

1 **High-resolution global peptide-protein docking using fragments-** 2 **based PIPER-FlexPepDock**

3 Nawsad Alam¹, Oriel Goldstein², Bing Xia³, Kathryn A. Porter³, Dima Kozakov^{4,5,6}, Ora
4 Schueler-Furman¹

5
6 ¹Department of Microbiology and Molecular Genetics, Institute for Medical Research Israel-Canada,
7 Hadassah Medical School, The Hebrew University, Jerusalem 91120, Israel

8 ²School of Computer Sciences and Engineering, The Hebrew University, Jerusalem 9190416,
9 Israel

10 ³Department of Biomedical Engineering, Boston University, 44 Cummington Street Boston, MA 02215,
11 USA

12 ⁴Department of Applied Mathematics and Statistics, Stony Brook University, Stony Brook, NY 11794

13 ⁵Laufer Center for Physical and Quantitative Biology, Stony Brook University, Stony Brook, NY 11794;

14 ⁶Institute for Advanced Computational Sciences, Stony Brook University, Stony Brook, NY 11794

15

16 **Abstract**

17 Peptide-protein interactions contribute a significant fraction of the protein-protein
18 interactome. Accurate modeling of these interactions is challenging due to the vast
19 conformational space associated with interactions of highly flexible peptides with large
20 receptor surfaces. To address this challenge we developed a fragment based high-
21 resolution peptide-protein docking protocol. By streamlining the Rosetta fragment picker
22 for accurate peptide fragment ensemble generation, the PIPER docking algorithm for
23 exhaustive fragment-receptor rigid-body docking and Rosetta FlexPepDock for flexible
24 full-atom refinement of PIPER docked models, we successfully addressed the challenge
25 of accurate and efficient global peptide-protein docking at high-resolution with
26 remarkable accuracy. Validation on a representative set of solved peptide-protein
27 complex structures demonstrates the accuracy and robustness of our approach, and
28 opens up the way to high-resolution modeling of many more peptide-protein interactions
29 and to the detailed study of peptide-protein association in general. PIPER-FlexPepDock
30 is freely available to the academic community as a server at
31 <http://piperfpd.furmanlab.cs.huji.ac.il>.

32

33 **Introduction**

34 Proteins are the workhorses inside living cells, and interactions among them are critical
35 for various important biological processes ¹. A significant fraction of these interactions
36 (15-40%) ² are peptide mediated, where a short stretch of residues from one partner
37 contributes most to its binding to the other. Such short peptidic regions, also termed
38 short linear interacting motifs (SLIMs) are often found embedded inside disordered
39 regions of intrinsically disordered proteins (IDPs) ^{2,3}, or appear as flexible linkers
40 connecting domains ⁴ and as flexible loops tethered to rigid segments ⁵.

41 The development of accurate structure based modeling tools is critical for atomic level
42 understanding of peptide-protein interactions, to allow the manipulation of known
43 interactions, to discover yet unknown peptide-protein interactions and networks, and to
44 provide starting points for the design of novel peptides and related molecules to target
45 specific systems of pharmacological interest ⁶. A number of computational tools have
46 been developed to assist the characterization of peptide-protein interactions, including
47 the prediction of peptide binding sites ⁷⁻⁹, refinement of coarse peptide-protein models
48 ¹⁰, folding and docking on a known binding site ¹¹ and most challenging of all, global
49 peptide-protein docking with no prior information about the peptide structure and the
50 binding site ¹²⁻¹⁷. The challenges associated with the global docking of flexible peptides
51 have been addressed in different ways, by reducing the conformational space to be
52 sampled both for the internal degrees of freedom of the peptide as well as its rigid-body
53 orientations on the receptor surface. For peptide docking within the HADDOCK docking
54 framework ¹², the peptide backbone is represented by idealized conformation(s), such
55 as alpha helix, beta strand and polyproline-II, followed by rigid-body, semi-flexible and
56 fully-flexible docking with explicit solvation ¹⁸. The pepATTRACT protocol ^{13,19} uses the
57 same approach to represent the peptide, followed by coarse-grained rigid-body docking
58 and flexible full-atom refinement. The AnchorDock protocol uses molecular dynamics
59 simulations to generate a set of plausible peptide conformations, which are then docked
60 using anchor-driven simulated annealing molecular dynamics around predicted
61 anchoring spots on the receptor ¹⁴. The CABS-dock protocol uses randomly generated
62 peptide conformations based on either predicted or known secondary structure,
63 randomly orients these peptides over the receptor surface, and refines them using

64 replica exchange Monte Carlo dynamics¹⁵. The MDockPep protocol¹⁶ uses peptide
65 sequence similarity to extract fragments from high resolution protein structures, which
66 are further refined using MODELLER²⁰ to generate plausible peptide conformations,
67 and then docked onto the receptor using rigid-body docking and flexible docking with
68 AutoDock Vina²¹. The recently published IDP-LZerD protocol models the binding of
69 long disordered segments to structured proteins using the Rosetta fragment picker
70 protocol²² to generate fragments of 9-residue overlapping windows followed by LZerD
71²³ rigid-body docking and molecular dynamics refinement¹⁷. Finally, we have recently
72 advanced a novel, global motif-based peptide fragment docking approach, PeptiDock²⁴,
73 in which peptide binding motif information rather than secondary structure propensity is
74 used to extract fragments from the Protein Data Bank (PDB²⁵), which are then docked
75 to the receptor using PIPER rigid body docking²⁶, followed by minimization using
76 CHARMM²⁷.

77 These significant recent advances in global peptide docking notwithstanding, present
78 approaches are still limited in their modeling quality and general applicability, and there
79 is ample room for improvements that would enable the detailed high-resolution study of
80 more peptide-protein interactions with higher accuracy. Here we describe PIPER-
81 FlexPepDock, a successful effort toward the development of such a robust, highly
82 accurate, global peptide-protein docking protocol. By integrating accurate peptide
83 fragment ensemble generation using the Rosetta fragment picker²², fast and exhaustive
84 fragment-receptor rigid-body docking using PIPER docking²⁸, and flexible full-atom
85 refinement of coarse PIPER models using Rosetta FlexPepDock¹⁰, we were able to
86 sample both the peptide backbone conformational states, as well as the landscape of
87 the peptide-receptor interactions very efficiently and with much higher accuracy than
88 current protocols: on a non-redundant dataset of peptide-protein complexes (**Table 1**
89 below), PIPER-FlexPepDock generates for about half models within 2.5 Å ligand RMSD
90 (2.0 Å, when restricted to motif regions where available), more than twice as many as
91 for existing peptide docking protocols such as pepATTRACT¹³ (among the 10 top-
92 ranked predictions; **Table 2** below).

93 Our results highlight the relevance of representing the peptide as a set of fragments that
94 can be exhaustively docked as rigid bodies onto the receptor structure and

95 subsequently refined using an accurate refinement protocol. They reinforce the
96 underlying biophysical model of a conformer ensemble of the free peptide that already
97 samples the bound conformation (at least in the encounter-complex, protein-like
98 environment) and involves only limited induced fit, not unlike to the classical association
99 between preformed protein domains. As a result, PIPER-FlexPepDock brings into reach
100 the study and targeted manipulation of a range of additional peptide-mediated
101 interactions not accessible before due to limitations in sampling and/or accuracy.

102

103 **Results**

104 **Overview of the PIPER-FlexPepDock protocol (Figure 1)**

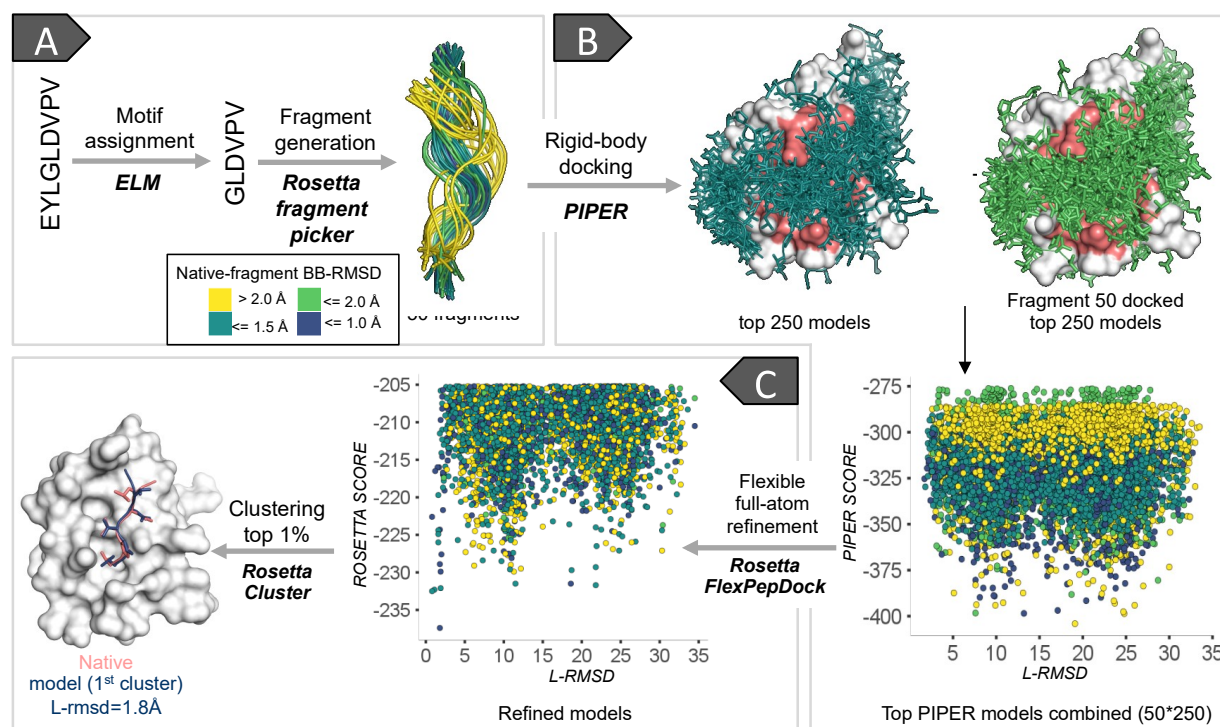
105 **Step A | Generation of fragment set to represent the peptide conformer ensemble:**

106 In a previous study we have shown that the bound peptide conformation can be well
107 represented by extraction of short fragments from the PDB based on information of
108 known binding sequence motifs²⁴. Here we have generalized this approach beyond
109 motifs, using fragment libraries selected by the Rosetta fragment picker protocol²²
110 based on sequence and secondary structure similarity (see **Methods**). The coordinates
111 of the top 50 mapped fragments are extracted from the PDB, including both backbone
112 and side-chain atoms, and non-identical residues in the extracted fragments are
113 mutated to the desired sequence. This set of fragments adequately represents the
114 peptide conformational ensemble, sampling also its receptor bound conformation (see
115 below). The peptide may be trimmed in cases where information is available about the
116 range of the binding segment (from motif databases such as the Eukaryotic Linear Motif
117 (ELM) resource^{29,30}, literature, or experiments such as alanine scanning), since
118 fragments generated for shorter peptide sequences are usually better representative
119 than longer fragments (as, e.g., for loop modeling³¹), and fraying ends beyond the motif
120 may contribute less to determine critical binding details.

