1 ## Tittle

2 Rosetta FunFolDes – a general framework for the computational design of

3 functional proteins

4

5 ## Authors

6

7 Jaume Bonet*, Sarah Wehrle*, Karen Schriever*, Che Yang*, Anne Billet, Fabian

8 Sesterhenn, Andreas Scheck, Freyr Sverrisson, Sabrina Vollers, Roxanne

9 Lourman, Melanie Villard, Stéphane Rosset, Bruno E. Correia.

10

11 * These authors contributed equally to this work.

12

13 ## Affiliations

14 Institute of Bioengineering, École Polytechnique Fédérale de Lausanne,

15 Lausanne, CH-1015, Switzerland

16 Swiss Institute of Bioinformatics (SIB), Lausanne, CH-1015, Switzerland

17

18

19 ## Abstract

20

21 The robust computational design of functional proteins has the potential to

22 deeply impact translational research and broaden our understanding of the

23 determinants of protein function, nevertheless, it remains a challenge for state-

24 of-the-art methodologies. Here, we present a computational design approach

25 that couples conformational folding with sequence design to embed functional

26 motifs into heterologous proteins. We performed extensive benchmarks, where

27 the most unexpected finding was that the design of function into proteins may

28 not necessarily reside in the global minimum of the energetic landscape, which

29 could have important implications in the field. We have computationally

30 designed and experimentally characterized a distant structural template and a *de*

31 *novo* "functionless" fold, two prototypical design challenges, to present

32    important viral epitopes. Overall, we present an accessible strategy to repurpose

33    old protein folds for new functions, which may lead to important improvements

34    on the computational design of functional proteins.

35

36

## Introduction

37

38

39 Proteins are one of the main functional building blocks of the cell. The ability to

40 create novel proteins outside of the natural realm has opened the path towards

41 innovative achievements, such as new protein pathways (Cross et al., 2017),

42 cellular functions (Joh et al., 2014), and therapeutic leads (Correia et al., 2010;

43 Correia et al., 2014; Kulkarni et al., 2015). Computational protein design is the

44 rational and structure-based approach to solve the inverse folding problem, i.e.

45 the search for the best putative sequence capable of fitting and stabilizing a given

46 protein's three-dimensional conformation (Coluzza, 2017). As such, a great deal

47 of effort has been placed into the understanding of the rules of protein folding

48 and stability (Koga et al., 2012; Marcos et al., 2017) and its relation to the

49 appropriate sequence space (Kuhlman & Baker, 2000).

50

51 Computational protein design focus on two main axes of search related to the

52 structural and sequence spaces that are explored. Fixed backbone approaches

53 work with a static protein backbone conformation, which greatly constrains the

54 sequence space that is explored by the computational algorithm (Kuhlman &

55 Baker, 2000). Following the same principles as naturally occurring homologs,

56 which often exhibit a certain degree of structural diversity, flexible backbone

57 approaches enhance the sequence diversity, adding the challenge of identifying

58 energetically favorable sequence variants that are correctly coupled to the

59 structural perturbations (Murphy et al., 2012).

60

61 Another variation for computational design approaches is *de novo* design, in

62 which protein backbones are assembled *in silico*, followed by sequence

63 optimization to fold into a pre-defined three-dimensional conformation without

64 being constrained by previous sequence information (Hill, Raleigh, Lombardi, &

65 DeGrado, 2000). This approach tests our understanding of the rules governing

66 the structure of different protein folds. The failures and successes of this

67 approach confirm and correct the principles used for the protein design process

68 (Koga et al., 2012; Marcos et al., 2017).

69

70    One of the main aims of computational protein design is the rational design of
71    functional proteins capable of carrying existing or novel functions into new
72    structural contexts (Street & Mayo, 1999). One can broadly classify three main
73    approaches for the design of functional proteins: redesigning of pre-existing
74    functions, grafting of functional sites onto heterologous proteins, and designing
75    of novel functions not found in the protein repertoire. The redesign of a pre-
76    existing function to alter its catalytic activity (Yu et al., 2014) or improve its
77    binding target recognition (Guntas, Purbeck, & Kuhlman, 2010) can be
78    considered the most conservative approach; as it is typically accomplished by
79    point mutations around the functional area of interest, it tends to have little
80    impact on global structure and stability of the designed protein. On the other
81    hand, the design of fully novel functions has most noticeably been achieved by
82    applying chemical principles that tested our fundamental knowledge of enzyme
83    catalysis (Jiang et al., 2008; Kries, Blomberg, & Hilvert, 2013).

84

85    Between these two approaches resides protein grafting. This method aims to
86    repurpose natural folds as carriers for exogenous known functions. It relies on
87    the strong relationship between protein structure and activity, to allocate a given
88    functionality from one protein to another by means of transferring the structural
89    motif responsible for the function (Azoitei et al., 2011; Correia et al., 2011;
90    Correia et al., 2010; Correia et al., 2014; Kulkarni et al., 2015; Procko et al., 2014;
91    Viana et al., 2013). The most successful grafting approaches are highly
92    dependent on structural similarity between the functional motif and the
93    insertion region in the protein scaffold. When the functional motif and the
94    insertion region are almost identical in backbone conformation, functional
95    transfer can be performed by side-chain grafting, i.e. mutating the target
96    residues into those of the functional motif (Correia et al., 2010; Kulkarni et al.,
97    2015). In much more challenging scenarios, full backbone grafting may be used
98    in conjunction with directed evolution to make the structure fully compatible
99    with the new function (Azoitei et al., 2011). Nevertheless, motif transfer is
100   limited between very similar structural regions, which greatly constrains the
101   subset of putative scaffolds that can be used for this purpose, especially as the
102   structural complexity of the functional motif grows.

103    Previously, we have demonstrated the possibility of expanding protein grafting
104    to scaffolds with segments that have low structural similarity . To accomplish
105    that task, we developed a prototype protocol named Rosetta Fold From Loops
106    (FFL) (Correia et al., 2014; Procko et al., 2014).

107

108    The distinctive feature of our protocol is the coupling of the folding and design
109    stages to bias the sampling towards structural conformations and sequences that
110    stabilize the grafted functional motif. In the past, FFL was used to obtain designs
111    that were functional (synthetic immunogens (Correia et al., 2014) and protein-
112    based inhibitors (Procko et al., 2014)) and where the experimentally determined
113    crystal structures closely resembled the computational models; however, the
114    structures of the functional sites were still very close to the insertion segments
115    of the hosting scaffolds.

116

117    Here, we present a complete re-implementation of the FFL protocol with
118    enhanced functionalities, simplified user interface and complete integration with
119    any other available Rosetta protocols. We have called this new, more generalist
120    protocol Rosetta Functional Folding and Design (FunFolDes), we have
121    benchmarked it in a number of scenarios providing important technical details
122    to better exploit and expand the capabilities of the original protocol.
123    Furthermore, we challenged FunFolDes with two design tasks to probe the
124    boundaries of applicability of the protocol. The design tasks were centered on
125    using distant structural template as hosting scaffold and functionalizing a *de*
126    *novo* designed protein – in both challenges, FunFolDes succeeded in
127    functionalizing the designed proteins. These results are encouraging and provide
128    a solid basis for the broad applicability of FunFolDes as a strategy for the robust
129    computational design of functionalized proteins.

130

131

132 ## Results

133

134 ### *Rosetta FunFolDes – a computational framework for design of functional*
135 *proteins*

136

137 The original prototype of the Rosetta Fold From Loops (FFL) protocol was
138 successfully used to transplant the structural motif of the Respiratory Syncytial
139 Virus (RSV) protein F site II neutralizing epitope into a protein scaffold in the
140 context of a vaccine design application (Correia et al., 2014).
141 FFL enabled the insertion and conformational stabilization of the structural
142 motif into a defined protein topology by using Rosetta's fragment insertion
143 machinery to fold the polypeptide chain to adopt the desired topology (Rohl,
144 Strauss, Misura, & Baker, 2004) which was then sequence designed. Information
145 content from the scaffold structure was used to guide the folding, ensuring an
146 overall similar topology while allowing for the conformational changes needed
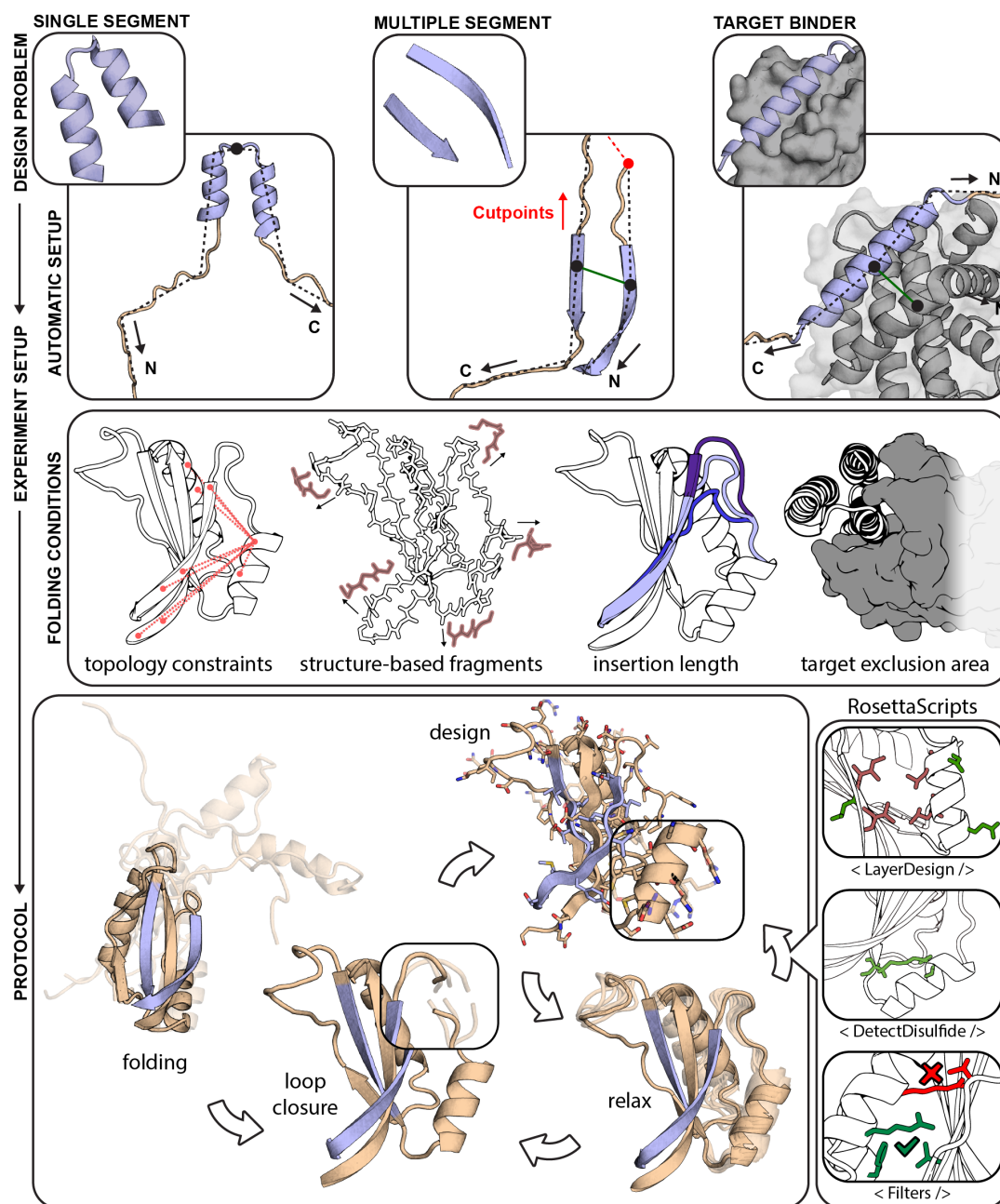147 to stabilize the inserted structural motif.

148

149 The final implementation of the protocol, referred to as FunFolDes, is
150 schematically represented in **Figure 1**, and fully described in Materials and
151 Methods. Our upgrades to FFL focused on three main aims: I) improve the
152 applicability of the system to allow handling of more complex structural motifs;
153 II) enhance the design of functional proteins by including binding partners in the
154 simulations; III) offer a higher degree of control over each stage of the simulation
155 while improving the usability for non-experts. These three aims were achieved
156 through the implementation of five core technical improvements described
157 below.

158

159 *Insertion of multi-segment functional sites.* The initial implementation of FFL was
160 limited to the insertion of a single-segment structural motif, which was sufficient
161 to demonstrate its potential at the time (Correia et al., 2014; Procko et al., 2014).
162 However, most functional sites in proteins typically entail, at the structural level,
163 multiple discontinuous segments; which is the case for protein-protein interfaces
164 or enzyme active-sites, among others (Aragues, Sali, Bonet, Marti-Renom, &

165    Oliva, 2007; Richter, Leaver-Fay, Khare, Bjelic, & Baker, 2011). FunFolDes can

166    now handle functional sites with any number of discontinuous segments,

167    ensuring the native orientations of each of the segments. Furthermore, it allows

168    control of the backbone flexibility of each of the insertion points and the order in

169    the protein scaffold sequence in which each segment is inserted. Finally, the

170    sequence length between the motif and the insertion region is not required to be

171    the same, allowing the user to search for protein scaffolds using alternative

172    metrics to the full backbone RMSDs between the motif and the protein scaffold

173    (Azoitei et al., 2011; Correia et al., 2010). These new features essentially allow

174    the replacement between completely different structural segments. Thus, they

175    greatly enhance the types of structural motifs that can be targeted with

176    FunFolDes, widening the applicability of the computational protocol.

177

178

**Figure 1. Rosetta FunFolDes - method overview.** FunFolDes was devised to tackle a wide range of functional protein design problems, combining a higher user control of the simulation parameters whilst simultaneously lowering the level of expertise required. FunFolDes is able to transfer single- and multi-segment motifs together with the target partner by exploiting Rosetta's FoldTree framework (top row). A wider range of information can be extracted from the template to shift the final conformation towards a more productive design space (middle row), including targeted distance constraints, generation of structure-based fragments, motif insertion in sites with different residue length and presence of the binding target to bias the folding stage. The bottom row showcases the most typical application of the FunFolDes protocol. Integration in RosettaScripts allows to tailor FunFolDes behavior and for a seamless integration with other

189   protocols, and complex selection logics can be added to address the different complexities in
190   each design task.

191

192   *Structural folding and sequence design in the presence of a binding partner.* Many
193   of the functional roles of proteins in cells require physical interaction with other
194   proteins, nucleic acids, or metabolites (Garcia-Garcia et al., 2012). Several
195   proposed mechanisms to regulate binding affinities and specificities in protein
196   interactions involve protein flexibility, such as induced fit and conformational
197   selection (Chakrabarti et al., 2016; Lange et al., 2008). Inspired by these
198   naturally occurring mechanisms, we devised a strategy to fold and design in the
199   presence of the desired binding partner. Including the binder in simulations has
200   a twofold benefit. On the one hand, is a way of explicitly representing functional
201   constraints to bias the designed protein towards a functional sequence space,
202   resolving putative clashes derived from the template scaffold and, thus,
203   significantly enlarging the number of usable templates. On the other hand, this
204   approach facilitates the design of new additional contact residues (outside of the
205   motif) that may afford enhanced affinity and/or specificity. Here, we tested
206   FunFolDes in a model system for which extensive experimental data has been
207   collected, and we show how this approach improves the sampling of productive
208   conformational and sequence space.

209

210   *Region-specific structural constraints.* FFL could exploit distance constraints from
211   the target scaffold to guide the folding stage. A simplified solution was
212   implemented in FFL with two possible simulation modes, where either
213   constraints are collected throughout the protein scaffold, or folding is
214   unconstrained. Currently, FunFolDes can collect from full-template to region-
215   specific constraints, allowing greater levels of flexibility in areas of the scaffold
216   that can be critical for function (e.g. segments close to the interface of a target
217   protein) and improving the sampling of conformations which otherwise could be
218   missed or highly underrepresented. Furthermore, FunFolDes is no longer limited
219   to atom-pair distance constraints (Rohl & Baker, 2002) and can incorporate
220   other types of kinematic constraints, such as angle and dihedral constraints

221  (Bowers, Strauss, & Baker, 2000), which have been used to improve success
222  rates when folding scaffolds rich in beta-strands (Marcos et al., 2017).

