



Universiteit  
Leiden  
The Netherlands

## Topological principles of protein folding

Scalvini, B.; Sheikhhassani, V.; Mashaghi, A.

### Citation

Scalvini, B., Sheikhhassani, V., & Mashaghi, A. (2021). Topological principles of protein folding. *Physical Chemistry Chemical Physics*, 23(37), 21316-21328.  
doi:10.1039/d1cp03390e

Version: Publisher's Version

License: [Licensed under Article 25fa Copyright Act/Law \(Amendment Taverne\)](#)

Downloaded from: <https://hdl.handle.net/1887/3277889>

**Note:** To cite this publication please use the final published version (if applicable).



Cite this: *Phys. Chem. Chem. Phys.*,  
2021, **23**, 21316

## Topological principles of protein folding†

Barbara Scalvini, Vahid Sheikhhassani and Alireza Mashaghi \*

What is the topology of a protein and what governs protein folding to a specific topology? This is a fundamental question in biology. The protein folding reaction is a critically important cellular process, which is failing in many prevalent diseases. Understanding protein folding is also key to the design of new proteins for applications. However, our ability to predict the folding of a protein chain is quite limited and much is still unknown about the topological principles of folding. Current predictors of folding kinetics, including the contact order and size, present a limited predictive power, suggesting that these models are fundamentally incomplete. Here, we use a newly developed mathematical framework to define and extract the topology of a native protein conformation beyond knot theory, and investigate the relationship between native topology and folding kinetics in experimentally characterized proteins. We show that not only the folding rate, but also the mechanistic insight into folding mechanisms can be inferred from topological parameters. We identify basic topological features that speed up or slow down the folding process. The approach enabled the decomposition of protein 3D conformation into topologically independent elementary folding units, called circuits. The number of circuits correlates significantly with the folding rate, offering not only an efficient kinetic predictor, but also a tool for a deeper understanding of theoretical folding models. This study contributes to recent work that reveals the critical relevance of topology to protein folding with a new, contact-based, mathematically rigorous perspective. We show that topology can predict folding kinetics when geometry-based predictors like contact order and size fail.

Received 24th July 2021,  
Accepted 1st September 2021

DOI: 10.1039/d1cp03390e

rsc.li/pccp

Over the last 20 years, it has been hypothesized that protein folding rates and mechanisms can be inferred from the native state topology.<sup>1</sup> The importance of local intra-chain contacts for small one-domain proteins emerged with the definition of Contact Order (CO), a “topological” parameter still widely used to date to predict protein folding rates.<sup>2</sup> This parameter was then coupled with size (length of the protein) with the introduction of absolute CO,<sup>3</sup> to allow for a better description of the folding kinetics of larger proteins. For such proteins, the folding pathway may be characterized by kinetic traps and escape from low free energy conformations.<sup>4,5</sup> In more recent years, other models have been suggested for folding rate prediction, based on total contact distance,<sup>6</sup> a small selection of contact information,<sup>7</sup> cumulative torsion angle<sup>8</sup> and other structural information.<sup>9–13</sup> Moreover, an evolution of the concept of contact order called *partial contact order* was envisioned in order to follow the progression of such topological descriptors from the unfolded to the folded state.<sup>14</sup> The partial contact order

(pCO) takes into account the likelihood that a certain contact is formed, and the associated reduction of loop entropy.<sup>14</sup> However, contact distance, contact order and protein length are not inherently topological properties, if topology is to be intended in the mathematical sense of the word. Topology is a mathematical concept characterizing the properties of objects which remain unaltered through continuous, invertible transformations such as stretching, shrinking and bending. A first step to introduce topology-based predictors for the quantification of entanglement was taken by Marco Baiesi *et al.*<sup>15–17</sup> Drawing from knot theory, the concept of Gaussian entanglement was first applied to the intertwined backbones of domain-swapped protein dimers,<sup>17</sup> and then to non-overlapping looping sub-chains of the same protein, where it proved to complement absolute CO in folding rate prediction on a set of 48 proteins.<sup>15</sup> However informative, these topologically inspired descriptors often concern a fairly limited portion of the available protein datasets, with about 15% of dimers displaying significant intertwining,<sup>17</sup> and 32% of proteins from the CATH database showing non-trivial Gaussian entanglement.<sup>16</sup> Topological concepts such as writhe and torsion were also applied to the protein backbone, yielding good results for folding rate prediction and revealing the role of handedness of proteins at both local and global organization levels.<sup>18</sup>

*Medical Systems Biophysics and Bioengineering, Leiden Academic Centre for Drug Research, Faculty of Science, Leiden University, Einsteinweg 55, 2333CC Leiden, The Netherlands. E-mail: a.mashaghi.tabari@lacdr.leidenuniv.nl*

† Electronic supplementary information (ESI) available. See DOI: 10.1039/d1cp03390e

Previous topological efforts to quantify the relationship between the native state three-dimensional arrangement and folding kinetics relied on the concept of entanglement as defined by knot theory, and focused on the entanglement of the backbone. However, knots are rare in proteins, and knotted proteins generally yield very low folding rates.<sup>19</sup> Other topologically inspired descriptors drawn from knot theory such as Gaussian entanglement also rely on the concept of backbone entanglement. The effect of entanglement of proteins with no knots and no slipknots on folding rates has been studied by Panagiotou and Plaxco<sup>18</sup> and Baiesi, Orlandini *et al.*<sup>15–17</sup> A mathematically rigorous topological concept, termed circuit topology (CT), has recently been proposed to describe the topology of unknots.<sup>20,21</sup> Circuit topology, in its first order definition, ignores possible backbone entanglement, and focuses only on the intra-chain contacts present in the native protein structure. Contacts are considered to be fixed. This allows circuit topology to provide a topological description of unknotted, yet folded linear chains,<sup>20–23</sup> a type of description which is complementary to that provided by Gaussian entanglement,<sup>15–17</sup> writhe and torsion.<sup>18</sup> Moreover, contact-based topological descriptors represent a very natural framework for proteins, since contacts often have not only geometrical but also biological relevance. The circuit topology framework allows us to readily combine our descriptors with information such as the energy of a contact, for example. The vast majority of proteins present intra-chain contacts, making our analysis applicable virtually to all proteins. Once contacts in a structure have been identified, they are classified based on their pairwise topological arrangement (Fig. 1A). According to CT, contacts can be in either one of three possible relationships with each other: series (S), parallel (P) and cross (X) (Fig. 1A). Series and parallel relationships also include a subset of relationships called *concerted relationships*, in which one of the two contact sites is shared between the two contacts. We call these concerted parallel (CP) and concerted series (CS) relationships. Here, we present a first order analysis; therefore, CP and CS will be included in the main sets and counted respectively as parallel and series. We note that CT was already suggested to have an impact on the folding dynamics of model polymers,<sup>22,24</sup> although its relevance to protein folding has not been evaluated.

Here, we show how the three fundamental topological relationships S, P, and X display differential patterns of correlation with folding rate, providing insight into which types of topological arrangements facilitate folding and which hinder it. We define the zipper effect as the mechanism with which a predominance of series arrangement slows folding, while parallel and cross arrangements (the so-called *entangled relationships*) yield higher folding rates. It is important to note that here the word ‘entangled’ is used in a broad sense, since we are dealing with unknots. Parallel and cross are designated as entangled because the two loops forming the relationship are not independent of each other. We show that both two-state and multi-state class proteins display statistical evidence of the zipper effect, if we only consider the topology of short range, attractive energy contacts. Lastly, we will show how proteins can be decomposed into *topological circuits*,<sup>25</sup> that is, topologically independent units. The number of these

