1 **Large-scale docking predicts that sORF-encoded peptides may function through protein-**

2 **peptide interactions in *Arabidopsis thaliana***

3

4 Rashmi R. Hazarika[1], Nikolina Sostaric[1], Yifeng Sun[1,2], Vera van Noort[1,3]

5

6 [1]KU Leuven, Department of Microbial and Molecular Systems, KU Leuven Kasteelpark Arenberg 22,

7 B-3001 Leuven, Belgium.

8 [2]KU Leuven, Faculty of Engineering Technology, Campus Group T, Andreas Vesaliusstraat 13, 3000

9 Leuven, Belgium.

10 [3]Leiden University, Institute of Biology Leiden, 2300 RA Leiden, The Netherlands.

11

12

13

14

15

16

17

18    **Abstract**

19    Several recent studies indicate that small Open Reading Frames (sORFs) embedded within multiple

20    eukaryotic non-coding RNAs can be translated into bioactive peptides of up to 100 amino acids in size.

21    However, the functional roles of the 607 Stress Induced Peptides (SIPs) previously identified from 189

22    Transcriptionally Active Regions (TARs) in *Arabidopsis thaliana* remain unclear. To provide a starting

23    point for function annotation of these peptides, we performed a large-scale prediction of peptide

24    binding sites on protein surfaces using and coarse-grained peptide docking. The docked models were

25    subjected to further atomistic refinement and binding energy calculations. A total of 530 peptide-

26    protein pairs were successfully docked. In cases where a peptide encoded by a TAR is predicted to

27    bind at a known ligand or cofactor-binding site within the protein, it can be assumed that the peptide

28    modulates the ligand or cofactor-binding. Moreover, we predict that several peptides bind at protein-

29    protein interfaces, which could therefore regulate the formation of the respective complexes. Protein-

30    peptide binding analysis further revealed that peptides employ both their backbone and side chain

31    atoms when binding to the protein, forming predominantly hydrophobic interactions and hydrogen

32    bonds. In this study, we have generated novel predictions on the potential protein-peptide interactions

33    in *A. thaliana*, which will help in further experimental validation.

34    **Author summary**

35    Due to their small size, short peptides are difficult to find and have been ignored in genome

36    annotations. Only recently, we have realized that these short peptides of less than 100 amino acids

37    may actually play an important role in the cell. Currently, there are no high-throughput methods to find

38    out what the functions of these peptides are in contrast with efforts that exist for 'normal' proteins. In

39    this work, we try to fill this gap by predicting with which larger proteins, the short peptides might

40    interact to exert their function. We find that many peptides bind to pockets where normally other

41    proteins or molecules bind. We thus think that these peptides that are induced by stress, may regulate

42    protein-protein and protein-molecule binding. We make this information available through our

43    database ARA-PEPs so that individual predictions can be followed up.

44    **Introduction**

45    Over the years, the functional importance of short plant signaling peptides has been overshadowed by

46    other groups of molecules. For instance, the phytohormone auxin was shown to be involved in

47    bidirectional polar transport across tissues, controlling plant growth-related processes (Grunewald &

48    Friml, 2010; Murphy *et al*, 2012). Furthermore, microRNAs are considered to be important signaling

49    molecules, regulating developmental processes in plants by moving from one cell to another over long

50    distances (Marín-González & Suárez-López, 2012). It was only within the last decade that the roles of

51    plant peptides in a wide variety of cellular functions were established by multiple studies

52    (Matsubayashi, 2011; Tavormina *et al*, 2015). Some of these peptides may be encoded by short Open

53    Reading Frames (sORFs), that were earlier assumed to be non-coding (Amor *et al*, 2009; Chen *et al*,

54    2015; Ladoukakis *et al*, 2011; Crappé *et al*, 2013; Ruiz-Orera & Messeguer, 2014; Andrews &

55    Rothnagel, 2014). Several recent studies clearly demonstrate that sORFs embedded within non-

56    coding RNAs (ncRNAs), intergenic regions and pseudogenes can indeed be translated into bioactive

57    peptides. In our previous work, we have identified several Transcriptionally Active Regions (TARs)

58    induced upon the application of biotic (*Botrytis cinerea*) and abiotic stress (Paraquat) in *Arabidopsis*

59    *thaliana*. These TARs could be translated into Stress-Induced Peptides (SIPs), which can be

60    specifically categorized depending on the applied stress condition into *Botrytis cinerea* Induced

61    Peptides (BIPs) and Oxidative Stress Induced Peptides (OSIPs) that we catalogued in a database

62    ARA-PEPs (Hazarika *et al*, 2017; De Coninck *et al*, 2013). Although some physiological effects of

63    sORF-encoded peptides have been discovered, the molecular mechanism by which they exert their

64    function through interaction with other molecules is largely unknown. We postulate that the peptides

65    could work through interactions with proteins, as protein-peptide interactions have previously been

66    well established as important mediators of protein-protein interactions, partaking in signal transduction,

67    cell-to-cell communication, protein trafficking and other regulatory pathways (London *et al*, 2010;

68    Schindler *et al*, 2015; Kilburg & Gallicchio, 2016; Petsalaki *et al*, 2009; Neduva & Russell, 2005;

69    Perkins *et al*, 2010; Pawson & Nash, 2003).

70

71    Peptide-mediated interactions constitute 15-40% of all protein-protein interactions (Petsalaki and

72    Russell 2008, Neduva and Russell, 2005). Most of the studies performed so far in order to understand

73    protein-peptide interactions focused on small peptides that may be short linear recognition motifs

74    originating from disordered protein regions (Kilburg & Gallicchio, 2016; London *et al*, 2010).

75    Investigating those interactions is experimentally challenging, and this has led to limited progress in

76    the field of protein-peptide interactions validation. On the other side, the successful modeling of such

77    complexes depends on prior structural knowledge of the protein that acts as a receptor. A number of

78    protein-peptide docking methods, such as Rosetta FlexPepDock (Raveh *et al*, 2010, 2011),

79    GalaxyPepDock (Ko *et al*, 2012; Heo *et al*, 2013; Lee *et al*, 2015), MedusaDock (Ding *et al*, 2010),

80    DynaDock (Antes, 2010), CABS-dock (Kurcinski *et al*, 2015; Wabik *et al*, 2015), pepATTRACT

81    (Schindler *et al*, 2015), HADDOCK (Dominguez *et al*, 2003; Trellet *et al*, 2013), and tools to predict

82    binding sites on proteins, such as PepSite2 (Trabuco *et al*, 2012; Petsalaki *et al*, 2009), have been

83    developed. Moreover, curated data also exists for characterization of protein-peptide interactions, e.g.

84    a non-redundant database of high-resolution peptide-protein complexes called the *peptiDB* (London *et*

85    *al*, 2010). Although docking strategies are the preferred methods for predicting protein-peptide

86    interactions, they are associated with certain limitations, such as difficulty in docking peptides longer

87    than 4 amino acids, owing to their high degree of conformational flexibility. In our current study, we

88    opted to combine the peptide-protein docking method pepATTRACT-local with binding site predictions

89    obtained from the PepSite2 server, that uses training data of known protein-peptide complexes from

90    Protein Data Bank (PDB) to define Spatial Position Specific Scoring Matrices (S-PSSMs).

91    Furthermore, as biological systems are not static, we also looked into dynamics of the obtained

92    docked models and calculated the energetics of binding based on multiple conformations that the

93    protein-peptide system can acquire in the solution.

94

95    We hypothesize that a SIP encoded by a TAR may bind on a protein at one of its pockets, or to a

96    known ligand or cofactor-binding site, and consequently affect the function of the protein as a whole.

97    Moreover, peptides may bind at the interfaces of multi-chain complexes and modulate their activity.

98    Protein-peptide interactions involve smaller interfaces whose affinity is usually weaker and are

99    transient as they can rapidly make and break interactions in response to sudden cellular perturbations,

100   for instance stress conditions (Stein & Aloy, 2008; Perkins *et al*, 2010).

101

102   In recent years, there has been a growing interest in developing protein-protein interaction inhibitors

103   based on peptides or peptide derivatives. Molecules that can mimic the binding or functional sites of

104   proteins are promising candidates for different types of biological applications. Synthetic peptides are

105   widely choiced molecules for mimicry of protein sites because they can be easily synthesized as exact

106   copies of protein fragments, or they may be generated by introducing diverse chemical modifications

107   to the peptide sequence, and/or by modifying the peptide backbone (Groß *et al*, 2016). Peptide mimics

4

108     have the potential to be developed as attractive targets for agriculture, especially plant disease control,

109     and for therapeutic interventions (Beekman & Howell, 2016).

110

111     As detailed above, when no information about the peptide-binding site on protein receptors is

112     available, there is need for computational approaches to predict peptide-binding sites on protein

113     surfaces, as these models can serve as starting points for experimental characterization of novel

114     protein-peptide interactions. This will be especially beneficial in studying the model plant *A. thaliana*, in

115     which peptides have multiple important roles, but have been understudied till now. Molecular docking

116     studies can be used effectively to explore the binding mode of putative peptides onto proteins, serving

117     as an excellent approach for *de novo* design of peptides targeting various other biosynthetic pathways

118     in major eukaryotes. In the current study, we investigate the potential roles of sORF-encoded stress

119     induced peptides in targeting the key regulatory enzymes, as this could further indicate their roles in

120     mediating the stress-response mechanisms.