223

224  *On-the-fly fragment picking.* Fragment insertion is a core algorithm in Rosetta
225  protocols exploring high degrees of freedom of the polypeptide chain, such as *ab*
226  *initio* protein prediction (Simons, Ruczinski, et al., 1999), loop modeling (Stein &
227  Kortemme, 2013), or more recently, FFL (Correia et al., 2014). Classically,
228  fragment libraries are generated through sequence-based predictions of
229  secondary structure and dihedral angles (Bowers et al., 2000). This information
230  is used in a Rosetta application to obtain three- and nine residue-long fragment
231  libraries from naturally occurring proteins, which are then provided to the
232  downstream protocols. Leveraging internal functionalities in Rosetta, FunFolDes
233  can assemble fragment sets automatically. Due to its particularities, secondary
234  structure, dihedral angles, and accessible solvent area can be automatically
235  computed from the protein scaffold's structure. Although sequence-based
236  fragments can still be provided, this removes the need for secondary applications
237  in the protocol pipeline, boosting the usability of FunFolDes by lowering the
238  barrier for non-experts. It also enables the assembly of protocols in which the
239  fragment sets are mutable along the procedure. The benchmark presented in this
240  paper evaluates the performance of such functionality.

241

242  *Compatibility with other Rosetta modules.* Finally, FunFolDes is compatible with
243  Rosetta's modular xml-interface: Rosetta Scripts (RS) (Fleishman et al., 2011).
244  This enables customization of the FunFolDes protocol and, more importantly,
245  connection with other protocols and filters available through the RS interface. In
246  order to obtain a full integration with this interface, the FunFolDes protocol is
247  divided in multiple *Movers* (i.e. modules capable of altering the information
248  content of a structure).

249

250  We devised two benchmark scenarios to test the performance of FunFolDes. One
251  of these aimed to capture conformational changes in small protein domains
252  caused by sequence insertions or deletions, and the second scenario assessed

10

253    protocol performance to fold and design a binder in the presence of the target

254    binding partner.

255

256    ***Capturing conformational and sequence changes in small protein domains***
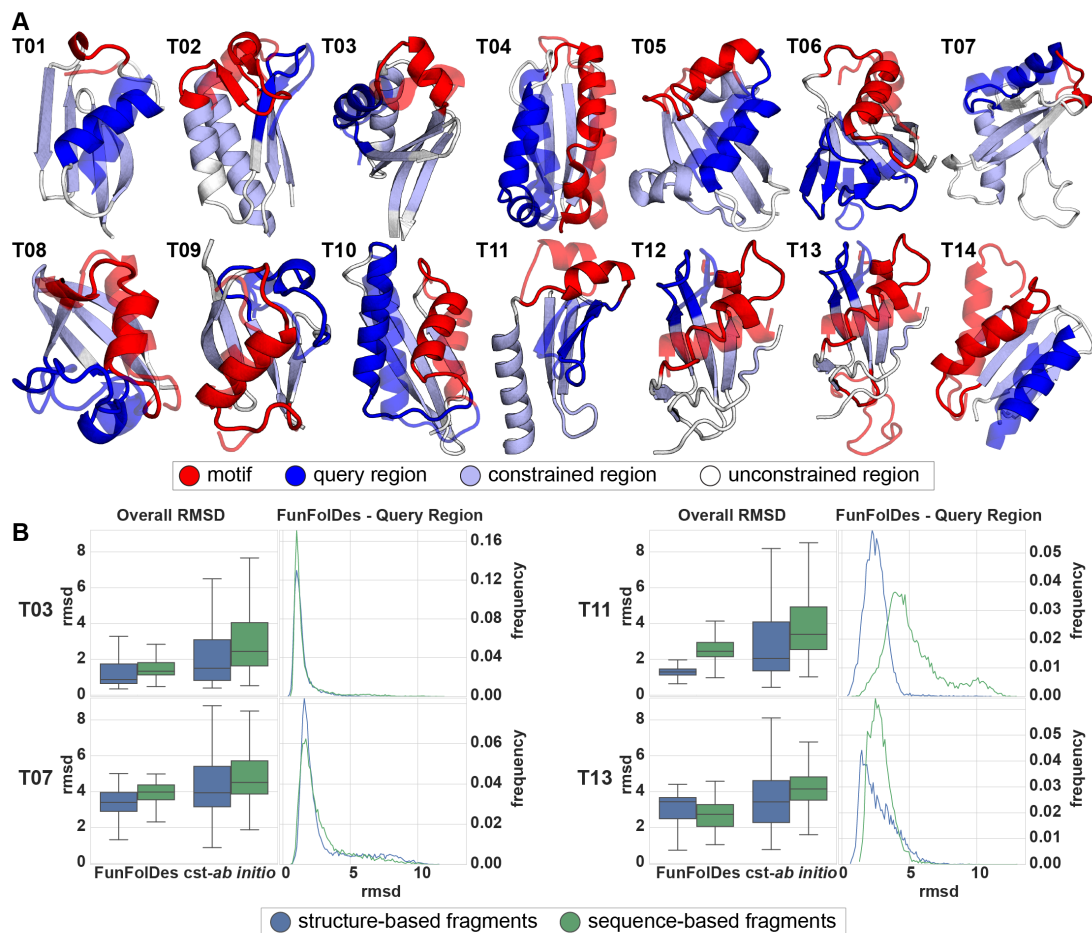
257

258    Typical protein design benchmarks are assembled by stripping native side

259    chains from known protein structures and evaluating the sequence recovery of

260    the design algorithm (Kuhlman & Baker, 2000). The main design aim of

261    FunFolDes is to insert structural motifs into protein folds while allowing

262    flexibility across the overall structure. This conformational freedom allows the

263    full protein scaffold to adapt and stabilize the functional motif's conformation.

264    This is a main distinctive point from other approaches to design functional

265    proteins that rely on a mostly rigid scaffold (Azoitei et al., 2011; Correia et al.,

266    2010; Fallas et al., 2017; Hill et al., 2000; Joh et al., 2014; Richter et al., 2011). For

267    many modeling problems, such as protein structure prediction, protein-protein

268    and protein-ligand docking, or protein design, standardized benchmark datasets

269    are available (Vreven et al., 2015) or easily accessible. Devising a benchmark for

270    designed proteins with propagating conformational changes across the structure

271    is challenging, as we are assessing both structural accuracy as well as sequence

272    recovery of the protocol.

273

274    To address this problem, we analyzed structural domains found repeatedly in

275    natural proteins and clustered them according to their definition in the CATH

276    database (Dawson et al., 2017). As a result, we were able to select a set of 14

277    benchmark targets labeled T01 through T14 (**Figure 2A**). A detailed description

278    on the construction of the benchmark can be found in the Materials and Methods

279    section.

280

11

**Figure 2. Benchmark test set to evaluate FunFolDes structural sampling.** A) Structural representation of the 14 targets used in the benchmark. Each target highlights the motif (red), the query region (blue), and the positions from which distance constraints were generated (light blue). Conformations of the motif and query regions, as found in the template structures, appear superimposed in a semi-transparent depiction. B) Full structure RMSD (Overall RMSD) and local RMSD for the query region (FunFolDes – Query Region) for four targets (full dataset presented in Supplementary Figure S1). Overall RMSD compares results for the two simulation modes (FunFolDes Vs. constrained–*ab initio (cst-ab initio)*) and the two fragment generation methods (structure-(blue) Vs. sequence-based fragments(green)). FunFolDes more frequently samples RMSDs closer to the conformation of the target structure. Generally, structure-based fragment also contribute to lower mean overall RMSDs. The FunFolDes – Query Region RMSD distributions show that the two fragments sets do not have a major importance in the structural recovery of the query region.

Briefly, for the benchmark we selected proteins with less than 100 residues, where each benchmark test case is composed of two proteins of the same CATH domain cluster. One of the proteins is dubbed template, and serves as a structural representative of the CATH domain. The second protein, dubbed

12

300  target, contains structural insertions or deletions (motif), to which a structural

301  change in a different segment of the same structure could be attributed (query

302  region). The motif and query regions for all the targets are shown in **Figure 2A**

303  and quantified in terms of the percentage of overall secondary structure in

304  **Figure 2 - Supplementary Figure 1A**. To a great extent, these structural

305  changes due to natural sequence insertions and deletions are analogous to those

306  occurring in the design scenarios for which FunFolDes was conceived.

307

308  Using FunFolDes, we folded and designed the target proteins while maintaining

309  the motif segment structurally fixed, mimicking a structural motif insertion.

310  Distance constraints between residues were extracted from the template in the

311  regions of shared structural elements of the template and the target, and were
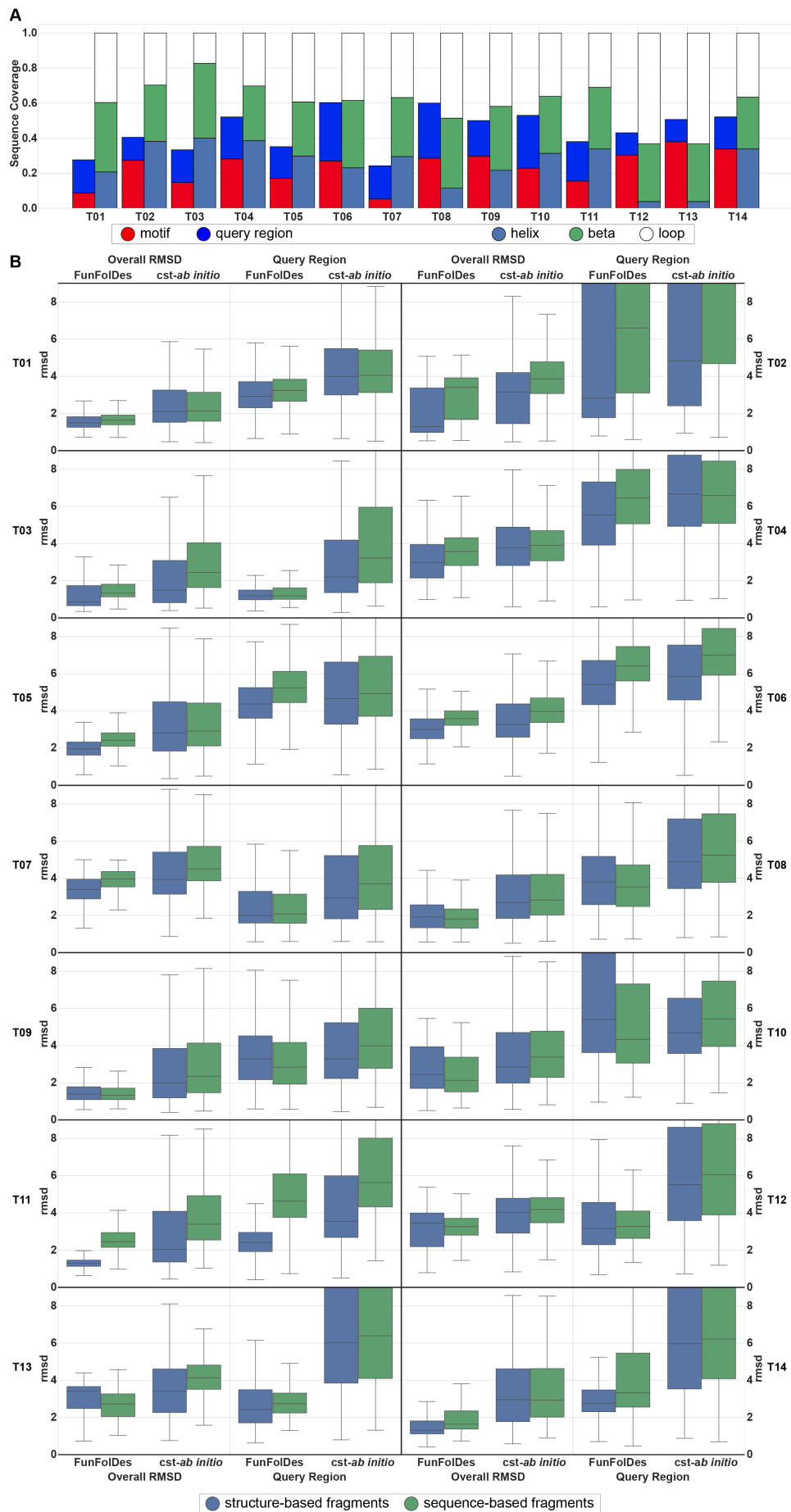
312  used to guide the folding simulations.

313

314  To check whether FunFolDes enhances sequence and structural sampling, we

315  compared the simulations to constrained *ab initio (*cst-*ab initio) simulations*

316  (Bowers et al., 2000). These simulations were performed using the same sets of

317  constraints but without the motif region as a static segment.

318

319  As Rosetta conformational sampling is highly dependent upon the fragment set

320  provided, in this benchmark we also tested the influence of structure- and

321  sequence-based fragments. The performance of the two protocols was broadly

322  analyzed by the global and local recovery of both structure and sequence.

323

324  Structural recovery was assessed through two main metrics: (a) global RMSD of

325  the full decoys against the target and (b) local RMSD of the query region. When

326  evaluating the distributions for global RMSD in the designed ensembles,

327  FunFolDes outperformed cst-*ab initio* by consistently producing populations of

328  decoys with lower mean (RMSDs mostly found below 5 Å), a result observed in

329  all 14 targets (**Figure 2B**, **Figure 2 – Supplementary Figure 1B**). This result is

330  especially reassuring considering that FFL simulations contain more structural

331  information of the target topology than the cst-*ab initio* simulations.

332

14

333 **Figure 2 - Supplementary Figure 1. Structural composition and overall results of the**
334 **benchmark targets.** A) Percentage of secondary structure type, motif and query region in the
335 overall structures. B) Full structure RMSD (Overall RMSD) and local RMSD for the query region
336 (Query Region) between the decoy populations and their respective targets. FunFolDes tends to
337 outperform cst–*ab initio* in all scenarios and the structure-based fragments yield decoy
338 population with lower mean RMSDs, albeit with small differences relative to the sequence-based
339 fragments.

340

341 Retrieval of the local RMSDs of the query unconstrained regions presented
342 mixed results across the benchmark set. In 13 targets, FunFolDes outperforms
343 cst-*ab initio*, showing lower mean RMSDs in the decoy population.
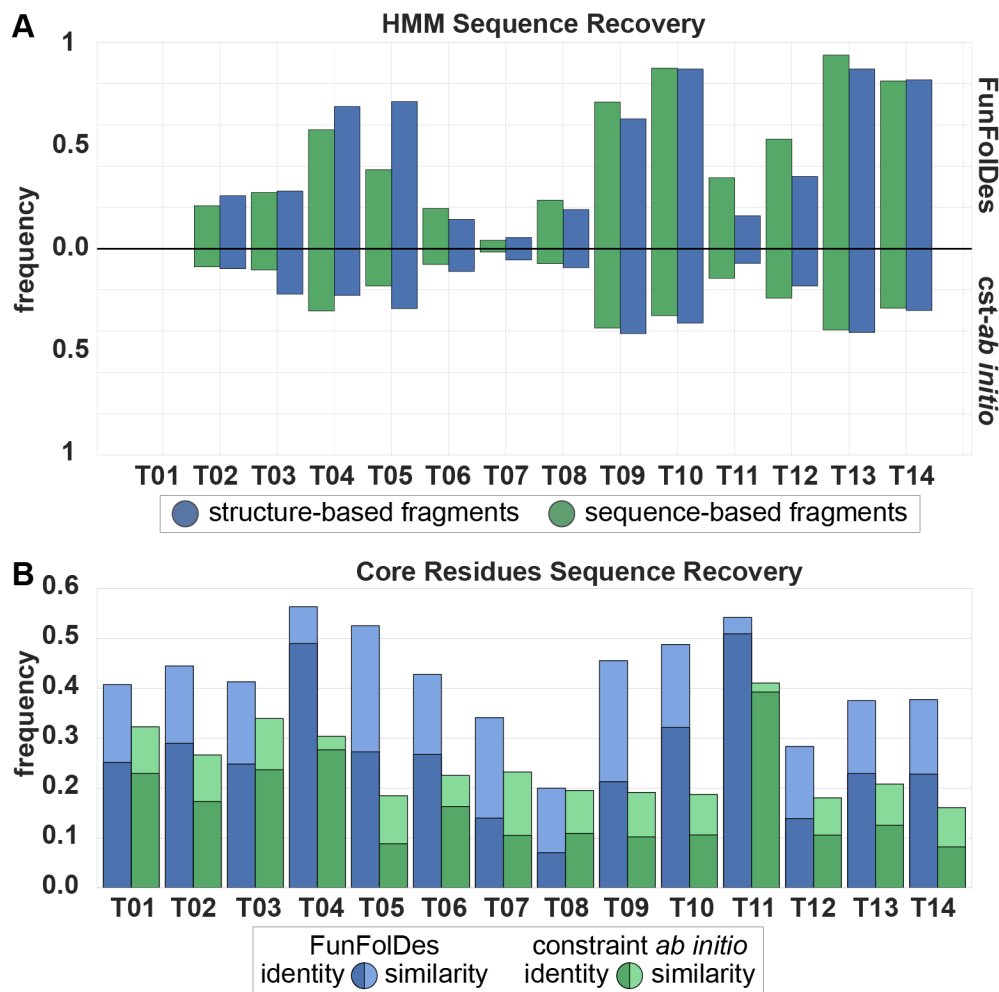
344

345 When comparing fragment sets (structure- vs sequence-based), both achieved
346 similar mean RMSDs in the decoy populations; nonetheless, the structure-based
347 fragments more often reach the lowest RMSDs for overall and query RMSDs
348 (**Figure 2B**, Fig**ure 2 – Supplementary Figure 1**). This is consistent with what
349 would be expected of structural information content within each set of
350 fragments. When paired with the technical simplicity of use, time-saving and
351 enhanced sampling of the desired topology, the structure-based fragments are an
352 added value for FunFolDes.