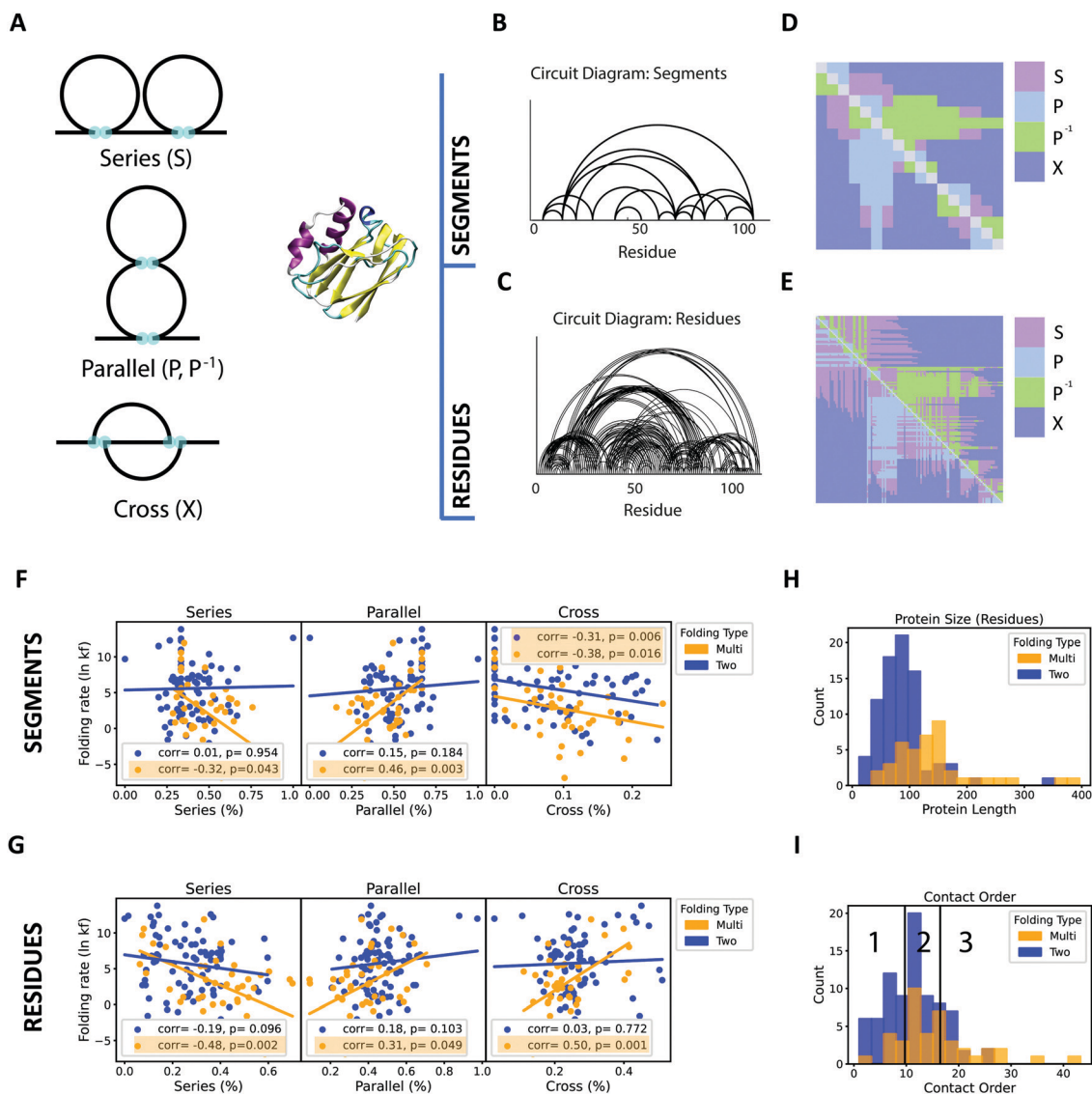
circuits normalized by size correlates positively with the logarithm of folding rate  $\ln(k_f)$ , suggesting that the localization of contacts inside topological circuits might play a role in facilitating folding efficiency.

## Results

### Topological parameters as kinetic predictors

Circuit topology utilizes contacts as basic elements for topological classification. However, a suitable definition of contacts is widely dependent on the purpose of the study. For folding rate prediction, contacts between residues have mostly been used for quantifying parameters such as CO.<sup>2</sup> Here, we will also consider contacts between residues. However, this is not the only choice; the flexibility of the CT framework allows us to consider other types of protein building blocks which can form contacts; one can identify *segments* of proteins which correspond to secondary structure elements, and perform CT analysis on the contacts created by these coarse-grained structures. In Fig. 1 we can see the CT diagram of segment–segment (Fig. 1B) and residue–residue (Fig. 1C) contacts, and their respective CT matrices, from which the frequencies of CT topological relationships can readily be extracted (Fig. 1D and E). Strikingly, these CT frequencies correlate with the logarithm of folding rate. The two choices of contact definition provide very different structural resolution, and we expectedly observe different degrees of correlation with folding rate. Contacts were retrieved from PDB structures, by defining a spatial cut-off for atom–atom distance (5.0 Å), and a threshold for the minimum number of atoms to be found in spatial proximity below the cut-off in order to consider the two residues/segments in contact (5 atoms for residues and 10 atoms for segments). Our main conclusions are robust with respect to the choice of parameters. For other cut-off choices, see the ESI.†

Next, we investigated whether the observed correlations depend on the complexity of the folding pathway. Many proteins fold and unfold with one main fast event, by a simple two-state transition. These “two-state folders”<sup>11,26,27</sup> have gathered much of the attention of scientific inquiry in the past, and their folding rates correlate with relative CO.<sup>2</sup> On the other hand, proteins with more intricate multi-state transitions – “multi-state folders” – have shown a strong dependency of their folding kinetics on the protein length (and not CO).<sup>28</sup> Notably, CT parameters also provide differential patterns of correlations for two- and multi-state folders, at first sight (Fig. 1F and G). For both segment and residue analyses we find statistically significant negative correlation between  $\ln(k_f)$  and series in multi-state folders (respectively  $r = -0.32$ ,  $p = 0.043$  and  $r = -0.48$ ,  $p = 0.002$ ). This result is understandable as series relationships favor delocalization along the chain, which seems to slow down the folding of multi-state folders, but leaves two-state folders unaffected. On the other hand, two-state proteins display moderate negative correlation with cross relationships ( $r = -0.31$  and  $p = 0.006$ ), in their segment representation. These differences might be due to the different average size of the two-state and multi-state proteins (Fig. 1H). Two-state proteins are



**Fig. 1** Segment and residue-based CT parameters correlate with folding rate. (A) Three pairwise arrangements of CT: series, parallel and cross. The inner contact is in parallel relation (P) with the outer contact, while the outer contact is in inverse parallel relation ( $P^{-1}$ ) with the inner contact. (B) Circuit diagram for segment-based contacts. (C) Circuit diagram for residue-based contacts. (D) CT matrix for segment-based contacts. (E) CT matrix for residue-based contacts. Segment and residue-based contacts offer very different resolution into protein topological arrangement, for the same protein (pseudoazurin, PDB code: 1ADW). F Scatterplot of topological fractions (series, parallel and cross) versus folding rate ( $\ln k_f$ ), for segment-based contacts. G Scatterplot of topological fractions (series, parallel and cross) versus folding rate ( $\ln k_f$ ), for residue-based contacts. H Size distribution (number of residues) for two-state and multi-state folders. I Contact Order distribution of the dataset. The dataset was divided into three sub-datasets (Lower, Average and Upper CO) by setting an upper (16.47) and lower (9.72) limit.

generally smaller; therefore, highly entangled topologies such as those favored by cross arrangement might be less likely to appear at the secondary structure level, for geometric and energetic constraints. The likelihood of finding such structures might increase for longer folding times. Therefore, it is not surprising to find a negative correlation between cross and folding rates in this instance. The folding rate in multi-state proteins displays a higher impact of topology, showing evidence of statistically relevant zipper effect at both residue and secondary structure levels, having negative correlation with series and positive correlations with at least one of the two entangled relationships – parallel for segments ( $r = 0.46$  and  $p = 0.003$ ),

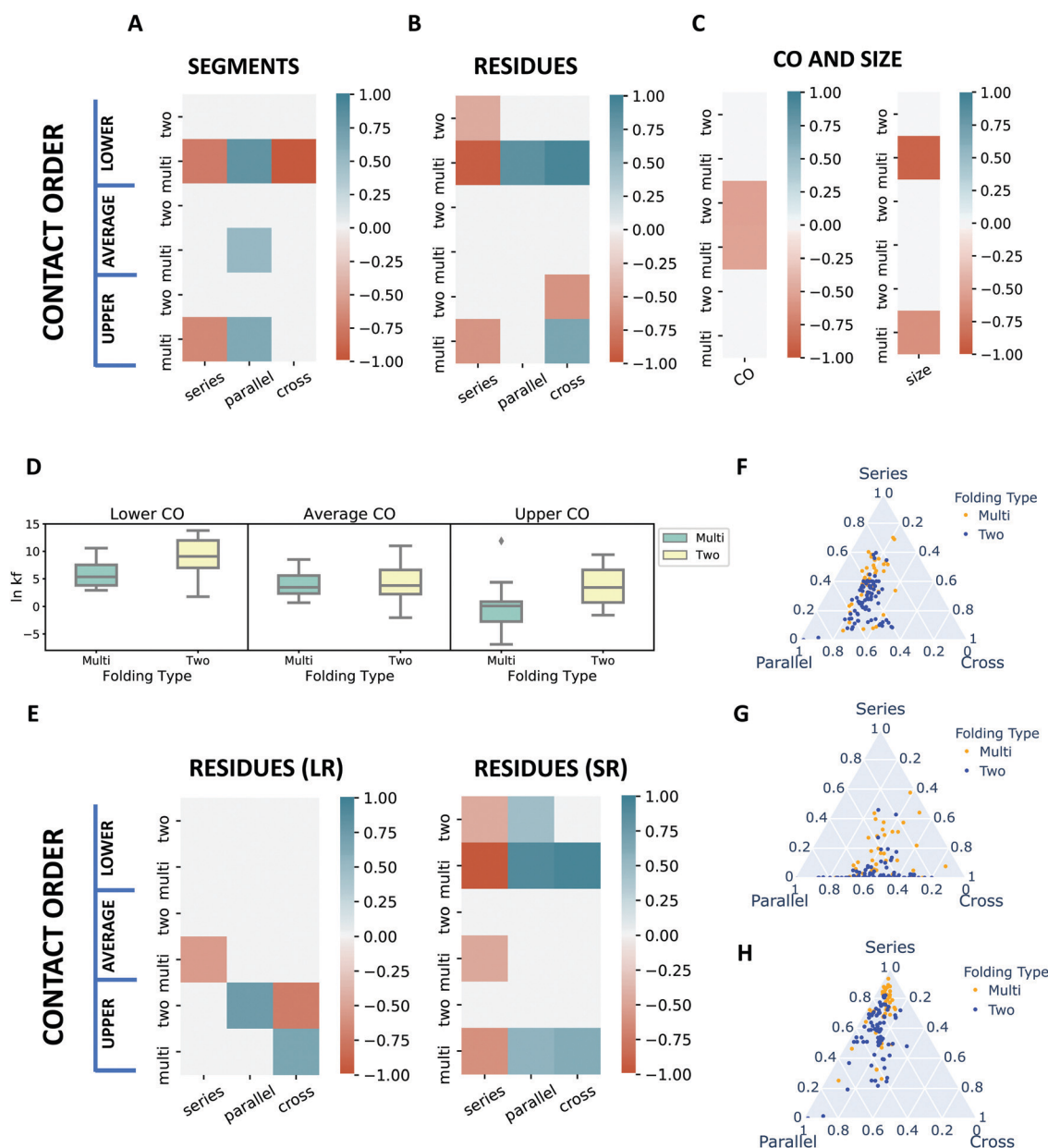
and both parallel and cross for residues ( $r = 0.31$  and  $p = 0.049$  for parallel,  $r = 0.50$  and  $p = 0.001$  for cross). Segment analysis yields correlation values for parallel (in multi-state folders) which are comparable to those obtained by Panagiotou and Plaxco for the writhe of the protein Primitive Path and the logarithm of the folding rate.<sup>18</sup> Chain writhing is a mechanism (possibly the main one in proteins) which can indeed create parallel contact topologies; thus, in this case, contact topology might be seen as a proxy for backbone topology. However, no correlation is found for parallel topology in two-state proteins, indicating possibly that the protein is too short to produce substantial writhe. CT parameters are normalized by the number of contacts in the chain, making it

