

This item is the archived peer-reviewed author-version of:

Nanoscale insight into silk-like protein self-assembly : effect of design and number of repeat units

Reference:

Razzokov Jamoliddin, Naderi Saber, van der Schoot Paul.- Nanoscale insight into silk-like protein self-assembly : effect of design and number of repeat units
Physical biology - ISSN 1478-3967 - 15:6(2018), 066010
Full text (Publisher's DOI): <https://doi.org/10.1088/1478-3975/AADB5E>
To cite this reference: <https://hdl.handle.net/10067/1535960151162165141>

ACCEPTED MANUSCRIPT

Nanoscale insight into silk-like protein self-assembly: effect of design and number of repeat units

To cite this article before publication: Jamoliddin Razzokov *et al* 2018 *Phys. Biol.* in press <https://doi.org/10.1088/1478-3975/aadb5e>

Manuscript version: Accepted Manuscript

Accepted Manuscript is “the version of the article accepted for publication including all changes made as a result of the peer review process, and which may also include the addition to the article by IOP Publishing of a header, an article ID, a cover sheet and/or an ‘Accepted Manuscript’ watermark, but excluding any other editing, typesetting or other changes made by IOP Publishing and/or its licensors”

This Accepted Manuscript is © 2018 IOP Publishing Ltd.

During the embargo period (the 12 month period from the publication of the Version of Record of this article), the Accepted Manuscript is fully protected by copyright and cannot be reused or reposted elsewhere.

As the Version of Record of this article is going to be / has been published on a subscription basis, this Accepted Manuscript is available for reuse under a CC BY-NC-ND 3.0 licence after the 12 month embargo period.

After the embargo period, everyone is permitted to use copy and redistribute this article for non-commercial purposes only, provided that they adhere to all the terms of the licence <https://creativecommons.org/licenses/by-nc-nd/3.0>

Although reasonable endeavours have been taken to obtain all necessary permissions from third parties to include their copyrighted content within this article, their full citation and copyright line may not be present in this Accepted Manuscript version. Before using any content from this article, please refer to the Version of Record on IOPscience once published for full citation and copyright details, as permissions will likely be required. All third party content is fully copyright protected, unless specifically stated otherwise in the figure caption in the Version of Record.

View the [article online](#) for updates and enhancements.

1
2
3 **Nanoscale insight into silk-like protein self-assembly: effect of design and number of**
4 **repeat units.**
5

6 Jamoliddin Razzokov,^{1,2} Saber Naderi,^{2,3} and Paul van der Schoot^{2,4}

7
8 ¹*Department of Chemistry, University of Antwerp, Universiteitsplein 1,*
9 *2610 Antwerp, Belgium*

10
11 ²*Faculteit Technische Natuurkunde, Technische Universiteit Eindhoven,*
12 *Postbus 513, 5600 MB Eindhoven, The Netherlands^{a)}*

13
14 ³*Dutch Polymer Institute, P.O. Box 902, 5600 AX Eindhoven,*
15 *The Netherlands*

16
17 ⁴*Instituut voor Theoretische Fysica, Universiteit Utrecht, Leuvenlaan 4,*
18 *3584 CE Utrecht, The Netherlands*

19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Accepted Manuscript

1
2
3 By means of replica exchange molecular dynamics simulations we investigate how
4 the length of a silk-like, alternating diblock oligopeptide influences its secondary and
5 quaternary structure. We carry out simulations for two protein sizes consisting of
6 three and five blocks, and study the stability of a single protein, a dimer, a trimer and
7 a tetramer. Initial configurations of our simulations are β -roll and β -sheet structures.
8 We find that for the triblock the secondary and quaternary structures upto and
9 including the tetramer are unstable: the proteins melt into random coil structures
10 and the aggregates disassemble either completely or partially. We attribute this to
11 the competition between conformational entropy of the proteins and the formation
12 of hydrogen bonds and hydrophobic interactions between proteins. This is confirmed
13 by our simulations on the pentablock proteins, where we find that, as the number
14 of monomers in the aggregate increases, individual monomers form more hydrogen
15 bonds whereas their solvent accessible surface area decreases. For the pentablock β -
16 sheet protein, the monomer and the dimer melt as well, although for the β -roll protein
17 only the monomer melts. For both trimers and tetramers remain stable. Apparently,
18 for these the entropy loss of forming β -rolls and β -sheets is compensated for, the
19 free-energy gain due to the hydrogen-bonding and hydrophobic interactions. We also
20 find that the middle monomers in the trimers and tetramers are conformationally
21 much more stable than the ones on the top and the bottom. Interestingly, the
22 latter are more stable on the tetramer than on the trimer, suggesting that as the
23 number of monomers increases protein-protein interactions cooperatively stabilize
24 the assembly. According to our simulations, the β -roll and β -sheet aggregates must
25 be approximately equally stable.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

43
44 Keywords: silk-like protein, implicit solvent, replica exchange molecular dynamics,
45 self-assembly, fiber.
46
47
48
49
50
51
52
53
54
55
56

57 ^{a)}Electronic mail: jamoliddin.razzokov@uantwerpen.be
58
59
60

I INTRODUCTION

A promising route to the rational design of functional nanomaterials is a bottom-up approach that makes use of the self-assembly of molecular building blocks and exploits the physical principles that govern inter- and intra-molecular interactions.¹⁻³ This is particularly true for biomimetic materials based on proteins and peptides as basic molecular units⁴⁻⁷ albeit that this does require an understanding of the relation between their molecular structure and self-organization properties. Recent advances in experimental, theoretical and simulation techniques have drawn attention to the role played by factors such as the protein concentration, hydrophobic, hydrophilic and electrostatic interactions, the salinity, temperature, pressure, pH and so on.⁸⁻¹¹ By controlling the primary sequence one can in principle control the way that the designer protein molecule responds to these factors and control its structure and function.¹²⁻¹⁶

An application of interest is gene therapy, relevant in the context of treatments of cancer, HIV (human immunodeficiency virus) and SCID (severe combined immunodeficiency), where a therapeutic gene is delivered into target cells.^{17,18} In order to successfully deliver the therapeutic gene into the target cells, a protective shell is required that shields the genetic material from degradation and attack by nucleases.^{19,20} One way to do this is to use a virus as a delivery vehicle but this apparently has disadvantages, toxic immunological responses.^{21,22} An alternative approach is to design bio-compatible proteins that are able to self-assemble on the genetic material and form an artificial virus-like particle.^{23,25}

Very recently a possible candidate for such an artificial DNA coat protein was proposed by de Vries and coworkers.²⁵ The protein is a triblock copolymer made up of a cationic DNA-binding motif, similar to the nucleic-acid binding domains in viruses, a collagen-like hydrophilic random coil motif that provides colloidal stability in aqueous solution and a silk-like protein block.

The design is based on the hypothesis that the latter self-assembles into β -sheets (at least in the aggregated state), that is, stack on top of each other via hydrogen bonding and hydrophobic interactions, and stabilize the protein aggregate around the DNA.

The work of de Vries et al., that the triblock copolymer is indeed capable of self-organizing into a capsid around the DNA, protect it against the action of nucleases and successfully transfect cells.²⁵ One of the important design parameters is the length of the silk-like block

1
2
3 that consists of, copies of a repeat unit. Below a critical length the sequence does not
4 seem to form a β -sheet, self-assemble and provide sufficiently dense coverage on the DNA
5 to protect it against attack by the nucleases. However, if the number of repeat units is too
6 large the binding of the protein to the DNA may well become too strong, interfering with
7 the expression of the target DNA in the cell.²⁵
8
9

10
11 To shed light on the influence of the number of repeat units in a silk-like protein similar
12 to that of the core of the protein of de Vries and collaborators, we carry out all-atom replica
13 exchange molecular dynamics simulations on proteins consisting of three and five repeat
14 units that we specify below. We restrict ourselves in this work to aggregates consisting of
15 one to four proteins in free solution and investigate the stability of the folded structure as
16 well as that of the aggregates themselves. In follow-up work we extend this to simulation
17 studies in the presence of DNA, which are computationally much more demanding.
18
19

20
21 Before presenting our simulations, we note that the influence of the length of β -sheet form-
22 ing silk-like proteins on their assembly into fibers has also been studied experimentally.^{12,24,26}
23
24

25
26 For instance, Davies and coworkers studied the self-assembly of peptides containing 7 and
27 9 amino acids with sequences $\text{CH}_3\text{CO-RLQLQLE-NH}_2$ and $\text{CH}_3\text{CO-QRLQLQLEQ-NH}_2$.¹²
28 They observe that above a critical concentration the 7mers and 9mers, which have a random
29 coil structure in the monomeric state, self-assemble into inter-molecular β -sheet tapes. The
30 critical concentration was lower for the case of the 9mers, indicating that the binding energy
31 for the 9mer must be larger than that of the 7mer. This is plausible because more hydrogen
32 bonds form between the longer peptides compared to the shorter ones.
33
34

35
36 From these experiments it is not possible to say whether dimers and trimers are also part
37 of the full fiber length distribution that might be dominated by long assemblies. The reason
38 is that the authors focus on circular dichroism spectroscopy that allows one to probe the
39 assembled fraction of proteins and not the size distribution. Whilst informative, measure-
40 ments like this cannot provide us with an insight at fully atomistic detail of the molecular
41 processes at the root of the stability of the very long assemblies. Computer simulations,
42 however, can provide such atomistic detail. For example, Schor and collaborators studied
43 by means of replica exchange molecular dynamics (REMD) simulations in explicit solvent
44 the most stable structure of a fiber forming silk-like protein, containing six repeat units of an
45 amino acid sequence $[(\text{GA})_3\text{-GE}]$ in aqueous solution. The authors performed simulations on
46 a single protein starting from β -sheet or β -roll structures as well as a dimer containing two
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Number of (GA) ₃ G blocks	Primary sequence
$n = 3$	K(GA) ₃ GQ-(GA) ₃ GK(GA) ₃ GQ
$n = 5$	[K(GA) ₃ GQ-(GA) ₃ G] ₂ K(GA) ₃ GQ

TABLE I. The primary sequence of our triblock ($n = 3$) and pentablock ($n = 5$) proteins.