353

354 In addition to structural metrics, we also quantified sequence recovery in the
355 decoy populations, both in terms of sequence identity as well as sequence
356 similarity according to the BLOSUM62 matrix (Henikoff & Henikoff, 1992)
357 (**Figure 3A**). In all targets, the sequence identities and similarities were higher
358 for FunFolDes populations than for cst-*ab initio*, and in line with sequence
359 recoveries presented for other design protocols (Murphy et al., 2012) (**Figure
360 3A**). This type of metrics has been shown to be highly dependent on the exact
361 backbone conformation used as input (Kuhlman & Baker, 2000; Murphy et al.,
362 2012). Given that FunFolDes is exploring larger conformational spaces, as a
363 proxy for the quality of the sequences generated, we used the target's Hidden
364 Markov Models (HMM) (Eddy, 2011) and quantified how many of the designed
365 sequences were identified as belonging to the target's CATH superfamily
366 according to its HMM definition (**Figure 3B**).

367



**Figure 3. Assessment of FunFolDes' sequence sampling quality of.** A) HMM Sequence Recovery measures the percentage of decoys generated that can be assigned to the original HMM from the CATH superfamily that the target belongs to. FunFolDes consistently outperforms cst–*ab initio*, which is consistent with the same behavior observed in the structural recovery. B) Core Residues Sequence Recovery reveals the conservation of core residues between each design set and its target. Conservation is measured in terms of sequence identity and sequence similarity (as assigned through BLOSUM62). Also according to this metric FunFolDes outperforms cst-*ab initio* in every instance, reaching, for some populations, levels of conservation similar to those found in more restrained flexible-backbone designs.

HMM recovery was computed as the percentage of decoys with an E-value under 10 and covering more than 50% of the full decoy sequence. FunFolDes decoy populations systematically outperformed those from cst-*ab initio* (**Figure 3B**). The performance of the two fragment sets shows no significant differences. Core

16

383    sequence identity and similarity was assessed over the structure-based fragment

384    set.

385

386    In summary, the results of this benchmark highlight the usability of FunFolDes to

387    generate close-to-native scaffold proteins to stabilize inserted structural motifs.

388    FunFolDes aims to refit protein scaffolds towards the requirements of a

389    functional motif. In this perspective, it is critical to explore, within certain

390    topological boundaries, structural variations around the original template. This

391    benchmark points to several variables in the protocol that resulted in enhanced

392    structural and sequence sampling.

393

394    ***Target-biased folding and design of protein binders***

395    The computational design of proteins that can bind with high affinity and

396    specificity to targets of interest remains a largely unsolved problem (Schreiber &

397    Fleishman, 2013). Within FunFolDes' conceptual approach of coupling folding

398    with sequence design, we sought to add the structure of the binding target

399    (**Figure 1**) to attempt to bias sampling towards functional structural and
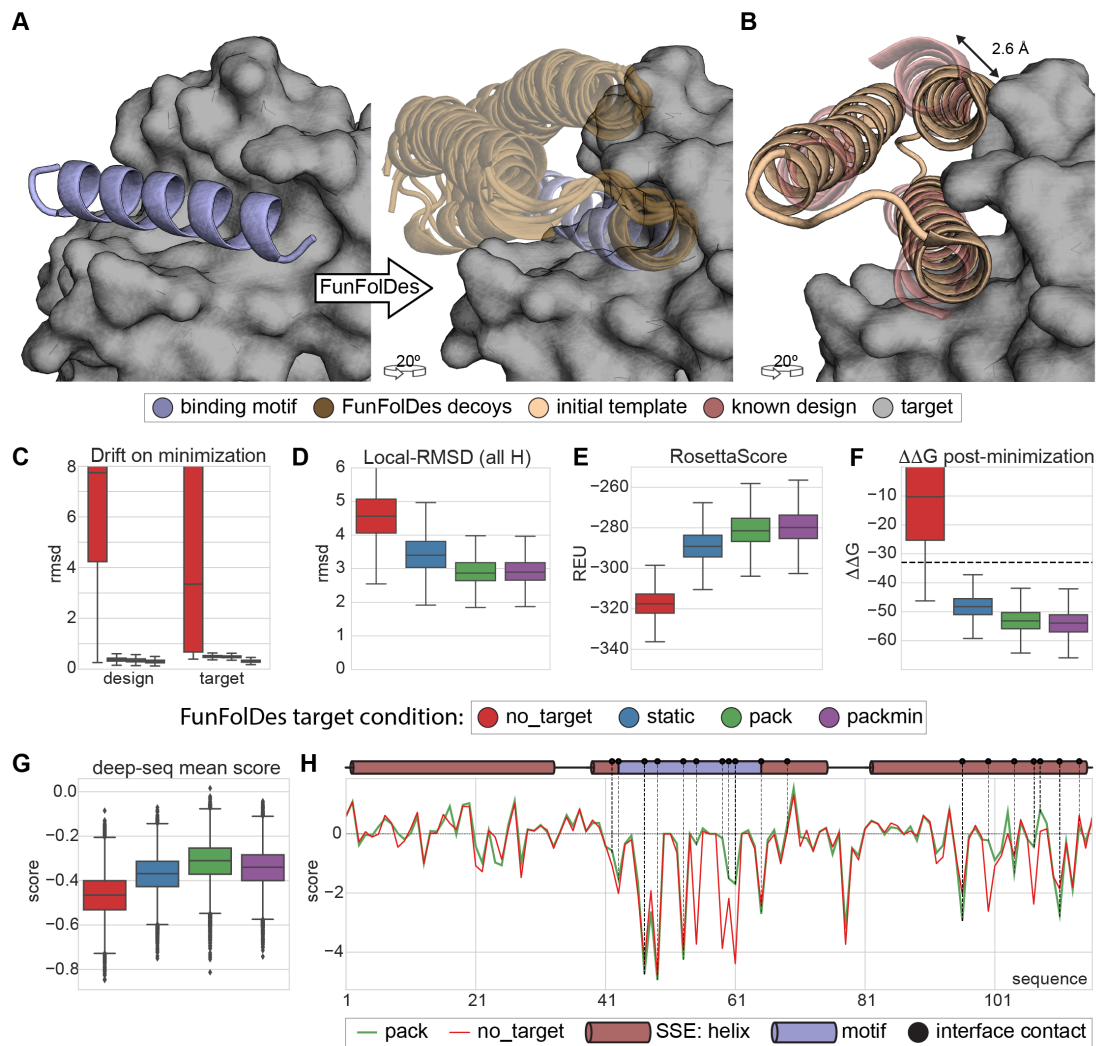
400    sequence spaces.

401

402    Previously, we used FFL to design a new binder (BINDI) to BHRF1 (**Figure 4A**),

403    an Epstein-Barr virus protein with anti-apoptotic properties directly linked to

404    the tumorigenic activity of EBV (Procko et al., 2014). FFL was used to generate

405    the initial designs that bound to BHRF1 with a dissociation constant ($K_D$) of 58-

406    60 nM, which were then affinity matured ($K_D$ = 220±50 pM) and showed

407    improved bacterial expression. BINDI was designed in the absence of the target

408    and then docked to BHRF1 through the known interaction motif. The BHRF1-

409    compatible models were further designed to ensure structural compatibility and

410    improve affinity. A striking observation from the overall approach was that the

411    FFL stage was highly inefficient, generating a large fraction of backbone

412    conformations incompatible with the binding mode of the complex.

413

414

415

**Figure 4. Target-biased design of a protein binder and assessment of performance based on saturation mutagenesis data.** A) Depiction of the initial design task, a single-segment binding motif (BIM-BH3), shown in purple cartoon representation, with its target (BHRF1), shown in gray surface is used by FunFolDes to generate an ensemble of designs compatible with the binding mode. B) Conformational difference between the initial template (PDB ID: 3LHP) and the previously designed binder (BINDI shown in pink cartoon representation), helix 3 requires a subtle but necessary shift (2.6 Å) to avoid steric clashes with the target. C-G) Scoring metrics for each design population according to the simulation mode: no_target - FunFolDes was used without the target protein; static - target present no flexibility allowed; pack - target allowed to repack the side-chains; packmin – side-chain repacking plus minimization and backbone minimization were allowed for the target. The target flexibility was allowed during the relax-design cycle of FunFolDes. C) Structural drift observed for design and target binder measured as the RMSD difference of each structure between pre- and post-minimization conformations. D) Structural recovery of the conformation observed in the BINDI-BHRF1 assess over the 3 helical segments of the bundle. E) Rosetta energy for the design populations generated with different simulation modes. F) Interaction energy (ΔΔG) between the designs and the target. G) Deep-sequence score distribution for each design population, computed as the mean score of each

18

434 sequence after applying a position score matrix based on the deep-sequence data. The pack
435 population slightly outperforms the other simulation modes. H) Per-residue scoring comparison
436 of the no_target and the pack populations according to the deep-sequence data. Although the
437 behavior is overall similar, pack outperforms no_target in multiple positions, several of which are
438 highlighted(black dots) as interfacial contacts or second shell residues close to the bind site
439 which were allowed to be designed throughout the simulations.

440

441 To test whether the presence of the target could improve structural and
442 sequence sampling, we leveraged the structural and sequence information
443 available for the BINDI-BHRF1 and benchmarked FunFolDes for this design
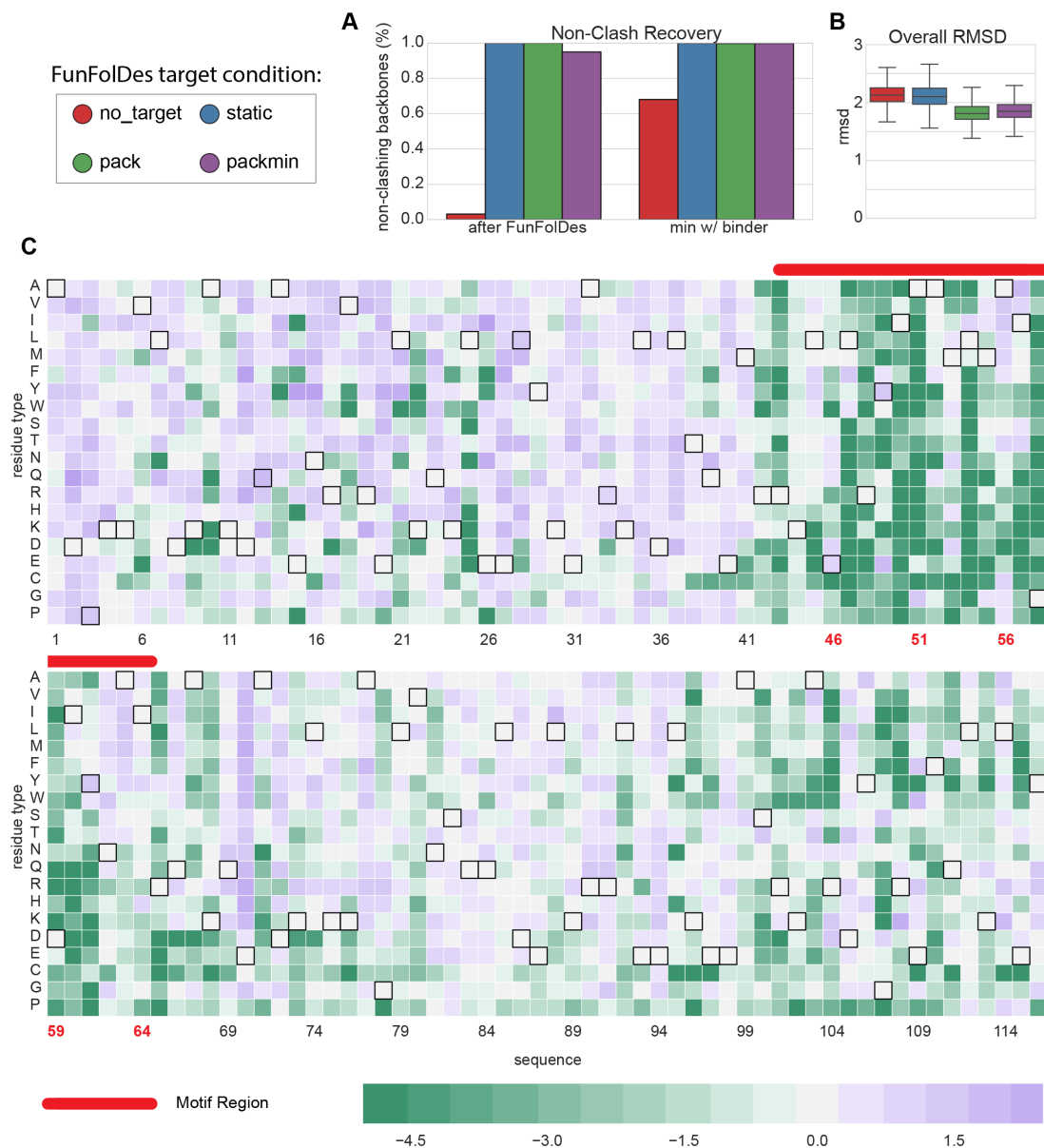444 problem.

445

446 As described by Procko and colleagues, when comparing the topological
447 template provided to FFL and the BINDI crystal structure, the last helix of the
448 bundle (helix 3) was shifted relative to the template to ensure structural
449 compatibility between BINDI and BHRF1 (**Figure 4B**). We used this case study to
450 assess the capabilities of FunFolDes to sample closer conformations to those
451 observed in the BINDI-BHRF1 crystal structure. In addition, we compared the
452 saturation mutagenesis data generated for BINDI (Procko et al., 2014) to
453 evaluate the sequence space sampled by FunFolDes.

454

455 A detailed description of this benchmark can be found in the Materials and
456 Methods section. Briefly, we performed four different FunFolDes simulations: I)
457 binding target absent (no_target); II) binding target present with no
458 conformational freedom (static); III) binding target present with side-chain
459 repacking (pack); IV) binding target present with side-chain repacking plus
460 minimization and backbone minimization (packmin). After the FunFolDes
461 simulations, the no_target set was docked to BHRF1 through the binding motif
462 and the remaining three simulations produced complexed structures. All the
463 complexes were globally minimized (both design and target) to assess the
464 conformational and energy changes as a proxy of the structural compatibility of
465 the designed binders.

466

467    Simulations performed with the target absent (no_target) very rarely produce

468    conformations compatible with the target (<10% of the total generated designs)

469    (**Figure 4 – Supplementary Figure 1A**). We observed an improvement on the

470    fraction of decoys compatible with the binding target (>60%) after global

471    minimization (**Figure 4 – Supplementary Figure 1A**). However, this was at the

472    cost of considerable structural drifts for both binder (mean RMSD 3.3 Å) and

473    target (mean RMSD 7.7 Å) (**Figure 4C**). These structural drifts are a reflection of

474    the energy optimization requirements by the relaxation algorithms but tat are

475    deemed biologically irrelevant due to the profound structural reconfigurations.

476    In contrast, simulations performed in the presence of the target clearly biased

477    the sampling to more productive conformational spaces. RMSD drifts upon

478    minimization were less than 1 Å for both designs and target (**Figure 4C**).

479

**Figure 4 - Supplementary Figure 1. Target-biased folding and design: structural features od the modeled designs and saturation mutagenesis data used for sequence recovery benchmark.** A) Quantification of the percentage of decoys compatible with a design-target binding conformation for the different simulation modes. The simulations performed without the target yield a very low percentage of binding compatible conformations. After minimization, this percentage increases with significant structural drifts. B) The initial template is a 3-helix bundle structure, the slight shift needed to adopt a binding-compatible conformation produces only a small global RMSD. C) Graphical representation of the deep-sequencing data as a position-specific score matrix. Black borders highlight the native BINDI residue type for each position. Mutations for which no data was obtained, likely reflect that these protein variants were unable to fold an display at the surface of yeast and were assigned the lowest score of -5.

