

# **B-SIDER: Computational algorithm for the design of complementary $\beta$ -sheet sequences**

Tae-Geun Yu<sup>1</sup>, Hak-Sung Kim<sup>1,\*</sup>, and Yoonjoo Choi<sup>1,\*</sup>

<sup>1</sup>Department of Biological Sciences, Korea Advanced Institute of Science and Technology (KAIST), Daejeon 34141, Republic of Korea

## **Correspondence**

Yoonjoo Choi

Department of Biological Sciences, Korea Advanced Institute of Science and Technology, Daejeon, 34141, Republic of Korea.

Tel: +82-42-350-2616

Email: [yunjoo.choi@kaist.ac.kr](mailto:yunjoo.choi@kaist.ac.kr)

Hak-Sung Kim

Department of Biological Sciences, Korea Advanced Institute of Science and Technology, Daejeon, 34141, Republic of Korea.

Tel: +82-42-350-2616

Email: [hskim76@kaist.ac.kr](mailto:hskim76@kaist.ac.kr)

## **Keywords**

Beta-sheet; Beta-strand; Complementary sequence; Protein design; Amyloid; Sequence design; Protein secondary structure;

## Abstract

**Motivation:** The  $\beta$ -sheet is an element of protein secondary structure, and intra- and inter-molecular  $\beta$ -sheet interactions play pivotal roles in biological regulatory processes including scaffolding, transporting, and oligomerization. In nature,  $\beta$ -sheet formation is tightly regulated, because dysregulated  $\beta$ -stacking often leads to severe diseases such as Alzheimer's, Parkinson's, systemic amyloidosis and diabetes. Thus, the identification of intrinsic  $\beta$ -sheet forming propensities could provide valuable insight into protein design for the development of novel therapeutics. However, structure-based design methods may not be generally applicable to such amyloidogenic peptides mainly due to high structural plasticity and complexity. Therefore, an alternative design strategy based on complementary sequence information is of great significance.

**Results:** We developed B-SIDER ( $\beta$ -Sheet Interaction DEsign for Reciprocity), a database search method for the design of complementary  $\beta$ -strands. The method makes use of the structural database information and generates a query-specific score matrix. The discriminatory power of the B-SIDER score function was tested on representative amyloidogenic peptide substructures against a sequence-based score matrix (PASTA2.0) and two popular *ab initio* protein design score functions (Rosetta and FoldX). B-SIDER was able to distinguish wild-type amyloidogenic  $\beta$ -strands as favored interactions in a more consistent manner than the other methods. B-SIDER is then prospectively applied to the design of complementary  $\beta$ -strands for the splitGFP scaffold. Three variants were identified to have stronger interactions than its original sequence selected by directed evolution, emitting higher fluorescence intensities. Our results support that B-SIDER can be applicable to the design of other  $\beta$ -strands, assisting in the development of therapeutics against disease-related amyloidogenic peptides.

**Availability:** B-SIDER is freely available at <http://bel.kaist.ac.kr/research/B-SIDER>.

# Introduction

The  $\beta$ -sheet is one of the major units of protein structure (Bhattacharjee and Biswas, 2010), and plays a variety of functional roles in transportation, recognition, scaffolding and enzymatic processes (Marcos, et al., 2018). Recently, the mechanism of  $\beta$ -sheet formation has received much attention because of its close relations with several critical diseases such as Alzheimer's disease, Parkinson's disease, type 2 diabetes and systemic amyloidosis (Chiti and Dobson, 2017; Richardson and Richardson, 2002). Such diseases are known to be linked to the precipitation of dysregulated  $\beta$ -stacking between neighboring  $\beta$ -strands (Colletier, et al., 2011; Liu, et al., 2012; Matthes, et al., 2014). In this regard, the information about the amino acid propensity of intrinsic  $\beta$ -sheet forming motifs and its use in the design of their complementary sequences are crucial for understanding the mechanism of  $\beta$ -sheet formation and developing potential therapeutics specifically targeting aggregation-prone regions (Giorgetti, et al., 2018).

While structure-based protein design approaches have shown notable successes in several cases (Huang, et al., 2016), their application to *de novo*  $\beta$ -sheet designs still remains challenging (Dou, et al., 2018; Marcos, et al., 2018). Structure-based design approaches require a well-defined protein structure, but amyloidogenic peptides usually have highly disordered structures (Jang, et al., 2016). Structural identification of such peptides has long been hindered by high degrees of structural plasticity, transiency, and complexity due to self-oligomerization (Dovidchenko and Galzitskaya, 2015; Zheng, et al., 2016). It is thus necessary to exploit the complementarity across neighboring  $\beta$ -strand pairs using sequence information.

Intriguingly, significant conservation and covariations of residue pairs between neighboring  $\beta$ -strands were identified in many protein families (Mandel-Gutfreund, et al., 2001). For instance, pairs of  $\beta$ -branched residues and cysteines are preferred at nonhydrogen-bonded positions.

Aromatic residues tend to be paired with valine or glycine (Steward and Thornton, 2002). Several computational algorithms have been developed to predict aggregation-prone regions based on the internal  $\beta$ -sheet forming patterns. While different in details, they make use of either statistical potentials such as Tango (Fernandez-Escamilla, et al., 2004), PASTA (Trovato, et al., 2006), SALSA (Zibae, et al., 2007), BETASCAN (Bryan Jr, et al., 2009), and Waltz (Maurer-Stroh, et al., 2010) or physicochemical properties of amino acids (Tartaglia and Vendruscolo, 2008). Additionally, consensus methods and machine-learning approaches have also been developed (Kim, et al., 2009; Tsolis, et al., 2013), showing fine agreements with experimental results.

It has been reported that  $\beta$ -strand interactions can be stabilized by introducing the  $\beta$ -sheet favored pairs (Kortemme, et al., 1998; Minor Jr and Kim, 1994; Quinn, et al., 1994; Stranges, et al., 2011) and charge pairing between neighboring  $\beta$ -strands (Shammas, et al., 2011; Wang and Hecht, 2002; West, et al., 1999). Recent studies showed that fragments derived from the amyloidogenic region can be used for  $\beta$ -stacking modeling (Gallardo, et al., 2016; Liu, et al., 2012). While the use of the amino acid pairing information in protein design has been attempted in elsewhere, practical applications of such patterns have been limited mainly owing to the lack of comprehensive quantification for residue pairing and noisy patterns of  $\beta$ -sheet forming residue pairs (Bhattacharjee and Biswas, 2010; Fujiwara, et al., 2012; Hutchinson, et al., 1998). The  $\beta$ -sheet forming peptides appear to have poor sequence commonalities and imperfect repeats (Bryan Jr, et al., 2009). Therefore, careful curation of meaningful patterns is required for the practical protein design strategy of complementary  $\beta$ -strands.