possible to compare proteins with very different geometrical properties. However, due to the assembly principles of proteins and geometrical and steric constraints, a non-trivial relationship between size and CT parameters exists (Fig. S1, ESI<sup>†</sup>).

### Disentangling the contributions of geometry and topology

We demonstrate that topology-based predictors complement CO and size, which are geometry-based predictors. In order to

do so, we divided the dataset into three sub-datasets based on their CO (Fig. 11): Upper, Average and Lower CO. CO values were retrieved from the ARCPro dataset<sup>29</sup> (cut-off value = 6 Å). Fig. 2A and B show the correlation coefficients for these three subsets, for two-state and multi-state proteins. Exact values can be seen in Tables S1 and S2 (ESI<sup>†</sup>). We also compare the CT correlation maps with those obtained by using CO and size on the same datasets (Fig. 2C and Table S3, ESI<sup>†</sup>). While CO is



**Fig. 2** Classification based on contact order and length filtering highlights differential patterns of correlation. (A) Folding rate correlation map for segment-based CT, with CO classification. (B) Folding rate correlation map for residue-based CT, with CO classification. (C) Folding rate correlation map for contact order and size, with CO classification. CT seems to be more informative than contact order for proteins with Lower and Upper Folding rates. (D) Boxplot of folding rates for different CO subsets. Slow folders populate the Upper CO sub-dataset, and display correlation between folding rate and long-range residue-based contacts. (E) Folding rate correlation map for residue-based CT, with CO classification. The two maps show only long-range contacts (on the left) and only short-range contacts (on the right). The threshold for range classification was set to 24 residues. (F) Triangular plot of the topological composition throughout the dataset, for residue-based CT. (G) Triangular plot of the topological composition throughout the dataset, for long-range residue-based CT. (H) Triangular plot of the topological composition throughout the dataset, for short-range residue-based CT.



moderately accurate in predicting  $\ln(k_f)$  for the Average CO dataset ( $r = -0.53$  and  $p = 4.5 \times 10^{-4}$  for two-state;  $r = -0.51$  and  $p = 0.03$  for multi-state), CT seems to obtain the best results for the two tails of the CO distribution, obtaining correlations as high as  $r = -0.93$  ( $p = 0.002$ ) for series and  $r = 0.94$  ( $p = 0.001$ ) for cross for multi-state proteins in the Lower CO range (Fig. 2B). These results imply that CO and CT give in fact complementary information about folding kinetics. Also, CT is able to provide resolution for those proteins that have a similar CO but present significant discrepancies in the folding rate. The size parameter also provides strong correlations for the Upper and Lower CO datasets (Fig. 2C and Table S4, ESI<sup>†</sup>), although only for multi-state proteins ( $r = -0.89$  and  $p = 0.01$  for Lower CO and  $r = -0.61$  and  $p = 0.01$ ), as expected. By combining the kinetic information on multi-state proteins provided by residue-based CT and size parameters, we see that not only the number of residues is impactful but also their topological arrangement, with contact delocalization favored by series relationships being as efficient as protein length in promoting a slow folding process. Interestingly, CT on the segment level displays an opposite trend for cross relationships for multi-state, Lower CO proteins, indicating that such a level of entanglement at the secondary structure level might actually be hindering folding for those members of the multi-state protein class which are smaller (Fig. 1H) and have higher folding rates (Fig. 2D). The fact that smaller multi-state proteins show similar correlations for the cross-fraction to two-state folders might suggest a rather continuous transition with respect to size between multi-state and two-state classes, rather than two binary distinct folding styles.<sup>5,30</sup>

### Arrangement of short-range attractive contacts as a topological driver of folding

Here, we investigate how topology, interaction energy, and interaction range work together to regulate the folding kinetics. The CO reflects the relative importance of local and non-local interactions in the molecule.<sup>2</sup> The conceptual background behind CO is that contacts between residues that are closer along the chain are less entropically costly, and therefore tend to happen early in the folding process. Therefore, simple protein structures which are rich in local contacts tend to fold faster.<sup>1,31</sup> Broglia and Tiana<sup>32</sup> highlighted the role of local contacts by identifying a specific hierarchy, which involves the formation of early local elementary structures (LESSs), followed by the assembly of the LES into a post critical folding nucleus at a later stage. Moreover, evidence exists that natural selection favors folds with a low contact order,<sup>33</sup> and therefore structures rich in local contacts. Indeed, off-lattice models of protein folding showed that the suppression of local interactions prevents the structure from reaching the native conformation.<sup>34</sup> However, the respective role of local *versus* non-local interactions is still a highly debated subject in the literature. *In silico* studies of three model 36-mers on a cubic lattice suggested that non-local interactions are the primary determinant of protein folding.<sup>35</sup>

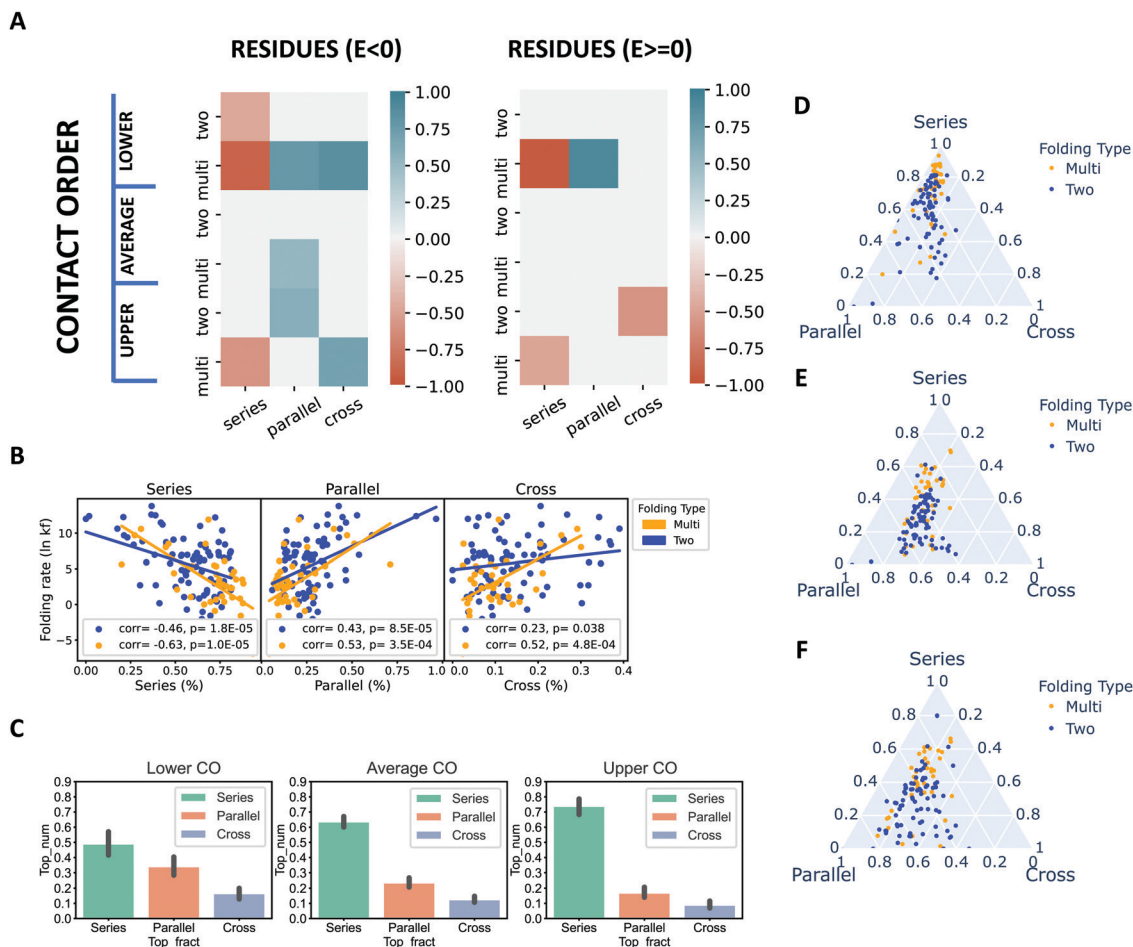
We can ask ourselves if not only the relative number of local *versus* non-local contacts, but also their topological