β -sheets or β -rolls that are stacked on top of each other to find which of the two structures is more stable. They find that in both monomeric and dimeric states the β -roll structure is the most stable structure.³⁷

Recently, we repeated the simulations of Schor and coworkers for a monomer and a dimer and in addition to this we carried out simulations with a trimer and a tetramer of their silk-like protein in both explicit and implicit solvent.²⁷ For the same amount of simulation time we obtain the same result as Schor et al. but by extending our simulations to longer time scales we find that in monomeric and dimeric states both the β -sheet and β -roll structures melt into random coil structures. This shows that more frequent sampling of the fewer replicas required for simulations with implicit solvent and probing longer time scales that can then be achieved in REMD simulations are crucial for sufficient sampling of the phase space of our protein. Here we present results of replica exchange computer simulations focusing on the self-assembly of a slightly different silk-like protein consisting of two types of repeat unit K[(GA)₃-G] and Q[(GA)₃-G]. See Table I for the primary sequence of our protein. In our simulations we make use of an implicit solvent model and consider two values of the number of repeat units, $n = 3$ and 5. We carry out simulations starting from one, two, three and four β -sheet or β -roll structures stacked on top of each other.

We find that for the case of the proteins with $n = 5$ repeat units, the monomer and dimer β -sheets melt into a disordered globular state whereas the dimer β -roll and trimers and tetramers of β -sheets and β -rolls are stable. In contrast, the proteins with $n = 3$ repeat units melt into random coil structures and disassemble within the simulation time. This indicates that the protein-protein binding free energy increases with increasing the number of repeat units. This is arguably caused by stronger inter-molecular interactions due to (i) a larger number of hydrogen bonds in the aggregates of the larger proteins and (ii) the fact that the solvent accessible surface area per monomer decreases more rapidly with aggregate size for the proteins with $n = 5$ than those with $n = 3$.

1
2
3 The remainder of this paper is organized as follows. In Section II, we present the methods
4 that we use in our simulations and the way we analyze our data. We discuss obtained results
5 in Section III. Section III A highlights our simulation results for the proteins with five repeat
6 units. These results are compared with our findings on the triblock proteins in Section III B.
7
8 Finally, Section IV is devoted to conclusions.
9

10 11 12 II SIMULATION METHOD AND ANALYSIS

13
14 We perform all-atom molecular dynamics simulations of a single protein, a dimer, a trimer
15 and a tetramer of a silk-like protein with an alternating amino acid sequence $K[(GA)_3-GX]_n$,
16 where X stands for glutamine and lysine, see Table I. We make use of the Amber simulation
17 package²⁸ and employ the ff99SB force field²⁹, in combination with the generalized Born (GB)
18 implicit solvent model³⁰ which represents solvent as a continuous medium without artifacts
19 from periodicity. In these simulations solvation energies become part of the total energy
20 of the system and the forces driving dynamics contain the derivatives of this term. The
21 GB is widely applied implicit-solvent model in atomic scale simulations, approximates long
22 range electrostatic interactions based on an analytical formula given in literatures³¹⁻³³. The
23 average energetic contribution of solute-solvent hydrogen bond is also included in implicit
24 solvent model.^{30,34} This also serves to enhance sampling and rapid convergence due to the
25 absence of solvent frictional forces which is not required re-organization of explicit water
26 molecules in response to protein conformational changes. The effective salt concentration
27 in our simulations equals 0.1 M, which is implemented in the simulations by including a
28 Debye-Hückel-type term³⁵ in the calculation of the electrostatic interaction energy of the
29 implicit solvent model.
30
31
32
33
34
35
36
37
38
39
40

41
42 By choosing an implicit solvent model instead of an explicit one we reduce the com-
43 putational cost of our simulations. This allows us to reach longer timescales compared to
44 simulations with explicit solvent because (i) we need less computing time for each replica
45 and (ii) we need fewer replicas. In a recent study on a similar protein with an amino acid
46 sequence of $[(GA)_3-GE]_5$ we have shown that (i) for this protein and for the same amount of
47 simulation time the results of implicit solvent simulations agree well with those of explicit
48 solvent ones and (ii) by extending the implicit solvent simulations to longer timescales we
49 have found that to study the self-assembly of our protein a sufficiently long simulation is
50 required.²⁷
51
52
53
54
55
56

57 We run simulations with $n = 3$ and 5 repeat units starting from β -sheet and β -roll
58
59
60

structures that are created using the Xleap program in AmberTools. The difference between β -sheet and β -roll structures is that in the β -sheet structure each repeat unit forms hydrogen bonds with its nearest neighbor in the protein sequence, whereas in β -rolls the hydrogen bonds are mostly formed between each repeat unit and the next-nearest neighbor. In order to create these structures we make linear protein chains using the Xleap program and pull the nearest (for the β -sheet) or the next-nearest (for the β -rolls) repeat units together during a short simulation run using the Amber simulation package. Cartoon view of monomers of a β -sheet and a β -roll structures is given in Fig. 1. Schematic images of proteins presented

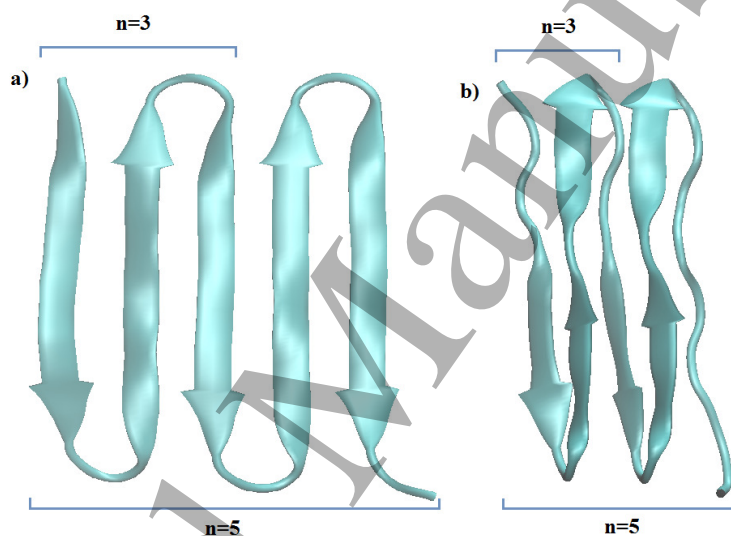


FIG. 1. Schematic representation of (a) a β -sheet and (b) a β -roll structure of our silk-like protein. Here each monomer contains $n = 5$ repeat units. $n = 3$ repeat units of protein is created taking 3 strands of each monomer.

in this paper are created using VMD software.³⁶

All simulations are performed using Langevin dynamics in implicit solvent, mimicking conditions where the lysine residues and the proteins as a whole are either charge neutral or the ionic strength is sufficiently large to effectively screen Coulomb interactions between charged groups. Partial positive and negative charges are invoked to model electro-negative and -positive atoms. Multimers of β -sheets and β -rolls are created by copying and translating monomers in the direction perpendicular to the β -sheet or β -roll plane over an arbitrary distance of 0.6 nm. From our simulations we find that for this choice individual proteins are able to find each other and bind within a reasonable simulation time. Note that our

1
2
3 procedure conserves the broken up-down symmetry of the proteins and maximizes the con-
4 tact area between them. This is necessary because the proteins in cross section are not
5 symmetric. All initial structures are energy minimized prior to our simulations.
6
7

8 For each of the protein aggregates we run replica exchange molecular dynamics (REMD)
9 simulations in the temperature range of 298 to 554 K. This temperature range we distribute
10 over 10 replicas for the monomers and dimers of the aggregates consisting of the protein
11 with $n = 3$ repeat units and for the monomers of proteins with five repeat units ($n = 5$).
12 For the trimers and tetramers of the shorter proteins with $n = 3$ and for the dimers, trimers
13 and tetramers of proteins with $n = 5$ the temperature range is somewhat smaller, 298 to
14 511 K, and spread over 12 replicas.
15
16
17
18
19
20

21 The temperature distribution is chosen in such a way that the replica exchange proba-
22 bilities are between 10 to 30 percent, as is customary.³⁷ A replica exchange attempt is done
23 every 500 MD steps, with each step representing 2 fs of real time. The exchange probability
24 is determined by a Metropolis algorithm, where the weight factors are given by the product
25 of the Boltzmann factors for each replica. The overall REMD simulation time is different
26 depending on the size of the aggregates, but varies between 20 ns and 40 ns per replica.
27 In our simulation the average time that it takes for a replica at the lowest temperature to
28 diffuse up in the temperature space, reach the highest temperature and diffuse back to the
29 lowest temperature is about 1 to 3 ns. This shows that the overall REMD simulation time
30 is sufficient for our replicas to explore the temperature space. Besides, substantial speedup
31 conformational changes is observed by using implicit solvent models in literatures.^{31,38}
32
33
34
35
36
37
38
39

40 To analyze our simulation data we calculate the number of hydrogen bonds within each
41 protein aggregate, the root mean-square deviation (RMSD) of each aggregate from its initial
42 structure and the solvent accessible surface area (SA). The number of hydrogen bonds in
43 each simulation frame is determined by measuring the distance between H-bond donors
44 and acceptors, using a cutoff of 0.3 nm, and the angle between the acceptor, hydrogen
45 and donor atoms for which the cutoff is 135° . For calculating the RMSD, we take the
46 initial configuration of each simulation as a reference structure, which consists of a single
47 or multiple β -sheets or β -rolls depending on the system we study, and compute the RMSD
48 between all atoms of each frame to the atoms in the reference configuration. To probe the
49 solvent accessible surface area (SA) of our protein systems, a “rolling-ball” algorithm is used
50 in which a spherical probe of radius of 0.14 nm moves on the surface of the protein.³⁹
51
52
53
54
55
56
57
58
59
60

