

Mechanistic insights into the evolution of DUF26-containing proteins in land plants

Aleksia Vaattovaara¹, Benjamin Brandt^{2#}, Sitaram Rajaraman^{1#}, Omid Safronov¹, Andres Veidenberg³, Markéta Luklová^{1‡}, Jaakko Kangasjärvi¹, Ari Löytynoja³, Michael Hothorn², Jarkko Salojärvi^{4,1*}, Michael Wrzaczek^{1*}

¹Organismal and Evolutionary Biology Research Programme, Viikki Plant Science Centre, VIPS, Faculty of Biological and Environmental Sciences, University of Helsinki, Viikinkaari 1 (POB65), FI-00014 Helsinki, Finland

²Structural Plant Biology Laboratory, Department of Botany and Plant Biology, University of Geneva, Geneva, Switzerland.

³Institute of Biotechnology, University of Helsinki, Viikinkaari 1 (POB65), FI-00014 Helsinki, Finland

⁴School of Biological Sciences, Nanyang Technological University, 60 Nanyang Drive, Singapore 637551, Singapore

[‡]Current address: Laboratory of Plant Molecular Biology, Institute of Biophysics AS CR, v.v.i. and CEITEC - Central European Institute of Technology, Mendel University in Brno, Zemědělská 1, CZ-613 00 Brno, Czech Republic

[#]These authors contributed equally to this manuscript.

ORCID IDs: 0000-0003-3452-0947 (AV), 0000-0001-5867-8760 (BB), 0000-0001-5171-3578 (SR), 0000-0003-2044-3258 (OS), 000-0001-7451-0626 (AVe), 0000-0003-4948-423X (ML), 000-0001-5389-6611 (AL), 0000-0002-8959-1809 (JK), 0000-0002-3597-5698 (MH), 0000-0002-4096-6278 (JS), 0000-0002-5946-9060 (MW)

* To whom correspondence should be addressed:

Michael Wrzaczek
Organismal and Evolutionary Biology Research Programme
Viikki Plant Science Centre, VIPS
Faculty of Biological and Environmental Sciences
Viikinkaari 1, PO Box 65
FIN-00014 Helsinki University
Finland
Email: michael.wrzaczek@helsinki.fi
Phone: +358 2941 57 773

Jarkko Salojärvi
School of Biological Sciences
Nanyang Technological University
60 Nanyang Drive, SBS 02s-88f
Singapore 637551
Singapore
Email: jarkko@ntu.edu.sg
Phone: +65 69047231

46 **Abstract**

47 Large protein families are a prominent feature of plant genomes and their size variation is a key element for
48 adaptation in plants. Here we infer the evolutionary history of a representative protein family, the DOMAIN
49 OF UNKNOWN FUNCTION (DUF) 26-containing proteins. The DUF26 first appeared in secreted proteins.
50 Domain duplications and rearrangements led to the emergence of CYSTEINE-RICH RECEPTOR-LIKE
51 PROTEIN KINASES (CRKs) and PLASMODESMATA-LOCALIZED PROTEINS (PDLPs). While the
52 DUF26 itself is specific to land plants, structural analyses of Arabidopsis PDLP5 and PDLP8 ectodomains
53 revealed strong similarity to fungal lectins. Therefore, we propose that DUF26-containing proteins constitute
54 a novel group of plant carbohydrate-binding proteins. Following their appearance, CRKs expanded both
55 through tandem duplications and preferential retention of duplicates in whole genome duplication events,
56 whereas PDLPs evolved according to the dosage balance hypothesis. Based on our findings, we suggest that
57 the main mechanism of expansion in new gene families is small-scale duplication, whereas genome
58 fractionation and genetic drift after whole genome multiplications drive families towards dosage balance.

59

60 **Keywords**

61 Receptor-like protein kinase, tandem repeat, cysteine-rich receptor-like protein kinase, plasmodesmata-
62 localized protein, lectin structure

63 Gene duplication and loss events constitute the main factor of gene family evolution¹. Duplications occur by
64 two major processes, whole genome multiplications (WGM) and small-scale duplications (SSD), including
65 tandem, segmental, and transposon-mediated duplications². There appears to be two distinct modes of
66 expansion, since the gene families that evolve through WGMs rarely experience SSD events³. The division is
67 visible also on the functional level, since genes duplicated in WGMs are enriched for transcriptional and
68 developmental regulation as well as signal transduction functions, whereas SSDs occur preferentially on
69 secondary metabolism and environmental response genes³. The prevailing explanation for the phenomenon is
70 dosage balance; in complex regulatory networks and protein complexes the stoichiometric balance between
71 the different components needs to be preserved, and therefore selection acts against losses after WGMs and
72 against duplications in SSDs⁴. In terms of sizes, the families retained after WGMs are stable across different
73 species whereas highly variable families evolve through SSDs⁵, suggesting high turnover rates. However,
74 these results have been obtained by analyzing the two extremes, the top families displaying pure WGM
75 retention or SSD characteristics³ while most of the gene families likely evolve in an intermediate manner.

76 Plants and other eukaryotes have developed a wide range of signal transduction mechanisms for controlling
77 cellular functions and to coordinate responses on cell, tissue, organ and organismal level. Plants in particular
78 encode large gene families of secreted proteins⁶⁻⁸ and proteins with extracellular domains to respond to
79 environmental and developmental cues but in most cases their functions are not known⁴. Signaling proteins
80 with extracellular domains include receptor-like protein kinases (RLKs)^{9,10} and receptor-like proteins
81 (RLPs)¹¹. In RLKs, extracellular domains are involved in signal perception and protein-protein interactions¹²
82 while the intracellular kinase domain transduces signals to intracellular substrate proteins. The RLKs are
83 involved in essential mechanisms including stress responses, hormone signaling, cell wall monitoring and
84 plant development¹². The large number of secreted proteins, RLKs and RLPs in plants may reflect their
85 sessile lifestyle and need for meticulous monitoring of signals from other cells, tissues, or the environment.
86 However, the large numbers make it difficult to dissect their conserved or specialized functions, and
87 therefore a detailed understanding of their evolution in different plant lineages is needed. Phylogenetic
88 relations between different groups of RLKs and RLPs have been described^{9,13-16} but only few have been
89 physiologically and biochemically characterized¹⁷.

90 Here we carry out an in-depth analysis of one protein family involved in signaling to explore the dynamics
91 and effect of the different duplication mechanisms on overall gene family evolution: the Domain of
92 Unknown Function 26 (DUF26; Gnk2 or stress-antifungal domain)-containing proteins^{18,19}. The DUF26 is an
93 extracellular domain harboring a conserved cysteine motif (C-8X-C-2X-C) in its core. It is present in three
94 types of plant proteins. The first class is CYSTEINE-RICH RECEPTOR-LIKE SECRETED PROTEINs
95 (CRRSPs). CRRSPs form large subgroups in *Arabidopsis thaliana* and rice (*Oryza sativa*) but in most plants
96 the size of the family has not been quantified. The best characterized CRRSP is Gnk2, a protein from *Ginkgo*
97 *biloba* with single DUF26 which exhibits antifungal activity and acts as mannose-binding lectin *in vitro*^{18,19}.
98 Two maize CRRSPs have been shown to also bind mannose and participate in defence against a fungal
99 pathogen²⁰. The second class, CYSTEINE-RICH RECEPTOR-LIKE PROTEIN KINASES (CRKs), has a
100 typical configuration of two DUF26 in the extracellular region and forms a large subgroup of RLKs in plants
101 with 44 members encoded in the *Arabidopsis thaliana* genome. CRKs participate in the control of stress
102 responses and development in *Arabidopsis* and in rice²¹⁻³¹. The third class of DUF26 domain-containing
103 proteins is the PLASMODESMATA-LOCALIZED PROTEINs (PDLPs). PDLPs contain two DUF26
104 domains in their extracellular region and a transmembrane helix, but lack a kinase domain. They associate
105 with plasmodesmata and regulate symplastic intercellular signaling³², are involved in pathogen responses³³,
106 systemic signaling³⁴, control of callose deposition³⁵ and are targets for viral movement proteins³⁶. However,
107 the precise biochemical functions of DUF26-containing proteins in plants remain unclear.

108 Tandem expansions are a driving force for diversification processes for example for F-Box proteins³⁷,
109 transcription factors³⁸, as well as RLKs¹⁶ and RLPs¹¹. These diversification processes include sub-
110 functionalization, where paralogs retain a subset of their original ancestral functions, and neo-
111 functionalization, where a protein acquires novel functions after duplication³⁸. CRKs and CRRSPs typically
112 exist in clusters on plant chromosomes²⁴, suggesting relatively recent tandem expansions. This makes the
113 DUF26-containing proteins perfect dataset for testing the power of sequence-based evolutionary
114 investigations. We propose that CRKs and CRRSPs experienced both ancestral and recent, lineage-specific
115 tandem duplications in different angiosperm lineages. In contrast to the general pattern of gene families
116 expanding by small-scale duplication events, these gene families experienced significant expansion also

117 during of after ancient whole genome duplication events. We combine phylogenetic analyses with
118 experimental structural biology to gain insight into the evolution of DUF26-containing proteins in plants.
119 While sequence analysis indicates that the DUF26 domain is specific to land plants, the domain shows strong
120 structural similarity to fungal carbohydrate-binding lectins. Our structural analyses suggest that DUF26-
121 containing proteins constitute a novel group of carbohydrate-binding proteins in plants. Consequently,
122 sequence similarity alone is not sufficient evidence of orthologs, and lineage-specific protein family
123 expansions can make translation of functional data between species difficult. Our results illustrate that a
124 detailed understanding of the evolution of large protein families is a prerequisite for translating findings from
125 model plants to different species and for dissecting conserved or specialized functions of protein family
126 members.

127 **Results**

128 **Identification and annotation of DUF26 genes**

129 We selected 32 plant species representing major lineages of the plant kingdom for which high-quality
130 genome assemblies are available and retrieved 1656 DUF26-containing gene models (Figure 1a, Table S1).
131 Manual curation identified 322 gene models that required correction, demonstrating the necessity of manual
132 validation of datasets for analysis of gene families (Figure S1). To further reduce the possible biases in
133 annotation quality, we searched and identified 268 gene models *de novo* from genomic sequences (see
134 Materials and Methods). Partial gene models and pseudogenes were excluded resulting in 1409 high-quality
135 models included in subsequent analyses.

136 According to the PFAM protein domain database³⁹, DUF26 is specific to embryophytes. We confirmed this
137 by querying the genomes of the diatom *Phaeodactylum tricornutum*, five algae species, the charophyte
138 *Klebsormidium flaccidum*, as well as fungi, insects and vertebrates (see Materials and Methods) and
139 identified no DUF26 or DUF26-like domain among the species (Figure 1a).

140 **DUF26-containing proteins have diverse domain compositions**

141 DUF26-containing proteins are grouped to three categories: CRRSPs, PDLPs and CRKs (Figure 1b).
142 CRRSPs consist of a signal peptide (SP) followed by one or more DUF26 domains, separated by a short
143 variable region. CRRSPs with a single DUF26 (sdCRRSPs) were identified from most land plants, including
144 the early-diverging liverwort (*Marchantia polymorpha*) and moss (*Physcomitrella patens*) lineages (Figure
145 1). CRRSPs with two DUF26 domains (ddCRRSPs) were identified from vascular plants including the early-
146 diverging lycophyte *Selaginella moellendorffii*; they represent the predominant type in all vascular plant
147 genomes (Figure 1). Rice as well as *Brassicaceae* display lineage-specific evolution with a large number of
148 ddCRRSPs while sdCRRSPs are absent (Figure 1a and S2).

149 CRKs contain a SP, two DUF26 domains, and a transmembrane region (TMR) followed by an intracellular
150 protein kinase domain. Similar to ddCRRSPs, CRKs were identified from vascular plants but not from
151 bryophytes (Figure 1a). The CRKs likely emerged as the result of a fusion of sdCRRSPs with TMR and
152 kinase domain from LRR_clade_3 RLKs in the common ancestor of vascular plants¹⁵, since the *Selaginella*

153 genome uniquely encodes single DUF26 CRKs (sdCRKs; Figure 1b). The two-domain configuration is
154 stable, since only few CRKs from eudicot plants contain more than two DUF26 domains.

155 Finally, PDLPs are composed of a SP, two DUF26, and a transmembrane region (TMR) followed by a 10-15
156 amino acid (AA)-long cytoplasmic extension and they were identified from all seed plants. Within the
157 angiosperms, we also identified several CRKs lacking SP, extracellular region and transmembrane domain.
158 These are subsequently referred to as CYSTEINE-RICH RECEPTOR-LIKE CYTOPLASMIC KINASES
159 (CRCKs).

160 **Evolution of CRKs, PDLPs and ddCRRSPs from small sdCRRSPs**

161 To investigate the relationships between CRRSPs, CRKs and PDLPs, we estimated phylogenetic trees using
162 full length AA sequences translated from gene models with intact DUF26 domains (Figure 2a and b). As a
163 result of the different domain compositions only the DUF26-containing region aligned across all sequences.
164 Due to their high sequence divergence CRCKs, DUF26-containing gene models from bryophytes and
165 monocot CRKs with a different intracellular protein kinase domain were excluded from the alignment (see
166 below). Overall, a phylogenetic tree for DUF26-containing proteins based on a filtered amino acid sequence
167 alignment split into two distinct groups, a basal group α and a variable group β (Figures 2a and b), where α is
168 paraphyletic with respect to β . In order to increase the number of informative sites and thus obtain better
169 resolution, we estimated separate phylogenetic trees for both groups (See Methods; Figures S2a and b); the
170 subgrouping observed within the basal α - and variable β -groups was present there as well as in the trees
171 estimated for each sub-family of DUF26-containing proteins (Figures S2c-e). To study gene family evolution
172 we reconciled the gene trees with the species tree, and estimated ancestral gene contents and duplication and
173 loss events for the sub-families in eleven species (see Materials and Methods; Figure S3). To identify
174 significant expansions we fitted birth-death rate models for DUF26-containing protein families and
175 compared the rates against different computationally derived gene families (orthogroups) for RLKs, all
176 protein kinases, and plasmodesmal proteins⁴⁰ using Badirate⁴¹ (see Materials and Methods). Finally, we
177 assessed selective pressure by estimating amino acid conservation patterns around the main cysteine-motif of
178 the DUF26 domains for major subclades within the α - and β -groups (Figures 2c). While most conserved

179 positions within DUF26-A and -B are either conserved in all DUF26-containing proteins or specific to
180 individual types, we were able to identify conserved sites specific to the α - or β -clades (Figures 2c).

181 The α -group is likely older, containing sequences from all vascular plants. Proteins in this group are
182 conserved in sequence level and identification of putative orthologs from different species is frequently
183 possible. Purifying selection, i.e. stabilizing selection by selective removal of (deleterious) variations, is
184 likely the main force acting on this clade, as suggested by low d_N/d_S values (one-rate model for whole
185 groups: bCRK-I 0.184, bCRK-II 0.192, PDLPs 0.267, sdCRRSPs 0.162, CRCK 0.134; more flexible model
186 with branch-specific d_N/d_S within each group yielded similar results). The subgroups within the basal α -group
187 have evolved independently but their DUF26 domains share a number of features which distinguish them
188 from the members of the variable β -group. These distinguishing features include a leucine or isoleucine
189 residue in the fourth position after the first cysteine in the DUF26-A and the position of the fourth cysteine in
190 the DUF26-B (Figure 2c). The sdCRRSPs appear to be the most ancient type of DUF26 genes in land plants,
191 since the sdCRRSPs are located close to the root of the α -group (Figure 2b) and form a monophyletic
192 subclade at the root of the CRRSP tree (Figure S2c). Furthermore, the sdCRRSPs are present in various early
193 diverging plant lineages such as the gymnosperm *Ginkgo biloba* (including Gnk2, the best studied
194 sdCRRSP^{18,19}) and the liverwort, *Marchantia polymorpha* (Figure S2f). The turnover rates of sdCRRSPs do
195 not differ from those of all gene families and show lineage-specific expansions in early diverging species
196 (Figure S3a).

197 The placement of *Selaginella* sdCRKs to the root of the CRK phylogeny (Figure S2d) and as sister to
198 sdCRRSPs in the α -group (Figure 2b) suggests an ancient origin. Our analyses indicate that the DUF26
199 domain likely has duplicated after fusing with TMR and kinase domain, thus establishing the typical double-
200 DUF26 CRK configuration found in seed plants (Figure S2d). Following the duplication, the two DUF26
201 domains diverged into distinct forms, DUF26-A and DUF26-B, which are evolutionarily conserved (Figure
202 2d). Overall, CRKs have expanded significantly in the branches leading to lycophytes and to angiosperms
203 compared to all RLKs (Figure 3a), and compared to all protein kinases they expanded significantly in the
204 branch from the ancestral node of lycophyte *Selaginella* to angiosperms (Figure S4a). In all of these branches
205 plants have experienced several WGDs⁴², suggesting that the ancestral CRKs have either been preferentially

206 retained after WGMs., or they have had a tandem birth rate that is higher than the death rate following
207 WGMs.