493   Given that global alignments of the designs do not emphasize the local

494   differences and the helical arrangement (**Figure 4 - Supplementary Figure 1**),

495   to analyze structural regions of particular interest, we aligned all the designs on

496   the conserved binding motif (**Figure 4A**) and measured the RMSD over the three

497   helices that compose the fold. The two key regions were helices 1 and 3, which

498   are in direct contact with the target.

499

500   According to this metric, FunFolDes simulations in the presence of the target

501   sampled mean RMSD of 3 Å with the BINDI structure as reference (**Figure 4D**),

502   and the closest designs were at approximately 2 Å. On the other hand,

503   simulations in the absence of the target showed a mean RMSD of 4.5 Å, and the

504   best designs around 2.5 Å.

505   While we acknowledge that these structural differences are modest, the data in

506   this benchmark suggest that these differences can be important in the sampling

507   of conformations and sequences competent for binding.

508

509   In addition to structural sampling, we also analyzed Rosetta Energies for the

510   different simulations. We observed noticeable differences in the overall energy

511   of the designed binders; in the absence of the binding target, the designs have an

512   mean energy of approximately -320 Rosetta Energy Units (REUs), while the

513   designs generated in the presence of the binding target showed an mean

514   between -280 and -290 REUs (**Figure 4E**). This difference is significant,

515   particularly for a protein of such small size (116 residues). Likewise, we also

516   observe considerable differences in terms of the binding energy ($\Delta\Delta G$) for the

517   designs folded in the absence or in the presence of the binding target,

518   corresponding to mean $\Delta\Delta G$s of -10 and -50 REUs, respectively (**Figure 4E**).

519

520   The energy metrics provide interesting insights regarding the design of

521   functional proteins. Although the sequence and structure optimization for the

522   designs in the absence of the target reaches lower energies, these designs are

523   structurally incompatible with the binding target and, even after refinement,

524   their functional potential (as assessed by the $\Delta\Delta G$) is not nearly as favorable as

525   those performed in the presence of the binding target (**Figure 4F**). These data

526   suggest that, in many cases, to optimize function it may be necessary to sacrifice

527   the overall computed energy of the protein which is often connected to the

528   experimental thermodynamic stability of the protein. Although stability is an

529   essential requirement for all functional proteins (Chevalier et al., 2017; Tokuriki,

530   Stricher, Serrano, & Tawfik, 2008), it may be necessary to design proteins that

531   are, *in silico*, less energetically favorable to ensure that the target functional

532   requirements can be accommodated. This observation provides a compelling

533   argument to perform biased simulations in the presence of the binding target,

534   which may broadly be defined as a "functional constraint".

535

536   To evaluate sequence sampling quality, we compared the computationally

537   designed sequences to a saturation mutagenesis dataset available for BINDI

538   (Procko et al., 2014). Briefly, this dataset was obtained by screening a saturation

539   mutagenesis library for binding interactions in a yeast-display setup coupled to a

540   deep sequencing readout. The impact on the binding affinity of each mutation

541   was assessed based on the relative frequencies of the mutants. Data from this

542   experiment were transformed into a positional scoring matrix (**Figure 4 –**

543   **Supplementary Figure 1C**). Point mutations that showed a beneficial effect on

544   the binding affinity to BHRF1 have a positive score, deleterious mutants a

545   negative score, and neutral score 0. Such a scoring scheme, will yield a score of 0

546   for the BINDI sequence.

547

548   When scoring the designs generated by the four different simulations, designs

549   performed in the presence of the binding target obtain higher mean scores as

550   compared to the no_target designs (**Figure 4G**). The pack simulation, where the

551   binding target is simply repacked, is the best performer with the highest

552   distribution mean, having one design that scores better than the BINDI sequence.

553   Furthermore, it is important to highlight that in some key positions at the

554   protein-protein interface, the pack designs clearly outperformed those generated

555   by the no_target simulation, when quantified in terms of a per-position score

556   (**Figure 4H**); meaning that across the design population, amino-acids that can be

557   conducive to productive binding interactions were sampled more often in the

558   presence of the binding target. This sequence sampling benchmark provides an

559    example of the benefits of using a "functional constraint" (binding target) to

560    improve the quality of the sequences obtained by computational design.

561    Overall, the BINDI benchmark provides important insights regarding the best

562    computational protocol within FunFolDes that can be utilized to improve the

563    outcome of design simulations in terms of frequency of functional proteins.

564

565    ***Repurposing a naturally occurring fold for a new function***
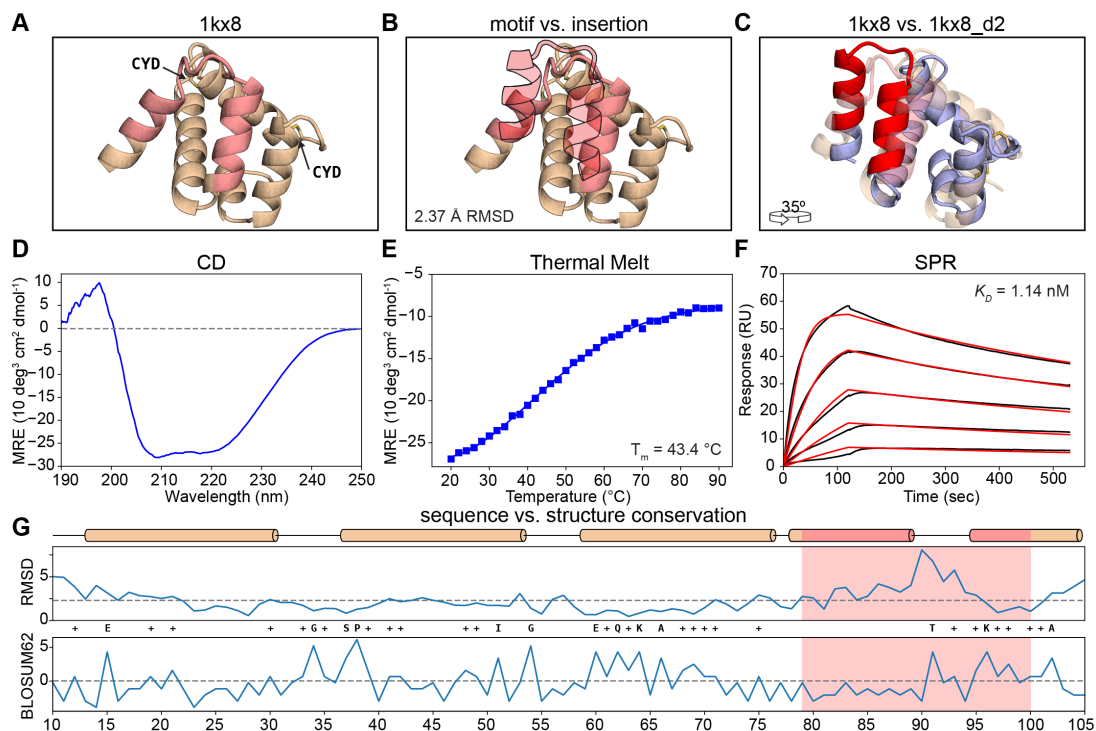
566

567    At its conception, FFL was primarily envisioned to aid the design of function into

568    proteins. To further test FunFolDes' design capabilities, we sought to transplant

569    a contiguous viral epitope that can be recognized by a monoclonal antibody with

570    high affinity (**Figure 5A**). The success of the designs was assessed by their

571    folding, thermal stability, and more importantly, binding affinities to the epitope-

572    specific antibody as the functional readout.

573

574    Specifically, we used as functional motif the RSVF site II epitope (PDB ID: 3IXT

575    (McLellan, Chen, Kim, et al., 2010)), a helix-loop-helix motif recognized by the

576    antibody motavizumab (mota). Previously, we have designed proteins with this

577    same epitope (Correia et al., 2014); however, we started from a structural

578    template with a similar conformation to that of the epitope, the RMSD between

579    the epitope and the scaffold segment was approximately 1 Å when measured

580    over the helical positions. Here, we sought to challenge FunFolDes by using a

581    distant structural template where the local RMSDs of the epitope structure and

582    the segment onto which the epitope was transplanted would be over 2 Å. We

583    used master (Zhou & Grigoryan, 2015) to perform a structural search of the site

584    II epitope over a subset of structures in the PDB. After filtering the results by

585    scaffold size (50-100 residues) and steric clashes using the structure of the

586    epitope in complex with mota, we selected as template scaffold the structure of

587    the A6 protein of the Antennal Chemosensory system from the moth *Mamestra*

588    *brassicae* (PDB ID: 1KX8 (Lartigue et al., 2002))(**Figure 5A**). The backbone

589    RMSD between the conformation of the epitope and the insertion region in 1kx8

590    is 2.37 Å (**Figure 5B**).

591

**Figure 5. Functional design of a distant structural template.** A) Structural representation of 1kx8. The insertion region is colored in light red and the two disulfide bonds are labeled (CYD). B) Structural comparison between the insertion region of 1kx8 and the mota epitope (light red-filled silhouette). Local RMSDs between the two segments reach 2.37 Å. C) Superimposition between 1kx8_d2 design model (blue with red motif) and the 1kx8 template (wheat and light red insertion site). Multiple conformational shifts are required throughout the structure to accommodate the site II epitope. D) CD spectrum of 1kx8_d2 showing a typical alpha-helical pattern with the ellipticity minimums at 208 nm and 220 nm. E) 1kx8_d2 shows a melting temperature ($T_m$) of 43.4°C. F) Binding affinity determined by SPR. 1kx8_d2 shows a $K_D$ of 1.14 nM. Experimental sensorgrams are shown in black and the fitted curves in red. G) Per-position evaluation of structural (top) and sequence (bottom) divergence between the design model 1kx8_d2 and the starting template 1kx8. The largest structural differences are observed in the region where the site II epitope was inserted, the overall difference of the two structures is 2.25 Å (dashed line). Sequence divergence is evaluated by applying the BLOSUM62 score matrix to the sequences, yielding a total of 13.5% identity and 38.5% similarity mostly in the structured regions. The epitope region is colored in light red. Identical positions between the 1kx8_d2 and 1kx8 are labeled with the residue one letter code while positively scored changes according to BLOSUM62 are labeled with a + symbol.
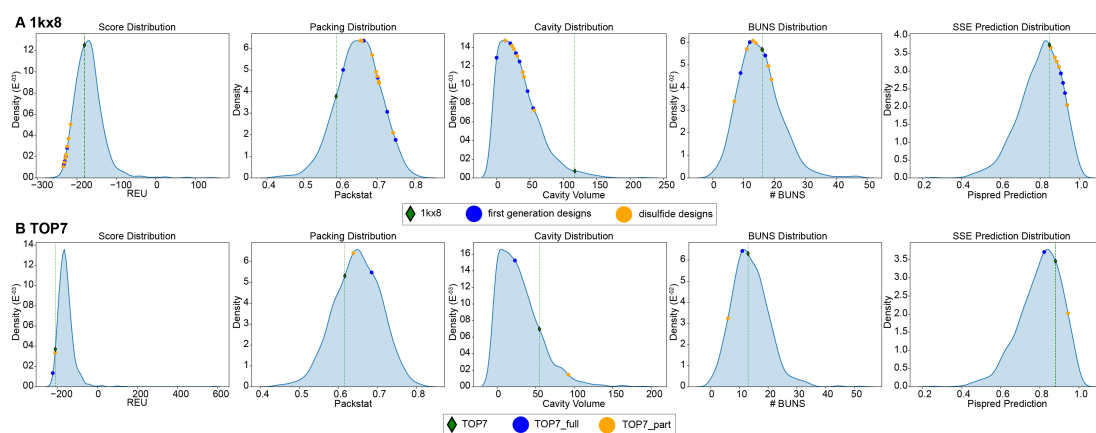
In terms of biological function, 1kx8 is involved in chemical communication and perception(Lartigue et al., 2002). Biochemically, it has been shown to bind to fatty-acid molecules with hydrophobic alkyl chains composed of 12-18 carbons.

25

615    Two prominent features are noticeable in the structure of 1kx8: two disulfide

616    bonds (**Figure 5A**) and a considerable void volume in the protein core, deemed

617    to be the binding site for fatty acid molecules. These features emphasize that the

618    initial design template is likely not a very stable protein.

619

620    In the design process we performed two stages of FunFolDes simulations; first

621    an exploratory stage to select properly folded designs with the functional motif

622    inserted (**Figure 5C**) that were fed into a second round of simulations, which

623    sampled much more extensively within the structural proximity of the 1st

624    generation template. For each stage, we generated 12'500 designs, eventually

625    selecting seven for initial experimental characterization according to several

626    structural features of the computational models, namely: Rosetta Energy,

627    packing score, and buried unsatisfied hydrogen bonds (**Figure 5 –**

628    **Supplementary Figure 1**). A detailed description of the process can be found in
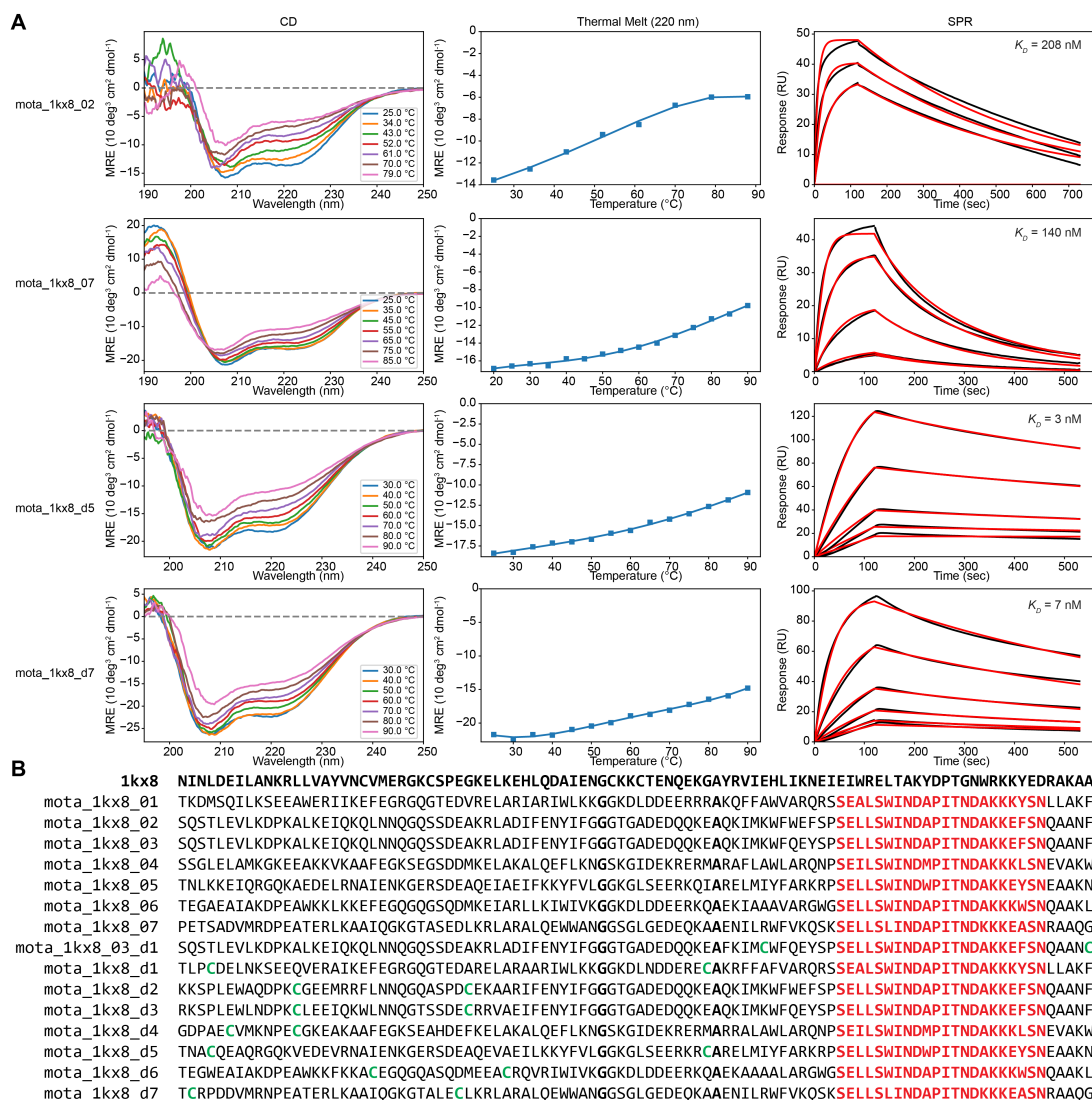
629    the Materials and Methods.

