



Analyzing the symmetrical arrangement of structural repeats in proteins with CE-Symm

Spencer E Bliven^{1,2,3}^{*}, Aleix Lafita^{1,4,5}^{*}, Peter W Rose^{6,7}, Guido Capitani^{1,8}[†], Andreas Prlić⁶, Philip E Bourne^{2,9}

1 Laboratory of Biomolecular Research, Paul Scherrer Institute, Villigen, Switzerland

2 National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA

3 Institute of Applied Simulation, Zurich University of Applied Science, Wädenswil, Switzerland

4 Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland

5 European Bioinformatics Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, Cambridgeshire, UK

6 RCSB Protein Data Bank, San Diego Supercomputing Center, University of California San Diego, La Jolla, CA, USA

7 Structural Bioinformatics Laboratory, San Diego Supercomputing Center, University of California San Diego, La Jolla, CA, USA

8 Department of Biology, ETH Zurich, Zurich, Switzerland

9 Department of Biomedical Engineering, University of Virginia, Charlottesville, VA, USA

 These authors contributed equally to this work.

[†]Deceased

* spencer.bliven@zhaw.ch, aleixlafita@ebi.ac.uk

Abstract

Many proteins fold into highly regular and repetitive three dimensional structures. The analysis of structural patterns and repeated elements is fundamental to understand protein function and evolution. We present recent improvements to the *CE-Symm* tool for systematically detecting and analyzing the internal symmetry and structural repeats in proteins. In addition to the accurate detection of internal symmetry, the tool is now capable of i) reporting the type of symmetry, ii) identifying the smallest repeating unit, iii) describing the arrangement of repeats with transformation operations and symmetry axes, and iv) comparing the similarity of all the internal repeats at the residue level. *CE-Symm 2.0* helps the user investigate proteins with a robust and intuitive sequence-to-structure analysis, with many applications in protein classification, functional annotation and evolutionary studies. We describe the algorithmic extensions of the method and demonstrate its applications to the study of interesting cases of protein evolution.

Availability: *CE-Symm* is an open source tool integrated into the BioJava library (www.biojava.org) and freely available at <https://github.com/rcsb/symmetry>.

Author summary

Many protein structures show a great deal of regularity. Even within single polypeptide chains, about 25% of proteins contain self-similar repeating structures, which can be organized in ring-like symmetric arrangements or linear open repeats. The repeats are often related, and thus comparing the sequence and structure of repeats can give an idea as to the early evolutionary history of a protein family. Additionally, the conservation and divergence of repeats can lead to insights about the function of the proteins.

This work describes *CE-Symm 2.0*, a tool for the analysis of protein repeats. It automatically detects repeats in protein structures and produces a single alignment of all repeats to the user. The algorithm is able to detect the geometric relationships between repeats, including cyclic, dihedral, and polyhedral symmetry, translational repeats, and cases where multiple symmetry operators are applicable in a hierarchical fashion. These complex relationships can then be visualized in a graphical interface as a complete structure, as a superposition of repeats, or as a multiple alignment of the protein sequence. *CE-Symm 2.0* can be used for the automatic detection of internal symmetry in protein structures, or as an interactive tool for the analysis of structural repeats.

Introduction

François Jacob described molecular evolution as a “tinkering” process, where pre-existing elements are combined and repurposed to solve new biological problems [1]. Traces of this “tinkerer evolution” can be seen in the widespread reuse of structural elements in proteins at different scales: small motifs [2], functional domains [3], and protein oligomerization [4]. One example is the repetition of structural elements within a protein chain, thought to arise from gene fusion and duplication events [5].

It is common for structural repeats in proteins to maintain a symmetric arrangement [6], which has been associated with many biological functions [7]. The internal symmetry of proteins is thought to arise from ancestral quaternary structures fused into a single polypeptide chain [8–10]. However, since symmetric protein folds theoretically have a folding thermodynamic advantage, their symmetry could have arisen by evolutionary convergence [11]. On the other hand, the evolution of functional patches is often symmetry breaking [12]. High-quality alignments of structural repeats are essential to resolve these opposing evolutionary explanations and understand the tension between conservation and divergence.

A number of tools have been developed for the detection of structural repeats and internal symmetry in proteins [6, 13–20]. Similarly to our *CE-Symm* method [21], a common approach to the problem is the identification of regions of similar structure within a protein chain, usually through the alignment of a protein structure to itself (self-alignment). However, the systematic extraction of repeating structural units and residue equivalencies among the repeats from the self-alignments is a nontrivial task. Here, we present an extension of *CE-Symm* (version 2.0) that accurately detects symmetry in proteins and defines the boundaries of the structural repeating elements. In addition, it reports the type of symmetry and describes the arrangement of repeats using the symmetry axes. Finally, the similarity of all the structural repeats can be compared at the residue level in a multiple structure alignment.

Types of symmetry

Several definitions of **internal symmetry** and **repeats** are possible, depending on the biological question of interest. For the purposes of this paper, we define it as the regular arrangement of a common repeating structural unit within a protein chain. Therefore, a

repeat is an asymmetric structural motif present multiple times in the same structure. We restrict our consideration of repeats to cases where the orientation between adjacent structural units is regular; that is, where a consistent geometric transformation can be applied to superimpose each repeat onto the next. In other words, *CE-Symm* focuses on identifying repeats which conserve both the structure and the interface between repeats.

Several types of internal symmetry can be derived from this broad definition. The most basic division is between closed symmetry and open symmetry. In proteins with **closed symmetry**, the repeats are arranged in a point group symmetry. This can be defined mathematically as a set of rotations that superimpose equivalent repeats while keeping at least one point at the center of rotation fixed. In contrast, repeats in proteins with **open symmetry** are related by transformations with a translational component. Examples of closed and open symmetry can be found in Fig. 1a-d and Fig. 1e-h, respectively.

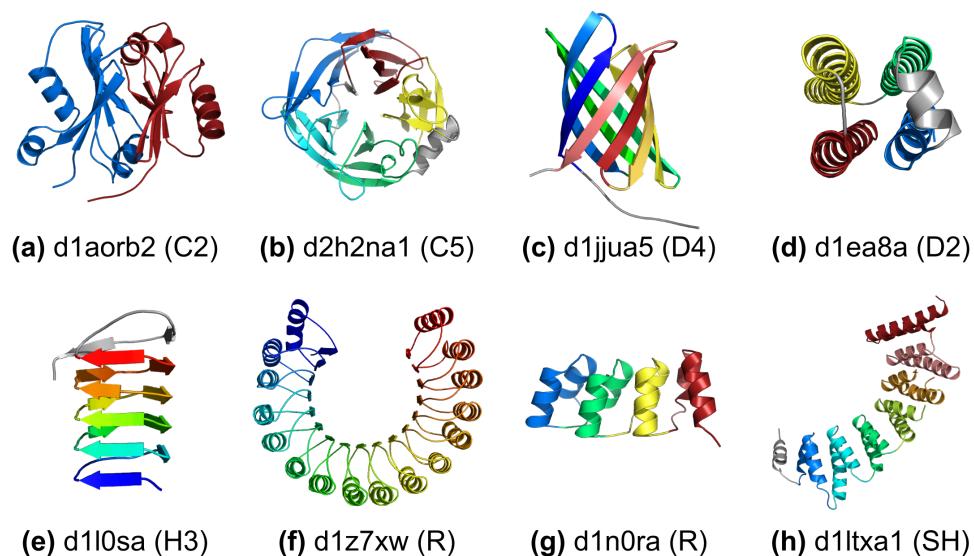


Fig 1. Examples of protein domains with internal symmetry. Protein domains are labeled with SCOP domain identifier [22]. a) N-terminal domain of aldehyde ferredoxin oxidoreductase with 2-fold rotational symmetry (C2) of an alpha+beta motif; b) 5-bladed beta propeller with a helical insertion between second and third blades with 5-fold rotational symmetry (C5) of a 4-stranded beta-sheet motif; c) Beta-barrel with 8 beta strands in a 4-fold dihedral symmetry (D4) of single stranded motifs; d) 4-helical bundle with dihedral symmetry connectivity (D2) of single helical motifs; e) Beta-helix with single stranded right-handed helical symmetry (H3); f) Leucine rich repeats with open rotational symmetry (R) of 16 up and down alpha-beta motifs; g) Designed Ankyrin repeat protein with 4 translational repeats (R) of double helical motifs; h) Alpha-alpha right-handed superhelix (SH) of double helical motifs. Repeats are colored from blue, N-terminal, to red, C-terminal. Non-repeating parts of the structure are colored in grey.