208 A monophyletic group of CRKs with representatives from gymnosperms and angiosperms is located near the
209 base of the CRK phylogeny (Figure S2d) and belongs to the α -group (Figure 2a and b). This group likely
210 represents the ancient CRKs in seed plants and will be subsequently referred to as basal CRKs (bCRKs).
211 Following the initial innovation in ancestral vascular plants the group has evolved at rates similar to
212 comparable orthogroups containing all protein kinases or all RLKs (Figure 3b, S3b and S4b). The bCRKs
213 split into two distinct subgroups, bCRK-I and bCRK-II (Figure 2b and S5), both of which are present in
214 gymnosperms and angiosperms, suggesting diverging duplicates in early seed plants. The larger bCRK-I
215 subclade further divides into distinct branches with tandemly duplicated *Amborella* bCRKs at their roots
216 (Figure S5, S6a-b) suggesting rapid differentiation after tandem duplication in ancestral angiosperms. The
217 lineage-specific size of the bCRK-I branch is conserved, except for an expansion specific to *Solanaceae*. The
218 small bCRK-II subclade is interestingly absent from *Brassicaceae*.

219 PDLPs are found in seed plants but not bryophytes or lycophytes. PDLPs belong to the α -group (Figure 2a
220 and b) and represent the most conserved class of DUF26-containing genes. As such, they do not display
221 different expansion rates compared to plasmodesmata-related orthogroups⁴⁰ (Figure 3c). PDLPs split into
222 two branches, PDLP-I and PDLP-II (Figure S2e), which both contain eudicot and monocot PDLPs,
223 suggesting that the divergence occurred already in common ancestral angiosperms. The PDLP-II branch
224 further divides into two angiosperm-specific branches with *Amborella trichopoda* sequences at their roots,
225 whereas the PDLP-I branch can be traced back to a single *Amborella trichopoda* PDLP. PDLPs and
226 ddCRRSPs originate from the loss of kinase domains and/or TMRs from CRKs. This two-step process is
227 supported by an atypical PDLP from *Amborella trichopoda* which is located at the root of the main
228 ddCRRSP clade (Figure S7). The timing of the event cannot be inferred, since it is unclear whether
229 gymnosperms also contain members of the PDLP-I branch. However, a database search identified one partial
230 gene model, a candidate PDLP from the fern *Marsilea quadrifolia*⁴³ lacking a transmembrane region (see
231 Materials and Methods). This putative fern PDLP shows high similarity to PDLPs and places to the root of
232 PDLPs in a phylogenetic tree estimated from PDLPs and CRKs (Figure S8).

233 A group of spruce-specific CRKs (spruce vCRKs) belongs to the α -group (Figure 2a and b) and are more
234 related to PDLPs than other CRKs. They form a distinct group between bCRKs and a large group of
235 angiosperm CRKs, the “variable CRK clade” (vCRKs; Figure 2a and b, S2d, S3d). These angiosperm
236 vCRKs form the β -group together with ddCRRSPs and atypical monocot sdCRRSPs (Figure 2a and b).
237 These CRRSPs likely evolved from vCRKs through the loss of TMR and kinase domains and, in case of
238 sdCRRSPs, also of the DUF26-B domains. The β -group is less conserved compared to the more ancient α -
239 group and branches into two eudicot-specific groups and one monocot-specific group with a small group of
240 *Amborella trichopoda* vCRKs at the root of the clade. Still, there are some conserved positions surrounding
241 the main cysteine motif that distinguish members of the β - from the α -group, for example a conserved
242 threonine following the first cysteine in DUF26-B (Figure 2c). Unlike proteins in the α -group, CRRSPs and
243 vCRKs in the β -group have undergone several independent tandem expansions in different plant taxa (Figure
244 3d, 3e, S3d, S4c, S6c) and expanded significantly during the diversification of monocots and dicots. CRRSPs
245 in the β -group are not monophyletic, suggesting several independent birth events resulting from partial
246 duplications of vCRKs. Hence, expansion rates and extrapolation of ancestral gene counts for ddCRRSPs
247 could not be reliably predicted (Figure S3e). Lineage-specific expansions in the β -group will make
248 identification of orthologs challenging.

249 **Plant DUF26 domains form conserved tandem assemblies and are structurally related to fungal lectins**

250 The high sequence divergence of the DUF26 proteins in different plant lineages and the strong lineage-
251 specific expansions raise the question whether their overall structure is conserved and what elements
252 distinguish the more closely related members of this protein family. The consensus DUF26 (PF01657)
253 domain as defined in PFAM comprises ~90-110 amino-acids and contains the conserved cysteine motif C-
254 8X-C-2X-C. Structural information is currently available only for the sdCRRSP Gnk2¹⁹ but not for proteins
255 with a double DUF26 configuration, such as ddCRRSPs, CRKs and PDLPs. Mechanistic constraints restrict
256 the evolution of protein structures, and therefore understanding structural conservation can provide essential
257 clues for protein function. Furthermore, selection patterns may differ between a young and lineage-specific
258 gene and an evolutionarily conserved gene.

259 Thus, we defined the structural relationship of tandem DUF26 domains by determining crystal structures of
260 the AtPDLP5 (residues 26-241) and AtPDLP8 (21-253) ectodomains to 1.25 and 1.95 Å resolution,
261 respectively (Table S2). Individual DUF26 domains feature two small α -helices folding on top of a central
262 anti-parallel β -sheet (Figure 4a). The PDLP5 DUF26-A domain is found to be N-glycosylated at positions
263 Asn69 and Asn132 in our crystals (Figure 4a). The secondary structure elements of DUF26 are covalently
264 linked by three disulfide bridges, formed by six conserved Cys residues, part of which belong to the C-8X-C-
265 2X-C motif (Figures 2d, 4a). We have previously suggested that tandem DUF26-domain containing proteins
266 could be involved in ROS or redox sensing^{24,26}. To assess the functional roles of the invariant disulfide
267 bridges in PDLPs, we mutated the partially solvent exposed PDLP5^{Cys101}, PDLP5^{Cys148} and PDLP5^{Cys191} to
268 alanine. While the wild-type PDLP5 ectodomain behaves as a monomer in solution (Figure S9), the mutant
269 proteins tend to aggregate in our biochemical preparations (Figure S9) and display reduced structural
270 stability in thermofluor assays (Figure S10, see Materials and Methods). These experiments and our
271 crystallographic data (Figure 4a) together suggest that the conserved disulfide bonds in PDLPs and
272 potentially in other DUF26-domain containing proteins are involved in structural stabilization rather than
273 redox signaling.

274 The N-terminal DUF26-A (PDLP5 residues 30-132) and the C-terminal DUF26-B (residues 143-236)
275 domains are connected by a structured loop (residues 133-142) and make extensive contacts with each other
276 (Figure 4a). The resulting ectodomain has a claw-like shape with the β -sheets of DUF26 A and B facing each
277 other (Figure 4a). The DUF26-A and B domains in PDLP5 and 8 closely align, with root mean square
278 deviations (r.m.s.d.s) of 1.6 and 1.2 Å when comparing 89 corresponding C $_{\alpha}$ atoms, respectively (Figure
279 S11a). Overall, DUF26-A is considerably more variable than DUF26-B on the sequence level (Figure 2d).
280 The DUF26-A and B domains in PDLP5 and PDLP8 have 24 % and 30 % of their residues in common, most
281 of which map to the hydrophobic core of the domain (including the six cysteine residues forming intra-
282 molecular disulfide bonds) and to the DUF26-A – DUF26-B interface (Figure 4b). This interface is formed
283 by a line of aromatic and hydrophobic residues originating from the proximal face of the β -sheet in DUF26-
284 A and B (Figure 4b, Supplementary Figure 12). Importantly, many of the interface residues are strongly
285 conserved among different PDLPs, but also among CRKs and ddCRRSPs (Figure S12). Consistently, the

286 ectodomains of PDLP5 and PDLP8 belonging to different phylogenetic clades (Figure S8) closely align with
287 an r.m.s.d. of ~ 1.6 Å when comparing 198 corresponding C $_{\alpha}$ atoms (Figure S11b). Together, these
288 observations suggest that evolutionarily distant DUF26 tandem proteins likely share the conserved three-
289 dimensional structure.

290 The physiological ligands for PDLPs are currently unknown. For this means, we performed structural
291 homology searches⁴⁴ to obtain insights into the biochemical function of plant DUF26 domains (See Materials
292 and Methods). The top hits include the single DUF26 domain protein ginkbilobin-2 (Gnk2) from *Ginkgo*
293 *biloba*¹⁹. Despite their moderate sequence similarity, the overall fold of Gnk2 and PDLP5 DUF26-A and B
294 as well as their disulfide-bond arrangement is fully conserved (Figure 4c). Notably, Glu130 and Arg132
295 implicated in mannose binding in Gnk2 are replaced by Asp131 and Lys133 in the DUF-A of PDLP5,
296 respectively (Figure 4d). A similar pocket is found in the DUF-A domain of PDLP8, but not in the DUF-B
297 domains of either PDLP5 or 8. Despite these structural homologies of Gnk2, PDLP5 DUF-A and PDLP8
298 DUF-A, we could not detect binding of mannose to the isolated PDLP5 ectodomain *in vitro* (Figure S13a).
299 We also tested other water soluble cell wall derived carbohydrates, but were not able to detect any binding to
300 the PDLP5 ectodomain (Figure S13b). The PDLP5 DUF26 domains share significant structural homology
301 not only with the plant Gnk2, but also with two fungal lectins, the α -galactosyl-binding *Lyophyllum decastes*
302 lectin (LDL)⁴⁵ and a glycan-binding Y3 lectin from *Coprinus comatus*⁴⁶. Both proteins closely align with the
303 plant DUF26 domain, and share one of the three disulfide bridges (Figure 4e-f). The surface areas involved
304 in globotriose and glycan binding, respectively, are not conserved in PDLPs, but the structural similarity of
305 plant DUF26 domains with different eukaryotic lectins could suggest a common evolutionary origin and a
306 potential role as carbohydrate recognition modules⁴⁵.

307 We next explored potential binding sites in the two molecules by identifying regions under positive or
308 purifying selection that could be indicative of domains involved in protein-protein interactions or ligand
309 perception. Analysis of site-wise selection for orthologs of PDLP5 and PDLP8 in their structural context
310 yielded low ω values, indicating strong conservation of residues buried inside the DUF26 domain fold, while
311 more variable residues (under more relaxed selection) appear on the surface of the structure (Figure 5a). The
312 high variability of the surface of the PDLP5 and PDLP8 DUF26 domains may be central to their ability to

313 interact with other proteins but also with potential ligands (Figure S13c). In case of PDLP5, the higher ω
314 values on the surface could indicate fast evolution events leading to sub- or neofunctionalization, as the
315 PDLP5 orthologs all originate from the more recent duplication in the lineage leading to Brassicaceae
316 species. The drastically different surface charge properties of related PDLPs from the same species (Figure
317 5b) suggest that different PDLPs and other DUF26 domain-containing proteins sense a rather diverse set of
318 ligands. While the nature of these molecules is currently unknown, cell-wall derived carbohydrates or small
319 extracellular molecules represent candidate ligands, but we were not able to identify any in our experiments
320 (Figure S13a-b). Notably, we observed typical lectin-dimers in our PDLP5 and PDLP8 crystals, in which
321 two lectin domains dimerized along an extended anti-parallel β -sheet (Figure 5c)⁴⁷. In principle, this mode of
322 dimerization could give rise to an extended binding cleft for a carbohydrate polymer, and presents an
323 attractive receptor activation mechanism for PDLPs and CRKs, in which a monomeric ground state forms
324 ligand-induced oligomers, as previously seen with LysM-domain containing carbohydrate receptors in
325 plants⁴⁸.

326 **The CRK kinase domain is related to LRR and S-locus RLKs**

327 Kinase domains transduce signals by phosphorylating substrate proteins and thereby are determining factors
328 for signal specificity. Typically, the intracellular kinase domain has been used to investigate phylogenetic
329 relationships between RLKs^{9,15,16}. The typical CRK kinase domain is similar to the kinase domain of S-locus
330 lectin and LRR RLKs from LRR_clade_3¹⁵ (Table S3). Based on the sequence of catalytic motifs in kinase
331 domains⁴⁹ most CRKs seem to be active protein kinases and the *in vitro* activity of several CRKs has been
332 experimentally confirmed^{25,28,30}. Most CRKs belong to the RD type^{50,51} which is considered to be capable of
333 auto-activation but a few non-RD CRKs are present in plant genomes⁴⁹.

334 Analyzing ectodomains and kinase domains of CRKs separately suggests that *Selaginella* ddCRKs share an
335 ancestor with bCRKs, while *Selaginella* sdCRKs share an ancestor with vCRKs (Figure 6a). The clear
336 separation of DUF26-A and DUF26-B (Figure 2d) and the timing of those events does not reveal whether the
337 duplication of the DUF26 domain in the extracellular region of CRKs has happened more than once or
338 whether functional constraints in the kinase domain led to the conserved similarity of *Selaginella* sdCRKs
339 and vCRKs. Juxtaposition of phylogenetic trees based on ectodomains and kinase domain suggests several

340 exchanges of kinase or extracellular regions among CRKs during evolution (Figure 6a). Most strikingly, a
341 group of monocot-specific CRKs separates from other CRKs in a phylogenetic tree based on the kinase
342 domain (Figure 6a). Those CRKs have an atypical gene model comprising a kinase domain with high
343 similarity to concanavalin A-like lectin protein kinase domains (Table S3), and a different exon-intron
344 structure (Figure 6b, S1b), altogether suggestive of chimeric gene formation following a tandem
345 duplication⁵². The switch of the kinase domain and the associated changes in exon-intron structure is specific
346 to grasses (*Poaceae*) and has likely resulted in a different set of target substrates. Exchange of kinase
347 domains is not the only alteration of domain composition within DUF26-containing proteins. Loss of
348 ectodomains and TMRs has established CRCKs at least three times; one group of CRCKs is specific to
349 angiosperms (CRCK-I clade), one is specific to *Brassicaceae* and one only to *Arabidopsis thaliana*.

350 **Mixed-mode evolution of large gene families**

351 In order to carry out more detailed analyses of gene family dynamics we analyzed the synteny, conservation
352 of the gene order between species, as well as tandem duplications in highly contiguous chromosome-level
353 assemblies of *Amborella trichopoda*, tomato (*Solanum lycopersicum*), Arabidopsis, rice and maize (*Zea*
354 *mays*; Figures 7a and S7), and estimated the timing of the duplication events by reconciliation of gene trees
355 with species trees (Figure S14a).

356 Within the young, rapidly diverging β -group the vCRKs show large lineage-specific expansions. The
357 ancestral origins for monocot and eudicot vCRKs differ, and neither synteny nor orthology can be identified
358 (Figure 7c and S14c). Altogether this suggests that this younger subfamily has a high birthrate and that it
359 expands rapidly by tandem duplications in all species. Additionally, many of the tandems are lost or
360 fractionated after WGMs. Similarly the CRRSPs demonstrate little synteny between different species (Figure
361 7a and S14d), and CRRSPs in rice and Arabidopsis experienced lineage-specific tandem duplications (Figure
362 S14d). In *Brassicaceae* this expansion can be traced to *Amborella* CRRSP (*AtrCRRSP2*), altogether
363 suggesting a tandem mode of expansion.

364 Tandem duplications evolve through unequal crossover or homologous recombination events⁵³. Unequal
365 crossover produces copy number variation, whereas homologous recombination such as gene conversion

366 plays a role in concerted evolution, which can maintain the similarity between gene copies over long
367 periods⁵⁴. Gene conversion is known to depend on the genomic distance as well as sequence homology.
368 Accordingly, we observed several events among the lineage-specific tandem vCRK expansions (Table S4),
369 whereas in case of bCRKs, events were observed only in the tandem expansion in *Amborella*. This suggests
370 that gene conversion is an important process maintaining the similarity between recent tandem duplicates but
371 as the sequences diverge over time the conversion events become increasingly rare.

372 The CRCK-I genes are present in most genomes as single copy genes within conserved syntenic genome
373 segments, suggesting that duplicates from WGD events have been lost during genome fractionation (Figure
374 7a). The evolution follows a specific dosage balance model where the maintenance of a single copy is critical
375 to the organism.

376 A hallmark for gene families evolving under dosage balance is that their overall numbers should be
377 conserved among species with similar WGM history. In the species tree (Figure 7a), most of the branches
378 contained one or two WGMs. Despite these events, the overall number of bCRKs is well conserved in
379 angiosperms (Figure 3b, 7b S3b and S7b). However, in *Amborella trichopoda* five bCRK genes appear in
380 tandem and these genes are at the roots of the respective orthologs (Figure S6b), indicating an ancestral SSD
381 origin still present in *Amborella*. The duplicate region experienced considerable fractionation during
382 evolution leading to *Brassicaceae* and *Solanaceae* lineages, resulting in scattered bCRK-I orthologs with
383 little conserved synteny, whereas in the two grasses the tandem duplicate was lost altogether. This indicates
384 rapid pseudogenization of the duplicated tandem blocks after WGMs, with, except for *Solanaceae*, no recent
385 tandem expansions. Altogether this suggests that a gene family that initially existed as a tandem duplicate
386 may have shifted towards a dosage balance mode of evolution. Dosage balance is observed in the second
387 subfamily of ancient origin, PDLPs, since they appear in genomic regions where synteny is conserved within
388 eudicots and monocots (Figure 7a), and no recent or ancient SSD events can be detected.

389 One of the predictions for the gene families evolving under dosage balance is that retained duplicates should
390 exhibit less functional divergence than other duplicates³. We explored functional conservation by analyzing
391 publicly available gene expression data on stress treatments (Table S5; Figure 7d, S15). In agreement with

392 studies in Arabidopsis and rice^{24,26,31,55}, pathogen treatments have the biggest impact on transcript abundance
393 of DUF26-containing genes, in particular CRKs and CRRSPs (Figure S15). Our analysis of gene expression
394 data suggests extensive lineage-specific functional diversification. This is visible in the correlation rank
395 between putative orthologs; in many cases higher correlation can be found with DUF26-containing genes
396 that have less similarity in sequence, indicating that the closely related genes have undergone sub- or neo-
397 functionalization following duplications^{56,57}.