630



631

632    **Figure 5 - Supplementary Figure 1. Structural and sequence evaluation of the**

633    **computational designs.** Assessment of structural and sequence features: Rosetta Energy,

634    packing score (packstat) (Alford et al., 2017), cavity volume, Buried UNSatisfied polar atoms and

635    secondary structure prediction (PSIPRED) for the template and the computational designs. Each

636    template (green diamond) and design (yellow and blue circles) are compared against a set of

637    non-redundant minimized structures of similar size (± 15 residues). A) Due to its natural

638    function, 1kx8 presents of a large cavity to bind its hydrophobic ligands. As such, the structure

639    presents generally low scores as compared to computationally designed proteins. B)

640    Distributions of the structural and sequence features of natural proteins and the TOP7 series of

641    designs.

642

643  We characterized experimentally the computationally designed; those expressed

644  in bacteria at good yields were further characterized using size exclusion

645  chromatography coupled to a multi-angle light scatter (SEC-MALS) to determine

646  the solution oligomerization state. To assess their folding and thermal

647  stability($T_m$) we used Circular Dichroism (CD) spectroscopy, and finally to assess

648  their functional properties we used surface plasmon resonance (SPR) to

649  determine binding dissociation constants ($K_D$s) to the mota antibody. We started

650  by expressing seven sequences from our first round of computational design; out

651  of these seven, six designs were purified and characterized further. While the

652  majority of the designs were monomers in solution and showed CD spectra

653  typical of helical proteins, in terms of stability we obtained both designs that

654  were not very stable nor did they exhibit cooperative unfolding (1kx8_02) and

655  also designs that were very stable and did not fully unfold under high

656  temperatures (1kx8_07) (**Figure 5 – Supplementary Figure 2**).

657

**Figure 5 - Supplementary Figure 2. Examples of experimental characterization performed for other variants on the 1kx8 design series.** A) CD wavelength spectra (left column), thermal denaturations (middle column) and SPR binding assays with the mota antibody (right column) were performed. B) Global sequence alignment of the wild-type protein 1kx8 and the computationally designed sequences. Red positions highlight the site II epitope insertion. Green positions highlight the cysteines performing the disulfide bridges. The two positions that consistently kept the original residue type of 1kx8 are highlighted in bold.

The determined binding affinities to mota ranged from 34 to 208 nM, which was an encouraging result. Nevertheless, comparing this affinity range to those of the peptide epitope ($K_D$ = 20 nM) and other designs with the same site grafted that were published previously ($K_D$ = 20 pM) (Correia et al., 2014), there was room for improvement. Therefore, we generated a second round of designs to attempt to improve stability and binding affinities. Driven by the observation that the

28

673 native fold has two disulfides bonds, our next set of designs included engineered
674 disulfide bonds.
675
676 In the second round, we tested eight designed variants with different disulfide
677 bonds and, if necessary, additional mutations to accommodate them. All eight
678 designs were soluble after purification and two were monomeric: 1kx8_d2 and
679 1kx8_3_d1, which also showed CD spectra typical of helical proteins (**Figure 5D**)
680 with melting temperatures ($T_m$s) of 43 and 48°C (**Figure 5E**), respectively.
681 Remarkably, 1kx8_d2 showed a $K_D$ of 1.14 nM (**Figure 5F**), an improvement of
682 approximately 30-fold compared to the best variants of the first round. 1kx8_d2
683 binds to the mota with approximately 20-fold higher affinity than the peptide-
684 epitope ($K_D \approx 20$ nM), and 50 fold lower compared to previously designed
685 synthetic scaffolds ($K_D = 20$ pM) (Correia et al., 2014). This difference in binding
686 is likely reflective of how challenging it can be to accomplish the repurposing of
687 protein structures with distant structural similarity.
688
689 Post-design analyses were performed to compare the sequence and structure of
690 the best design model with the initial template. **Figure 5G**, shows a per-residue
691 RMSD measurement upon a global alignment of the 1kx8 structure with the
692 designed model. The global RMSD between the two structures is 2.25 Å. Much of
693 the structural variability arises from the inserted motif, while the surrounding
694 segments adopt a configuration similar to the original template scaffold. The
695 sequence identity of 1kx8_d2 as compared to the native protein is approximately
696 13%. The sequence conservation per-position (**Figure 5G**) was evaluated
697 through the BLOSUM62 matrix, where positive scores are attributed if the
698 original residue is not mutated or if the substitution is deemed favorable
699 according the scoring matrix, and negative if unfavorable. Overall, 38.5% of the
700 residues in 1kx8_d2 scored positively, and 61.4% of the residues had a score
701 equal to or lower than 0. This is particularly interesting in the perspective that
702 multiple mutations deemed unfavorable according the statistics condensed in
703 the BLOSUM62 matrix are still able to yield well folded and, in this case,
704 functional proteins.
705

706    The successful design of this protein is a relevant demonstration of both the

707    broad usability of the FFL algorithm and of the overall strategy of designing

708    functional proteins by coupling the folding and design process to incorporate

709    functional motifs in unrelated protein folds. In a subsequent design challenge, we

710    sought to functionalize a *de novo* design fold, which unlike natural proteins, did

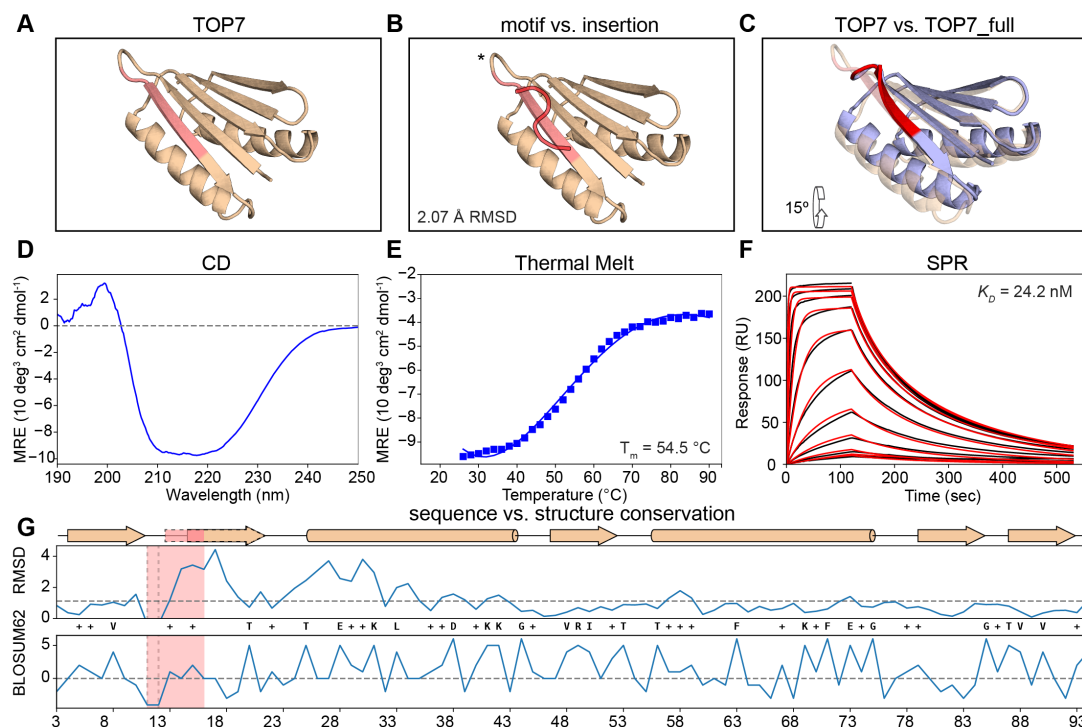711    not evolve under any sort of functional pressure.

712

713    ***Functionalization of a functionless fold***

714

715    Advances in computational design methodologies have achieved remarkable

716    results in the design of *de novo* protein sequences and structures (Hill et al.,

717    2000; Koga et al., 2012; Marcos et al., 2017). However, the majority of the

718    designed proteins are "functionless" and were designed to test the performance

719    of computational algorithms in predicting structural accuracy. Here, we sought

720    to use one of the hallmark proteins of *de novo* design efforts – TOP7 (Kuhlman et

721    al., 2003) (**Figure 6A**) – and functionalize it using FunFolDes. To do so, we

722    leveraged several of the newly implemented features in FunFolDes. The

723    functional site selected to insert into TOP7 was another viral epitope from RSVF,

724    commonly referred to as site IV, which is recognized by the 101F antibody

725    (McLellan, Chen, Chang, et al., 2010). When bound to the 101F antibody, site IV

726    adopts a β-strand-like conformation (**Figure 6B**), which in terms of secondary

727    structure content is compatible with one of the edge strands of the TOP7

728    topology (**Figure 6C**). Despite the secondary structure similarity, the RMSD of

729    the site IV backbone in comparison with that of TOP7 is 2.1 Å, and upon

730    alignment of the antibody in this particular orientation, clashes arise between

731    TOP7's helix 1 and the antibody interface. Therefore, this design challenge is yet

732    another prototypical application for FunFolDes. In this design challenge we

733    followed two distinct routes: I) a conservative approach where we fixed the

734    amino-acid identities of roughly half of the core of TOP7 and allowed mutations

735    mostly on the contacting shell of the epitope insertion site; and II) a sequence

736    unconstrained design where all the positions of the scaffold were allowed to

737    mutate. We attempted five designs for recombinant expression in *E. coli* and two

738    (TOP7_full and TOP7_partial) were selected for further biochemical and

739    biophysical characterization, one from each of the two design strategies
740    mentioned above. According to SEC-MALS, both behaved as monomers in
741    solution, with TOP7_partial being a less well-behaved protein with higher
742    aggregation propensity. Both TOP7_full and TOP7_partial (**Supplementary Figure 5**)
743    were folded according to CD measurements, with the TOP7_full showing a CD
744    spectrum (**Figure 6D**) which very closely resembles that of the native TOP7
745    (Kuhlman et al., 2003). TOP7_full was subjected to thermal denaturation
746    monitored by CD, where we observed that the newly designed protein is much
747    less stable than the original TOP7 (**Figure 6E**). To quantify the functional
748    component of TOP7_full, we determined the $K_D$ of its interaction with 101F to be
749    24.2 nM (**Figure 6F**), which is within the range measured for the native viral
750    protein RSVF (3.6 nM) (McLellan, Chen, Chang, et al., 2010). Importantly, the $K_D$
751    for TOP7_full is 2400 fold higher than that of the peptide-epitope (58.4 μM)
752    (McLellan, Chen, Chang, et al., 2010), suggesting that productive conformational
753    stabilization and/or extra contacts to the rest of the protein were successfully
754    designed.

755



757    **Figure 6. Functionalization of the functionless de novo fold TOP7.** A) Structure of TOP7 with
758    the insertion region highlighted in light red. B) Structural comparison between 101F and the
759    insertion region of TOP7 reveals a 2.07 Å RMSD. C) TOP7_full model (in blue and red for the

31

760    motif) superimposed over the TOP7 crystal structure. 101F's insertion is structurally

761    compensated mostly by the first pairing beta strand and a shift of the first alpha helix D) CD

762    spectrum shows a broad ellipticity signal between 210 nm and 222 nm as a representative of

763    mixed secondary structural propensities. E) The melting temperature ($T_m$) for TOP7_full was

764    54.5 °C. F) Binding affinity determined by SPR. TOP7_full shows a $K_D$ of 24.2 nM. Experimental

765    sensorgrams are shown in black and the fitted curves in red. G) Per-position evaluation of

766    structural (top) and sequence (bottom) divergence between the design model TOP7_full and the

767    starting template TOP7. The largest structural differences are observed in the region

768    downstream of the site IV epitope , the overall difference of the two structures is 1.5 Å (dashed

769    line). Sequence divergence is evaluated by applying the BLOSUM62 score matrix to the

770    sequences, yielding a total of 27.7% identity and 52.2% similarity. The epitope region is colored

771    in light red. Identical positions between the TOP7_full and TOP7 are displayed as their residue

772    types while positively scored changes according to BLOSUM62 are labeled with a + symbol.
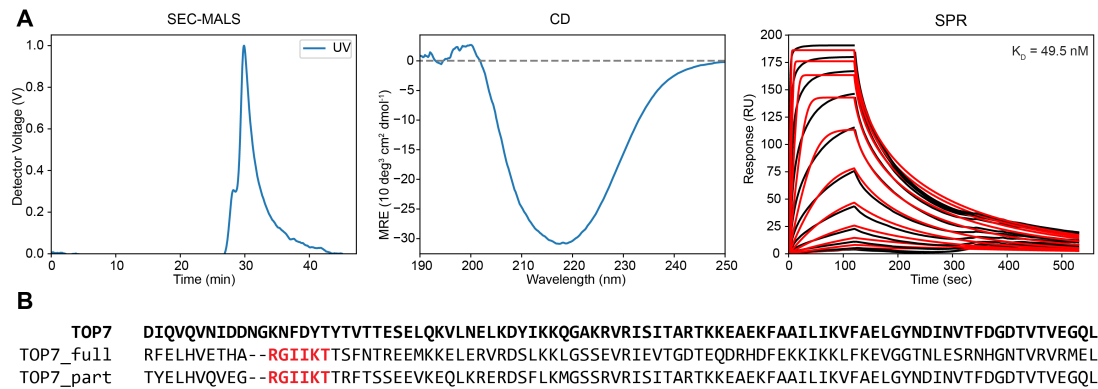
773

774    Given the successful functionalization of TOP7, we sought to understand the

775    levels of sequence and structural change (**Figure 6G**). Per-residue sequence

776    recovery and structural similarity were evaluated for TOP7_full against TOP7

777    (**Figure 6G**). We compared the per-residue RMSD between the TOP7_full model

778    and the crystal structure of TOP7, revealing that most conformational changes

779    occur from the site IV insertion region and displacement of the neighboring

780    alpha-helix, with the overall backbone RMSD between both structures being 1.5

781    Å. The connecting loop between the strand that holds the epitope and the

782    adjacent strand was also shortened to obtain a tighter connection between the 2

783    strands (**Figure 6C**).

784

785    Remarkably, the sequence identity of the most aggressive design (TOP7_full) is

786    only 28%, and using the BLOSUM62 based scoring system, we observe that most

787    of the TOP7_full residues were actually favorable, obtaining positive scores. This

788    low conservation is especially relevant considering that intensive studies on

789    TOP7 have revealed the importance of beta-sheet conservation in order to keep

790    its foldability (Boschek et al., 2009; Soares, Boschek, Apiyo, Baird, & Straatsma,

791    2010; Viana et al., 2013).

792

**Figure 6 - Supplementary Figure 1. Experimental characterization of TOP7_variants.** A) Experimental characterization for the TOP7_partial design: SEC-MALS elution profile (left column); CD wavelength scan spectrum; SPR binding assays with the 101F antibody (right column). Noticeably the TOP7_partia show a CD spectrum notoriously different from WT TOP7 and the TOP7_full design. B) Global sequence alignment of the wild-type protein TOP7 and the computationally designed sequences. Red positions highlight the site IV epitope insertion.

In summary, our results show that FunFolDes was able to repurpose a functionless protein by folding and designing its structure to harbor a functional site, which in this case was a viral epitope. Previously, these computationally designed proteins with embedded viral epitopes were dubbed epitope-scaffolds and showed their biomedical applicability as immunogens that were able to elicit viral neutralizing antibodies(Correia et al., 2014).

## Discussion and Conclusions

809

810

811 The robust computational design of proteins that bear a biochemical function

812 remains an important challenge for present methodologies. The ability to

813 consistently repurpose old folds for new functions or the *de novo* design of

814 functional proteins could bring new insights into the determinants necessary to

815 encode function into proteins (e.g. dynamics, stability, etc.) as well as important

816 advances in translational applications (e.g. biotechnology, biomedical,

817 biomaterials, etc.).

818 Here, we present the second-generation computational design protocol Rosetta

819 FunFolDes, which was conceived to embed functional motifs into protein

820 topologies, allowing for a global retrofitting of the overall protein topology to

821 favorably host the functional motif. FunFolDes has evolved to incorporate two

822 types of constraints to guide the design process: topological and functional. The

823 former entail the fragments to assemble the protein structure and sets of

824 distance constraints that bias the folding trajectories towards a desired topology;

825 and the latter are the structure of the functional motif inserted and the binding

826 target, if used.