Closed symmetries can be further characterized according to the possible chiral point groups: cyclic (C_n), generated by a single n -fold rotational operator (Fig. 1a-b); dihedral (D_n), which requires an n -fold rotation and n perpendicular 2-fold operators (Fig. 1c-d); and polyhedral point-groups (T, O, and I), which feature non-perpendicular rotation operators. Both cyclic and dihedral internal symmetries are common in proteins, but, although common at the quaternary structure level, polyhedral symmetries have not yet been observed within a single polypeptide chain.

Open symmetry can be further subdivided into special cases of helical, translational, and superhelical repeats. Helical symmetry consists of repeats arranged around a screw axis, where each repeat is related to the next by a fixed linear translation combined by a rotation around the central axis (Fig. 1e). In cases where the rotation angle is close to an fraction of a turn, we indicate the approximate number of subunits needed per turn (Hn). Proteins with open symmetry that have negligible translation are called rotational repeats (Fig. 1f), and those with negligible rotation between repeats are called translational repeats (Fig. 1g), both annotated as R. Superhelical symmetry (SH) provides the most general description of repeats with open symmetry, and is reserved for cases which cannot be expressed as a single fixed operator relating each repeat to the next. Instead, the rotation axis between adjacent repeats precesses along a helical path (Fig. 1h). Proteins with open symmetry are sometimes referred to as solenoid proteins [23].

Methods

CE-Symm analyzes the symmetry in a protein structure and produces a multiple alignment of all repeats, as well as ancillary information about the type and order of symmetry in the structure. An overview of *CE-Symm* alignment steps is shown in Fig. 2. These are described in detail below, but consist of (1) structural self-alignment, (2) order detection, (3) refinement to a multiple alignment, (4) Monte Carlo optimization of the multiple alignment, and (5) point group symmetry detection. These steps are repeated iteratively to detect multiple levels of symmetry (hierarchical symmetry) and higher-order point groups.

Self-alignment

CE-Symm begins with a structural self-alignment (other than the identity alignment) of the input protein structure using the Combinatorial Extension (CE) algorithm [24]. Identifying significant self-alignments was the primary focus of the first version of the algorithm [21]. In the self-alignment of structures with closed symmetry the first and last repeats are aligned, forming a circular permutation (CP) of the structure. This is why the structure alignment method used in *CE-Symm* shares algorithmic primitives with *CE-CP* [25]. For proteins with open symmetry, the initial self-alignment will always be missing one of the repeats due to the translation component of the symmetry operator.

The alignment quality is quantified using TM-Score [26]. Both inconsistently arranged repeats and large asymmetric regions in a structure will reduce the score of the self-alignment. In addition, open symmetry will generally have lower scores than closed symmetry, because the terminal repeats will be missing from the self-alignment.

Order of symmetry detection

The order of symmetry is defined as the number of symmetric units (repeats) in a structure. Extracting the order of symmetry is a key part of symmetry detection, and posterior analyses heavily depend on it. Several algorithms were evaluated for automatically determining the order [27].

The most straightforward method for determining the order in cases of closed symmetry is based on the angle of rotation, which can be calculated based on the superposition operator [18]. The distance between a measured angle of rotation, θ , and the closest theoretical angle of rotation for order k is given by a triangle wave of frequency k :

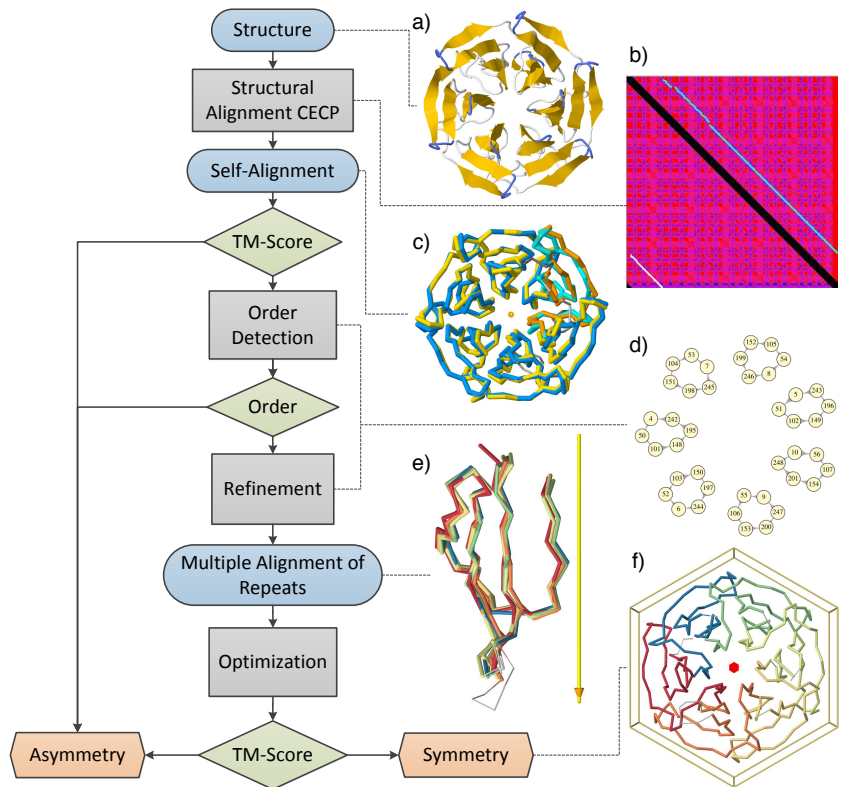


Fig 2. Flowchart of one iteration of the *CE-Symm* algorithm. Algorithm steps are grey rectangles, inputs and outputs are blue rounded rectangles, decision rules are green rhomboid boxes and final classifications are orange hexagonal rectangles. Additional iterations on the resulting repeats may be performed to detect further symmetry axes or hierarchical symmetry. The images on the right represent, from top to bottom: a) initial structure, colored by secondary structure elements; b) self-alignment dot-plot matrix, where similarity score is a range from blue (high similarity) to magenta (low similarity), the identity alignment is blacked out and the optimal self-alignment path is in white; c) superposition of the structure against itself based on the optimal self-alignment, where the original structure is in blue and cyan and a copy of the structure is in yellow and orange (orange and cyan correspond to the regions of the alignment involving a circular permutation); d) subset of the alignment graph with seven connected components of six aligned residues each; e) superposition of the six internally symmetric repeats according to the symmetry axis (yellow bar) and their residue equivalencies; and f) structure inside the six-fold cyclic symmetry (C_6) box, with repeats colored differently.

$$\delta(\theta, k) = \frac{2\pi}{k} \left| \left(\left\lfloor \frac{\theta k}{2\pi} - \frac{1}{2} \right\rfloor \bmod 1 \right) - \frac{1}{2} \right| \quad (1)$$

The best-fit order is then the k that minimizes this distance, up to some maximum order.

