

# Detecting Repetitions and Periodicities in Proteins by Tiling the Structural Space

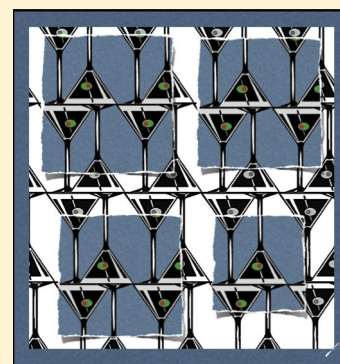
R. Gonzalo Parra,<sup>†</sup> Rocío Espada,<sup>†</sup> Ignacio E. Sánchez,<sup>†</sup> Manfred J. Sippl,<sup>‡</sup> and Diego U. Ferreiro<sup>\*†</sup>

<sup>†</sup>Protein Physiology Lab, Departamento de Química Biológica, Facultad de Ciencias Exactas y Naturales, UBA-CONICET-IQUIBICEN, Buenos Aires, Argentina

<sup>‡</sup>Center of Applied Molecular Engineering, Division of Bioinformatics, Department of Molecular Biology, University of Salzburg, Salzburg, Austria

## Supporting Information

**ABSTRACT:** The notion of energy landscapes provides conceptual tools for understanding the complexities of protein folding and function. Energy landscape theory indicates that it is much easier to find sequences that satisfy the “Principle of Minimal Frustration” when the folded structure is symmetric (Wolynes, P. G. Symmetry and the Energy Landscapes of Biomolecules. *Proc. Natl. Acad. Sci. U.S.A.* **1996**, *93*, 14249–14255). Similarly, repeats and structural mosaics may be fundamentally related to landscapes with multiple embedded funnels. Here we present analytical tools to detect and compare structural repetitions in protein molecules. By an exhaustive analysis of the distribution of structural repeats using a robust metric, we define those portions of a protein molecule that best describe the overall structure as a tessellation of basic units. The patterns produced by such tessellations provide intuitive representations of the repeating regions and their association toward higher order arrangements. We find that some protein architectures can be described as nearly periodic, while in others clear separations between repetitions exist. Since the method is independent of amino acid sequence information, we can identify structural units that can be encoded by a variety of distinct amino acid sequences.



## INTRODUCTION

“There is something breathtaking about the basic forms of crystals. They are in no sense a discovery of the human mind; they just “are”, existing quite independently of us. The most that man can do is become aware, in a moment of clarity, that they are there, and take cognizance of them.” M.C. Escher

Natural protein molecules are peculiar polymers. Unlike most of the random amino acid sequences, natural protein chains spontaneously find functional states by folding to a discrete collection of structures constituting a *native* state. Our deepest understanding of this phenomenon is grounded in the energy landscape theory of protein folding, which simplifies the complexity of folding to a few general descriptors of the configurational space.<sup>1,2</sup> These abstractions provide conceptual tools to infer reliable energy functions<sup>3</sup> and to build simple and powerful predictive models,<sup>4,5</sup> and most importantly, they provide a common language for the development (and healthy discussion!) of ideas.<sup>6,7</sup> The basic notion underlying these developments is the *principle of minimal frustration*:<sup>8</sup> in order to fold to a stable structure, a polymer must possess a funneled energy landscape.

According to energy landscape theory, proteins are information-bearing molecules that evolved to funneled energy surfaces, contrasting them to random heteropolymeric chains that have rugged energy landscapes.<sup>1</sup> Since amino acids in natural proteins generally appear to be distributed at random,<sup>9</sup> higher order correlations must be present in sequences that

result in stable folds. Energy landscape theory predicts that funneled landscapes and low energy structures are much easier to realize in the presence of symmetry as compared to asymmetric arrangements.<sup>10</sup> The identification of funneled energy landscapes as a general requirement for stable folds implies that patterns can form in different parts of the molecule with relative independence which subsequently assemble to higher order structures. This greatly reduces the search problem by efficiently arranging relatively small fundamental building blocks or “foldons”<sup>11</sup> in a repetitive fashion. The mere existence of repetitions or fundamental units does not guarantee that the system will be symmetric, but these units should arrange in particular ways and coalesce into higher order patterns. Hence, a periodicity guarantees a certain symmetry, but there can be repetitions without symmetry. Therefore, detecting repeated units and patterns is a first step toward an understanding of their assembly to complete structures and the emergence of symmetry. Such structural mosaics are accompanied by energy landscapes with multiple funnels embedded within each other.<sup>12–14</sup>

**Special Issue:** Peter G. Wolynes Festschrift

**Received:** February 28, 2013

**Revised:** June 9, 2013

**Published:** June 11, 2013

Several algorithms have been used to characterize repetitions in protein sequences.<sup>15,16</sup> Most methods are based on the self-alignment of the primary structure, while others implement spectral analysis of pseudochemical characteristics of the amino acids.<sup>15</sup> Since the same structural motif can be encoded by sequences that appear completely unrelated, it is not surprising that sequence-based methods fail to infer true structural repetitions when the sequence similarity is low. In contrast to sequence based methods, only a few methods for the detection of repetitions in protein structures are available. These usually search for repeats by aligning the structure against itself.<sup>17,18</sup> Some methods add sophisticated transformations of the alignment matrices that enhance the detection and characterization of structural repeats,<sup>19,20</sup> and machine learning provided a fast method to recognize repeat regions in solenoid structures.<sup>21</sup> Although many families of proteins with repeating motifs can be identified,<sup>16,22</sup> there is still no consensus on how to reconcile the often conflicting characterizations of repetitions in proteins<sup>15,23</sup> even for basic parameters such as the size of the repeating elements, the number and location of the occurrences, and the grouping of these into higher order patterns.

Here we develop basic concepts and methods for the detection and analysis of repeats in protein structures. Using a fast and robust structural alignment protocol and a proper metric,<sup>24</sup> we exhaustively analyze the repetition of every possible continuous fragment of a protein structure and define the portions that best describe the overall structure when this fragment is repeated, translated, and rotated exhaustively with respect to the complete molecule. The result is a tessellation of the whole protein in terms of a set of basic tiles. The tessellation lends itself to an intuitive visualization of the repeating units and their association into higher order patterns. We find that some architectures can be described as nearly periodic, while in some others clear separations between repetitions exist. Since this method is independent of sequence, it allows for comparison of recurring structures and tiles that represent a common structural motif that can be encoded by a variety of distinct sequence elements.

## METHODS

**Structural Alignments and Tiles.** For the characterization of repetitions and the identification of tiles in protein structures, we use the TopMatch tool.<sup>24,25</sup> Given a pair of protein structures, this algorithm generates an exhaustive list of partial alignments along with the transformations (rotations and translations) that maximize the superposition of equivalent C<sup>α</sup> atoms. The alignments are ranked according to the TopMatch score

$$S = \sum_i^L e^{-r_i^2/\sigma^2}$$

which provides a metric for structural similarity.<sup>26</sup> Here  $L$  is the length of the alignment and  $r_i$  is the euclidean distance between equivalent C<sup>α</sup> atoms. Basically,  $S$  is a function of the alignment length  $L$  and the structural deviation of the superimposed structural fragments, where the scaling factor  $\sigma$  determines the rate of reduction of  $L$  as a function of the structural deviation. Here we used  $\sigma = 6.35 \text{ \AA}$  as reported previously.<sup>24</sup> Proteins often contain recurrent structural motifs that can be considered as repetitions and variations of a basic structural unit. In order to detect this kind of structural repetition, we treat the structure

as a mosaic and try to decompose it into smaller units or *tiles* with the constraint that these tiles are all structurally similar to each other. In a protein, the possible tiles are not necessarily unique nor are they required to cover a chain completely. However, in any case, it is certainly possible to identify those tiles that, when repeated in a non-overlapping fashion, cover a maximum fraction of the structure.

Given a protein structure, every continuous fragment of the polypeptide is a possible tile. Hence, the length of tiles ranges from the sequence length  $N$  down to a single residue. Since the C<sup>α</sup> traces of tiles of one or a few residues are too small for meaningful comparisons, we use a lower tile length of six amino acid residues. In a protein of length  $N$ , there is one tile of length  $N$ , two tiles of length  $N - 1$ , and so on, and hence, the total number of tiles is  $N_T = \sum_{L=6}^N (N - L + 1)$ . Each of these tiles  $T_i$  is then used as a query in TopMatch to identify all other tiles  $T_k$  that are structurally similar to  $T_i$ . Each match is uniquely identified by its length  $L_{ik}$ , the location of its center  $Z_{ik}$ , and the associated score  $S_{ik}$ . The matches are then sorted by  $S_{ik}$ , where the self-alignment ( $i \equiv k$ ) necessarily has the highest score, since the respective alignment length is maximal and the structural deviation is zero. Hence,  $L_{ii} = S_{ii}$  i.e., the score obtained from an alignment of a tile with itself evaluates to the length  $L_{ii}$  of the alignment.