121

122     **Results**

123         **Short peptides may exert their function by interacting with proteins**

124     In a previous study, we identified 189 TARs in response to plant oxidative stress by the herbicide

125     Paraquat and the fungus *Botrytis cinerea*, which could be translated into 607 SIPs (Hazarika *et al*,

126     2017; De Coninck *et al*, 2013). A peptide fragment library consisting of 23,113 k-mers, ranging from 4

127     to 10 amino acid residues, was generated and searched for potential binding sites on *A. thaliana*

128     proteins from the PDB repository, using PepSite2. We screened 12,540,140 protein-peptide pairs and

129     found 3,769,393 significant matches at PepSite2 score > 60 and *p*-value <= 0.1. We additionally

130     screened for short peptide motif matches on *A. thaliana* proteins using BLASTP and found 302

131     matches. We performed initial docking analysis using the pepATTRACT protocol, and a larger subset

132     of 576 protein-peptide pairs was devised by pooling together the above 302 protein-peptide pairs as

133     well as others with significant PepSite2 score (Figure 1A). The list of docked pairs can be accessed

134     through the url

135     (https://www.biw.kuleuven.be/CSB/ARA-PEPs/SIP_PDB_interactions.php). In our study, 46 protein-

136     peptide complexes failed to dock, and the reason could be the large conformational changes of the

137     protein caused by binding of a flexible peptide; this indeed remains a big problem of docking methods

5

138    (Trellet *et al*, 2013). From the docked models, we filtered out the 104 top protein-peptide pairs, and for

139    each of them characterized the protein-peptide binding, and determined the free energy of binding.

140    Our results show that there exists a huge repertoire of potential peptide binding sites on all available

141    *A. thaliana* proteins out of which we explored only 0.015%. Large-scale predictions of potential protein-

142    peptide pairs can aid in future experimental validations for understanding cell-to-cell communication

143    during plant development or stress-tolerance mechanisms.

144

145    Specific inhibitors may mimic portions of protein interfaces and can bind to a peptide binding pocket

146    located at the interface between two monomers. In our study we found that 15 peptides bind at the

147    interface between subunits of protein complexes, indicating the ability to modulate complex's activity

148    (Figure 1A, Supplementary table 2). Among the 15 models, in three of the cases the peptides bind in a

149    similar way to known characterized short peptides or portion of a full-length protein (Supplementary

150    Figure 3). We also observed that 30 peptides bind to a known ligand/cofactor binding site (Figure 1A,

151    Supplementary table 1). Ligand and protein binding sites may often overlap within protein families as it

152    has been shown that a peptide may compete with the ligand for the binding site or non-competitively

153    bind to the pocket along with the ligand molecule. The structural models obtained in the current study

154    will facilitate future validations as we can directly compare the binding modes of peptides on proteins.

155    While we observe that both pepATTRACT-local and pepATTRACT-blind produce similar results for

156    56% of the docked pairs, the other 44% of the pairs were not docked at the same site of the protein

157    using the two methods. Under assumption that the PepSite2 correctly indicated the binding regions,

158    these results show how the use of restraints in the docking protocol can help in concentrating the

159    search around relevant regions of the protein-peptide interaction space. Other reports have previously

160    shown that restraint-based dockings yield better results as compared to blind docking methods (Vazda

161    and Kozakov, 2009, de Vries et al., 2007).

162    In addition, we investigated if peptide binding pockets on proteins could bind multiple peptides, or in

163    other words, whether a peptide prefers to bind to one specific pocket on a protein. To test this, we

164    generated a randomly shuffled list of peptides while keeping the list of PDB structures intact, followed

165    by scanning for binding sites using PepSite2. In 95.15% of the cases the peptides preferred to bind to

166    the same pocket, while in only 4.84% of the cases the peptides were bound to different pockets on the

167    same receptor (Supplementary Figure 1A). It is possible that the S-PSSMs capture the binding modes

168    of amino acids in such a way that amino acids in a peptide sequence may prefer to bind to chemically

169    similar binding sites on proteins, e.g. hydrophobic amino acids from the peptide tend to bind

170    hydrophobic protein regions (Petsalaki *et al*, 2009) of appropriate sizes. While some reports suggest

171    that peptides often look for a large enough pocket to bind followed by latching onto it with the help of a

172    few hotspot residues (London et al, 2010, 2012), other reports suggest that several different peptides

173    are able to bind to the same protein domain by exhibiting special properties such as promiscuity

174    (Bhattacherjee & Wallin, 2013). Furthermore, the seemingly more important role of peptide backbone

175    compared to side chain atoms (detailed below) in protein binding provides another explanation for the

176    observed promiscuity, as backbone atoms are the same independent of the amino acid sequence of

177    the peptide.

178

179    We mapped 835 unique *A. thaliana* protein chains from the initial screening to Uniprot IDs using

180    annotations from SIFTS database (Velankar *et al*, 2013) and performed Gene Ontology

181    (GO) enrichment analysis using REACTOME_Pathways and GO_BiologicalProcess ontology. GO

182    analysis of the *A. thaliana* proteins with significant scores revealed that they may be categorized into 5

183    main groups *viz.* defense response, cellular response to organic cyclic compound, organonitrogen

184    compound metabolic process, regulation of stomatal movement and cellular response to chemical

185    stimulus (Figure 1B).

186

187    **Figure 1 (A)** Computational pipeline for prediction of protein-peptide pairs. 23,113 k-mers were

188    screened for binding sites on 1009 *A. thaliana* proteins. Out of 3,769,393 significant matches, 576

189    pairs were docked, and 104 pairs were further studied in detail. 30 peptides may bind to a ligand

190    binding/catalytic site on a protein and 15 peptides may bind at the dimer interface between 2 chains of

191    a protein complex. The peptide binding pocket is highlighted in yellow. **(B)** Histogram showing specific

192    GO terms related to the associated proteins from protein-peptide screening analysis. The bars

193    represent the number of proteins from the analyzed cluster associated with the term, and the label

194    displayed on the bars is the percentage of proteins compared to all proteins associated with the term.

195    The overview pie-chart presents functional groups for the proteins where the name of the group is

196    given by the most significant term in the group. GO enrichment analysis revealed 5 main groups and

197    each group section in the pie-chart correlates with the number of terms in each group.

198

### Characterization of protein-peptide binding interactions

200 We carried out a general characterization of protein-peptide interactions in the 104 docked pairs. We

201 extracted top 10 models from the docking results after iATTRACT refinement and calculated the

202 average number of interactions within each docked pair. Each receptor atom that comes within 4 Å of

203 any ligand atom is considered as a close contact. We determined the mean number of close contacts

204 per docked pair to be 167±49. All protein-peptide pairs interacted with each other using hydrogen

205 bonds and hydrophobic interactions. The mean number of hydrophobic interactions per system is

206 24±8, and the mean number of hydrogen bonds per docked pair is 4.5±2, where donors of hydrogen

207 bonds are localized on protein in 51% of the cases (Figure 2E). While hydrogen bonds and

208 hydrophobic interactions are omnipresent, not all pairs formed salt bridges and pi-pi stackings (Figure

209 2E and 2D). Within the peptides in our top 104 models, 22.5% of the total amino acid residues are

210 charged (Arg, Lys, His, Asp, Glu) and 14.4% residues are aromatic (His, Phe, Tyr, Trp). In agreement

211 with the amino acid composition, we also found salt bridges to be more prevalent than pi-pi stackings

212 in the docked systems: 52% of protein-peptide pairs contain salt bridges with the mean number per

213 system 1.72±1, and 7.7% form pi-pi stackings with the mean of 1.63±0.7. An additional 5.8% of pairs

214 have both salt bridges and pi-pi stackings, while the remaining 34.6% do not form any salt bridges or

215 pi-pi stackings (Figure 2D).

216 In total, 23.8% of peptide side chain and 38.7% of backbone atoms participate in close contacts, with

217 an average number of contacts per interacting atom being 3.6 and 2.8, respectively (Figure 2E). In

218 average, the ratio of the unique peptide side chain:backbone atoms involved in close contacts is 2:1

219 for the top model in 104 docked systems, while the overall side chain:backbone ratio of atoms is 3.2:1.

220 If side chain and backbone atoms of the peptide were equally important in protein binding, we would

221 expect the ratio of atoms involved in close contacts to also be 3.2:1. Instead, its lower value suggests

222 that peptide backbone atoms might be more important in protein-peptide interactions than the side

223 chain ones.

224 We determined the overall hydropathy index for all the peptide fragments (576) and found that 60.2%

225 of the peptides are hydrophilic, while the remaining 39.8% are hydrophobic in nature (Figure 2C). A

226 larger fraction of SIPs in our dataset have high hydrophilicity or a lower GRAVY index score,

227 suggesting that they may mainly interact with globular proteins rather than with hydrophobic regions

228    that spans membranes. Peptides with fewer ionic/charged groups are generally less soluble in water

229    and are therefore prone to aggregation and interacting with hydrophobic pockets of larger proteins.

230    After analysis of the docked models, we took a further look into the dynamics of the top model from

231    each of the 104 protein-peptide docked pairs and calculated the free energy of binding based on 100

232    conformational snapshots from molecular dynamics for each system (Figure 2A, B). Per-residue

233    decomposition of protein-peptide binding energies also allowed identification of amino acid types that

234    frequently (in multiple systems) have significant binding contribution (Supplementary Figure 2A). For

235    instance, arginine residue, located in proteins at the peptide binding interface, stands out as a

236    recurring amino acid with significant stabilizing effect on the binding (negative value of the $\Delta G_{bind}$

237    contribution). Interestingly, a prevalent contribution of negatively charged peptide amino acids is

238    seemingly lacking. Visual investigation of trajectories obtained by molecular dynamics shows that

239    arginines make salt bridges with negatively charged carboxyl groups of peptide C-terminal amino

240    acids in 70% of cases (31/44), making this interaction independent of amino acid type present in the

241    peptide. Other prominent protein residues that predominantly stabilize interactions with the peptides

242    are the charged (Glu, Asp) and aromatic ones (Trp, Phe, Tyr).

243    Local destabilizing effect on binding is shown by different peptide amino acids, containing side chains

244    of largely different properties (Supplementary Figure 2A). However, a more detailed view reveals that

245    this is a consequence of amino acid location within a peptide, rather than its chemical composition

246    (Supplementary Figure 2B). In average, non-terminally located amino acids contribute to the binding in

247    a stabilizing manner, N-termini destabilize protein-peptide interaction, while C-terminal amino acids

248    have different average effect depending on amino acid type, and rarely have significant contribution

249    (Supplementary Figure 2B, C). Positively charged arginine residue in peptide is an interesting

250    example: its overall contribution is stabilizing (Supplementary Figure 2A) but depending on its position

251    within the peptide it shows effects that range from stabilizing to destabilizing (Supplementary Figure

252    2B). The same holds true for several other amino acids.