Detecting the order based on the rotation angle works only for cases of closed symmetry. Another method, called graph-component order, can be used for both open and closed symmetry. Conceptually, the self-alignment is treated as a directed graph over the set of aligned residues. Residues that are aligned in all k repeats will form a path with k nodes. For open symmetry these paths tend to be disjoint, so simply finding the most frequent size of the connected components in the graph can accurately determine the order for open symmetry. For well-aligned cases of closed symmetry, the aligned residues form a cycle of k nodes, so the same method can also work in the general case. Those residues which participate in a path or cycle of the most frequent size form the refined alignment discussed in the following section.

For cases of closed symmetry, small alignment errors can lead to a situation where repeated applications of the alignment do not eventually form a closed cycle back to the original residue, but rather to a residue at a small offset. This can lead to failures of the graph-component order detector due to the merging of multiple alignment paths. This case can be handled by the sequence-function order detector. For each potential order k under consideration, all paths of length k are considered in that alignment graph. The number of residues separating the start and end of the path becomes a good indicator for the agreement between the graph and order k . Orders between 1 and some k_{max} (8 by default) are considered, with a pronounced decrease in average distance (40% by default) indicating that the correct order has been discovered. This method was already described in more detail in our previous publication [21].

By default, *CE-Symm* first runs the graph component order detector and, if the symmetry is determined to be closed, the sequence-function method is then used to improve the order determination.

Refinement to a multiple alignment

The refinement procedure takes as input the self-alignment of the structure and the order of symmetry and returns a multiple alignment of the repeats. *CE-Symm* has two implementations of the refinement procedure: graph-component and sequence-function, which are closely related to their respective order detectors.

The graph-component refiner combines the maximally connected components of the self-alignment graph with size equal to the order of symmetry. Each connected component contributes one column to the refined alignment of repeats, taking special care that the repeat sequences preserve the sequence order of the domain. A heuristic method is used to decide which components to include in the refined multiple alignment.

The sequence-function refiner uses all the path of length k (the order of symmetry) to construct a multiple alignment of the symmetric repeats. The residues in each path are set as one column in the alignment, sorted in increasing order, resulting in a multiple alignment of the repeats. Note that, like the graph-component refiner, the multiple alignments obtained at the end of this stage do not contain gaps, so all repeats are of the same size.

Optimization

The multiple alignment obtained from the refinement is sometimes far from optimal, and depends very much on the goodness and consistency of the self-alignment. In addition, the refinement process prioritizes precision over coverage, which means that

only the best residue equivalencies will be included, resulting in a shorter multiple alignment. The goal of the optimization is to increase the multiple alignment length while keeping the RMSD low. Furthermore, the optimization procedure can improve parts of the alignment that were not fully represented in the self-alignment, and thus not captured in the refinement result.

The optimization process uses a similar approach to the Combinatorial Extension Monte Carlo (CEMC) multiple structure alignment algorithm [28]. The multiple alignment can be described by a matrix, where the rows represent aligned structures and the columns represent aligned positions (residue equivalencies). Rearranging and modifying the entries of the matrix results in changes of the multiple alignment. There are four possible moves (changes in the multiple alignment):

1. **Expand**: increase the alignment length by extending the boundary of a group of sequential residues, chosen randomly. This move requires the addition of an alignment column.
2. **Shrink**: decrease the alignment length by decreasing the boundary of a group of sequential residues, chosen randomly. This move requires a deletion of an alignment column.
3. **Shift**: move a group of sequential residues, chosen randomly, of one structure (row), chosen randomly, one position to the right or to the left.
4. **Insert gap**: delete one entry of the matrix, chosen randomly. This is equivalent to inserting a gap in one residue position (column) of one structure (row).

The insertion of gaps allows for partial repeat similarities in the alignment. All moves take into consideration that rows of the alignment occur sequentially in the protein sequence, so unaligned residues between repeats can be considered either at the end of a repeat or the beginning of the following one. In addition, the shrink and insert gap moves have been biased, so that the probability of choosing an alignment column or an equivalent residue, respectively, is proportional to the average inter-residue distance of the given column or the given residue, respectively. A geometric distribution with parameter 0.5 is chosen to allocate the probability among alignment columns. A schematic representation of the steps and how they affect the multiple alignment is provided in supplementary figure S1.

After each optimization step, an alignment score is calculated. The score function to be optimized has also been smoothed with respect to the original CEMC score to remove discontinuities:

$$S = \sum_{i=0}^N \sum_{k=0}^L \left[\frac{C}{1 + \left(\frac{d_{ik}}{d_0}\right)^2} - A \right] - G \quad (2)$$

N is the number of structures (rows) in the alignment; L is the number of equivalent positions (columns) in the alignment, including gaps; C is the maximum score of an alignment position (by default set to 20); d_{ik} is the average distance from aligned residue k in structure i to all its equivalent residues; d_0 is the structural similarity function parameter, as defined by the TM-score [26]; A is the distance cutoff penalization, which shifts the function to negative values when the maximum allowed average distance of an aligned position (d_c) is higher than d_{ik} ; and G is a linear gap penalty term. Calculation of A using a distance cutoff parameter d_c (by default set to 7Å) is straightforward from the condition that the score S has to be 0 when $d_{ik} = d_c$. The shape of the score function for different values of d_c is shown in supplementary figure S2.

The acceptance probability of a move is proportional to the score difference and decreases proportional to the number of optimization steps. 188
189

$$p = \left[\frac{C - \Delta S}{C\sqrt{m}} \right] \left(1 - \frac{m}{M} \right) \quad (3)$$

m is the current iteration number, ΔS is the change in alignment score and M is the maximum number of iterations. The maximum number of optimization iterations is proportional to the length of the domain, by default a hundred times the number of residues in the protein. Optimization finishes either because it reaches the maximum number of iterations or in the case that no moves are accepted for a fraction of the total number of iterations (by default M divided by 50). 190
191
192
193
194
195

Recursive symmetry detection 196

So far, the procedure described can only identify symmetry operations that require a single axis. However, some structures present symmetries represented with more than one axis. This is the case for point groups other than cyclic, like dihedral symmetry, or structures with more than one level of symmetry, what we define as hierarchical symmetries. Multiple *CE-Symm* iterations are run in a recursive manner, i.e. repeats found in previous rounds are recursively fed into the next run until a non-significant result (no symmetry) is found. The goal is to find all the significant symmetry levels of a structure. 197
198
199
200
201
202
203
204

At the end of an iteration, repeats are extracted from the internal symmetry result and one of them is chosen as the representative, by default the N-terminal repeat. Results of successive iterations are merged by combining the symmetry axes and multiple alignments, generating a unique result for the query structure. 205
206
207
208

The recursive symmetry detection allows better order determination for difficult cases (e.g., TIM barrels), because usually fractions of the order of symmetry are initially found (e.g., 2-fold instead of 8-fold). Continuing the analysis recursively breaks the structure down to the true asymmetric repeating units (e.g., with three levels of symmetry: 2-fold, 4-fold and finally 8-fold). 209
210
211
212
213

Significance 214

There are three decision checkpoints in the algorithm flowchart in Fig. 2. The first significance criterion for a symmetry result is the self-alignment TM-score. Like in the previous version, the default threshold value is set to 0.4. The second significance criterion is the order of symmetry. A symmetric structure must have symmetry order greater than 1 and the refinement of the self-alignment into a multiple repeat alignment has to be successful. The third significance criterion is the average TM-score of the multiple alignment of repeats, defined as the average TM-score of all pairwise repeat alignments. The default threshold value for the average TM-score is set to 0.36, because a 10% decrease from the original TM-score is allowed after refinement due to the restrictive conditions imposed on it. In addition the number secondary structure elements (SSE) of the final asymmetric repeating unit is considered. If the the number of SSE of each repeat is lower than the threshold, the result will be considered non-significant. For many applications it may be desirable to exclude simple repeat units (e.g. helical bundle proteins), but these are included in *CE-Symm* analysis by default in order to find the highest possible symmetry in a structure. 215
216
217
218
219
220
221
222
223
224
225
226
227
228
229

Symmetry type determination

The recursive symmetry detection identifies a collection of symmetry axes that describe the arrangement of repeats in the query structure. In many cases, several of these axes can be combined to form higher-order symmetries. For example, a two-fold rotation axis can be combined with another orthogonal axis to form dihedral symmetry. Near-identical rotation axes can also be combined to form higher-order rotational symmetry.