121

122 **Step B | Fragment rigid-body docking using PIPER:** Each of the fragments is docked
123 onto the receptor structure using PIPER, an exhaustive Fast Fourier Transform (FFT)-
124 based rigid body docking algorithm²⁸, as implemented previously for PeptiDock²⁴ (see
125 **Methods**), and top ranking fragment orientations from each docking run are collected

126 and combined together. These models are of low resolution as no flexibility is included
127 in the PIPER algorithm, and therefore ranked using a soft potential that allows a certain
128 degree of steric clashes to overcome the limitations of rigid-body only docking.
129



130 **Figure 1. Overview of the PIPER-FlexPepDock peptide docking protocol.** Example shown:
131 PDZ domain-peptide interaction [PDB IDs of receptor structure 1MFG (bound) and 2H3L (free)].
132 For a given receptor structure and peptide sequence, the divide and conquer strategy involves
133 first the description of the peptide as an **ensemble of fragments (A)**, their fast and exhaustive
134 **rigid body docking** (using PIPER) onto the whole receptor (binding site region is shaded
135 salmon) **(B)**, and subsequent **high-resolution refinement** (using Rosetta FlexPepDock; the top
136 5000 models are included in the plot) **(C)**, followed by clustering and selection of top ranking
137 representatives. Fragments are colored according to their similarity to the native bound peptide
138 conformation. L-RMSD: Ligand root mean square deviation from crystal structure; see text for
139 more details.

140
141 **Step C | FlexPepDock refinement of PIPER models and selection of final models:**

142 Each of the PIPER models is refined by a single fully flexible refinement run using the
143 Rosetta FlexPepDock Refinement algorithm¹⁰ (see **Methods**). The top ranking refined
144 models are clustered (as in Gray *et al.*³²), clusters are ranked based on the reweighted
145 score of the best scoring model in each cluster (as in Raveh *et al.*¹¹), and the top 10
146 ranked cluster representatives are selected as prediction (following the CAPRI scheme
147 that accepts 10 models³³).

148

149 **Initial calibration of the PIPER-FlexPepDock on a small set of protein-peptide**
 150 **complexes**

151 Motivated by our recent advance in global peptide docking using a motif-focused
 152 approach ²⁴ we ventured into the development of a more generalized protocol. We
 153 initially calibrated our docking approach on a small representative set of nine peptide–
 154 protein complexes (highlighted in bold in **Table 1**; see also **Supplementary Table**
 155 **S1A**). We trimmed the peptide based on the motif defined in ELM, where available. For
 156 all complexes impressive modeling accuracy was achieved for this new global docking
 157 approach (within $\leq 2.5\text{\AA}$ Ligand RMSD models among the top 10 ranking clusters; **Table**
 158 **1**). For the full length peptides modeling near-native models were obtained for 5/9
 159 cases, highlighting the benefits for motif (or shorter peptide sequence) focused
 160 modeling, due to better fragment quality compared to the corresponding full-length
 161 peptides (**Table 1**). Encouraged by these initial results, we proceeded to the validation
 162 of our protocol on a larger and representative set of peptide-protein complexes (**Table 1**
 163 and **Supplementary Table S1B**).

164

165 **Table 1.** Benchmark of peptide-protein complexes used in this study (non-redundant
 166 set; see **Table S1C** for full set). PDB ids of the initial calibration set are highlighted in
 167 bold.

PDB ID	Peptide sequence ^a / secondary structure	Fragment similarity ^b	PIPER ^c			PIPER-FPD (<i>motif</i>)			PIPER-FPD (<i>full peptide</i>)		
			L ^d	L	I ^e	Fnat	L	I	Fnat		
Known binding motif (n=12)											
1CZY:CE	<u>POQATDD</u>	2.2(2.5)	17.6	1.6	0.6	0.86	2.4	0.9	0.76		
1CA4:A	CEECCCC										
1EG4:AP	NMTPYRSPPPYVP	0.6(3.0)	21.6	12.9	4.1	0.19	29.3	11.5	0.00		
1EG3:A	TTTTTTCCECCCC										
1ELW:AC	<u>GPTIEEVD</u>	1.0(2.4)	3.2	0.8	2.5	0.75	2.6	3.0	0.71		
1A17:A	CCCCCCCC										
1JD5:AB	<u>AIAFYIPD</u>	0.7(2.3)	2.8	1.2	0.5	0.88	8.2	2.9	0.19		
1JD4:A	CEEEETCC										
1JWG:BD	<u>DEDLLHI</u>	2.8(2.9)	3.4	2.2	0.8	0.90	2.2	0.8	0.90		
1JWF:A	CCCCCCC										
1MFG:AB	<u>EYLGLDVPV</u>	1.5(2.7)	3.1	1.8	0.8	0.73	8.7	3.1	0.29		
2H3L:A	CCCCCEEC										
1NTV:AB	<u>NFDNPVYRKT</u>	2.7(3.3)	5.6	5.2	1.7	0.43	3.8 ^f	1.5	0.52		
1P3R:B	CEETTTTCCC										
1RXZ:AB	<u>KSTQATLERWF</u>	2.4(4.0)	7.4	5.0	1.9	0.31	3.2	1.7	0.39		

1RWZ:A	CEEECTTTTTC								
1SSH:AB	GPPPAMPARPT	1.1(2.3)	2.6	7.6	2.9	0.56	1.9^f	1.1	0.87
1OOT:A	CCCCCCCCCCC								
1X2R:AB	LDEETGEFL	0.2(0.5)	1.1	1.3	0.5	0.74	1.7	0.6	0.72
1X2J:A	CTTTTTCCC								
2A3I:AB	QQKSLLOQLLTE	0.3(3.5)	1.6	1.0	0.4	0.93	4.8	2.1	0.72
2AA2:A	CCCCHHHHHHHC								
2CCH:DF	HTLKGRRLVFDN	1.8(4.6)	3.4	1.0	0.4	0.91	3.9	1.6	0.67
1H1R:B	TTTTCCCCCCCC								
No known binding motif (n=15)									
1AWR:CI	HAGPIA	1.6	6.5				1.3	0.5	0.97
2ALF:A	CCCCC								
1ER8:EI	HPFHLLVY	1.9	4.1				1.2	0.8	0.80
4PAE:A	CCCBCCBC								
1LVM:AE	ENLYFQ	1.9	2.3				1.4	0.6	0.91
1LVB:B	CCEEEC								
1NVR:AB	ASVSA	2.0	5.5				7.1	2.2	0.56
2QHN:A	CEEEC								
1NX1:AC	DAIDALSSDFT	1.8	1.9				1.3	0.9	0.80
1ALV:A	HHHHHHHHHCC								
1OU8:BD	GAANDENY	3.1	6.1				6.4	2.4	0.39
1OU9:A	CCCCCCCC								
1U00:AP	ELPPVKIHC	2.3	6.0				2.1	1.7	0.71
2V7Y:A	CCCCEECCC								
2B9H:AC	RRNLKGLNLSLH	3.2	16.4				15.8	6.3	0.18
2B9F:A	CCTTTTCCCCC								
2C3I:BA	KRRRHPSG	2.5	3.5				8.7	2.5	0.26
2J2I:B	CCCCCCCC								
2DS8:BP	ALRVVK	1.4	4.7				1.2	0.6	0.86
2DS7:A	CCEECC								
2FMF:AB	QDQVDDLLDSLGF	1.1	5.2				1.3	0.7	1.00
1JBE:A	HHHHHHHHHHHCC								
2H9M:CD	ARTKQ	2.2	4.7				4.0	1.1	0.37
2H14:A	TTTTTC								
2HPL:AB	DDLYG	1.6	3.0				1.4	0.5	0.91
2HPJ:A	CCCCC								
2O02:AP	GLLDALDLAS	2.4	2.8				3.4	1.2	0.74
2BQ0:A	THHHHHCCCC								
3D1E:AP	GQLGLF	2.5	4.1				6.7	2.0	0.51
3D1G:A	CBCCCC								

168 ^a Motif (as defined by the ELM database ²⁹) is underlined (motif details are provided in the
169 **Supplementary Table S1B**).

170 ^b Similarity between fragments and bound peptide conformation: Median backbone RMSD (Å) (in
171 parentheses: results for full peptide).

172 ^c Results for PIPER simulations are given for the motif / full peptide for known/unknown motifs,
173 respectively. The models are selected as in PeptiDock ²⁴ (without the minimization step).

174 ^{d,e} Modeling accuracy: L – Ligand RMSD (**models within 2.5 Å are highlighted in bold**)^c; I –
175 Interface RMSD^d. Defined as in the CAPRI experiment ^{34,35}.

176 ^f Complexes for which docking of the full peptide provides better models than docking of the
177 motif only.

178

179

180

181 **Assessment of peptide docking performance**

182 We assessed the performance of PIPER-FlexPepDock on a larger, non-redundant set
183 consisting of 27 complexes (non-redundant at the domain level, as defined by CATH³⁶;
184 see **Methods**), among them 12 with reported binding motif. The benchmark is
185 summarized in **Table 1 (Supplementary Table S1C)** provides results for a redundant
186 set of 42 complexes used in previous studies, as well as additional details, including
187 performance of other approaches for comparison).