253    Overall, the largest destabilizing factor in binding across the 104 top protein-peptide models is the

254    inability of protein to stabilize the N-terminal positively charged amino group of the peptide. However,

255    the negative $\Delta G_{bind}$ values for almost all systems (101 out of 104; Figure 2A, B) show that this factor is

256    insufficient to destabilize the overall binding of the peptide to the predicted part of the protein.

257

258    **Figure 2**. Characterization of peptides and protein-peptide binding interactions. **(A)** Histogram

259    showing $\Delta G_{bind}$ values for 104 top models **(B)** $\Delta G_{bind}$ values as individual data points **(C)** Hydropathy

260    index for all the 576 peptides predicted to interact with *A. thaliana* proteins **(D)** Total number of

261    charged and aromatic residues in the peptides that interact with proteins. The plot also shows the total

262    number of salt-bridges and pi-pi stackings formed in the top models. **(E)** Different types of interactions

263    formed by the protein-peptide pairs

264

265

266    **Peptide BIP142_3 (LAEDTFGEIS)  binds to CRYD protein**

267    The 10-mer peptide fragment LAEDTFGEIS from BIP142_3/OSIP134_3 can be translated from

268    BcTAR142/PQTAR134, expressed under stress conditions involving either *B. cinerea* or Paraquat.

269    The above TAR is expressed solely under stress conditions and shows no expression under mock

270    treatments (Figure 3A). We split the entire 41 AA long peptide sequence of BIP142_3/OSIP134_3 into

271    10-mer fragments using a sliding window and scanned against all *A. thaliana* PDB structures. High

272    confidence peptide bindings with *p*-value less than 0.001 were retained (Figure 3B). We picked the

273    protein-peptide model LAEDTFGEIS-2VTB(D) because the mentioned 10-mer shows sequence

274    similarity to cryptochrome DASH (CRYD) protein (chain D in 2VTB PDB structure), hence might act as

275    a peptide mimic or affect the activity of the protein by binding to one of its pockets. Members of the

276    cryptochrome DASH subclade are involved in the DNA repair of cyclobutane pyrimidine dimers in

277    single stranded DNA (Selby & Sancar, 2006). Cryptochromes in general are photolyase-like

278    flavoproteins that mediate blue-light regulation of gene expression and photomorphogenic responses,

279    including abiotic stress responses in *Arabidopsis*, as well as in all kingdoms of life (Yu et al., 2010).

280

281    For the LAEDTFGEIS-2VTB(D) model, binding sites were predicted using PepSite2, and a coarse

282    model of the interacting residues between peptide and protein was built. The peptide bound to *A.*

283    *thaliana* CRYD shows one of the strongest bindings among the 104 top models, with $\Delta G_{bind}$ value of -

284    48.02 kcal mol$^{-1}$. Several residues in this complex have large contribution to this binding (represented

285    as sticks at Figure 3C), with all except N-terminal leucine of the peptide contributing in a stabilizing

286    manner. Protein and peptide are bound via various types of interactions: salt bridges (Arg 436 and Asp

287    4; Arg 487, Arg 490 and Ser 10), stacking interactions (Trp 365 and Phe 6) and hydrogen bonds (Trp

10

288    365 and Thr 5). Destabilizing effect of N-terminal peptide residue, found in multiple other systems as

289    well (Supplementary Figure 2B), is likely the consequence of the lack of a negatively charged residue

290    at the corresponding position in the protein, which would make favourable interactions with N-terminal

291    amine group.

292

293    Two cofactors can bind to CRYD: flavin adenine dinucleotide (FAD) and 5,10-methenyltetrahydrofolate

294    (MTHF), out of which the first one is necessary for catalytic activity. According to UniProt database,

295    amino acids Arg 436 and Asp 485, which coincide with LAEDTFGEIS binding site, are involved in ATP

296    binding. If the peptide indeed binds CRYD in a way predicted in this study, it could block FAD binding

297    or even bind simultaneously with it, therefore having an effect on the activity of this enzyme, and

298    consequently on the aforementioned type of DNA repair (Figure 3D).

299

300    **Figure 3**. **(A)** Overview of Transcriptionally Active Region BcTAR142/PQTAR134 induced under

301    stress conditions which might encode a short peptide. The TAR shows mRNA expression levels under

302    treated (PQ and BC) and mock conditions (mock_BC and mock_PQ). **(B)** Screening of all possible

303    peptide fragments from BIP142_3/OSIP134_3 against all *A. thaliana* proteins in PDB. **(C)** A coarse

304    model of peptide LAEDTFGEIS bound to CRYD protein (chain D of 2VTB). Restraint based docking

305    was performed, followed by surrounding of the 3D model with explicit water. The solvated structure

306    was optimized and then used for MD simulation. The conformational snapshots from the MD were

307    used to calculate $\Delta G_{bind}$ value for protein-peptide binding, and for visual inspection of the mode of

308    binding (details in the text). Finally, superposition of the docked model (CRYD as grey, and docked

309    LAEDTFGEIS as magenta surface) and the original PDB structure, which has FAD (green) and MTHF

310    (yellow) co-factors bound, suggests that peptide might have an effect on FAD binding, and

311    consequently on CRYD's function.

312

313    **Discussion**

314    Steroids, peptides and other small bioactive compounds mainly regulate cellular communication in

315    eukaryotes, including plants. Over the last decade, an increasing number of secreted peptides have

316    been shown to influence a variety of developmental processes in plants, such as meristem size, root

317    growth, stomatal differentiation, and organ abscission (Butenko & Aalen, 2012). sORFs that might

11

318   encode peptides have been overlooked in gene prediction programs owing to their small size.

319   Moreover, there exist only a handful of publicly available T-DNA insertion collections of peptide

320   encoding genes (Butenko *et al*, 2014; Lease & Walker, 2006). In this study, we predict that a large

321   number of SIPs in *A. thaliana* may exert their function through protein-peptide interactions, by binding

322   on protein surfaces. We found 30 peptides that may bind at known ligand/cofactor binding sites on

323   proteins. The identification of ligand/cofactor binding sites in protein structures can aid in determination

324   of peptide ligand types and experimental validation of the function of the receptor (Glaser *et al*, 2006).

325   A peptide may compete for the ligand-binding site on the receptor, or it may non-competitively bind to

326   the pocket together with the ligand molecule and play a role in modulating the receptor. Additionally,

327   we screened 15 peptides that may bind to a pocket at the interface between two monomers of a multi-

328   chain complex. The design of peptides and peptidomimetics that mimic portions of dimeric/multimeric

329   protein interfaces have been shown to be an useful approach for the discovery of inhibitors that bind at

330   protein-protein interfaces (Cardinale *et al*, 2011). Currently, there is a lot of interest in drugs that can

331   inhibit dimerization of a functionally obligate homodimeric enzyme. However, design of peptides that

332   may disrupt protein-protein interactions is far more challenging than designing enzyme active site

333   inhibitors, due to factors such as the large interfacial areas involved, and flat and featureless

334   topologies that these binding surfaces may exhibit (Fletcher and Hamilton, 2006).

335

336   The characterization of protein-peptide interactions can be used to evaluate the binding affinity of the

337   model. One of the major factors determining the binding of a peptide to a protein is the size of the

338   pocket (Laskowski *et al*, 1996). In our study, binding analysis of the predicted protein-peptide pairs

339   revealed that different peptides tend to bind in the same pocket of one protein. We also observed that

340   the peptides mostly interact with proteins with the help of side chains, but this is due to the reason that

341   we have 3.2 times more of side chain atoms than backbone atoms. However, the peptide backbone

342   atoms participate in more unique protein-peptide interactions as compared to the side chains. This is

343   in agreement with another finding where peptides use more H-bonds in binding to their protein partner

344   involving the peptide backbone. In the *PeptiDB* dataset comprising of 103 protein-peptide complexes,

345   19 peptides bind as β-strands, which use far more H-bonds on average, while 18 peptides were bound

346   as α-helices, which form less H-bonds with proteins and contain more nonpolar atoms at the interface

347   (London *et al*, 2010).

348    The pepATTRACT-local docking method has advantages over other protein-peptide docking methods.

349    First, pepATTRACT-local docking outperforms blind docking whose performance is similar with other

350    local docking methods (Schindler *et al*, 2015). Second, this approach completes a run in about one

351    hour for each pair, which is beneficial for large-scale prediction of protein-peptide interactions. In our

352    study, 46 protein-peptide pairs failed to dock. The reason may be due to large conformational changes

353    upon peptide binding onto the receptor, which still remains a huge problem while trying to accurately

354    predict interactions (Trellet *et al*, 2013). Some failed cases reveal that the peptide is deeply buried into

355    the protein surface. These failed pairs can be docked using other local docking methods. However, for

356    most of protein-peptide interactions, only very small conformational changes upon peptide binding

357    have been observed on the protein surface.

358    The *A. thaliana* genome may encode thousands of small proteins that could function as peptide

359    signals and more than 600 plasma membrane-bound receptor-type proteins that could act as

360    receptors for peptide ligands (Shiu & Bleecker, 2001). Several sORF-encoded peptides may target

361    regulatory enzymes involved in metabolic pathways by downregulating or upregulating the activity of

362    these key enzymes. Predicting potential protein-peptide pairs and confirmation of physical interaction

363    between these pairs is crucial to advance our understanding of cell-to-cell communication during plant

364    development or stress-tolerance mechanisms (Murphy *et al*, 2012). Nevertheless, there are a few

365    protein-peptide pairs that have been quite comprehensively studied such as the CLE (CLV3/ESR;

366    CLAVATA3/EMBRYO SURROUNDING REGION-related) family peptides, which are plant-specific

367    peptide hormones that mediate cellular communication and are involved in meristem maintenance,

368    vascular development and nematode feeding cell formation. Apart from these CLE peptides there may

369    be many more peptides that remain to be identified. In general, our study aims at predicting potential

370    protein-peptide interactions on protein surfaces which can be experimentally validated by researchers

371    in the future.