To determine the point group symmetry, we build on the algorithm described by Levy *et al.* [29]. The symmetry axes can be found efficiently by first considering only the centroids of each repeat, since they must be in a symmetric configuration if the entire complex is symmetric. To find all possible symmetry axes, the centroids are rotated around axes that go through the centroid of the whole structure using an orientation grid search in quaternion space [30]. For each orientation, the RMSD of the aligned centroids is calculated. If the centroids align within a threshold, then all $C\alpha$ atoms are superimposed. The symmetry axis is then defined by the rotation matrix of this superposition. If the RMSD is less than a threshold value (i.e., 5 Å), the symmetry operation is considered valid. Since symmetry operations form a group, only a few are needed to complete the full point group

This procedure allows the combination of axes that have been considered separately by *CE-Symm*. The point group is included in the final symmetry output and displayed to the user as a polyhedron box around the protein structure.

Internal symmetry dataset

For evaluation purposes we used the manually curated dataset of 1,007 domains selected randomly from the set of *SCOP* superfamilies, introduced in our previous study [21]. A small number of classifications were updated to be more consistent with the new symmetry definitions, specially for the cases of open symmetry. The updated version of the internal symmetry dataset (v2.0), together with the reasons of the modified annotations, is available at <https://github.com/rcsb/symmetry-benchmark> and summarized in supplementary table S1. An important note is that in the evaluation of the previous version the open symmetry cases in the benchmark were part of the asymmetric (negative) set, while they are part of the symmetric (positive) set in this evaluation.

Results

Method evaluation

In our previous article we compared the performance of *CE-Symm* against other internal symmetry detection methods. In the new version, the detection performance has only been affected by the additional order detection and alignment refinement steps. The ROC curves of both versions are very similar, with a slight reduction of the false positives in the new one (supplementary figure S3). At the default TM-score threshold values for result significance, the false positive (FP) rate has decreased from 5.5% to 2.5%, while the true positive (TP) rate has been reduced from 81% to 76% in the benchmarking dataset. The bottleneck in symmetry detection continues to be finding a significant self-alignment, as we previously suggested.

The different methods for order detection perform similarly for closed symmetry cases in the benchmark (supplementary table S2). The simpler graph-component method performs worse than the others, but it is the only one that can be used for open symmetries, while the sequence-function detector performs better than the

rotation-angle method, particularly for difficult cases. The evaluation of the default order detection strategy of *CE-Symm* is shown in supplementary figure S4. The strategy achieves an overall 89% precision, although high order open symmetries are still an open challenge.

On average for symmetric entries in the benchmark where *CE-Symm* could find symmetry, the optimization step extended the repeat length by 43%, reduced the RMSD by 1.8% and increased the average TM-Score of the repeat alignment by 19.6%. Furthermore, using optimization an additional 23 cases (9% of the symmetric structures in the benchmark) were correctly identified as symmetric (209 with optimization, 186 without), which is a 12% improvement in symmetry detection. Because the highest scoring alignment of the simulation trajectory is taken as the result, optimization can only improve the initial alignment.

Sequence-structure analysis

This new version of *CE-Symm* is the first tool capable of presenting internal symmetry as a multiple structural alignment between repeats. This feature provides a direct correlation between sequence and structure and can be used for comparative and evolutionary analysis of a variety of protein folds and families.

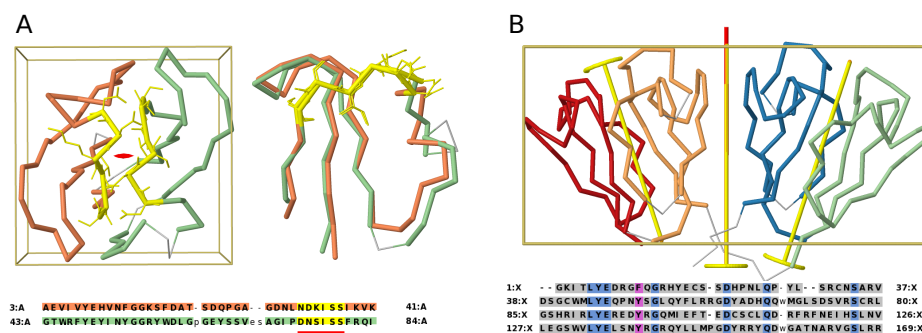


Fig 3. Internal symmetry in crystallin proteins A) An archaean $\beta\gamma$ -crystallin with two repeats per chain (3HZ2). The full chain is displayed along its 2-fold axis, followed by a superposition of the repeats. The conserved calcium binding motif N/D-N/D-#-I-S/T-S is highlighted in yellow throughout. B) Human γ -D crystallin structure with four repeats per chain (1HK0). Two levels of symmetry exist: a C2 symmetry within each domain, and an additional C2 axis relating the two domains. The calcium binding motif has been lost (red bar below sequence), but other conserved positions (blue and magenta in the sequence) show the homology between the repeats.

Structural alignment of the repeats can reveal conserved motifs that have persisted since the duplication event. One example is the $\beta\gamma$ -crystallin superfamily, which occurs in a variety of repeat arrangements. Many $\beta\gamma$ -crystallins contain a calcium binding site motif [31]. As shown in Fig. 3A, the calcium binding motif is structurally conserved after a 2-fold rotation around the symmetry axis, and the residue side-chains preserve their orientation. Furthermore, calcium coordinates residues from both repeats, making the two-fold symmetry an essential feature of the binding site.

On the other hand, duplication events allow the appearance of asymmetry by independent sequence and structural divergence of the repeats. An example is the MaoC-like thioesterase/thiol ester dehydrase-isomerase superfamily (SCOP: d.38.1.4). Members of this family fold into a characteristic 'hot dog fold' which binds coenzyme A and catalyzes the dehydration of various bound fatty acids. Typically the MaoC-like proteins contain one hot dog domain per chain and assemble into dimers, tetramers, or