Herein, we present a database search method, B-SIDER ( $\beta$ -Stacking Interaction DEsign for Reciprocity), to design complementary  $\beta$ -strands. The method generates a query-specific score matrix from the structure database. To utilize the pairing information and overcome the pattern

noise, we hypothesized that significant complementary pairs can be amplified by superposing a subset of sequence fragments. Moreover, the recent growth boom of  $\beta$ -sheet structures (Marcos and Silva, 2018) allows the solid statistics of  $\beta$ -sheet forming residue pairings (Sormanni, et al., 2015). Based on the hypothesis and statistics of  $\beta$ -sheet forming residue pairings, we developed a fast and reliable computational method for the design of complementary  $\beta$ -strand sequences. The methodology augments  $\beta$ -sheet forming residue preferences through overlaying complementary fragment sequences (**Fig. 1**). We retrospectively validated our approach using a set of curated amyloidogenic targets and compared it with two popularly used structure-based methods (Rosetta and FoldX) and a sequence-based aggregation prediction algorithm (PASTA2.0). Our algorithm was shown to clearly distinguish favorable  $\beta$ -sheet forming sequences entirely based on the query sequence, whereas the structure-based energy functions exhibited inconsistent results depending on targets. The utility and potential of our method were demonstrated by designing novel complementary peptides for splitGFP. The designed sequences showed stronger interactions with neighboring strands of the scaffold and consequently higher fluorescence emissions than the original peptide selected by directed mutagenesis (Cabantous, et al., 2005).

## Methods

### Computational algorithm for the design of complementary sequences

#### *Collection of $\beta$ -strand information*

Non-redundant structures determined by high resolution X-ray crystallography were collected from the PDB: < 90 % sequence identity, < 3 Å resolution. Given the query target sequence, the non-redundant structure database was used to extract pairing information from

matched sequences. Initially, the target sequence is divided into linear moving-windows whose residues in length range from 3 to the entire target sequence length. Any structures with identical target sequence fragments to the split queries were collected, followed by further filtering based on the definition of  $\beta$ -sheet secondary structure (the distance between backbone nitrogen-oxygen atom pairs  $< 5 \text{ \AA}$ ). In order to remove redundancy, protein structures that contain the matches were compared using TMalign (Zhang and Skolnick, 2005). If TM-score  $> 0.7$  and sequence identity  $> 90 \%$ , one of the matched sequences was removed.

While the method is applicable to both parallel and anti-parallel  $\beta$ -sheets in theory, we mainly focused on anti-parallel  $\beta$ -sheets in this study since anti-parallel cases are more frequently observed compared to parallel ones (Hubbard, 1994). Disease-related amyloidogenesis is also known to be initiated with anti-parallel  $\beta$ -sheets and soluble oligomeric amyloid species mainly exist as anti-parallel (Cerf, et al., 2009; Gordon, et al., 2004).

### ***Complementary sequence score***

The  $\beta$ -sheet complementarity score function is derived from the environment-specific substitution score (Choi and Deane, 2010). We hypothesized that each position of a  $\beta$ -strand is independent of one another, and their complementarities are determined by residue pairs from neighboring strands. Given the query sequence, all of the identified neighboring sequences are pooled together as described in the previous section. The amino acid frequency at each complementary position is counted as

$$A_{i,p} = O_{i,p} / \sum_i O_{i,p}. \quad (1)$$

where  $A_{i,p}$  is the frequency of a certain amino acid  $i$  at a specific complementary position  $p$ .  $O_{i,p}$  is the total count of the amino acid at  $p$ . The background frequency of the certain amino acid ( $B_i$ ) is counted from the HOMSTRAD database (Mizuguchi, et al., 1998) and calculated as

$$B_i = \sum_i O_{i,p} / \sum_{i,p} O_{i,p}. \quad (2)$$

The complementary sequence score of the amino acid at the position ( $S_{i,p}$ ) is calculated as

$$S_{i,p} = -\log(A_{i,p}/B_i). \quad (3)$$

It should be noted that the complementarity score is completely data-driven, i.e. if an amino acid never appears at a certain position, a high penalty score is imposed. We only consider complementary amino acids which are found at least once in the entire identified sequences. The final score is the sum of the scores at all of the positions.

## Protein expression and complementation assay

### *Gene construction*

The gene coding for splitGFP (Cabantous, et al., 2005) consists of the 1-10<sup>th</sup> strands (GFP1-10) and 11<sup>th</sup> strand (GFP11) template (**Table S1**). They were cloned into pET-28a (Novagen) vector between the Nde-I and Xho-I restriction sites. We introduced additional mutations to GFP1-10 to inhibit aggregation and convenient expression (Kim, et al., 2015). The GFP11 strand was fused with a P22 virus-like particle scaffolding protein (McCoy and Douglas, 2018) for soluble and stable expression. Mutations on GFP11 were introduced by PCR using the mutagenic primers (**Table S2**), and the resulting genes were cloned into the pET-28a vector. Six histidine residues were fused to the N-terminal of the GFP1-10 and GFP11 genes as an affinity purification tag.

158

## 159 ***Protein expression and purification***

160 All of the constructs were transformed into BL21 (DE3) *E. coli* strains. The transformed  
 161 cells were grown overnight and inoculated into a Luria-Bertani media containing 50 µg/ml of  
 162 kanamycin at 37 °C. Then, cells were grown until an optical density of the cells reached 0.6-0.8 at  
 163 600 nm, followed by addition of 0.7 mM of IPTG (isopropyl β-D-1-thiogalactopyranoside) for  
 164 induction. After incubation for 16-18 hours at 18 °C, the cells were harvested and suspended in a  
 165 lysis buffer containing 50 mM Tris (pH 8.0), 150 mM NaCl, and 5 mM imidazole. The suspended  
 166 cells were disrupted by sonication, and insoluble fractions were removed by centrifugation at  
 167 18,000 g for 1 hour. The supernatants were filtered using 0.22 µm syringe filters and purified  
 168 through affinity chromatography with Ni-NTA agarose Superflow (Qiagen). The solutions were  
 169 applied to the resin-packed columns and washed with a buffer containing 50 mM Tris (pH 8.0),  
 170 150 mM NaCl, and 10 mM imidazole, until no protein was detected by Bradford assay. Then, an  
 171 elution buffer (50 mM Tris (pH 7.4), 150 mM NaCl, 300 mM imidazole) was applied to the column.  
 172 The buffer exchange was performed by PD-10 column (GE health-care) to PBS (phosphate  
 173 buffered saline, pH 7.4). The concentrations of the proteins were determined by measuring the  
 174 absorbance at 280 nm. All the purification processes were performed at 4 °C. The purities of  
 175 proteins were then evaluated by SDS-PAGE.