372

373    **Materials and Methods**

374        **Screening of peptide binding pockets on protein surfaces**

375    We generated a peptide fragment library consisting of 23,113 k-mers ranging in size from 4 to 10

376    amino acids, using sequences of SIPs (Hazarika *et al*, 2017; De Coninck *et al*, 2013) and following a

377    sliding window approach. We extracted 2,561 structures corresponding to 1,009 *A. thaliana* proteins

378    from Protein Data Bank, PDB (www.rcsb.org) (Berman *et al*, 2000). 996 structures were retained after

379    filtering out redundant ones. We carried out an all-vs-all screening of potential peptide binding sites on

380    *A. thaliana* proteins using PepSite2 (Petsalaki *et al*, 2009; Trabuco *et al*, 2012), and retained motif

381    matches with score > 60 and *p*-value <= 0.1. The reliability of the PepSite2 method is based on the

382    measure of positive predictive value. For *p*-values below 0.003, false positive rate was reported to be

383    0.01, and true positive rate was 0.1 representing a positive predictive value of 89.9%.

384    We used the default settings of BLASTP2.2.28+ algorithm (Altschul *et al*, 1990) to screen out peptides

385    that show sequence similarity with a protein chain, as entire or a part of a SIP may mimic a specific

386    binding motif on the protein, or resemble a loop from a large structured protein, a disordered region in

387    protein termini or interfaces between defined domains (London *et al*, 2010; Kilburg & Gallicchio, 2016).

388    **Building peptide models, docking and structure refinement**

389    We shortlisted 576 protein-peptide pairs with *p*-value < 0.1 from the PepSite2 output in order to build

390    atomistic models and perform docking studies, using the protein-peptide coarse-grained *ab initio*

391    docking protocol pepATTRACT (Schindler *et al*, 2015). For each peptide, three idealized peptide

392    conformations (extended, α-helical and polyproline) were built using the Python library PeptideBuilder

393    (Tien *et al*, 2013). The backbone dihedral angles used to represent the three peptide conformations

394    were α-helical ($\Phi$= –57°, $\Psi$ = –47°), extended ($\Phi$ = –139°, $\Psi$ = –135°), and polyproline conformations

395    ($\Phi$ = –78°, $\Psi$ = 149°) (Trellet *et al*, 2013).

396    In the current study, the rigid body docking models were ranked by ATTRACT score, and the top-

397    ranked 100 structures were subjected to atomistic refinement using the flexible interface refinement

398    method iATTRACT. We used the distance restraint based local docking protocol of pepATTRACT to

399    restrict the sampling during rigid body sampling stage and flexible refinement stage towards the

400    PepSite2 predicted interface residues. The placement of peptide and protein was optimized during

401    iATTRACT refinement. At this stage, the interface region of the peptide and the protein were treated

402    as fully flexible, while simultaneously optimizing the center of mass position and orientation of the

403    peptide.

404    **Molecular dynamics**

405    **Preparation and parametrization**

406    We used *pdb4amber* from Amber16 (Case *et al*, 2017) to make the pdb files of 104 high confidence

407    protein-peptide complexes top models suitable for using this software package. All disulfide bonds

408    detected by *pdb4amber* were retained in the system. Detected protein gaps were treated by addition

409    of *N*-methyl group (NME) to the carbon of the backbone amide group of the C-terminal, and acetyl

410    group (ACE) to the backbone nitrogen of the N-terminal amino acid, using PyMol Molecular Graphics

411    System, Version v1.7.4.4, Schrodinger, LLC (Delano, 2002). Capping prevents the amino acids that

412    are flanking the gap from being recognized as protein termini, and therefore charged.

413        Parametrization of the systems was done using *teLeap* from Amber16. Counter-charged ions

414    ($Na^+$ or $Cl^-$) were added to the non-neutral systems, and each protein-peptide complex was

415    surrounded by a rectangular box of explicit TIP3P water spanning 10 Å from the system. Force field

416    *ff14SB* was used for parametrization of proteins and peptides, and *tip3p* for parametrization of water.

417    Joung/Cheatham parameters were employed for monovalent ions in the chosen water type.

**Optimization**

419    Systems were optimized in 25,000 steps divided in five cycles, using *sander* from Amber16. First

420    1,000 steps of each cycle were performed by steepest descent method, while conjugated gradient was

421    used for the remaining steps. In first three cycles, the constraint was applied to 1. the entire protein, 2.

422    heavy protein:peptide atoms, and 3. backbone atoms, using force constant 100 kcal $mol^{-1}$ $Å^{-2}$.

423    Constraint on backbone atoms was reduced to 50 kcal $mol^{-1}$ $Å^{-2}$ in the fourth cycle, and no constraints

424    were applied in the fifth.

**Molecular dynamics simulations**

426    After optimization, each system was equilibrated during the initial 500 ps, using *pmemd* from Amber16

427    package. In the first 300 ps, the canonical NVT ensemble was simulated, with constraint applied to

428    atoms in the protein:peptide complex using the force constant 25 kcal $mol^{-1}$ $Å^{-2}$. Temperature was

429    increasing from 0 to 300 K during the first 250 ps. In the last 200 ps of equilibration, isothermal-

430    isobaric ensemble NpT was simulated, with temperature held constant at 300 K and pressure at 1.0

431    bar, with no constraints applied to the system. Throughout equilibration, the SHAKE algorithm was

432    used to apply constraints on bonds containing hydrogen atoms, and time step of 2 fs was used. The

433    cutoff distance for non-bonded interactions was set to 15 Å, and the neighbor list was updated each

434    20 steps.

435    Production phase was done as a 4.5 ns continuation of the 500 ps long equilibration, using Gromacs 5

436    software (Lindahl *et al*, 2001; Hess *et al*, 2008; Van Der Spoel, 2005; Berendsen *et al*, 1995;

437    Essmann, 1995; Hess, 2007; Miyamoto & Kollman, 1992; Bussi *et al*, 2007). Conversion from Amber

438    to Gromacs file formats was performed with the help of ParmEd 2.7 tool (Swails *et al*). Constraint on

439    bonds that contain hydrogen atoms was applied using LINCS algorithm, the time step was 2 fs, and

440    the coordinates were written each picosecond. The temperature and pressure were kept at 300 K and

441    1.0 bar using modified Berendsen thermostat for temperature, and Parrinello-Rahman barostat for

442    pressure coupling. Particle mesh Ewald method☐ was used for electrostatic interactions, the cutoff

443    distance for non-bonded interactions was 12 Å, and the neighbor list was updated each 20 steps.

444    Periodic boundary conditions were applied throughout equilibration and production phase.

445    **Analysis and binding energy calculation**

446    The obtained trajectories were visualized by Visual Molecular Dynamics VMD program (Humphrey *et*

447    *al*, 1996), and tools from Gromacs package were used to correct for periodic boundary conditions and

448    calculate root mean square deviation (RMSD) of complexes' backbones. Matplotlib (Hunter, 2007) ☐

449    was   used   to   visualize   the   results   of   analyses.

450    Molecular Mechanics energies with Generalized Born and Surface Area continuum solvation

451    (MM/GBSA) method was used to calculate the Gibbs energy of protein:peptide binding in the 104 top

452    docking models. The binding energy is calculated as the following average:

453    $$\Delta G_{bind} = \overline{G_{complex} - G_{free} - G_{peptide}} \qquad (1)$$

454    with each Gibbs energy term being the following sum:

455    $$G = E_b + E_{el} + E_{vdw} + G_{GB} + G_{SA} - TS \qquad (2)$$

456    where the bonded, electrostatic and van der Waals interaction energies terms are obtained by

457    molecular mechanics, the polar solvation term by generalized Born, the non-polar solvation term from

458    linear relation to the solvent accessible surface area, while the entropy term is often omitted

459    (Genheden & Ryde, 2015), as is in this study.

460    Amber *MMPBSA.py.MPI* was used here to calculate $\Delta G_{bind}$ for protein:peptide systems by

461    MM/GBSA method, using a single trajectory of the complex. The topology files of dry complexes, as

462    well as ligand (peptide) and receptor (protein), were prepared with Amber *ante-MMPBSA.py*. The

463    Gibbs energy terms in equation (1) were calculated for 100 conformational snapshots from the last 2.5

464    ns of the production phase for each system, using salt concentration of 0.15 mol dm$^{-3}$. During the

465    MM/GBSA calculations, per-residue binding energy decomposition was also performed in order to get

466    insight into contributions of specific protein and peptide residues to binding.

467    **Amino acids contribution to binding**

16

468   The output files of the MM/GBSA per-residue energy decomposition were used to analyze the

469   characteristics of protein:peptide binding. In each of the 104 systems, the residue with the largest

470   contribution to binding, either in stabilizing or destabilizing manner, was detected. The threshold was

471   then set to 40 % of its binding energy contribution value, and all residues that contributed more than

472   the threshold in a given system were taken for the analysis, with taking into account whether amino

473   acid belongs to protein or peptide. The number of appearances of individual amino acid was then

474   calculated, as well as average binding contribution of different amino acids, separately for proteins and

475   peptides, using Python.

476        **Characteristics of SIPs and *A. thaliana* proteins**

477   We scanned SIPs for hydrophobicity using the grand average of hydropathy (GRAVY) number, which

478   is a measure of the hydrophobicity/hydrophilicity of a protein based on Kyte and Doolittle equation.

479   The hydropathy values range from -2 to +2 for most proteins, with the positively ranked proteins being

480   more hydrophobic.