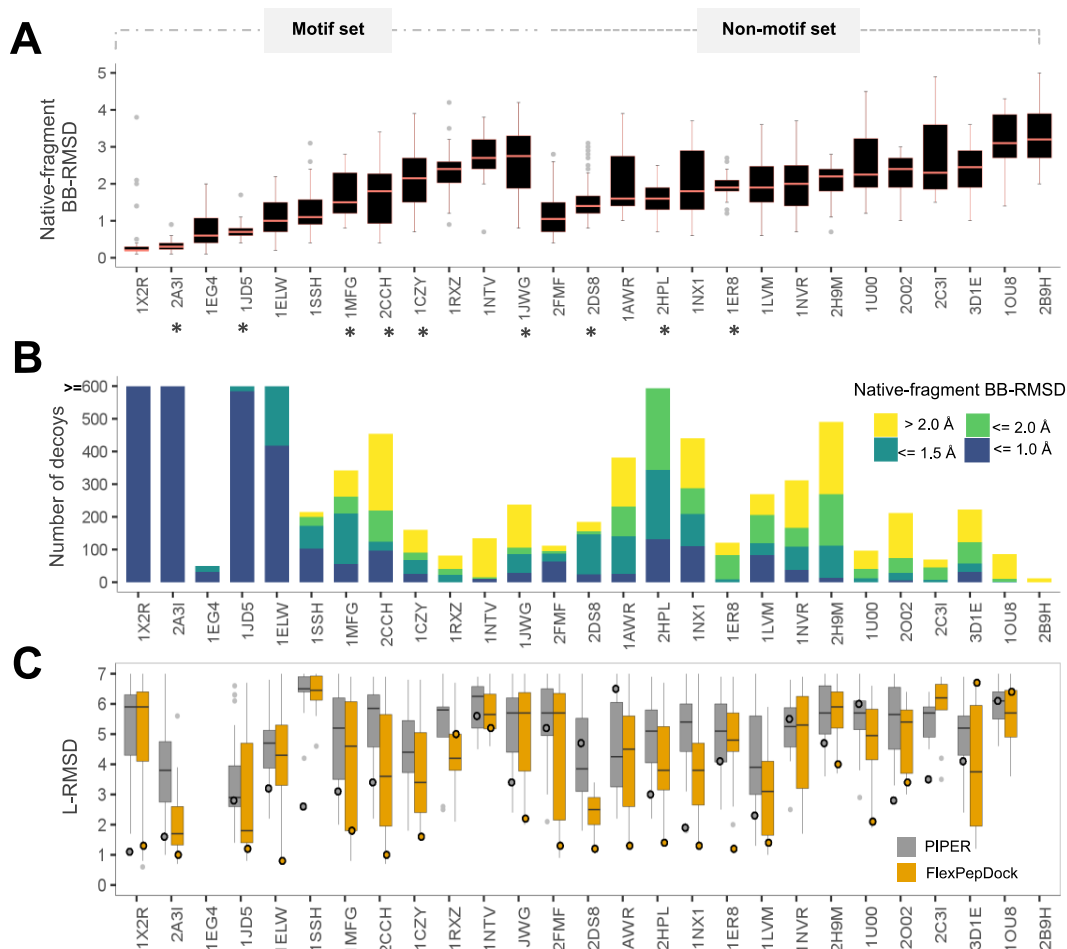
188

189 **Representing the peptide conformational states using fragments**

190 Fragments derived from solved protein structures contain valuable information about the
191 local structural context that can be used to efficiently reduce sampling space for various
192 modeling applications, including e.g. *ab initio* protein folding³⁷ and loop building^{31,38,39}.
193 In our protocol we use the Rosetta fragment picker protocol²² to generate fragments
194 consistent with both the peptide sequence and the (predicted) secondary structure (See
195 **Methods**).

196 How accurately do the fragments represent the peptide conformational states? Most
197 importantly, how similar are peptide conformations to the one adapted when bound to
198 their receptor? A significant representation of similar fragments could guarantee that,
199 when docked with high density in the binding site using an exhaustive but accurate
200 rigid-body docking algorithm, they could efficiently be refined to high resolution using an
201 accurate refinement algorithm such as Rosetta FlexPepDock. To assess the quality of
202 the fragments (*i.e.*, their coverage of the bound conformation) we analyzed the
203 distribution of backbone RMSDs of the fragments relative to the bound peptide
204 conformation. Reassuringly, the fragment pool generated using the Rosetta fragment
205 picker protocol represents the bound like peptide conformation with high accuracy in our
206 benchmark of 27 peptide-protein complexes (**Figure 2A**: median backbone RMSD
207 within 2.0 Å for 15 out of the 27 cases, with average backbone RMSD of the best ten
208 fragments within 1.0 and 1.5 Å for 14 and 21 cases, respectively). The best accuracy is
209 achieved for helical peptide motifs (e.g., the helical nuclear receptor box motif in 2A3I⁴⁰;
210 for helical peptides with coiled terminus segments such as 2FMF⁴¹ and 1NX1⁴² the
211 median backbone RMSD is slightly higher). Even for the remaining peptides the
212 fragment ensemble is often composed of a significant portion of bound like

213 representatives. The worst representation of bound-like peptides is obtained for few
 214 longer coil peptides, such as 2B9H⁴³, which defines the limitation of the fragment picker
 215 protocol for longer sequences. In such cases, trimming the peptide might improve the
 216 quality significantly.



217
 218 **Figure 2. Assessment of performance of the different steps of PIPER-FlexPepDock (A)**
 219 **Fragment quality:** distribution of fragment backbone RMSDs relative to the native bound
 220 peptide conformation (defined as fragment quality). PDBs with and without motif information are
 221 grouped separately. The initial calibration set is marked with asterisks (*). **(B) PIPER rigid body**
 222 **docking:** distribution of the number of models within 5Å ligand (L)-RMSD from the native,
 223 colored according to fragment quality. **(C) Improvement after FlexPepDock refinement:**
 224 distribution of the L-RMSDs of the top 1% FlexPepDock refined models (in orange) and
 225 corresponding PIPER models (in gray). Shown are the results of runs starting from the unbound
 226 receptor structure and including receptor minimizations (see also **Figure 3**). The circles
 227 represent the L-RMSD values of the best model among the top 10 ranking clusters. The Y-axis
 228 has been trimmed to 7Å. (Note that the gray circles are taken from the top-ranked models of the
 229 *full* PIPER run based on density clustering [see **Methods** and Porter *et al.*²⁴]), while the
 230 distributions represent the subset of models that served as starting structures for the models
 231 selected after FlexPepDock refinement).

232
233 We previously showed that extracting fragments based on sequence motif information
234 allows identification of bound peptide conformations that reflect the structural pattern of
235 these motifs²⁴. We demonstrate here that representative fragments are not restricted to
236 peptides with known motifs. In fact, a comparison to the fragments extracted based on
237 sequence motif (for the dataset analyzed in the PeptiDock study, using the motif
238 definition therein²⁴) shows that the fragment picker approach produces overall
239 ensembles that contain structures more similar to the bound peptide conformation (see
240 **Supplementary Table S2**).

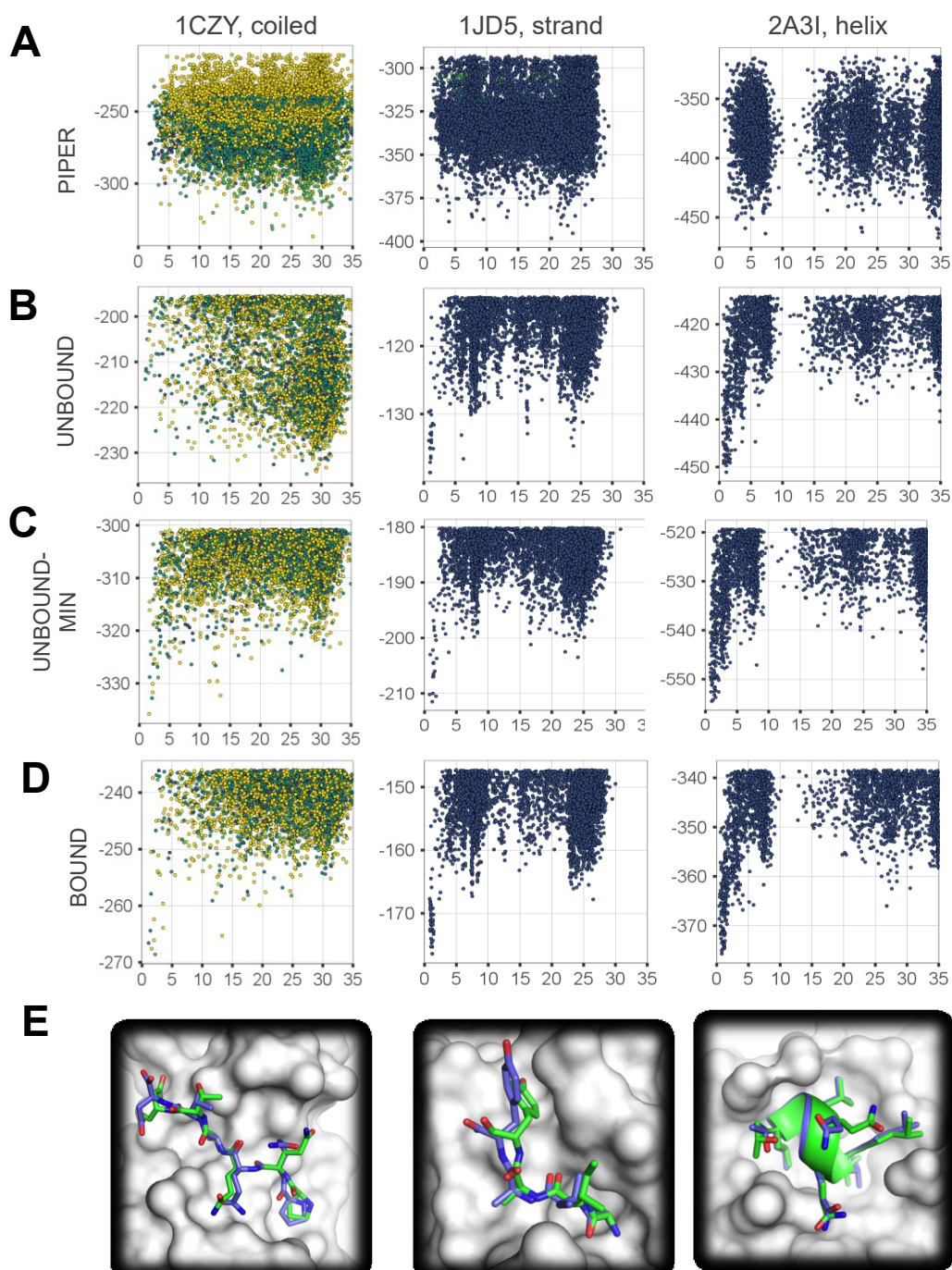
241
242 **Rigid-body docking: fragment quality and PIPER performance**
243 The fact that the fragment ensembles include bound-like conformations justifies
244 proceeding to the next step, namely their docking onto the receptor. The PIPER rigid-
245 body docking protocol allows fast and exhaustive sampling to provide coarse models of
246 fragment-receptor interactions, of which the top-scoring can be followed up by
247 subsequent refinement to allow for conformational changes upon binding. The effective
248 range for successful refinement using the FlexPepDock protocol was previously found
249 to be within up to 5Å in terms of Cartesian RMSD, and up to 50 degrees in terms of ϕ - ψ
250 RMSD (distance of fragment from the bound peptide conformation in ϕ - ψ dihedral
251 space)¹⁰. It is thus important for the PIPER docking stage to identify a large pool of
252 fragments that are densely docked in close proximity (within effective Cartesian RMSD
253 range) of the native peptide binding mode, involving docked fragments that are similar
254 to the native peptide bound conformation (within effective phi-psi RMSD range). Indeed,
255 analysis of the top ranking PIPER models shows presence of good quality fragments at
256 the binding site (in fact, most complexes include <1.0Å bb RMSD fragments; **Figure**
257 **2B**).