From the set of matches, we extract that subset of fragments that maximizes the sum over the scores  $C_i = \max \sum_k S_{ik}$  where any two tiles  $T_{k_1}$  and  $T_{k_2}$  that occur in the sum must not overlap. This sum defines the coverage  $C_i$  of tile  $T_i$  which was used to generate the matches. We define the associated tile score as

$$\Theta_i = \frac{C_i - L_{ii}}{N - L_{ii}} \quad (1)$$

which represents the fraction of the structural space that can be covered by repetitions of a given tile. When considering the ranked list of hits, there are several ways to define the set of non-overlapping alignments. In the most restrictive variant, we include only those repeats  $T_k$  for which the aligned region comprises the whole tile, i.e.,  $L_{ik} \equiv L_{ii}$ . A more flexible variant is to include all alignments where  $L_{ii}/2 < L_{ik} \leq L_{ii}$ , that is, when more than half of  $T_k$  matches  $T_i$ . In the latter case, we use the additional restriction that the first and last residues of any two tiles  $T_h$  and  $T_k$  in the optimal subset must not overlap.

**Homogeneous Model.** To evaluate the upper limits of the tiling scoring functions, we calculated the tile score  $\Theta_i$  expected for a homogeneous model, where the protein is represented as a finite linear string of amino acids. In this case, every alignment of tile  $T_i$  and repeat  $T_k$  has a perfect match, and thus, the alignment score  $S_{ik}$  will be equal to  $L_{ii}$ . Then, the coverage  $C_i$  is the product of  $L_{ii}$  and the number of tile copies  $n_c$  that can be accommodated which, depending on the tile center  $Z_{ii}$  is  $n_c = \lfloor N/L_{ii} \rfloor$  if the chain ends are covered or  $n_c = \lfloor N/L_{ii} \rfloor - 1$  if they are not.

When alignments with  $L_{ii}/2 < L_{ik} \leq L_{ii}$  are permitted, then  $\Theta_i$  has an extra term that takes into account the coverage at the chain ends:

$$\Theta_i = \frac{(n_c - 1) \cdot L_{ii} + C_{\text{beg}} \cdot \chi(C_{\text{beg}} - L_{ii}/2) + C_{\text{end}} \cdot \chi(C_{\text{end}} - L_{ii}/2)}{N - L_{ii}} \quad (2)$$

where  $\chi(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ ,  $n_c$  is the number of full length tile copies that can be accommodated along the protein, and  $C_{\text{beg}}$  and  $C_{\text{end}}$  are the maximum number of amino acids left uncovered by the copies at the limits of the protein, and can be calculated as

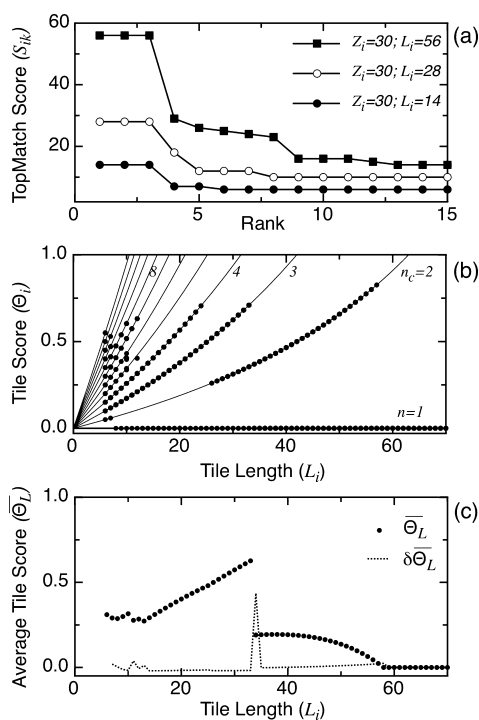
$$C_{\text{beg}} = Z_i + \left[ \frac{1}{2} - \frac{Z_i}{L_i} \right] \cdot L_i - \frac{L_i}{2} \quad (3)$$

$$C_{\text{end}} = N - \left[ Z_i + \left[ \frac{N}{L_i} - \frac{1}{2} - \frac{Z_i}{L_i} \right] \cdot L_i + \frac{L_i}{2} \right] \quad (4)$$

Further details of this model can be found in the Supporting Information.

## RESULTS AND DISCUSSION

To illustrate the characteristic properties of tessellations of protein structures, we use the protein “4ank” (pdb: 1n0r,  $N = 126$  residues) which is a synthetic construct of canonical ankyrin repeats.<sup>27</sup> Figure 1a shows the scores of the top 15 hits for three different fragments of the structure used as the query tile  $T_i$ . In all instances, the highest ranking tile corresponds to the self-alignment ( $i \equiv k$ ) and in each of these cases there are



**Figure 1.** Scoring the tiles: continuous portions of a protein model (pdb: 1n0r, A) are selected and structurally aligned to the whole protein. A ranked list of alignments is generated for every fragment according to TopMatch score ( $S_{ik}$ ), three of which are shown in panel a.  $L_i$  is the length of the fragment in amino acid units and  $Z_i$  is the center, according to the numbering scheme of the  $C^\alpha$  atoms of the pdb. (b) Distributions of tile score  $\Theta_i$  for every tile length ( $L_i$ ). Each point corresponds to the experimental values obtained when perfect matching ( $S_{ik} = L_i$ ) is restricted. The lines correspond to the expression  $\Theta_i = (n_c - 1)L_i/(N - L_i)$ , with  $N = 126$  and the number of tile copies that can be repeated is  $n_c = 1, 2, 3, \dots, 12$  as indicated. (c) The points correspond to the average  $\bar{\Theta}_L$  calculated for every  $L_i$ . The dotted line is the difference between consecutive points  $\delta\bar{\Theta}_L$ .

two tiles ( $i \neq k$ ) that yield nearly perfect matches. For the subsequent tiles, the score drops rapidly.

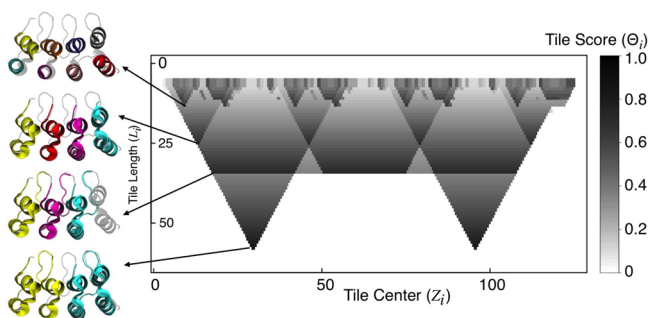
Next we use the ranked list to pick out non-overlapping fragments in order to cover the protein structure as repeats of tile  $T_i$ . For each possible tile  $T_j$ , the tile score  $\Theta_i$  is calculated as described above. Clearly, the tile score  $\Theta_i$  for tiles with  $L_{ij} > N/2$  is always zero, as no repetitions of such fragments are possible (Figure 1b). The largest tile that can be repeated twice has  $L_i = 57$  amino acids. Tiles nested within these largest tiles necessarily have smaller scores. Three repeats are observed for  $L_i = 33$ , and four for  $L_i = 24$  (Figure 1b). The peaks in Figure 1b correspond to the largest fragments that occur more than once and for which each of its extensions occurs fewer times; that is, they are *maximal elements*. The steady decrease in  $\Theta_i$  results from fragments that are nested within the maximal ones. This can be inferred from the homogeneous model where a group of tiles that occurs  $n_c$  times yields the tile score  $\Theta_i = (n_c - 1)L_i/(N - L_i)$ .

The fact that there is a number of tiles of similar score  $\Theta_i$  but varying length  $L_i$  implies that the overall protein architecture can be covered by a set of nested tiles. Hence, the question arises which of the possible tile lengths yields a tessellation of maximum coverage. In the case of real proteins, copies of individual tiles generally exhibit structural variations with respect to a basic tile. Such variations reduce the score  $S_{ik}$  of the respective structural matches. The relative reduction is generally much more pronounced for small tiles as compared to larger tiles which may result in a relatively large decrease of the overall tile score  $\Theta_i$ . In short, if the various copies of small tiles have relatively large structural deviations, then the associated tile score  $\Theta_i$  may appear suboptimal with respect to tile scores obtained from larger tiles. It is therefore convenient to take the average  $\bar{\Theta}_L$  over all tile scores  $\Theta_i$  that have the same tile length  $L_i$  (Figure 1c). In the example, it is evident that the maximum occurs at  $\bar{\Theta}_L = 33$  residues, indicating that tiles of this size tessellate the structure in an optimal way. Formally, the optimal length is obtained as a root of the derivative  $d\bar{\Theta}_L/dL_i$ ; i.e., it can be obtained from the finite differences  $\Delta\bar{\Theta}_L = \bar{\Theta}_L - \bar{\Theta}_{L-1}$ . Note that this identifies the optimal tile length  $L$  but not the particular tile  $T_i$  that optimizes the tessellation.