481   Gene ontology (GO) analysis was performed using the ClueGo Cytoscape plugin (Bindea *et al.*, 2009).

482   Lists of 835 unique proteins from the initial screening analysis were mapped to corresponding Uniprot

483   IDs   using   mappings   from   SIFTS   database   (Velankar   *et   al*,   2013)

484   (www.ebi.ac.uk/pdbe/docs/sifts/index.html).   The   list   of   proteins   was   used   to   query

485   REACTOME_Pathways and GO_BiologicalProcess ontology and the type of evidence set was

486   All_experimental. Pathways with *p-values* ≤ 0.05 were displayed, the minimum GO tree interval was

487   set as 3 and the maximum level was set as 8, the GO term/pathway selection was set as a threshold

488   of 4% of genes per pathway and the kappa score was set as 0.4.

489        **Analysis of interactions at the protein-peptide interface**

490   We manually inspected the top 10 models for each docked protein-peptide pair predicted by

491   pepATTRACT using molecular visualization softwares UCSF Chimera (Pettersen *et al*, 2004) and

492   PyMOL Molecular Graphics System. PDBeMotif, a web server for checking the PDB structure for

493   ligands and binding sites (Gutmanas *et al*, 2014) and Catalytic Site Atlas, a database of enzyme active

494   sites and catalytic residues on enzymes (Porter, 2004) was used for finding ligand/cofactor binding

495   sites and enzyme active sites respectively. We analyzed if a specific peptide binding site lies at the

496   interface of multi-chain proteins and assumed that residues on the 2 chains less than 6.0 Å apart were

17

497    interacting residues. The distance between Cα atoms located in chains A and B, with coordinates $A(x_1, y_1, z_1)$

498    and $B(x_2, y_2, z_2)$, was calculated according to the Euclidean distance equation $D(A,B)$: $\sqrt{\{(x_1-x_2)^2}$

499    $+ (y_1-y_2)^2 + (z_1-z_2)^2\}}$. All calculations were performed using the Biopython package from Python.

500    Protein-peptide bindings were characterized using BINding ANAlyzer (BINANA) (Durrant &

501    McCammon, 2011), HBPLUS (McDonald & Thornton, 1994) and Protein-Ligand Interaction Profiler

502    (PLIP) (Salentin *et al*, 2015) tools. BINANA was used to characterize important protein-ligand

503    interactions such as close contacts (any receptor atom within 4.0 Å of the ligand atoms), hydrogen

504    bonds (distance cutoff = 4.0 Å and angle cutoff <= 40°), hydrophobic contacts (ligand carbon atom

505    within 4.0 Å of a receptor carbon atom), salt bridges and pi-pi interactions.

506

**References:**

Altschul SF, Gish W, Miller W, Myers EW & Lipman DJ (1990) Basic local alignment search tool. *J. Mol. Biol.* **215:** 403–10

Amor B Ben, Wirth S, Merchan F, Laporte P, D'Aubenton-Carafa Y, Hirsch J, Maizel A, Mallory A, Lucas A, Deragon JM, Vaucheret H, Thermes C & Crespi M (2009) Novel long non-protein coding RNAs involved in Arabidopsis differentiation and stress responses. *Genome Res.* **19:** 57–69

Andrews SJ & Rothnagel JA (2014) Emerging evidence for functional peptides encoded by short open reading frames. *Nat. Rev. Genet.* **15:** 193–204

Antes I (2010) DynaDock: A now molecular dynamics-based algorithm for protein-peptide docking including receptor flexibility. *Proteins Struct. Funct. Bioinforma.* **78:** 1084–1104

Baker EN & Hubbard RE (1984) Hydrogen bonding in globular proteins. *Prog. Biophys. Mol. Biol.* **44:** 97–179

Beekman AM & Howell LA (2016) Small-Molecule and Peptide Inhibitors of the Pro-Survival Protein Mcl-1. *ChemMedChem* **11:** 802–813

Belkhadir Y, Wang X & Chory J (2006) Arabidopsis Brassinosteroid Signaling Pathway. *Sci. STKE* **2006:** cm5-cm5

Belkhadir Y, Yang L, Hetzel J, Dangl JL & Chory J (2014) The growth-defense pivot: Crisis management in plants mediated by LRR-RK surface receptors. *Trends Biochem. Sci.* **39:** 447–456

Berendsen, H. J. C., van der Spoel, D. & van Drunen, R. GROMACS: A message-passing parallel molecular dynamics implementation. *Comput. Phys. Commun.* **91,** 43–56 (1995).

Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN & Bourne PE (2000) The protein data bank. *Nucleic Acids Res.* **28:** 235–242

538     Bhattacherjee A & Wallin S (2013) Exploring Protein-Peptide Binding Specificity through

539        Computational Peptide Screening. *PLoS Comput. Biol.* **9:**

540     Bindea,G. *et al.* (2009) ClueGO: A Cytoscape plug-in to decipher functionally grouped gene ontology

541        and pathway annotation networks. *Bioinformatics*, **25**, 1091–1093.

542

543     Bosshard HR, Marti DN & Jelesarov I (2004) Protein stabilization by salt bridges: concepts,

544        experimental approaches and clari cation of some misunderstandings. *J. Mol. Recognit.* **17:** 1–

545        16

546     Bussi, G., Donadio, D. & Parrinello, M. Canonical sampling through velocity rescaling. *J. Chem. Phys.*

547     **126,** 14101 (2007).

548     Butenko MA & Aalen RB (2012) Receptor Ligands in Development. In *Receptor-like Kinases in Plants:*

549        *From Development to Defense*, Tax F & Kemmerling B (eds) pp 195–226. Berlin, Heidelberg:

550        Springer Berlin Heidelberg

551     Butenko M a, Wildhagen M, Albert M, Jehle A, Kalbacher H, Aalen RB & Felix G (2014) Tools and

552        Strategies to Match Peptide-Ligand Receptor Pairs. *Plant Cell* **26:** 1838–1847

553     Cardinale D, Guaitoli G, Tondi D, Luciani R, Henrich S, Salo-Ahen OMH, Ferrari S, Marverti G,

554        Guerrieri D, Ligabue A, Frassineti C, Pozzi C, Mangani S, Fessas D, Guerrini R, Ponterini G,

555        Wade RC & Costi MP (2011) Protein–protein interface-binding peptides inhibit the cancer

556        therapy target human thymidylate synthase. *Proc. Natl. Acad. Sci.* **108:** E542–E549

557

558     D.A. Case, D.S. Cerutti, T.E. Cheatham, III, T.A. Darden, R.E. Duke, T.J. Giese, H. Gohlke, A.W.

559     Goetz, D. Greene, N. Homeyer, S. Izadi, A. Kovalenko, T.S. Lee, S. LeGrand, P. Li, C. Lin, J. Liu, T.

560     Luchko, R. Luo, D. Mermelstein, K.M. Merz, G. Monard, H., D. M. Y. and P. A. K. AMBER 2017.

561     (2017)

562     Chen M, Chen J & Zhang D (2015) Exploring the secrets of long noncoding RNAs. *Int. J. Mol. Sci.* **16:**

563        5467–5496

564     Clark S, Running M & Meyerowitz E (1995) *CLAVATA3* is a specific regulator of shoot and floral

565        meristem development affecting the same processes as *CLAVATA1*. *Development* **121:** 2057–

566        2067

567     Clark        SE,        Williams        RW        &        Meyerowitz        EM        (2017)        The

568  <em><strong>CLAVATA1</strong></em><strong>Gene Encodes a Putative Receptor Kinase

569  That Controls Shoot and Floral Meristem Size in Arabidopsis</strong>. *Cell* **89:** 575–585

570  De Coninck B, Carron D, Tavormina P, Willem L, Craik DJ, Vos C, Thevissen K, Mathys J & Cammue

571  BP a (2013) Mining the genome of Arabidopsis thaliana as a basis for the identification of novel

572  bioactive peptides involved in oxidative stress tolerance. *J. Exp. Bot.* **64:** 5297–5307

573  Crappé J, Van Criekinge W, Trooskens G, Hayakawa E, Luyten W, Baggerman G & Menschaert G

574  (2013) Combining in silico prediction and ribosome profiling in a genome-wide search for novel

575  putatively coding sORFs. *BMC Genomics* **14:** 648

576  Dagliyan O, Proctor EA, D'Auria KM, Ding F & Dokholyan N V. (2011) Structural and dynamic

577  determinants of protein-peptide recognition. *Structure* **19:** 1837–1845

578  Delano WL (2002) The PyMOL Molecular Graphics System.

579  Ding F, Yin S & Dokholyan N V (2010) Rapid Flexible Docking Using a Stochastic Rotamer Library of

580  Ligands. *J. Chem. Inf. Model.* **50:** 1623–1632

581  Dominguez C, Boelens R & Bonvin AMJJ (2003) HADDOCK: A protein-protein docking approach

582  based on biochemical or biophysical information. *J. Am. Chem. Soc.* **125:** 1731–1737

583  Durrant JD & McCammon JA (2011) BINANA: A Novel Algorithm for Ligand-Binding Characterization.

584  *J. Mol. Graph. Model.* **29:** 888–893

585  Essmann, U. *et al.* A smooth particle mesh Ewald method. *J. Chem. Phys.* **103,** 8577–8593 (1995).

586  Fletcher S & Hamilton AD (2006) Targeting protein-protein interactions by rational design: mimicry of

587  protein surfaces. *J. R. Soc. Interface* **3:** 215–33

588  Genheden, S. & Ryde, U. The MM/PBSA and MM/GBSA methods to estimate ligand-binding affinities.

589  Expert Opin. Drug Discov. 10, 0 (2015).

590  Glaser F, Morris RJ, Najmanovich RJ, Laskowski RA & Thornton JM (2006) A method for localizing

591  ligand binding pockets in protein structures. *Proteins Struct. Funct. Genet.* **62:** 479–488

592  Groß A, Hashimoto C, Sticht H & Eichler J (2016) Synthetic Peptides as Protein Mimics. *Front.*

593  *Bioeng. Biotechnol.* **3:**

594  Grunewald W & Friml J (2010) The march of the PINs: developmental plasticity by dynamic polar

595  targeting in plant cells. *EMBO J.* **29:** 2700–2714

596  Gutmanas A, Alhroub Y, Battle GM, Berrisford JM, Bochet E, Conroy MJ, Dana JM, Fernandez

597  Montecelo MA, Van Ginkel G, Gore SP, Haslam P, Hatherley R, Hendrickx PMS, Hirshberg M,

598    Lagerstedt I, Mir S, Mukhopadhyay A, Oldfield TJ, Patwardhan A, Rinaldi L, et al (2014) PDBe:

599    Protein data bank in Europe. *Nucleic Acids Res.* **42:** 308–317

600    Hazarika RR, Coninck B De, Yamamoto LR, Martin LR, Cammue BPA & Noort V Van (2017) ARA-

601    PEPs: a repository of putative sORF- encoded peptides in Arabidopsis thaliana. *BMC*

602    *Bioinformatics***:** 1–9

603    Heo L, Park H & Seok C (2013) GalaxyRefine: Protein structure refinement driven by side-chain

604    repacking. *Nucleic Acids Res.* **41:** 384–388

605    Hess, B. P-LINCS: A Parallel Linear Constraint Solver for Molecular Simulation. (2007).

606    doi:10.1021/CT700200B

607    Hess, B., Kutzner, C., van der Spoel, D. & Lindahl, E. GROMACS 4: Algorithms for Highly Efficient,

608    Load-Balanced, and Scalable Molecular Simulation. J. Chem. Theory Comput. 4, 435–447

609    (2008).