258
259 **Improvement of PIPER models by FlexPepDock refinement**
260 The FFT algorithmic implementation of rigid-body sampling in PIPER makes exhaustive
261 orientation search possible with significant computational efficiency, but is defined on a
262 grid. Consequently, the scoring function can successfully isolate the best few hundreds

263 from the vast pool of billions of positions of the peptide fragment relative to the receptor,
264 but not discriminate the top rigid-body docked models further (**Figure 3A &**
265 **Supplementary Figure S1**). In turn, the Rosetta scoring function used in the
266 FlexPepDock Refinement protocol (currently Talaris 2014 ⁴⁴) is highly accurate, but this
267 flexible docking protocol lacks the ability for fast and exhaustive sampling. Thus, to
268 address the problem of exhaustive sampling with high accuracy, we combine the fast
269 and exhaustive rigid-body sampling of PIPER with accurate flexible refinement by
270 FlexPepDock of the top ranking few hundred best models. Indeed, the FlexPepDock
271 refinement stage significantly improves the model quality, as well as better model
272 ranking (See **Figures 2C, 3C & Supplementary Figure S1**). This includes the
273 identification of a near-native funnel missed before (e.g. 1CZY in **Figure 3 – compare A**
274 **to C**), or significant enhancement of a near-native funnel (e.g. 1JD5 and 2A3I). More
275 examples can be found in **Supplementary Figure S1**.

276 We performed three runs to assess protocol performance (Summarized in **Figure 4A**
277 and **Supplementary Table S1B**; specific examples are shown in **Figure 3**): First, we
278 applied the protocol to bound receptor structures. For these runs a near-native peptide
279 conformation (L-RMSD $\leq 2.0\text{\AA}$, see **Methods** section) was found among the top 10
280 ranked clusters for 19 out of 27 complexes (success rate=70%, **Figure 3D**). We then
281 proceeded to the real-world scenario, in which the free receptor structure was provided
282 as starting point (unbound run), leading to worse performance, as expected (10
283 complexes successfully modeled - success rate=37%, **Figure 3B**). Importantly however,
284 when including also receptor flexibility during the refinement stage (unbound-min run),
285 these results improved, in particular if 10 best models are considered (14/27 complexes
286 successfully modeled - success rate=52%, **Figure 3C**).

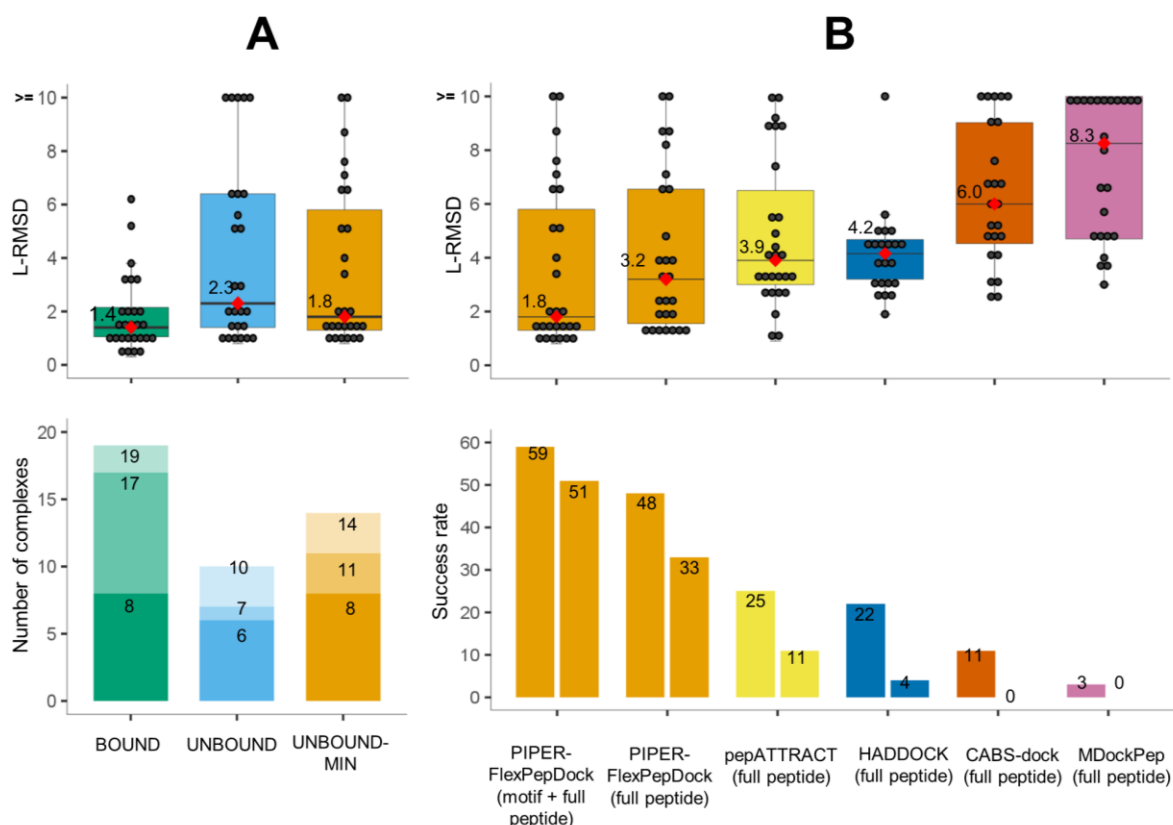
287



288
289
290
291
292
293
294
295
296

Figure 3. Examples of global peptide docking energy landscapes: Left: PDB id 1CZY (coiled peptide); Center: 1JD5 (extended peptide); Right: 2A3I (helical peptide). (A) Energy landscape as sampled in the first docking step of the protocol by PIPER rigid body docking of peptide fragments onto the *unbound* receptor structure. **(B-D)** Energy landscapes for the PIPER-FPD scheme, starting from the unbound receptor structure **(B)**, the unbound receptor structure including receptor flexibility **(C)**, and the corresponding bound receptor for comparison **(D)**. Models are colored according to fragment quality, as in previous Figures. **(E)** Comparison of the modeled to the native structure (shown in blue and green, respectively).

297



298 **Figure 4. PIPER-FlexPepDock peptide docking performance. (A) Overall performance on**
 299 **a non-redundant set of 27 peptide-protein complexes. Top:** Distribution of best model L-
 300 **RMSDs (among top 10 ranking clusters) for runs using the bound (BOUND) and free**
 301 **(UNBOUND & UNBOUND-MIN) receptor structures, the latter including also receptor flexibility**
 302 **in the final refinement step (only the motif region was modeled for the 12 complexes with known**
 303 **motif). The median values are shown as red diamonds and printed alongside. Bottom:**
 304 **Distribution of the ranks of the first near-native cluster (defined as L-RMSD ≤ 2.0 Å), shown**
 305 **using different shades (for corresponding results among the top1, top3 and top10 ranked**
 306 **predictions). (B) Comparison to performance by other algorithms. Top:** Box plots of best L-
 307 **RMSDs among top 10 ranking clusters, including results for the motif part where the motif is**
 308 **known (as in A), as well as for the full peptide, for comparison. Bottom:** Performance is shown
 309 **for different cutoffs (2.0 Å and 3.0 Å L-RMSD in left and right boxes, respectively) (See**
 310 **Supplementary Table S1B for more details).**

311

312

313 Comparison with other global docking protocols

314 We compared the results of PIPER-FlexPepDock (unbound-min run) with other existing
 315 global peptide-protein docking protocols such as HADDOCK¹², pepATTRACT¹³,
 316 CABS-dock¹⁵, and MDockPep¹⁶ on our non-redundant set of 27 complexes, as well as
 317 on the set of 42 complexes used by these protocols in previous studies (34 complexes

318 were compared with HADDOCK as other 8 cases were not included in their unbound
 319 run set). Since full length peptides were modeled using the other protocols, we modeled
 320 full length peptides for the motif set cases for valid comparison. The success rate for
 321 generating near-native models (*i.e.*, L-RMSD within 2.0Å, or 3.0Å) was significantly
 322 better for PIPER-FlexPepDock than any other protocol, even for models of the full
 323 peptides (see **Figure 4B** and **Table 2**).

324
 325 **Table 2.** Summary of performance of PIPER-FlexPepDock, and comparison to other peptide
 326 docking protocols. Results are shown for PIPER-FlexPepDock runs on unbound receptor
 327 structures, including receptor minimization.