1
2
3 In REMD simulations, trajectories are in a way not continuous in time, because they are
4 exchanged between replicas at given intervals with a certain probability. Therefore, instead
5 of calculating the number of hydrogen bonds, the RMSD and the SA as a function of time,
6 we obtain the free energy of our protein systems, ΔG , as a function of each of the quantities
7 mentioned above. To do this, first, from the data of the lowest temperature replica, we
8 count the number of occurrences in which our system is in state X , which corresponds to a
9 value of, e.g., the RMSD. Next, the probability, $\Pi(X)$, of finding our system at a state X is
10 computed from the number of occurrences. Finally, the free-energy is calculated using the
11 Boltzmann relation $\Pi(X) \propto \exp -\Delta G(X)/k_B T$, where k_B is the Boltzmann constant and
12 T is the temperature. The global minimum of the free energy we arbitrarily set to zero as
13 we can only probe free energy differences within a single simulation run.
14
15
16
17
18
19
20
21

22 In the following, we shall first discuss our results for the protein consisting of five repeat
23 units, and after that contrast that to what we find for the triblock protein.
24
25

26 III RESULTS AND DISCUSSION

27 III A. Cooperative stabilization of folded pentablock proteins

28 Focusing first on a discussion of our simulation results on the assemblies of our five-repeat
29 unit protein, we remind the reader that for each of the assemblies the initial configurations
30 start off from all β -sheet and all β -roll structures. Shown in Fig. 2 are representative
31 snapshots of a monomer, a dimer, a trimer and a tetramer taken from the simulations with
32 initial β -sheet structures, and similar ones for the initial β -roll structures in Fig. 9. Note
33 that these snapshots are not necessarily the final structures of the simulations, they represent
34 the most stable ones.
35
36
37
38
39
40
41

42 In our simulations, the monomers and the dimers of the β -sheets melt into a random coil
43 whereas their trimeric and tetrameric states are much more stable. For the β -rolls only the
44 monomers melt within our simulation time, and dimers, trimers and tetramers are stable
45 albeit one of the outer monomers of the tetramer melts. To quantify these observations we
46 measure the root mean-square deviation (RMSD) of our proteins using snapshots that we
47 take every 500 MD steps starting from their initial structures. Shown in Fig. 3 are results
48 for the β -sheets, we have very similar ones for the β -rolls (not shown).
49
50
51
52
53

54 Fig. 3, shows that the minimum of free-energy for the monomer and one of the monomers
55 in the dimer occur almost at the same value of the RMSD. This is because in our simulations
56 the monomers in the dimer melt into random-coil-like structures, separate and turn into two
57
58
59
60

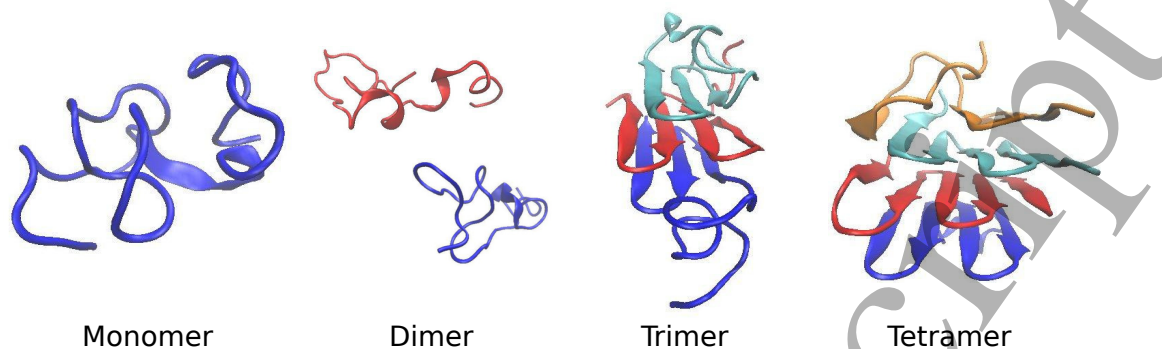


FIG. 2. Representative snapshots of simulations at the temperature $T = 298K$ starting from a monomer, a dimer, a trimer and a tetramer of β -sheets with five repeat units. Here each monomer contains $n = 5$ repeat units.

individual random-coil monomers. The minimum of free energy moves to lower values of the RMSD for the monomers that reside in the center of the trimer and the tetramer; the RMSD curves corresponding to these monomers are in Fig. 3 referred to as *trimer(middle)* and *tetramer(middle(1))*.

There are two minima in the free-energy landscape of the *trimer(middle)* that are indicated in the figure as 5S1 and 5S2. One of these minima, is associated with a β -sheet structure and the other one corresponding to the deepest minimum, 5S2 is a β -sheet with one of the strands molten. Renderings of these structures are given in Fig. 4. For the *tetramer(middle(1))*, there are also two minima approximately on the same values of the RMSD but the deepest one corresponds to the 5S1 structure. The width of the potential well in this case is smaller compared to that of the *trimer(middle)* protein indicating that the *tetramer(middle(1))* protein has less conformational freedom.

Within the trimer and the tetramer structures, the conformations of the monomers in the center fluctuate much less than those at both ends of the aggregates. The fluctuations can be quantified by measuring the RMSD for the individual monomers in the trimer and tetramer. Shown in Fig. 5 is the RMSD for the individual monomers in the tetramer. The monomers in the middle of the tetramer are very stable and their free-energies exhibit a deep minimum at 0.22 nm. The standard deviation of RMSD for these monomers is about 0.04 nm. In contrast, the free-energy landscape of the ones on the top and bottom of the aggregate, which are in contact with water, is much broader. In this case the standard deviation of RMSD is about 0.2 nm. Their free-energy minimum corresponds to the 5S2

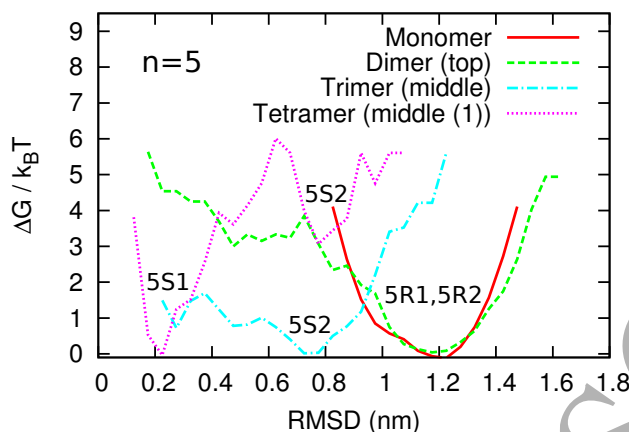


FIG. 3. Free energy, ΔG , obtained from simulations starting from a monomer, a dimer, a trimer and a tetramer of β -sheets consisting of five repeat units, as a function of the root-mean square deviation of snapshots of the monomer, one of the monomers in the dimer and monomers in the middle of the trimer and tetramer from their initial (β -sheet) structures. The structures corresponding to the minima of the free energy are denoted as 5S1, 5S2, 5R1 and 5R2 and shown in Fig. 4.

structure that can turn into the 5S3 structure by overcoming a relatively small free-energy barrier as is clear from Fig. 5.

Our results on structures starting from β -rolls mirror these findings. The only difference between the initial β -roll and β -sheet structures is in the conformational asymmetry of the top and bottom monomers. For the β -sheet structure our results point at a much more symmetric configuration, as evident from Fig. 5. Since both β -rolls and β -sheets are inherently chemically not mirror symmetric, we expect both types of assembly to be symmetry-broken. Obviously we cannot exclude the possibility that in a longer simulation run both structures would exhibit this structural asymmetry.

The monomers in the center of the trimers and tetramers are more stable than the ones on the top and the bottom of the aggregate. However, the free-energy landscapes of the monomers in the center of the trimers are broader than those in the center of the tetramers. This suggests that the trimer is not as stable a structure as a tetramer. Interestingly, the monomers at the top and the bottom of the tetramer are also structurally more stable than the outer monomers of the trimer (data not shown). This indicates that the more stable β -sheets in the center of the aggregates restrict the fluctuations of the outer monomers via inter-molecular interactions and hence stabilize the structure as a whole. Considering that

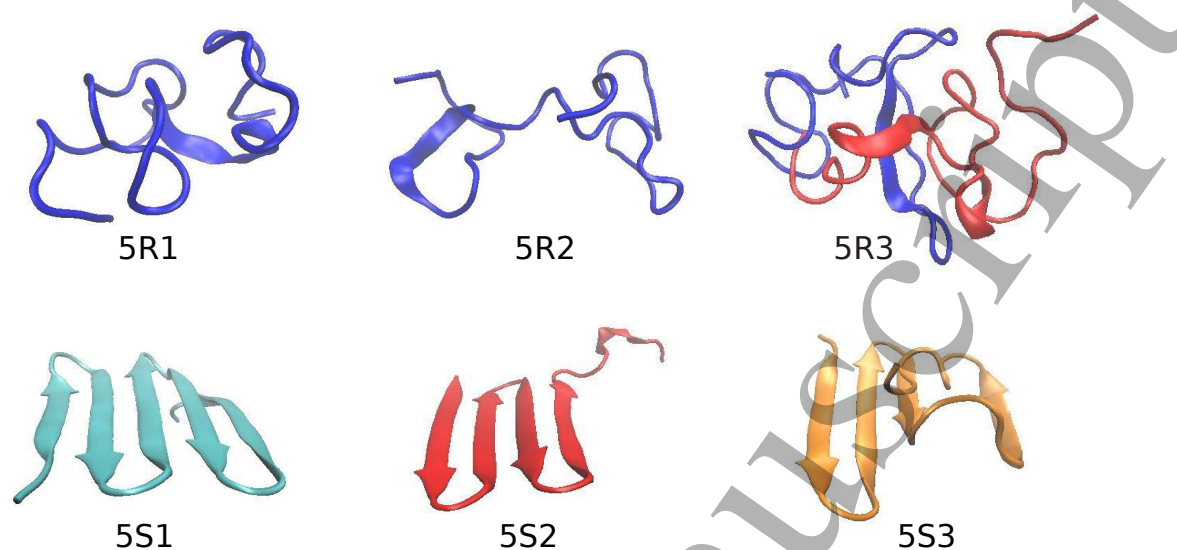


FIG. 4. Snapshots of minimum free energy configurations representing a (5R1) random-coil structure, (5R2) a random coil with more surface exposed to water, (5R3) two random coils bound to each other, (5S1) a β -sheet, (5S2) a β -sheet with a molten strand, (5S3) a β -sheet with two molten strands. See also Fig. 5. These snapshots are taken from the lowest temperature replicas corresponding to a temperature $T = 298 \text{ K}$. 5R1, 5R2 are snapshots of simulations started from a monomer β -sheet and 5R3 is taken from a simulation on a dimer starting from β -sheet structures. 5S1, 5S2 and 5S3 structures are taken from simulation on tetramers consisting of β -sheets.