176

## 177 ***Complementation assay***

178 The assembly of splitGFP variants was monitored and measured by fluorescence  
 179 complementation assay. Excessive amount of GFP1-10 (50 pmol) in 180 µl and 20 µl of equal



molar concentration of each GFP 11 strand (3 pmol) were co-incubated in PBS buffer (pH 7.4). Fluorescence kinetics ( $\lambda_{ex}$  = 488 nm /  $\lambda_{em}$  = 530 nm) were monitored for 12 hours at 25 °C by TECAN infinite M200 microplate reader at 5 minutes intervals (Cabantous, et al., 2005) with shaking for 2 seconds between intervals. Each experiment was performed in triplicate with Nunc F 96 Micro-well black plate, blocked with a solution of PBS containing 0.5 % of Bovine serum albumin (BSA) for 30 minutes before the assay.

## Results and Discussion

### Overview of the design process

We hypothesized that repetitively observed amino acid pairing patterns indicate the “smoking-gun” of strong preferences to  $\beta$ -sheet. It was also assumed that the sequence with most frequent patterns would directly form a  $\beta$ -sheet without considering other environmental contributions.

There are two major steps in the algorithm: 1) The extraction of  $\beta$ -sheet complementarity information and 2) the construction of scoring matrix (See **Fig. 1** and **Fig. S1**). When a query sequence is given, it is fragmented into several pieces of short peptides longer than three residues in length and matched neighboring strands are collected. This fragmentation and overlaying processes naturally impose weights on complementary-prone positions and amplify pattern signals (**Fig. S1**). After the collection of matched sequences, a position-specific complementarity scoring matrix is constructed. The obtained scoring matrix is used to evaluate and design the complementarity of  $\beta$ -strand interactions.

## Validation of the score function on retrospective cases

In an effort to validate the complementarity score, we manually curated a test set of naturally occurring  $\beta$ -strand pairs whose environmental effects are minimal. It is known that  $\beta$ -strand pairing is in general greatly hindered by local environments (Zaremba and Gregoret, 1999), but amyloidogenic peptide segments are known to form natural  $\beta$ -sheets in themselves (Trovato, et al., 2006). We thus selected a set of widely known amyloidogenic structures whose aggregation-prone regions have been identified (**Table 1** and **Table S3**).

To assess the complementarity of the native sequences, we compared their scores with those of random sequences. The natural amyloidogenic segments are known to be highly aggregation-prone, so they are expected to be highly preferred, i.e., having fairly low scores in the random sequence score distributions. **Figure 2** shows that all the native sequence scores are ranked extremely low in all of the distributions. On average, the native sequences are within 4.1 % of the distributions (**Fig. 2**). The results indicate that the scoring function is extremely useful in detecting favorable  $\beta$ -strand counterparts.

We also compared the B-SIDER score with two structure-based all-atom energy functions and a sequence-based score matrix. For structure-based methods, we picked Rosetta (Talaris 13) (Alford, et al., 2017; Kuhlman, et al., 2003) and FoldX (Schymkowitz, et al., 2005), which have been popularly used in *de novo* protein designs (Fleishman, et al., 2011; Rocklin, et al., 2017). PASTA2.0 (Walsh, et al., 2014) is a method to predict aggregation-prone regions using the scoring matrix derived from residue pairing patterns of  $\beta$ -sheets. To avoid any biases, 1,000 random sequences were newly prepared per target. “FastRelax” protocol (Tyka, et al., 2011) from Rosetta (Ver. 3.7), “BuildModel” command from FoldX (Ver. 4.0) and the scoring matrix from PASTA2.0 were used against the native and the random sequences. The predictive power of a score function

was assessed by the percentile value of the native sequence score against the random sequence score distribution.

**Figure 3** shows that structure-based score functions are in general worse than the sequence-based scoring matrices. The results indicate that the Rosetta energy score function is not sufficiently accurate for ranking complementary  $\beta$ -strands (35.8<sup>th</sup> percentile on average), whereas the predictive powers of PASTA2.0 and FoldX were moderate, showing 10.8<sup>th</sup> and 14.7<sup>th</sup>, respectively. B-SIDER was shown to be the most accurate in an extremely consistent manner. While the assessment of PASTA2.0 is also fairly consistent, the query-specific nature of B-SIDER may give better results.

Considering the Rosetta relax protocol performs flexible backbone refinements, the use of the fixed-backbone calculation seems to be better for the evaluation of  $\beta$ -sheet complementarity. It should be noted that the inconsistent results of Rosetta imply that FoldX prediction would be also highly driven by structure preparation, i.e. design with ill-defined models may not be generally successful. On the other hand, B-SIDER and PASTA2.0 do not depend on query structures, and thus, it can be applied to general cases such as  $\beta$ -sheet interactions with high structural plasticity and poor structural integrity, which are the common features of amyloidogenic peptides. Furthermore, the process of collecting complementary motifs of B-SIDER also appeared to be powerful, making it possible to distinguish favorable complementary sequences not easily detected by one-to-one residue pairing.

### **Prospective application of the algorithm to splitGFP**

As shown in the retrospective test, B-SIDER is extremely useful in discriminating naturally  $\beta$ -strand forming sequences. As a proof of concept, we prospectively designed novel

complementary  $\beta$ -strands for splitGFP. SplitGFP is a fragmented protein pair derived from superfolderGFP (Cabantous, et al., 2005), comprising a scaffold containing 10  $\beta$ -strands (GFP1-10) and its complementary  $\beta$ -strand peptide (GFP11). GFP11 specifically interacts with GFP1-10 and the strand tightly forms a stable  $\beta$ -sheet structure, which facilitates the chromophore maturation in an irreversible manner (Köcker, et al., 2018). This assembly process results in the emission of the green fluorescence. Because GFP11 is known to be disordered in solution, its conformational transition from the disordered to induced  $\beta$ -sheet is similar to amyloidogenic peptides (Ito, et al., 2013; Xu, et al., 2005). This model system thus efficiently assesses whether designed sequences by B-SIDER have favorable  $\beta$ -sheet interactions.