398 Overall bCRKs show elevated transcript levels in response to stress treatments, while many vCRKs have
399 elevated transcript levels in pathogen treated samples. Rice PDLPs display altered transcript levels in some
400 specific stress treatments. Despite re-arrangements and lineage-specific expansions the data provide support
401 for seven putative orthologs, including three PDLP and three CRRSP relationships (Figures 7d; Table S5).
402 Even though the synteny of bCRKs (and PDLPs) is more conserved compared to CRRSPs, bCRKs
403 demonstrate varying responses to stimuli, whereas in CRRSPs synteny is associated with similar functions.

404 The second prediction from the dosage balance model is that since the protein products of the genes are
405 highly connected and thus interact with many other proteins, disturbances in the dosage balance should have
406 large effects on an organism's phenotype⁵⁸. Reanalysis of phenotyping data of T-DNA mutant insertion
407 lines²⁴ confirms that bCRKs indeed demonstrate a larger variance in phenotypes than vCRKs ($p=0.03$;
408 Wilcox test; Figure 7e). Altogether the analysis suggests that PDLPs and bCRKs are evolving according to
409 the dosage balance model, whereas the vCRKs and CRRSPs evolve by SSD mechanisms.

410 **Discussion**

411 Compared to animal genomes, plant genomes encode a large number of large gene families⁵⁹. In particular
412 signal transduction components including transcription factors, protein kinases and phosphatases have
413 experienced drastic expansions in plants⁵⁹. This might reflect the adaptation to a sessile lifestyle but also
414 could indicate a different strategy for signal transduction and integration at the cellular level. The large, in
415 part lineage-specific expansions and conversions between different domain arrangements seriously hamper
416 the identification of orthologous proteins in different plant species. Here we studied the evolution of a large
417 plant protein family which is hallmarked by heterogenous domain architecture and drastic lineage-specific

418 expansions of subgroups, the DUF26-containing proteins. We identified 1409 high-quality gene models
419 representing CRRSPs, CRKs and PDLPs from major plant lineages. Our analyses suggest that sdCRRSPs are
420 the ancestral type of DUF26-containing proteins. CRKs originated from a fusion of CRRSPs with TMR and
421 kinase domain of LRR_clade_3 RLKs¹⁵ in the lineage leading to lycophytes. PDLPs and ddCRRSPs
422 emerged subsequently through the loss of the kinase domain or the TMR and kinase domain. Our results
423 reveal an ancient split into two distinct groups. The α -group is strongly conserved in size and sequence
424 throughout embryophytes. This facilitates identification of functional orthologs and extrapolation of
425 functional information from model plant species to crops. The β -group evolved before the split of monocots
426 and eudicots and contains CRKs and CRRSPs that expanded through WGDs followed by lineage-specific
427 tandem duplications. Domain re-arrangements in the β -clade led to secondary groups of ddCRRSPs and
428 sdCRRSPs while the recruitment of a different kinase domain in grasses suggests the re-routing of signaling
429 pathways towards novel phosphorylation substrates. Thus, it is likely that members of the β -group have been
430 subject to sub- and neo-functionalization, which is a challenge for functional analyses. The domain
431 exchanges in DUF26-containing proteins highlight the importance of comparative analysis of phylogenies
432 inferred from full-length protein sequences with those inferred from individual domains. WGDs have been
433 associated with periods of environmental upheaval and an increase in biological complexity^{2,60}. Accordingly,
434 the appearance and radiation of DUF26-containing proteins with different domain structures as well as CRK
435 and CRRSP expansions co-occur with the evolution of novel physiological characteristics, such as
436 vasculature, and with the adaptation to new habitats and lifestyles (Figure 1b).

437 Sequence analysis suggested that DUF26 proteins could be specific to embryophytes. Crystallographic
438 analysis of two PDLP ectodomains reveals that the structure of the DUF26 domains closely matches the fold
439 of the sdCRRSP Gnk-2, which is evolutionarily distant from the PDLPs. PDLPs contain two DUF26 domain
440 and the structure of Gnk-2 is more similar to the DUF26-A. However, despite the high structural similarity
441 the mannose-binding function of Gnk-2 is not conserved in the PDLP DUF26-A domain. Intriguingly, plant
442 DUF26 domains share significant structural similarity to fungal carbohydrate-binding modules. Notably, the
443 tandem arrangement of two lectin-like DUF26 domains appears to be plant-specific. Rapid sequence
444 divergence⁶¹ is a limiting factor in detection of homology at the amino acid sequence level, seen e.g. in the

445 marked differences between DUF26 from *Marchantia polymorpha* and *Physcomitrella patens* and those
446 from other plants. This may obscure identification of DUF26 domains in charophytes and other algal species.
447 The physiological ligands of ddCRRSP, CRKs and PDLPs remain to be discovered and our work suggests
448 that different tandem DUF26 domains likely recognize diverse sets of ligands. Similar to plant malectin
449 receptors⁶², DUF26 domains may have evolved novel or additional functions which might include mediation
450 of protein-protein interactions at the cell surface^{20,35}. The strong structural similarity between DUF26
451 domains and fungal lectins suggests a common origin, and DUF26 proteins represent novel carbohydrate-
452 binding domains in plants. Identification of ligands for different DUF26-domains will provide novel insights
453 in to perception of cell wall status or environmental signals. However, this may be challenging since plant
454 cells and their cell walls contain a large number of carbohydrates and related compounds.

455 From the evolutionary analysis, an overall model emerges (Figure 8). The young gene families initially
456 expand through tandem duplications and therefore experience more relaxed selection⁶³. This is supported by
457 the fact that the tandem genes function in processes that require fast responses such as adaptation to
458 environment, pathogen responses and secondary metabolism^{2,64}, and that these gene families show high
459 variation across species and have high kn/ks rates⁵. In tandems, main evolutionary forces are unequal
460 crossover and concerted evolution through gene conversion, but over time the genes evolve into their
461 specific functions. This process may be interrupted by WGM events. Since the tandem genes are not
462 evolving under dosage balance, there is no compensatory drift⁶⁵, and thus drift and selection by dosage
463 eventually drives one of the duplicates into fixation while others turn into pseudogenes. Assuming that the
464 elements driving tandem duplications are still present after fractionation, the remaining duplicates may in
465 turn expand. In case of a tandem where all genes have established a unique functional role in the system,
466 drift may drive the duplicated tandem into scattered orthologs. These orthologs may eventually assume a
467 fixed syntenic position in the genome and switch to a dosage balance mode of evolution. The evolutionary
468 mode of the gene family would depend on the balance between the death rate after WGMs and the birth rate
469 of the tandem duplications.

470 Our study of DUF26-containing proteins demonstrates both the challenges of analyses of large protein
471 families and the power of combining advanced evolutionary and structural methods. Our analysis will

472 provide a model for future studies of similarly large protein families and will facilitate the forthcoming
473 detailed biochemical and physiological investigation of the mechanistic functions of CRKs, PDLPs and
474 CRRSPs in different plant species.

475 **Materials and methods**

476 **Gene identification and annotation**

477 Altogether 32 plant and algae genomes (Table S1) covering the major plant lineages were selected for
478 analyses. For 27 species protein annotations (primary transcripts) and genome sequence data was retrieved
479 from Phytozome⁶⁶, and Barley (*Hordeum vulgare*) from Gramene (<http://www.gramene.org>) with the latest
480 names for gene models from IPK server (http://webblast.ipk-gatersleben.de/barley_ibsc/)⁶⁷. Silver birch
481 (*Betula pendula*) was sequenced at the University of Helsinki⁵⁶. Eggplant (*Solanum melongena*) data was
482 retrieved from Eggplant Genome DataBase (<http://eggplant.kazusa.or.jp/>). *Klebsormidium flaccium* and
483 Sacred lotus (*Nelumbo nucifera*) genome data were from NCBI (<https://www.ncbi.nlm.nih.gov>).
484 Additionally the FungiDB⁶⁸ (www.fungidb.org), InsectBase⁶⁹ (<http://www.insect-genome.com>) and human
485 (*Homo sapiens*), chicken (*Gallus gallus*) and zebrafish (*Danio rerio*) genomes were screened for DUF26.
486 Detailed information of the genome versions and references are given in the Table S6.

487 HMMER (version 3.1b2) search⁷⁰ for PFAM domain with ID PF01657 (stress-antifungal domain) was
488 carried out among AA sequences representing gene models from different species⁷¹. Genome sequences were
489 checked with Wise2 (version 2.4.1) software^{72,73}. All gene models found with HMMER were manually
490 curated, and new genes found with Wise2 were manually annotated using Fgenesh+⁷⁴. Birch (*Betula*
491 *pendula*)⁵⁶ and Sacred lotus (*Nelumbo nucifera*) were fully manually annotated as they did not have gene
492 models *a priori*. High rates of manual annotation and curation were needed for *Selaginella moellendorffii*,
493 grapevine (*Vitis vinifera*; version Genoscope.12X⁷⁵) and potato (*Solanum tuberosum*). Sequences from each
494 species were further checked by carrying out a multiple sequence alignment and phylogenetic tree estimation
495 with PASTA⁷⁶. Partial gene models were identified by checking sequences individually. Genes were defined
496 as pseudogenes if the genomic sequence was available but no full domain structure could be predicted. In
497 cases where the prediction problem was caused by the length of the contig or a gap in the genome sequence
498 the gene model was marked as partial. Pseudogenes and partial gene models were not included in the
499 subsequent analyses.

500 For domain analyses and phylogenetic trees containing only domain sequences, the domain borders were
501 defined with HMMER using the PFAM domain PF01657 for DUF26, and PF07714 for the kinase domain

502 from curated dataset. The ectodomain region was defined to end at the border of the transmembrane region
503 in the PDLPs and CRKs. The partial PDLP from *Marsilea quadrifolia* was identified by using pBLAST
504 search against sequences in the NCBI database.

505 **Phylogenetic trees**

506 Only full gene models were used to infer phylogenetic trees. Sequence quality in alignments was checked
507 using Guidance (version 2.01) and alignments were built using the MAFFT option⁷⁷. Sequences with low
508 quality score were removed from datasets and alignments were built again with PASTA. For phylogenetic
509 trees, alignments were filtered in in Wasabi⁷⁸ to remove residues with less than 10 percent coverage.
510 Filtering was required due to the high sequence diversity (on less conserved regions) resulting in a high
511 number of gaps in multiple sequence alignments. Maximum likelihood (ML) phylogenetic trees were
512 inferred for filtered and also unfiltered data using RAxML (version 8.1.3)⁷⁹.

513 ML phylogenetic trees were bootstrapped using RAxML (version 8.1.3) for 1000 bootstrap replicates. For
514 phylogenetic trees containing full length sequences with all domain structures bootstrapping was also carried
515 out with partitioning (both DUF26 and kinase domains defined separately). The PROTGAMMAJTT model
516 was used in phylogenetic analyses using RAxML. Model selection was based on a Perl script for identifying
517 the optimal protein substitution model (available in RAxML webpage, provided by Alexandros Stamatakis).
518 Bootstrapped trees are available on Wasabi⁷⁸ (see figure legends). Comparison on phylogenetic trees based
519 on CRK ectodomain and kinase domain regions was visualized in R using the "dendextend" package.

520 **Exon intron structure**

521 The number of exons for all genes was estimated using Scipio (version 1.4.1)⁸⁰ using default parameters
522 (minimum identity of 90% and coverage of 60%). It internally uses BLAT to perform the initial alignment of
523 the protein sequences against the genome followed by refinement of hits to determine the exact splicing
524 borders and to obtain the final gene structure. The number of exons per gene was extracted from the final
525 result.

526 **Orthogroup generation**

527 11 representative species from different clades (*Arabidopsis thaliana*, *Amborella trichopoda*, *Oryza sativa*,
528 *Zea Mays*, *Vitis vinifera*, *Populus trichocarpa*, *Aquilegia coerulea*, *Brachypodium distachyon*,

529 *Physcomitrella patens*, *Selaginella moellendorffii* and *Spirodela polyrhiza*) were chosen to study the
530 evolution of the DUF26-containing proteins. Primary protein sequences of these 11 species were
531 downloaded from Phytozome (version 11.0). An all-against-all BLAST was run for all the protein sequences
532 followed by generation of orthogroups using the software OrthoMCL (version 2.0.9)⁸¹ with an inflation
533 parameter of 1.5 for the clustering phase. Clustering yielded 34,535 orthogroups.

534 **Species tree generation**

535 Orthogroups containing one representative protein for each of the 11 species were chosen to generate the
536 species tree. Multiple sequence alignment was carried out on the single copy orthogroups using PRANK⁸²
537 and the output was used to infer a species tree using RAxML⁷⁹.

538 **Evolutionary rate and ancestral size estimation**

539 The evolutionary rate and ancestral size of the orthogroups were modelled using Badirate software (version
540 1.35)⁴¹. The species tree and orthogroups generated from the previous steps were used as input for Badirate.
541 The BDI (Birth, Death, Innovation) rate model was used. The Free Rates (FR) branch model was chosen
542 which would assume every branch of the species tree to have its own turnover rates. Turnover rates of
543 orthogroups were estimated using the maximum likelihood fitting. Orthogroups were defined as protein
544 kinases if they included sequences with PFAM domain PF00069. Orthogroups containing RLKs were
545 defined based on known *Arabidopsis* RLKs¹⁵. Plasmodesmata-related orthogroups were defined based on
546 *Arabidopsis thaliana* genes related to plasmodesmata⁴⁰.

547 **Nucleotide CDS sequence generation from protein sequence for PAML**

548 The GFF file output from Scipio⁸⁰ was pre-processed by an in-house script and processed with the gff3
549 module of the GenomeTools (version 1.5.4)⁸³ software. The final GFF file along with the corresponding
550 species genome in fasta formatted file was passed as an input to the extractfeat module of the GenomeTools
551 software to extract the final nucleotide CDS sequences.

552 **PAML analyses**

553 We estimated d_N/d_S ratios (ratio of non-synonymous and synonymous sites, ω) for conserved clades (bCRK-
554 I, bCRK-II, CRCKs (orthologs of AtCRK43), PDLPs and sdCRRSPs) from eleven species (*Arabidopsis*
555 *thaliana*, *Amborella trichopoda*, *Oryza sativa*, *Zea Mays*, *Vitis vinifera*, *Populus trichocarpa*, *Aquilegia*

556 *coerulea*, *Brachypodium distachyon*, *Physcomitrella patens*, *Selaginella moellendorffii* and *Spirodela*
557 *polyrhiza*) by using the codeml program from PAML (version 4.9)⁸⁴. We applied the one-ratio model (M0)
558 to estimate overall d_N/d_S ratios for each conserved group separately and free ratios neutral model (M1) to
559 estimate d_N/d_S ratios for each branch within conserved clades⁸⁵. To study the evolution of PDL5 and
560 PDL8, sitewise-analyses of their homologs was carried out. As PDL5 is specific to *Brassicaceae*, we
561 additional nucleotide sequences for orthologs of AtPDL5 from NCBI, Phytozome and CoGe databases.
562 Furthermore, additional sequences for orthologs of AtPDL8 were included in the alignment to improve
563 depth and reliability of the analysis. Multiple sequence alignments of coding nucleotide sequences were
564 constructed with PRANK⁸² and phylogenetic trees were estimated using RAxML⁷⁹ for codeml.

565 **Syntenic vs tandem duplications**

566 Syntenic and tandem duplications were analysed using Synmap application in CoGe⁸⁶, using default settings.
567 Tandem duplications were defined as genome regions with at least three to five duplicate genes (Table S4).
568 Synteny comparisons were done between *Arabidopsis thaliana* and *Solanum lycopersicum*, *S. lycopersicum*
569 and *Amborella trichopoda*, *A. trichopoda* and *Oryza sativa* and *Zea mays* and *Oryza sativa*. Tandem
570 duplication results from DAGchainer were collected for each species. The results were filtered based on
571 annotated gene models from selected species. The currently available *Amborella trichopoda* genome is
572 presented only as scaffolds, and the genes were placed to chromosomes based on physical mapping⁸⁷.
573 Scaffolds not assigned to any chromosome were added separately. Thus the location of the *Amborella*
574 *trichopoda* genes in the genome is only a rough estimate (Figure 7a).

575 **Gene conversion analyses**

576 Gene conversion events were estimated from nucleotide sequences for the same eleven species that were
577 analyzed for d_N/d_S ratios with GENECONV (version 1.81a)⁸⁸. Analyses were carried out for the main clades
578 of the eleven species. For bCRKs and vCRKs separate analyses were carried out using sequences from the
579 five species used in synteny analyses (*Arabidopsis thaliana*, *Amborella trichopoda*, *Oryza sativa*, *Solanum*
580 *lycopersicum* and *Zea mays*). The largest tandem region of vCRKs in *A. thaliana* chromosome 4 was
581 analyzed separately to validate the results from the analysis with all vCRKs from *A. thaliana*.

582

583 **Gene tree reconciliation**

584 Gene tree reconciliation was carried out using DLCpar (version 1.0)⁸⁹ downloaded from
585 <https://www.cs.hmc.edu/~yjw/software/dlcpar/>. NCBI taxonomy was used as the species tree, downloaded in
586 newick format from PhyloT website, <http://phylot.biobyte.de/>. Reconciliation was carried out using DLCpar
587 search with 20 prescreening iterations, followed by 1000 search iterations. The solution was visualized in R,
588 using custom scripts and ‘ape’ package.