827

828 We have extensively benchmarked the protocol, leveraging natural structural

829 and sequence variation of proteins within the same fold, as well as deep

830 mutational scanning data for the computationally designed protein BINDI. In our

831 first benchmark, we observed that with FunFolDes we can efficiently bias the

832 sampling towards improved structural and sequence spaces. Protocol features

833 that enable higher quality sampling in design simulations are extremely

834 important. Improved sampling may contribute to solving some of the major

835 limitations in protein design, related to "junk" sampling, where most of the

836 generated designs are not physically realistic, exhibiting obvious flaws according

837 to general principles of protein structure. Importantly, higher quality sampling

838 will likely contribute to improve the success rate of designs that are tested

839 experimentally. The BINDI benchmark allowed us to test FunFolDes in a system

840 with a large amount of experimental data, which included both sequences and

841 structures. Perhaps the most enlightening observation was that designs that

34

842    were theoretically within a sequence/structure space productive for binding to

843    the target were rather far from the energetic minimum that the protein fold can

844    achieve in the absence of the binding target. Considering that the large majority

845    of the design algorithms are energy "greedy" and the sequence/structure

846    searches are performed with the central objective of finding the global minimum

847    of the energetic landscape, by introducing functional constrains into the

848    simulations, FunFolDes presents an alternative way of designing functional

849    molecules and efficiently skewing the searches towards off-minima regions of

850    the global landscape. We anticipate that such finding will be more relevant for

851    protein scaffolds that need to undergo a large degree of structural adaptation to

852    perform the desired function. If confirmed that this finding is generalized across

853    multiple design problems, it could be an important contribution for the field of

854    computational protein design.

855

856    Furthermore, we used FunFolDes to tackle two design challenges and

857    functionalized two proteins with two distinct viral epitopes. These design

858    challenges were devised to test the applicability of FunFolDes. Importantly, in

859    previous applications FFL always used three-helix bundles as design templates,

860    here we diversified the template folds and used an all-helical protein that is not a

861    bundle (1kx8) and a mixed alpha-beta protein (TOP7), clearly showing the

862    applicability to other folds. For the 1kx8 design series, we evaluated the

863    capability of using distant structural templates as starting topologies as a

864    demonstration of how one can use the many naturally occurring protein

865    structures available and repurpose their function even when the initial template

866    and the target structures are quite different. We obtained stable proteins that

867    where recognized by an anti-RSV antibody with high affinity, showing that in this

868    case, we successfully repurposed a distant structural template for a different

869    function, a task for which other computational approaches (Silva, Correia, &

870    Procko, 2016) would have limited applicability. We see this result as an exciting

871    step forward towards using the wealth of the natural structural repertoire for

872    the design of novel functional proteins.

873

874    In a last effort, we functionalized a "functionless" fold, based on one of the first *de*
875    *novo* designed proteins – TOP7. For us, this challenge has important implications
876    in order to understand the design determinants and biochemical consequences
877    of inserting a functional motif into a protein that was mainly optimized for
878    thermodynamic stability. We were successful in functionalizing TOP7 differently
879    than previous published efforts, where TOP7 was mostly used as a carrier
880    protein with functional motifs fused onto loop regions or side chains grafted in
881    the helical regions, while our functional motif was embedded in the beta sheet of
882    the protein template (Boschek et al., 2009; Soares et al., 2010; Viana et al., 2013).
883    Exciting advances in the area of *de novo* protein design are also yielding many
884    new proteins, which could then be functionalized with FunFolDes, highlighting
885    the usefulness of this approach. Interestingly, we observed that the
886    functionalized version of TOP7 showed a dramatic decrease in thermodynamic
887    stability as compared to the parent protein. While this observation can be the
888    result of many different factors, it is compelling to interpret it as the "price of
889    function", meaning that to harbor function, the TOP7 protein was penalized in
890    terms of stability, which would be consistent with our findings in the BINDI
891    benchmark example.
892    Recently, there have also been several *de novo* proteins which were designed for
893    functional purposes (Chevalier et al., 2017); however, these efforts were limited
894    to linear motifs that carried the functions, and the functionalization was mainly
895    accomplished by side-chain grafting (Correia et al., 2010; Kulkarni et al., 2015)
896    and relied on screening of a much larger number of designed proteins.

897

898    From our perspective, and considering all the technical improvements,
899    FunFolDes has matured to become a valuable resource for the robust
900    functionalization of proteins using computational design. Here, we present a
901    number of important findings provided by the detailed benchmarks performed
902    and used the protocol to functionalize proteins in design tasks which are
903    representative of some of the common challenges that the broad scientific
904    community faces when using computational design approaches.

905

## Materials and Methods

906

907

908    *Computational protocol description*

909

910    Rosetta Functional Folding and Design (FunFolDes) is a general approach for
911    grafting functional motifs into protein scaffolds. It's main purpose is to provide
912    an accessible tool to tackle specifically those cases in which structural similarity
913    between the functional motif and the insertion region is low, thus expanding the
914    pool of structural templates that can be considered useful scaffolds. This
915    objective is achieved by folding the scaffold after motif insertion while keeping
916    the structural motif static. This process allows the scaffold's conformation to
917    change and properly adapt to the three dimensional restrictions enforced by the
918    functional motif. The pipeline of the protocol (summarized in **Figure 1**) proceeds
919    as follows:

920

921    I) *Selection of the functional motif.* A single or multi-segment motif must be
922    selected and provided as an input. In the most common mode of the protocol
923    dihedral angles, side chain identities and conformations are kept fixed
924    throughout the whole protocol.

925

926    II) *Selection of the protein scaffold.* Searches for starting protein scaffolds can be
927    achieved, but are not limited to, RMSD similarity matches to the Protein Data
928    Bank (PDB) (Rose et al., 2017). The ability of FunFolDes to adapt the scaffold to
929    the needs of the motif widens the structural space of what can be considered as a
930    suitable template. Thus, this step requires human intervention and has to be
931    performed outside of the main protocol.

932

933    III) *Generation of fragment databases.* The usage of fragments lies at the core of
934    many Rosetta protocols, particularly those that perform large explorations of the
935    conformational space required for structure prediction and design. The most
936    standard way of assembling fragment sets is to generate sequence-based
937    fragments using the FragmentPicker application (Kim, Blum, Bradley, & Baker,
938    2009). Despite the usefulness of the sequence-based fragments in typical design

939    and structure prediction problems, FunFolDes-derived designs depend on the
940    structural content of the template rather than its sequence. Thus we
941    implemented the *StructFragmentMover*, a mover that performs on-the-fly
942    fragment picking based on secondary structure, dihedral angles and solvent
943    accessibility, calculated from the template's structural information. The typical
944    three- and nine residue-long fragment sets are generated from the fragment
945    database included in the Rosetta tools release.

946

947    IV) *Generation of constraints.* Residue-pair distance and backbone dihedral angle
948    constraints can be extracted from the protein scaffold to guide the folding
949    process. These constraints may include the full-length protein or focus in specific
950    segments while allowing a wider flexibility in other regions. Although not
951    required, the use of constraints greatly increases the quality of the sampling. The
952    protocol can be also made aware of other constraint types (such as cartesian
953    constraints) by properly modifying the score functions applied to the *ab initio*
954    stage (Simons, Bonneau, Ruczinski, & Baker, 1999).

955

956    V) *Construction of the extended pose.* The extended structure is composed of all
957    the segments of the target motif maintain their native backbone conformation
958    and internal rigid body orientation. The scaffold residues are linearly attached to
959    previously defined insertion points. In multi-segment motif scenarios, the
960    construct will present a chain break between each of the motif composing
961    segments. The number of chain-breaks in the pose scales with the number of
962    segments(n) within a motif always resulting in n-1 chain-breaks. Once the
963    extended pose is assembled, it is represented at the centroid level (all side-chain
964    atoms in a single virtual atom) to reduce the computational cost of the
965    simulation.

966

967    VI) *Folding the extended pose.* Fragment insertion is performed to accomplish the
968    folding stage. Kinematics of the pose are controlled through the FoldTree (Wang,
969    Bradley, & Baker, 2007), a system to control the propagation of the torsion
970    angles applied to a structure. In single-segment motif structures, the FoldTree
971    starts in the center of the motif and propagates in opposite directions towards

972 the N- and C-terminal of the protein. In multi-segment motifs, in which the pose

973 bears chain-breaks between each pair of motif segments, the FoldTree has a

974 fixed node in the center of each segment and expands towards both sides

975 (**Figure 1**). The chain-breaks in the structure are marked as cut-points, which

976 avoid further propagation of the kinematic movement throughout the

977 polypeptide chain, and are subjected to a score term to promote their spatial

978 proximity. All the nodes of the FoldTree are placed in the motif segments are

979 kept fixed relative to each other in three-dimensional space; this setup allows for

980 the folding of the protein while maintaining the relative position of all the motif

981 segments.

982

983 VII) *Inclusion of the binding target*. If a binding target (protein, nucleic acid or

984 small molecule ligand) is provided, a new FoldTree node is added to the closest

985 residue between the first motif segment and each binding element. Similarly to

986 the multi-segment kinematics, this ensures that the rigid-body orientation

987 between the motif and its target is maintained. FunFolDes can handle

988 simulations with both multi-segment and binding targets simultaneously.

989

990 VIII) *Folding post-processing*. Folding trajectories are considered successful if

991 they generate structures under a user-defined RMSD threshold of the starting

992 scaffold. In case of a multi-segment motif, a preliminary loop closure will be

993 executed to generate a continuous polypeptide chain, and the kinematic setup

994 maintained to avoid segment displacement during the design step. After the

995 folding stage performed at the centroid level, full atom information is recovered.

996 All the steps necessary to perform the setup of the extended pose (kinematic

997 setup, folding, post-processing) are carried out by a newly implemented mover

998 called *NubInitioMover*.

999

1000 IX) *Protein design and conformational relaxation*. The folded structure is

1001 subjected to iterative cycles of sequence design (Hu, Wang, Ke, & Kuhlman,

1002 2007) and structural relaxation (Tyka et al., 2011) in which the sequence search

1003 is coupled with confined conformational sampling (Kuhlman & Baker, 2004). A

1004 MoveMap is defined to control backbone dihedrals and side chain conformations

1005    of the motif segments and the binding target while allowing for backbone and

1006    side-chain exploration of the movable residues. TaskOperations are used to

1007    avoid undesired mutations in the functional motif.

1008

1009    X) *Loop closure*. If multi-segment motifs are used, a final loop closure step is

1010    required in order to obtain a polypeptide chain without breaks. The

1011    *NubInitioLoopClosureMover* performs this last step using the Cyclic Coordinate

1012    Descend (CCD) protocol (Wang et al., 2007), while ensuring that the original

1013    conformation and rigid-body orientation of the motifs is maintained. After the

1014    closure of each cut-point, a final round of fixed backbone design is performed on

1015    the residues of the cut-points and surroundings.

1016

1017    XI) *Selection, scoring and ranking*. Finally, the decoys are ranked and selected

1018    according to Rosetta energy, structural metrics (core packing, buried unsatisfied

1019    polar atoms, etc) (Alford et al., 2017), sequence-based predictions such as

1020    secondary structure propensity (Jones, 1999) and folding propensity (Simons,

1021    Bonneau, et al., 1999) or any other metrics accessible through RosettaScripts

1022    (RS).

1023

1024    The pipeline components described here represent the most standardized

1025    version of the FunFolDes protocol. By means of its integration in RS, different

1026    stages can be added, removed or modified to tailor the protocol to the specific

1027    needs of the design problem at hand.

1028

1029    *Capturing conformational and sequence changes in small protein domains*

1030

1031    To test the ability of FunFolDes to recover the required conformational changes

1032    to stabilize a given structural motif, we created a benchmark of 14 target cases of

1033    proteins with less than 100 residues, named T01 to T14. Each target case was

1034    composed of two structures of the same CATH superfamily(Dawson et al., 2017).

1035    One of the structures was representative of the shared structural features of the

1036    CATH family; we called this structure the reference. The second protein within

1037    each target case presents two structural variations with respect to the reference:

1038    I) an insertion or deletion (indel) region and II) a conformational change. Direct

1039    structural contacts between these two regions make it so that the indel region is

1040    the cause for the conformational change. We called this second structure the

1041    target (**Figure 2**, **Table 1**).

1042

1043

| ID | CATH | # | reference | target | motif range |
|----|------|---|-----------|--------|-------------|
| **T01** | CATH.3.40.140.10 | 1 | 1pgxA | 2pw9C | 69-73 |
| **T02** | CATH.3.30.310.50 | 1 | 3i3wA | 4bjuA | 464-486 |
| **T03** | CATH.3.30.70.980 | 1 | 1lfpA | 1mw7A | 140-150 |
| **T04** | CATH.3.30.70.100 | 1 | 1rjjA | 1lq9A | 19-45 |
| **T05** | CATH.3.10.20.30 | 1 | 2q5wD | 2pkoA | 49-64 |
| **T06** | CATH.2.30.29.30 | 1 | 1c1yB | 1h4rA | 39-59 |
| **T07** | CATH.3.10.20.90 | 1 | 2bkfA | 2al6B | 115-119 |
| **T08** | CATH.3.10.20.90 | 1 | 1wj4a | 1wiaA | 181-200 |
| **T09** | CATH.3.10.20.90 | 1 | 3ny5B | 3phxB | 100-121 |
| **T10** | CATH.3.10.20.310 | 1 | 2x8xX | 2qdfA | 103-121 |
| **T11** | CATH.3.10.320.10 | 1 | 4p5mA | 2bc4C | 56-66 |
| **T12** | CATH.2.40.40.20 | 1 | 1cr5B | 2pjhB | 119-142 |
| **T13** | CATH.2.40.40.20 | 2 | 1cr5B | 2pjhB | 119-142, 168-173 |
| **T14** | CATH.3.30.110.40 | 1 | 1jdqA | 3lvjC | 14-37 |

1044    **Table 1. Targets included in the conformational and sequence recovery benchmark.** For

1045    each of the benchmark target is indicated the CATH superfamily and representatives used in the

1046    simulations. (#) indicates the number of segments in the target protein that are considered motif.

1047    Motif range indicates the residues considered motif according to the PDB numbering.

1048

1049    For each template protein we generated approximately 10000 decoys with

1050    FunFolDes by folding the target with the following conditions: 1) the indel region

1051    was considered as the motif, meaning that its structural conformation was kept

1052    fixed and no mutations allowed; 2) residue-pair distance constraints were

1053    derived from the secondary structure elements conserved between reference

1054    and the target (constrained region); 3) the region of the protein which showed

1055 the largest structural variations (query region) was constraint-free throughout
1056 the simulation.

1057

1058 FunFolDes simulations were compared with constrained *ab initio (*cst-*ab initio)*
1059 simulations, the key difference between them being that the cst-*ab initio*
1060 simulations allowed for backbone flexibility in the motif region. The comparison
1061 between both approaches provides insights on the effects of a static segment in
1062 the folding trajectory of the polypeptide chain. In both scenarios a threshold was
1063 set after the folding stage where only decoys that had less than 5 Å RMSD from
1064 the template were carried to the design stage.

1065

1066 The importance of the input fragments was assessed with our benchmark. Both
1067 protocols were tested with sequence-based fragments and structure-based
1068 fragments generated on-the-fly by FunFolDes. Comparison between the two
1069 types of fragments provides insight into how to utilize FunFolDes in the most
1070 productive manner.

1071

1072 Structural recovery was evaluated by RMSD with the target structure. Global
1073 RMSD, understood as the minimum possible RMSD given the most optimal
1074 structural alignment, was used to assess the overall structural recovery of each
1075 decoy population. Local RMSD, was evaluated for the unconstrained (query)
1076 region and the motif by aligning each decoy to the template through the
1077 constrained segments (excluding the motif). This metric aimed to capture the
1078 specific conformational changes required to accommodate the motif into the
1079 structure (**Figure 2B**, **Figure 2 - Supplementary Figure S1B**).

1080

1081 Sequence recovery was evaluated through two different criteria, sequence
1082 associated statistics and Hidden Markov Model (HMM) (Eddy, 2011). For the
1083 sequence associated statistics, we quantified sequence identity and similarity
1084 according to BLOSUM62 for the core residues of each protein, as defined by
1085 Rosetta's *LayerSelector* (Koga et al., 2012). Motif residues, that were not allowed
1086 to mutate, were excluded from the statistics. In the second criteria, position
1087 specific scoring matrices with inter-position dependency known as Hidden

1088     Markov Model (HMM) were used to evaluate fold specific sequence signatures. In

1089     this case, the closest HMM to the template structure provided by CATH was used

1090     to query the decoys and identify those that matched the HMM under two

1091     conditions: I) an e-value under 10 and II) a sequence coverage over 50%.