Original GFP11 was designed by directed mutagenesis, and it shows a high intrinsic propensity to form hydrogen bonds with the neighboring  $\beta$ -strands of GFP1-10 (Miller, et al., 2015). In our case, the queries are the neighboring strands of GFP1-10 (**Fig. 4A**). It is known that the residues pointing inward (1, 3, 5, and 7<sup>th</sup> positions) directly interact with the chromophore and thus they were not subject to mutation. B-SIDER identified 2,637 non-redundant sequences from the structure database. The native sequence is ranked at a modest score among randomly chosen 1,000 possible sequence variants (46<sup>th</sup> percentile), indicating that there could be room for complementary sequences with stronger interactions than the original one (**Fig. 4B**). We then selected 10 sequences with the lowest B-SIDER scores (top\_vars; **Table 2**). Amino acid compositions of the 10 variants are mostly hydrophobic or branched amino acids (**Fig. S2**). Additionally, one sequence with a high score (> 75<sup>th</sup> percentile) was randomly selected as a negative control (neg\_var, 77<sup>th</sup>).

The selected variants were successfully expressed and purified (**Fig.S3**) except for four clones (top\_var6, 7, 8, and 10) which were observed to be insoluble, perhaps due to aggregation.

Among those expressed, three variants (top\_var1, 2, and 9) showed faster assembly patterns and higher signals compared to the original GFP11 (**Fig. 5**). No functional aberrance with excitation and emission was observed (**Fig. S4**). All the successful variants, which emitted stronger fluorescence levels than the original one, were shown to have the pair of phenylalanine and threonine at the positions 6 and 8, respectively. These results demonstrate that the designed variants indeed formed complementary  $\beta$ -strands in a more favorable fashion than its original peptide as predicted. The other variants showed slightly lower signals than the original one, but still gave rise to clear fluorescence signals (**Fig. 5**). The negative control (neg\_var) barely emitted any signal, suggesting that the score indeed indicates the complementarity of  $\beta$ -stacking interactions. We also assigned scores of the GFP11 variants using other scoring methods. As shown in the retrospective test set, Rosetta and FoldX were not able to discriminate top\_vars as favorable (**Fig. S6**). However, PASTA2.0 was again fairly accurate in this case.

## Conclusion

$\beta$ -sheet forming patterns are crucial for understanding the aggregation mechanism of disease-related  $\beta$ -sheets and developing potential therapeutics against them. Unlike  $\alpha$ -helices, however, there has been no established design principle for the complementarity of  $\beta$ -sheets. In this study, we developed B-SIDER, a database search method for the design of complementary  $\beta$ -strands based on the intrinsic  $\beta$ -sheet forming propensities. Statistical patterns of interacting residue pairs between neighboring  $\beta$ -strands enable to quantify the complementary interaction. We demonstrated that the statistical potential can be directly applied to the design of complementary  $\beta$ -strand sequences. Using splitGFP as a model system, we successfully designed fragment variants, which led to stronger fluorescence emissions than the native one originally identified by directed

mutagenesis. The results clearly indicate that B-SIDER can be useful for the detection and design of  $\beta$ -stacking interactions between unstructured fragments. Therefore, our approach can find wide applications to protein designs where structure-based methods are not effective, including the development of protein binders specifically against disease-related intrinsically disordered proteins.

## Acknowledgements

This work was performed using Alphacom high-performance computing cluster in the department of biological sciences at the Korea Advanced Institute of Science and Technology (KAIST).

## Funding

This work was supported by the Korea Research Fellowship Program [2016H1D3A1938246 to Y.C.], Global Research Laboratory (NRF-2015K1A1A2033346), and Mid-Career Researcher Program (NRF-2017R1A2A 1A05001091) of the National Research Foundation (NRF) funded by the Ministry of Science and ICT.

## References

1. Alford, R.F., *et al.* The Rosetta all-atom energy function for macromolecular modeling and design. *J. Chem. Theory Comput.* 2017;13(6):3031-3048.
2. Allen, M., *et al.* Raincloud plots: a multi-platform tool for robust data visualization. *PeerJ* 2018;6:27137.
3. Bhattacharjee, N. and Biswas, P. Position-specific propensities of amino acids in the  $\beta$ -strand. *BMC Struct. Biol.* 2010;10(1):29.
4. Bryan Jr, A.W., *et al.* BETASCAN: probable  $\beta$ -amyloids identified by pairwise probabilistic analysis. *PLOS Comput. Biol.* 2009;5(3):1000333.
5. Cabantous, S., Terwilliger, T.C. and Waldo, G.S. Protein tagging and detection with engineered self-assembling fragments of green fluorescent protein. *Nat. Biotechnol.* 2005;23(1):102.