589 **Phenomics data analysis**

590 Phenotyping data of T-DNA mutant insertion lines was normalized against the Col-0 data by calculating Z-
591 scores, see Bourdais *et al.*²⁴ The standard deviation (SD) over all experiments was calculated for each allele,
592 and in case of several insertion alleles the one with maximum SD was selected. The residuals of the bCRK
593 vs vCRK split in the data were tested for normality using Shapiro’s test. Since the null hypothesis
594 (normality) was rejected with $p < 0.05$ the difference between groups was tested with Wilcox test.

595 **Transcriptomic analyses**

596 Paired end RNAseq data was collected from the publicly available sequence read archive (SRA) database by
597 fastq-dump.2 (version 2.5.7) for *Arabidopsis thaliana*, *Oryza sativa*, *Solanum lycopersicum* and *Zea mays*.
598 FastQC (version 0.11.4) (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) was used to check the
599 quality of the samples. Low quality reads and bases were removed by Trimmomatic (version 0.36)⁹⁰ with the
600 following options: phred33, TRAILING: 20, and MINLEN: 30. Filtered reads were mapped to gene models
601 from Phytozome version 12, by Kallisto, run in paired end mode (version 0.43.1, --bias and --bootstrap:
602 200)⁹¹. Bootstrap samples were averaged (custom R code) and gene expression abundance (transcript per
603 million [TPM]) was estimated by tximport (version 1.2.0)⁹² followed by averaging over biological replicates.
604 Ortholog comparison between species was carried out by grouping the experiments into seven categories,
605 with maximum TPM among experiments representing gene response. Pearson correlation was calculated
606 among orthologs and all other possible pairs.

607 **Protein expression and purification**

608 An expression construct coding for the *PDLP5* ectodomain (amino acids 1-241) was codon optimized for
609 *Spodoptera frugiperda* and synthesized by Geneart (Thermo Fisher). Using the PfuX7 polymerase⁹³, the

610 gene for the *PDLP8*-ECD (1-253) was amplified from *Arabidopsis thaliana* cDNA. The Gibson assembly
611 method⁹⁴ was employed to insert the *PDLP5* and *PDLP8* ectodomain coding sequences into an adapted
612 pFAST-BAC1 vector (Geneva Biotech), providing a C-terminal 2x-STREP-9xHIS tag. *PDLP5* point
613 mutations (C101A, C148A and C191A) were then introduced as described⁹⁵. Bacmids were generated by
614 transforming the plasmids (confirmed by sequencing) into *Escherichia coli* DH10MultiBac (Geneva
615 Biotech). Virus particles were created by transfecting (Profectin, AB Vector) the bacmids into *Spodoptera*
616 *frugiperda* SF9 cells. For secreted protein production, *Trichoplusia ni* Tnao38 cells were infected with a viral
617 multiplicity of 1, incubated for 3 days at 22 °C. The protein-containing supernatant was separated from the
618 intact cells by centrifugation and subjected to Ni²⁺-affinity chromatography (HisTrap Excel; GE Healthcare)
619 in buffer A (10 mM Hepes 7.5, 500 mM NaCl). Bound proteins eluted in buffer A supplemented with 500
620 mM imidazole. The elution fractions were pooled and further purified by StrepII-affinity purification (Strep-
621 Tactin XT Superflow high capacity, IBA) in buffer B (20 mM Tris pH 8.0, 250 mM NaCl, 1 mM EDTA).
622 The column was washed with 5-10 column volumes of buffer B and eluted in buffer B supplemented with 50
623 mM biotin. The C-terminal 2x-STREP-9xHIS tag was subsequently removed by adding tobacco etch virus
624 (TEV)-protease to the StrepII elution in a 1:100 ratio for 16h at 4 °C. The 2x-STREP-9xHIS-tag and the HIS-
625 tagged TEV-protease were then separated from the respective ectodomain by an additional Ni²⁺-affinity
626 chromatography step (HisTrap Excel; GE Healthcare). Cleaved *PDLP5*, *PDLP5*^{C101A}, *PDLP5*^{C148A},
627 *PDLP5*^{C191A} and *PDLP8* ectodomains were next subjected to preparative size exclusion chromatography
628 using either a HiLoad 26/600 Superdex 200 pg (*PDLP5* and *PDLP8*) or HiLoad 16/600 Superdex 200 pg
629 (*PDLP5*^{C101A}, *PDLP5*^{C148A} and *PDLP5*^{C191A}) column, equilibrated in 20 mM sodium citrate pH 5.0 and 150
630 mM NaCl. Monomeric peak fractions were collected and concentrated using an Amicon Ultra (Milipore)
631 filter device. The concentrated monomeric peak fractions of *PDLP5*, *PDLP5*^{C101A}, *PDLP5*^{C148A} and
632 *PDLP5*^{C191A} were additionally subjected to analytical size exclusion chromatography on a Superdex 200
633 Increase 10/300 GL column (GE healthcare) equilibrated in 20 mM citrate pH 5.0 and 150 mM NaCl.

634 **Thermostability assay**

635 20 µl reactions consisted of either *PDLP5*, *PDLP5*^{C101A}, *PDLP5*^{C148A} and *PDLP5*^{C191A} ectodomains at a
636 concentration of 1.5 mg/ml in 20 mM citrate pH 5.0, 150 mM NaCl, 10x SYPRO Orange dye (Thermo

637 Fisher), and were mixed in a 384-well ABI PRISM plate (Applied Biosystems). Using a 7900HT Fast Real-
638 Time PCR system SYPRO Orange fluorescence was measured. The reactions were initially incubated for 2
639 min at 25 °C and then the temperature was increased to 95 °C at a heating rate of 0.5 °C/min. Resulting
640 melting curves were fitted with a Boltzman function using GraphPad Prism and the melting temperatures,
641 T_m , correspond to the first inflection point of the Boltzman fit.

642 **Isothermal titration calorimetry**

643 ITC experiments were performed at 25°C using a Nano ITC (TA Instruments, New Castle, USA) with a 1.0
644 mL standard cell and a 250 µl titration syringe. The PDLP5 ectodomain was gelfiltrated into ITC buffer (20
645 mM sodium citrate pH 5.0, 150 mM NaCl) and all carbohydrates were resuspended into ITC buffer. The
646 experiments were carried out by injecting 24 times 10 µl of D-+-Mannose (1 mM; Sigma), Pectic Galactan
647 (2mg/ml; Megazyme), Rhamnogalacturonan (2mg/ml; Megazyme), Polygalacturonic Acid (2mg/ml;
648 Megazyme), Cellohexaose (1 mM; Megazyme) or Arabinohexaose (1 mM; Megazyme) aliquots into PDLP5
649 (~100 µM) in the cell at 150 s intervals. ITC data for the D-+-Mannose experiment were corrected for the
650 heat of dilution by subtracting the mixing enthalpies for titrant solution injections into protein free ITC
651 buffer. Data were analyzed using the NanoAnalyze program (version3.5) as provided by the manufacturer.

652 **Protein crystallization and crystallographic data collection**

653 The PDLP5 ectodomain formed crystals in hanging drops composed of 1 µl of protein solution (70 mg/ml in
654 20 mM citrate pH 5.0 and 150 mM NaCl) and 1 µl of crystallization buffer (17.5 % [w/v] polyethylene
655 glycol 4,000, 250 mM (NH₄)₂SO₄) suspended over 800 µl of the latter as reservoir solution. Protein crystals
656 were transferred into crystallization buffer supplemented with 25% (v/v) ethylene glycol, which served as
657 cryoprotectant, and snap frozen in liquid N₂. PDLP8 crystals (52 mg/ml in 20 mM citrate pH 5.0, 150 mM
658 NaCl) developed in hanging drops containing 17.5 % (w/v) polyethylene glycol 4,000, 0.1 M citrate pH 5.5,
659 20 % (v/v) 2-propanol. Crystals were frozen directly in liquid N₂. For PDLP5 native ($\lambda= 1.0 \text{ \AA}$) and
660 redundant sulfur SAD ($\lambda= 2.079 \text{ \AA}$) data were collected to 1.29 Å resolution at beam line PX-III of the Swiss
661 Light Source (SLS), Villigen, Switzerland. A 1.95 Å native data set of PDLP8 was acquired at the same
662 beam line. Data processing and reduction was done with XDS (version: Jan, 2018)⁹⁶.

663 **Structure solution and refinement**

664 The structure of PDLP5 was solved using the single-anomalous diffraction (SAD) method. 24 S sites
665 corresponding to the 12 disulfide bonds in the PDLP5 crystallographic dimer were located with the program
666 SHELXD⁹⁷, site-refinement and phasing was done in SHARP⁹⁸ and the starting phases were used for
667 automated model building in BUCCANEER⁹⁹ and ARP/wARP¹⁰⁰. The model was completed in alternating
668 cycles of model correction in COOT¹⁰¹ and restrained refinement in Refmac5¹⁰². The structure of PDLP8 was
669 solved using the molecular replacement methods as implemented in the program PHASER¹⁰³, and using the
670 refined PDLP5 tandem ectodomain as search model. Inspection with MolProbity¹⁰⁴ revealed excellent
671 stereochemistry for the final models. Structural and surface representations were done in Pymol
672 (<http://pymol.sourceforge.org>) and Chimera¹⁰⁵.

673 **Data availability**

674 Materials used in this study and data generated are available from the corresponding author upon request.
675 Phylogenetic trees with bootstrap information for 1000 replicates and corresponding sequence alignments
676 have been deposited on Wasabi (<http://wasabiapp.org>); identifiers are available in the figure legends as web
677 links. Information on used genomic data is available in Table S5. Publically available gene expression data
678 was taken from the Sequence Read Archive (SRA) database; identifiers are listed in Table S4.
679 Crystallographic coordinates and structure factors have been deposited with the Protein Data Bank
680 (<http://rcsb.org>) with accession codes 6GRE (PDLP5) and 6GRF (PDLP8).

681 **Code availability**

682 All R scripts developed for parsing the data and visualizing the results are available upon request.

683 **Authors' contributions**

684 AV, BB, JK, JS, MH and MW conceived and designed the project. AL, BB, OS, SR, ML, AVe, AL, MH,
685 and JS carried out the analyses. AV, BB, AL, MH, JS, and MW analyzed the data. AV, BB, MH, JS and MW
686 wrote the manuscript. All authors read and contributed to the final manuscript.

687 **Acknowledgments**

688 The authors would like to thank Ms. Kerri Hunter, Drs. Julia Krasensky-Wrzaczek, Sachie Kimura and
689 Alexey Shapiguzov for critical comments on the manuscript. We thank Prof. David Collinge and Prof.

690 Michael Lyngkjær (University of Copenhagen, Denmark) for help with the barley genome and Prof. Eric
691 Schranz (Wageningen University, Netherlands) for helpful comments. This work was supported by the
692 Doctoral Programme in Plant Sciences (DPPS) of the University of Helsinki and by the Finnish Cultural
693 Foundation Suomen Kulttuurirahasto (to AV), the Academy of Finland (grant numbers #275632, #283139
694 and #312498 to MW) and the University of Helsinki (Three-year fund allocation to MW) and by grant
695 31003A_176237 from the Swiss National Science Foundation and by an International Research Scholar
696 Award from the Howard Hughes Medical Institute (to MH). AV, OS, SR, JK, JS and MW are members of
697 the Centre of Excellence in the Molecular Biology of Primary Producers (2014-2019) funded by the
698 Academy of Finland (grant numbers #271832 and #307335). BB was supported by an EMBO long-term
699 fellowship. We would also like to thank the staff at beam line PXIII of the Swiss Light Source (Villigen,
700 Switzerland) for technical assistance during data collection.

701 References

- 702 1 Demuth, J. P. & Hahn, M. W. The life and death of gene families. *Bioessays* **31**, 29-39 (2009).
703 2 Panchy, N., Lehti-Shiu, M. & Shiu, S. H. Evolution of gene duplication in plants. *Plant Physiol* **171**,
704 2294-2316 (2016).
705 3 Tasdighian, S. *et al.* Reciprocally retained genes in the angiosperm lineage show the hallmarks of
706 dosage balance sensitivity. *Plant Cell* **29**, 2766-2785 (2017).
707 4 Birchler, J. A., Bhadra, U., Bhadra, M. P. & Auger, D. L. Dosage-dependent gene regulation in
708 multicellular eukaryotes: implications for dosage compensation, aneuploid syndromes, and
709 quantitative traits. *Dev Biol* **234**, 275-288 (2001).
710 5 Wang, P. *et al.* Factors influencing gene family size variation among related species in a plant
711 family, Solanaceae. *Genome Biol Evol* **10**, 2596-2613 (2018).
712 6 Nakamura, S., Suzuki, T., Kawamukai, M. & Nakagawa, T. Expression analysis of *Arabidopsis*
713 *thaliana* small secreted protein genes. *Biosci Biotechnol Biochem* **76**, 436-446 (2012).
714 7 Agrawal, G. K., Jwa, N. S., Lebrun, M. H., Job, D. & Rakwal, R. Plant secretome: unlocking secrets
715 of the secreted proteins. *Proteomics* **10**, 799-827 (2010).
716 8 Tavormina, P., De Coninck, B., Nikonorova, N., De Smet, I. & Cammue, B. P. The plant peptidome:
717 An expanding repertoire of structural features and biological functions. *Plant Cell* **27**, 2095-2118
718 (2015).
719 9 Shiu, S. H. & Bleecker, A. B. Receptor-like kinases from *Arabidopsis* form a monophyletic gene
720 family related to animal receptor kinases. *Proc Natl Acad Sci USA* **98**, 10763-10768 (2001).
721 10 Shiu, S. H. & Bleecker, A. B. Expansion of the receptor-like kinase/Pelle gene family and receptor-
722 like proteins in *Arabidopsis*. *Plant Physiol* **132**, 530-543 (2003).
723 11 Fritz-Laylin, L. K., Krishnamurthy, N., Tor, M., Sjolander, K. V. & Jones, J. D. Phylogenomic
724 analysis of the receptor-like proteins of rice and *Arabidopsis*. *Plant Physiol* **138**, 611-623 (2005).
725 12 Smakowska-Luzan, E. *et al.* An extracellular network of *Arabidopsis* leucine-rich repeat receptor
726 kinases. *Nature* **553**, 342-346 (2018).
727 13 Shiu, S. H. & Bleecker, A. B. Plant receptor-like kinase gene family: diversity, function, and
728 signaling. *Sci STKE* **2001**, re22 (2001).
729 14 Shiu, S. H. *et al.* Comparative analysis of the receptor-like kinase family in *Arabidopsis* and rice.
730 *Plant Cell* **16**, 1220-1234 (2004).
731 15 Zulawski, M., Schulze, G., Braginets, R., Hartmann, S. & Schulze, W. X. The *Arabidopsis* Kinome:
732 phylogeny and evolutionary insights into functional diversification. *BMC Genomics* **15**, 548 (2014).

- 733 16 Fischer, I., Dievart, A., Droc, G., Dufayard, J. F. & Chantret, N. Evolutionary dynamics of the
734 leucine-rich repeat receptor-like Kinase (LRR-RLK) subfamily in angiosperms. *Plant Physiol* **170**,
735 1595-1610 (2016).
- 736 17 Kimura, S., Waszczak, C., Hunter, K. & Wrzaczek, M. Bound by fate: The role of reactive oxygen
737 species in receptor-like kinase signaling. *Plant Cell* **29**, 638-654 (2017).
- 738 18 Miyakawa, T. *et al.* A secreted protein with plant-specific cysteine-rich motif functions as a
739 mannose-binding lectin that exhibits antifungal activity. *Plant Physiol* **166**, 766-778 (2014).
- 740 19 Miyakawa, T., Miyazono, K., Sawano, Y., Hatano, K. & Tanokura, M. Crystal structure of
741 ginkbilobin-2 with homology to the extracellular domain of plant cysteine-rich receptor-like kinases.
742 *Proteins* **77**, 247-251 (2009).
- 743 20 Ma, L. S. *et al.* The *Ustilago maydis* repetitive effector Rsp3 blocks the antifungal activity of
744 mannose-binding maize proteins. *Nat Commun* **9**, 1711 (2018).
- 745 21 Acharya, B. R. *et al.* Overexpression of CRK13, an *Arabidopsis* cysteine-rich receptor-like kinase,
746 results in enhanced resistance to *Pseudomonas syringae*. *Plant J* **50**, 488-499 (2007).
- 747 22 Chen, K., Du, L. & Chen, Z. Sensitization of defense responses and activation of programmed cell
748 death by a pathogen-induced receptor-like protein kinase in *Arabidopsis*. *Plant Mol Biol* **53**, 61-74
749 (2003).
- 750 23 Chen, K., Fan, B., Du, L. & Chen, Z. Activation of hypersensitive cell death by pathogen-induced
751 receptor-like protein kinases from *Arabidopsis*. *Plant Mol Biol* **56**, 271-283 (2004).
- 752 24 Bourdais, G. *et al.* Large-scale phenomics identifies primary and fine-tuning roles for CRKs in
753 responses related to oxidative stress. *PLOS Genetics* **11**, e1005373 (2015).
- 754 25 Idänheimo, N. *et al.* The *Arabidopsis thaliana* cysteine-rich receptor-like kinases CRK6 and CRK7
755 protect against apoplastic oxidative stress. *Biochem Biophys Res Commun* **445**, 457-462 (2014).
- 756 26 Wrzaczek, M. *et al.* Transcriptional regulation of the CRK/DUF26 group of receptor-like protein
757 kinases by ozone and plant hormones in *Arabidopsis*. *BMC Plant Biol* **10**, 95 (2010).
- 758 27 Yeh, Y. H., Chang, Y. H., Huang, P. Y., Huang, J. B. & Zimmerli, L. Enhanced *Arabidopsis* pattern-
759 triggered immunity by overexpression of cysteine-rich receptor-like kinases. *Front Plant Sci* **6**, 322
760 (2015).
- 761 28 Yadeta, K. A. *et al.* A cysteine-rich protein kinase associates with a membrane immune complex and
762 the cysteine residues are required for cell death. *Plant Physiol* **173**, 771-787 (2017).
- 763 29 Lee, D. S. K., Young Cheon; Kwon, Sun Jae; Ryu, Choong-Min; Park, Ohkmae K. The *Arabidopsis*
764 Cysteine-Rich Receptor-Like Kinase CRK36 regulates immunity through interaction with the
765 cytoplasmic kinase BIK1. *Frontiers in Plant Science* **8**, 1856 (2017).
- 766 30 Tanaka, H. *et al.* Abiotic stress-inducible receptor-like kinases negatively control ABA signaling in
767 *Arabidopsis*. *Plant J* **70**, 599-613 (2012).
- 768 31 Chern, M. *et al.* A genetic screen identifies a requirement for cysteine-rich-receptor-like kinases in
769 Rice NH1 (OsNPR1)-mediated immunity. *PLoS Genet* **12**, e1006049 (2016).
- 770 32 Brunkard, J. O. & Zambryski, P. C. Plasmodesmata enable multicellularity: new insights into their
771 evolution, biogenesis, and functions in development and immunity. *Curr Opin Plant Biol* **35**, 76-83
772 (2017).
- 773 33 Caillaud, M. C. *et al.* The plasmodesmal protein PDL1 localises to haustoria-associated membranes
774 during downy mildew infection and regulates callose deposition. *PLoS Pathog* **10**, e1004496 (2014).
- 775 34 Lim, G. H. *et al.* Plasmodesmata localizing proteins regulate transport and signaling during systemic
776 acquired immunity in plants. *Cell Host Microbe* **19**, 541-549 (2016).
- 777 35 Cui, W. & Lee, J. Y. *Arabidopsis* callose synthases CalS1/8 regulate plasmodesmal permeability
778 during stress. *Nat Plants* **2**, 16034 (2016).
- 779 36 Amari, K. *et al.* A family of plasmodesmal proteins with receptor-like properties for plant viral
780 movement proteins. *PLoS Pathog* **6**, e1001119 (2010).
- 781 37 Xu, G., Ma, H., Nei, M. & Kong, H. Evolution of F-box genes in plants: different modes of sequence
782 divergence and their relationships with functional diversification. *Proc Natl Acad Sci USA* **106**, 835-
783 840 (2009).
- 784 38 Rody, H. V., Baute, G. J., Rieseberg, L. H. & Oliveira, L. O. Both mechanism and age of
785 duplications contribute to biased gene retention patterns in plants. *BMC Genomics* **18**, 46 (2017).
- 786 39 Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable future. *Nucleic
787 Acids Res* **44**, D279-285 (2016).