arrangement has an impact on folding kinetics. To address this question, we applied a 24 residue threshold in order to discriminate between short-range and long-range contacts prior to CT analysis (Fig. 2E and Tables S5, S6, ESI<sup>†</sup>). It is apparent to see that the topology of short-range contacts displays correlations which are higher in magnitude and also more widespread over the whole CO range, as opposed to long range contacts. Multi-state proteins in the Lower CO range still display the highest correlations between the topological content and  $\ln(k_f)$ :  $r = -0.97$ ,  $p = 1.9 \times 10^{-4}$  for series,  $r = 0.89$ ,  $p = 0.007$  for parallel, and  $r = 0.94$ ,  $p = 0.002$  for cross. The zipper effect appears to be confirmed in the results from the short-range correlation panel (Fig. 2E): once local contacts are uncoupled from non-local contacts in CT analysis, negative correlations with folding rates are only seen for series relationships, and positive correlations are observed with the *entangled relationships*, cross and parallel. Short-range contacts appear as the main topological folding drivers. This is compatible with the findings of Adesh Kumar and co-workers,<sup>36</sup> who theorized that local contacts might be fundamental for the differentiation between the native-like conformations during folding, by Monte Carlo simulation of three protein structures. However, correlations with long-range contacts also appear for the 'slow folding' Upper CO proteins (Fig. 2D and E). Since non-local contacts along the chain are generally formed at a later stage during folding,<sup>32</sup> they can only affect the folding process after longer characteristic times. This finding suggests that, for very fast folders, the impact of the topology of long-range contacts might be negligible.

Moreover, we find that short and long-range contacts are also qualitatively different with respect to the topological content. Fig. 2F portrays in a triangular plot the percentages of series, cross and parallel for all residues. We can compare it to the topological content in long-range (Fig. 2G) and short-range (Fig. 2H) contacts; we see that non-local contacts are actually much richer in cross relationships with respect to local contacts. This finding indicates that on the short-range contacts, high levels of entanglement such as those promoted by cross relationships are unfavorable.

Contacts can also be discriminated by assigning energy-like quantities based on the statistical potential suggested by Paul Thomas and Ken Dill.<sup>37,38</sup> This is a first order attempt to add bio-chemical information to contact topology: contacts can be filtered based on the sign of the potential matrix element associated with residue-residue interaction, resulting in 'repulsive energy contacts' ( $E > 0$ ) and 'attractive energy contacts' ( $E < 0$ ). The correlation map for energy filtering (Fig. 3A and Tables S7, S8, ESI<sup>†</sup>) clearly highlights how the topology of attractive energy contacts plays the biggest role in folding kinetics. However, repulsive contacts can still correlate with slower folding processes, as in the case of Lower CO multi-state ( $r = -0.95$ ,  $p = 0.001$  for series), for Upper CO two-state proteins (cross,  $r = -0.58$ ,  $p = 0.05$ ) and Upper CO multi-state proteins (series,  $r = -0.50$ ,  $p = 0.048$ ).

Considering the results for length and energy-based contact filtering, it becomes clear that not all contact topologies are equally impactful when it comes to folding. Local, negative



**Fig. 3** Classification based on contact order and energy filtering highlights the kinetic role of the topology of short-range attractive contacts. (A) Folding rate correlation maps for residue-based CT, with CO classification. The two maps show only negative (attractive) energy contacts (on the left) and only positive (repulsive) energy contacts (on the right). (B) Scatterplot for residue-based CT fractions and folding rate: only short-range attractive energy contacts were included. With this type of filtering, both folding types display the zipping effect, and all correlations are significant ( $p$  value  $\leq 0.05$ ). (C) Bar plot of topological fractions with respect to contact order classification, for attractive energy short-range residue-based CT. With increasing CO, we observe an increase in series fraction and a decrease in entangled fraction (parallel, cross). (D) Triangular plot of the topological composition throughout the dataset, for attractive energy short-range residue-based CT. (E) Triangular plot of the topological composition throughout the dataset, for attractive energy residue-based CT. (F) Triangular plot of the topological composition throughout the dataset, for repulsive energy residue-based CT.

(attractive) energy contacts seem to be the topological drivers of the folding process. It is therefore natural to reconsider correlations for the whole dataset while considering short-range, attractive energy contacts exclusively (Fig. 3B). Interestingly, this type of filtering yields statistically significant correlations for both two-state and multi-state proteins, for all three topological relationships. Even more notably, now both two-state and multi-state proteins show evidence of zipper effect, making the distinction between the two classes more quantitative than qualitative; correlations seem to be still more pronounced in the case of multi-state folders, but correlation trends are the same for the two classes. Fig. 3C shows another evidence of zipper effect; with decreasing contact order (higher folding rate), there is a gradual increase in entangled relationships. However, the triangular plot of the energy/length filtered dataset (Fig. 3D) is a closer match to the short-range triangular plot (Fig. 2H) rather than to the attractive energy plot (Fig. 3E), indicating that the best predictor for the topological

content is the distance between contacts, and not the energy of the contacts. Moreover, the topological content for repulsive energy contacts (Fig. 3F) does not look significantly different from the one for attractive energy contacts (Fig. 3E).

### Linear combination of CO and CT parameters as an improved folding rate predictor

The analysis outlined so far suggests complementarity between folding rate descriptors such as CT parameters and Contact Order. We see, for example, how the pre-filtering of data based on Contact Order is useful to uncover differential patterns of correlation for CT parameters (Fig. 2A–C, E and 3A). CO pre-filtering highlights also how proteins belonging to different CO ranges might be best described by CO, CT parameters or size, when it comes to folding rate prediction. It is then natural to ask whether these folding rate descriptors could be combined to produce more accurate folding rate predictions. In order to

test this hypothesis, we envisioned a multilinear regression analysis of the dataset, using CT parameters, CT parameters combined with CO, and CT parameters combined with size as independent variables. Folding rate predictions yielded by using only CO and size are also reported for comparison. For the analysis, we use CT fractions derived from attractive energy short-range contacts, since this unifies two- and multi-state folders for what concerns their correlation patterns with respect to CT (Fig. 3B). All CT parameter values reported in this paper were previously normalized by the total sum of S, P and X relationships in the protein. This normalization implies that, once we provide two CT fractions, the third is automatically determined, as the sum of all three fractions needs to yield 1. Thus, we can compare proteins with very different number of contacts. However, one of the three CT parameters is actually redundant, when it comes to multilinear regression analysis (MLR). Therefore, we decided to discard one and only use two CT parameters for folding rate prediction. Since the independent variables used for MLR should not be too highly correlated, we chose the two CT parameters which presented the lowest correlation coefficient when confronted with each other, that is, parallel and cross contacts ( $r = 0.23$ ,  $p = 0.011$ ). The CT-based folding rate predictor is therefore defined as:

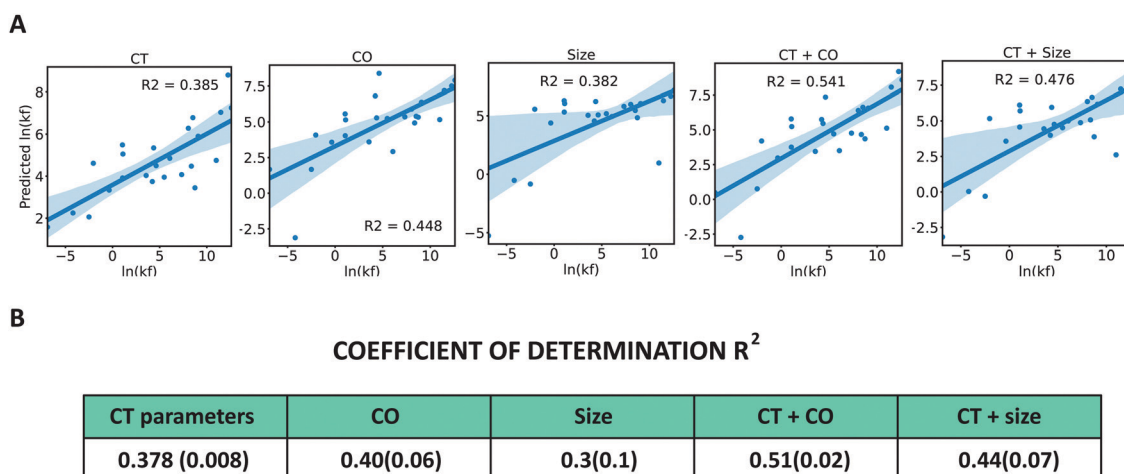
$$K_{CT} = c_P P + c_X X$$

where  $K_{CT}$  is the predicted logarithm of the folding rate,  $P$  and  $X$  are the parallel and cross fractions and  $c_P$  and  $c_X$  are coefficients which are calculated by the MLR model over the training set. Following the same reasoning, CT parameters can be combined with the CO and size to obtain new predictors:

$$K_{CT+CO} = c_P P + c_X X + c_{CO} CO$$

$$K_{CT+L} = c_P P + c_X X + c_L L$$

where  $L$  is the size of the protein (number of residues), and  $c_L$ ,  $c_{CO}$  are the coefficients calculated by the MLR model. In order to perform this analysis we relied on a freely available Python tool for machine learning and predictive data analysis, scikit-learn 0.24.2.<sup>39</sup> Thanks to the scikit-learn cross-validator module, we divided the dataset into 5 consecutive folds (sub-sets). Iteratively, 4 of these 5 datasets were used as training sets for the model, and the remaining one as the test set for folding rate prediction. Folding rate predictions on one of these test sets can be seen in Fig. 4A, for all predictors. Predictions for all test sets can be found in Fig. S2 (ESI<sup>†</sup>). A useful parameter to quantify the goodness of our prediction (how well the MLR model is representative of our dataset) is the coefficient of determination  $R^2$ .<sup>40</sup> The table in Fig. 4B presents the average  $R^2$  over the predictions from the 5 test sets; it is clear to see that both CO and size have a higher predictive power when combined with CT than when they are used on their own, with  $K_{CT+CO}$  representing the best folding rate predictor. Folding rate predictions from the first and last test set were excluded from the comparison, as the residual (predicted folding rate – experimental folding rate) distribution from CO prediction did not satisfy the normality requirement (Table S9, ESI<sup>†</sup>). However,  $R^2$  values and adjusted  $R^2$  values from all test sets can be seen in Tables S10 and S11 (ESI<sup>†</sup>) respectively. The adjusted determination coefficient  $R_{adj}^2$  is a modified version of  $R^2$  which takes into account the number of independent variables in the model. It discriminates whether the added variables provide an improvement to the prediction which is higher than what would be expected by the addition of random parameters. The  $R_{adj}^2$  coefficient confirms our general conclusions which identify  $K_{CT+CO}$  as the best predictor (Table S11, ESI<sup>†</sup>). This result proves the complementarity of CT parameters and CO for folding rate prediction, which we already hypothesized from the analysis presented in Fig. 2B and C.



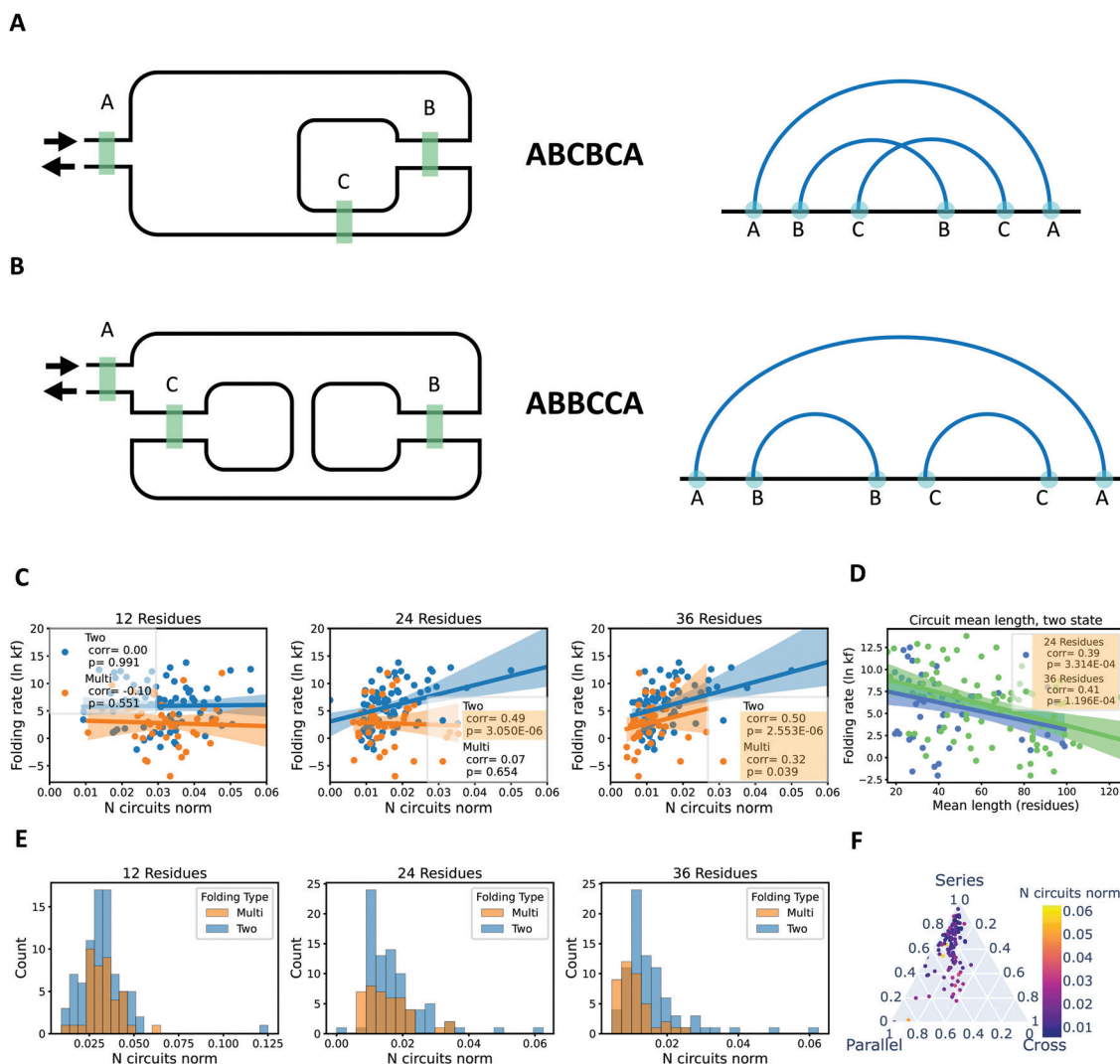
**Fig. 4** The linear combination of CT parameters and CO allows for folding rate prediction with increased statistical significance. (A) Scatterplots of predicted folding rate (obtained with multilinear regression over CT fractions, CO, protein length and a combination of these parameters) and experimental folding rate ( $\ln k_f$ ), calculated over one of the 5 training/test set combinations. (B) Average  $R^2$  score for CT parameters (parallel, cross), CO, size (protein length expressed in number of residues) and linear combination of CT parameters and CO, CT parameters and size. Numbers between parentheses indicate the standard deviation. The average was performed over 3 different choices of training/test subsets. Predictions obtained over test sets 1 and 5 were discarded by residual analysis, as their residual distribution did not satisfy the normality requirement.



### Circuits as elementary folding units

The circuit topology of a chain enables bottom-up analysis of a fold architecture. We investigate whether higher order topological features are seen in proteins and whether they contribute to folding kinetics. It was previously suggested that protein folding might proceed in a step-wise manner from separately cooperative elementary units of about 20 residues, called foldons.<sup>41</sup> Analogously, we can look for the topological equivalent of folding sub-units by exploiting a string notation of contacts, such as that utilized by generalized circuit topology.<sup>25</sup> The string notation allows for the identification of *circuits* in the chain, formally defined as a segment of a string that consists only of pairs of letters. Circuits represent well-separated topological structures

within a complex topology. Let us have a look at Fig. 5A and B to clarify the notation. Letters are assigned to contacts in the order in which they appear along the chain. Each contact site will then be represented in the string by that letter; consequently, each letter will appear in the sequence twice, as each contact is formed by two contact sites (residues, in this case). Thus, if we take the diagram shown in Fig. 5A, and we follow the chain from beginning to end, we first encounter contact A, then contact B, then C, B, C and finally A. Therefore, the string notation is **ABCBCA**. The choice of the letter (or general symbol to identify the contact) is arbitrary; the notation is valid as long as the symbol used is unique to that contact in the string. Each segment of the chain which consists of full pairs of letters represents a circuit.



**Fig. 5** The number of topological circuits normalized by the size of the protein correlates positively with folding rate. (A) Example of a circuit, with string and diagram representation. This circuit can be further decomposed, as BCBC is itself a circuit. (B) Example of a circuit, with string and diagram representation. Also in this case, the circuit can be further decomposed. If we remove contact B, we would obtain circuit ACCA, leaving the topology of contact C and A unaffected. The same goes for contact C and circuit ABBA. Contacts C and B together also form a circuit, BBCC. (C) Scatterplot for number of circuits normalized by size and folding rate. Legends display Spearman correlation coefficients. (D) Scatterplot for the circuit mean length and folding rate, for 24 and 36 residue thresholds for long-range exclusion. No correlation was detected for the 12 residue threshold. (E) Histograms of the number of circuits, normalized by protein length, for all long-range exclusion thresholds. (F) Triangular plot of CT fractions for residue-based CT. The color code indicates the number of circuits normalized by size, calculated with 36 residues long-range exclusion threshold.