Cutoff L ^a	PIPER-FlexPepDock		pepATTRACT ¹³	HADDOCK ¹²	CABSDOCK ¹⁵	MDockPep ¹⁶
	motif	full				
non-redundant set (n=27)						
1.5	12 (44%)	7 (26%)	2 (7%)	0	0	0
2	14 (52%)	9 (33%)	3 (11%)	1 (5%)	0	0
2.5	16 (59%)	12 (44%)	4 (15%)	2 (9%)	1 (4%)	0
3	16 (59%)	13 (48%)	7 (26%)	5 (23%)	3 (11%)	1 (4%)
calibration set (n=9)						
1.5	6 (67%)	3 (33%)	1 (11%)	0	0	0
2	8 (89%)	3 (33%)	2 (22%)	0	0	0
2.5	9 (100%)	5 (56%)	2 (22%)	1 (11%)	0	0
3	9 (100%)	5 (56%)	4 (44%)	1 (11%)	1 (11%)	0
redundant set (n=42) ^b						
1.5	20 (48%)	10 (24%)	4 (10%)	1 (3%)	1 (2%)	0
2	26 (62%)	17 (40%)	7 (17%)	3 (9%)	2 (5%)	1 (2%)
2.5	29 (69%)	22 (52%)	8 (19%)	4 (12%)	5 (12%)	2 (5%)
3	29 (69%)	25 (60%)	13 (31%)	11 (32%)	7 (17%)	4 (10%)

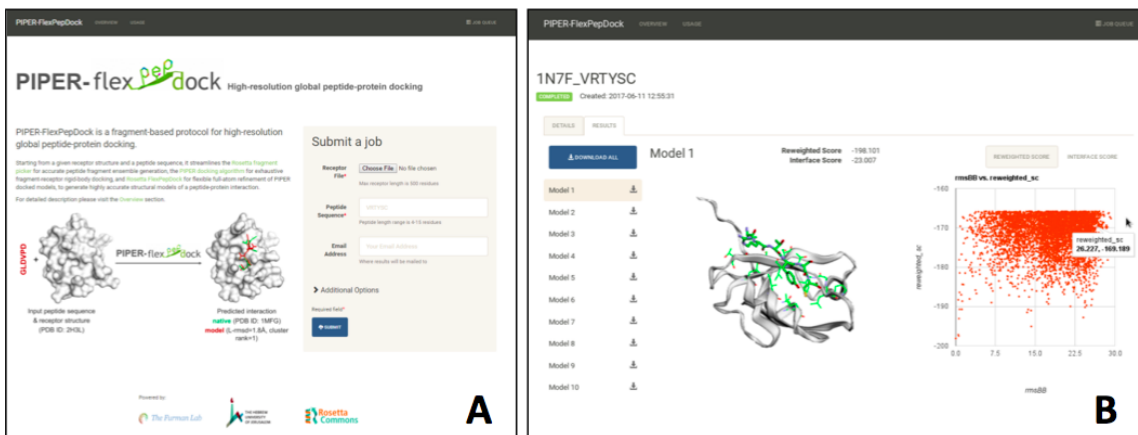
328 ^a number of models within the mentioned ligand RMSD value (Å)

329 ^b n=34 for HADDOCK results

330
 331 **The PIPER-FlexPepDock server for the high-resolution modeling of peptide-**
 332 **protein interactions**

333 In order to maximize the impact of our new protocol for global peptide-protein docking
 334 and to make it accessible to the modeling of many new peptide-protein complexes, we
 335 have set up a user-friendly server open to the scientific community (**Figure 5**). All that is
 336 needed is a structure of the receptor and a sequence of the peptide, but additional
 337 information about peptide secondary structure can also be included to narrow the

338 search. The top-ranking resulting models can be downloaded, or inspected by an
339 interactive viewer using the 3Dmol.js libraries⁴⁵.



340 **Figure 5. The PIPER-FlexPepDock server: (A)** Job submission page: the required input
341 includes the structure of the receptor and the sequence of the peptide; advanced options are
342 accessible via a button. The tabs at the top provide links to detailed descriptions of the server,
343 as well as to the Queue (upper right). **(B)** Results of an example peptide docking run: The liprin
344 C-terminal peptide sequence VRTYSC docked onto the PDZ domain of GRIP1 (free receptor
345 PDB id 1N7E). The top10 ranking models can be downloaded, and links to the individual models
346 are provided to the left for inspection using an interactive viewer. In this case, Model 1 is an
347 accurate prediction (L-RMSD=1.0Å from solved structure PDB id 1N7F). On the right side a
348 scatter plot shows the sampled energy landscape (relative to an arbitrary reference structure).
349
350

351 Discussion

352 **A new approach for global peptide docking with excellent performance** - With the
353 presentation of our new PIPER-FlexPepDock algorithm, we have demonstrated that
354 combining fast and exhaustive rigid-body docking (using the FFT-based PIPER docking
355 algorithm) of a representative peptide conformer ensemble (approximated by fragments
356 extracted from solved structures, based on local similarity of sequence and secondary
357 structure), with high-resolution refinement (using Rosetta FlexPepDock) is a widely-
358 applicable approach for the generation of models of peptide-receptor structures of
359 remarkable accuracy – significantly better than any other current protocol - starting from
360 the sequence of the peptide and the structure of the receptor. The performance on a
361 large benchmark of solved peptide-protein complex structures demonstrates both
362 accuracy and robustness of our modeling approach, and opens up the way of modeling
363 many more peptide-protein interactions at much higher resolution and accuracy than
364 any existing global peptide-protein global docking protocol.

365 **Receptor-bound peptide conformations are adequately represented by fragments**
366 **extracted from protein monomer structures** – This study demonstrates that
367 fragments derived from solved protein structures, based on secondary structure and
368 sequence similarity (rather than on sequence binding motifs which are not always
369 available) represent the peptide conformational states with high accuracy, in particular
370 the bound state. Interestingly, it is this same observation regarding the representation of
371 local conformational preference that provided originally the platform for the
372 breakthrough of Rosetta *ab initio* protein structure prediction⁴⁶. This indicates that while
373 isolated peptides in solution rarely show significant conformational preferences⁴⁷, in the
374 encounter complex regime in vicinity of other proteins, their conformational freedom
375 seems to be restricted significantly (similar to local peptide regions within a full protein)
376 and can be represented by fragment libraries, in concordance with previous reports that
377 show similar arrangements of fragments within monomers and peptide-protein
378 interactions⁴⁸.

379
380 **Effective sampling of the energy landscape:** The simplified scoring function and
381 exhaustive sampling with PIPER allows uniform sampling of the fragments onto the
382 receptor on a smoothed energy landscape. The top scoring PIPER models represent
383 the dense sampling into wider energy basins. Though the ranking of models might lack
384 the accuracy at this stage, the following refinement stage performs local sampling to
385 efficiently locate the minimum. Interestingly, this approach is much more effective than
386 the local refinement starting from one representative model (only one FlexPepDock
387 optimization run is necessary starting from each PIPER model, compared to several
388 hundred to thousand runs starting from a representative (defined, e.g. from a PIPER run
389 as implemented in the PeptiDock peptide motif docking algorithm²⁴). This is most
390 probably due to the fact that these starting coarse models are trapped in many distinct
391 states, each near a distinct local minimum, simplifying sampling during optimization.

392
393 **Mapping encounter complexes and more:** The peptide-receptor binding energy
394 landscape can provide a broader understanding of the binding mechanism itself. The
395 exhaustive sampling with accurate refinement provides a high-resolution map of the

396 energy landscape and helps us understand the energetic of the encounter between the
397 peptide and the receptor. In a previous study, we were able to demonstrate that
398 experimentally observed encounter complexes are well reproduced from a global
399 protein docking energy landscape ⁴⁹, and we anticipate that the corresponding peptide-
400 protein docking energy landscape will provide similar information.

401 Our novel global peptide docking pipeline allows high-resolution modeling of peptide-
402 protein interaction with much higher accuracy than ever before. The robustness of our
403 approach can be seen in the validation of known complexes and further improvement
404 such as clever incorporation of increased receptor flexibility and development of
405 accurate peptide conformational ensemble generation will further enhance the range of
406 applications.

407

408

409 **Materials and Methods**

410 **Data set (Table 1 and Supplementary Table S1)**

411 Docking performance and analysis was calibrated and assessed on a benchmark of
412 peptide-protein complexes derived from the PeptiDB database ⁵⁰, filtered according to
413 the following criteria:

414 (1) *Availability of both the complex and the free receptor structure*, solved by X-ray
415 crystallography (resolution of the complex $\leq 2.0\text{\AA}$).

416 (2) *Absence of crystal contacts that could influence the peptide conformation*. In certain
417 cases this further interaction is of biological relevance, leading to receptor
418 multimerization and clustering (e.g. PeptiDB entries involving some of the SH3 domain-
419 peptide interactions, 2AK5 ⁵¹, and 2J6F ⁵²). Since for these cases, obtaining high-
420 resolution models might be challenging without including the symmetry mate, such
421 examples were removed from the dataset.

422 (3) *Absence of large receptor rearrangement upon peptide binding*. Even though the
423 present implementation of PIPER-FlexPepDock does allow for local conformational
424 changes in the receptor (backbone as well as side chains), accurate modeling of more
425 significant movement of the receptor upon peptide binding (e.g. significant loop
426 movement at the binding interface in PeptiDB entry 1D4T ⁵³) require the development of

427 algorithms for efficient modeling of more significant receptor flexibility, which is beyond
428 the scope of the present study.

429 (4) *Non-redundant dataset*. The criteria above result in a dataset of 42 complexes
430 (**Supplementary Table S1C**) that is very similar to the one used in previous studies by
431 different groups^{12,13,15,16}. To ensure that no bias towards a certain peptide-receptor
432 would be introduced, we extracted a domain non-redundant set (defined by CATH
433 classification³⁶), resulting in the 27 complexes described in this study in detail (**Table 1**
434 and **Supplementary Table S1B**).

435 The dataset was further divided into two subsets, based on available information about
436 a peptide binding motif (defined in this study based on ELM²⁹, <http://elm.eu.org>): For
437 the **motif set** (12 complexes) we modeled only the motif part, since it contributes most
438 to binding, and shorter peptides are easier to model. To enable comparison to
439 performance of other protocols, we subsequently also docked the full peptide. For the
440 **non-motif set** (15 complexes), the full peptide was docked.

441 **Initial calibration set:** For initial calibration, we selected a smaller subset of 9
442 complexes (**Supplementary Table S1A**). The established protocol was then validated
443 on the remaining complexes, to ensure similar performance and thereby prevent
444 overfitting of the modeling protocol.

445

446 **The steps of the PIPER-FlexPepDock protocol**

447 In the following we provide specific details of the different steps of the PIPER-
448 FlexPepDock protocol. For runline commands, see the **Supplementary Materials**
449 section.