- 788 40 Fernandez-Calvino, L. *et al.* Arabidopsis plasmodesmal proteome. *PLoS One* **6**, e18880 (2011).
- 789 41 Librado, P., Vieira, F. G. & Rozas, J. BadiRate: estimating family turnover rates by likelihood-based
790 methods. *Bioinformatics* **28**, 279-281 (2012).
- 791 42 Van de Peer, Y., Maere, S. & Meyer, A. The evolutionary significance of ancient genome
792 duplications. *Nat Rev Genet* **10**, 725-732 (2009).
- 793 43 Hsu, T. C. *et al.* Early genes responsive to abscisic acid during heterophyllous induction in *Marsilea*
794 *quadrifolia*. *Plant Mol Biol* **47**, 703-715 (2001).
- 795 44 Holm, L. & Rosenstrom, P. Dali server: conservation mapping in 3D. *Nucleic Acids Res* **38**, W545-
796 549 (2010).
- 797 45 van Eerde, A., Grahn, E. M., Winter, H. C., Goldstein, I. J. & Kregel, U. Atomic-resolution
798 structure of the alpha-galactosyl binding *Lyophyllum decastes* lectin reveals a new protein family
799 found in both fungi and plants. *Glycobiology* **25**, 492-501 (2015).
- 800 46 Zhang, P. *et al.* Cytotoxic protein from the mushroom *Coprinus comatus* possesses a unique mode
801 for glycan binding and specificity. *Proc Natl Acad Sci USA* **114**, 8980-8985 (2017).
- 802 47 Vijayan, M. & Chandra, N. Lectins. *Curr Opin Struct Biol* **9**, 707-714 (1999).
- 803 48 Hohmann, U., Lau, K. & Hothorn, M. The structural basis of ligand perception and signal activation
804 by receptor kinases. *Annu Rev Plant Biol* **68**, 109-137 (2017).
- 805 49 Berrabah, F. *et al.* A nonRD receptor-like kinase prevents nodule early senescence and defense-like
806 reactions during symbiosis. *New Phytol* **203**, 1305-1314 (2014).
- 807 50 Dardick, C. & Ronald, P. Plant and animal pathogen recognition receptors signal through non-RD
808 kinases. *PLoS Pathog* **2**, e2 (2006).
- 809 51 Dardick, C., Schwessinger, B. & Ronald, P. Non-arginine-aspartate (non-RD) kinases are associated
810 with innate immune receptors that recognize conserved microbial signatures. *Curr Opin Plant Biol*
811 **15**, 358-366 (2012).
- 812 52 Rogers, R. L., Shao, L. & Thornton, K. R. Tandem duplications lead to novel expression patterns
813 through exon shuffling in *Drosophila yakuba*. *PLoS Genet* **13**, e1006795 (2017).
- 814 53 Hastings, P. J., Lupski, J. R., Rosenberg, S. M. & Ira, G. Mechanisms of change in gene copy
815 number. *Nat Rev Genet* **10**, 551-564 (2009).
- 816 54 Chen, J. M., Cooper, D. N., Chuzhanova, N., Ferec, C. & Patrinos, G. P. Gene conversion:
817 mechanisms, evolution and human disease. *Nat Rev Genet* **8**, 762-775 (2007).
- 818 55 Zou, C., Lehti-Shiu, M. D., Thomashow, M. & Shiu, S. H. Evolution of stress-regulated gene
819 expression in duplicate genes of *Arabidopsis thaliana*. *PLoS Genet* **5**, e1000581 (2009).
- 820 56 Salojärvi, J. *et al.* Genome sequencing and population genomic analyses provide insights into the
821 adaptive landscape of silver birch. *Nat Genet* **49**, 904-912 (2017).
- 822 57 Fischer, I. *et al.* Impact of recurrent gene duplication on adaptation of plant genomes. *BMC Plant*
823 *Biol* **14**, 151 (2014).
- 824 58 Veitia, R. A., Bottani, S. & Birchler, J. A. Cellular reactions to gene dosage imbalance: genomic,
825 transcriptomic and proteomic effects. *Trends Genet* **24**, 390-397, doi:10.1016/j.tig.2008.05.005
826 (2008).
- 827 59 Guo, Y. L. Gene family evolution in green plants with emphasis on the origination and evolution of
828 *Arabidopsis thaliana* genes. *Plant J* **73**, 941-951 (2013).
- 829 60 Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of polyploidy. *Nat Rev*
830 *Genet* **18**, 411-424 (2017).
- 831 61 Copley, R. R., Goodstadt, L. & Ponting, C. Eukaryotic domain evolution inferred from genome
832 comparisons. *Curr Opin Genet Dev* **13**, 623-628 (2003).
- 833 62 Franck, C. M., Westermann, J. & Boisson-Dernier, A. Plant malectin-like receptor kinases: From
834 cell wall integrity to immunity and beyond. *Annu Rev Plant Biol* **69**, 301-328 (2018).
- 835 63 Innan, H. & Kondrashov, F. The evolution of gene duplications: classifying and distinguishing
836 between models. *Nat Rev Genet* **11**, 97-108 (2010).
- 837 64 Moghe, G. D. & Last, R. L. Something old, something new: Conserved enzymes and the evolution
838 of novelty in plant specialized metabolism. *Plant Physiol* **169**, 1512-1523 (2015).
- 839 65 Thompson, A., Zakon, H. H. & Kirkpatrick, M. Compensatory drift and the evolutionary dynamics
840 of dosage-sensitive duplicate genes. *Genetics* **202**, 765-774 (2016).
- 841 66 Goodstein, D. M. *et al.* Phytozome: a comparative platform for green plant genomics. *Nucleic Acids*
842 *Res* **40**, D1178-1186 (2012).

- 843 67 Deng, W., Nickle, D. C., Learn, G. H., Maust, B. & Mullins, J. I. ViroBLAST: a stand-alone BLAST
844 web server for flexible queries of multiple databases and user's datasets. *Bioinformatics* **23**, 2334-
845 2336 (2007).
- 846 68 Stajich, J. E. *et al.* FungiDB: an integrated functional genomics database for fungi. *Nucleic Acids*
847 *Res* **40**, D675-681 (2012).
- 848 69 Yin, C. *et al.* InsectBase: a resource for insect genomes and transcriptomes. *Nucleic Acids Res* **44**,
849 D801-807 (2016).
- 850 70 Eddy, S. R. Profile hidden Markov models. *Bioinformatics* **14**, 755-763 (1998).
- 851 71 Vaattovaara, A., Salojärvi, J. & Wrzaczek, M. Extraction and curation of gene models for plant
852 receptor kinases for phylogenetic analysis. *Methods Mol Biol* **1621**, 79-91 (2017).
- 853 72 Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. *Genome Res* **14**, 988-995 (2004).
- 854 73 Birney, E. & Durbin, R. Using GeneWise in the *Drosophila* annotation experiment. *Genome Res* **10**,
855 547-548 (2000).
- 856 74 Solovyev, V. in *Handbook of Statistical Genetics*. 97-159 (John Wiley & Sons, Ltd, 2008).
- 857 75 Jaillon, O. *et al.* The grapevine genome sequence suggests ancestral hexaploidization in major
858 angiosperm phyla. *Nature* **449**, 463-467 (2007).
- 859 76 Mirarab, S. *et al.* PASTA: Ultra-large multiple sequence alignment for nucleotide and amino-acid
860 sequences. *J Comput Biol* **22**, 377-386 (2015).
- 861 77 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
862 improvements in performance and usability. *Mol Biol Evol* **30**, 772-780 (2013).
- 863 78 Veidenberg, A., Medlar, A. & Löytynoja, A. Wasabi: An integrated platform for evolutionary
864 sequence analysis and data visualization. *Mol Biol Evol* **33**, 1126-1130 (2016).
- 865 79 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
866 phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).
- 867 80 Keller, O., Odronitz, F., Stanke, M., Kollmar, M. & Waack, S. Scipio: using protein sequences to
868 determine the precise exon/intron structures of genes and their orthologs in closely related species.
869 *BMC Bioinformatics* **9**, 278 (2008).
- 870 81 Li, L., Stoeckert, C. J., Jr. & Roos, D. S. OrthoMCL: identification of ortholog groups for eukaryotic
871 genomes. *Genome Res* **13**, 2178-2189 (2003).
- 872 82 Löytynoja, A. & Goldman, N. An algorithm for progressive multiple alignment of sequences with
873 insertions. *Proc Natl Acad Sci USA* **102**, 10557-10562 (2005).
- 874 83 Gremme, G., Steinbiss, S. & Kurtz, S. GenomeTools: a comprehensive software library for efficient
875 processing of structured genome annotations. *IEEE/ACM Trans Comput Biol Bioinform* **10**, 645-656
876 (2013).
- 877 84 Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**, 1586-1591
878 (2007).
- 879 85 Yang, Z. N., R. Synonymous and nonsynonymous rate variation in nuclear genes of mammals. *J Mol*
880 *Evol* **46**, 409-418 (1998).
- 881 86 Lyons, E., Pedersen, B., Kane, J. & Freeling, M. The value of nonmodel genomes and an example
882 using SynMap within CoGe to dissect the hexaploidy that predates the rosids. *Tropical Plant Biology*
883 **1**, 181-190 (2008).
- 884 87 Chamala, S. *et al.* Assembly and validation of the genome of the nonmodel basal angiosperm
885 *Amborella*. *Science* **342**, 1516-1517 (2013).
- 886 88 Sawyer S. Statistical tests for detecting gene conversion. *Mol Biol Evol* **6**, 526-538 (1989).
- 887 89 Wu, Y. An algorithm for constructing parsimonious hybridization networks with multiple
888 phylogenetic trees. *J Comput Biol* **20**, 792-804 (2013).
- 889 90 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data.
890 *Bioinformatics* **30**, 2114-2120 (2014).
- 891 91 Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq
892 quantification. *Nat Biotechnol* **34**, 525-527 (2016).
- 893 92 Sonesson, C., Love, M. I. & Robinson, M. D. Differential analyses for RNA-seq: transcript-level
894 estimates improve gene-level inferences. *F1000Res* **4**, 1521 (2015).
- 895 93 Norholm, M. H. A mutant Pfu DNA polymerase designed for advanced uracil-excision DNA
896 engineering. *BMC Biotechnol* **10**, 21 (2010).

- 897 94 Gibson, D. G. *et al.* Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat*
898 *Methods* **6**, 343-345 (2009).
- 899 95 Liu, H. & Naismith, J. H. An efficient one-step site-directed deletion, insertion, single and multiple-
900 site plasmid mutagenesis protocol. *BMC Biotechnol* **8**, 91 (2008).
- 901 96 Kabsch, W. Automatic processing of rotation diffraction data from crystals of initially unknown
902 symmetry and cell constants. *Journal of Applied Crystallography* **26**, 795-800 (1993).
- 903 97 Sheldrick, G. M. A short history of SHELX. *Acta Crystallogr A* **64**, 112-122 (2008).
- 904 98 Bricogne, G., Vonrhein, C., Flensburg, C., Schiltz, M. & Paciorek, W. Generation, representation
905 and flow of phase information in structure determination: recent developments in and around
906 SHARP 2.0. *Acta Crystallogr D Biol Crystallogr* **59**, 2023-2030 (2003).
- 907 99 Cowtan, K. The Buccaneer software for automated model building. 1. Tracing protein chains. *Acta*
908 *Crystallogr D Biol Crystallogr* **62**, 1002-1011 (2006).
- 909 100 Cohen, S. X. *et al.* ARP/wARP and molecular replacement: the next generation. *Acta Crystallogr D*
910 *Biol Crystallogr* **64**, 49-60 (2008).
- 911 101 Emsley, P. & Cowtan, K. Coot: model-building tools for molecular graphics. *Acta Crystallogr D*
912 *Biol Crystallogr* **60**, 2126-2132 (2004).
- 913 102 Murshudov, G. N., Vagin, A. A. & Dodson, E. J. Refinement of macromolecular structures by the
914 maximum-likelihood method. *Acta Crystallogr D Biol Crystallogr* **53**, 240-255 (1997).
- 915 103 McCoy, A. J. *et al.* Phaser crystallographic software. *J Appl Crystallogr* **40**, 658-674 (2007).
- 916 104 Davis, I. W. *et al.* MolProbity: all-atom contacts and structure validation for proteins and nucleic
917 acids. *Nucleic Acids Res* **35**, W375-383 (2007).
- 918 105 Pettersen, E. F. *et al.* UCSF Chimera--a visualization system for exploratory research and analysis. *J*
919 *Comput Chem* **25**, 1605-1612 (2004).
- 920 106 Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
921 *Nucleic Acids Res* **32**, 1792-1797 (2004).
- 922 107 Kabsch, W. & Sander, C. Dictionary of protein secondary structure: pattern recognition of hydrogen-
923 bonded and geometrical features. *Biopolymers* **22**, 2577-2637 (1983).

924

925 **Figure legends**

926 **Figure 1. Overview and distribution of DUF26-containing genes in plants.** a) DUF26-containing genes
927 are absent from algae and charophytes but present in land plants. *Marchantia polymorpha* and
928 *Physcomitrella patens* genomes encode sdCRRSPs. *Selaginella moellendorffii* possesses sdCRRSPs,
929 sdCRKs and canonical CRKs. Seed plant (gymnosperm and angiosperm) genomes encode the whole set of
930 DUF26-containing genes. CRKs were defined as basal group CRKs (bCRKs) or variable group CRKs
931 (vCRKs) based on their phylogenetic positions. Whole genome duplication (WGD) events are presented with
932 green circle and whole genome triplication (WGT) events with dark blue circle. Ferns were omitted from
933 analyses due to lack of available genome assemblies. b) Overview of different domain compositions of
934 proteins containing DUF26 in different plant lineages. The number of representative species in the analyses
935 is given in brackets after the name of the group. Numbers in the table present the number of species in each
936 lineage in which the domain structure was found. In abbreviations sd (single domain), dd (double domain), td
937 (triple domain) and qd (quadruple domain) refers to the number of the DUF26 domains.

938 **Figure 2. Phylogenetic tree of CRRSPs, CRKs and PDLPs.** a) The phylogenetic tree was estimated with
939 the maximum-likelihood method using all high quality full-length DUF26-containing sequences from
940 lycophytes onwards. CRCKs and concA-CRKs were excluded while GNK2 from *Ginkgo biloba* was
941 included. Overall, DUF26-containing genes split into basal and variable group. Detailed phylogenetic trees
942 with bootstrap support (1000 replicates) and filtered sequence alignments are available at
943 <http://was.bi?id=IaroPa> (full tree), <http://was.bi?id=wpEHGt> (basal group separately) and
944 http://was.bi?id=aJJe_D (variable group separately). b) The same phylogenetic tree as in panel a rooted to
945 ancestral sdCRRSPs and sdCRKs from *Selaginella moellendorffii* showing that the variable group branches
946 out from the basal group. c) The MEME figures present the conservation pattern of amino acid positions
947 around the main cysteine motif within the DUF26 domains for sdCRRSPs, bCRKs and PDLPs from the
948 basal group and CRRSPs and vCRKs from the variable group. The features specific only to genes either in
949 the basal group or in the variable group are highlighted. d) The DUF26-A and DUF26-B domains are clearly
950 separated in an unrooted phylogenetic tree containing DUF26 domain sequences. The MEME figures present

951 differences in the conservation of the AA sequence surrounding the conserved cysteines in DUF26-A and
952 DUF26-B.