Since a particular tile  $T_i$  is characterized by the position of its center  $Z_i$  along the amino acid sequence and a match between two tiles  $T_i$  and  $T_k$  by the respective alignment length  $L_{ik}$ , the multitude of tessellations of a particular structure is representable in two dimensions and the associated score  $\Theta_i(L_{ij}, Z_i)$  can be indicated by shades of gray (Figure 2). Such representations show how copies of each of the possible tiles cover the whole structure. In the case of 1n0r, the structure is covered by two repeats of 57 amino acids, centered at residues 30 and 96. These repeats decay into two smaller repeats of 24 amino acids, where the decomposition results in a loss of approximately 12% coverage. These tiles in turn consist of two smaller tiles of 8 and 10 amino acids. The latter correspond to two  $\alpha$ -helices that are part of the canonical ankyrin motif (Figure 2).

A peculiar phenomenon is apparent for tiles of length  $L_i = 33$ . Any tile of this size provides a nearly complete tessellation of the structure. Moreover, at this length scale, the tiles are separated by a distance that is equal to the size of the tile itself. Hence, the whole structure has the characteristics of a wave. The characteristic wavelength is  $L = 33$ , and the structure can be completely covered starting with any phase  $\phi = 0, 1, \dots, L - 1$ . Taken together, these observations imply that those tiles





**Figure 2.** Tiling a highly symmetric protein: a designed ankyrin-repeat protein (pdb: 1n0r, A) was fragmented in 7381 different tiles. These are ordered according to their size (vertical axis) and their center (horizontal axis) in amino acid units. The tile score  $\Theta_i$  of each one is displayed in grayscale. The structures of the protein and the respective tiling at different  $L_i$  are shown on the left. The native structure is colored gray, and superimposed to it is the selected tile (yellow); the copies of it are colored cyan, magenta, red, etc.

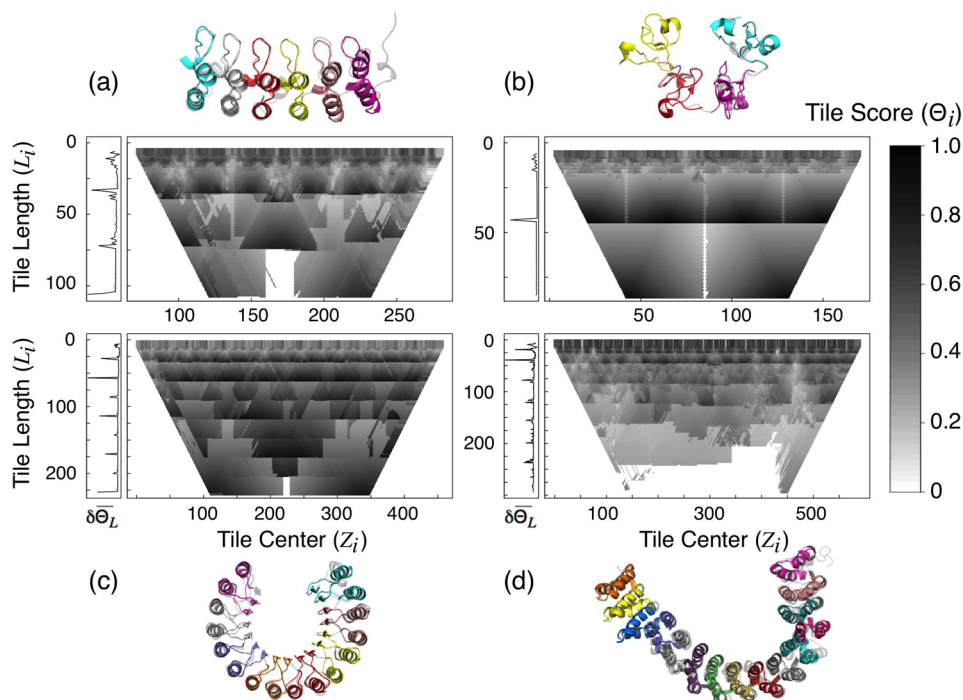
optimally cover a repetitive protein structure whose average score  $\overline{\Theta}_L$  is a maximum, and it seems that such maxima are accompanied by a large value of  $\Delta\overline{\Theta}_L$  (Figure 1b). From the set of tiles that contribute to  $\overline{\Theta}_L$ , we may define the most typical tile as that particular tile  $T_i$  that has the largest score  $\Theta_i(L_p, Z_i)$  with respect to all other tiles  $T_k(L_p)$  in this set.

Repeats in protein structures are thought to be the result of duplication of amino acid sequences. In general, a duplication results in an exact copy of the duplicated material. On the level of amino acid sequences, the similarity among the copies decays in time due to the accumulation of amino acid substitutions,

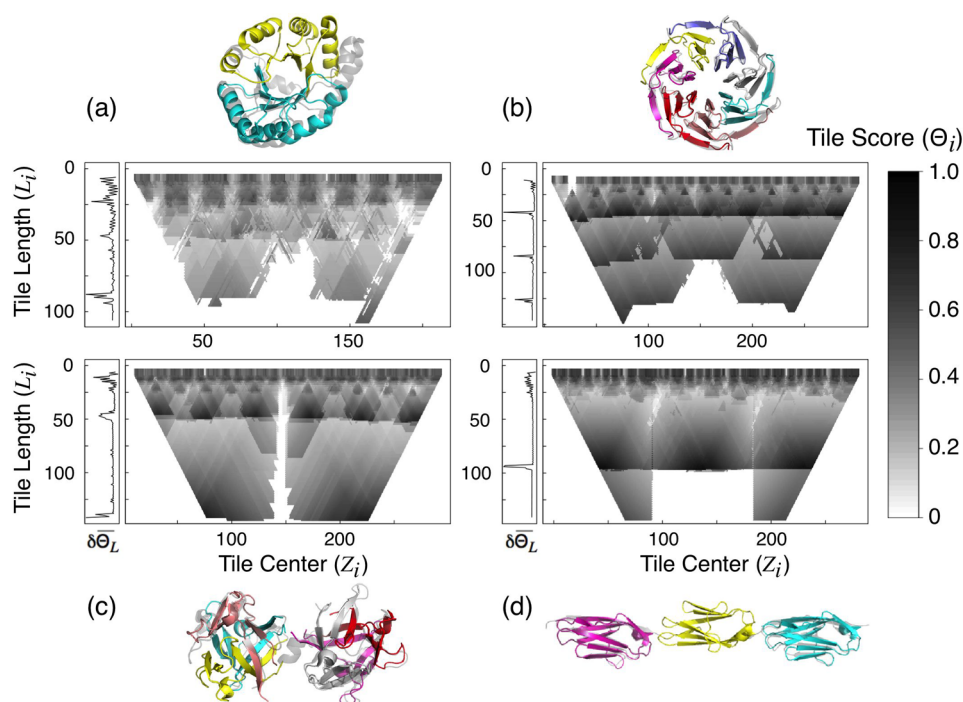
insertions, and deletions. The respective structures are more robust in the sense that the similarity among the sequences decays much faster as compared to the similarity among the polypeptide backbone. Nevertheless, insertions, deletions, and other events also affect the three-dimensional structures of the individual copies, and therefore, in natural proteins, structural repeats are rarely exact and they are often interspersed by nonrepetitive regions. In what follows, we discuss tessellations obtained for a broad variety of protein structures. This method does not rely on visual inspection. We define the characteristic frequency at the highest peak in  $\Delta\overline{\Theta}_L$ , and the basic tile-unit as the one that scores the highest  $\Theta_i$  at this  $L_i$ . The nonrepetitive regions found in these tessellations are marked as insertions (Table S1, Supporting Information).

**Tessellations of Classical Repetitive Proteins.** Many natural proteins contain tandem repeats of similar amino acid stretches. They are broadly classified in groups according to the length of the minimal repeating unit. Short repeats of up to five residues usually form fibrillar structures such as collagen or silk, while repeats longer than about 100 residues frequently fold independently as globular domains.<sup>16,28</sup> There is a class of repeat proteins that lies in between these for which folding of the repeating units is coupled and “domains” are not obvious to define.<sup>29</sup> Since defects in the regularities of the repeating array are likely to affect the folding transitions and the biological function, we aimed at defining these from a purely geometrical perspective using the tiling approach described above.