450 **(1) Generation of peptide conformations using Rosetta fragment picker and** 451 **Rosetta fixbb design**

452 The Rosetta fragment picker²² uses a scoring measure composed of a weighted
453 combination of secondary structure propensity, sequence profile similarity and residue
454 propensities for local regions in the Ramachandran plot⁵⁴ to map fragments to *vall*, a
455 database of solved high-resolution protein structures. Consequently, the mapped
456 fragments are consistent with the peptide sequence (as defined by a sequence
457 similarity profile generated with PSI-BLAST⁵⁵) and secondary structure (as predicted

458 using PSIPRED⁵⁶; even though PSIPRED was shown to perform quite well for shorter
459 sequences¹¹, we use the full protein sequence from which the peptide was derived for
460 PSIPRED and PSIBLAST runs, where available). If the preferred secondary structure is
461 already known (e.g. the alpha helical nuclear receptor box motif) it can be provided
462 instead of PSIPRED predictions. Secondary structural information can also be obtained
463 from experimental techniques such as Circular dichroism (CD) spectroscopy, or
464 approximated by residue Ramachandran local region propensities (derived from
465 statistical analysis of high-resolution protein structure⁵⁷). The coordinates of the top fifty
466 assigned fragments are extracted from the PDB, and side chains for non-identical
467 residues are modeled using the Rosetta fixbb design algorithm⁵⁸. The whole process
468 results in an ensemble of 50 fragments for the query peptide sequence.

469

470 **(2) Rigid body docking using PIPER**

471 Each of the fifty fragments is globally docked onto the receptor using the PIPER Fourier
472 transform (FFT) docking algorithm, as detailed before²⁴, decomposing the free receptor
473 into independent binding units (either a single domain or repeated, non-decomposable
474 domains; as in Lavi *et al.*⁹). The calculations are performed for each of 70,000
475 rotations, and one lowest-energy translation for each rotation is retained. For each
476 fragment docking run the top ranked 250 solutions (total 50x250 = 12500 models) are
477 collected for refinement in the next step (see **Supplementary Figure S3** for a
478 comparison of performance using different numbers of top-ranked solutions).

479

480 **Selection of final model from a PIPER simulation:** In order to compare performance
481 of a protocol involving only the first PIPER rigid body docking step (in **Table 1**), we
482 selected the final models as reported previously (similar to the PeptiDock
483 implementation²⁴, but without minimization). In short, the models collected are clustered
484 (with radius of 3.5Å C_α RMSD), and cluster density is used for ranking and selection of
485 representatives.

486

487 **(3) The Rosetta FlexPepDock refinement algorithm**

488 The FlexPepDock Refinement protocol refines all of the peptide's degrees of freedom
489 (*i.e.* its rigid body orientation as well as backbone dihedral angles), as well as the
490 receptor side chain conformations. Rosetta FlexPepDock refinement was performed as
491 described previously ¹⁰, with slight changes: (1) *Sampling*: In our present
492 implementation, we also allowed the receptor backbone to move during minimization
493 steps, to allow for slight readjustment upon binding (compare e.g. **Figures 3B** and **3C**).
494 (2) *Scoring*: Rosetta energy function Talaris2014 ⁴⁴ was used. Clustering of models was
495 performed as previously described, using a threshold of 2.0Å ³². The top-scoring
496 member of each cluster (according to reweighted score) was selected as the
497 representative member, and clusters were ranked based on the reweighted score of the
498 representative members (as in Raveh *et al.* ¹¹).

499

500 **Model Evaluation Criteria**

501 For each global docking run the 10 top ranking clusters were selected as prediction and
502 evaluated for quality based on ligand RMSD (L-RMSD), calculated between the native
503 and model peptide backbone atoms after optimal superimposition of the receptor, as
504 done in the CAPRI assessment ^{34,35}. L-RMSD and other measures, such as Fnat and I-
505 RMSD, were calculated using DockQ ⁵⁹.

506

507 **Rosetta release version**

508 The protocol and tests described in this manuscript follow the FlexPepDock protocol, as
509 implemented within the Rosetta weekly release version 2016.20.58704.

510

511 **Simulation running time**

512 The processing time for the different stages of the protocol depends on both the length
513 of receptor and the peptide sequence. For example the global docking the carboxy-
514 terminal tail of the ErbB2 Receptor GLDVPV onto the free ERBIN PDZ domain (103
515 residues) the generation of 50 fragments takes ~8 CPU minutes over an AMD Sun
516 cluster with 300 cores. For the same complex a single PIPER fragment docking
517 simulation takes ~2 minutes and a single refinement run of the PIPER docked model
518 takes ~1 minutes on the same system architecture (~ 1.5 hours to refine all models).

519

520 **Protocol availability**

521 The runline commands are provided in the **Supplementary Material Section**. The
522 Rosetta software is available for free to the academic community. The details regarding
523 downloading and installation is available at <https://www.rosettacommons.org>. PIPER
524 FFT rigid body docking is available as part of the protein-protein docking server ClusPro
525 (PeptiDock at <https://peptidock.cluspro.org>).

526

527 **Acknowledgements**

528 We thank Dr. Barak Raveh for insightful discussions. We also thank Dr. Christina
529 Schindler for providing the pepATTRACT models as reported in Schindler et al. ¹³, and
530 Dr. Mikael Trellet and Prof. Alexandre Bonvin for providing the link for the SBGrid
531 deposited HADDOCK models (<https://data.sbgrid.org/dataset/131/>) as reported in Trellet
532 et al. ¹², for the comparison of performance.

533

534 **Author contributions**

535 NA conceived the study, designed the project, developed and implemented the protocol,
536 generated and analyzed data, wrote the manuscript; OG, NA, and BX developed the
537 server; KP developed the PIPER peptide docking protocol; DK conceived the study,
538 analyzed the data; OSF conceived the study, designed and coordinated the project,
539 analyzed the data, wrote the manuscript.

540

541 **Competing interests**

542 No competing interests exist.

543

544 **References**

- 545 1 Pawson & Nash. Assembly of cell regulatory systems through protein interaction
546 domains. *Science* **300**, 445-452, (2003).
547 2 Petsalaki & Russell. Peptide-mediated interactions in biological systems: new
548 discoveries and applications. *Curr Opin Biotechnol* **19**, 344-350, (2008).

- 549 3 Neduva, Linding, Su-Angrand, Stark, de Masi, Gibson, Lewis, Serrano & Russell.
550 Systematic discovery of new recognition peptides mediating protein interaction networks.
551 *PLoS Biol* **3**, e405, (2005).
- 552 4 Vacic, Oldfield, Mohan, Radivojac, Cortese, Uversky & Dunker. Characterization of
553 molecular recognition features, MoRFs, and their binding partners. *J Proteome Res* **6**,
554 2351-2366, (2007).
- 555 5 Gamble, Vajdos, Yoo, Worthylake, Houseweart, Sundquist & Hill. Crystal structure of
556 human cyclophilin A bound to the amino-terminal domain of HIV-1 capsid. *Cell* **87**, 1285-
557 1294, (1996).
- 558 6 London, Raveh & Schueler-Furman. Druggable protein-protein interactions--from hot
559 spots to hot segments. *Curr Opin Chem Biol* **17**, 952-959, (2013).
- 560 7 Trabuco, Lise, Petsalaki & Russell. PepSite: prediction of peptide-binding sites from
561 protein surfaces. *Nucleic Acids Res* **40**, W423-427, (2012).
- 562 8 Saladin, Rey, Thevenet, Zacharias, Moroy & Tuffery. PEP-SiteFinder: a tool for the blind
563 identification of peptide binding sites on protein surfaces. *Nucleic Acids Res* **42**, W221-
564 226, (2014).
- 565 9 Lavi, Ngan, Movshovitz-Attias, Bohnuud, Yueh, Beglov, Schueler-Furman & Kozakov.
566 Detection of peptide-binding sites on protein surfaces: the first step toward the modeling
567 and targeting of peptide-mediated interactions. *Proteins* **81**, 2096-2105, (2013).
- 568 10 Raveh, London & Schueler-Furman. Sub-angstrom modeling of complexes between
569 flexible peptides and globular proteins. *Proteins* **78**, 2029-2040, (2010).
- 570 11 Raveh, London, Zimmerman & Schueler-Furman. Rosetta FlexPepDock ab-initio:
571 simultaneous folding, docking and refinement of peptides onto their receptors. *PLoS*
572 *One* **6**, e18934, (2011).
- 573 12 Trellet, Melquiond & Bonvin. A unified conformational selection and induced fit approach
574 to protein-peptide docking. *PLoS One* **8**, e58769, (2013).
- 575 13 Schindler, de Vries & Zacharias. Fully Blind Peptide-Protein Docking with pepATTRACT.
576 *Structure* **23**, 1507-1515, (2015).
- 577 14 Ben-Shimon & Niv. AnchorDock: Blind and Flexible Anchor-Driven Peptide Docking.
578 *Structure* **23**, 929-940, (2015).
- 579 15 Kurcinski, Jamroz, Blaszczyk, Kolinski & Kmiecik. CABS-dock web server for the flexible
580 docking of peptides to proteins without prior knowledge of the binding site. *Nucleic Acids*
581 *Res* **43**, W419-424, (2015).
- 582 16 Yan, Xu & Zou. Fully Blind Docking at the Atomic Level for Protein-Peptide Complex
583 Structure Prediction. *Structure* **24**, 1842-1853, (2016).
- 584 17 Peterson, Roy, Christoffer, Terashi & Kihara. Modeling disordered protein interactions
585 from biophysical principles. *PLoS Comput Biol* **13**, e1005485, (2017).
- 586 18 Dominguez, Boelens & Bonvin. HADDOCK: a protein-protein docking approach based
587 on biochemical or biophysical information. *J Am Chem Soc* **125**, 1731-1737, (2003).
- 588 19 de Vries, Rey, Schindler, Zacharias & Tuffery. The pepATTRACT web server for blind,
589 large-scale peptide-protein docking. *Nucleic Acids Res*, (2017).
- 590 20 Webb & Sali. Comparative Protein Structure Modeling Using MODELLER. *Curr Protoc*
591 *Bioinformatics* **47**, 5 6 1-32, (2014).
- 592 21 Trott & Olson. AutoDock Vina: improving the speed and accuracy of docking with a new
593 scoring function, efficient optimization, and multithreading. *J Comput Chem* **31**, 455-461,
594 (2010).
- 595 22 Gront, Kulp, Vernon, Strauss & Baker. Generalized fragment picking in Rosetta: design,
596 protocols and applications. *PLoS One* **6**, e23294, (2011).
- 597 23 Venkatraman, Yang, Sael & Kihara. Protein-protein docking using region-based 3D
598 Zernike descriptors. *BMC Bioinformatics* **10**, 407, (2009).