953 **Figure 3. Comparison of evolutionary rates between gene families.** Analyses were carried out with
954 Badirate for eleven species (*Physcomitrella patens*, *Selaginella moellendorffii*, *Amborella trichopoda*,
955 *Arabidopsis thaliana*, *Populus trichocarpa*, *Vitis vinifera*, *Aquilegia coerulea*, *Spirodela polyrhiza*, *Zea*
956 *mays*, *Oryza sativa* and *Brachypodium distachyon*). Neutral branches are reported as bold black lines;
957 branches involving gene family expansion are reported as bold purple lines and branches with contraction as
958 blue dashed lines. Branches with a significant differences (false discovery rate adjusted $p < 0.05$) to birth-
959 death rate model estimates are marked with arrows. Node labels present the ancestral gene family sizes
960 estimated by Badirate. Tip labels contain species abbreviations and the change in numbers compared to the
961 most recent ancestral node. **a)** All CRKs compared to other receptor like kinases (RLKs). **b)** bCRKs
962 compared to RLKs. **c)** PDLPs compared to other plasmodesmata related orthogroups. **d)** vCRKs compared to
963 RLKs. **e)** Phylogenetic maximum-likelihood tree showing differences in lineage specific expansions in
964 monocot and dicot vCRKs following the split of *Amborella trichopoda*. Species-specific expansions (at least
965 two genes from same species) are marked with red and clades including sequences from only *Brassicaceae*
966 or *Solanaceae* are marked with blue.

967 **Figure 4: The crystals structures of the PDLP5 and PDLP8 ectodomains reveal a conserved tandem**
968 **architecture of two lectin-like domains.** **a)** Overview of the PDLP5 ectodomain. The two DUF26 domains
969 are shown as ribbon diagrams, colored in blue (DUF26-A) and orange (DUF26-B), respectively. N-glycans
970 are located at Asn69 and Asn132 of DUF26-A and are depicted in bonds representation (in cyan). The
971 DUF26-A and DUF26-B domains each contain 3 disulfide bridges labeled 1 (Cys89-Cys98), 2 (Cys101-
972 Cys126), 3 (Cys36-Cys113), 4 (Cys191-Cys200), 5 (Cys203-Cys228) and 6 (Cys148-Cys215). **b)** Close-up
973 view of the DUF26-A – DUF26-B interface in PDLP5 (orange) and PDLP8 (blue), shown in bonds
974 representation. **c)** Superimposition of the GnK2 extracellular DUF26 domain (PDB-ID 4XRE) with either
975 PDLP5 DUF-26A (r.m.s.d. is ~ 1.4 Å comparing 100 aligned C_{α} atoms) or PDLP5 DUF26-B (r.m.s.d. is ~ 2.0
976 Å comparing 93 corresponding C_{α} atoms). Corresponding disulfide bridges shown in bonds representation
977 (PDLP5 in green, GnK2 in yellow) are highlighted in grey. GnK2-bound mannose is shown in magenta (in

978 bonds representation). **d)** Close-up view of the residues involved in the binding of mannose of Gnk2 (bonds
979 representation, in blue and magenta, respectively) and putative residues involved in substrate binding of
980 PDLP5 DUF26-A (in orange). **e)** The fungal LDL DUF26 domain (C_{α} trace, in blue; PDB-ID 4NDV) and
981 PDLP5 DUF26-A (in orange) superimposed with an r.m.s.d. of ~ 2.4 Å comparing 75 aligned C_{α} atoms).
982 Disulfide bridges (LDL in yellow and PDLP5 in green; aligned disulfide bridges highlighted in grey) and the
983 LDL bound globotriose (magenta) are shown in bonds representation. **f)** C_{α} traces of the structural
984 superimposition of the fungal Y3 protein (PDB-ID 5V6I) and PDLP5 DUF26-A (r.m.s.d. is ~ 2.6 Å
985 comparing 78 corresponding C_{α} atoms). Disulfide bridges of Y3 (yellow) and PDLP5 DUF26-A (green) are
986 shown alongside, one corresponding disulfide pair is highlighted in gray.

987 **Figure 5: PDLP5 and PDLP8 may have drastically different oligomerisation modes, surface charge**
988 **distributions and surface exposed residues are not widely conserved.** **a)** The conservation of amino acid
989 residues illustrated on the molecular surface of the PDLP5 or PDPL8 crystallization dimers, respectively.
990 Site-wise ω (dN/dS) values, indicating the intensity and direction of selection on amino acid changing
991 mutations, illustrated on the molecular surfaces (upper) and in ribbon diagrams (lower) of PDLP5 or PDPL8.
992 The ω values range from 0.15 (green) to slightly over 1.0 (magenta), reflecting conserved sites under
993 purifying selection and sites evolving close a neutral process, respectively. **b)** Electrostatic potential mapped
994 onto molecular surfaces of the putative PDLP5 and PDLP8, orientation as in c) dimer, respectively. **c)**
995 Ribbon diagrams of PDLP5 (orange) and PDLP8 (blue) crystallographic dimers. In both dimers large,
996 antiparallel β -sheets are formed, using different protein-protein interaction surfaces.

997 **Figure 6. CRKs experienced domain rearrangements.** **a)** Comparison of phylogenetic trees based on
998 ectodomain region and kinase domain of 880 CRKs. Phylogenetic maximum-likelihood trees are presented
999 as tanglegram where the tree of the CRK ectodomain region is plotted against the tree of the kinase domain.
1000 The kinase tree is rooted to atypical monocot CRKs with a Concanavalin-A type kinase domain and the
1001 ectodomain tree is rooted to CRKs from *Selaginella moellendorffii*. The ectodomain tree was detangled
1002 based on the kinase domain tree. Lines connect the ectodomain and kinase domain belonging to same gene,
1003 and connection are drawn in different colors for better visibility. Juxtaposition of the trees shows
1004 rearrangements and domain swaps of ecto- and kinase domains. Black circles highlight the difference

1005 between the ectodomains and kinase domains of the *Selaginella* sdCRKs and ddCRKs and also the group of
1006 the atypical monocot CRKs which have exchanged the kinase domain. **b)** The exon-intron structure of the
1007 CRKs. Usually CRKs contain seven exons: one encoding DUF26 domains, one encoding transmembrane
1008 region (TMR) and five exons encoding the kinase domain. In atypical monocot CRKs with exchanged kinase
1009 domain, whole gene is encoded by one or two exons. The scale bar for each gene represents 100 bases.
1010 Regions encoding the DUF26-A are colored with blue, the DUF26-B with orange, the transmembrane region
1011 (TMR) with pink and the kinase domain with green.

1012 **Figure 7. Identification of the modes of gene family evolution in DUF26-containing genes in**
1013 ***Arabidopsis thaliana*, tomato, rice, maize and *Amborella trichopoda*.** **a)** Gene families that are
1014 preferentially retained after whole genome multiplications (WGMs) are typically identified by synteny
1015 analysis. The figure illustrates syntenic regions containing DUF26 genes from *Amborella trichopoda* to
1016 monocots *Oryza sativa* and *Zea mays* (to left from middle) and to eudicots *Solanum lycopersicum* and
1017 *Arabidopsis thaliana* (right from the middle). In the synteny analysis within monocots and dicots, segments
1018 with at least 5 syntenic genes were included, whereas in comparisons to *Amborella* the minimum threshold
1019 was 3 syntenic genes. Analyses were carried out with Synmap software within CoGe. For *Amborella*
1020 *trichopoda* genomic locations of DUF26-containing genes are only known on chromosome/scaffold level
1021 based on physical mapping. **b and c)** Gene families with a preferential retention pattern after WGMs show
1022 conserved gene counts over species. Phylogenetic tree of the five species shown in the panel was used to
1023 reconcile the gene trees and estimate gene counts in ancestral nodes for **b)** bCRKs and **c)** vCRKs, using
1024 *Selaginella moellendorffii* as outgroup. The gains are highlighted with red and losses with blue. **d)** Gene
1025 families with preferential retention pattern should have many orthologs. Heatmaps of the normalized
1026 transcriptional expression counts (Transcript per million [TPM]) of candidate DUF26 orthologs from four of
1027 the species: *Solanum lycopersicum*, *Arabidopsis thaliana*, *Zea mays*, and *Oryza sativa*. Coloring in heatmaps
1028 is proportional to \log_2 (TPM) value that represents the gene expression level. The corresponding \log_2 (TPM)
1029 value is displayed next to the color key. The rows represent gene models and the columns show the
1030 experiments, collected from publicly available Sequence Read Archive (SRA) database. SRA accessions are
1031 annotated to relevant stress conditions (descriptions are presented in Table S4). Solid lines connect putative

1032 orthologs based on evidence from phylogenetic and synteny analyses; dashed lines connect putative
1033 orthologs based on evidence from either phylogenetic or synteny analyses. **e)** Final prediction of gene
1034 families evolving under dosage balance is that their knockouts demonstrate a high phenotypic effect. This
1035 can be seen by reanalysis of phenotype data from (Bourdais *et al.*²⁴); the bCRK T-DNA insertion mutants
1036 display a significantly larger standard deviation (Y-axis) over different phenotyping experiments than vCRK
1037 mutants.

1038 Pathogens: *Agrobacterium tumefaciens*, *Alternaria brassicicola*, *Botrytis cinerea*, *Cercospora zeina*,
1039 *Cladosporium fulvum*, *Colleotrichum graminicola*, *Magnaporthe grisei*, *Pseudomonas putida*, *Pseudomonas*
1040 *fluorescens*, *Pseudomonas syringae* pv. tomato DC3000, *Rizoctonia solani*, *Ustilago maydis*, *Xanthomonas*
1041 *oryzae*.

1042 **Figure 8. Model of mixed-type gene family evolution.** Gene families evolve through two major events,
1043 whole genome multiplications (WGM) and small-scale duplications (SSD). Genes related to environmental
1044 responses and secondary metabolism experience SSDs in the form of tandems, whereas highly connected
1045 genes associated with transcriptional and developmental regulation or signal transduction functions are
1046 preferentially retained after WGMs. **a)** Prevailing hypothesis for the retention pattern is dosage-balance; in
1047 case of highly connected genes the stoichiometric balance needs to be maintained, and therefore selection
1048 acts against gene losses after WGMs and against duplications by SSDs. **b)** On the other hand, gene family
1049 evolving through tandem duplications (**b**; evolution before the speciation node) has a high birth rate and
1050 therefore the number of duplicates between species can vary. After duplications the homogeneity of the
1051 duplicates is maintained through gene conversion events, which has a high probability with near-by
1052 homologous sequences. This can be maintained for long periods, but eventually over time the sequences
1053 diverge by drift and selection based on dosage. Our data suggests that a tandemly expanding gene family
1054 may evolve into a dosage balance mode as a result of WGMs (**b**; evolution after speciation node). Following
1055 WGMs, the duplicated tandems may experience extensive fractionation due to drift and selection by dosage
1056 which fragments the tandem structure. At the same time, the connectivity of the gene family has been
1057 accumulating through sub- and neofunctionalization, increasing pressure for retention of the gene models.
1058 These phenomena together may result into a dosage balance model of evolution (top branch after speciation

1059 node). This does not necessarily occur across all WGM events and depends on the tandem duplication rate,
1060 as was observed for bCRKs in Solanaceae (bottom branch), where there exist both single copies and a later
1061 tandem expansion in the genome. Different subfamilies can be in different states of this process. c) CRRSPs
1062 and PDLPs follow dosage balance mode after the paleohexaploid event, whereas bCRKs have assumed the
1063 mode in later WGM events. The overall numbers of the bCRKs are preserved but identification of orthologs
1064 between species that have experienced independent WGMs is difficult, suggesting that convergent
1065 functionality of the members is recent. Gene families expanding through tandem duplications such as vCRKs
1066 and CRRSPs have high birthrate and demonstrate several lineage-specific expansions.

1067

1068 **Supplementary figure legends**

1069 **Figure S1. Summary of manual gene annotation and correction.** a) The number of corrected, manually
1070 annotated and partial/pseudo gene models in the studied species. Percentage of corrected gene models is
1071 marked with light gray, manually annotated genes with black and genes classified as partial or pseudogenes
1072 with dark gray. Silver birch (*Betula pendula*) and sacred lotus (*Nelumbo nucifera*) genes were fully manually
1073 annotated, as the gene models were not available when the study was initiated. *Selaginella moellendorffii* and
1074 *Vitis vinifera* required highest percentage of manual corrections. The high percentage of pseudogenes in
1075 *Physcomitrella patens* is explained by low gene number (two out of three gene models are likely
1076 pseudogenes). b) Average exon numbers of CRRSPs, PDLPs and CRKs. Average exon numbers were
1077 calculated for sdCRRSPs, ddCRRSPs, PDLPs and CRKs in *Amborella trichopoda*, *Arabidopsis thaliana* and
1078 *Oryza sativa*. c) The amount of curated and manually annotated gene models in basal and variable groups.
1079 Corrected (red) and manually annotated (green: species with pre-existing annotations; blue: species without
1080 previous annotations) gene models marked in both groups. Corrected or annotated genes can be found in all
1081 subgroups within these groups. There are several examples of corrected or previously non-annotated genes
1082 that are basal for subgroups, indicating the importance of gene model validation for correct tree topology.

1083 **Figure S2. Phylogenies of DUF26-containing proteins.** a) A phylogenetic maximum-likelihood tree was
1084 estimated with full-length sequences for the basal group containing *Selaginella* sdCRRSPs and CRKs,
1085 Norway spruce CRRSPs and CRKs, monocot and eudicot bCRKs and PDLPs. Detailed phylogenetic trees

1086 with bootstrap support (1000 replicates) and filtered sequence alignment can be found at
1087 <http://was.bi?id=wpEHGt>. **b)** The phylogenetic maximum-likelihood tree for the variable group contains
1088 angiosperm CRRSPs and vCRKs. Tree was estimated using the full-length sequences. Detailed phylogenetic
1089 trees with bootstrap support (1000 replicates) and filtered sequence alignment can be found at
1090 http://was.bi?id=aJJe_D. Phylogenetic maximum likelihood trees of **c)** CRRSPs **d)** CRKs and **e)** PDLPs.
1091 Detailed phylogenetic trees containing gene identifiers as well as bootstrap support (1000 replicates) and
1092 filtered sequence alignment can be found at <http://was.bi?id=zbH7i> (CRRSPs), <http://was.bi?id=i9To8q>
1093 (CRKs) and <http://was.bi?id=Fe1A3A> (PDLPs). **f)** Phylogenetic maximum-likelihood tree of all DUF26
1094 genes in *Marchantia polymorpha*, *Selaginella moellendorffii* and *Amborella trichopoda*. Tree is estimated
1095 from sequence alignment of full length gene models where the sites with coverage less than 10% have been
1096 filtered out. Tree is rooted to sdCRRSPs from *Marchantia polymorpha*. A detailed phylogenetic tree with
1097 gene identifiers as well as bootstrap support (1000 replicates) and filtered sequence alignment can be found
1098 at <http://was.bi?id=VeeQZ6>.

1099 **Figure S3. Ancestral gene counts for DUF26-containing genes.** DLCpar was used for inferring the most
1100 parsimonious history of protein groups in the presence of duplications, losses, and incomplete lineage
1101 sorting. The panels illustrate ancestral gene counts and lineage-specific expansions in **a)** sdCRRSPs in the
1102 basal group, **b)** basal CRKs, **c)** PDLPs, **d)** variable group CRKs, and **e)** ddCRRSPs in the variable group.
1103 Numbers with black color show the gene counts in the species and their most recent common ancestor.
1104 Estimated gene gains are marked with red and losses with blue.

1105 **Figure S4. Badirate comparisons for evolutionary rates.** Analyses were carried out with Badirate for
1106 eleven species (*Physcomitrella patens*, *Selaginella moellendorffii*, *Amborella trichopoda*, *Arabidopsis*
1107 *thaliana*, *Populus trichocarpa*, *Vitis vinifera*, *Aquilegia coerulea*, *Spirodela polyrhiza*, *Zea mays*, *Oryza*
1108 *sativa* and *Brachypodium distachyon*). Neutral branches: bold black lines; gene family expansion: bold
1109 purple lines; gene family contraction: blue dashed lines. Branches with a significant difference to birth-death
1110 model estimated from orthogroup data are marked with arrows. Node labels present the gene family size in
1111 ancestral nodes as estimated by Badirate. Tip labels contain species abbreviation and the change in number

1112 compared to the most recent ancestral node. **a)** All CRKs compared to all kinases. **b)** bCRKs compared to all
1113 kinases. **c)** vCRKs compared to all kinases.

1114 **Figure S5. Phylogenetic maximum-likelihood tree of bCRKs.** The full length sequences belonging to this
1115 clade were re-aligned and the alignment was filtered to exclude sites with less than 10% coverage. Bootstrap
1116 support is calculated with 1000 replicates. A detailed phylogenetic tree and filtered sequence alignment can
1117 be found at <http://was.bi?id=6Z7yhQ>.