6. Cerf, E., *et al.* Antiparallel  $\beta$ -sheet: a signature structure of the oligomeric amyloid- $\beta$  peptide. *Biochem. J.* 2009;421(3):415-423.
7. Chiti, F. and Dobson, C.M. Protein misfolding, amyloid formation, and human disease: a summary of progress over the last decade. *Annu. Rev. Biochem.* 2017;86:27-68.
8. Choi, Y. and Deane, C.M. FREAD revisited: accurate loop structure prediction using a database search algorithm. *Proteins* 2010;78(6):1431-1440.
9. Colletier, J.-P., *et al.* Molecular basis for amyloid- $\beta$  polymorphism. *Proc. Natl. Acad. Sci. U.S.A.* 2011;108(41):16938-16943.
10. Dou, J., *et al.* De novo design of a fluorescence-activating  $\beta$ -barrel. *Nature* 2018;561(7724):485.
11. Dovidchenko, N.V. and Galzitskaya, O.V. Computational approaches to identification of aggregation sites and the mechanism of amyloid growth. In, *Lipids in Protein Misfolding*. Springer; 2015. p. 213-239.
12. Fernandez-Escamilla, A.-M., *et al.* Prediction of sequence-dependent and mutational effects on the aggregation of peptides and proteins. *Nat. Biotechnol.* 2004;22(10):1302.
13. Fleishman, S.J., *et al.* Computational design of proteins targeting the conserved stem region of influenza hemagglutinin. *Science* 2011;332(6031):816-821.
14. Fujiwara, K., Toda, H. and Ikeguchi, M. Dependence of  $\alpha$ -helical and  $\beta$ -sheet amino acid propensities on the overall protein fold type. *BMC Struct. Biol.* 2012;12(1):18.
15. Gallardo, R., *et al.* De novo design of a biologically active amyloid. *Science* 2016;354(6313):4949.
16. Giorgetti, S., *et al.* Targeting amyloid aggregation: an overview of strategies and mechanisms. *Int. J. Mol. Sci.* 2018;19(9):2677.
17. Gordon, D.J., *et al.* Increasing the amphiphilicity of an amyloidogenic peptide changes the  $\beta$ -sheet structure in the fibrils from antiparallel to parallel. *Biophys. J.* 2004;86(1):428-434.
18. Huang, P.-S., Boyken, S.E. and Baker, D. The coming of age of de novo protein design. *Nature* 2016;537(7620):320.
19. Hubbard, T.J. Use of  $\beta$ -strand Interaction Pseudo-Potentials in Protein Structure Prediction and Modeling. In, *Proceedings of the Hawaii International Conference on System Sciences*. 1994. p. 336-336.
20. Hutchinson, E.G., *et al.* Determinants of strand register in antiparallel  $\beta$ -sheets of proteins. *Protein Sci.* 1998;7(11):2287-2300.
21. Ito, M., Ozawa, T. and Takada, S. Folding Coupled with Assembly in Split Green Fluorescent Proteins Studied by Structure-based Molecular Simulations. *J. Phys. Chem. B* 2013;117(42):13212-13218.
22. Jang, H., *et al.* Computational methods for structural and functional studies of Alzheimer's amyloid ion channels. In, *Protein Amyloid Aggregation*. Springer; 2016. p. 251-268.
23. Kim, C., *et al.* NetCSSP: web application for predicting chameleon sequences and amyloid fibril formation. *Nucleic Acids Res.* 2009;37(suppl\_2):469.
24. Kim, Y.E., *et al.* Green fluorescent protein nanopolygons as monodisperse supramolecular assemblies of functional proteins with defined valency. *Nat. Commun.* 2015;6:7134.
25. Köker, T., Fernandez, A. and Pinaud, F. Characterization of Split Fluorescent Protein Variants and Quantitative Analyses of Their Self-Assembly Process. *Sci. Rep.* 2018;8(1):5344.

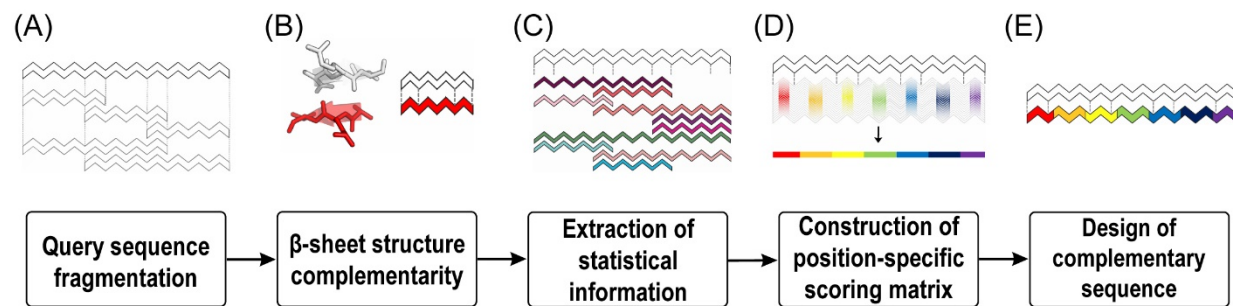


26. Kortemme, T., Ramírez-Alvarado, M. and Serrano, L. Design of a 20-amino acid, three-stranded  $\beta$ -sheet protein. *Science* 1998;281(5374):253-256.
27. Kuhlman, B., *et al.* Design of a novel globular protein fold with atomic-level accuracy. *Science* 2003;302(5649):1364-1368.
28. Liu, C., *et al.* Out-of-register  $\beta$ -sheets suggest a pathway to toxic amyloid aggregates. *Proc. Natl. Acad. Sci.* 2012;109(51):20913-20918.
29. Mandel-Gutfreund, Y., Zaremba, S.M. and Gregoret, L.M. Contributions of residue pairing to  $\beta$ -sheet formation: conservation and covariation of amino acid residue pairs on antiparallel  $\beta$ -strands. *J. Mol. Biol.* 2001;305(5):1145-1159.
30. Marcos, E., *et al.* De novo design of a non-local  $\beta$ -sheet protein with high stability and accuracy. *Nat. Struct. Mol. Biol.* 2018;25(11):1028.
31. Marcos, E. and Silva, D.-A. Essentials of de novo protein design: Methods and applications. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* 2018;8(6):1374.
32. Matthes, D., *et al.* Spontaneous aggregation of the insulin-derived steric zipper peptide VEALYL results in different aggregation forms with common features. *J. Mol. Biol.* 2014;426(2):362-376.
33. Maurer-Stroh, S., *et al.* Exploring the sequence determinants of amyloid structure using position-specific scoring matrices. *Nat. Methods* 2010;7(3):237.
34. McCoy, K. and Douglas, T. In Vivo Packaging of Protein Cargo Inside of Virus-Like Particle P22. In, *Virus-Derived Nanoparticles for Advanced Technologies*. Springer; 2018. p. 295-302.
35. Miller, K.E., *et al.* Bimolecular fluorescence complementation (BiFC) analysis: advances and recent applications for genome-wide interaction studies. *J. Mol. Biol.* 2015;427(11):2039-2055.
36. Minor Jr, D.L. and Kim, P.S. Measurement of the  $\beta$ -sheet-forming propensities of amino acids. *Nature* 1994;367(6464):660.
37. Mizuguchi, K., *et al.* HOMSTRAD: a database of protein structure alignments for homologous families. *Protein Sci.* 1998;7(11):2469-2471.
38. Quinn, T.P., *et al.* Betadoublet: de novo design, synthesis, and characterization of a beta-sandwich protein. *Proc. Natl. Acad. Sci.* 1994;91(19):8747-8751.
39. Richardson, J.S. and Richardson, D.C. Natural  $\beta$ -sheet proteins use negative design to avoid edge-to-edge aggregation. *Proc. Natl. Acad. Sci. U.S.A.* 2002;99(5):2754-2759.
40. Rocklin, G.J., *et al.* Global analysis of protein folding using massively parallel design, synthesis, and testing. *Science* 2017;357(6347):168-175.
41. Schymkowitz, J.W., *et al.* Prediction of water and metal binding sites and their affinities by using the Fold-X force field. *Proc. Natl. Acad. Sci.* 2005;102(29):10147-10152.
42. Shammass, S.L., *et al.* Perturbation of the stability of amyloid fibrils through alteration of electrostatic interactions. *Biophys. J.* 2011;100(11):2783-2791.
43. Sormanni, P., Aprile, F.A. and Vendruscolo, M. Rational design of antibodies targeting specific epitopes within intrinsically disordered proteins. *Proc. Natl. Acad. Sci.* 2015;112(32):9902-9907.
44. Steward, R.E. and Thornton, J.M. Prediction of strand pairing in antiparallel and parallel  $\beta$ -sheets using information theory. *Proteins* 2002;48(2):178-191.
45. Stranges, P.B., *et al.* Computational design of a symmetric homodimer using  $\beta$ -strand assembly. *Proc. Natl. Acad. Sci.* 2011;108(51):20562-20567.