Therefore, the chain identified by **ABCBCA** is itself a circuit. **BCBC** is also a circuit, while **ABCB** is not. In Fig. 5B, **ABBCCA** is a circuit, as also **BB** and **CC**. These are topologically independent units; the circuit **CC** could be removed, and the topology of **ABBA** would be unaffected; **ABBA** would still be a circuit. However, which circuit should we consider, when decomposing a longer chain? **ABBCCA** or the shorter **BB** and **CC**? This depends on the threshold we impose for the exclusion of long-range contacts. Three thresholds were tested on our dataset, 12, 24 and 36 residues. Imposing a threshold implies, for example, that the contacts which are formed by residues that are more than 12 residues apart along the chain are erased, in order to reveal the self-contained topological sub-structures of this length range.

The retrieved number of circuits is related to the size of the protein, but not in a trivial way.<sup>25</sup> The number of circuits in a protein can be considered its topological size. Topological and geometrical size are clearly two closely related concepts. Here we show, however, that the information provided by these two parameters is not redundant. Correlations between protein size and number of circuits decrease as we increase the threshold for long-range contact exclusion in circuit calculations. Correlations go from being as high as  $r = 0.91$  ( $p = 1.85 \times 10^{-47}$ ) for  $t_{lr} = 12$ , to  $r = 0.64$  ( $p = 2.98 \times 10^{-15}$ ) for  $t_{lr} = 36$ . This consideration sets the lower boundary for our analysis, as for thresholds which are as low as 12 residues, the detected topological length size coincides with the geometrical size. This becomes apparent when we normalize the number of circuits by protein length, and use the normalized number of circuits as folding rate predictor (Fig. 5C). While we observe no correlation between the normalized number of circuits and  $\ln(k_f)$  for  $t_{lr} = 12$ , the correlation increases as we go towards a higher threshold. The correlation is particularly pronounced for two-state folders, for which we observe significant correlations for both  $t_{lr} = 24$  ( $r = 0.49$ ,  $p = 3.05 \times 10^{-6}$ ) and  $t_{lr} = 36$  ( $r = 0.50$ ,  $p = 2.55 \times 10^{-6}$ ). Multi-state folders, on the other hand, only display correlation for  $t_{lr} = 36$ , which is also weaker in magnitude as opposed to that of two-state proteins ( $r = 0.32$ ,  $p = 0.039$ ). This is interesting especially if we consider how traditionally size as a folding rate predictor was particularly successful when applied to multi-state folders. Observing a significant, albeit weak correlation for multi-state folders for the normalized number of circuits indicates that topological and geometrical sizes are not always equivalent concepts. This consideration is particularly true when we consider two-state folders, where size generally provides only modest correlations. Indeed, for this dataset the correlation between the protein length and  $\ln(k_f)$  for two-state folders is  $r = -0.28$  and  $p = 0.010$  (Fig. S3, ESI†); this finding suggests that the topological length might be a better descriptor for folding kinetics than the geometrical size, for two-state proteins. In general, the correlations in Fig. 5C suggest that, for proteins of comparable length, a subdivision in a higher number of topologically independent units might facilitate folding. Moreover, the size of the circuits also seems to matter for two-state folders, with proteins made up by smaller circuits folding faster (Fig. 5D).

The distribution of the normalized number of circuits for two and multi-state folders also contains crucial information

concerning the topological makeup of the two folder types (Fig. 5E). While for  $t_{lr} = 12$  and  $t_{lr} = 24$  we do not observe any particular difference between the two distributions, for  $t_{lr} = 36$ , we actually observe a shift between the two, with two-state folders having a longer distribution tail towards high values of normalized number of circuits. For  $t_{lr} = 36$  residues, the multi and two-state distributions for normalized number of circuits are statistically different, as quantified by the Mann-Whitney  $U$  test ( $p = 5.05 \times 10^{-4}$ ). There is still significant overlap between the two distributions for low values of normalized number of circuits, indicating, again, that the difference between two and multi-state folders is not binary. Nevertheless, the results suggest that topology might be informative not only of the speed but also of the quality of the folding process.

Concerning the topological content of the circuits, we do not observe a clear trend between the normalized number of circuits and topological fractions (Fig. 5F). While short-range contacts contained inside one circuit tend to be in series with local contacts present in other circuits, we also find that a relatively high number of normalized circuits can also be compatible with high percentages of entangled relationships. This enrichments in cross and parallel fractions can be due to the fact that circuits favor tight knit local interaction and tend to bring protein strands closer together. Moreover, circuits can also create long-range entangled relationships with each other, which are generally ignored in the computation of circuits, if they happen for residues which are more distant along the chain than the threshold for long-range exclusion.

## Discussion

Thanks to the theoretical framework of CT, we were able to draw a correlation between topological properties and folding kinetics, disentangling the role of topology from that of geometry. A significant step in the direction of topological description of folding phenomena was undertaken by Nikolay V. Dokholyan *et al.*, who demonstrated that average graph connectivity was a determinant of folding probability for pre-transition and post-transition states in the protein folding pathway.<sup>42</sup> Different approaches drawn from knot theory were also used to describe the entanglement, torsion and writhe of the protein backbone,<sup>15–18</sup> devising topologically inspired descriptors which yielded fairly good correlations with the logarithm of the protein folding rate.<sup>15,18</sup> Here, we have taken a fundamentally new step forward, by showing how folding rate can be predicted by CT parameters. Circuit topology (as presented in this study) only focuses on contacts, therefore ignoring the entanglement of the backbone. Moreover, this method does not require cumbersome mathematical and computational operations to connect the ends of the chain, such as those applied by Sulkowska *et al.*<sup>43</sup> CT not only considers the number of contacts in the protein, but also shows that there are differential patterns of correlations with respect to the topological arrangement of the contacts. Series, parallel and cross contacts are invariant with respect to shrinking,

bending, stretching and other continuous transformations,<sup>20</sup> and thus present true topological features of protein folds. Our analysis reveals that CT and CO have complementary ranges of validity and can be coupled to predict with accuracy the folding rate of a protein within a wide range of sizes and folding complexity. Moreover, CT offers invaluable information about what type of topological arrangements favor or hinder folding, therefore adding a mechanistic insight into folding rate prediction. The evidence of the zipper effect for short-range, attractive energy contacts offers a generalized model for folding which resolves the qualitative discrepancy between two-state and multi-state proteins. This unified view is beneficial since often attribution to two-state or multi-state classes is somewhat arbitrary,<sup>30</sup> and the folding state of a protein might also not be known *a priori*. Moreover, we found that the zipper effect yields a particularly high correlation for multi-state folders, which were previously found to mainly correlate with protein length.<sup>28</sup>

Although the presented implementation of CT ignores backbone entanglement, one can consider a comparison between the correlation scores obtained by CT analysis and those extracted by other topologically inspired descriptors such as torsion, writhe, Gaussian linking number and its linear combinations with relative and absolute contact order.<sup>15,18</sup> It is natural to compare the results obtained in Fig. 1F for segments to the analysis carried out by Panagiotou and Plaxco, about torsion and writhe of the protein backbone. They obtained correlation scores as high as 0.48 and 0.45 for writhe and torsion, respectively, with respect to the logarithm of the folding rate. We obtain comparable results when considering the parallel relationship. However, we only obtain it for multi-state folders, while writhe and torsion correlation values were only provided for two-state folders.<sup>18</sup> A combination of the two approaches might provide a more complete description for protein folding kinetics at the secondary structure level. For what concerns Gaussian entanglement, correlations as high as  $-0.64$  and  $-0.74$  were obtained for two and multi-state proteins, respectively,<sup>15</sup> with correlations increasing when these scores were combined with RCO and ACO. However, these results were obtained on small datasets (26 two- and 22 multi-state proteins); it is important to take into account that this type of analysis is sensitive to the size and characteristics of the dataset.<sup>15</sup> In fact, CT provides comparable scores when applied to smaller subsections of the datasets, with sizes comparable to the ones in these studies (Fig. 2 and 3). Moreover, combining CT with traditionally used descriptors such as CO and protein length allows for an increase in the predictive power of both parameters (Fig. 4). One might also consider the advantages of combining contact- and entanglement-based descriptors for protein folding prediction. Generalized CT<sup>25,44</sup> expands CT concepts to entangled subloops of a chain (the so-called *soft contacts*), therefore offering the opportunity for such complete description in future research.