1092     Although these conditions are wide, they were within the variability found

1093     between members of CATH superfamilies with high structural and sequence

1094     variability like the ones used in the benchmark.

1095

1096     *Target-biased design of protein binders*

1097

1098     To assess the performance of FunFolDes in the presence of a binder target we

1099     recreated the design of BINDI as a binder for BHRF1 (Procko et al., 2014), the

1100     BHRF1 binding motif from the BIM BH3 protein (PDB ID:2WH6 (Kvansakul et al.,

1101     2010)) was inserted into a previously described 3-helix bundle scaffold (PDB

1102     ID:3LHP (Correia et al., 2010)).

1103

1104     On that account, four different design simulations were performed: one without

1105     the binder (no_target) and three in the presence of the binder (static, pack and

1106     packmin). The difference between the last three relates to how the binding target

1107     was handled. In the static simulations the binding target was kept fixed and no

1108     conformational movement in the side chains was allowed throughout the

1109     protocol. In the pack simulations the side chains of the binding target were

1110     repacked during the binder design stage. Finally, in the packmin simulations the

1111     binding target side-chains were allowed to repack and both side-chains and

1112     backbone were subjected to minimization. These three target configurations are

1113     easily obtained by altering MoveMap definitions, demonstrating the flexibility of

1114     the protocol to include variable conditions. In all cases, the two terminal residues

1115     on each termini of the binding motif were allowed backbone movement to

1116     optimize the insertion in the 3-helix bundle scaffold. For each one of these

1117     scenarios, approximately 20000 decoys were generated.

1118

1119     For the no_target simulations the FunFolDes designs were docked to BHRF1

1120     using the inserted motif as guide to assess their complementary and interface

1121 metrics. In all the simulations, a final round of global minimization was
1122 performed where both proteins of the complex were allowed backbone
1123 flexibility. During this minimization, the jump between the design and target was
1124 kept fixed to maintain the binding motif and target in place. The final ΔΔG of the
1125 complexes was measured after the minimization step to enable comparasions
1126 between the no_target decoys and the remaining simulation modes. Structural
1127 changes related to this minimization step were evaluated as the global RMSD
1128 between each structure before and after the process, this measure is referred to
1129 as RMSD drift.

1130

1131 Structural evaluation includes global RMSD against the BINDI design crystal
1132 structure (PDB ID: 4OYD (Procko et al., 2014)) as well as local RMSD against
1133 regions of interest of BINDI. For the Local-RMSD the structures were aligned
1134 through the inserted motif, as it was kept throughout all simulations and with
1135 respect to BINDI. The local RMSD analysis was performed over all the helical
1136 segments contained in the structures (all H), which provides a measurement of
1137 the structural shifts on the secondary structure regions of the designs.

1138

1139 To evaluate the sequence recovery of our simulations we leveraged BINDI's
1140 saturation mutagenesis data analyzed by deep sequencing performed in the by
1141 Procko et al (Procko et al., 2014). The experimental fitness of each mutation was
1142 summarized in a score matrix where a score was assigned for each amino-acid
1143 substitution for the 116 positions of the protein (**Figure 4 – Supplementary
1144 Figure 1**). In summary, point mutations that improved BINDI's binding to
1145 BHRF1 are assigned positive scores while deleterious mutations present
1146 negative values. These scores are computed based on experimental data where
1147 the relative populations of each mutant were compared between a positive
1148 population of cells displaying the designs (binders) and negative populations
1149 (mutants that display but don't bind), these experiments have been described in
1150 detail elsewhere(Procko et al., 2014). After normalization by the score of the
1151 final BINDI sequence in each position, a position sequence specific matrix
1152 (PSSM) was created. Like the original data, this matrix also assigns a positive
1153 score to each point mutation if it resulted in an improved binding for the design.

1154   This normalization provides a score of 0 for the BINDI sequence, which is useful

1155   as a reference score.

1156

1157   *Repurposing naturally occurring folds for a new functions*

1158

1159   To experimentally validate the capabilities of FunFolDes and insert functional

1160   sites in structurally distant templates, we decided to transfer the site II epitope

1161   from the Respiratory Syncytial Virus (RSV) protein F (PDB ID:3IXT (McLellan,

1162   Chen, Kim, et al., 2010)) into heterologous scaffolds. This is a continuous, single

1163   segment, helix-loop-helix conformation epitope. The main objective was to

1164   challenge the capabilities of FunFolDes to reshape the structure of the scaffold to

1165   the requirements of the functional motif, we aimed to search for insertion

1166   segments with RMSDs towards the site II structure higher than 2 Å.

1167

1168   We searched for host scaffolds using MASTER (Zhou & Grigoryan, 2015) where

1169   we used the full-length site II segment as a query against a subset of 17539

1170   protein structures from the PDB composed of 30% non-redundant sequences

1171   included in the MASTER distribution. The RMSD between the query and

1172   segments on the scaffolds were assessed using backbone $C_\alpha$s. All matches with

1173   $RMSD_{C\alpha} < 5.5$ Å relative to site II were recovered and further filtered by protein

1174   size, where only proteins between 50 and 100 residues were kept. These

1175   scaffolds were then ranked regarding antibody-binding compatibility, where

1176   each match was realigned to the antibody–epitope complex and steric clashes

1177   between all glycine versions of the scaffold and antibody were quantified using

1178   Rosetta. All matching scaffolds with ΔΔG values above 100 REU were discarded

1179   under the assumption that their compatibility with the antibody binding mode

1180   was to low. The remaining scaffolds were visually inspected and PDB ID: 1kx8

1181   (Lartigue et al., 2002) ($RMSD_{C\alpha} = 2.37$ Å) was selected for design with FunFolDes.

1182

1183   The twenty-one residues from the site II epitope (motif) as present in 3IXT were

1184   grafted into a same sized segment (residues 79-100) of 1kx8 using the

1185   *NubInitioMover*. Up to three residues in each insertion region of the motif were

1186   allowed backbone flexibility in order to allow proper conformational transitions

1187     in the insertion points . Atom pair constraints with a standard deviation of 3 Å

1188     were defined for all template residues, leaving the motif residues free of

1189     constraints. The generous standard deviation was set up to favour necessary

1190     conformational changes to allow the optimal fitting of the motif within the

1191     topology . Regardless, the total allowed deviation for template was limited at 5 Å

1192     to ensure the retrieval of the same topology. In this design series we used

1193     sequence-based fragments generated with the 1kx8 native sequence. Three

1194     cycles of design/relax were performed on the template residues with the

1195     *FastDesignMover*.

1196

1197     A first generation of 12500 designs was ranked according to Rosetta energy.

1198     From the top 50 decoys, only one presented the motif without distortions on the

1199     edges derived from the allowed terminal flexibility. This decoy was used as

1200     template on the second generation of FunFolDes to enhance the sampling of

1201     properly folded conformations, with the same input conditions as the previous

1202     one.

1203

1204     In the second generation, the top 50 decoys according to Rosetta energy were

1205     further optimized through additional cycles of design/relax. After a selection

1206     based on Rosetta energy, buried unsatisfied polars and secondary structure

1207     prediction using PSIPRED (Jones, 1999), a total of 7 designs were manually

1208     optimized and selected for experimental characterization. After the initial

1209     characterization, designs with added disulphide bridges were generated in order

1210     to improve protein stability and affinity (**Figure 4 – Supplementary Figures 1**

1211     **and 2**).

1212

1213     *Functionalization of a functionless fold*

1214

1215     In a second effort to test the design capabilities of FunFolDes we sought to insert

1216     a functional motif in one of the first de novo designed proteins – TOP7 (PDB ID:

1217     1QYS) (Kuhlman et al., Science, 2003).

1218

1219 Six residues from the complex between the antibody 101F and the peptide-
1220 epitope, corresponding to residues 429-434 (RGIIKT) on the full-length RSV F
1221 protein (McLellan, Chen, Chang, et al., 2010), were grafted into the edge strand of
1222 the TOP7 backbone using FunFolDes. The choice between epitope and hosting
1223 scaffold was made based on the secondary structure adopted by the epitope and
1224 the acceptable structural compatibility of the TOP7 structure, the $RMSD_{C\alpha}$
1225 between the epitope an the insertion segment was 2.07 Å.

1226

1227 To ensure that the majority of the β-strand secondary structure was maintained
1228 throughout the grafting protocol, the epitope motif was extended by one residue
1229 and a designed 4-residue β-strand (KVTV) pairing with the backbone of the C-
1230 terminal epitope residues was co-grafted as a discontinuous segment into the
1231 adjacent strand in the TOP7 backbone. With this strategy we circumvented a
1232 Rosetta sampling limitation, where often times is necessary an extensive set of
1233 constrains to achieve the desired backbone hydrogen-bond pairing on beta-
1234 strands (Marcos et al., 2017). After defining the motif consisting of the epitope
1235 plus the pairing strand and the sites of insertion on the TOP7 scaffold, FunFolDes
1236 was used to graft the motif.

1237

1238 Backbone flexibility was allowed for the terminal residues of the functional motif
1239 and a β-turn connection between the two strands was modelled during the
1240 folding process (*NubInitioMover*). During the folding process, 101F antibody was
1241 added to the simulation in order to limit the explored conformational space
1242 towards binding productive designs. Finally, the *NubInitioLoopClosureMover* was
1243 applied to ensure that a proper polypeptide chain was modelled and no chain-
1244 brakes remained, a total of 800 centroid models were generated after this stage.
1245 Next, we applied an RMSD filter to select scaffolds with similar topology to TOP7
1246 (< 1.5 Å) and a hydrogen bond long-range backbone score (HB_LR term) to
1247 favour the selection of proteins with proper beta-sheet pairing. The top 100
1248 models according the HB_LR score and that also fulfilled the RMSD filter, were
1249 then subjected to an iterative sequence-design relax protocol, alternating fixed
1250 backbone side-chain design and backbone relaxation using the *FastDesignMover*.
1251 Designable positions were limited to a subset of residues according to their

47

1252  position in the core or surface of the protein and secondary structure identity.
1253  Two different design strategies were pursued: I) partial design - amino acid
1254  identities of the C-terminal half of the protein (residues 45 through 92) were
1255  retained from TOP7 while allowing repacking of the side chains and backbone
1256  relaxation; II) full-design - the full sequence space in all residues of the structure
1257  (with the exception of the 101F epitope) was explored. No backbone or side
1258  chain movements were allowed in the 6-residue epitope segment whereas the
1259  adjacently paired β-strand was allowed to both mutate and relax. Tight Cα atom-
1260  pair distance constraints (standard deviation of 0.5 Å) were used to restrain
1261  movements of the entire sheet throughout the structural relaxation iterations.
1262
1263  From the 100 designs generated, only those that passed a structural filter
1264  requiring that 80% of secondary structure composition of the β-sheet after
1265  backbone relaxation were selected for further analysis. The 93 designs passing
1266  this filter were evaluated based on several metrics such as: REU, hydrogen-bond
1267  long-range backbone interactions and core packing. The best-scored designs
1268  were finally submitted to human-guided optimisation, mutations of single
1269  surface residues (1-3) and shortening of the connecting loop between the two
1270  inserted strands using the Rosetta Remodel application.
1271
1272  Interestingly, in an attempt to reproduce the same grafting exercise with
1273  *MotifGraftMover* (Silva et al., 2016), this resulted in non-resolvable chain breaks
1274  when trying to graft either the two segment-motif or the epitope alone into the
1275  TOP7 scaffold.
1276
1277  *Protein Expression and Purification*
1278
1279  DNA sequences of the designs were purchased from Twist Bioscience. For
1280  bacterial expression the DNA fragments were cloned via Gibson cloning into a
1281  pET21b vector containing a C-terminal His-tag and transformed into *E. coli*
1282  BL21(DE3). Expression was conducted in Terrific Broth supplemented with
1283  ampicillin (100 μg/ml). Cultures were inoculated at an $OD_{600}$ of 0.1 from an
1284  overnight culture and incubated at 37°C with a shaking speed of 220 rpm. After
1285  reaching $OD_{600}$ of 0.7, expression was induced by the addition of 1 mM IPTG and

cells were further incubated for 4-5h at 37°C. Cells were harvested by centrifugation and pellets were resuspended in lysis buffer (50 mM TRIS, pH 7.5, 500 mM NaCl, 5% Glycerol, 1 mg/ml lysozyme, 1 mM PMSF, 1 μg/ml DNase). Resuspended cells were sonicated and clarified by centrifugation. Ni-NTA purification of sterile-filtered (0.22 μm) supernatant was performed using a 1 ml His-Trap™ FF column on an ÄKTA pure system (GE healthcare). Bound proteins were eluted using an imidazole concentration of 300 mM. Concentrated proteins were further purified by size exclusion chromatography on a Superdex™ 75 300/10 or a Superdex™ Hiload 16/600 75 pg column (GE Healthcare) using PBS buffer (pH 7.4) as mobile phase.

For IgG expression, heavy and light chain DNA sequences were cloned separately into pFUSE-CHIg-hG1 (InvivoGen) mammalian expression vectors. Expression plasmids were co-transfected into HEK293-F cells in FreeStyle™ medium (Gibco™) using polyethylenimine (Polysciences) transfection. Supernatants were harvested after 1 week by centrifugation and purified using a 5 ml HiTrap™ Protein A HP column (GE Healthcare). Elution of bound proteins was accomplished using a 0.1 M glycine buffer (pH 2.7) and eluents were immediately neutralized by the addition of 1 M TRIS ethylamine (pH 9). The eluted IgGs were further purified by size exclusion chromatography on a Superdex 200 10/300 GL column (GE Healthcare) in PBS buffer (pH 7.4). Protein concentrations were determined by measuring the absorbance at 280 nm using the sequence calculated extinction coefficient on a Nanodrop (Thermo Scientific).

*Circular Dichroism (CD)*

Far-UV circular dichroism spectra of designed scaffolds were collected between a wavelength of 190 nm to 250 nm on a Jasco J-815 CD spectrometer in a 1 mm path-length quartz cuvette. Proteins were dissolved in PBS buffer (pH 7.4) at concentrations between 20 μM and 40 μM. Wavelength spectra were averaged from two scans with a scanning speed of 20 nm min$^{-1}$ and a response time of 0.125 sec. The thermal denaturation curves were collected by measuring the change in ellipticity at 220 nm from 20 to 95°C with 2 or 5 °C increments.

1319

1320 *Size-exclusion Chromatography combined with Multi-Angle Light-Scattering (SEC-*
1321 *MALS)*

1322

1323 Multi-angle light scattering was used to assess the monodispersity and molecular
1324 weight of the proteins. Samples containing between 50 -100 µg of protein in PBS
1325 buffer (pH 7.4) were injected into a Superdex™ 75 300/10 GL column (GE
1326 Healthcare) using an HPLC system (Ultimate 3000, Thermo Scientific) at a flow
1327 rate of 0.5 ml min$^{-1}$ coupled in-line to a multi-angle light scattering device
1328 (miniDAWN TREOS, Wyatt). Static light-scattering signal was recorded from
1329 three different scattering angles. The scatter data were analysed by ASTRA
1330 software (version 6.1, Wyatt)

1331

1332 *Surface Plasmon Resonance (SPR)*

1333

1334 To determine the dissociation constants of the designs to the mota or 101F
1335 antibody, surface plasmon resonance was used. Experiments were performed on
1336 a Biacore 8K at room temperature with HBS-EP+ running buffer (10 mM HEPES
1337 pH 7.4, 150 mM NaCl, 3mM EDTA, 0.005% v/v Surfactant P20) (GE Healthcare).
1338 Approximately 1200 response units of mota or 101F antibody were immobilized
1339 via amine coupling on the methylcarboxyl dextran surface of a CM5 chip (GE
1340 Healthcare). Varying protein concentrations were injected over the surface at a
1341 flow rate of 30 µl/min with a contact time of 120 sec and a following dissociation
1342 period of 400 sec. Following each injection cycle, ligand regeneration was
1343 performed using 3M MgCl$_2$ (GE Healthcare). Data analysis was performed using
1344 1:1 Langmuir binding kinetic fits within the Biacore evaluation software (GE
1345 Healthcare).

1346

1347 *Availability*

1348

1349 FunFolDes is available as part of the Rosetta software suite and is fully
1350 documented in the Rosetta Commons documentation website as one of the
1351 Composite Protocols. All data and scripts necessary to recreate the analysis and

1352 design simulations described in this work are available at

1353 https://github.com/lpdi-epfl/FunFolDesData.