$\text{I}\kappa\text{B}\alpha$  is an ankyrin repeat containing protein that binds to and inhibits the transcription factor NF- $\kappa\text{B}$ .<sup>30</sup> The fragmentation and tiling procedure correctly identifies a characteristic 33 amino acids length corresponding to the canonical ankyrin repeat size (Figure 3a). We found deviations from this



**Figure 3.** Tiling classical repeat-containing proteins. The tiling profile is shown in grayscale, together with the  $\delta\overline{\Theta}_L$  projected on the left. The structures of the native protein and the highest scoring tiling at the characteristic frequency are shown, using the same coloring scheme of Figure 2. The length ( $L_i$ ) and center ( $Z_i$ ) of the selected tile are (a) ankyrin repeat:  $\text{I}\kappa\text{B}\alpha$  (pdb: 1nfi, E)  $L_i = 33$ ,  $Z_i = 191.5$ ; (b) hevein: wheat-germ agglutinin (pdb: 1k7u, A)  $L_i = 43$ ,  $Z_i = 150.5$ ; (c) leucine-rich: porcine ribonuclease inhibitor (pdb: 2bnh, A)  $L_i = 57$ ,  $Z_i = 139.5$ ; and (d) HEAT: PR65/A (pdb: 1b3u, A)  $L_i = 39$ ,  $Z_i = 530.5$ .



**Figure 4.** Tiling globular repeat-containing topologies. The tiling profile is shown in grayscale, together with the  $\overline{\delta\Theta_L}$  projected on the left. The structures of the native protein and the highest scoring tiling at the characteristic frequency are shown, using the same coloring scheme of Figure 2. The length ( $L_i$ ) and center ( $Z_i$ ) of the selected tile is are (a) TIM barrel (pdb: 1fq0, A)  $L_i = 88$ ,  $Z_i = 60$ ; (b)  $\beta$ -propeller (pdb: 3ow8, A)  $L_i = 42$ ,  $Z_i = 200$ ; (c) trefoil (pdb: 1ybi, A)  $L_i = 46$ ,  $Z_i = 176$ ; and (d) Ig-repeats (pdb: 2rik, A)  $L_i = 94$ ,  $Z_i = 140$ .

canonical size ranging from 30 to 39 residues, indicating that not all the ankyrin repeats are geometrically equivalent. Fragments with highest scores can be placed six times, covering about 92% of the structure (Table S1, Supporting Information). It is apparent that the most C-terminal repetition is distorted relative to the others, as the  $\Theta_i$  corresponding to this region are lower. The grouping of consecutive repeats at bigger  $L_i$  segregate pairs where the central one scores best, indicating that the insertions detected at length 33 distort the symmetry of the array at a higher length scale. Maybe it is no coincidence that this protein was shown to fold *in vitro* in three consecutive transitions roughly corresponding with the pairing of repeats at  $L_i = 70$ .<sup>31,32</sup>

The monomeric chain of wheat-germ agglutinin has been described to contain four hevein subdomains.<sup>33</sup> The tiling approach detects that this structure can be composed with two tiles of  $L_i = 86$  amino acids, as well as four repetitions of  $L_i = 43$ , both covering 100% of the structure (Figure 3b). Taking the average of the  $\Theta_i$  at each  $L_i$  points that a discontinuity occurs at size 43, defining a characteristic frequency. At this size, most tiles are equally good in covering the structural space with three repetitions. The highly symmetric disposition of the four best tiles at this length scale makes the whole structure appear nearly periodic, and a preferred phase is determined by the N and C termini of the chain.

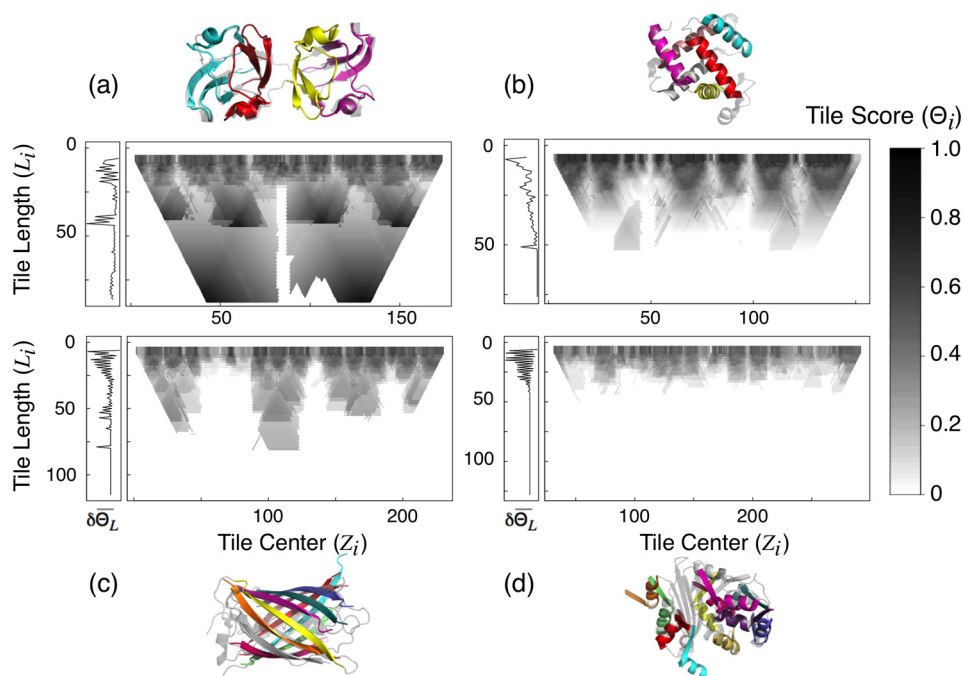
Porcine ribonuclease inhibitor is a leucine-rich repeat protein for which 16 consecutive repetitions were defined in its sequence. Although very similar at the primary level, these repeats are not structurally equivalent. We detect that there are two different types of tiles, each consisting of 28 and 29 amino acids (Figure 3c). Moreover, we found that these are alternated along the structure, appearing as a square-tooth pattern at this length scale (Figure S3, Supporting Information). Since these

units are arranged in a symmetric fashion, the structure can be represented as well by bigger fragments (Figure 3c). At the length of  $L_i = 57$  residues, almost every fragment repetition is as good as others in explaining the overall structure. Thus, the repeating length is better described with two canonical leucine-rich repeats. It is striking to note that Haigis et al. previously identified a 57-residue repeat as the evolutionary unit of this protein by analyzing the exon boundaries of the primary transcripts.<sup>34</sup>

The scaffolding subunit of protein phosphatase 2A, PR65/A, is a large repeat-protein of the HEAT class.<sup>35</sup> The tiling procedure detects the best tile at size 39 amino acids and identifies 15 copies of it in the structure, coincident with the detection in amino acid sequence patterns of the HEAT motif (Figure 3d). This protein exhibits an overall superhelical structure, yet irregularities in the array cause unevenness in the grouping of consecutive repeats at higher length scales. The periodic packing of HEAT repeats is interrupted between repeats 3 and 4 ( $Z_i = 117$ ) and between 12 and 13 ( $Z_i = 471$ ).<sup>35</sup> This is reflected at higher  $L_i$  where the tiles centered around amino acid 300 display consistent higher scores, indicating that the central repeats are more symmetrically arranged than the terminal ones (Figure 3d).

**Tessellations of Globular Proteins.** In contrast to the solenoidal architectures usually acquired by classical repeat-proteins, some protein folds display point rotational symmetries. Often the N and C terminal repetitions come in contact, closing up the structure in polyhedral-like forms. We investigated how the tiling procedure identifies structural repetitions and tessellation patterns in some of the most common topologies of this kind.

The TIM barrel is one of the most common folds among monomeric enzymes.<sup>36</sup> This is typically described as a



**Figure 5.** Tiling classical globular proteins. The tiling profile is shown in grayscale, together with the  $\delta\overline{\Theta}_L$  projected on the left. The structures of the native protein and example tilings are shown, using the same coloring scheme of Figure 2. The length ( $L_i$ ) and center ( $Z_i$ ) of the selected tile are (a)  $\beta\gamma$ -crystallin (pdb: 1h4a, X)  $L_i = 43$ ,  $Z_i = 149.5$ ; (b) myoglobin (pdb: 1mbd, A)  $L_i = 18$ ,  $Z_i = 29$ ; (c) green fluorescent protein (pdb: 1gfl, A)  $L_i = 15$ ,  $Z_i = 182.5$ ; and (d)  $\beta$ -lactamase (pdb: 4blm, A)  $L_i = 15$ ,  $Z_i = 185.5$ .

collection of  $\beta$ - $\alpha$  motifs linked by variable loops that close up a cylinder of parallel  $\beta$ -strands surrounded by a layer of  $\alpha$ -helices. There is a relatively high structural conservation among proteins of this type, yet their sequences can appear unrelated, opening room for discussion about the nature of the repeating units and their arrangement.<sup>37</sup> We applied the tiling procedure on some of the most discussed cases, and for most, we detect signals for 2, 4, and 8 repeats (Table S1 and Figure S4, Supporting Information). Not all the TIM barrels showed the same characteristic frequency. Some of the structures are best described with fragments that correspond to half barrel (Figure 4a), while others displayed comparable signals at sizes corresponding to half or quarter barrel (Figure S4, Supporting Information). The most irregular examples have characteristic frequencies at even lower length scales (Table S1, Supporting Information). On the basis of amino acid sequence patterns, Soding et al. annotated equivalent deviations in this topological family.<sup>37</sup>