- 599 24 Porter, Bing, Beglov, Bohnuud, Alam, Schueler-Furman & Kozakov. ClusPro PeptiDock:
600 Efficient global docking of peptide recognition motifs using FFT. *Bioinformatics*
601 **10.1093/bioinformatics/btx216.**, (2017).
- 602 25 Berman, Westbrook, Feng, Gilliland, Bhat, Weissig, Shindyalov & Bourne. The Protein
603 Data Bank. *Nucleic Acids Res* **28**, 235-242, (2000).
- 604 26 Kozakov, Beglov, Bohnuud, Mottarella, Xia, Hall & Vajda. How good is automated
605 protein docking? *Proteins* **81**, 2159-2166, (2013).
- 606 27 Brooks, Brooks, Mackerell, Nilsson, Petrella, Roux, Won, Archontis, Bartels, Boresch,
607 Caflisch, Caves, Cui, Dinner, Feig, Fischer, Gao, Hodoscek, Im, Kuczera, Lazaridis, Ma,
608 Ovchinnikov, Paci, Pastor, Post, Pu, Schaefer, Tidor, Venable, Woodcock, Wu, Yang,
609 York & Karplus. CHARMM: the biomolecular simulation program. *J Comput Chem* **30**,
610 1545-1614, (2009).
- 611 28 Kozakov, Brenke, Comeau & Vajda. PIPER: an FFT-based protein docking program with
612 pairwise potentials. *Proteins* **65**, 392-406, (2006).
- 613 29 Dinkel, Van Roey, Michael, Kumar, Uyar, Altenberg, Milchevskaya, Schneider, Kuhn,
614 Behrendt, Dahl, Damerell, Diebel, Kalman, Klein, Knudsen, Mader, Merrill, Staudt, Thiel,
615 Welti, Davey, Diella & Gibson. ELM 2016-data update and new functionality of the
616 eukaryotic linear motif resource. *Nucleic Acids Res* **44**, D294-300, (2016).
- 617 30 Puntervoll, Linding, Gemund, Chabanis-Davidson, Mattingsdal, Cameron, Martin,
618 Ausiello, Brannetti, Costantini, Ferre, Maselli, Via, Cesareni, Diella, Superti-Furga,
619 Wyrwicz, Ramu, McGuigan, Gudavalli, Letunic, Bork, Rychlewski, Kuster, Helmer-
620 Citterich, Hunter, Aasland & Gibson. ELM server: A new resource for investigating short
621 functional sites in modular eukaryotic proteins. *Nucleic Acids Res* **31**, 3625-3630,
622 (2003).
- 623 31 Messih, Lepore & Tramontano. Loopng: a template-based tool for predicting the
624 structure of protein loops. *Bioinformatics* **31**, 3767-3772, (2015).
- 625 32 Gray, Moughon, Wang, Schueler-Furman, Kuhlman, Rohl & Baker. Protein-protein
626 docking with simultaneous optimization of rigid-body displacement and side-chain
627 conformations. *J Mol Biol* **331**, 281-299, (2003).
- 628 33 Lensink, Velankar & Wodak. Modeling protein-protein and protein-peptide complexes:
629 CAPRI 6th edition. *Proteins* **85**, 359-377, (2017).
- 630 34 Mendez, Leplae, De Maria & Wodak. Assessment of blind predictions of protein-protein
631 interactions: current status of docking methods. *Proteins* **52**, 51-67, (2003).
- 632 35 Mendez, Leplae, Lensink & Wodak. Assessment of CAPRI predictions in rounds 3-5
633 shows progress in docking procedures. *Proteins* **60**, 150-169, (2005).
- 634 36 Pearl, Bennett, Bray, Harrison, Martin, Shepherd, Sillitoe, Thornton & Orengo. The
635 CATH database: an extended protein family resource for structural and functional
636 genomics. *Nucleic Acids Res* **31**, 452-455, (2003).
- 637 37 Rohl, Strauss, Misura & Baker. Protein structure prediction using Rosetta. *Methods*
638 *Enzymol* **383**, 66-93, (2004).
- 639 38 Park, Lee, Heo & Seok. Protein loop modeling using a new hybrid energy function and
640 its application to modeling in inaccurate structural environments. *PLoS One* **9**, e113811,
641 (2014).
- 642 39 Vanhee, Verschueren, Baeten, Stricher, Serrano, Rousseau & Schymkowitz. BriX: a
643 database of protein building blocks for structural analysis, modeling and design. *Nucleic*
644 *Acids Res* **39**, D435-442, (2011).
- 645 40 Li, Suino, Daugherty & Xu. Structural and biochemical mechanisms for the specificity of
646 hormone binding and coactivator assembly by mineralocorticoid receptor. *Mol Cell* **19**,
647 367-380, (2005).

- 648 41 Guhaniyogi, Robinson & Stock. Crystal structures of beryllium fluoride-free and beryllium
649 fluoride-bound CheY in complex with the conserved C-terminal peptide of CheZ reveal
650 dual binding modes specific to CheY conformation. *J Mol Biol* **359**, 624-645, (2006).
- 651 42 Todd, Moore, Deivanayagam, Lin, Chattopadhyay, Maki, Wang & Narayana. A structural
652 model for the inhibition of calpain by calpastatin: crystal structures of the native domain
653 VI of calpain and its complexes with calpastatin peptide and a small molecule inhibitor. *J*
654 *Mol Biol* **328**, 131-146, (2003).
- 655 43 Remenyi, Good, Bhattacharyya & Lim. The role of docking interactions in mediating
656 signaling input, output, and discrimination in the yeast MAPK network. *Mol Cell* **20**, 951-
657 962, (2005).
- 658 44 Leaver-Fay, O'Meara, Tyka, Jacak, Song, Kellogg, Thompson, Davis, Pache, Lyskov,
659 Gray, Kortemme, Richardson, Havranek, Snoeyink, Baker & Kuhlman. Scientific
660 benchmarks for guiding macromolecular energy function improvement. *Methods*
661 *Enzymol* **523**, 109-143, (2013).
- 662 45 Rego & Koes. 3Dmol.js: molecular visualization with WebGL. *Bioinformatics* **31**, 1322-
663 1324, (2015).
- 664 46 Simons, Kooperberg, Huang & Baker. Assembly of protein tertiary structures from
665 fragments with similar local sequences using simulated annealing and Bayesian scoring
666 functions. *J Mol Biol* **268**, 209-225, (1997).
- 667 47 Ho & Dill. Folding very short peptides using molecular dynamics. *PLoS Comput Biol* **2**,
668 e27, (2006).
- 669 48 Vanhee, Stricher, Baeten, Verschueren, Lenaerts, Serrano, Rousseau & Schymkowitz.
670 Protein-peptide interactions adopt the same structural motifs as monomeric protein folds.
671 *Structure* **17**, 1128-1136, (2009).
- 672 49 Kozakov, Li, Hall, Beglov, Zheng, Vakili, Schueler-Furman, Paschalidis, Clore & Vajda.
673 Encounter complexes and dimensionality reduction in protein-protein association. *Elife* **3**,
674 e01370, (2014).
- 675 50 London, Movshovitz-Attias & Schueler-Furman. The structural basis of peptide-protein
676 binding strategies. *Structure* **18**, 188-199, (2010).
- 677 51 Jozic, Cardenes, Deribe, Moncalian, Hoeller, Groemping, Dikic, Rittinger & Bravo. Cbl
678 promotes clustering of endocytic adaptor proteins. *Nat Struct Mol Biol* **12**, 972-979,
679 (2005).
- 680 52 Moncalian, Cardenes, Deribe, Spinola-Amilibia, Dikic & Bravo. Atypical polyproline
681 recognition by the CMS N-terminal Src homology 3 domain. *J Biol Chem* **281**, 38845-
682 38853, (2006).
- 683 53 Poy, Yaffe, Sayos, Saxena, Morra, Sumegi, Cantley, Terhorst & Eck. Crystal structures
684 of the XLP protein SAP reveal a class of SH2 domains with extended, phosphotyrosine-
685 independent sequence recognition. *Mol Cell* **4**, 555-561, (1999).
- 686 54 Ramachandran, Ramakrishnan & Sasisekharan. Stereochemistry of polypeptide chain
687 configurations. *J Mol Biol* **7**, 95-99, (1963).
- 688 55 Altschul, Madden, Schaffer, Zhang, Zhang, Miller & Lipman. Gapped BLAST and PSI-
689 BLAST: a new generation of protein database search programs. *Nucleic Acids Res* **25**,
690 3389-3402, (1997).
- 691 56 Ward, McGuffin, Buxton & Jones. Secondary structure prediction with support vector
692 machines. *Bioinformatics* **19**, 1650-1655, (2003).
- 693 57 Berman, Battistuz, Bhat, Bluhm, Bourne, Burkhardt, Feng, Gilliland, Iype, Jain, Fagan,
694 Marvin, Padilla, Ravichandran, Schneider, Thanki, Weissig, Westbrook & Zardecki. The
695 Protein Data Bank. *Acta Crystallogr D Biol Crystallogr* **58**, 899-907, (2002).
- 696 58 Kuhlman, Dantas, Ireton, Varani, Stoddard & Baker. Design of a novel globular protein
697 fold with atomic-level accuracy. *Science* **302**, 1364-1368, (2003).