β -sheet dimers are not stable at all, this suggests a gradual increase in stability of the assemblies with increasing degree of polymerization. This is in-line with density functional theory calculations of Filot et al., where the authors find that the absolute value of average interaction energy between self-assembling C_3 -symmetrical trialkylbenzene-1,3,5-tricarboxamide-based polymers increases as a function of the number of monomers.⁴⁰

To verify that the stability of our protein increases with increasing degree of polymerization, we calculate the number of hydrogen bonds and the solvent accessible surface area as a function of the number of monomers in the aggregates. Shown in Fig. 6 is the free energy landscape of all of our assemblies as a function of number of hydrogen bonds per monomer, starting off from β -sheet structures. In Fig. 7, only simulation results for the β -roll structures for the monomer and tetramer are shown for clarity. As expected, the location of the minimum of the free energy of the monomer and the dimer configurations is approximately at the same value, corresponding to about four hydrogen bonds per monomer. Again, this is

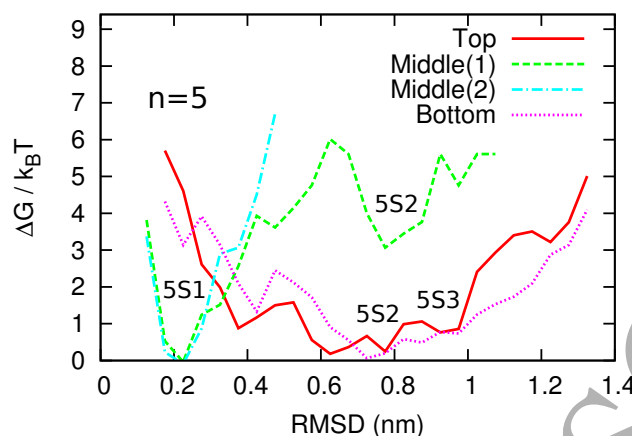


FIG. 5. Free energy, ΔG , obtained from simulations starting from a tetramer of β -sheets with five repeat units as a function of the root-mean square deviation of snapshots of the individual monomers in the tetramer. The structures corresponding to the minima of the free energy that are denoted as 5S1, 5S2, 5R1 and 5R2 shown in Fig. 4. The average of RMSD of the monomers is $RMSD \approx 0.22$ nm with a standard deviation of about 0.04 nm. The average and standard deviation of the top and bottom monomers are about 0.7 nm and 0.2 nm.

because the β -sheet monomers in the dimer melt, separate and form individual monomers. As the number of monomers increases the average number of hydrogen bonds per monomer shifts to higher values and the corresponding standard deviation decreases. This confirms that our protein structures become more stable as the number of monomers in the aggregates increases.

For the trimer and tetramer configurations, we find that the location of the free-energy minimum shifts to larger number of hydrogen bonds as a function of the number of monomers, because of the formation of inter-monomer hydrogen bonds. So, the free energy gain associated with additional hydrogen bonding increases with the number of monomers, arguably offsetting the entropy loss of formation of the β -sheet structures from random coils (molten globules) that also increases with the number of monomers. In fact, hydrogen bonds are also involved in the bonding between the proteins. An even deeper analysis is recurrent.

Assemblies are presumably not only stabilized by inter-molecular hydrogen bonding but also by hydrophobic interactions between the proteins. This we can probe by calculating the solvent accessible surface area (SA) for each of the assemblies because our proteins are

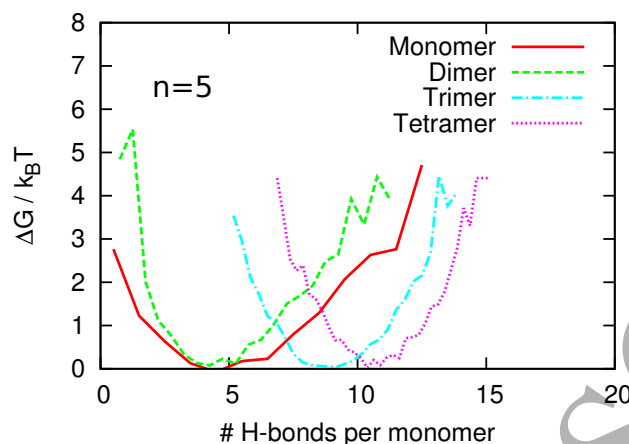


FIG. 6. Free energy, ΔG , obtained from simulations at temperature $T = 298$ K starting from a monomer, a dimer, a trimer and a tetramer of β -sheets of our protein consisting of five repeat units, as a function of the number of hydrogen bonds divided by the number of monomers. The values of the average and standard deviation for a monomer, a dimer, a trimer and a tetramer are (4.6, 2.19), (4.5, 1.7), (8.7, 1.5) and (10.5, 1.4), respectively.

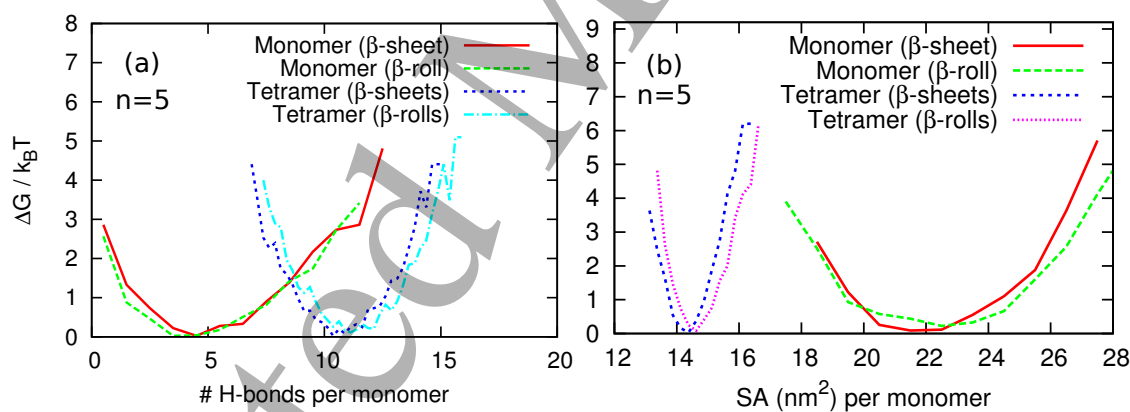


FIG. 7. Free energy, ΔG , as a function of (a) the number of hydrogen bonds per monomer and (b) the solvent accessible surface area, SA, for the monomers and tetramers computed from simulations that start from β -roll or β -sheet structures.

almost entirely made up of hydrophobic amino acid residues (only Glutamine is polar). For all of the structures starting from β -sheets this is shown in Fig. 8, and for the monomer and tetramer of the initial β -roll configurations in Fig. 7. The figures clearly demonstrate that as the assemblies grow in size more of the hydrophobic surface is buried inside the assembly. This means that the overall free-energy cost of the contact of the hydrophobic

residues of the proteins with water decreases as a function of the degree of polymerization of the assemblies. This is herein bound to be utilized to estimate effective binding free energies and a critical polymerization concentration, that, interestingly, are virtually the same for both the β -sheets and the β -rolls.

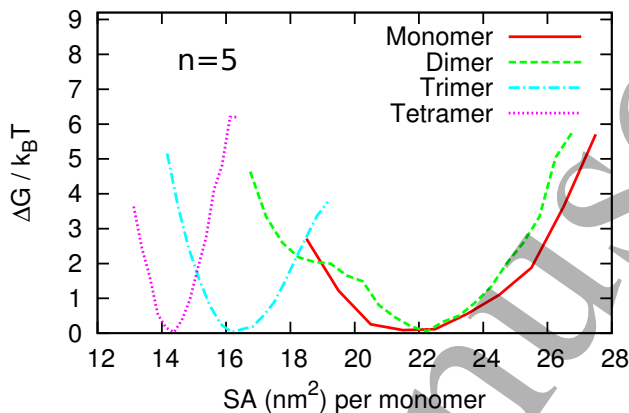


FIG. 8. Free-energy, ΔG , obtained from simulations at temperature $T = 298$ K starting from a monomer, a dimer, a trimer and a tetramer of β -sheets consisting of five repeat units as a function of the solvent accessible surface area divided by the number of monomers.

All of this indicates that β -sheet and β -roll assemblies must be equally stable, at least within our simulations. This is rather surprising, considering that the large-scale structure of the both types of assembly appear very different. Indeed, as shown in Fig. 9b, the width of a tetramer consisting of the β -rolls is smaller than that of the β -sheets. The former is about 1.3 nm and the latter equals approximately 1.9 nm. On the other hand, if we consider the internal structure of the β -roll tetramer then it becomes evident that the β -rolls form inter-molecular β -sheets, that is, hydrogen bonds are formed between strands of two or more of the β -rolls. See Fig. 9. The β -roll tetramer actually resembles a tetramer consisting of tilted β -sheets. This presumably explains the small difference in stability. An in-depth reconsideration follows for the same issue.

First, however, we discuss stability of the structure of assemblies formed by the triblock protein and compare this with what we found for the pentablock.