Several proteins can be grouped into the  $\beta$ -propeller class. These contain a variable number of radially arranged antiparallel  $\beta$ -sheets appropriately named “blades”.<sup>38</sup> We identify that in most cases the best tiles distinguish this motif and annotate 4-, 5-, 6-, and 7-bladed propellers (Figure S5, Supporting Information), even when a nonpropeller domain is present in the same polypeptide chain (Figure S5d, Supporting Information). An interesting exception occurs in the subclass of WD-repeat propellers where the selected tile does not correspond with a blade (Figure 4b). In this case, we detect a characteristic frequency of  $L_i = 42$  amino acids, with tiles repeated seven times and contributing three strands to one blade and one strand to the next one (Figure 4b). Notably, this particular phase was the one originally described when no structure of members of this class was known.<sup>39</sup>

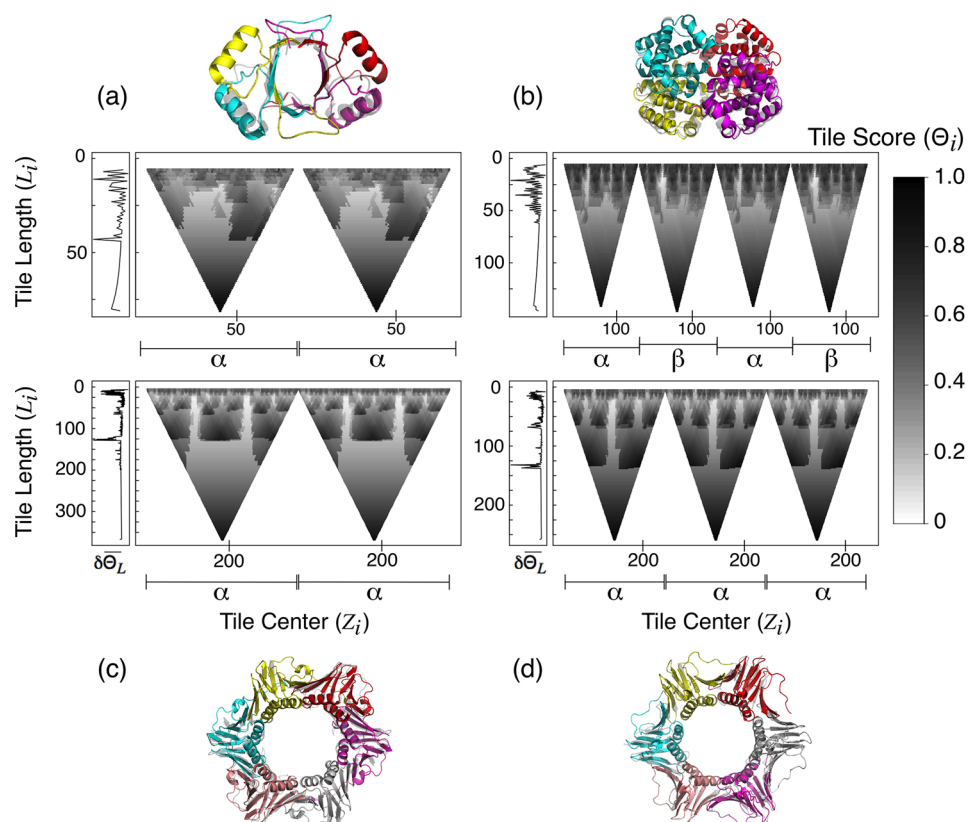
The hemagglutinating protein HA33 from *Clostridium botulinum* is a neurotoxin-associated protein that folds in an appealing topology of two consecutive  $\beta$ -trefoil subdomains (Figure 4c). The characteristic frequency ( $L_i = 142$ ) points to two fragments that have the highest  $\Theta_i$  and correspond to the tiles of each subdomain. The best phase at the second peak ( $L_i = 46$ ) corresponds to tiles that can be fitted three times in each subdomain and match the annotated foil of the  $\beta$ -trefoil architecture.

Surveying other architectures with repeating motifs, we noted that in some cases the highest scoring tiles are at the characteristic frequency. Figure 4d shows the results for a fragment of titin that contains three tandem immunoglobulin-like (Ig) domains. At  $L_i = 94$  amino acids, the best phase coincides with the Ig domains. The fact that other phases also score high at this length scale is indicative that the arrangement between the Ig domains is regular, as if this were not the case, those fragments would not display that high  $\Theta_i$ .

At some level, all proteins are formed by repetitions of amino acids. The symmetry of the backbone interactions in secondary structures was key to the Pauling and Corey proposal of these arising from the regular repetition of planar peptide bonds.<sup>40,41</sup> Recurrent secondary structure motifs were once candidates for fundamental building blocks of globular domains, in line with the success of structure prediction by fragment assembly.<sup>42–44</sup> Since repetitions can be confidently found by tiling the structural space, we explored to what extent any given protein structure can be said to be composed with tiles, illustrating with some classical examples.

Synthesized at embryonic stages and hopefully lasting soluble for a lifetime,<sup>45</sup>  $\beta\gamma$ -crystallins lens proteins increase the refractory index and maintain transparency of the vertebrates’ eyes. Since its initial description, it has been a clear example of structural motifs coalescing into higher order patterns.





**Figure 6.** Tiling quaternary complexes. The tiling profile is shown in grayscale, together with the  $\delta\bar{\Theta}_i$  projected on the left. The structures of the native protein and example tilings are shown, using the same coloring scheme of Figure 2. The length ( $L_i$ ) and center ( $Z_i$ ) of the selected tile are (a) homodimeric HPV-16 E2c (pdb: 1r8p)  $L_i = 43$ ,  $Z_i = 58.5$ ; (b) deoxy-hemoglobin (pdb: 2hhb)  $L_i = 141$ ,  $Z_i = 71.5$  from chain A; (c) the  $\beta$ -subunit of *Thermotoga maritima* DNA polymerase III (pdb: 1vpk)  $L_i = 128$ ,  $Z_i = 297$ ; and (d) the processivity factor of *Saccharomyces cerevisiae* DNA polymerase- $\delta$  (pdb: 1plq)  $L_i = 132$ ,  $Z_i = 190$ .

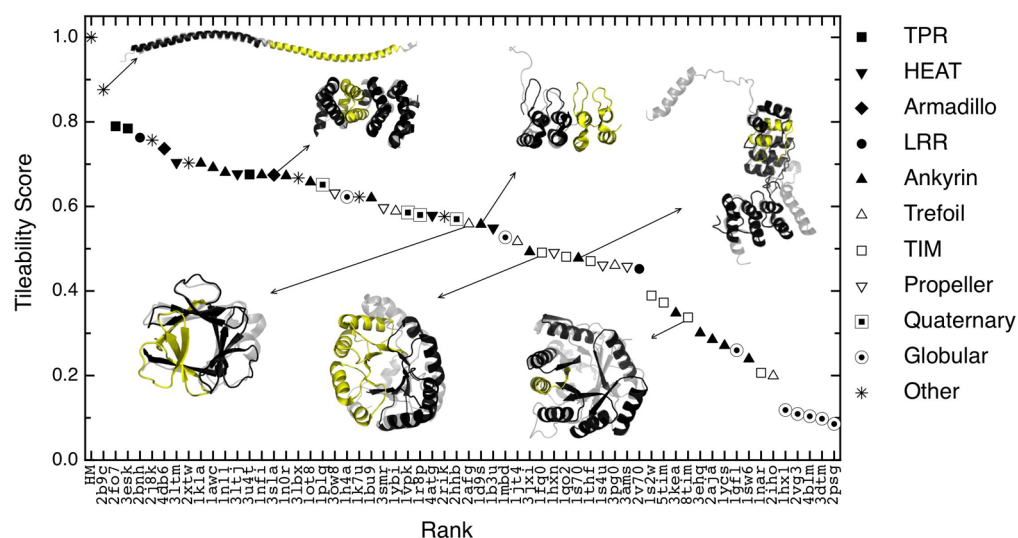
Coincident with the classical descriptions of these folds, tiling the structural space detects that this protein can be very well described with two repetitions of an eight-stranded  $\beta$ -barrel of  $L_i = 87$  amino acids centered at positions  $Z_i = 44.5$  and  $Z_i = 133.5$  (Figure 5a). In turn, each of these can be composed with two units of about 40 residues that correspond to the Greek-key motif, that can be further decomposed into three 10-residue  $\beta$ -strands. The characteristic frequency is at  $L_i = 43$  amino acids, selecting out the Greek-key as the repetition we annotate. It is apparent that there are irregularities in the structure that make the second and fourth Greek-keys have a higher  $\Theta_i$  than the others and indeed different maximal  $L_i$ .