1118 **Figure S6. Species trees and reconciled phylogenetic trees for DCLpar analyses.** **a)** Species tree for the
1119 24 species where all DUF26-domain genes were comprehensively analyzed. The tree was downloaded from
1120 PhyloT. The node labels indicate the speciation event IDs that are used in panels b and c. **b)** Reconciled gene
1121 tree for the bCRKs from DCLpar. The node labels provide the timing of the event by referring to the
1122 speciation event ID in the species tree. **c)** Reconciled gene tree for the variable group CRRSPs from DLCpar.
1123 The node labels provide the timing of the event by referring to the speciation event ID in the species tree.

1124 **Figure S7. Phylogenetic maximum-likelihood tree of 5 species used in segmental duplication analyses**
1125 **and *Selaginella moellendorffii* as outgroup.** The tree includes DUF26 genes from *Amborella trichopoda*,
1126 *Solanum lycopersicum*, *Arabidopsis thaliana*, *Oryza sativa*, *Zea mays* and *Selaginella moellendorffii*. The
1127 full length gene models were used for the sequence alignment and the sites with less than 10% coverage
1128 were filtered out. Bootstrap support is calculated with 1000 replicates. A detailed phylogenetic tree and
1129 filtered sequence alignment can be found at <http://was.bi?id=2NeJCb>.

1130 **Figure S8. Phylogenetic maximum-likelihood tree of PDLPs with possible partial PDLP from *Marsilea***
1131 ***quadrifolia*.** The phylogenetic tree is based on the sequence covering the part of ectodomain that is present in
1132 the partial gene model from *Marsilea quadrifolia*. The ddCRKs from *Selaginella moellendorffii*, *Picea*
1133 *abies* and *Amborella trichopoda* were used as outgroup for PDLPs. The partial gene model from fern
1134 *Marsilea quadrifolia* is placed close to the root of PDLP clade and thus could be a PDLP. Bootstrap support
1135 is calculated with 1000 replicates. A detailed phylogenetic tree and filtered sequence alignment can be found
1136 at <http://was.bi?id=usJEbx>.

1137 **Figure S9: Mutation of disulfide bridge-forming cysteines in PDLP5 results in protein aggregation.**
1138 PDLP5, PDLP5^{C101A}, PDLP5^{C148A} and PDLP5^{C191A} ectodomains were subjected to preparative size exclusion
1139 chromatography (left). Non-aggregated fractions were combined and subjected to analytical size exclusion
1140 chromatography (right). Molecular mass standards: A = Thyroglobulin, 669 kDa B = Aldolase, 158 kDa; C =
1141 Conalbumin, 75 kDa; D = Ovalbumin, 44 kDa; E = Ribonuclease A, 13.7 kDa.

1142 **Figure S10: Mutations in disulfide bridge forming residues in PDLP5 result in lower protein stability:**
1143 Melting curves (4 replicates in green, brown, red and blue) of PDLP5, PDLP5^{C101A}, PDLP5^{C148A}, PDLP5^{C191A}
1144 ectodomains and of the blank without protein (blank measurements for PDLP5, PDLP5^{C101A}, PDLP5^{C148A} are
1145 the same as the experiments were carried out together). For PDLP5, PDLP5^{C101A}, PDLP5^{C148A} ectodomains
1146 average melting temperatures are given +/- SDM (n=4). PDLP5^{C191A} was unstable at the given conditions and
1147 no melting curve could be acquired.

1148 **Figure S11: Structural comparisons of PDLP5 and PDLP8 DUF26 domains reveal a high degree of**
1149 **structural similarity (a)** Superimposition of the DUF26-A (orange; C_α trace) and the DUF26-B (blue; C_α
1150 trace) domains of PDLP5 (left; r.m.s.d. is ~1.6 Å comparing 89 corresponding C_α atoms) and PDLP8 (right;
1151 r.m.s.d. is ~1.2 Å comparing 89 corresponding C_α atoms) demonstrate the structural similarity of DUF26-A
1152 and DUF26-B domains. Glycosylated asparagines are indicated by an arrow **(b)** Structural superposition of
1153 PDLP5 (orange, shown as C_α trace) and PDLP8 (blue) reveals a high degree of overall structural similarity
1154 (r.m.s.d. is ~1.6 Å comparing 198 corresponding C_α atoms), and a conserved pattern of disulfide bridges
1155 (grey highlights). The disulfide bridges in PDLP8 are: 1 (Cys89-Cys98), 2 (Cys101-Cys126), 3 (Cys34-
1156 Cys113), 4 (Cys191- Cys200), 5 (Cys203-Cys228) and 6 (Cys148-Cys215). Disulfide bridges are depicted in
1157 bonds representation (PDLP5 in yellow, PDLP8 in green).

1158 **Figure S12: Cysteines forming disulfide bonds and residues involved in the interaction of DUF26-A**
1159 **and DUF26-B domains are conserved in bCRKs, vCRKs, CRSPs and PDLPs.** A set of PDLPs, bCRKs,
1160 vCRKs and CRRSPs were selected based on the structure and their sequences were aligned with
1161 MUSCLE¹⁰⁶. The result shows the conservation of amino acids present in the interaction patch of DUF26-A
1162 and DUF26-B in either PDLP5s (red highlight) or all double DUF26 containing proteins (highlight in blue).

1163 Cysteines and disulfide bridges are highlighted in yellow. A secondary structure assignment of the DUF26-A
1164 (blue) and DUF26-B domains¹⁰⁷ is given above the sequences.

1165 **Figure S13: The PDLP5 ectodomain does not bind mannose or other cell wall derived sugars. a)**
1166 Mannose was titrated into a cell containing the PDLP5 ectodomain in an isothermal titration calorimetry
1167 (ITC) assay (n.d., no binding detected). **b)** ITC experiments were carried out to test binding of plant cell wall
1168 sugars to the isolated PDLP5 ectodomain.

1169 **Figure S14. Gene tree reconciliation of the five species used in segmental duplication analyses. (a)**
1170 Species tree of the five species in the analyses. The phylogeny was downloaded from PhyloT, and the node
1171 labels indicate the speciation event ID. These IDs are used in Figures S15b-d). The reconciled gene trees
1172 were estimated with DLCpar for **(b)** bCRKs, **(c)** vCRKs, **(d)** ddCRRSPs. The node labels provide the timing
1173 of the split by referring to the speciation event ID in the species tree (Figure S15a). *Selaginella*
1174 *moellendorffii* was used as outgroup.

1175 **Figure S15. The DUF26 genes show transcriptional response to several stress treatments.** Heatmap
1176 illustrating transcriptional response of DUF26 genes from *Arabidopsis thaliana*, *Oryza sativa*, *Zea*
1177 *mays* and *Solanum lycopersicum*. The dendrogram shows a phylogenetic tree of the 253 DUF26-containing
1178 genes (rows) in the four species. The columns represent the RNAseq experiments from Sequence Read
1179 Archive (see Table S4; accession numbers not shown here for clarity), categorized into pathogen defence
1180 (red highlight) and miscellaneous (blue). The heatmap colors represent the log₂(TPM) values, as illustrated
1181 by the color key. The NA values are displayed with white color.

1182 **Supplementary file 1: Sequences used for analyses in fasta format**

1183 **Supplementary file 2: wwPDB X-ray Structure Validation Summary Report 6GRE (PDLP5)**

1184 **Supplementary file 3: wwPDB X-ray Structure Validation Summary Report 6GRF (PDLP8)**

1185 **Table S1: Information of DUF26 proteins included in this study.** Information of DUF26 protein
1186 sequences found in study species.

1187 **Table S2: Data collection, phasing and refinement statistics for structural analyses**

1188 **Table S3: The related kinase domains for the CRK kinase domains.** PBLAST results for selected CRKs
1189 from *Amborella trichopoda*, *Arabidopsis thaliana*, *Oryza sativa* and *Selaginella moellendorffii*. Only amino
1190 acid sequence of the kinase domain of each CRK was used as query. Best hit outside the CRKs was marked
1191 in the table.

1192 **Table S4: Gene conversion analyses results.**

1193 **Table S5: Identified orthologs and information of transcriptome data used in analyses.**

1194 **Table S6: Genome version information and references for plant genomes used in phylogenetic**
1195 **analyses.**

1196

1197

1198

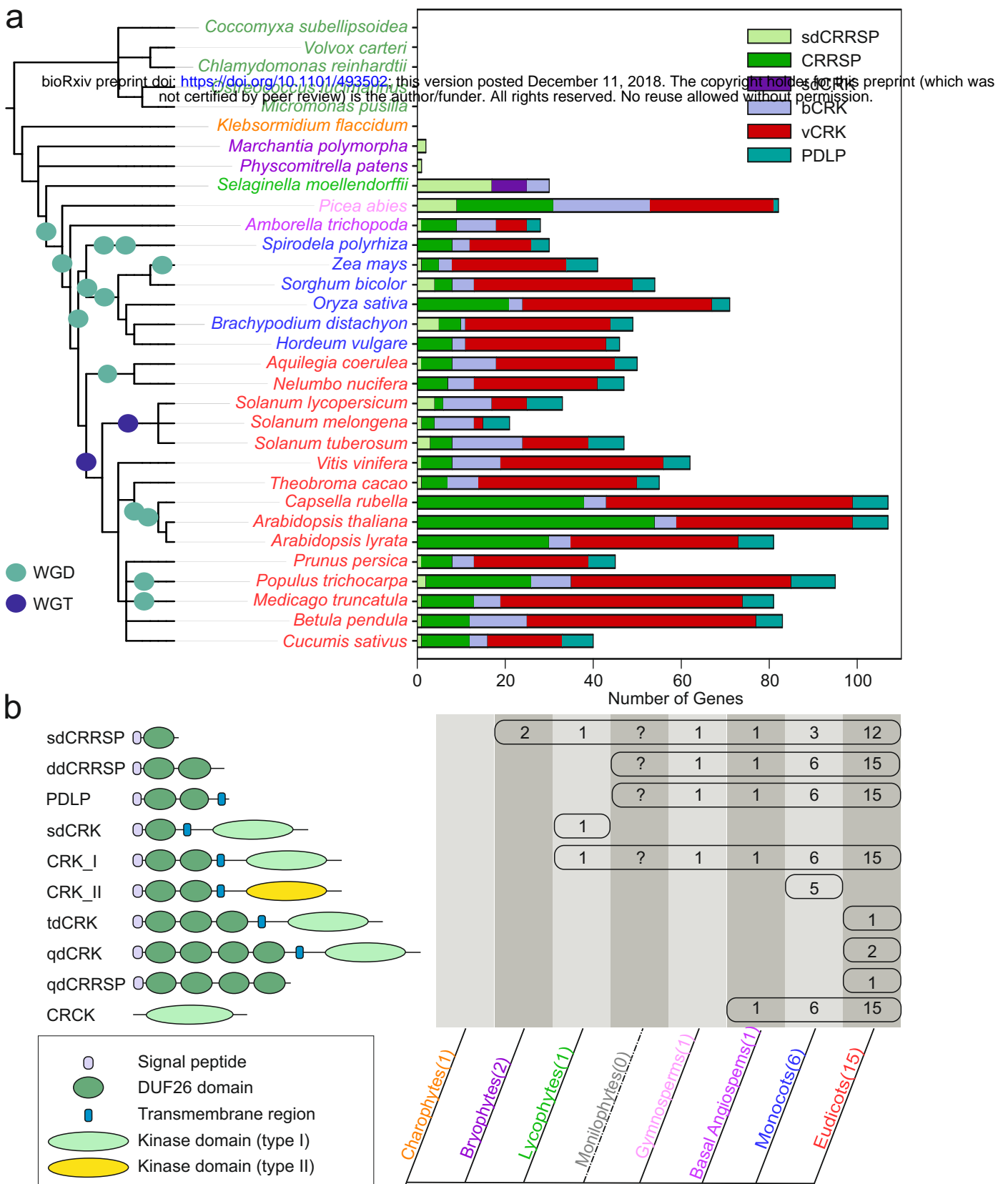


Figure 1. Overview and distribution of DUF26-containing genes in plants. a DUF26-containing genes are absent from algae and charophytes but present in land plants. *Marchantia polymorpha* and *Physcomitrella patens* genomes encode sdCRRSPs. *Selaginella moellendorffii* possesses sdCRRSPs, sdCRKs and canonical CRKs. Seed plant (gymnosperm and angiosperm) genomes encode the whole set of DUF26-containing genes. CRKs were defined as basal group CRKs (bCRKs) or variable group CRKs (vCRKs) based on their phylogenetic positions. Whole genome duplication (WGD) events are presented with green circle and whole genome triplication (WGT) events with dark blue circle. Ferns were omitted from analyses due to lack of available genome assemblies. **b** Overview of different domain compositions of proteins containing DUF26 in different plant lineages. The number of representative species in the analyses is given in brackets after the name of the group. Numbers in the table present the number of species in each lineage in which the domain structure was found. In abbreviations sd (single domain), dd (double domain), td (triple domain) and qd (quadruple domain) refers to the number of the DUF26 domains.

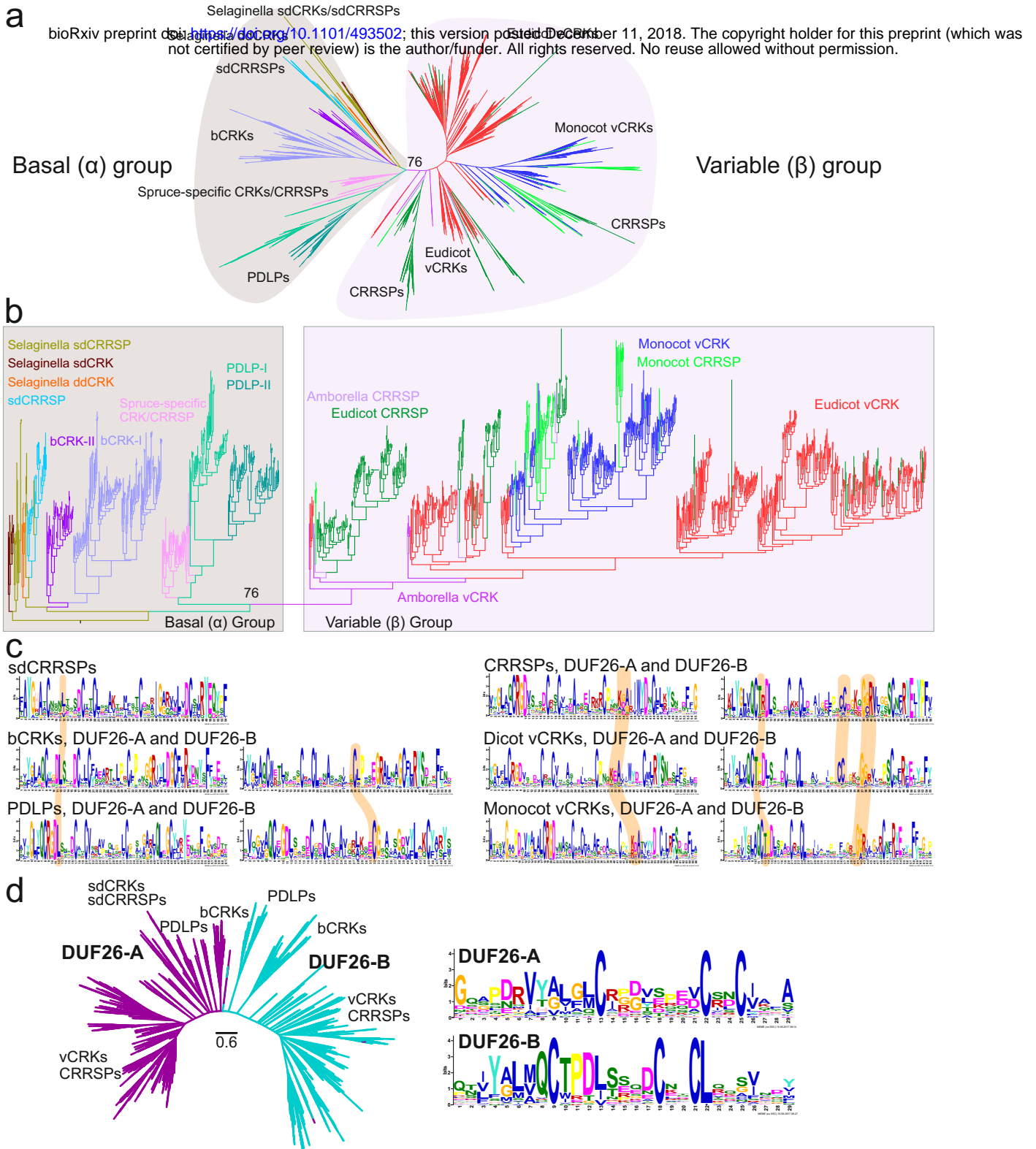


Figure 2. Phylogenetic tree of CRRSPs, CRKs and PDLPs. **a**) The phylogenetic tree was estimated with the maximum-likelihood method using all high quality full-length DUF26-containing sequences from lycophytes onwards. CRCKs and concA-CRCKs were excluded while GNK2 from *Ginkgo biloba* was included. Overall, DUF26-containing genes split into basal and variable group. Detailed phylogenetic trees with bootstrap support (1000 replicates) and filtered sequence alignments can be found at <http://was.bi?id=IaroPa> (full tree), <http://was.bi?id=wpEHGt> (basal group separately) and http://was.bi?id=aIJe_D (variable group separately). **b**) The same phylogenetic tree rooted to ancestral sdCRRSPs and sdCRKs from *Selaginella moellendorffii* showing that the variable group branches out from the basal group. **c**) The MEME figures present the conservation pattern of amino acid positions around to the main cysteine motif within the DUF26 domains for sdCRRSPs, bCRKs and PDLPs from the basal group and CRRSPs and vCRKs from the variable group. The features specific only to genes either in the basal group or in the variable group are highlighted. **d**) The DUF26-A and DUF26-B domains are clearly separated in an unrooted phylogenetic tree containing DUF26 domain sequences. The MEME figures present differences in the conservation of the AA sequence surrounding the conserved cysteines in DUF26-A and DUF26-B.