The statistically significant correlations found between folding rate and the topology of short-range contacts, as well as the number of circuits, suggest that folding happens primarily at

the circuit level. We might find parallels between the concept of topological circuits and the one of local elementary structures (LESs) envisioned by a hierarchical model of protein folding.<sup>32,45</sup> Following this reasoning, one would envision a folding model in which folding occurs early on inside the circuits, and at a later stage the circuits are arranged with respect to each other, forming inter-circuit contacts. This type of folding model also matches the 'zipping and assembly' mechanism theorized by S. Bano Ozkan and co-workers.<sup>46</sup> This folding mechanism would be compatible with our observations that the topology of long-range contacts correlates with folding rate only for slow-folding proteins (Fig. 2E). Circuits presumably represent the elementary topological units of folding. The correlations between the normalized number of circuits and folding rate for two-state folders (and to a lesser extent, multi-state folders) indicate that, for proteins of comparable sizes, the ones that present multiple, small folding elementary units will fold faster. The high correlations obtained by two-state folders shed light on the nature of the different mechanisms experienced by two and multi-state proteins during folding. In particular, it would seem that the topological length, as opposed to the geometrical length, might play a role in folding rate prediction for two-state proteins.

These insights into the role of native topology offer not only new tools for the theoretical understanding of protein kinetics, but also powerful principles for protein design. The framework of CT has already proved to be effective in the field of molecular engineering.<sup>47</sup> The zipper effect and circuit decomposition might provide an easily applicable topological prescription for obtaining proteins with the desired kinetic properties.

## Methods

All proteins, CO and kinetic information were retrieved from the ARCPro dataset.<sup>29</sup> Contact order retrieved from the ARCPro dataset was computed as the absolute contact order (ACO) based on a 6 Å cut-off value for determining contacts by a multiple contact all-heavy atom method. Four proteins were excluded from analysis: 1FMK, 1M9S and 2BLM for the incompleteness of structural information in the PDB files, and 1RA9 for incompleteness in kinetic information in the dataset. Therefore, the whole dataset for analysis comprised 122 proteins. The sub-datasets contained the following number of proteins: 36 proteins for Lower CO (multi-state: 7 and two-state: 27), 58 proteins for Average CO (multi-state: 18 and two-state: 40) and 28 for Upper CO (multi-state: 16 and two-state: 12). The partitioning of the dataset into CO ranges was made by calculating mean  $x$  and standard deviation  $\sigma_x$  of the CO distribution and defining the following thresholds:

$$t_{\text{Upper}} = x + \frac{\sigma_x}{2}$$

$$t_{\text{Lower}} = x - \frac{\sigma_x}{2}$$

Circuit topology parameters were retrieved using our custom-made Python code, which allows for energy, length filtering and circuit decomposition options. All PDBs are

pre-processed automatically before analysis, in order to remove water molecules, hydrogen atoms and various binders. Only one chain (the first contained in the PDB) is selected.

Contacts between segments were calculated based on a distance cut-off of 5.0 Å and a cut-off in number of atoms equal to 10. See the ESI† for distance cut-off values equal to 3.5, 4.0, 5.0, 5.5 and 6 Å (Fig. S4, ESI†). For the definition of segments, the secondary structure files of the proteins as produced by the free web service STRIDE<sup>48</sup> were used. Each secondary structural element as defined in the STRIDE file represents a segment, to which contacts formed by atoms included in the segment are assigned.

Contacts between residues were calculated based on a distance cut-off of 5.0 Å. Residues were deemed to be in contact when more than  $n_a = 5$  atoms were found to be closer than the distance cut-off. We repeated the analysis for cut-off values equal to 3.5, 4.0, 4.5, 5.5 and 6.0 Å (Fig. S5, ESI†) and for  $n_a = 1, 2, 3, 4, 5$  and 6 (Fig. S6, ESI†). The four closest neighbors of each residue were excluded from analysis. Each retrieved contact site in the protein structure was given an index. Indexes were given based on the order in which contact sites appeared along the protein chain, from the left end to the right end of the chain. In this way, each contact was characterized by the two indexes ( $i, j$ ) of its constituent contact sites. In order to define the CT relationship between two contacts, their contact indexes ( $i, j$ ) and ( $r, s$ ) were compared. CT relationships were assigned based on the mathematical relationships summarized below:

$$C_{i,j}SC_{r,s} \Leftrightarrow [i,j] \cap [r,s] = \emptyset$$

$$C_{i,j}PC_{r,s} \Leftrightarrow [i,j] \subset (r,s)$$

$$C_{i,j}XC_{r,s} \Leftrightarrow [i,j] \cap [r,s] \notin \{[i,j], [r,s]\} \cup \mathcal{P}(\{i,j,r,s\})$$

$$C_{i,j}CSC_{r,s} \Leftrightarrow (([i,j] \cap [r,s] = \{i\}) \vee ([i,j] \cap [r,s] = \{j\}))$$

$$C_{i,j}CPC_{r,s} \Leftrightarrow (([i,j] \subset [r,s]) \wedge (i = r \vee j = s))$$

where  $\mathcal{P}$  denotes the powerset *i.e.*, all subsets of a set including the null set ( $\emptyset$ ). The topological relationships introduced above are sufficient and necessary to describe the topology of any folded linear chain with binary contacts.<sup>20</sup> For simplicity, CP and CS relationships were counted respectively as parallel and series in the analysis presented in this paper. One can readily adjust the set theory definition to reduce the relation set  $\{P, S, X, CP, CS\}$  to  $\{P, S, X\}$  and to make the parallel relation symmetric so that  $P = P^{-1}$ :

$$\text{Series: } C_{i,j}SC_{r,s} \Leftrightarrow [i,j] \cap [r,s] \subset \{i,j,r,s\}$$

$$\text{Parallel } C_{i,j}PC_{r,s} \Leftrightarrow [i,j] \subset [r,s] \vee [r,s] \subset [i,j]$$

$$\text{Cross: } C_{i,j}XC_{r,s} \Leftrightarrow [i,j] \cap [r,s] \notin \{[i,j], [r,s]\} \cup \mathcal{P}(\{i,j,r,s\})$$

Correlation analysis for segments and residues subdivided in CO subgroups, for different distance cut-off values, can be seen in Fig. S7 and S8 (ESI†). Distance filtering (short-range *versus* long-range contacts) was carried out with a threshold for long-range exclusion of 24 residues. The analysis was also

repeated for thresholds equal to 12 and 36 residues (Fig. S9, ESI†). Energy filtering was carried out by exploiting the statistical potential matrix calculated by P. Thomas and K. Dill.<sup>37</sup> The Pearson correlation coefficient and two-tailed  $p$  values were calculated using custom-made data analysis Jupyter lab files. All correlation maps shown in the paper display correlations with  $p$  value  $\leq 0.05$ .

Multilinear regression was performed by using an ordinary least squares linear regression from the linear\_model module in scikit-learn 0.24.2. The subdivision into subsequent training and test sets was performed by the model\_selection module, with the K Fold function. The five sets are formed respectively by protein 1 to 25, protein 26 to 49, protein 50 to 73, protein 74 to 97 and 98 to 121. Indexes refer to those assigned to proteins in the ARCPro database.

Residuals from folding rate prediction were tested for normality with the Shapiro test. Distributions with  $P$  values  $< 0.05$  were considered not normal. In order to evaluate the quality of the prediction, the determination coefficient  $R^2$  was used, as calculated by the metrics.r2\_score function:

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

where  $\hat{y}_i$  is the predicted value of the  $i$ -th data point,  $y_i$  is the corresponding true value,  $n$  is the total number of samples and  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . The adjusted determination coefficient is defined as:

$$R_{\text{adj}}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - p - 1}$$

where  $n$  is the number of samples and  $p$  is the number of predictors (independent variables).

Circuit decomposition and counting were performed by setting a threshold on the length of the circuits. Given the average length  $l$  of the circuits in a protein, and  $\sigma_l$  their standard deviation, the circuits with a length below a threshold  $t_1 = l - \frac{\sigma_l}{2}$  were discarded. Exclusion of smaller circuits was done under the assumption that the folding speed of bigger circuits represents the bottleneck for the folding rate. Results without application of threshold  $t_1$  are displayed in Fig. S10 (ESI†).

## Data availability

Data available on request from the authors.

## Code availability

All codes are available through the following links: CT analysis for proteins: [https://github.com/circuittopology/circuit\\_topology](https://github.com/circuittopology/circuit_topology); Correlation analysis for CT and folding rate: [https://github.com/Barbaraleidenuniv/Topology\\_analysis](https://github.com/Barbaraleidenuniv/Topology_analysis). For permissions, please contact the corresponding author.



## Author contributions

A. M. conceived the project and supervised the research. B. S. performed the research. B. S., V. S. and A. M. analyzed the data and wrote the paper.

## Conflicts of interest

There are no conflicts to declare.

## Acknowledgements

The authors thank Jaie Woodard and Duane Moes for help with the codes, Elnaz Bani-Jamali for insight into the protein structure analysis, and Enrico Liscio for useful suggestions about predictive data analysis. The research in Mashaghi lab is supported by funding from Muscular Dystrophy Association (USA), Grant Number MDA628071, and Dutch Research Council (Nederlandse Organisatie voor Wetenschappelijk Onderzoek) through Open Competition XS (OCENW.XS.076).