About 70% of the mean structure of myoglobin, the hydrogen atom of biology,<sup>46</sup> can be described with six copies of an 18 amino acid fragment. This corresponds to the “B”  $\alpha$ -helix, and constitutes a maximal fragment. The score at higher length scales decreases rapidly (Figure 5b). In this case, we could not detect a relevant frequency above the  $\alpha$ -helical segments, indicating that these do not contiguously repeat in a highly symmetrical way, a fact that strongly surprised Kendrew et al. when they solved the crystal structure.<sup>47</sup>

Green fluorescent protein folds as a  $\beta$ -barrel with a coaxial helix, with the fluorophore forming from the central helix.<sup>48</sup> We identify fragments of  $L_i = 15$  that can cover about 71% of the structural space with 11 repetitions, corresponding with  $\beta$ -strands (Figure 5c). At higher length scales, no fragment significantly raises the signal.

*Bacillus licheniformis*  $\beta$ -lactamase illustrates an example of a mixed  $\alpha\beta$  topology, composed of two discontinuous subdomains.<sup>49</sup> Here again, there is no particular length scale at which a useful characteristic frequency can be defined (Figure 5d). The best tiling occurs at  $\Theta_i = 15$  where the fragment corresponds to 1 of 10  $\alpha$ -helices and covers 74% of the structural space when repeated.

**Tessellations of Oligomers.** In their natural environment, most of the polypeptide chains of living organisms are not found folded as spheroidal monomers but typically come together, forming oligomeric complexes with two or more subunits. Most frequently, they form homodimeric complexes, but hetero-oligomers are not uncommon and even thousands are to be found. The symmetrical basis of this phenomenon has been explored even before the first protein structures were solved.<sup>50</sup> A recent survey estimates that over 95% of the homodimeric complexes crystallized are symmetric,<sup>51</sup> and it is expected that small insertions and deletions can have profound effects on protein functionality, modulating oligomer stability, specificity, and aggregation.<sup>52</sup> To analyze the details of symmetry in multichain complexes, we can first define the elementary blocks that constitute the array. To explore this, we applied the same procedure of fragmenting and tiling described above now using the quaternary arrangements of subunits as the target structure to cover. If the monomers that form a homo-oligomer cannot be decomposed into significant tiles, we expect the best tile to correspond to the monomeric chain itself. Indeed, we find this is the case for the majority of



**Figure 7.** Tileability of protein structures. The tiling procedure was applied to the protein models (indicated with their PDB code, HM: homogeneous model) and ranked according to their tileability score  $\Xi$ . Example tessellations are shown with the tile-unit colored yellow and the copies colored black, superimposed to the native structure in gray. Filled symbols: solenoidal repeat proteins. Empty symbols: globular repeat proteins.

the oligomers we evaluated. We noted however interesting cases in which the subunits can be decomposed into significant tiles.

Papillomavirus E2c-DNA binding protein is a remarkable model to study sequence specific recognition.<sup>53</sup> This domain is composed of two identical chains that come together, forming in a  $\beta$ -barrel architecture that exposes four  $\alpha$ -helices. The tiling procedure identifies an 81-residue fragment as the best scoring fragments, corresponding to the monomeric chains (Figure 6a). However, these can be further decomposed in tiles of  $L_i = 43$ , covering about 90% of the structural space. The best tile at this frequency corresponds to a  $\beta\alpha\beta$  motif that intertwines in each monomer and together contributes half  $\beta$ -barrel (Figure 6a).

Hemoglobin (the helium atom of biology?) is the prime example of a symmetrical quaternary arrangement, a tetramer of  $\alpha_2\beta_2$  chains. Figure 6b shows a regular tiling pattern in which four nearly identical regions can be distinguished. This highlights the long-established structural identity of the  $\alpha$  and  $\beta$  chains. As in the case of myoglobin, no significant decomposition of the structure can be made with continuous fragments.

On occasion, protein structures reveal geometrical chances and necessities of their history. Figure 6 shows the structures of the  $\beta$ -subunit of an archaeal DNA polymerase III (a homodimer, Figure 6c), together with the processivity factor of eukaryotic DNA polymerase- $\delta$  (a homotrimer, Figure 6d). Tiling these quaternary complexes identifies the subunits and further points to similar characteristic frequencies of  $L_i = 128$  and  $L_i = 132$ . In both cases, the chosen tiles at their respective  $L_i$  cover about 94% of the structure of the complexes. It is apparent that a DNA clamp of this kind can be constructed with either two or three polypeptide chains, each containing three or two tiles, that pack in a 6-fold rotational fashion.<sup>24</sup> This common tile can be further decomposed into two tiles of  $L_i = 65$  amino acids yet compromising about 10% coverage. It is interesting to note that these smaller fragments get intertwined when forming a higher order structure, unlike any other of the maximal fragments identified.

## CONCLUSIONS

Foldable sequences with funneled landscapes are easier to find if the low energy structure is symmetric.<sup>10</sup> Modern natural philosophers appreciate the existence of symmetry as an emergent feature of the parsimony of nature, resulting from the limited modes of interaction between a small number of elementary parts assembling into higher order structures<sup>50,54–56</sup> It is the inexact symmetries of biological molecules that are most striking.<sup>10,54</sup> Subtle aperiodicities can give rise to big biological effects,<sup>57</sup> and thus, their modulation can be at the core of the physiological workings of these “frozen accidents”.

In order to detect and characterize repetitions in protein structures, we presented a simple scheme based on analyzing the distribution of suboptimal structural alignments of continuous fragments. The procedure identifies maximal fragments, those for which any extension occurs fewer times in the ensemble of solutions (Figure 1B). By counting the number of occurrences of non-overlapping fragments and having a good metric for the overall coverage, we defined a score that ranks how a structure can be tessellated with similar, though not identical, fragments. We found that in most cases there is a defined fragment length at which the coverage gained by the repetitions is highest, defining a characteristic frequency. In some cases, there is a discrete collection of fragments that allows one to unequivocally define a best phase. In these cases, the repeat unit, the number of occurrences, and their boundaries can be confidently defined (Table S1, Supporting Information). In other cases, there are several equivalent phases at the characteristic frequency, pointing to structures that can be considered almost periodic and where the definition of a basic tile must remain arbitrary (Figure 2 and Table S1, Supporting Information). This is a common theme in the cases of solenoidal proteins where different researchers have defined the repeat unit at distinct frequencies and phases.<sup>23</sup> Including other information beyond geometry could indicate if there is a *biologically* preferred phase, such as the characterization of insertion sites, the variability in orthologous sequences, exon boundaries, or folding mechanisms.<sup>58</sup>



Proteins in which the repeats pack symmetrically against each other but do not translate along an axis can form closed structures. The fragmenting and tiling approach can be readily applied to such topologies like barrels, propellers, trefoils, and so on. Within these, we can distinguish nested repeating units and even resolve fine geometrical differences (Figure 4 and Table S1, Supporting Information). If the fundamental tiles are arranged symmetrically, then there must be larger tiles which are multiples of these basic tiles. These higher order tiles appear as additional maxima of  $\Theta_{L_i}$  toward larger  $L_i$  as compared to the basic tile. This hierarchical nesting of tiles can be captured by a tessellation score that is computed in the following way. For each tile length  $L_i$ , take the maximum tile score  $\Theta_i$  (e.g., the maximum score for a particular  $L_i$  in Figure 1b) and take the average over all  $L$ . This *tileability* score ( $\Xi$ ) is 1.0 for the homogeneous model, approaches 1 for highly regular structures like  $\alpha$ -helices, and goes to zero for nonrepetitive structures. In Figure 7, a variety of proteins are ranked by their respective tileability score  $\Xi$  (Table S1, Supporting Information). The largest value of  $\Xi$  is obtained for a long  $\alpha$ -helix from a coiled coil. The helix is followed by several solenoidal proteins with the most regular designed proteins ranking higher than the more irregular natural ones. These structures are followed by repetitive proteins with an overall globular shape. At the end of this scale, we find typical globular domains that do not have any periodicity larger than a few residues. We note that not all members of a particular topology group together; they rather get segregated according to the irregularities they display (Figure 7).

The same tiling procedure can be applied at the level of protein complexes, analyzing the details of how fragment copies between chains cover the structural space. At this level, we found that the best tiles often correspond with the monomeric chains or classical globular domains within them. However, interesting exceptions can mark chains that can be further decomposed into smaller units (Figure 6). It will be appealing to extend this now limited survey and characterize how frequently the distribution of geometric tiles coincides with the polypeptide chains, globular domains, exons boundaries, foldons, or motifs.