III B. Stability of the triblock proteins

The question arises how the results that we obtained in the previous sections depend on the length of the silk-like protein sequence, i.e., on the number of repeat units that make up this protein. As alluded to in the introduction, this is relevant to the experiments of

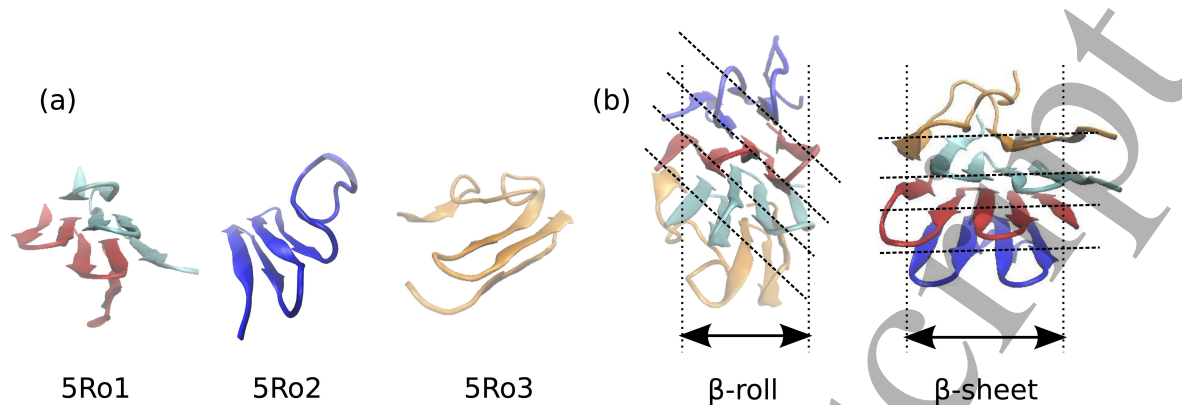


FIG. 9. (a) From the left to the right: snapshots of simulations representing two β -rolls forming inter-molecular β -sheet (5Ro1), a β -roll with a molten strand (5Ro2), a β -roll with a molten strand and a half molten strand (5Ro3). (b) Snapshots of tetramers containing the β -rolls (on the left) and β -sheets (on the right). For the former the approximate width of the aggregate shown by arrows is about 1.3 nm and for the latter it is about 1.9 nm. Orientation of each intra- and inter-molecular β -sheet plane is shown by a dashed line.

de Vries et al., where increasing the number of repeat units of the silk-like block enhances the self-assembly of their proteins on a DNA molecule.²⁵ To address this question, we again rely on REMD simulations of monomer, dimer, trimer and tetramer of β -sheet and β -roll structures of a smaller sequence consisting of $n = 3$ repeat units. This protein is sufficiently large for formation of a β -sheet and a β -roll structures and is small enough to allow us to study the differences between its self-assembly and that of the protein with $n = 5$ repeat units. For compactness, we only discuss the results we obtained for the β -sheets. Again, as was the case for the pentablock protein, our findings for the β -sheet and β -roll assemblies of the triblock are consistent.

In our simulations, the monomers melt into random coil structures. The monomers in the dimers also melt and subsequently separate. For the trimer, first, the top and bottom monomers of the assembly melt. Next, one of these two monomer detaches itself from the trimer. Following that, the remaining folded monomer of the dimer also melts and eventually the dimer disassembles. This indicates that the protein-protein interactions between the triblock proteins are weaker than those of the pentablocks: trimers of the latter species are indeed stable. The tetramer largely follows the route of the trimer, that is, first the outer two monomers melt and one of them detaches towards the end of the simulation run.

1
2
3 Presumably, the remaining trimer would also completely melt and disassemble in a longer
4 simulation run. However, within our simulation time of 20 ns per replica we are not able
5 to observe the completion of this process. It seems that none of the assemblies upto and
6 including the tetramer are stable under the conditions tested.
7
8
9

10 To illustrate how we obtained this conclusion, apart from studying snapshots, we analyze
11 how the free energy curves that we obtain depend on the RMSD, the number of hydrogen
12 bonds and the solvent accessible surface area of the protein configurations sampled. Because
13 we sample over the entire simulation time, which includes the actual melting process, our free
14 energy is not a true free energy and contains information on non-equilibrium configurations.
15 This can be seen in Fig. 10 that shows the free-energy of the aggregates as a function of the
16 RMSD of the simulations starting from the β -sheet structures.
17
18
19
20
21

22 The obtained free-energy curves of the monomer and the monomers in the dimer are
23 similar. They all have a minimum at a relatively large value of the RMSD, confirming that
24 the monomers of the dimer separate and form random coil structures. Here, the melting
25 process is so quick that the folded structure does not significantly contribute to the free
26 energy. This is not so for the trimer. The free energy of this structure has two minima, one
27 of which is associated with the original β -sheet and the other one with the molten state.
28 The former is indicated in the figure with the acronym $3S1$ and the latter with $3R2$ that
29 correspond to structures illustrated in Fig. 11.
30
31
32
33
34
35

36 The melting and breakup of the tetramer is so slow that the free energy of the proteins in
37 the center retain most of their original configuration. This expresses itself in the free energy
38 landscape of the monomer in the center of the tetramer shown in Fig. 10 in the existence of
39 a minimum at an RMSD between 0.3 and 0.4 nm. This is consistent with an almost perfect
40 β -sheet structure such as the $3S1$ structure shown in Fig. 11. A similar picture emerges
41 if we compare the free energy landscapes of all four monomers in the tetramer, shown in
42 Fig. 12. All the monomers in the tetramer retain at least some of their original structures,
43 associated with the β -sheet and β -roll. Conformationally more flexible ones, i.e., ones that
44 have to cross a lower barrier to the actual free energy minimum, develop additional minima
45 associated with different structures (denoted $3S2$, $3R1$ and $3R2$ in Fig. 11).
46
47
48
49
50
51
52
53

54 More information can be extracted from Fig. 12 allowing us to compare the structures of
55 the triblock and pentablock proteins. For instance, the width of the free energy well for the
56 monomer below the top monomer, which is denoted as $middle(1)$ in Fig. 12, is wider than
57
58
59
60

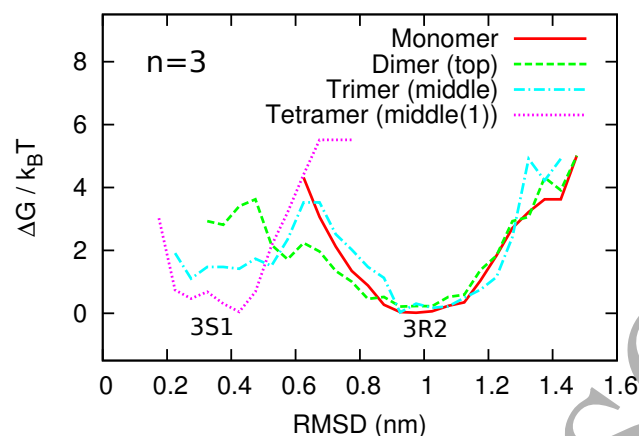


FIG. 10. Free energy, ΔG , obtained from simulations starting from a monomer, a dimer, a trimer and a tetramer of β -sheets with three repeat units as a function of the root-mean square deviation of snapshots of the monomer, one of the monomers in the dimer, the second monomer in the middle of the trimer and tetramer from their initial β -sheet structures. The structures corresponding to the minima of the free energy that are denoted as 3S1 and 3R2 are shown in Fig. 11.

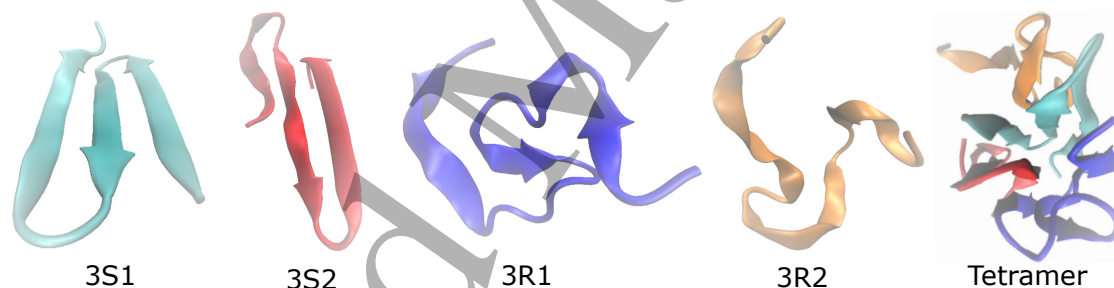


FIG. 11. From the left to the right: snapshots of simulations representing a β -sheet (3S1), a β -sheet with a molten strand (3S2), a half molten β -sheet (3R1), a random coil structure (3R2) and a tetramer containing proteins with three repeat units. The snapshots are taken from simulations at the lowest temperature ($T = 298$ K) replica.

both monomers in the middle of the tetramer consisting of pentablocks. Hence, the middle triblock protein fluctuates more than the corresponding pentablock ones. Also, the global minimum in the free-energy landscape of the second middle monomer, named *middle(2)* in Fig. 12, is associated with a β -sheet with a molten strand (see 3S2 in Fig. 11).

Moreover, the monomer at the bottom of the aggregate melts into a disordered structure, similar to that of 3R1 in Fig. 11. The free-energy landscape of the monomer on the top has three minima that correspond to three different structures. This includes the initial

3S1 structure that morphs into a partially molten, two-stranded 3S2 structure and the final random coil structure 3R2. The latter eventually separates from the rest of the aggregate. Apparently, the time evolution of the top and bottom proteins in the tetramer is not the same. Presumably this reflects the up-down asymmetry of the initial β -sheet structures and hence of the entire assembly.