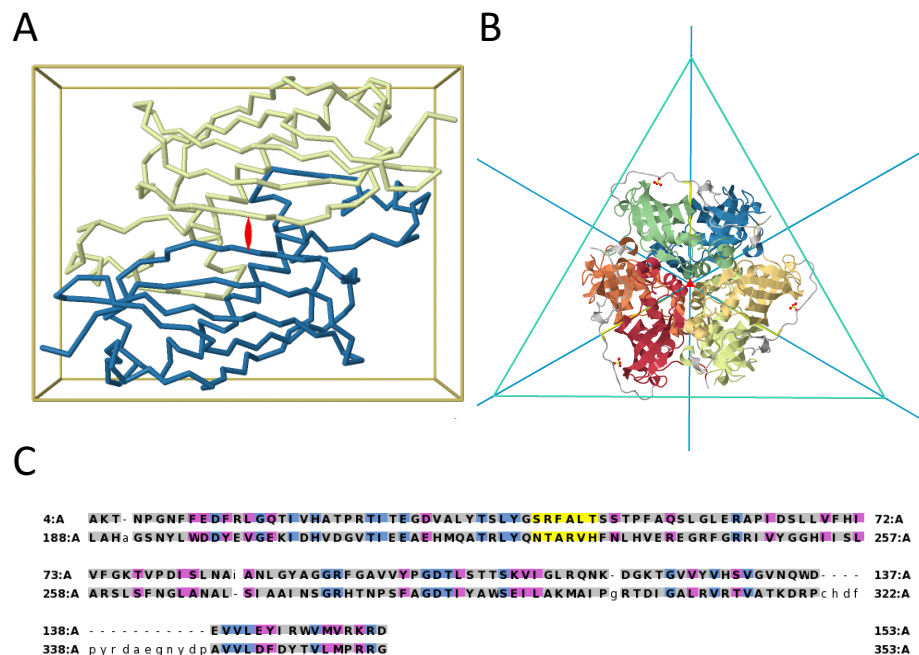


Fig 4. Hot dog fold duplication. A) Internal C2 symmetry in one chain of a MaoC domain protein dehydratase from *Chloroflexus aurantiacus* (4E3E) displaying a “double hot dog” fold. B) Full trimeric assembly, with the six individual hot dog domains colored. The quaternary structure has a threefold cyclic symmetry that combines with the twofold internal symmetry into a dihedral D3 symmetry equivalent to the quaternary structure of homologs without the internal domain duplication. C) Sequence alignment showing that the catalytic R/N-####-H motif (yellow) is lost in the first domain but retained in the second. Amino acid identity is shown in blue and similarity in magenta.

hexamers [32]. Some members of the family contain a duplication of the hot dog fold [33], accompanied by the loss of the catalytic motif R/N-####-H in one of the domains, in order to accommodate bulkier substrates which would otherwise not fit in a single domain [32]. The structural divergence of the catalytic site in one of the repeats of the double hot dog subunit can be easily observed with *CE-Symm* (Fig. 4).

Multiple levels of symmetry

Some proteins contain more than one axis of symmetry. In those cases, the axes of symmetry can be collinear, orthogonal or independent to each other. If the axes are collinear, they can be combined into a single axis with higher symmetry order. If the axes are orthogonal, they can be combined into a point group of higher symmetry order.

If the axes are independent to each other, multiple levels of symmetry exist in the structure in a hierarchical organization. This can be an indication of multiple independent duplication events, like in the case of γ -crystallins (Fig. 3B), where four repeats are related by two independent 2-fold axes corresponding to two successive duplication events.

Internal symmetry and assembly stoichiometry

Additionally, the internal symmetry axes can also combine with the quaternary symmetry axes. Therefore, internal symmetry can increase the order of symmetry of a protein complex. Returning to the previous MaoC-like protein example, the internal two-fold axis of the double hot dog domain in Fig. 4B is orthogonal to the three-fold quaternary symmetry axis, combining for an overall dihedral symmetry. This arrangement is structurally similar to the D3 quaternary symmetry of hexameric single hot dog proteins (e.g. 1YLI). Accounting for internal symmetry when comparing the symmetry two protein assemblies is therefore important, because proteins can have a similar overall structure despite their different subunit compositions. It would be misleading to say that the structure of the trimeric and hexameric MaoC-like proteins are substantially different. Another well-know example of similar overall arrangement with different subunit composition are DNA clamps, which promote processivity in DNA replication. In archaea and eukaryotes, the clamp is a trimer, while in bacteria it is a dimer [34]. Furthermore, all DNA clamps have further internal symmetry axes leading to an overall D6 symmetry. As a historical note, the homology between bacterial and eukaryotic DNA clamps was only acknowledged when the structures were solved and the similarity of their complexes was identified [35].

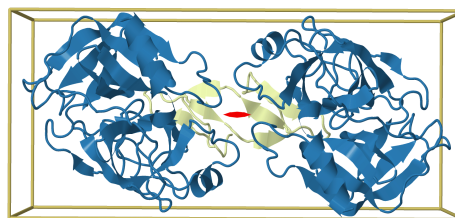


Fig 5. Uneven stoichiometry as a consequence of internal symmetry. The Bowman-Birk inhibitor from snail medic seeds (yellow) forms a complex with two bovine trypsin subunits (blue) in an uneven (A2B) stoichiometry and asymmetric assembly (2ILN). However, a 2-fold symmetry axis (red) can be identified in the complex when the internal symmetry of the inhibitor is taken into account, showing that the complex is equivalent to one with even (A2B2) stoichiometry.

Furthermore, internal symmetry is important in understanding the stoichiometry of protein assemblies. Uneven stoichiometry assemblies are those with an unbalanced number of each entity type in the complex and occur rarely in the biological environment. It was previously reported that up to 40% of all protein assemblies with uneven stoichiometry in the PDB can be explained by the presence of internal symmetry in one or multiple of the subunits in the complex [36]. One such example is the artificial complex of Bowman-Birk inhibitor from snail medic seeds with bovine trypsin, which has an A2B stoichiometry (Fig. 5). Although the complex is asymmetric, considering the internal symmetry of the inhibitor shows that the assembly is structurally comparable to an even A2B2 assembly with C2 overall symmetry. This property has also functional consequences, since the binding of two trypsin proteins symmetrically allows the inhibitor to efficiently induce dimerization and block the peptidase activity. Symmetry is characteristic of biological assemblies and can be considered by methods, like EPPIC, in order to predict the biological assembly in the context of crystal lattices [37]. Including internal symmetry in these methods could further improve predictions for some known cases like, for example, uneven stoichiometries.

Open Symmetry

The majority of proteins with internal symmetry have closed symmetry. In the case of quaternary symmetry, this is expected since homooligomers with open symmetry are disfavored due to their aggregation potential [38]. However, this is not the case for open internal symmetry due to the ability for terminal repeats to diverge to avoid undesirable homotypic interactions. Open symmetry is common in internal repeats.

The most general formulation of open repeats in the literature is that of superhelical symmetry, where the repeated subunit is simultaneously translated along a helical path (curvature) and rotated around this path (twist) [23]. *CE-Symm* cannot in principle identify superhelical symmetries, where both curvature and twist are relevant, because of the fundamental limit of the method to find a single symmetry axis (or multiple independent axes). However, we observed that the majority of structures containing tandem repeats that are classified as superhelical in the literature (solenoids) can be described by a single axis of symmetry. They fall in one of the following four conditions: i) the twist is negligible; ii) the curvature is negligible; iii) both twist and curvature are negligible; or iv) the twist is much larger than the curvature. In all those cases, *CE-Symm* can identify the symmetry in the structures and annotate them as helical, translational or open rotational symmetries.