It is tempting to speculate about the functional consequences that the symmetrical distribution of similar fragments can have at different length scales. Energy landscape theory *modus operandi* appreciates that packing subunits in symmetrically equivalent ways give rise to structures with similar free energies, allowing multiple funnels to coexist in the energy landscape<sup>59</sup> and small perturbations to switch between these states.<sup>60</sup> Symmetry has been pointed out as being the key in other functional phenomena such as folding cooperativity, multiple ligand binding, thermodynamic stability, coding compression, and finite assembly.<sup>55,50</sup> Symmetric organization is an easy (and perhaps unavoidable) way for allostery to emerge.<sup>61,62</sup> Repetitions with point symmetries give rise to closed arrays such as barrels and the like at the tertiary level, and rings or polyhedra at the quaternary level. Helical symmetries form solenoids at the tertiary level that correspond with tubular organizations at the quaternary level. Nucleation and capping of these repeating arrays is often pointed to be critical to their physiological behavior both at the tertiary and quaternary levels. Potentially unbounded periodicity may require other mechanisms to terminate growth. It is thus not surprising that physiological workings and pathological states are the result of

aggregation of similar fragments, such as cytoskeleton dynamics,<sup>63</sup> epigenetic phenomena,<sup>64</sup> sickle-cell anemia,<sup>65</sup> and amyloid-related processes.<sup>66</sup>

The organization of protein molecules can be appreciated at many levels, from amino acid sequence motifs to dynamic interacting networks of thousands of components.<sup>63</sup> As the relevant contributions of the physical forces change at different length and time scales, the organizational agencies at each level will necessarily differ, but some common principles may underlie. The concepts postulated by energy landscape theory can be a guide in such a search.<sup>67–69</sup>

## ■ ASSOCIATED CONTENT

### 📄 Supporting Information

The tile and tessellation parameters for all the surveyed structures (Table S1), the derivation of the homogeneous model, and figures of tilings and tessellations of example proteins. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [ferreiro@qb.fcen.uba.ar](mailto:ferreiro@qb.fcen.uba.ar).

### Notes

The authors declare no competing financial interest.

## ■ ACKNOWLEDGMENTS

The works of Peter G. Wolynes transcended fields and provided us with deep impressions and equations to appreciate nature's beauty and comprehend its rich complexity. Enumerating his vast production and scientific achievements does not come close to the illuminating experience the lucky of us had in investigating with him. To Peter then we raise our Martini in *funneled glasses* and repeat "Salud!". This work was supported by the Consejo Nacional de Investigaciones Científicas y Técnicas de Argentina (CONICET), the Agencia Nacional de Promoción Científica y Tecnológica (ANPCyT), and by FWF Austria, grant number P21294-B12. R.G.P. and R.E. hold fellowships from CONICET, and I.E.S. and D.U.F. are Career Investigators.

## ■ REFERENCES

- (1) Bryngelson, J. D.; Onuchic, J. N.; Socci, N. D.; Wolynes, P. G. Funnel, Pathways, and the Energy Landscape of Protein Folding: a Synthesis. *Proteins* **1995**, *21*, 167–195.
- (2) Wolynes, P. G. Recent Successes of the Energy Landscape Theory of Protein Folding and Function. *Q. Rev. Biophys.* **2005**, *38*, 405–410.
- (3) Wolynes, P. G. Energy Landscapes and Solved Protein-folding Problems. *Philos. Trans. R. Soc., A* **2005**, *363*, 453–464.
- (4) Davtyan, A.; Schafer, N. P.; Zheng, W.; Clementi, C.; Wolynes, P. G.; Papoian, G. A. AWSEM-MD: Protein Structure Prediction Using Coarse-Grained Physical Potentials and Bioinformatically Based Local Structure Biasing. *J. Phys. Chem. B* **2012**, *116*, 8494–8503.
- (5) Zheng, W.; Schafer, N. P.; Davtyan, A.; Papoian, G. A.; Wolynes, P. G. Predictive Energy Landscapes for Protein-Protein Association. *Proc. Natl. Acad. Sci. U.S.A.* **2012**, *109*, 19244–19249.
- (6) Wolynes, P. G.; Eaton, W. A.; Fersht, A. R. Chemical Physics of Protein Folding. *Proc. Natl. Acad. Sci. U.S.A.* **2012**, *109*, 17770–17771.
- (7) Oliveberg, M.; Wolynes, P. G. The Experimental Survey of Protein-Folding Energy Landscapes. *Q. Rev. Biophys.* **2005**, *38*, 245–288.