610    Hothorn M, Belkhadir Y, Dreux M, Dabi T, Noel JP, Wilson IA & Chory J (2011) Structural basis of

611    steroid hormone perception by the receptor kinase BRI1. *Nature* **474:** 467–471

612

613    Humphrey, W., Dalke, A. & Schulten, K. VMD: visual molecular dynamics. *J. Mol. Graph.* **14,** 33–8,

614    27–8 (1996).

615    Hunter, J. D. Matplotlib: A 2D graphics environement. *Comput. Sci. Eng.* **9,** 90–95 (2007).

616    Kilburg D & Gallicchio E (2016) Recent Advances in Computational Models for the Study of Protein-

617    Peptide Interactions. *Adv. Protein Chem. Struct. Biol.* **105:** 27–57

618    Ko J, Park H & Seok C (2012) GalaxyTBM: template-based modeling by building a reliable core and

619    refining unreliable local regions. *BMC Bioinformatics* **13:** 198

620    Kucera M, Isserlin R, Arkhangorodsky A & Bader GD (2016) AutoAnnotate: A Cytoscape app for

621    summarizing networks with semantic annotations. *F1000Research* **5:** 1717

622    Kurcinski M, Jamroz M, Blaszczyk M, Kolinski A & Kmiecik S (2015) CABS-dock web server for the

623    flexible docking of peptides to proteins without prior knowledge of the binding site. *Nucleic Acids*

624    *Res.* **43:** W419–W424

625    Ladoukakis E, Pereira V, Magny EG, Eyre-Walker A & Couso JP (2011) Hundreds of putatively

626    functional small open reading frames in Drosophila. *Genome Biol.* **12:** R118

627    Laskowski RA, Luscombe NM, Swindells MB & Thornton JM (1996) Protein clefts in molecular

628    recognition and function. *Protein Sci.* **5:** 2438–52

629    Lease K a & Walker JC (2006) The Arabidopsis unannotated secreted peptide database, a resource

630        for plant peptidomics. *Plant Physiol.* **142:** 831–838

631    Lee H, Heo L, Lee MS & Seok C (2015) GalaxyPepDock: A protein-peptide docking tool based on

632        interaction similarity and energy optimization. *Nucleic Acids Res.* **43:** W431–W435

633    Lindahl, E., Hess, B. & van der Spoel, D. GROMACS 3.0: a package for molecular simulation and

634    trajectory analysis. *J. Mol. Model.* **7,** 306–317 (2001).

635    London N, Movshovitz-attias D & Schueler-furman O (2010) The Structural Basis of Peptide-Protein

636        Binding Strategies. *Struct. Des.* **18:** 188–199

637    Maere S, Heymans K & Kuiper M (2005) BiNGO: A Cytoscape plugin to assess overrepresentation of

638        Gene Ontology categories in Biological Networks. *Bioinformatics* **21:** 3448–3449

639    Marín-González E & Suárez-López P (2012) 'And yet it moves': Cell-to-cell and long-distance

640        signaling by plant microRNAs. *Plant Sci.* **196:** 18–30

641    Martinez CR & Iverson BL (2012) Rethinking the term 'pi-stacking'. *Chem. Sci.* **3:** 2191

642    Matsubayashi Y (2011) Post-translational modifications in secreted peptide hormones in plants. *Plant

643        Cell Physiol.* **52:** 5–13

644    McDonald IK & Thornton JM (1994) Satisfying Hydrogen Bonding Potential in Proteins. *J. Mol. Biol.*

645        **238:** 777–793

646    Merico D, Isserlin R, Stueker O, Emili A & Bader GD (2010) Enrichment map: A network-based

647        method for gene-set enrichment visualization and interpretation. *PLoS One* **5:**

648    Miyamoto, S. & Kollman, P. A. Settle: An analytical version of the SHAKE and RATTLE algorithm for

649    rigid water models. *J. Comput. Chem.* **13,** 952–962 (1992).

650    Morris JH, Apeltsin L, Newman AM, Baumbach J, Wittkop T, Su G, Bader GD & Ferrin TE (2011)

651        ClusterMaker: A multi-algorithm clustering plugin for Cytoscape. *BMC Bioinformatics* **12:** 1–14

652    Murphy E, Smith S & De Smet I (2012) Small Signaling Peptides in Arabidopsis Development: How

653        Cells Communicate Over a Short Distance. *Plant Cell* **24:** 3198–3217

654    Neduva V & Russell RB (2005) Linear motifs: Evolutionary interaction switches. *FEBS Lett.* **579:**

655        3342–3345

656    Pace CN, Fu H, Fryar KL, Landua J, Trevino SR, Schell D, Thurlkill RL, Imura S, Scholtz JM, Gajiwala

657        K, Sevcik J, Urbanikova L, Myers JK, Takano K, Hebert EJ, Shirley BA & Grimsley GR (2014)

658    Contribution of hydrogen bonds to protein stability. *Protein Sci.* **23:** 652–661

659    Pace CN, Fu H, Fryar KL, Landua J, Trevino SR, Shirley BA, Hendricks MM, Iimura S, Gajiwala K,

660    Scholtz JM & Grimsley GR (2011) Contribution of hydrophobic interactions to protein stability. *J.*

661    *Mol. Biol.* **408:** 514–528

662    Pawson T & Nash P (2003) Assembly of Cell Regulatory Systems Through Protein Interaction

663    Domains. *Science (80-. ).* **300:** 445 LP-452

664    Perkins JR, Diboun I, Dessailly BH, Lees JG & Orengo C (2010) Transient Protein-Protein

665    Interactions: Structural, Functional, and Network Properties. *Structure* **18:** 1233–1243

666    Petsalaki E, Stark A, García-Urdiales E & Russell RB (2009) Accurate prediction of peptide binding

667    sites on protein surfaces. *PLoS Comput. Biol.* **5:**

668    Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC & Ferrin TE (2004)

669    UCSF Chimera - A visualization system for exploratory research and analysis. *J. Comput. Chem.*

670    **25:** 1605–1612

671    Porter CT (2004) The Catalytic Site Atlas: a resource of catalytic sites and residues identified in

672    enzymes using structural data. *Nucleic Acids Res.* **32:** 129D–133

673    Raveh B, London N & Schueler-Furman O (2010) Sub-angstrom modeling of complexes between

674    flexible peptides and globular proteins. *Proteins Struct. Funct. Bioinforma.* **78:** 2029–2040

675    Raveh B, London N, Zimmerman L & Schueler-Furman O (2011) Rosetta FlexPepDock ab-initio:

676    Simultaneous Folding, Docking and Refinement of Peptides onto Their Receptors. *PLoS One* **6:**

677    e18934

678    Rodrigues JPGLM, Trellet M, Schmitz C, Kastritis P, Karaca E, Melquiond ASJ & Bonvin AMJJ (2012)

679    Clustering biomolecular complexes by residue contacts similarity. *Proteins Struct. Funct.*

680    *Bioinforma.* **80:** 1810–1817

681    Rojas CM, Senthil-Kumar M, Tzin V & Mysore KS (2014) Regulation of primary plant metabolism

682    during plant-pathogen interactions and its contribution to plant defense. *Front. Plant Sci.* **5:** 17

683    Ruiz-Orera J & Messeguer X (2014) Long non-coding RNAs as a source of new peptides. *arXiv Prepr.*

684    *arXiv …:* 1–40

685    Salentin S, Schreiber S, Haupt VJ, Adasme MF & Schroeder M (2015) PLIP: Fully automated protein-

686    ligand interaction profiler. *Nucleic Acids Res.* **43:** W443–W447

687    Schindler CEM, De Vries SJ & Zacharias M (2015) Fully Blind Peptide-Protein Docking with

688 pepATTRACT. *Structure* **23:** 1507–1515

689 Selby CP & Sancar A (2006) A cryptochrome/photolyase class of enzymes with single-stranded DNA-
690 specific photolyase activity. *Proc. Natl. Acad. Sci. U.S.A.* **103**: 17696-700

691 Shiu S-H & Bleecker  a B (2001) Receptor-like kinases from Arabidopsis form a monophyletic gene
692 family related to animal receptor kinases. *Proc. Natl. Acad. Sci.* **98:** 10763–10768

693 Swails, J. *et al.* ParmEd. Available at: https://github.com/ParmEd/ParmEd.

694 Stein A & Aloy P (2008) Contextual specificity in peptide-mediated protein interactions. *PLoS One* **3:**
695 1–10

696 Tavormina P, De Coninck B, Nikonorova N, De Smet I & Cammue BPA (2015) The Plant Peptidome:
697 An Expanding Repertoire of Structural Features and Biological Functions. *Plant Cell* **27:** 2095–
698 118