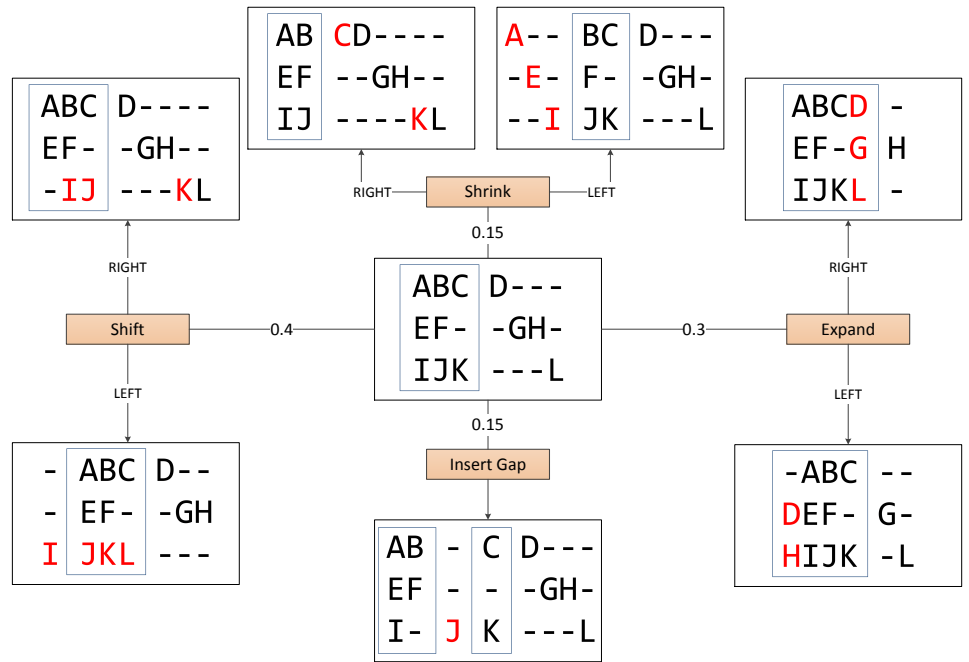
For instance, from the 18 solenoid protein representatives from table 1 in Kobe and Kajava [23], in 10 either the twist or the curvature are reported to be small (helical symmetry applies), in 5 both the twist and curvature are annotated as small (translational symmetry applies), and the remaining 3 structure representatives have irregular twist (asymmetric applies). Although many folds are classified as superhelical, only a small number have regular repeats but do not fit into one of the above categories. Therefore, in practice *CE-Symm* can also be a good method for identifying, classifying and characterizing solenoid and other repeat proteins with open symmetry. We hypothesize that the low prevalence of actual superhelical symmetry in proteins could be a consequence of the benefit in conserving interfaces between adjacent repeats.

Conclusion

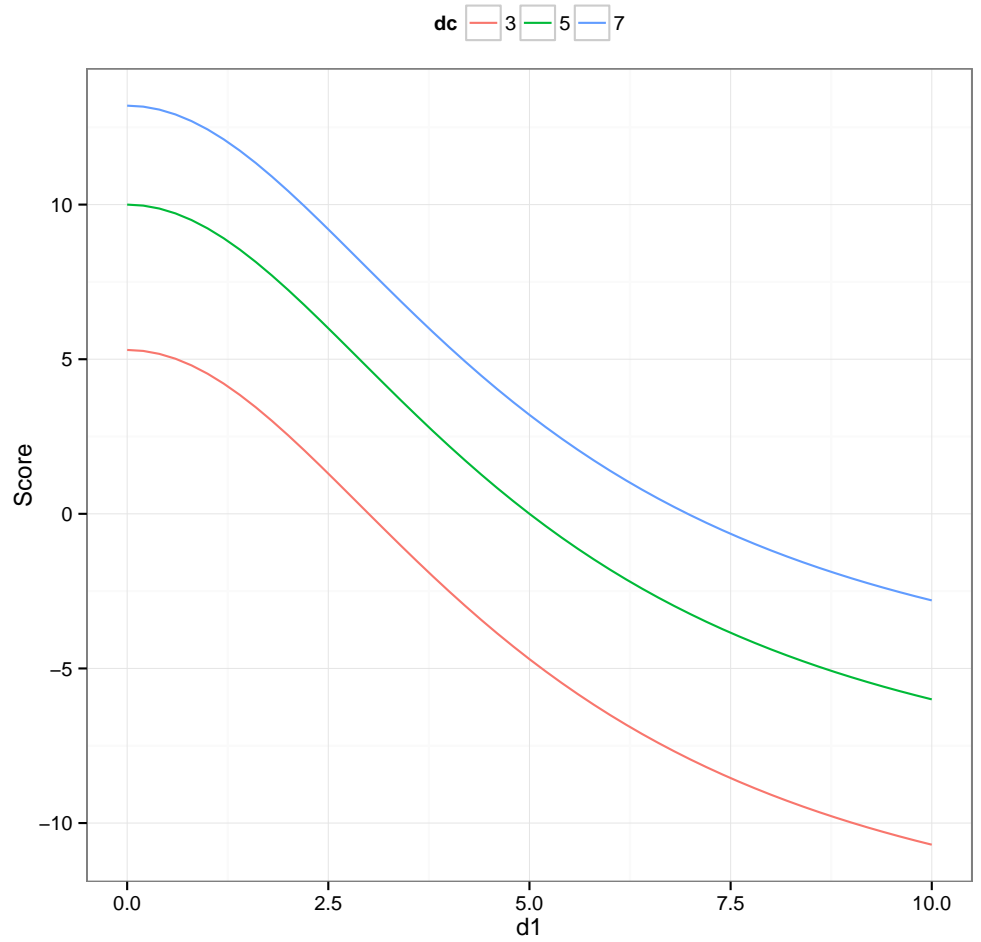
We have extended our internal symmetry analysis tool in order to improve its usability, capabilities and the interpretability of results. In addition to detecting symmetry in protein structures, the tool can identify corresponding residues of the protein from each repeating element and the symmetry operations between them. *CE-Symm 2.0* adds broad capabilities for the detection of all types of internal symmetry, providing information about the type and order of symmetry and the repeat boundaries. The alignments between the repeats are eminently useful in identifying conserved and differential features between repeats, and can be applied to understanding protein function and evolution.

Determining whether the high prevalence of internal symmetry in protein structures is predominantly a consequence of thermodynamic selection or an indication of the history of protein evolution remains an open question. Here, we have presented examples where internal symmetry is a result of evolution and tied to functional consequences, and how our tool can help researchers in the protein evolution, classification and annotation fields. We have made *CE-Symm 2.0* and its source code freely available on GitHub, and we are working to integrate it into leading bioinformatics resources for protein analysis.

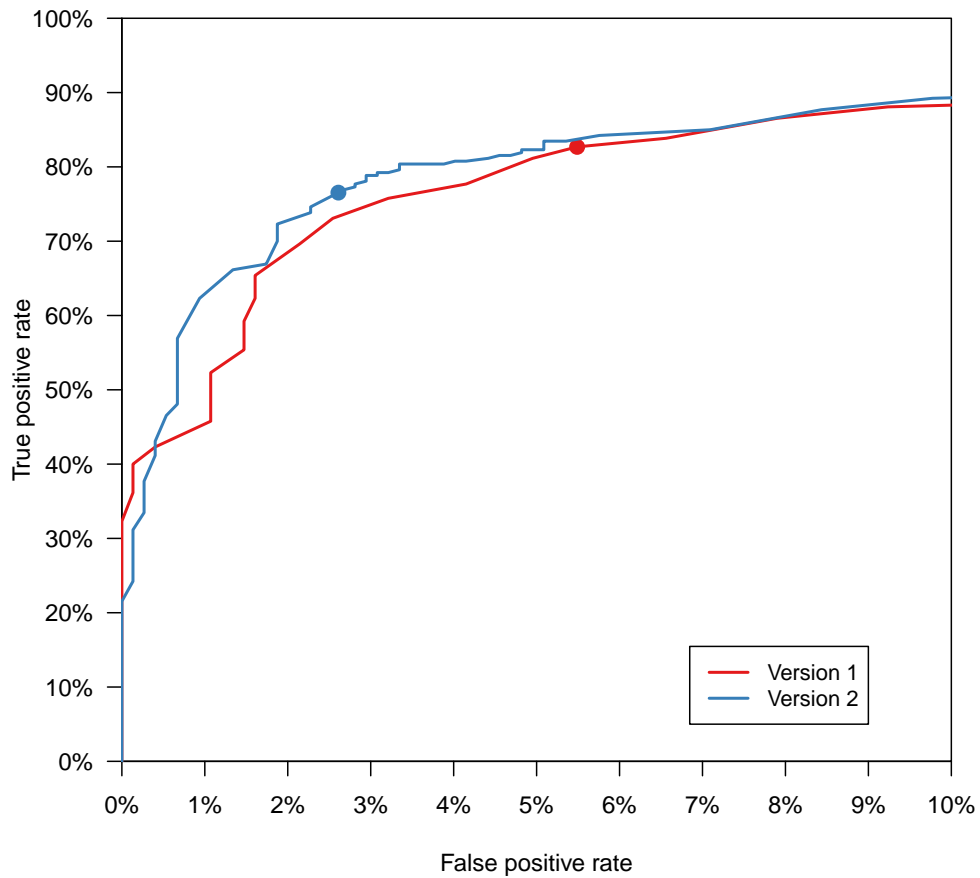
Supporting information



S1 Fig. Schematic representation of the Monte Carlo optimization moves. The starting alignment is shown in the center. The probability of each of the moves are indicated along the edges.