- (8) Bryngelson, J. D.; Wolynes, P. G. Spin Glasses and the Statistical Mechanics of Protein Folding. *Proc. Natl. Acad. Sci. U.S.A.* **1987**, *84*, 7524–7528.
- (9) Weiss, O.; Jimenez-Montano, M. A.; Herzog, H. Information Content of Protein Sequences. *J. Theor. Biol.* **2000**, *206*, 379–386.
- (10) Wolynes, P. G. Symmetry and the Energy Landscapes of Biomolecules. *Proc. Natl. Acad. Sci. U.S.A.* **1996**, *93*, 14249–14255.
- (11) Panchenko, A. R.; Luthey-Schulten, Z.; Wolynes, P. G. Folds, Protein Structural Modules, and Exons. *Proc. Natl. Acad. Sci. U.S.A.* **1996**, *93*, 2008–2013.
- (12) Wales, D. J. Symmetry, Near-Symmetry and Energetics. *Chem. Phys. Lett.* **1998**, *285*, 330–336.
- (13) Ferreira, D. U.; Wolynes, P. G. The Capillarity Picture and the Kinetics of One-Dimensional Protein Folding. *Proc. Natl. Acad. Sci. U.S.A.* **2008**, *105*, 9853–9854.
- (14) Itoh, K.; Sasai, M. Multidimensional Theory of Protein Folding. *J. Chem. Phys.* **2009**, *130*, 145104.
- (15) Luo, H.; Nijveen, H. Understanding and Identifying Amino Acid Repeats. *Briefings Bioinf.* **2013**, DOI: 10.1093/bib/bbt003.
- (16) Kajava, A. V. Tandem Repeats in Proteins: from Sequence to Structure. *J. Struct. Biol.* **2012**, *179*, 279–288.
- (17) Shih, E. S.; Hwang, M. J. Alternative Alignments from Comparison of Protein Structures. *Proteins* **2004**, *56*, 519–527.
- (18) Abraham, A. L.; Rocha, E. P.; Pothier, J. Swelpe: a Detector of Internal Repeats in Sequences and Structures. *Bioinformatics* **2008**, *24*, 1536–1537.
- (19) Murray, K. B.; Taylor, W. R.; Thornton, J. M. Toward the Detection and Validation of Repeats in Protein Structure. *Proteins* **2004**, *57*, 365–380.
- (20) Taylor, W. R.; Heringa, J.; Baud, F.; Flores, T. P. A Fourier Analysis of Symmetry in Protein Structure. *Protein Eng.* **2002**, *15*, 79–89.
- (21) Walsh, I.; Sirocco, F. G.; Minervini, G.; Di Domenico, T.; Ferrari, C.; Tosatto, S. C. RAPHAEL: Recognition, Periodicity and Insertion Assignment of Solenoid Protein Structures. *Bioinformatics* **2012**, *28*, 3257–3264.
- (22) Marcotte, E. M.; Pellegrini, M.; Yeates, T. O.; Eisenberg, D. A Census of Protein Repeats. *J. Mol. Biol.* **1999**, *293*, 151–160.
- (23) Schaper, E.; Kajava, A. V.; Hauser, A.; Anisimova, M. Repeat or Not Repeat?—Statistical Validation of Tandem Repeat Prediction in Genomic Sequences. *Nucleic Acids Res.* **2012**, *40*, 10005–10017.
- (24) Sippl, M. J.; Wiederstein, M. Detection of Spatial Correlations in Protein Structures and Molecular Complexes. *Structure* **2012**, *20*, 718–728.
- (25) Sippl, M. J.; Wiederstein, M. A Note on Difficult Structure Alignment Problems. *Bioinformatics* **2008**, *24*, 426–427.
- (26) Sippl, M. J. On Distance and Similarity in Fold Space. *Bioinformatics* **2008**, *24*, 872–873.
- (27) Mosavi, L. K.; Minor, D. L.; Peng, Z. Y. Consensus-Derived Structural Determinants of the Ankyrin Repeat Motif. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, *99*, 16029–16034.
- (28) Wolynes, P. G. Folding Funnels and Energy Landscapes of Larger Proteins within the Capillarity Approximation. *Proc. Natl. Acad. Sci. U.S.A.* **1997**, *94*, 6170–6175.
- (29) Ferreira, D. U.; Walczak, A. M.; Komives, E. A.; Wolynes, P. G. The Energy Landscapes of Repeat-Containing Proteins: Topology, Cooperativity, and the Folding Funnels of One-Dimensional Architectures. *PLoS Comput. Biol.* **2008**, *4*, e1000070.
- (30) Ferreira, D. U.; Komives, E. A. Molecular Mechanisms of System Control of NF- $\kappa$ B Signaling by IkappaBalpha. *Biochemistry* **2010**, *49*, 1560–1567.
- (31) DeVries, I.; Ferreira, D. U.; Sanchez, I. E.; Komives, E. A. Folding Kinetics of the Cooperatively Folded Subdomain of the IkappaBalpha Ankyrin Repeat Domain. *J. Mol. Biol.* **2011**, *408*, 163–176.
- (32) Ferreira, D. U.; Cervantes, C. F.; Truhlar, S. M.; Cho, S. S.; Wolynes, P. G.; Komives, E. A. Stabilizing IkappaBalpha by “Consensus” Design. *J. Mol. Biol.* **2007**, *365*, 1201–1216.
- (33) Muraki, M.; Ishimura, M.; Harata, K. Interactions of Wheat-Germ Agglutinin with GlcNAc Beta 1,6Gal Sequence. *Biochim. Biophys. Acta* **2002**, *1569*, 10–20.
- (34) Haigis, M. C.; Haag, E. S.; Raines, R. T. Evolution of Ribonuclease Inhibitor by Exon Duplication. *Mol. Biol. Evol.* **2002**, *19*, 959–963.
- (35) Groves, M. R.; Hanlon, N.; Turowski, P.; Hemmings, B. A.; Barford, D. The Structure of the Protein Phosphatase 2A PR65/A Subunit Reveals the Conformation of Its 15 Tandemly Repeated HEAT Motifs. *Cell* **1999**, *96*, 99–110.
- (36) Nagano, N.; Orengo, C. A.; Thornton, J. M. One Fold with Many Functions: the Evolutionary Relationships between TIM Barrel Families Based on Their Sequences, Structures and Functions. *J. Mol. Biol.* **2002**, *321*, 741–765.
- (37) Soding, J.; Remmert, M.; Biegert, A. HHrep: De Novo Protein Repeat Detection and the Origin of TIM Barrels. *Nucleic Acids Res.* **2006**, *34*, W137–W142.
- (38) Fulop, V.; Jones, D. T. Beta Propellers: Structural Rigidity and Functional Diversity. *Curr. Opin. Struct. Biol.* **1999**, *9*, 715–721.
- (39) Neer, E. J.; Schmidt, C. J.; Nambudripad, R.; Smith, T. F. The Ancient Regulatory-Protein Family of WD-Repeat Proteins. *Nature* **1994**, *371*, 297–300.
- (40) Pauling, L.; Corey, R. B.; Branson, H. R. The Structure of Proteins; Two Hydrogen-Bonded Helical Configurations of the Polypeptide Chain. *Proc. Natl. Acad. Sci. U.S.A.* **1951**, *37*, 205–211.
- (41) Pauling, L.; Corey, R. B. The Pleated Sheet, a New Layer Configuration of Polypeptide Chains. *Proc. Natl. Acad. Sci. U.S.A.* **1951**, *37*, 251–256.
- (42) Moulton, J.; Fidelis, K.; Krysztofowicz, A.; Tramontano, A. Critical Assessment of Methods of Protein Structure Prediction (CASP)—Round IX. *Proteins* **2011**, *79* (Suppl. 10), 1–5.
- (43) Hegler, J. A.; Lätzer, J.; Shehu, A.; Clementi, C.; Wolynes, P. G. Restriction versus Guidance in Protein Structure Prediction. *Proc. Natl. Acad. Sci. U.S.A.* **2009**, *106*, 15302–15307.
- (44) Simons, K. T.; Kooperberg, C.; Huang, E.; Baker, D. Assembly of Protein Tertiary Structures from Fragments with Similar Local Sequences Using Simulated Annealing and Bayesian Scoring Functions. *J. Mol. Biol.* **1997**, *268*, 209–225.
- (45) Truscott, R. J. W. Macromolecular Deterioration as the Ultimate Constraint on Human Lifespan. *Ageing Res. Rev.* **2011**, *10*, 397–403.
- (46) Frauenfelder, H.; McMahon, B. H.; Fenimore, P. W. Myoglobin: the Hydrogen Atom of Biology and a Paradigm of Complexity. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 8615–8617.
- (47) Kendrew, J.; Bodo, G.; Dintzis, H. M.; Parrish, R. G.; Wyckoff, H.; Phillips, D. C. A Three-Dimensional Model of the Myoglobin Molecule Obtained by X-ray Analysis. *Nature* **1958**, *181*, 662–666.
- (48) Ormö, M.; Cubitt, A. B.; Kallio, K.; Gross, L. A.; Tsien, R. Y.; Remington, S. J. Crystal Structure of the Aequorea Victoria Green Fluorescent Protein. *Science* **1996**, *273*, 1392–1395.
- (49) Santos, J.; Gebhard, L. G.; Risso, V. A.; Ferreyra, R. G.; Rossi, J. P. F. C.; Ermácora, M. R. Folding of an Abridged Beta-Lactamase. *Biochemistry* **2004**, *43*, 1715–1723.
- (50) Goodsell, D. S.; Olson, A. J. Structural Symmetry and Protein Function. *Annu. Rev. Biophys. Biomol. Struct.* **2000**, *29*, 105–153.
- (51) Swapna, L. S.; Srikeerthana, K.; Srinivasan, N. Extent of Structural Asymmetry in Homodimeric Proteins: Prevalence and Relevance. *PLoS One* **2012**, *7*, e36688.
- (52) Hashimoto, K.; Panchenko, A. R. Mechanisms of Protein Oligomerization, the Critical Role of Insertions and Deletions in Maintaining Different Oligomeric States. *Proc. Natl. Acad. Sci. U.S.A.* **2010**, *107*, 20352–20357.
- (53) Sánchez, I. E.; Ferreira, D. U.; Dellarole, M.; de Prat-Gay, G. Experimental Snapshots of a Protein-DNA Binding Landscape. *Proc. Natl. Acad. Sci. U.S.A.* **2010**, *107*, 7751–7756.
- (54) Wolynes, P. G. Aperiodic Crystals: Biology, Chemistry and Physics in a Fugue with Stretto. *AIP Conf. Proc.* **1988**, *180*, 39–65.
- (55) Wales, D. J. Decoding the Energy Landscape: Extracting Structure, Dynamics and Thermodynamics. *Philos. Trans. R. Soc., A* **2012**, *370*, 2877–2899.

- (56) Denton, M. J.; Marshall, C. J.; Legge, M. The Protein Folds as Platonic Forms: New Support for the Pre-Darwinian Conception of Evolution by Natural Law. *J. Theor. Biol.* **2002**, *219*, 325–342.
- (57) Schrödinger, E. *What Is Life?*; Cambridge University Press: Cambridge, UK, 1944.
- (58) Schafer, N. P.; Hoffman, R. M.; Burger, A.; Craig, P. O.; Komives, E. A.; Wolynes, P. G. Discrete Kinetic Models from Funneled Energy Landscape Simulations. *PLoS One* **2012**, *7*, e50635.
- (59) Levy, Y.; Cho, S. S.; Shen, T.; Onuchic, J. N.; Wolynes, P. G. Symmetry and Frustration in Protein Energy Landscapes: a near Degeneracy Resolves the Rop Dimer-Folding Mystery. *Proc. Natl. Acad. Sci. U.S.A.* **2005**, *102*, 2373–2378.
- (60) Hegler, J. A.; Weinkam, P.; Wolynes, P. G. The Spectrum of Biomolecular States and Motions. *HFSP J.* **2008**, *2*, 307–313.
- (61) Monod, J.; Wyman, J.; Changeux, J. P. On the Nature of Allosteric Transitions: a Plausible Model. *J. Mol. Biol.* **1965**, *12*, 88–118.
- (62) Kuriyan, J.; Eisenberg, D. The Origin of Protein Interactions and Allostery in Colocalization. *Nature* **2007**, *450*, 983–990.
- (63) Wang, S.; Wolynes, P. G. On the Spontaneous Collective Motion of Active Matter. *Proc. Natl. Acad. Sci. U.S.A.* **2011**, *108*, 15184–15189.
- (64) Jablonka, E.; Raz, G. Transgenerational Epigenetic Inheritance: Prevalence, Mechanisms, and Implications for the Study of Heredity and Evolution. *Q. Rev. Biol.* **2009**, *84*, 131–176.
- (65) Pauling, L.; Itano, H. A. Sickle Cell Anemia a Molecular Disease. *Science* **1949**, *110*, 543–548.
- (66) Treusch, S.; Cyr, D. M.; Lindquist, S. Amyloid Deposits: Protection Against Toxic Protein Species? *Cell Cycle* **2009**, *8*, 1668–1674.
- (67) Frauenfelder, H.; Sligar, S. G.; Wolynes, P. G. The Energy Landscapes and Motions of Proteins. *Science* **1991**, *254*, 1598–1603.
- (68) Frauenfelder, H. Proteins: Paradigms of Complexity. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, *99* (Suppl. 1), 2479–2480.
- (69) Zhuravlev, P. I.; Papoian, G. A. Protein Functional Landscapes, Dynamics, Allostery: a Tortuous Path towards a Universal Theoretical Framework. *Q. Rev. Biophys.* **2010**, *43*, 295–332.