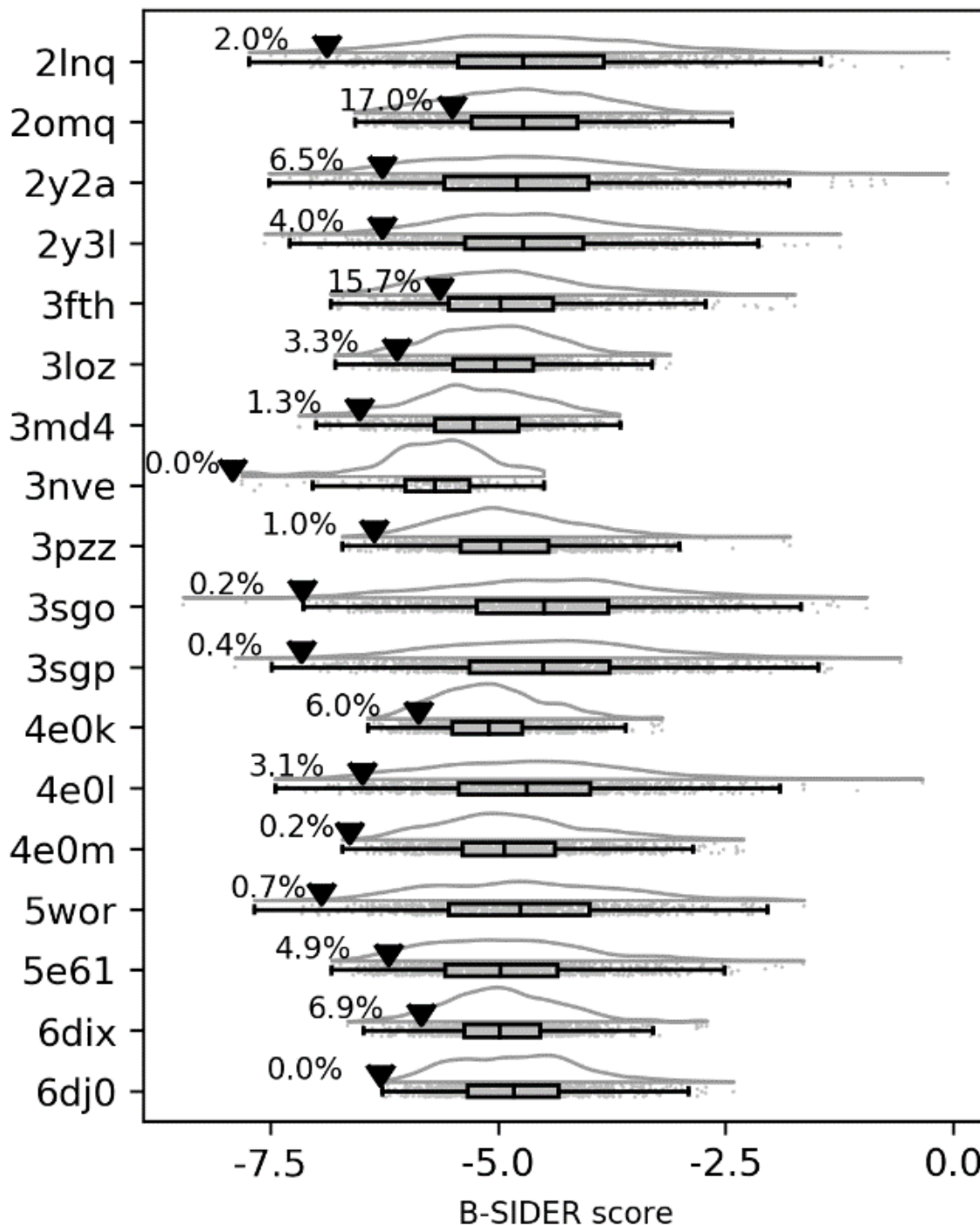


46. Tartaglia, G.G. and Vendruscolo, M. The Zyggregator method for predicting protein aggregation propensities. *Chem. Soc. Rev.* 2008;37(7):1395-1401.
47. Trovato, A., *et al.* Insight into the structure of amyloid fibrils from the analysis of globular proteins. *PLOS Comput. Biol.* 2006;2(12):170.
48. Tsolis, A.C., *et al.* A consensus method for the prediction of 'aggregation-prone' peptides in globular proteins. *PLoS One* 2013;8(1):54175.
49. Tyka, M.D., *et al.* Alternate states of proteins revealed by detailed energy landscape mapping. *J. Mol. Biol.* 2011;405(2):607-618.
50. Walsh, I., *et al.* PASTA 2.0: an improved server for protein aggregation prediction. *Nucleic Acids Res.* 2014;42(W1):301.
51. Wang, W. and Hecht, M.H. Rationally designed mutations convert de novo amyloid-like fibrils into monomeric  $\beta$ -sheet proteins. *Proc. Natl. Acad. Sci.* 2002;99(5):2760-2765.
52. West, M.W., *et al.* De novo amyloid proteins from designed combinatorial libraries. *Proc. Natl. Acad. Sci.* 1999;96(20):11211-11216.
53. Xu, Y., *et al.* Conformational transition of amyloid  $\beta$ -peptide. *Proc. Natl. Acad. Sci.* 2005;102(15):5403-5407.
54. Zaremba, S.M. and Gregoret, L.M. Context-dependence of Amino Acid Residue Pairing in Antiparallel  $\beta$ -Sheets. *J. Mol. Biol.* 1999;291(2):463-479.
55. Zhang, Y. and Skolnick, J. TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids Res.* 2005;33(7):2302-2309.
56. Zheng, W., *et al.* Exploring the aggregation free energy landscape of the amyloid- $\beta$  protein (1-40). *Proc. Natl. Acad. Sci.* 2016;113(42):11835-11840.
57. Zibae, S., *et al.* A simple algorithm locates  $\beta$ -strands in the amyloid fibril core of  $\alpha$ -synuclein, A $\beta$ , and tau using the amino acid sequence alone. *Protein Sci.* 2007;16(5):906-918.