S2 Fig. Score function for the Monte Carlo optimization procedure.



S3 Fig. Comparison of the ROC curves of the symmetry detection for the old (Version 1) and new (Version 2) versions of *CE-Symm*. Differences in the ROC curves are not significant. The dots indicate the sensitivity and specificity at the default TM-score threshold (0.4).

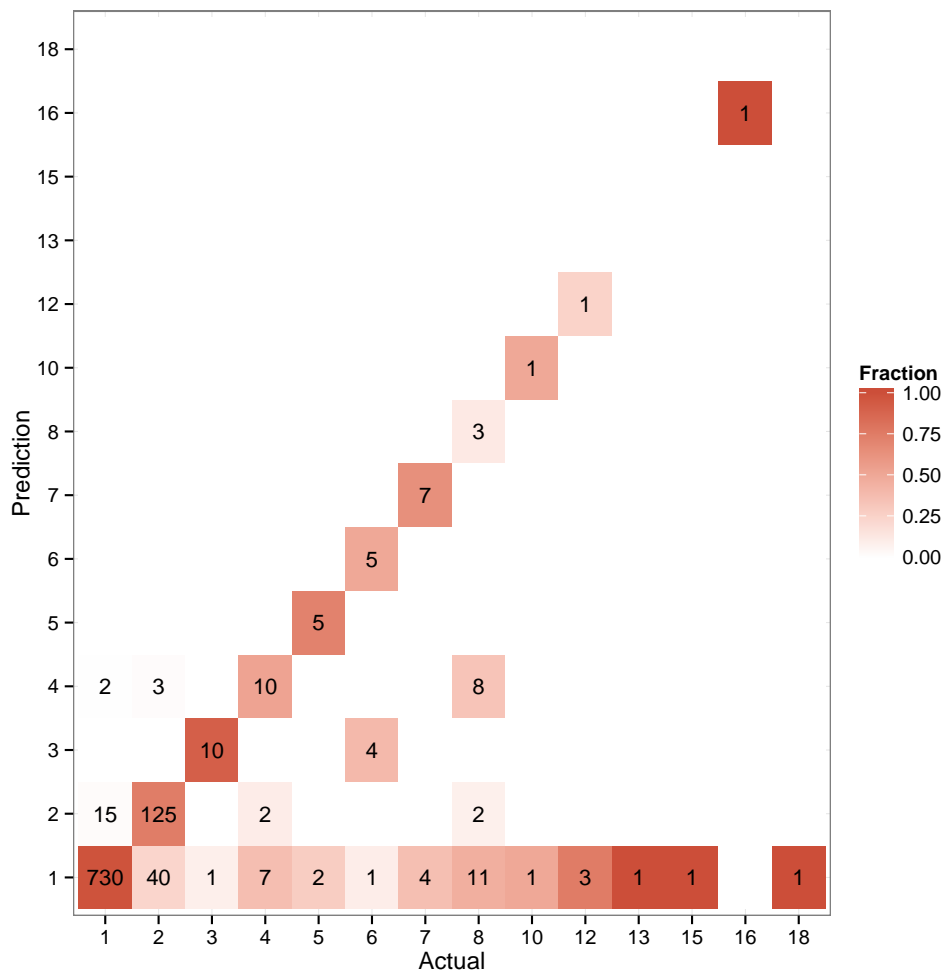
409

410

411

412

413



S4 Fig. Confusion matrix of actual and predicted symmetry orders of the structures in the benchmark. Entries of the matrix are colored by the recall of each symmetry order (columns).

414
415
416
417

Type	Count	Percentage
Asymmetric	747	74.2%
Rotational	214	21.2%
C2	160	74.8%
C3	10	4.7%
C4	2	0.9%
C5	3	1.4%
C6	9	4.2%
C7	10	4.7%
C8	20	9.3%
Dihedral	18	1.8%
D2	14	77.8%
D3	1	5.6%
D4	2	11.1%
D5	1	5.6%
Helical	11	1.1%
H2	9	81.8%
H3	2	18.2%
H10	1	9.1%
Superhelical	2	0.2%
Repeats	15	1.5%

S1 Tab. Summary of the updated annotations in the internal symmetry benchmarking dataset.

Method	Precision	Cramer V
Graph Component	0.598	0.652
Sequence Function	0.783	0.728
Rotation Angle	0.754	0.642

S2 Tab. Performance measures of the symmetry order detection methods for closed symmetry domains in the benchmark dataset. Precision measures the total fraction of correct predictions and Cramer V measures the correlation between actual and predicted classes. Both measures have values in the [0,1] interval, where 1 means perfect precision and correlation.

Competing interests

The authors declare that they have no competing interests.

Author's contributions

P.B. initiated the CE-Symm project and provided direction. S.B., A.L, P.R. and A.P. developed code for the *CE-Symm* tool. S.B. and A.L. designed, implemented and benchmarked the algorithmic extensions of the method and performed the analyses presented in this article. S.B., A.L. and G.C. wrote the manuscript. All authors reviewed and approved the final version of the manuscript.

Acknowledgements

We dedicate this article to Guido, catalyst and indispensable part of this project, who sadly left us before its completion.

We thank Philippe Youkharibache for testing and suggesting new features for *CE-Symm* and Jose Duarte for useful discussions during the development of the method. We also thank the developers that contributed to the *BioJava* library, which has been fundamental throughout the development of *CE-Symm*.

This research was supported in part by the Intramural Research Program of the National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health (support to S.B and P.B), and by the National Science Foundation, National Institutes of Health, and U.S. Department of Energy [DBI-1338415] (support to A.P., P.R., and P.B.). Financial support to G.C. from the Swiss National Science Foundation and the Research Committee of the Paul Scherrer Institute is gratefully acknowledged. Additional support for S.B from COST Action BM1405 and COST Switzerland SEFRI project IZCNZ0-174836.

References

1. Jacob F. Evolution and tinkering. *Science*. 1977;196(4295):1161–1166. doi:10.1126/science.860134.
2. Lupas AN, Ponting CP, Russell RB. On the evolution of protein folds: are similar motifs in different protein folds the result of convergence, insertion, or relics of an ancient peptide world? *Journal of structural biology*. 2001;134(2-3):191–203. doi:10.1006/jsbi.2001.4393.
3. Han JH, Batey S, Nickson AA, Teichmann SA, Clarke J. The folding and evolution of multidomain proteins. *Nature Reviews Molecular Cell Biology*. 2007;8(4):319–330. doi:10.1038/nrm2144.
4. Levy ED, Teichmann S. Structural, evolutionary, and assembly principles of protein oligomerization. *Progress in Molecular Biology and Translational Science*. 2013;117:25–51. doi:10.1016/B978-0-12-386931-9.00002-7.
5. Andrade MA, Perez-Iratxeta C, Ponting CP. Protein repeats: structures, functions, and evolution. *Journal of Structural Biology*. 2001;134(2-3):117–131. doi:10.1006/jsbi.2001.4392.
6. Guerler A, Wang C, Knapp EW. Symmetric structures in the universe of protein folds. *Journal of Chemical Information and Modeling*. 2009;49(9):2147–2151. doi:10.1021/ci900185z.
7. Goodsell DS, Olson AJ. Structural Symmetry and Protein Function. *Annu Rev Biophys Biomol Struct*. 2000;29(1):105–53. doi:10.1146/annurev.biophys.29.1.105.
8. Abraham AL, Pothier J, Rocha EPC. Alternative to Homo-oligomerisation: The Creation of Local Symmetry in Proteins by Internal Amplification. *Journal of Molecular Biology*. 2009;394(3):522–534. doi:10.1016/j.jmb.2009.09.031.
9. Lee J, Blaber M. Experimental support for the evolution of symmetric protein architecture from a simple peptide motif. *Proceedings of the National Academy of Sciences of the United States of America*. 2011;108:126–130. doi:10.1073/pnas.1015032108.
10. Broom A, Doxey AC, Lobsanov YD, Berthin LG, Rose DR, Howell PL, et al. Modular evolution and the origins of symmetry: Reconstruction of a three-fold symmetric globular protein. *Structure*. 2012;20(1):161–171. doi:10.1016/j.str.2011.10.021.