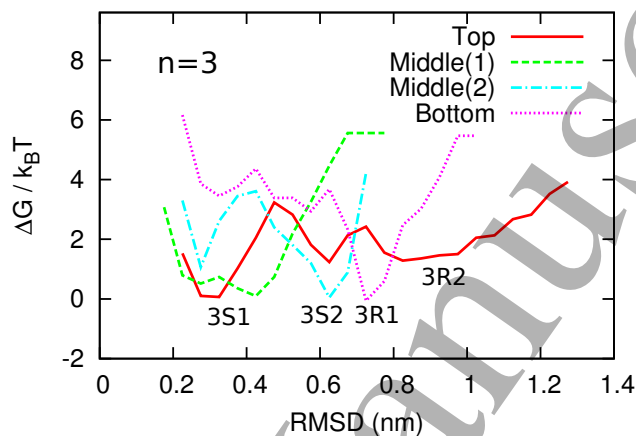


FIG. 12. Free energy, ΔG , obtained from simulations starting from a tetramer of β -sheets with three repeat units as a function of the root-mean square deviation, RMSD, of the individual monomers in the tetramer from their initial β -sheet structure. The structures corresponding to the minima of the free energy that are denoted as 3S1, 3S2, 3R1 and 3R2 are shown in Fig. 11.

To investigate more quantitatively the reason for the structural differences between the aggregates that consist of the two proteins with three and five repeat units, we also compute the free energy of the configurations as a function of the number of hydrogen bonds per monomer. Results are shown in Fig. 13. Again, the free-energy curves of the monomer and the dimer approximately match, because both melt in a short time, but they are different from that of the trimer. This is because the monomers in the trimer melt into random coil structures but they only separate towards the end of the simulation. This causes the shift in the minimum and the asymmetry of the free energy landscape. The shift we observe for the tetramer has the same root cause, amplified by the even slower conformational relaxation processes that take place in this assembly.

For the maximum number of intra- and inter-molecular hydrogen bonds in a fully folded (β -sheet or β -roll) tetramer of the pentablock, we find from our simulations a value of about 15 per monomer. Given the size of the triblock we would expect this to go down to about

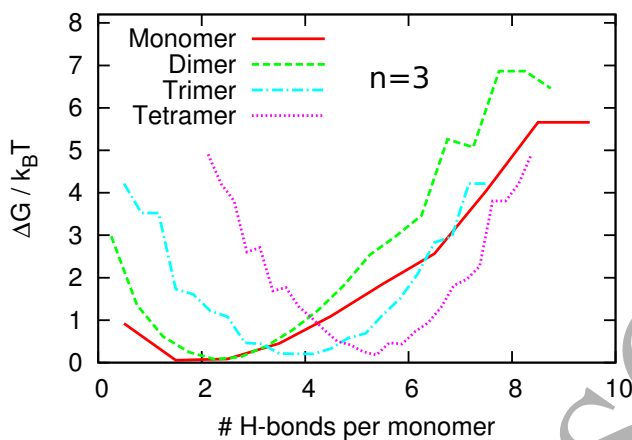


FIG. 13. Free energy, ΔG , obtained from simulations starting from a monomer, a dimer, a trimer and a tetramer of β -sheets with three repeat units as a function of the number of hydrogen bonds divided by the number of monomers.

9. We read off from Fig. 13 that the optimal values are much smaller than that, caused by the melting of the folded structures during the course of the simulation. The melting is partially caused by a reduction in the maximum solvent accessible surface area that can be shielded from water due to the smaller size of the protein.

Fig. 14 confirms this expectation. Comparing with the solvent accessible surface area of the pentablock proteins, we see a fifty per cent reduction in the gain of that quantity upon assembly. Of course, this is not really surprising given the differences in protein size. Also, the ratio of the average solvent accessible surface area, SA, of the molten monomer for the case of triblock to that of the pentablock is about 0.68, which is slightly smaller than the ratio that is expected for collapsed polymers: $(5/3)^{(2/3)} = 0.71$.

For the triblock and pentablock proteins we find that a minimum number of proteins is required to produce a stable aggregate. For the triblock this number is probably larger than four, whilst that for the pentablock it is either two or three depending on whether we are dealing with a β -roll or β -sheet structure. The existence of a difference in configuration of free monomers and the smallest stable folded aggregate implies (i) that this minimum number of disordered free monomers need to come together and (ii) that they have to undergo a conformational change in order to form the smallest thermodynamically stable assembly. Consequently, the linear polymerization of this kind of protein must be nucleated, be it kinetically and/or thermodynamically.⁴¹ Kinetic nucleation occurs through high-free

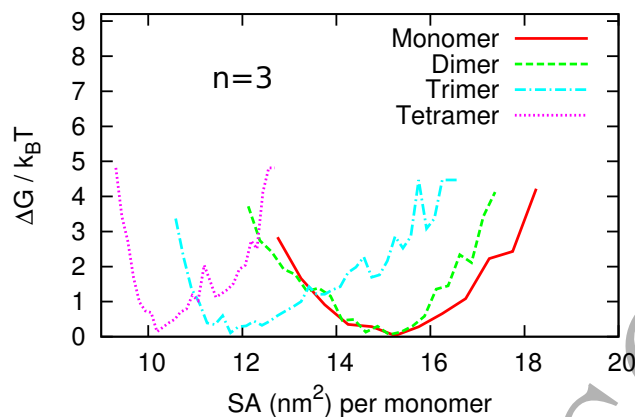


FIG. 14. Free energy, ΔG , obtained from simulations starting from a monomer, a dimer, a trimer and a tetramer of β -sheets with three repeat units as a function of the solvent accessible surface area divided by the number of monomers.

energy intermediate conformational states of monomers between de- and attached states, thermodynamic nucleation through high-free energy conformers in bound state. The latter are stabilized by a binding free energy.

This could explain why the long self-assembled fibers of a triblock construct of two collagen-like disordered sequences sandwiching a silk-like central block take such a long time to form in the experiments of Martens and coworkers.⁴² Incidentally, these authors find that the width of the fibers that form in solutions of their protein is almost half of the value that is expected for fibers made out of β -sheets.⁴² This suggests that the β -roll structure in the assembled state of the protein is the most stable one albeit that we cannot exclude the possibility that this is due to kinetics rather than thermodynamics.

As already alluded to, from our simulations we have to conclude that β -roll and β -sheet assemblies must be virtually equally stable or very nearly that, at least up to the tetramer level. We can make this statement more quantitative. For this we make use of insights based on a simple coarse-grained model that describes the cooperative binding of monomers into supramolecular polymers.⁴¹ Two free energies describe activation and elongation of these self-assembled polymers. The first is a free energy penalty associated with assuming the polymerization-active state of the molecules and the second a binding free energy between polymerization-active molecules in the assembly. The presumption is that the polymerization-active state represents an excited state or conformation of the molecule that in the assembly is stabilized by the binding to its neighbors. Hence, the model describes

thermodynamic nucleation.

These two free energies we now attempt to estimate from our simulations for the two proteins. Within the nucleated assembly model the total binding free energy, ΔG , depends on the number of monomers in an assembly, N , as,

$$\Delta G(N) = (N - 1)\Delta G_b + N\Delta G_c. \quad (1)$$

Here, ΔG_b is the free energy of the formation of a bond between a protein and the assembly, representing elongation of the chain. ΔG_c is the free energy of the conformational change required for the protein to be able to bind. To calculate ΔG_b and ΔG_c , we first compute the $\Delta G(N)$ for our pentablock protein aggregates.

As a first estimate, we presume that binding is driven by hydrophobic interactions that we associate with the reduction of solvent accessible surface area upon binding. We compute the $\Delta G(N)$ by (i) calculating the SA for folded structures upto $N = 4$, (ii) subtract from that the SA of respectively three and four molten monomers, (iii) multiply this by the surface tension of the hydrophobic parts of the protein, γ , (iv) fit Eq. 1 to these two data points. The surface tension, γ , is used in our implicit solvent simulations to calculate the non-polar contribution of the solvation free energy, and has a value of $0.84 k_B T/nm^2$ at room temperature ($T = 298$ K).⁴³

Applying this procedure for the $N = 3$ and $N = 4$ β -sheets of the pentablock protein gives values of $\Delta G_c \simeq +8 k_B T$ and $\Delta G_b \simeq -20 k_B T$. From this admittedly crude estimate we can indeed conclude that the folded structure represent a high-energy state, stabilized by protein-protein interactions. Interestingly, these values are not all that far off from those obtained by Aggeli and coworkers who analyzed their experimental data on the self-assembly of a β -sheet forming oligopeptide using the same nucleated assembly model: $\Delta G_c \simeq +7 k_B T$ and $\Delta G_b \simeq -31 k_B T$.⁴ Note that our estimate is based entirely on the contribution from hydrophobic interactions and hence should be considered as a lower estimate. We return to this issue below.

If we repeat this exercise for the $N = 2, 3$ and 4 β -rolls of the pentablock proteins, we find $\Delta G_c \simeq +3 \pm 2 k_B T$ and $\Delta G_b \simeq -13 \pm 5 k_B T$. For the β -roll we have error bars because we have three data points. This is not so for the β -sheets for which we only have two data points and hence by definition obtain a perfect fit. See Fig. 16. The estimate of the errors for both protein types we expect to be comparable, hence we do not think there

is a real statistically significance between the different values of the free energies that we obtain. Indeed, the critical assembly concentration, X_c (a mole fraction), we obtain this way is virtually identical for the two conformers⁴¹: $k_B T \ln X_c = \Delta G_b + \Delta G_c$. For our β -sheets $\ln X_c = -12$ and for the β -rolls $\ln X_c = -10 \pm 6$. That ΔG_c acts as a thermodynamic nucleation free energy is shown in Fig. 15, where we show the mean degree of polymerization, N , as a function of the protein concentration for $\Delta G_b + \Delta G_c = -12 k_B T$ and two values of $\Delta G_c = +8 k_B T$ and $0 k_B T$, obtained from the two constant assembly models.⁴¹ For $\Delta G_c = +8 k_B T$ the polymerization is sharper than for $\Delta G_c = 0 k_B T$. For a description of the model the reader is referred to the work of van der Schoot.⁴¹

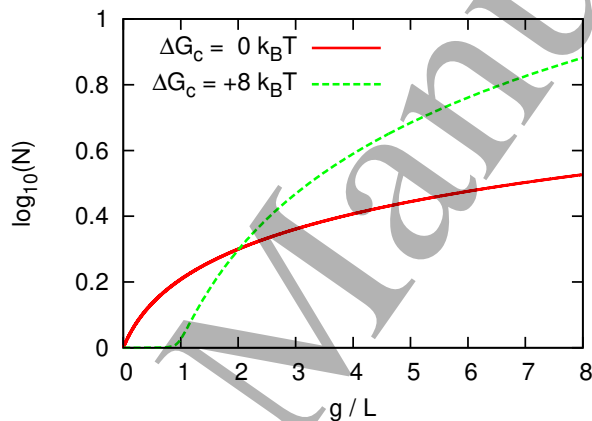


FIG. 15. Mean degree of polymerization, N , as a function of concentration of our pentablock protein for two values of the free energy of conformational change of the protein upon binding, $\Delta G_c = +8 k_B T$ and $0 k_B T$. For both curves $\Delta G_b + \Delta G_c = -12$.