1354

1355 *Contributions*

1356

1357 J.B. coded the algorithm described. A.S coded the *StructFragMover*. K.S., A.B., A.S.

1358 and F.S. performed computational design simulations. S.W., K.S, C.Y., A.B., F.S.,

1359 S.V., R. L., M. V. and S.R. contributed to experimental characterization of the

1360 designed proteins. J.B. and B.E.C. designed the study and wrote manuscript.

1361

1379

1380 *References*

1381

1382 Alford, R. F., Leaver-Fay, A., Jeliazkov, J. R., O'Meara, M. J., DiMaio, F. P., Park, H., . .
1383 . Gray, J. J. (2017). The Rosetta All-Atom Energy Function for
1384 Macromolecular Modeling and Design. *J Chem Theory Comput, 13*(6),
1385 3031-3048. doi:10.1021/acs.jctc.7b00125

1386  Aragues, R., Sali, A., Bonet, J., Marti-Renom, M. A., & Oliva, B. (2007).
1387        Characterization of protein hubs by inferring interacting motifs from
1388        protein interactions. *PLoS Comput Biol, 3*(9), 1761-1771.
1389        doi:10.1371/journal.pcbi.0030178
1390  Azoitei, M. L., Correia, B. E., Ban, Y. E., Carrico, C., Kalyuzhniy, O., Chen, L., . . .
1391        Schief, W. R. (2011). Computation-guided backbone grafting of a
1392        discontinuous motif onto a protein scaffold. *Science, 334*(6054), 373-376.
1393        doi:10.1126/science.1209368
1394  Boschek, C. B., Apiyo, D. O., Soares, T. A., Engelmann, H. E., Pefaur, N. B.,
1395        Straatsma, T. P., & Baird, C. L. (2009). Engineering an ultra-stable affinity
1396        reagent based on Top7. *Protein Engineering Design & Selection, 22*(5),
1397        325-332. doi:10.1093/protein/gzp007
1398  Bowers, P. M., Strauss, C. E., & Baker, D. (2000). De novo protein structure
1399        determination using sparse NMR data. *J Biomol NMR, 18*(4), 311-318.
1400  Chakrabarti, K. S., Agafonov, R. V., Pontiggia, F., Otten, R., Higgins, M. K., Schertler,
1401        G. F. X., . . . Kern, D. (2016). Conformational Selection in a Protein-Protein
1402        Interaction Revealed by Dynamic Pathway Analysis. *Cell Rep, 14*(1), 32-
1403        42. doi:10.1016/j.celrep.2015.12.010
1404  Chevalier, A., Silva, D. A., Rocklin, G. J., Hicks, D. R., Vergara, R., Murapa, P., . . .
1405        Baker, D. (2017). Massively parallel de novo protein design for targeted
1406        therapeutics. *Nature, 550*(7674), 74-79. doi:10.1038/nature23912
1407  Coluzza, I. (2017). Computational protein design: a review. *J Phys Condens
1408        Matter, 29*(14), 143001. doi:10.1088/1361-648X/aa5c76
1409  Correia, B. E., Ban, Y. E., Friend, D. J., Ellingson, K., Xu, H., Boni, E., . . . Schief, W. R.
1410        (2011). Computational protein design using flexible backbone remodeling
1411        and resurfacing: case studies in structure-based antigen design. *J Mol Biol,
1412        405*(1), 284-297. doi:10.1016/j.jmb.2010.09.061
1413  Correia, B. E., Ban, Y. E., Holmes, M. A., Xu, H., Ellingson, K., Kraft, Z., . . . Schief, W.
1414        R. (2010). Computational design of epitope-scaffolds allows induction of
1415        antibodies specific for a poorly immunogenic HIV vaccine epitope.
1416        *Structure, 18*(9), 1116-1126. doi:10.1016/j.str.2010.06.010
1417  Correia, B. E., Bates, J. T., Loomis, R. J., Baneyx, G., Carrico, C., Jardine, J. G., . . .
1418        Schief, W. R. (2014). Proof of principle for epitope-focused vaccine design.
1419        *Nature, 507*(7491), 201-206. doi:10.1038/nature12966
1420  Cross, L. L., Paudyal, R., Kamisugi, Y., Berry, A., Cuming, A. C., Baker, A., &
1421        Warriner, S. L. (2017). Towards designer organelles by subverting the
1422        peroxisomal import pathway. *Nat Commun, 8*(1), 454.
1423        doi:10.1038/s41467-017-00487-7
1424  Dawson, N. L., Lewis, T. E., Das, S., Lees, J. G., Lee, D., Ashford, P., . . . Sillitoe, I.
1425        (2017). CATH: an expanded resource to predict protein function through
1426        structure and sequence. *Nucleic Acids Res, 45*(D1), D289-D295.
1427        doi:10.1093/nar/gkw1098
1428  Eddy, S. R. (2011). Accelerated Profile HMM Searches. *PLoS Comput Biol, 7*(10),
1429        e1002195. doi:10.1371/journal.pcbi.1002195
1430  Fallas, J. A., Ueda, G., Sheffler, W., Nguyen, V., McNamara, D. E., Sankaran, B., . . .
1431        Baker, D. (2017). Computational design of self-assembling cyclic protein
1432        homo-oligomers. *Nat Chem, 9*(4), 353-360. doi:10.1038/nchem.2673
1433  Fleishman, S. J., Leaver-Fay, A., Corn, J. E., Strauch, E. M., Khare, S. D., Koga, N., . . .
1434        Baker, D. (2011). RosettaScripts: a scripting language interface to the

Rosetta macromolecular modeling suite. *PLoS One, 6*(6), e20161. doi:10.1371/journal.pone.0020161

Garcia-Garcia, J., Bonet, J., Guney, E., Fornes, O., Planas, J., & Oliva, B. (2012). Networks of ProteinProtein Interactions: From Uncertainty to Molecular Details. *Mol Inform, 31*(5), 342-362. doi:10.1002/minf.201200005

Guntas, G., Purbeck, C., & Kuhlman, B. (2010). Engineering a protein-protein interface using a computationally designed library. *Proc Natl Acad Sci U S A, 107*(45), 19296-19301. doi:10.1073/pnas.1006528107

Henikoff, S., & Henikoff, J. G. (1992). Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci U S A, 89*(22), 10915-10919.

Hill, R. B., Raleigh, D. P., Lombardi, A., & DeGrado, W. F. (2000). De novo design of helical bundles as models for understanding protein folding and function. *Acc Chem Res, 33*(11), 745-754.

Hu, X., Wang, H., Ke, H., & Kuhlman, B. (2007). High-resolution design of a protein loop. *Proc Natl Acad Sci U S A, 104*(45), 17668-17673. doi:10.1073/pnas.0707977104

Jiang, L., Althoff, E. A., Clemente, F. R., Doyle, L., Rothlisberger, D., Zanghellini, A., . . . Baker, D. (2008). De novo computational design of retro-aldol enzymes. *Science, 319*(5868), 1387-1391. doi:10.1126/science.1152692

Joh, N. H., Wang, T., Bhate, M. P., Acharya, R., Wu, Y., Grabe, M., . . . DeGrado, W. F. (2014). De novo design of a transmembrane Zn(2)(+)-transporting four-helix bundle. *Science, 346*(6216), 1520-1524. doi:10.1126/science.1261172

Jones, D. T. (1999). Protein secondary structure prediction based on position-specific scoring matrices. *J Mol Biol, 292*(2), 195-202. doi:10.1006/jmbi.1999.3091

Kim, D. E., Blum, B., Bradley, P., & Baker, D. (2009). Sampling bottlenecks in de novo protein structure prediction. *J Mol Biol, 393*(1), 249-260. doi:10.1016/j.jmb.2009.07.063

Koga, N., Tatsumi-Koga, R., Liu, G., Xiao, R., Acton, T. B., Montelione, G. T., & Baker, D. (2012). Principles for designing ideal protein structures. *Nature, 491*(7423), 222-227. doi:10.1038/nature11600

Kries, H., Blomberg, R., & Hilvert, D. (2013). De novo enzymes by computational design. *Curr Opin Chem Biol, 17*(2), 221-228. doi:10.1016/j.cbpa.2013.02.012

Kuhlman, B., & Baker, D. (2000). Native protein sequences are close to optimal for their structures. *Proc Natl Acad Sci U S A, 97*(19), 10383-10388.

Kuhlman, B., & Baker, D. (2004). Exploring folding free energy landscapes using computational protein design. *Curr Opin Struct Biol, 14*(1), 89-95. doi:10.1016/j.sbi.2004.01.002

Kuhlman, B., Dantas, G., Ireton, G. C., Varani, G., Stoddard, B. L., & Baker, D. (2003). Design of a novel globular protein fold with atomic-level accuracy. *Science, 302*(5649), 1364-1368. doi:10.1126/science.1089427

Kulkarni, M. R., Islam, M. M., Numoto, N., Elahi, M., Mahib, M. R., Ito, N., & Kuroda, Y. (2015). Structural and biophysical analysis of sero-specific immune responses using epitope grafted Dengue ED3 mutants. *Biochim Biophys Acta, 1854*(10 Pt A), 1438-1443. doi:10.1016/j.bbapap.2015.07.004

Kvansakul, M., Wei, A. H., Fletcher, J. I., Willis, S. N., Chen, L., Roberts, A. W., . . . Colman, P. M. (2010). Structural basis for apoptosis inhibition by Epstein-

Barr virus BHRF1. *PLoS Pathog, 6*(12), e1001236. doi:10.1371/journal.ppat.1001236

Lange, O. F., Lakomek, N. A., Fares, C., Schroder, G. F., Walter, K. F., Becker, S., . . . de Groot, B. L. (2008). Recognition dynamics up to microseconds revealed from an RDC-derived ubiquitin ensemble in solution. *Science, 320*(5882), 1471-1475. doi:10.1126/science.1157092

Lartigue, A., Campanacci, V., Roussel, A., Larsson, A. M., Jones, T. A., Tegoni, M., & Cambillau, C. (2002). X-ray structure and ligand binding study of a moth chemosensory protein. *J Biol Chem, 277*(35), 32094-32098. doi:10.1074/jbc.M204371200

Marcos, E., Basanta, B., Chidyausiku, T. M., Tang, Y., Oberdorfer, G., Liu, G., . . . Baker, D. (2017). Principles for designing proteins with cavities formed by curved beta sheets. *Science, 355*(6321), 201-206. doi:10.1126/science.aah7389

McLellan, J. S., Chen, M., Chang, J. S., Yang, Y., Kim, A., Graham, B. S., & Kwong, P. D. (2010). Structure of a major antigenic site on the respiratory syncytial virus fusion glycoprotein in complex with neutralizing antibody 101F. *J Virol, 84*(23), 12236-12244. doi:10.1128/JVI.01579-10

McLellan, J. S., Chen, M., Kim, A., Yang, Y., Graham, B. S., & Kwong, P. D. (2010). Structural basis of respiratory syncytial virus neutralization by motavizumab. *Nat Struct Mol Biol, 17*(2), 248-250. doi:10.1038/nsmb.1723

Murphy, G. S., Mills, J. L., Miley, M. J., Machius, M., Szyperski, T., & Kuhlman, B. (2012). Increasing sequence diversity with flexible backbone protein design: the complete redesign of a protein hydrophobic core. *Structure, 20*(6), 1086-1096. doi:10.1016/j.str.2012.03.026

Procko, E., Berguig, G. Y., Shen, B. W., Song, Y., Frayo, S., Convertine, A. J., . . . Baker, D. (2014). A computationally designed inhibitor of an Epstein-Barr viral Bcl-2 protein induces apoptosis in infected cells. *Cell, 157*(7), 1644-1656. doi:10.1016/j.cell.2014.04.034

Richter, F., Leaver-Fay, A., Khare, S. D., Bjelic, S., & Baker, D. (2011). De novo enzyme design using Rosetta3. *PLoS One, 6*(5), e19230. doi:10.1371/journal.pone.0019230

Rohl, C. A., & Baker, D. (2002). De novo determination of protein backbone structure from residual dipolar couplings using Rosetta. *J Am Chem Soc, 124*(11), 2723-2729.

Rohl, C. A., Strauss, C. E., Misura, K. M., & Baker, D. (2004). Protein structure prediction using Rosetta. *Methods Enzymol, 383*, 66-93. doi:10.1016/S0076-6879(04)83004-0

Rose, P. W., Prlic, A., Altunkaya, A., Bi, C., Bradley, A. R., Christie, C. H., . . . Burley, S. K. (2017). The RCSB protein data bank: integrative view of protein, gene and 3D structural information. *Nucleic Acids Res, 45*(D1), D271-D281. doi:10.1093/nar/gkw1000

Schreiber, G., & Fleishman, S. J. (2013). Computational design of protein-protein interactions. *Curr Opin Struct Biol, 23*(6), 903-910. doi:10.1016/j.sbi.2013.08.003

Silva, D. A., Correia, B. E., & Procko, E. (2016). Motif-Driven Design of Protein-Protein Interfaces. *Methods Mol Biol, 1414*, 285-304. doi:10.1007/978-1-4939-3569-7_17

1533   Simons, K. T., Bonneau, R., Ruczinski, I., & Baker, D. (1999). Ab initio protein
1534        structure prediction of CASP III targets using ROSETTA. *Proteins, Suppl 3*,
1535        171-176.
1536   Simons, K. T., Ruczinski, I., Kooperberg, C., Fox, B. A., Bystroff, C., & Baker, D.
1537        (1999). Improved recognition of native-like protein structures using a
1538        combination of sequence-dependent and sequence-independent features
1539        of proteins. *Proteins, 34*(1), 82-95.
1540   Soares, T. A., Boschek, C. B., Apiyo, D., Baird, C., & Straatsma, T. P. (2010).
1541        Molecular basis of the structural stability of a Top7-based scaffold at
1542        extreme pH and temperature conditions. *Journal of Molecular Graphics &*
1543        *Modelling, 28*(8), 755-765. doi:10.1016/j.jmgm.2010.01.013
1544   Stein, A., & Kortemme, T. (2013). Improvements to robotics-inspired
1545        conformational sampling in rosetta. *PLoS One, 8*(5), e63090.
1546        doi:10.1371/journal.pone.0063090
1547   Street, A. G., & Mayo, S. L. (1999). Computational protein design. *Structure, 7*(5),
1548        R105-109.
1549   Tokuriki, N., Stricher, F., Serrano, L., & Tawfik, D. S. (2008). How protein stability
1550        and new functions trade off. *PLoS Comput Biol, 4*(2), e1000002.
1551        doi:10.1371/journal.pcbi.1000002
1552   Tyka, M. D., Keedy, D. A., Andre, I., Dimaio, F., Song, Y., Richardson, D. C., . . . Baker,
1553        D. (2011). Alternate states of proteins revealed by detailed energy
1554        landscape mapping. *J Mol Biol, 405*(2), 607-618.
1555        doi:10.1016/j.jmb.2010.11.008
1556   Viana, I. F. T., Soares, T. A., Lima, L. F. O., Marques, E. T. A., Krieger, M. A., Dhalia,
1557        R., & Lins, R. D. (2013). De novo design of immunoreactive conformation-
1558        specific HIV-1 epitopes based on Top7 scaffold. *Rsc Advances, 3*(29),
1559        11790-11800. doi:10.1039/c3ra41562g
1560   Vreven, T., Moal, I. H., Vangone, A., Pierce, B. G., Kastritis, P. L., Torchala, M., . . .
1561        Weng, Z. (2015). Updates to the Integrated Protein-Protein Interaction
1562        Benchmarks: Docking Benchmark Version 5 and Affinity Benchmark
1563        Version 2. *J Mol Biol, 427*(19), 3031-3041. doi:10.1016/j.jmb.2015.07.016
1564   Wang, C., Bradley, P., & Baker, D. (2007). Protein-protein docking with backbone
1565        flexibility. *J Mol Biol, 373*(2), 503-519. doi:10.1016/j.jmb.2007.07.050
1566   Yu, F., Cangelosi, V. M., Zastrow, M. L., Tegoni, M., Plegaria, J. S., Tebo, A. G., . . .
1567        Pecoraro, V. L. (2014). Protein design: toward functional metalloenzymes.
1568        *Chem Rev, 114*(7), 3495-3578. doi:10.1021/cr400458x
1569   Zhou, J., & Grigoryan, G. (2015). Rapid search for tertiary fragments reveals
1570        protein sequence-structure relationships. *Protein Sci, 24*(4), 508-524.
1571        doi:10.1002/pro.2610
1572