11. Wolynes PG. Symmetry and the energy landscapes of biomolecules. *Proceedings of the National Academy of Sciences of the United States of America*. 1996;93(25):14249. doi:10.1073/pnas.93.25.14249.
12. Bonjack-Shterengartz M, Avnir D. The near-symmetry of proteins. *Proteins: Structure, Function and Bioinformatics*. 2015;83(4):722–734. doi:10.1002/prot.24706.
13. Kinoshita K, Kidera a, Go N. Diversity of functions of proteins with internal symmetry in spatial arrangement of secondary structural elements. *Protein science : a publication of the Protein Society*. 1999;8:1210–1217. doi:10.1110/ps.8.6.1210.
14. Murray KB, Taylor WR, Thornton JM. Toward the detection and validation of repeats in protein structure. *Proteins: Structure, Function and Genetics*. 2004;57(2):365–380. doi:10.1002/prot.20202.
15. Shih ESC, Gan RCR, Hwang MJ. OPAAS: A web server for optimal, permuted, and other alternative alignments of protein structures. *Nucleic Acids Research*. 2006;34(WEB. SERV. ISS.). doi:10.1093/nar/gkl264.
16. Abraham AL, Rocha EPC, Pothier J. Swelpe: A detector of internal repeats in sequences and structures. *Bioinformatics*. 2008;24(13):1536–1537. doi:10.1093/bioinformatics/btn234.
17. Marsella L, Sirocco F, Trovato A, Seno F, Tosatto SCE. REPETITA: detection and discrimination of the periodicity of protein solenoid repeats by discrete Fourier transform. *Bioinformatics*. 2009;25(12):i289–i295.
18. Kim C, Basner J, Lee B. Detecting internally symmetric protein structures. *BMC bioinformatics*. 2010;11:303. doi:10.1186/1471-2105-11-303.
19. Do Viet P, Roche DB, Kajava AV. TAPO: A combined method for the identification of tandem repeats in protein structures. *FEBS Letters*. 2015;589(19):2611–2619. doi:10.1016/j.febslet.2015.08.025.
20. Paladin L, Hirsh L, Piovesan D, Andrade-Navarro MA, Kajava AV, Tosatto SCE. RepeatsDB 2.0: improved annotation, classification, search and visualization of repeat protein structures. *Nucleic Acids Res*. 2017;45(6):3613–3613.
21. Myers-Turnbull D, Bliven SE, Rose PW, Aziz ZK, Youkharibache P, Bourne PE, et al. Systematic detection of internal symmetry in proteins using CE-symm. *Journal of Molecular Biology*. 2014;426(11):2255–2268. doi:10.1016/j.jmb.2014.03.010.
22. Murzin AG, Brenner SE, Hubbard T, Chothia C. SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J Mol Biol*. 1995;247(4):536–540.
23. Kobe B, Kajava AV. When protein folding is simplified to protein coiling: The continuum of solenoid protein structures. *Trends in Biochemical Sciences*. 2000;25(10):509–515. doi:10.1016/S0968-0004(00)01667-4.
24. Shindyalov IN, Bourne PE. Protein structure alignment by incremental combinatorial extension (CE) of the optimal path. *Protein Engineering Design and Selection*. 1998;11(9):739–747. doi:10.1093/protein/11.9.739.

25. Bliven SE, Bourne PE, Prlić A. Detection of circular permutations within protein structures using CE-CP. *Bioinformatics*. 2015;31(8):1316–1318. doi:10.1093/bioinformatics/btu823.
26. Zhang Y, Skolnick J. Scoring function for automated assessment of protein structure template quality. *Proteins: Structure, Function and Genetics*. 2004;57(4):702–710. doi:10.1002/prot.20264.
27. Bliven SE. *Structure-Preserving Rearrangements: Algorithms for Structural Comparison and Protein Analysis*. University of California San Diego; 2015. Available from: <https://escholarship.org/uc/item/0t54p4gj>.
28. Guda C, Scheeff ED, Bourne PE, Shindyalov IN. A New Algorithm for the Alignment of Multiple Protein Structures Using Monte Carlo Optimization. *Pacific Symposium on biocomputing*. 2001;6:275–286.
29. Levy ED, Pereira-Leal JB, Chothia C, Teichmann SA. 3D complex: A structural classification of protein complexes. *PLoS Computational Biology*. 2006;2(11):1395–1406. doi:10.1371/journal.pcbi.0020155.
30. Karney CFF. Quaternions in molecular modeling. *Journal of Molecular Graphics and Modelling*. 2007;25(5):595–604. doi:10.1016/j.jmkgm.2006.04.002.
31. Aravind P, Mishra A, Suman SK, Jobby MK, Sankaranarayanan R, Sharma Y. The gamma-crystallin superfamily contains a universal motif for binding calcium. *Biochemistry*. 2009;48(51):12180–12190. doi:10.1021/bi9017076.
32. Pidugu LS, Maity K, Ramaswamy K, Suroliya N, Suguna K. Analysis of proteins with the 'hot dog' fold: Prediction of function and identification of catalytic residues of hypothetical proteins. *BMC Structural Biology*. 2009;9. doi:10.1186/1472-6807-9-37.
33. Qin YM, Haapalainen AM, Kilpeläinen SH, Marttila MS, Koski MK, Glumoff T, et al. Human peroxisomal multifunctional enzyme type 2. Site-directed mutagenesis studies show the importance of two protic residues for 2-enoyl-CoA hydratase 2 activity. *Journal of Biological Chemistry*. 2000;275(7):4965–4972.
34. Kelman Z, O'Donnell M. Structural and functional similarities of prokaryotic and eukaryotic DNA polymerase sliding clamps. *Nucleic Acids Res*. 1995;23(18):3613–3620.
35. Leipe DD, Aravind L, Koonin EV. Did DNA replication evolve twice independently? *Nucleic Acids Res*. 1999;27(17):3389–3401.
36. Marsh Ja, Rees Ha, Ahnert SE, Teichmann Sa. Structural and evolutionary versatility in protein complexes with uneven stoichiometry. *Nature communications*. 2015;6:6394. doi:10.1038/ncomms7394.
37. Bliven SE, Lafita A, Parker A, Capitani G, Duarte JM. Automated evaluation of quaternary structures from protein crystals. *bioRxiv*. 2017; p. 224717.
38. Capitani G, Duarte JM, Baskaran K, Bliven S, Somody JC. Understanding the fabric of protein crystals: Computational classification of biological interfaces and crystal contacts. *Bioinformatics*. 2015;32(4):481–489. doi:10.1093/bioinformatics/btv622.