698 59 Basu & Wallner. DockQ: A Quality Measure for Protein-Protein Docking Models. *PLoS*
699 *One* **11**, e0161879, (2016).
700
701

702 **Supplementary Information**

703

704 **Supplementary Material**

705 **Runline commands**

706 The runline commands for the different stages are given below:

707

708 **1. Fragment generation using fragment picker:**

709

710 The `make_fragments.pl` script is used to run PSI-BLAST and PSIPRED to generate the
711 peptide secondary structure and the sequence similarity profile:

712

```
713 $make_fragments.pl -verbose peptide.fasta
```

714

715 Rosetta fragment picker is used to assign fragments consistent with the predicted secondary
716 structure and sequence profile to the `vall` database of high-resolution protein fragments:

717

```
718 $fragment_picker.linuxgccrelease -database rosetta_database -  
719 in:file:vall vall.jul19.2011 -in:file:checkpoint pep_seq.checkpoint -  
720 frags:frag_sizes 6 -frags:n_candidates2000 -frags:n_fragments 50 -  
721 frags:ss_pred pep_seq.psipred_ss2 psipred -frags:scoring:config  
722 psi_L1.cfg -frags:bounded_protocol true
```

723

724 These assigned fragments are extracted from the Protein Data Bank (including the side-chains).
725 The non-identical residues are mutated using the Rosetta `fixbb` design protocol:

726

```
727 $fixbb.linuxgccrelease -database rosetta_database -in:file:s  
728 fragment_1.pdb -resfile mutation_resfile -ex1 -ex2 -use_input_sc
```

729

730

731 **2. PIPER Docking:**

732

733 *Step 1:* preprocessing the input receptor and fragments using `pdbprep.pl` and
734 `pdbnmd.pl`:

735

```
736 $perl pdbprep.pl receptor.pdb  
737 $perl pdbnmd.pl receptor.pdb '?'
```

738

739 Each of the 50 fragments is similarly processed.

740

741 **Step II: Running PIPER FFT docking:**

742

```
743 $piper.acpharis.omp.20120803 -vv -c1.0 -k4 --msur_k=1.0 --maskr=1.0 -T  
744 FFTW_EXHAUSTIVE -R 70000 -t 1 -p atoms.0.0.4.prm.ms.3cap+0.5ace.Hr0rec  
745 -f coeffs.0.0.4.motif -r rot70k.0.0.4.prm receptor_nmin.pdb  
746 fragment1_nmin.pdb >piper.log
```

747

748 Each of the fragments is docked onto the receptor.

749

750 **Step III: Top scoring 250 PIPER models are extracted:**

751

```
752 $python apply_ftresult.py -i PIPER_model_ID ft.000.00 rot70k.0.0.4.prm  
753 fragment1_nmin.pdb --out-prefix PIPER_model_ID
```

754

755 Where PIPER_model_ID is an integer value assigned to each transformation for a fragment.

756

757

758 **3. FlexPepDock Refinement:**

759

760 *Step I: prepacking the PIPER model*

761 In the PIPER docked model the receptor is replaced with a prepacked receptor. A single
762 prepacked receptor is used.

763

764 *Step II: Running refinement*

765

```
766 $FlexPepDocking.mpi.linuxgccrelease -database rosetta_database -  
767 in:file:s PIPER_model_1.pdb -scorefile score.sc -min_receptor_bb -  
768 lowres_preoptimize -pep_refine -flexpep_score_only -ex1 -ex2aro -  
769 use_input_sc -unboundrot free_receptor.pdb
```

770

771 Where PIPER_model_1.pdb is the prepacked model.

772

773 **4. Clustering:**

774

775 The top scoring 1% refined models (125) are clustered using the Rosetta cluster application:

776

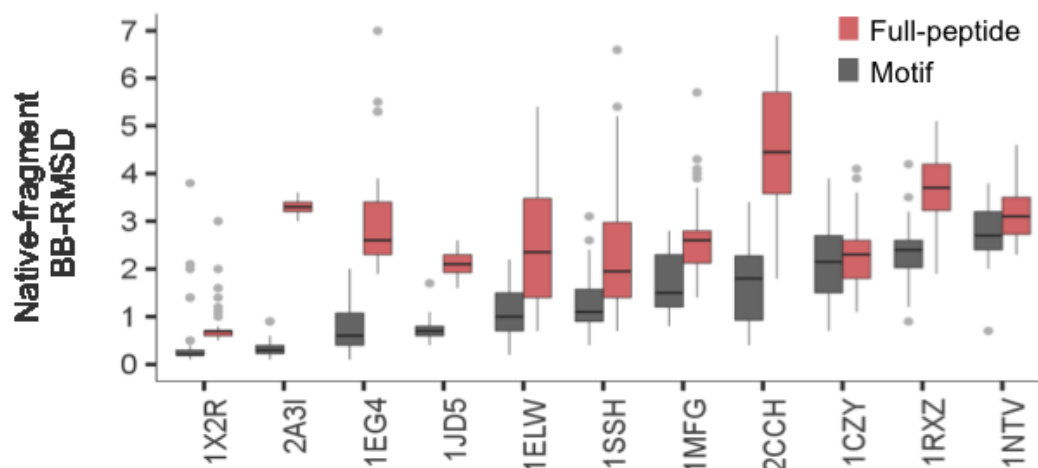
```
777 $cluster.linuxgccrelease -in:file:silent decoys.silent top_model_list -  
778 in:file:silent_struct_type binary -database $PATH_TO_DB -cluster:radius  
779 2.0 -in:file:fullatom -tags `cat top_refined_list` -  
780 silent_read_through_errors
```

781

782 The clusters are ranked based on the top scoring decoys from each clusters, based on
783 reweighted score, and top ranking 10 clusters are selected as putative models.

784 **Supplementary Figures**

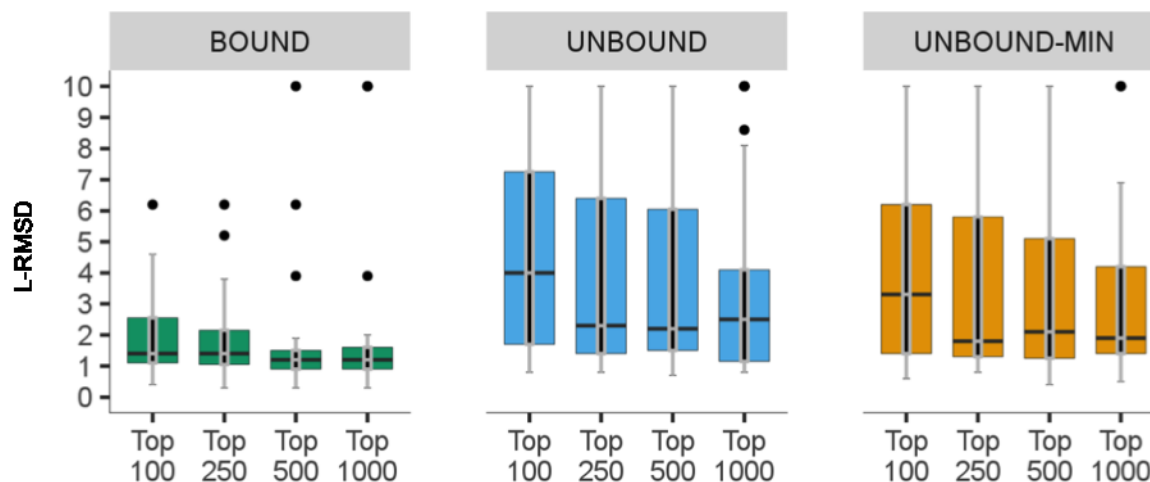
785 **Figure S1.** Global peptide docking energy landscapes for the full dataset (accompanies
786 **Figure 3**; provided as separate file). X-axis: L-RMSD, Y-axis: reweighted score. **Top**
787 **line:** PIPER rigid body docking of peptide fragments onto the unbound receptor structure;
788 **Middle lines:** FlexPepDock refinement of the PIPER docked fragments on the unbound
789 rigid (**second line**) and flexible (**third line**) receptor structure; **Bottom line:** PIPER-
790 FlexPepDock results starting from a bound receptor structure.



791 **Figure S2.** Fragment quality is significantly better for shorter, motif-defined peptide
792 segments (accompanies Figure 2A): Distributions of fragments backbone RMSD values
793 relative to the bound peptide conformations for the motif segments and corresponding
794 full length peptides. The motif set complexes 1JWG and 1TP5 are not added as in these
795 cases the motif covers the whole peptide.

797

798



799 **Figure S3.** The performance of the PIPER-FPD with different number of top PIPER
800 models selected for the refinement stage: Distributions of L-RMSDs of the best models
801 among top 10 ranking clusters for runs using the bound receptor structure (BOUND)
802 and the free receptor structure (UNBOUND & UNBOUND-MIN), the latter including also
803 receptor flexibility in the final refinement step (only the motif region was modeled for the
804 12 complexes with known motif). The number of PIPER models taken for the
805 FlexPepDock refinement step is shown below each boxplot. Based on these results, we
806 determined a cutoff of 250 models for optimal tradeoff between performance and
807 running time.

808

809

810 Supplementary Tables

811
 812 **Table S1.** Details of the datasets of peptide-protein complexes, including modeling
 813 results for PIPER-FlexPepDock and other peptide docking protocols (accompanies
 814 **Table 1**; provided as separate xls file). **(A)** calibration set (n=9 complexes); **(B)** non-
 815 redundant set (n=27 complexes); **(C)** redundant set (n=42 complexes).

816
 817 **Table S2.** Median fragment-native Backbone-RMSD values for the PeptiDock set
 818 complexes obtained using Rosetta fragment picker and the motif-based fragment
 819 generation approach (used in PeptiDock²⁴)

PDB ID		Peptide sequence (defined by PeptiDock)	Fragment quality ^a	
Complex	Free Receptor		Rosetta fragment picker	PeptiDock sequence – motif based fragment picker
1D4T:A	1D1ZA	TIYAQV	2.4	1.9^b
1SSH:A	1OOTA	PAMPAR	2.0	2.1
1MFG:A	2H3LA	LDVPV	1.4	1.4
2H9M:A	2H14A	ARTKQ	2.1	2.4
2FOJ:A	2FAWA	RAHSS	1.6	2.0
2HPL:A	2HPJA	DDLYG	1.7	2.2
1CZY:A	1CA4A	PQQATDD	2.4	2.5
1JD5:A	1JDA	AIAYF	1.4	1.5
2VJ0:A	1B9KA_1	WVTFE	1.1	2.1
2VJ0:A	1B9KA_2	FEDNF	2.0	2.5
2C3I:B	2J2IB	RRRHPS	2.5	2.4
2CCH:B	1H1RB	KGRRL	1.5	1.9
1EG4:A	1EG3A	RSPPPY	1.6	1.8
1RXZ:A	1RWZA	QATLERWF	2.9	2.8
1ER8:E	4PAEA	HLLVY	1.8	2.0
1JWG:A	1JWFA	DEDLL	2.1	1.6

820 ^a Similarity between fragments and bound peptide conformation: Median backbone RMSD (Å)

821 ^b In bold: significantly better fragments ($\Delta \geq 0.3$ Å)

822