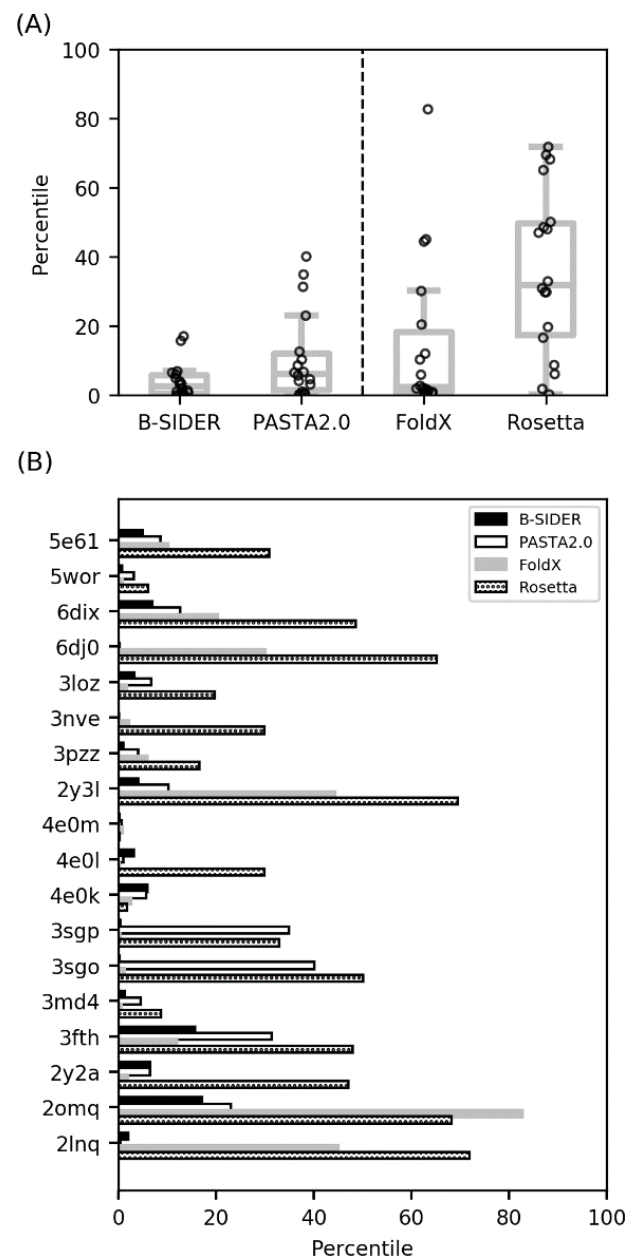
**Fig. 1. Overview of the B-SIDER algorithm** (A) When a query sequence is given, the query is divided into smaller linear peptides ranging from 3 residues to its full length, and exactly matched sequences are identified. (B) The matched sequences are checked against the structure database so that the matches indeed form  $\beta$ -sheets. (C) If the matches form  $\beta$ -sheets, their complementary sequences are extracted. (D) A position-specific score matrix is constructed based on the complementary sequence information. (E) Final complementary sequences are designed using the score matrix.



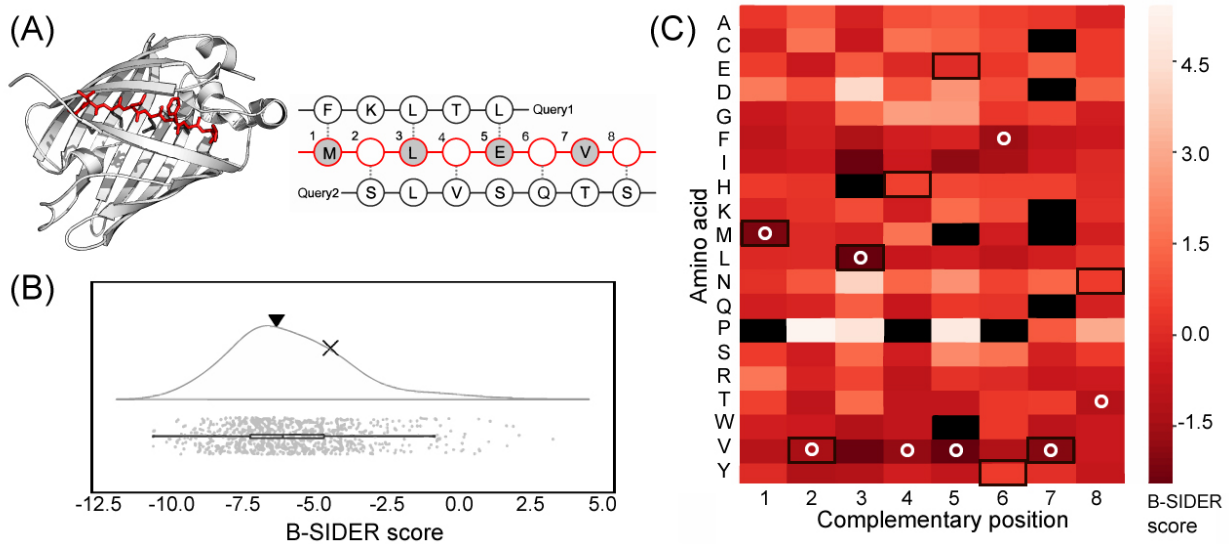
**Fig. 2. The predictive power of the B-SIDER score.** In this test set, native sequences were compared against 1,000 random sequences. The lower, the more favorable. The native complementarity scores are marked as (▼) and their percentile values are displayed next to the marks. This plot was generated using the Raincloud Python package (Allen, et al., 2018)