Our analysis shows that according to our simulations self-assembled fibrils consisting of β -sheets and those of β -rolls of our protein must be approximately equally stable. From the experiments of Martens et al.⁴² on their silk-like protein we obtain a lower bound for $\Delta G_b + \Delta G_c$ of $-18 k_B T$. This number is much more negative than what we find but this is not surprising because their silk-like protein is about five times larger than ours ($n = 24$). On the other hand, their block is connected to a very large disordered protein blocks that apparently significantly reduce the net binding free energy of that protein. Interestingly, binding energies close to $20 k_B T$ are typically found in the context of protein assembly of, e.g., viruses.^{44,45}

As discussed in Sec. III A, the fibers that Martens and collaborators find must, because

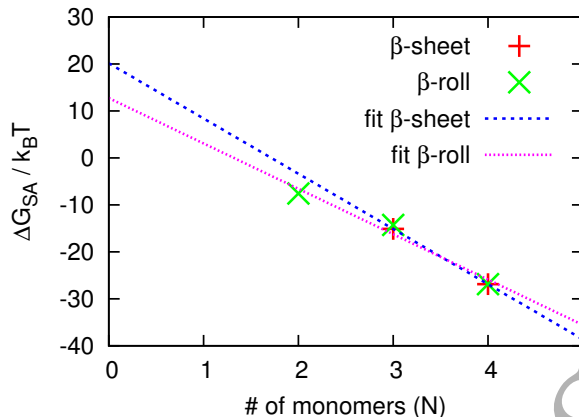


FIG. 16. Free energy contribution of hydrophobic interactions, ΔG_{SA} , obtained from simulations starting from N β -sheets and β -rolls with five repeat units as a function of the number of monomers N . The lines are plotted by fitting our data to eq. 1.

of their dimensions, be of the β -roll variety. This is confirmed by the REMD simulations of Schor and collaborators on a much smaller silk-like protein with $n = 6$.³⁷ The authors find that in monomeric and the dimeric states the β -roll structures are more stable than the others. However, in a recent study where we repeat the simulations of Schor et al. and we find that the most stable monomer structure is a disordered globule not a β -roll.²⁷ This is also true for the dimers.

Note that the repeat unit in Schor's work that was inspired by that of Martens et al. is slightly different from the one studied in this work: $(GA)_3$ -GE instead of our $(GA)_3$ -GX with X alternatingly Q and K. The difference between our previous results and those of Schor on the same sequence we described to differences in simulation time.²⁷ Our explicit solvent simulations were twice and our implicit solvent ones ten times as long. This highlights the importance of full equilibration of the simulation, in this context meaning that the proteins must visit all the replicas sufficiently often.

In this work, by the same REMD equilibration criterion, our simulations should have completely equilibrated. This is because, as mentioned in Sec. II, the average simulation time that it takes for a replica at the lowest temperature to diffuse up to the highest temperature and again down to the lowest one is about 1 to 3 ns and our simulations are 20 to 40 ns long. Therefore during a simulation run replicas sufficiently explore the temperature space. However, we do observe that one of the monomers of the tetramer consisting of our triblock

1
2
3 proteins melts and detaches to form a trimer. We find that starting from a folded trimer,
4 the trimer melts and dissociates. This implies that the simulation on the tetramer at least
5 as far as disassembly is concerned does not reach its thermal equilibrium. This does not
6 mean that the folded tetramer of the pentablock has not equilibrated but does suggest that
7 it would be useful in the future study to revisit this problem using much longer simulation
8 times analyzing results with a specific tool Wordom⁴⁷ that currently are out of our reach.
9

10
11
12 Finally, we have in our discussion above ignored any contribution from intra- and inter-
13 molecular hydrogen bonding to the stability of the proteins. Their contribution is actually
14 not so trivial to assess because one would need to measure the difference in free energy
15 between the hydrogen bonds between the donors and acceptors on the protein and those
16 between moieties on the protein and water. If we ignore this complication and assign a net
17 binding free energy to every hydrogen bond, G_h , then we find that the contribution to ΔG_c
18 and ΔG_b for the pentablock β -sheet structures amount to approximately $-12 G_h$ and $+24$
19 G_h , respectively. A reasonable estimate for G_h is between -2.6 and $-3.2 k_B T$.⁴⁶
20
21
22
23
24
25
26
27

28 This suggests that hydrogen bonding could in principle significantly enhance the stability
29 of the folded aggregate structure but also that hydrogen bonding increases the free energy
30 of the folded state. The latter conclusion is somewhat surprising, because one would naively
31 expect that the folded structure engages in more intra-molecular hydrogen bonds than the
32 molten structure does. Hence, it appears that inter-molecular hydrogen bonds stabilize the
33 β -sheet and β -roll structures of our protein. On the other hand, a naive uncoupling of the
34 effects of burying hydrophobic surface area and the formation of hydrogen bonds ignores
35 the exchange of intra-molecular hydrogen bonds in the globular structure for inter-molecular
36 ones.
37
38
39
40
41
42

43 IV CONCLUSIONS

44
45 We applied replica exchange molecular dynamics simulations to study the influence of the
46 length of the primary sequence of a “blocky” silk-like protein on its quaternary structure.
47 The stability of aggregates containing one, two, three and four of the protein monomers
48 consisting of three or five repeat units was investigated. All simulations were started from
49 two initial configurations: β -sheets and β -rolls. We find that monomers in the assemblies
50 containing the smaller-length proteins melt into random coil structures and separate. Single
51 monomers of both β -sheets and -rolls of the pentablock protein also melt. Dimers of the
52 β -sheet proteins also melt into random coil structures, whereas those of the β -roll seem
53
54
55
56
57
58
59
60

1
2
3 stable. Trimers and tetramers of both folds remain folded within our simulation time.
4

5 This shows that the inter-molecular interactions stabilizing the protein aggregates must
6 be stronger for the longer proteins. We find that the monomers in the center of the trimers
7 and tetramers are conformationally more stable than the ones on the top and the bottom the
8 assemblies. Interestingly, the ones in the middle of the tetramers are more stable than the
9 ones of the trimers. This results in more conformationally stable folded monomers on the
10 top and bottom of the tetramer due to the inter-molecular interactions with the middle ones
11 and suggest that as the assemblies become longer their stability increases. Presumably this
12 effect levels off for sufficiently long assemblies, as is seen in theoretical studies of assemblies
13 of small molecules.⁴⁰
14
15
16
17
18
19

20 Our findings allow an in-depth understanding of self-assembly of protein building blocks
21 and provide a new approach for the preparation of bio-inspired nanoscale materials with a
22 unique morphologies.
23
24

25 **ACNOWLEDGMENTS**

26
27 The work of J. Razzokov is supported by Jepa-Limmat Foundation. We thank Sarah
28 Harris (University of Leeds) and Alexey Lyulin (Eindhoven University of Technology), for
29 useful discussions and advice on the simulations. Eindhoven University of Technology is
30 thanked by J. Razzokov for their hospitality. We are grateful for computer time provided
31 by the Dutch National Computing Facilities at the LISA facility at SURFsara. The work of
32 S. Naderi forms part of the research program of the Dutch Polymer Institute (DPI, Project
33 No. 698). This work was supported by NWO Exacte Wetenschappen (Physical Sciences) for
34 the use of supercomputer facilities, with financial support from the Nederlandse Organisatie
35 voor, Onderzoek (Netherlands Organization for Scientific Research, NWO).
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

- ¹Cheng, J.Y., Ross, C.A., Smith, H.I., Thomas, E.L. Templated Self-Assembly of Block Copolymers: Top-Down Helps Bottom-Up. *Adv. Mater.*, **2006**, 18, 2505-2521.
- ²Masatsugu Shimomura, Tetsuro Sawadaishi. Bottom-up strategy of materials fabrication: a new trend in nanotechnology of soft materials. *Curr. Opin. Colloid Interface Sci.*, **2001**, 6, 11 - 16.
- ³Bai, F., Wang, D., Huo, Z., Chen, W., Liu, L., Liang, X., Chen, C., Wang, X., Peng, Q., Li, Y. A Versatile Bottom-up Assembly Approach to Colloidal Spheres from Nanocrystals. *Angew. Chem. Int. Ed.*, **2007**, 46, 6650-6653
- ⁴Aggeli, A., Nyrkova, I. A., Bell, M., Harding, R., Carrick, L., McLeish, T. C. B., Boden, N. Hierarchical self-assembly of chiral rod-like molecules as a model for peptide β -sheet tapes, ribbons, fibrils, and fibers. *P. Natl. Acad. Sci.*, **2001**, 98, 11857-11862.
- ⁵Vauthey, S., Santoso, S., Gong, H., Watson, N., Zhang, S. *P. Natl. Acad. Sci.*, Molecular self-assembly of surfactant-like peptides to form nanotubes and nanovesicles. **2002**, 99, 5355-5360.
- ⁶Cui, H., Webber, M. J., Stupp, S. I. Self-assembly of peptide amphiphiles: From molecules to nanostructures to biomaterials. *Peptide Science*, **2010**, 94, 1-18.
- ⁷Ulijn, R. V., Smith, A. M. Designing peptide based nanomaterials. *Chem. Soc. Rev.*, **2008**, 37, 664-675.
- ⁸Chen, C. L., Bromley, K. M., Moradian-Oldak, J., DeYoreo, J. J. In situ AFM study of amelogenin assembly and disassembly dynamics on charged surfaces provides insights on matrix protein self-assembly. *J. Am. Chem. Soc.*, **2011**, 133, 17406-17413.
- ⁹Graveland-Bikker, J. F., Schaap, I. A. T., Schmidt, C. F., De Kruif, C. G. Structural and mechanical study of a self-assembling protein nanotube. *Nano Letters*, **2006**, 6, 616-621.
- ¹⁰Klein, M. L., Shinoda, W. Large-scale molecular dynamics simulations of self-assembling systems. *Science*, **2008**, 321, 798-800.
- ¹¹Colombo, G., Soto, P., Gazit, E. Peptide self-assembly at the nanoscale: a challenging target for computational and experimental biotechnology *Trends Biotechnol.*, **2007**, 25, 211-218.
- ¹²Davies, R. P. W., Aggeli, A. Self-assembly of amphiphilic sheet peptide tapes based on aliphatic side chains. *J. Pept. Sci.*, **2011**, 17, 107-114.