699 The PyMOL Molecular Graphics System, Version 1.8 Schrödinger, LLC.

700 Tien MZ, Sydykova DK, Meyer AG & Wilke CO (2013) PeptideBuilder: A simple Python library to
701 generate model peptides. *PeerJ* **1:** e80

702 Torii KU (2004) Leucine-Rich Repeat Receptor Kinases in Plants: Structure, Function, and Signal
703 Transduction Pathways. *Int. Rev. Cytol.* **234:** 1–46

704 Trabuco LG, Lise S, Petsalaki E & Russell RB (2012) PepSite: Prediction of peptide-binding sites from
705 protein surfaces. *Nucleic Acids Res.* **40:** 423–427

706 Trellet M, Melquiond ASJ & Bonvin AMJJ (2013) A Unified Conformational Selection and Induced Fit
707 Approach to Protein-Peptide Docking. *PLoS One* **8:**

708 Van Der Spoel, D. *et al.* GROMACS: fast, flexible, and free. *J. Comput. Chem.* 26, 1701–18 (2005).

709 Velankar S, Dana JM, Jacobsen J, Van Ginkel G, Gane PJ, Luo J, Oldfield TJ, O'Donovan C, Martin
710 MJ & Kleywegt GJ (2013) SIFTS: Structure Integration with Function, Taxonomy and Sequences
711 resource. *Nucleic Acids Res.* **41:** 483–489

712 Wabik J, Kurcinski M & Kolinski A (2015) Coarse-Grained Modeling of Peptide Docking Associated
713 with Large Conformation Transitions of the Binding Protein: Troponin I Fragment–Troponin C
714 System. *Molecules* **20:** 10763–10780

715

716

717

718 **Supporting Information Captions**

719

720 **Supplementary Figure 1 (A)** Effect of random shuffling on the binding of peptides to pockets **(B)**
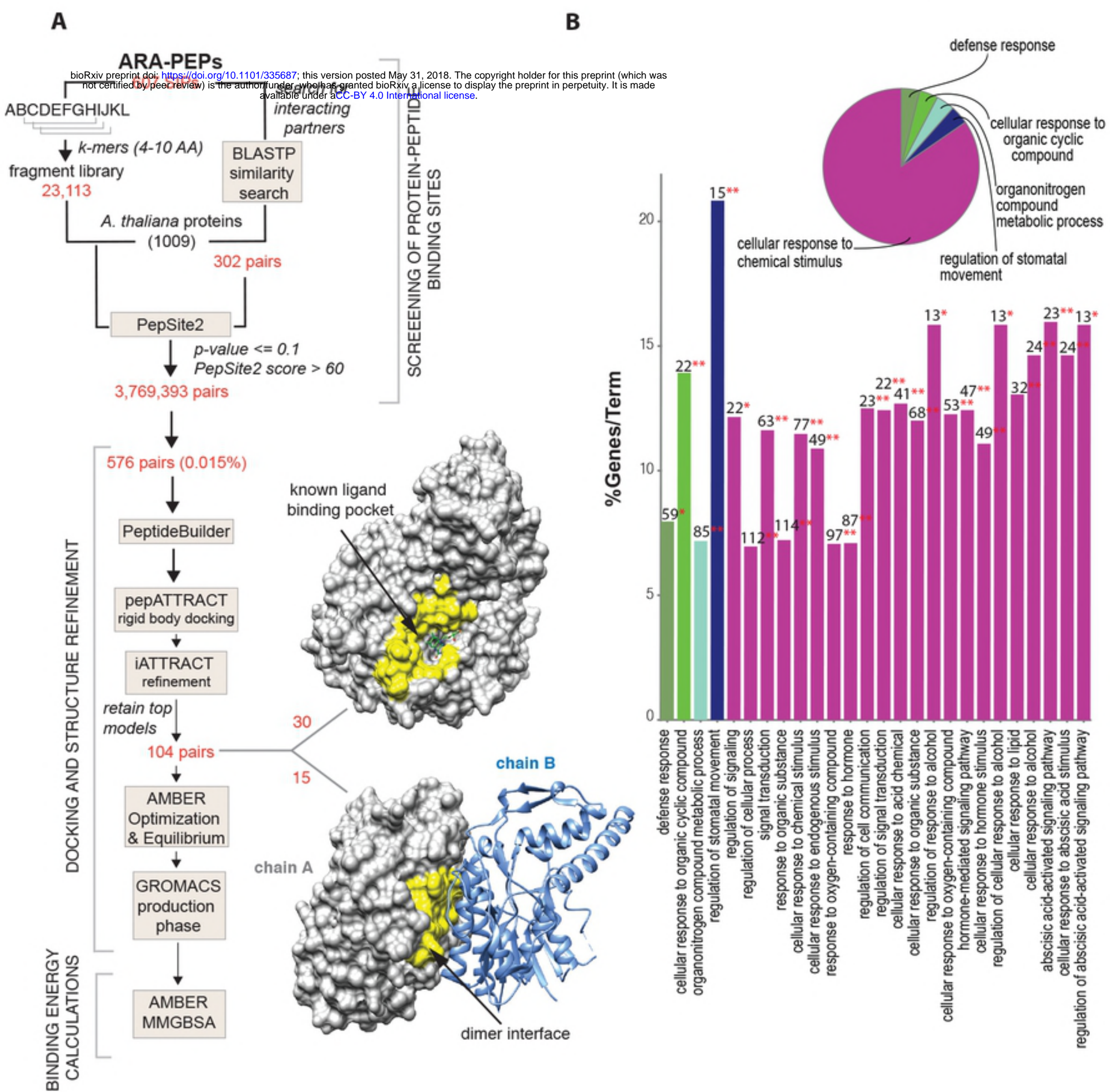
721 Comparison of pepATTRACT-local and blind docking protocols

722

723 **Supplementary Figure 2 (A)** Average contributions to the binding energy for each amino acid type,

724 for peptide and protein amino acids separately, and **(B)** for peptide amino acids at different locations

725 within the peptides. **(C)** Individual data points for all amino acids, from which the averages were made,

726 with red lines representing the average values. The represented data includes only amino acids whose

727 binding contribution is at least 40 % of the maximal contribution value within the respective system.

728

729 **Supplementary Figure 3**.: Examples of protein-peptide models showing binding modes of