## References

- 1 D. Baker, A surprising simplicity to protein folding, *Nature*, 2000, **405**, 39–42.
- 2 K. W. Plaxco, K. T. Simons and D. Baker, Contact order, transition state placement and the refolding rates of single domain proteins, *J. Mol. Biol.*, 1998, **277**, 985–994.
- 3 D. N. Ivankov, *et al.*, Contact order revisited: Influence of protein size on the folding rate, *Protein Sci.*, 2003, **12**, 2057–2062.
- 4 P. Sormanni, *et al.*, Simultaneous quantification of protein order and disorder, *Nat. Chem. Biol.*, 2017, **13**, 339–342.
- 5 H. Kaya and H. S. Chan, Solvation Effects and Driving Forces for Protein Thermodynamic and Kinetic Cooperativity: How Adequate is Native-centric Topological Modeling?, *J. Mol. Biol.*, 2003, **326**, 911–931.
- 6 H. Zhou and Y. Zhou, Folding rate prediction using total contact distance, *Biophys. J.*, 2002, **82**, 458–463.
- 7 L. Censoni and L. Martínez, Prediction of kinetics of protein folding with non-redundant contact information, *Bioinformatics*, 2018, **34**, 4034–4038.
- 8 Y. Li, Y. Zhang and J. Lv, An Effective Cumulative Torsion Angles Model for Prediction of Protein Folding Rates, *Protein Pept. Lett.*, 2020, **27**, 321–328.
- 9 D. N. Ivankov and A. V. Finkelstein, Prediction of protein folding rates from the amino acid sequence-predicted secondary structure, *Proc. Natl. Acad. Sci. U. S. A.*, 2004, **101**, 8942–8944.
- 10 H. Gong, D. G. Isom, R. Srinivasan and G. D. Rose, Local secondary structure content predicts folding rates for simple, two-state proteins, *J. Mol. Biol.*, 2003, **327**, 1149–1154.
- 11 M. M. Gromiha and S. Selvaraj, Comparison between long-range interactions and contact order in determining the folding rate of two-state proteins: application of long-range order to folding rate prediction, *J. Mol. Biol.*, 2001, **310**, 27–32.
- 12 Z. Ouyang and J. Liang, Predicting protein folding rates from geometric contact and amino acid sequence, *Protein Sci.*, 2008, **17**, 1256–1263.
- 13 M. M. Gromiha, Multiple Contact Network Is a Key Determinant to Protein Folding Rates, *J. Chem. Inf. Model.*, 2009, **49**, 1130–1135.
- 14 L. L. Chavez, J. N. Onuchic and C. Clementi, Quantifying the roughness on the free energy landscape: Entropic bottlenecks and protein folding rates, *J. Am. Chem. Soc.*, 2004, **126**, 8426–8432.
- 15 M. Baiesi, E. Orlandini, F. Seno and A. Trovato, Exploring the correlation between the folding rates of proteins and the entanglement of their native states, *J. Phys. A: Math. Theor.*, 2017, **50**, 504001.
- 16 M. Baiesi, E. Orlandini, F. Seno and A. Trovato, Sequence and structural patterns detected in entangled proteins reveal the importance of co-translational folding, *Sci. Rep.*, 2019, **9**, 1–12.
- 17 M. Baiesi, E. Orlandini, A. Trovato and F. Seno, Linking in domain-swapped protein dimers, *Sci. Rep.*, 2016, **6**, 1–11.
- 18 E. Panagiotou and K. W. Plaxco, A topological study of protein folding kinetics, *arXiv*, 2018, 1–13, DOI: 10.1090/conm/746/15010.
- 19 E. Shakhnovich, To knot or not to knot?, *Nat. Mater.*, 2011, **10**, 84–86.
- 20 A. Mashaghi, R. J. Van Wijk and S. J. Tans, Circuit topology of proteins and nucleic acids, *Structure*, 2014, **22**, 1227–1237.
- 21 A. Mashaghi, Circuit Topology of Folded Chains, *Notices Amer. Math. Soc.*, 2021, **68**(3), 420–423.
- 22 M. Heidari, H. Schiessel and A. Mashaghi, Circuit Topology Analysis of Polymer Folding Reactions, *ACS Cent. Sci.*, 2020, **6**(6), 839–847.
- 23 B. Scalvini, *et al.*, Topology of Folded Molecular Chains: From Single Biomolecules to Engineered Origami, *Trends Chem.*, 2020, **2**, 609–622.
- 24 A. Mugler, S. J. Tans and A. Mashaghi, Circuit topology of self-interacting chains: implications for folding and unfolding dynamics, *Phys. Chem. Chem. Phys.*, 2014, **16**, 22537–22544.
- 25 A. Golovnev and A. Mashaghi, Generalized Circuit Topology of Folded Linear Chains, *iScience*, 2020, **23**, 101492.
- 26 K. L. Maxwell, *et al.*, Protein folding: Defining a “standard” set of experimental conditions and a preliminary kinetic data set of two-state proteins, *Protein Sci.*, 2005, **14**, 602–616.
- 27 D. Barrick, What have we learned from the studies of two-state folders, and what are the unanswered questions about two-state protein folding?, *Phys. Biol.*, 2009, **6**, 015001.
- 28 O. V. Galzitskaya, S. O. Garbuzynskiy, D. N. Ivankov and A. V. Finkelstein, Chain length is the main determinant of the folding rate for proteins with three-state folding kinetics, *Proteins: Struct., Funct., Genet.*, 2003, **51**, 162–166.
- 29 A. S. Wagaman, A. Coburn, I. Brand-Thomas, B. Dash and S. S. Jaswal, A comprehensive database of verified experimental data on protein folding kinetics, *Protein Sci.*, 2014, **23**, 1808–1812.

- 30 S. E. Jackson, How do small single-domain proteins fold?, *Fold. Des.*, 1998, **3**, 81–91.
- 31 R. Doyle, K. Simons, H. Qian and D. Baker, Local Interactions and the Optimization of Protein Folding, *Proteins*, 1997, **291**, 282–291.
- 32 R. A. Broglia and G. Tiana, Hierarchy of events in the folding of model proteins, *J. Chem. Phys.*, 2001, **114**, 7267–7273.
- 33 P. Cossio, *et al.*, Exploring the universe of protein structures beyond the protein data bank, *PLoS Comput. Biol.*, 2010, **6**(11), e1000957.
- 34 A. Irbäck, C. Peterson, F. Potthast and O. Sommelius, Local interactions and protein folding: A three-dimensional off-lattice approach, *J. Chem. Phys.*, 1997, **107**(1), 273–282.
- 35 V. I. Abkevich, A. M. Gutin and E. I. Shakhnovich, Impact of local and non-local interactions on thermodynamics and kinetics of protein folding, *J. Mol. Biol.*, 1995, **252**, 460–471.
- 36 A. Kumar, A. Baruah and P. Biswas, Role of local and nonlocal interactions in folding and misfolding of globular proteins, *J. Chem. Phys.*, 2017, **146**(6), 065102.
- 37 P. D. Thomastt and K. E. N. A. Dill, An iterative method for extracting energy-like quantities from protein structures, *Proc. Natl. Acad. Sci. U. S. A.*, 1996, **93**, 11628–11633.
- 38 P. D. Thomas and K. A. Dill, Statistical potentials extracted from protein structures: how accurate are they?, *J. Mol. Biol.*, 1996, **257**, 457–469.
- 39 O. Kramer, Scikit-Learn, *Machine Learning for Evolution Strategies*, Springer International Publishing, 2016, pp. 45–53, DOI: 10.1007/978-3-319-33383-0\_5.
- 40 A. Di Bucchianico, Coefficient of Determination ( $R^2$ ), in *Encyclopedia of Statistics in Quality and Reliability*, John Wiley & Sons, Ltd, 2008, DOI: 10.1002/9780470061572.eqr173.
- 41 S. W. Englander and L. Mayne, The nature of protein folding pathways, *Proc. Natl. Acad. Sci. U. S. A.*, 2014, **111**, 15873–15880.
- 42 N. V. Dokholyan, L. Li, F. Ding and E. I. Shakhnovich, Topological determinants of protein folding, *Proc. Natl. Acad. Sci. U. S. A.*, 2002, **99**, 8637–8641.
- 43 J. I. Sulkowska, E. J. Rawdon, K. C. Millett, J. N. Onuchic and A. Stasiak, Conservation of complex knotting and slipknotting patterns in proteins, *Proc. Natl. Acad. Sci. U. S. A.*, 2012, **109**, E1715–E1723.
- 44 J. Ceniceros, M. Elhamdadi and A. Mashaghi, Coloring Invariant for Topological Circuits in Folded Linear Chains, *Symmetry*, 2021, **13**(6), 919.
- 45 R. L. Baldwin and G. D. Rose, Is protein folding hierarchic? I. Local structure and peptide folding, *Trends Biochem. Sci.*, 1999, **0004**, 26–33.
- 46 S. B. Ozkan, G. A. Wu, J. D. Chodera and K. A. Dill, Protein folding by zipping and assembly, *Proc. Natl. Acad. Sci. U. S. A.*, 2007, **104**, 11987–11992.
- 47 V. Kočar, *et al.*, Design principles for rapid folding of knotted DNA nanostructures, *Nat. Commun.*, 2016, **7**, 10803.
- 48 M. Heinig and D. Frishman, STRIDE: a web server for secondary structure assignment from known atomic coordinates of proteins, *Nucleic Acids Res.*, 2004, **32**, W500–W502.