- 1
2
3 ¹³Scharnagl, C., Reif, M., Friedrich, J. Stability of proteins: temperature, pressure and the
4 role of the solvent. *BBA-Proteins. Proteom.*, **2005**, 1749, 187-213.
5
6 ¹⁴Shukla, D., Schneider, C. P., Trout, B. L. Complex interactions between molecular ions in
7 solution and their effect on protein stability. *J. Am. Chem. Soc.*, **2011**, 133, 18713-18718.
8
9 ¹⁵Miklos, A. C., Sarkar, M., Wang, Y., Pielak, G. J. Protein crowding tunes protein stability.
10 *J. Am. Chem. Soc.*, **2011**, 133, 7116-7120.
11
12 ¹⁶Pace, C. N., Fu, H., Fryar, K. L., Landua, J., Trevino, S. R., Shirley, B. A., ... Grimsley,
13 G. R. Contribution of hydrophobic interactions to protein stability. *J. Mol. Biol.*, **2011**,
14 408, 514-528.
15
16 ¹⁷Cavazzana-Calvo, M., Fischer, A. Gene therapy for severe combined immunodeficiency:
17 are we there yet?. *J. Clin. Invest.*, **2007**, 6, 1456-1465.
18
19 ¹⁸Cavazzana-Calvo, M., Lagresle, C., Hacein-Bey-Abina, S., Fischer, A. Gene therapy for
20 severe combined immunodeficiency. *Annu. Rev. Med.*, **2005**, 56, 585-602
21
22 ¹⁹Roos, W. H., Bruinsma, R., Wuite, G. J. L. Physical virology. *Nat. Phys.*, **2010**, 6, 733-743.
23
24 ²⁰Giacca, M., Zacchigna, S. Virus-mediated gene delivery for human gene therapy. *J. Con-*
25 *trol. Release.*, **2012**, 161, 377-388.
26
27 ²¹Marshall, E. *Science*, **1999**, 286, 2244-2245.
28
29 ²²Check, E. Gene therapy: a tragic setback. *Nature*, **2002**, 420, 116-118.
30
31 ²³Thomas, M., Klivanov, A.M. Non-viral gene therapy: polycation-mediated DNA delivery.
32 *Appl. Microbiol. Biotechnol.*, **2003**, 62, 27-34.
33
34 ²⁴Smeenck, J. M., Schn, P., Otten, M. B., Speller, S., Stunnenberg, H. G., van Hest, J. C.
35 Fibril formation by triblock copolymers of silklike β -sheet polypeptides and poly (ethylene
36 glycol). *Macromolecules*, **2006**, 39, 2989-2997.
37
38 ²⁵Hernandez-Garcia, A., Kraft, D. J., Janssen, A. F., Bomans, P. H., Sommerdijk, N. A.,
39 Thies-Weesie, D. M., ... & de Vries, R. Design and self-assembly of simple coat proteins
40 for artificial viruses. *Nature nanotechnology*, **2014**, 9, 698-702.
41
42 ²⁶Smeenck, J. M., Otten, M. B., Thies, J., Tirrell, D. A., Stunnenberg, H. G., van Hest, J.
43 Controlled Assembly of Macromolecular Sheet Fibrils. *Angew. Chem. Int. Ed.*, **2005**, 44,
44 1968-1971.
45
46 ²⁷Razzokov, J., Naderi, S., van der Schoot, P. Prediction of the structure of a silk-like
47 protein in oligomeric states using explicit and implicit solvent models. *Soft matter*, **2014**,
48 10, 5362-5374
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3²⁸SalomonFerrer, R., Case, D.A. Walker, R.C. An overview of the Amber biomolecular
4 simulation package. *Wiley. Interdiscip. Rev. Comput. Mol. Sci*, **2013**, 3, 198-210.
- 5
6²⁹Hornak, V., Abel, R., Okur, A., Strockbine, B., Roitberg, A., Simmerling, C. Comparison
7 of multiple Amber force fields and development of improved protein backbone parameters.
8 *Proteins*, **2006**, 65, 712-725.
- 9
10
11
12³⁰Tsui, V., Case, D. A. Theory and applications of the generalized Born solvation model in
13 macromolecular simulations. *Biopolymers*, **2000**, 56, 275-291.
- 14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- ³¹Anandakrishnan, R., Aleksander, Drozdetski., Ross, Walker., Alexey, Onufriev. Speed
of conformational change: comparing explicit and implicit solvent molecular dynamics
simulations. *Biophys. J.*, **2015**, 5, 1153-1164.
- ³²Feig, M., Im, W., Brooks III, C.L. Implicit solvation based on generalized Born theory in
different dielectric environments. *J. Chem. Phys.*, **2004**, 120, 903-911.
- ³³Zhang, L.Y., Gallicchio, E., Friesner, R.A., Levy, R.M. Solvent models for proteinligand
binding: Comparison of implicit solvent Poisson and surface generalized Born models with
explicit solvent simulations. *J. Com. Chem.*, **2001**, 22, 591-607.
- ³⁴Lomize, A.L., Pogozheva, I.D., Mosberg, H.I., Anisotropic solvent model of the lipid bi-
layer. 1. Parameterization of long-range electrostatics and first solvation shell effects. *J.*
Chem. Inf. Model., **2011**, 51, 918-929.
- ³⁵Srinivasan, J., Trevathan, M.W., Beroza, P. and Case, D.A., 1999. Application of a pair-
wise generalized Born model to proteins and nucleic acids: inclusion of salt effects. *Theor.*
Chem. Acco., **1999**, 101, 426-434.
- ³⁶Humphrey, W., Dalke, A., Schulten, K. VMD: visual molecular dynamics. *J. Mol. Graph-*
ics., **1996**, 14, 33-48.
- ³⁷Schor, M., Martens, A. A., Stuart, M. A. C., Bolhuis, P. G. Prediction of solvent dependent
 β -roll formation of a self-assembling silk-like protein domain. *Soft Matter*, **2009**, 5, 2658-
2665.
- ³⁸Kleinjung, J., Fraternali, F. Design and application of implicit solvent models in biomolec-
ular simulations. *Curr. Opin. Struct. Biol*, **2014**, 25, 126-134.
- ³⁹Shrake, A., Rupley, J. A. Environment and exposure to solvent of protein atoms. Lysozyme
and insulin. *J. Mol. Biol*, **1973**, 79, 351 - 371.
- ⁴⁰Filot, I. A., Palmans, A. R., Hilbers, P. A., van Santen, R. A., Pidko, E. A., de Greef, T. F.
Understanding cooperativity in hydrogen-bond-induced supramolecular polymerization: a

- 1
2
3 density functional theory study. *J. Phys. Chem. B.*, **2010**, 114, 13667-13674.
4
5 ⁴¹van der Schoot, P. *Nucleation and Co-Operativity in Supramolecular Polymers. Advances*
6 *In Chemical Engineering*, **2009**, 35, 45-77.
7
8 ⁴²Martens, A. A., Portale, G., Werten, M. W., de Vries, R. J., Eggink, G., Cohen Stuart,
9 M. A., de Wolf, F. A. Triblock protein copolymers forming supramolecular nanotapes and
10 pH-responsive gels. *Macromolecules*, **2009**, 42, 1002-1009.
11
12 ⁴³Sitkoff, D., Sharp, K. A., Honig, B. Accurate calculation of hydration free-energies using
13 macroscopic solvent models. *J. Phys. Chem.*, **1994**, 98, 1978-1988.
14
15 ⁴⁴Kegel, W. K., van der Schoot, P. Competing hydrophobic and screened-Coulomb interac-
16 tions in hepatitis B virus capsid assembly. *Biophys. J.*, **2004**, 86, 3905-3913.
17
18 ⁴⁵Kegel, W. K., van der Schoot, P. Physical regulation of the self-assembly of tobacco mosaic
19 virus coat protein. *Biophys. J.*, **2006**, 91, 1501 - 1512.
20
21 ⁴⁶Sheu, S. Y., Yang, D. Y., Selzle, H. L., Schlag, E. W. Energetics of hydrogen bonds in
22 peptides. *Proc. Nat. Acad. Sci.*, **2003**, 100, 12683-12687.
23
24 ⁴⁷Seeber, M., Cecchini, M., Rao, F., Settanni, G., Caffisch, A. Wordom: a program for
25 efficient analysis of molecular dynamics simulations. *Bioinformatics*, **2007**, 23, 2625-2627.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

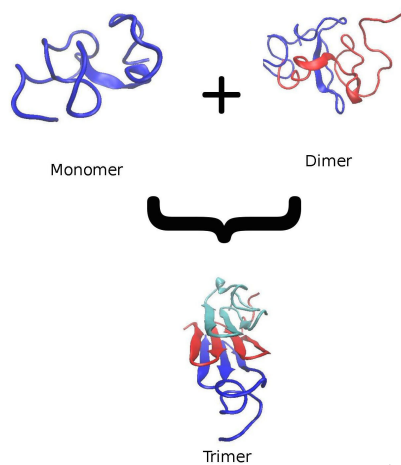


FIG. 17. For Table of Contents Only

Accepted Manuscript