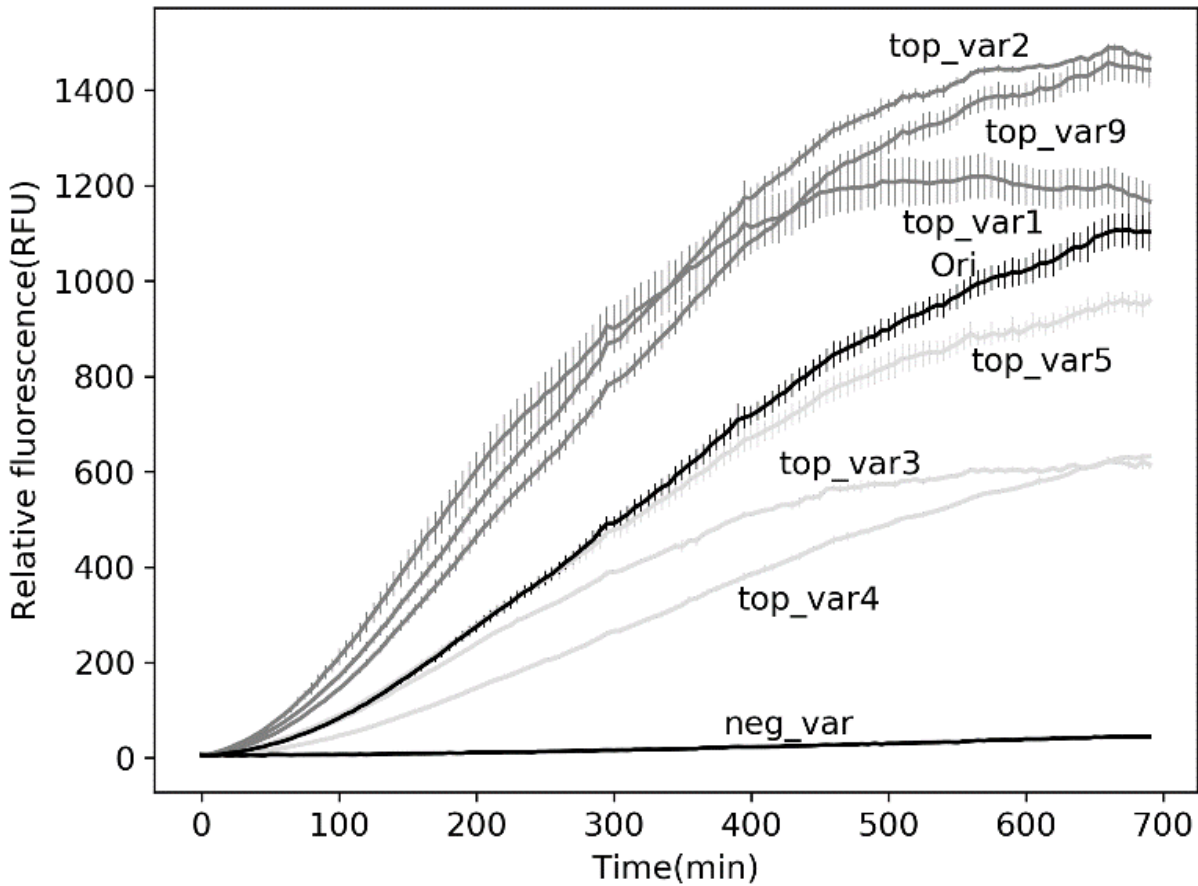
**Fig. 3. The comparison of the B-SIDER score with the structure-based energy scores and existing sequence-based method.** The B-SIDER score is compared against two popularly used structure-based protein design methods (Rosetta and FoldX) and the scoring matrix for the prediction of amyloid formation (PASTA2.0). For a fair comparison, we generated new sets of 1,000 random sequences per target and assessed their scores using each energy function. (A) The B-SIDER score gave the most consistent assessment. (B) The percentile values of native complementary sequences against the score distributions of the randomly generated sequences are presented. Overall, the sequence-based approaches showed higher accuracies. The better performance of FoldX than Rosetta may indicate that the fixed-backbone design strategy may be more useful in the design and prediction of  $\beta$ -sheet complementarity.



**Fig. 4. The design process of GFP11 variants.** (A) The strand to design (GFP11) is highlighted in red and query sequences to consider are shown on the right. Dotted lines represent hydrogen bonds. The residues pointing inward, which are not subject to mutation, are colored in gray. (B) The B-SIDER score of the original GFP11 (▼) is compared against a score distribution of randomly generated sequences. The negative variant (neg\_var) is marked as (×, the 77<sup>th</sup> percentile). (C) The position-specific scoring matrix for the query sequences. Amino acids which never appeared in each position are colored in black. The amino acids of the native sequence are highlighted with black boxes. The white circles represent the lowest scores at each position.



**Fig. 5. Complementation assay with the designed GFP11 variants.** Among the six successfully expressed variants, three variants exhibited stronger fluorescence emissions than the original peptide identified by directed mutagenesis.



**Table 1. Retrospective test set**

Source	PDB ID
Amyloid- $\beta$	2lnq, 2y2a, 2y3l, 3pzz
Insulin	2omq
IAPP(amylin)	3fth, 5e6l
Prion	3md4, 3nve
Tau	4e0m
Alpha B crystalline	3sgo, 4sgp
$\beta$ -2 microglobulin	3loz, 4e0k, 4e0l
Immunoglobulin	6dj0, 6dix
SOD1	5wor

**Table 2. GFP 11 variants**

Variant	Sequence	Score	Relative Assembly*
top_var1	M <u>V</u> L <u>V</u> E <u>F</u> V <u>T</u>	-12.34	1.17
top_var2	M <u>Y</u> L <u>V</u> E <u>F</u> V <u>T</u>	-12.14	1.4
top_var9	M <u>T</u> L <u>V</u> E <u>F</u> V <u>T</u>	-11.79	1.36
top_var3	M <u>V</u> L <u>V</u> E <u>I</u> V <u>T</u>	-12.11	0.59
top_var4	M <u>V</u> L <u>V</u> E <u>F</u> V <u>Y</u>	-11.95	0.56
top_var5	M <u>Y</u> L <u>V</u> E <u>I</u> V <u>T</u>	-11.92	0.88
top_var6	M <u>V</u> L <u>V</u> E <u>V</u> V <u>T</u>	-11.90	N.D.
top_var7	M <u>V</u> L <u>V</u> E <u>F</u> V <u>V</u>	-11.89	N.D.
top_var8	M <u>V</u> L <u>V</u> E <u>F</u> V <u>W</u>	-11.83	N.D.
top_var10	M <u>V</u> L <u>V</u> E <u>F</u> V <u>F</u>	-11.77	N.D.
neg_var	M <u>V</u> L <u>G</u> E <u>K</u> V <u>E</u>	-4.60	0.04
Ori	MVLHEYVN	-6.50	1

\*Normalized values compared to the fluorescence level of the original GFP11 strand.