730 characterized peptides and SIPs. The peptide binding pocket is highlighted in yellow.
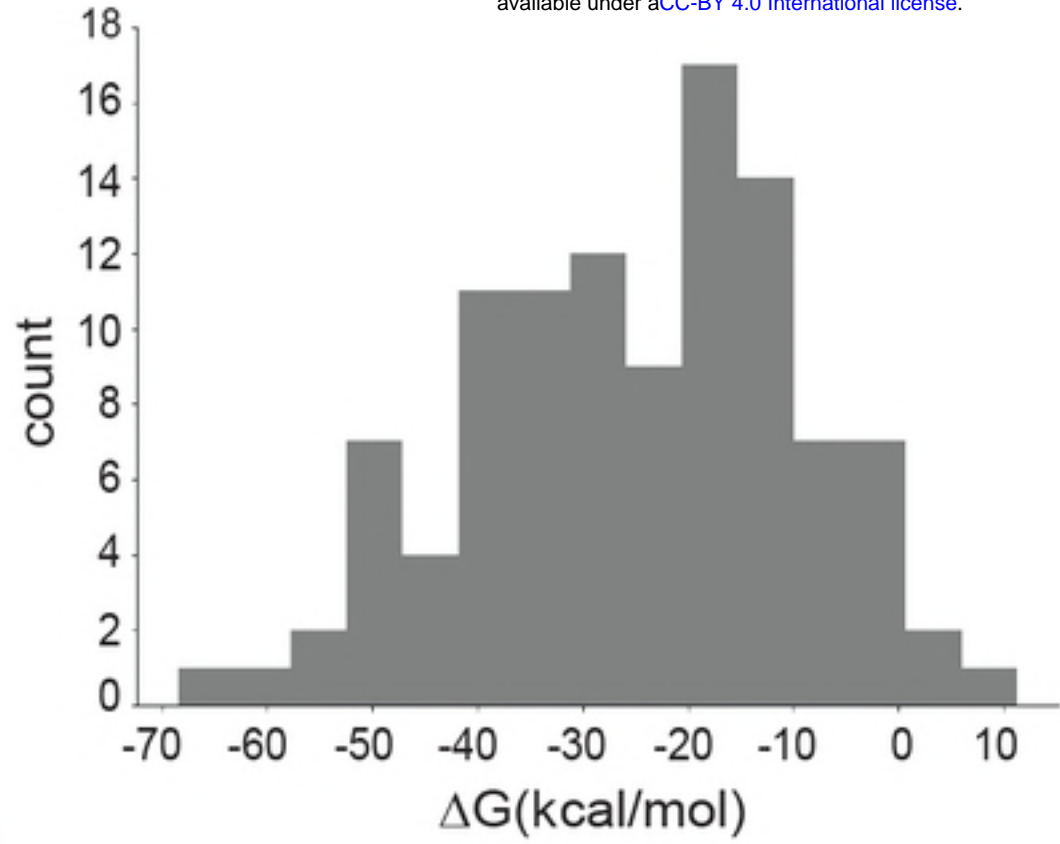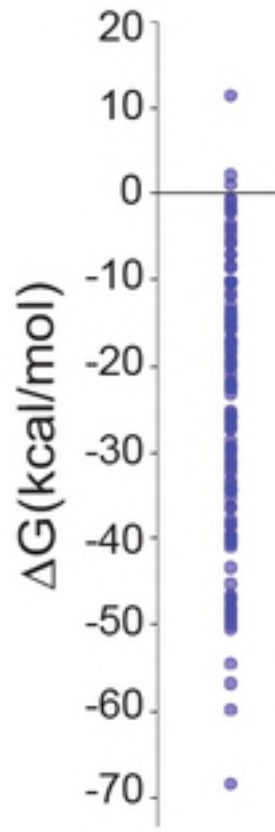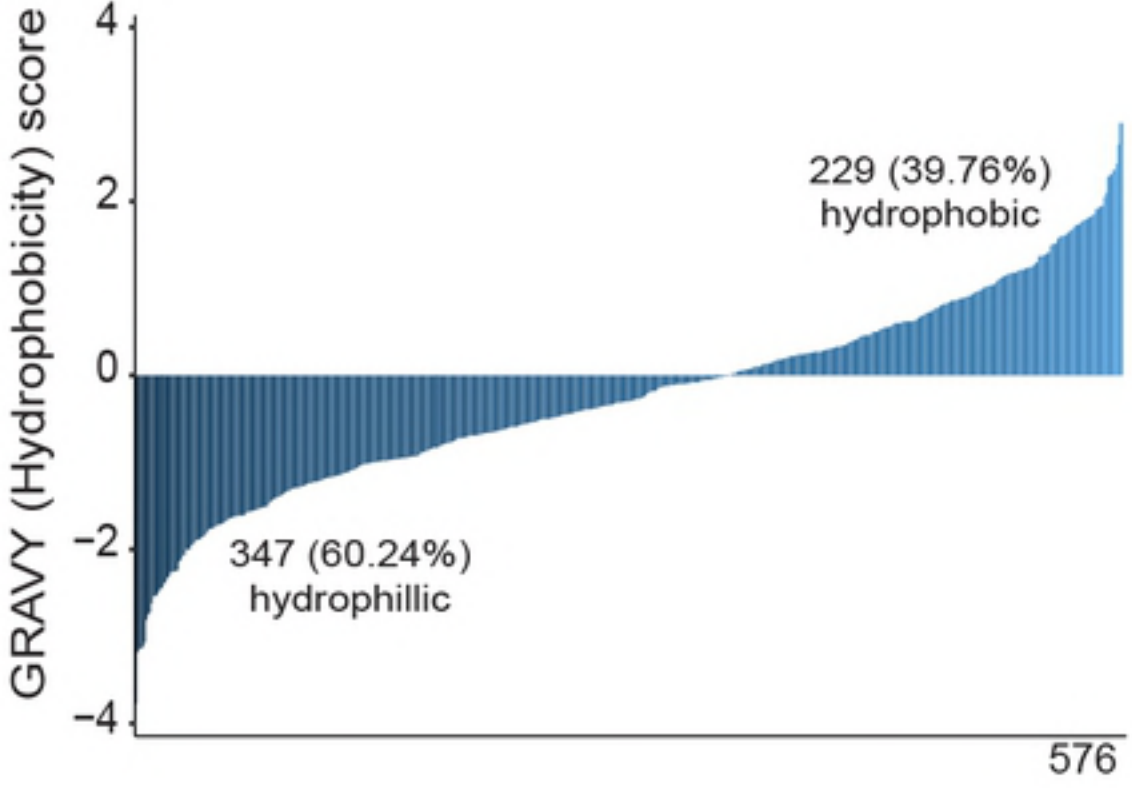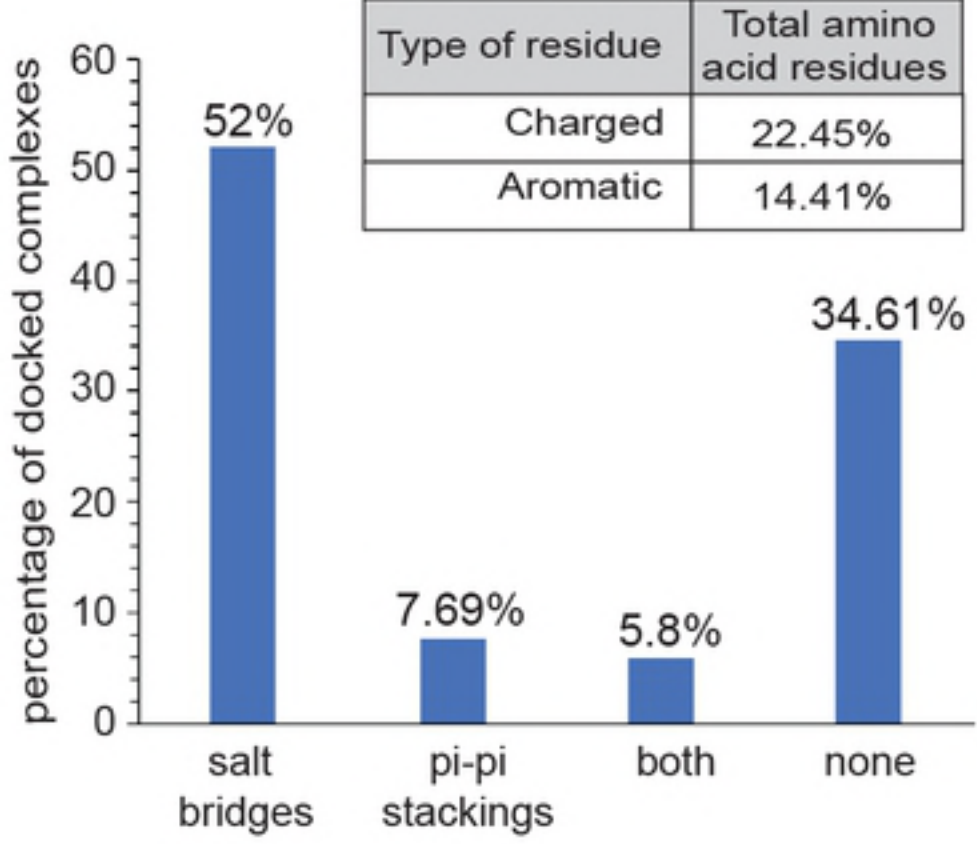
731

732 **Supplementary table 1**: List of protein-peptide models where the peptide-binding pocket overlaps

733 with known ligand binding or catalytic sites. Empty fields in the ligand column indicate only catalytic

734 sites and no known ligand is known to bind at the respective sites.

735

736

737 **Supplementary table 2**: List of protein-peptide models where the peptide binding pocket lies at the

738 subunits interface of a protein complex. For fields that are indicated as monomers in Protein

739 stoichiometry, the other chain in the structure is either a characterized peptide or the monomers may
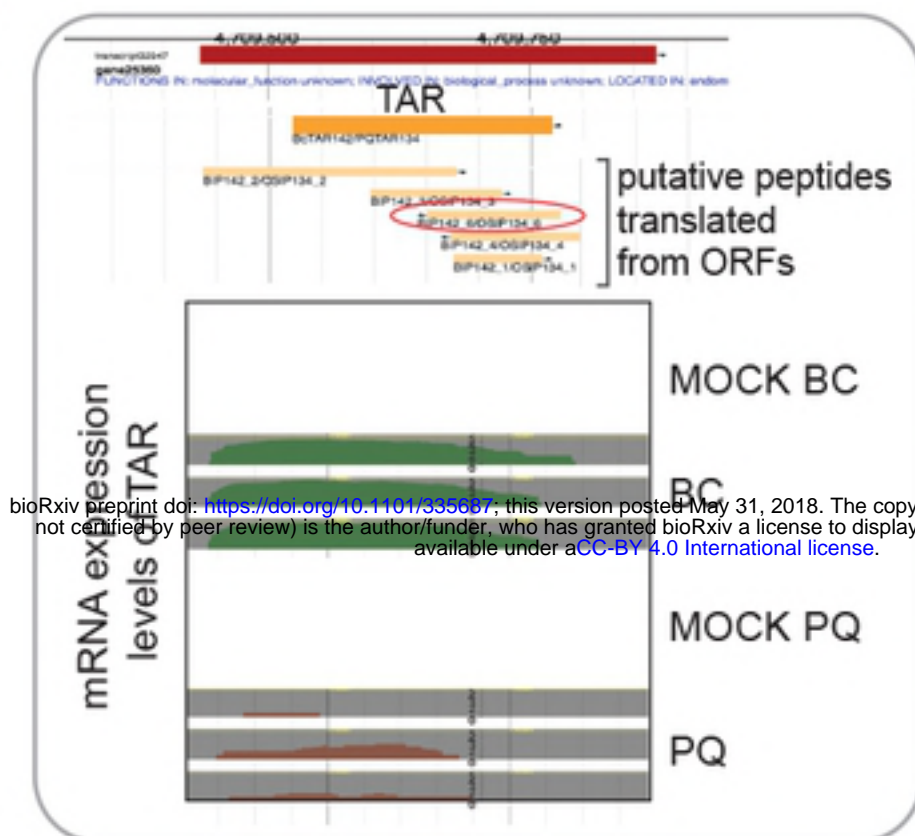
740 biologically aggregate to form dimers.

741

742

**A**



**B**



**C**



347 (60.24%) hydrophillic

229 (39.76%) hydrophobic

576

**D**



| Type of residue | Total amino acid residues |
| --- | --- |
| Charged | 22.45% |
| Aromatic | 14.41% |

**E**

| | side chain atoms | backbone atoms |
| --- | --- | --- |
| Total number in peptides | 13181 | 4080 |
| Unique involved in close contacts | 3142 | 1579 |
| Total involved in close contacts | 11305 | 4436 |

| Type of contacts | Mean±SD |
| --- | --- |
| Close contacts (4 Å) | 167±49 |
| Hydrophobic | 24±8 |
| H-bonds | 4.5±2 |

**A**

TAR

putative peptides
translated
from ORFs

mRNA expression
levels of TAR

MOCK BC

BC

MOCK PQ

PQ

**B**

## screening of protein-peptide pairs

PDB structures

peptide fragments

*p-value*

0.0016

0.0004

**C**

LAEDTFGEIS -2VTB(D)

coarse
model

LEU
ALA
ASP
GLU
THR
PHE

3D model in explicit water

molecular dynamics

protein-peptide binding analysis

Arg 490

Arg 436

Arg 487

Asp 4

Leu 1

Ser 10

Thr 5

Trp 365

Phe 6