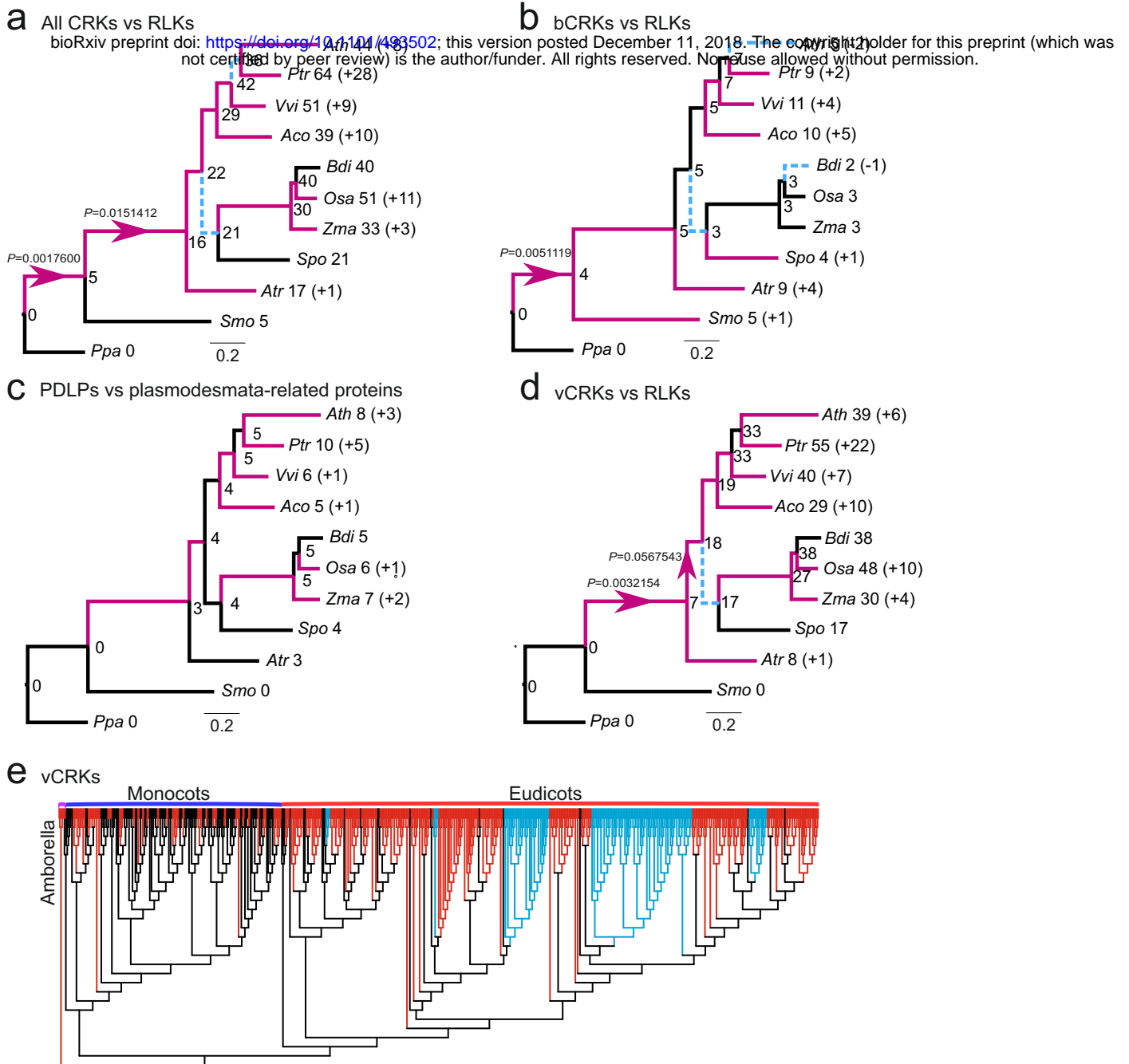


Figure 3. Comparison of evolutionary rates between gene families. Analyses were carried out with Badirate for eleven species (*Physcomitrella patens*, *Selaginella moellendorffii*, *Amborella trichopoda*, *Arabidopsis thaliana*, *Populus trichocarpa*, *Vitis vinifera*, *Aquilegia coerulea*, *Spirodela polyrhiza*, *Zea mays*, *Oryza sativa* and *Brachypodium distachyon*). Neutral branches are reported as bold black lines; branches involving gene family expansion are reported as bold purple lines and branches with contraction as blue dashed lines. Branches with a significant difference to birth-death rate model estimates are marked with arrows. Node labels present the gene family sizes in ancestors as estimated by Badirate. Tip labels contain species abbreviation and the change in numbers compared to the most recent ancestral node. **a**) All CRKs compared to other receptor like kinases (RLKs). **b**) bCRKs compared to RLKs. **c**) PDLPs compared to other plasmodesmata related orthogroups. **d**) vCRKs compared to RLKs. **e**) Phylogenetic maximum-likelihood tree showing differences in lineage specific expansions in monocot and dicot vCRKs following the split of *Amborella trichopoda*. Species-specific expansions (at least two genes from same species) are marked with red and clades including sequences from only Brassicaceae or Solanaceae are marked with blue.

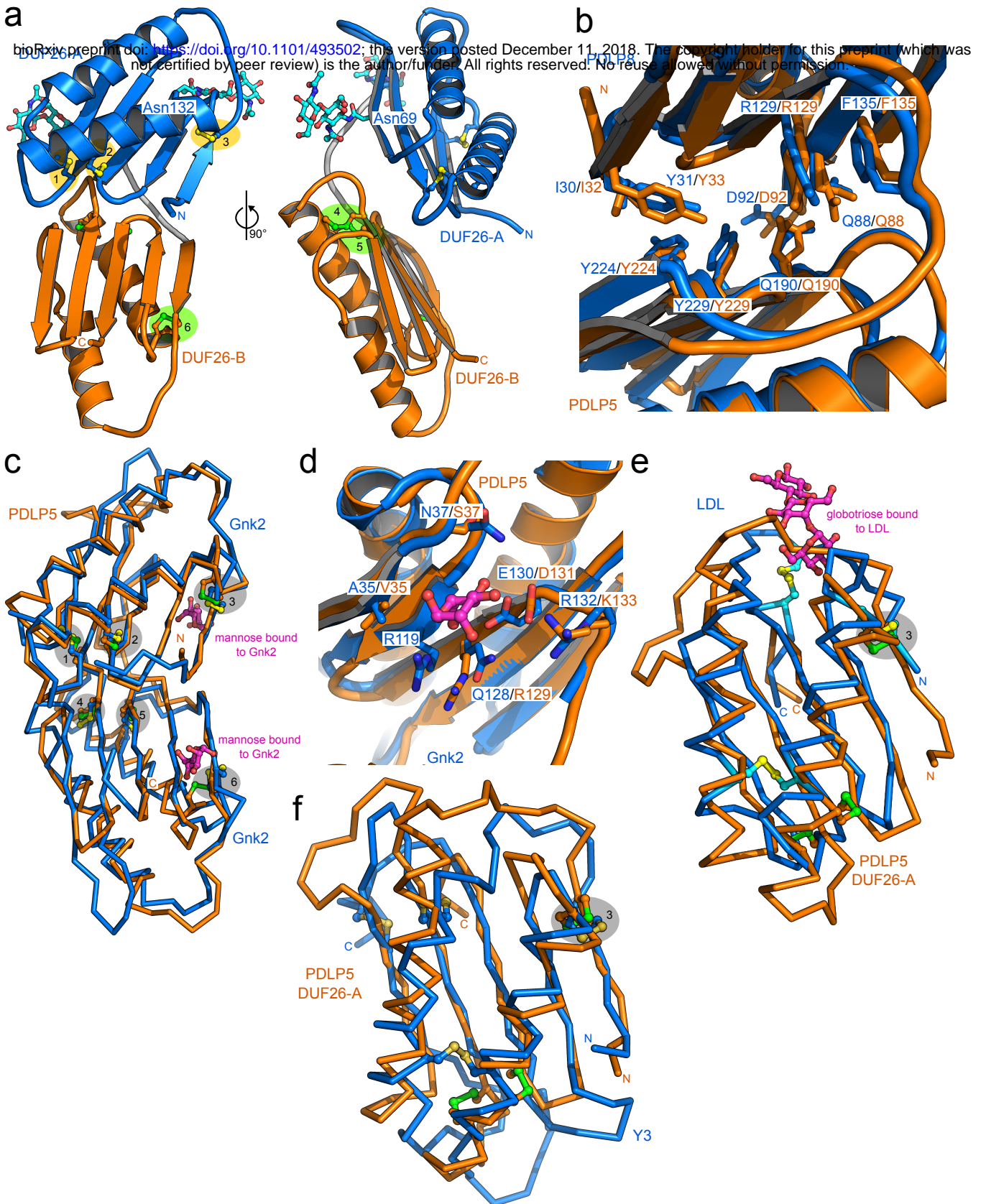


Figure 4: The crystals structures of the PDLP5 and PDLP8 ectodomains reveal a conserved tandem architecture of two lectin-like domains.

a) Overview of the PDLP5 ectodomain. The two DUF26 domains are shown as ribbon diagrams, colored in blue (DUF26-A) and orange (DUF26-B), respectively. N-glycans are located at Asn69 and Asn132 of DUF26-A and are depicted in bonds representation (in cyan). The DUF26-A and DUF26-B domains each contain 3 disulfide bridges labeled 1 (Cys89-Cys98), 2 (Cys101-Cys126), 3 (Cys36-Cys113), 4 (Cys191-Cys200), 5 (Cys203-Cys228) and 6 (Cys148-Cys215). **b)** Close-up view of the DUF26-A – DUF26-B interface in PDLP5 (orange) and PDLP8 (blue), shown in bonds representation. **c)** Superimposition of the Gnk2 extracellular DUF26 domain (PDB-ID 4XRE) with either PDLP5 DUF-26A (r.m.s.d. is ~1.4 Å comparing 100 aligned C_α atoms) or PDLP5 DUF26-B (r.m.s.d. is ~2.0 Å comparing 93 corresponding C_α atoms). Corresponding disulfide bridges shown in bonds representation (PDLP5 in green, Gnk2 in yellow) are highlighted in grey. Gnk2-bound mannose is shown in magenta (in bonds representation). **d)** Close-up view of the residues involved in the binding of mannose of Gnk2 (bonds representation, in blue and magenta, respectively) and putative residues involved in substrate binding of PDLP5 DUF26-A (in orange). **e)** The fungal LDL DUF26 domain (C_α trace, in blue; PDB-ID 4NDV) and PDLP5 DUF26-A (in orange) superimposed with an r.m.s.d. of ~2.4 Å comparing 75 aligned C_α atoms). Disulfide bridges (LDL in yellow and PDLP5 in green; aligned disulfide bridges highlighted in grey) and the LDL bound globotriose (magenta) are shown in bonds representation. **f)** C_α traces of the structural superimposition of the fungal Y3 protein (PDB-ID 5V6I) and PDLP5 DUF26-A (r.m.s.d. is ~2.6 Å comparing 78 corresponding C_α atoms). Disulfide bridges of Y3 (yellow) and PDLP5 DUF26-A (green) are shown alongside, one corresponding disulfide pair is highlighted in grey.

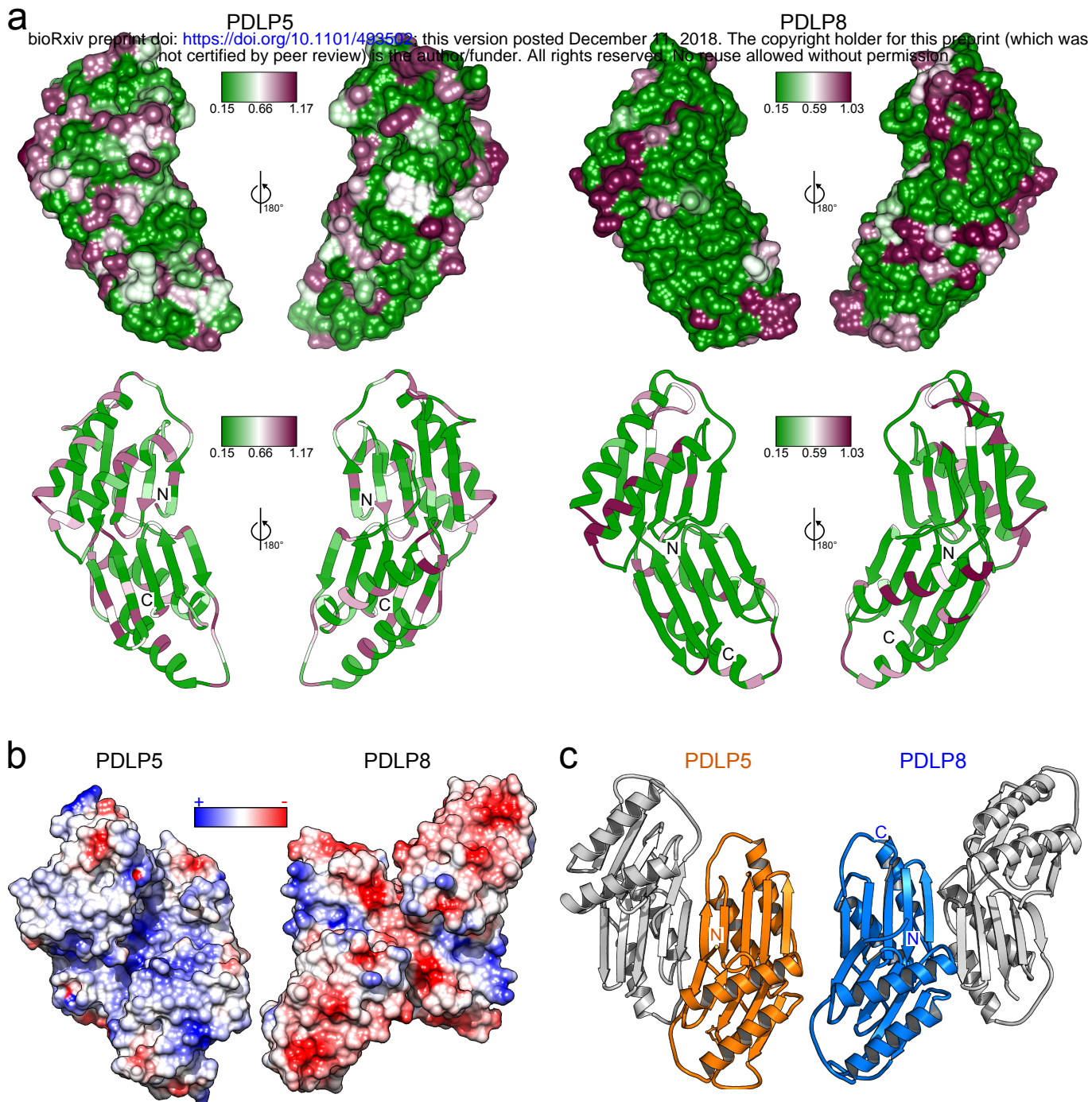


Fig 5: PDLP5 and PDLP8 may have drastically different oligomerisation modes, surface charge distributions and rates of evolution. a) Site-wise ω (d_s/d_s) values of amino acid positions ranging from 0.15 (green) reflecting conserved sites under purifying selection to 1.17 or 1.03 for PDLP5 or PDLP8 (magenta), respectively, reflecting variable sites possibly containing advantageous mutations illustrated on the molecular surfaces (upper) and in ribbon diagrams (lower) of PDLP5 (left) or PDLP8 (right), respectively. **b)** Electrostatic potential mapped onto molecular surfaces of the putative PDLP5 (left) and PDLP8 (right), orientation as in c) dimer, respectively. **c)** Ribbon diagrams of PDLP5 (orange, left) and PDLP8 (blue, right) crystallographic dimers. Note that in both dimers large, antiparallel β -sheets are formed, using different protein-protein interaction surfaces.

a Ectodomain Kinase domain

bioRxiv preprint doi: <https://doi.org/10.1101/493502>; this version posted December 11, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

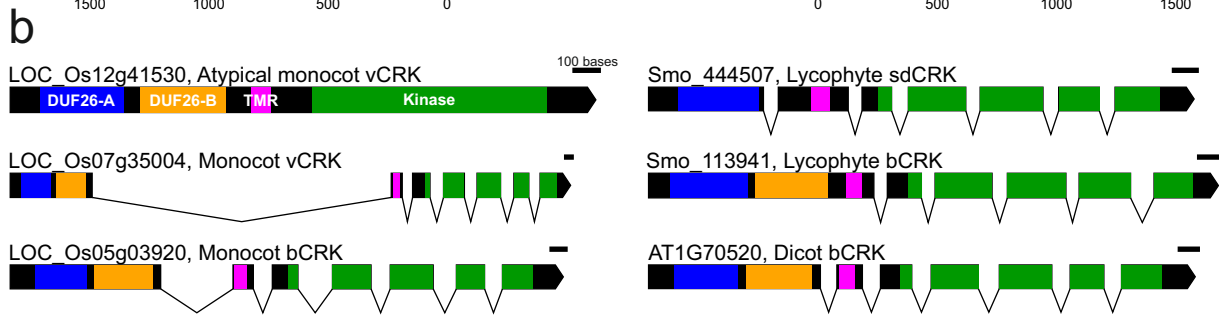
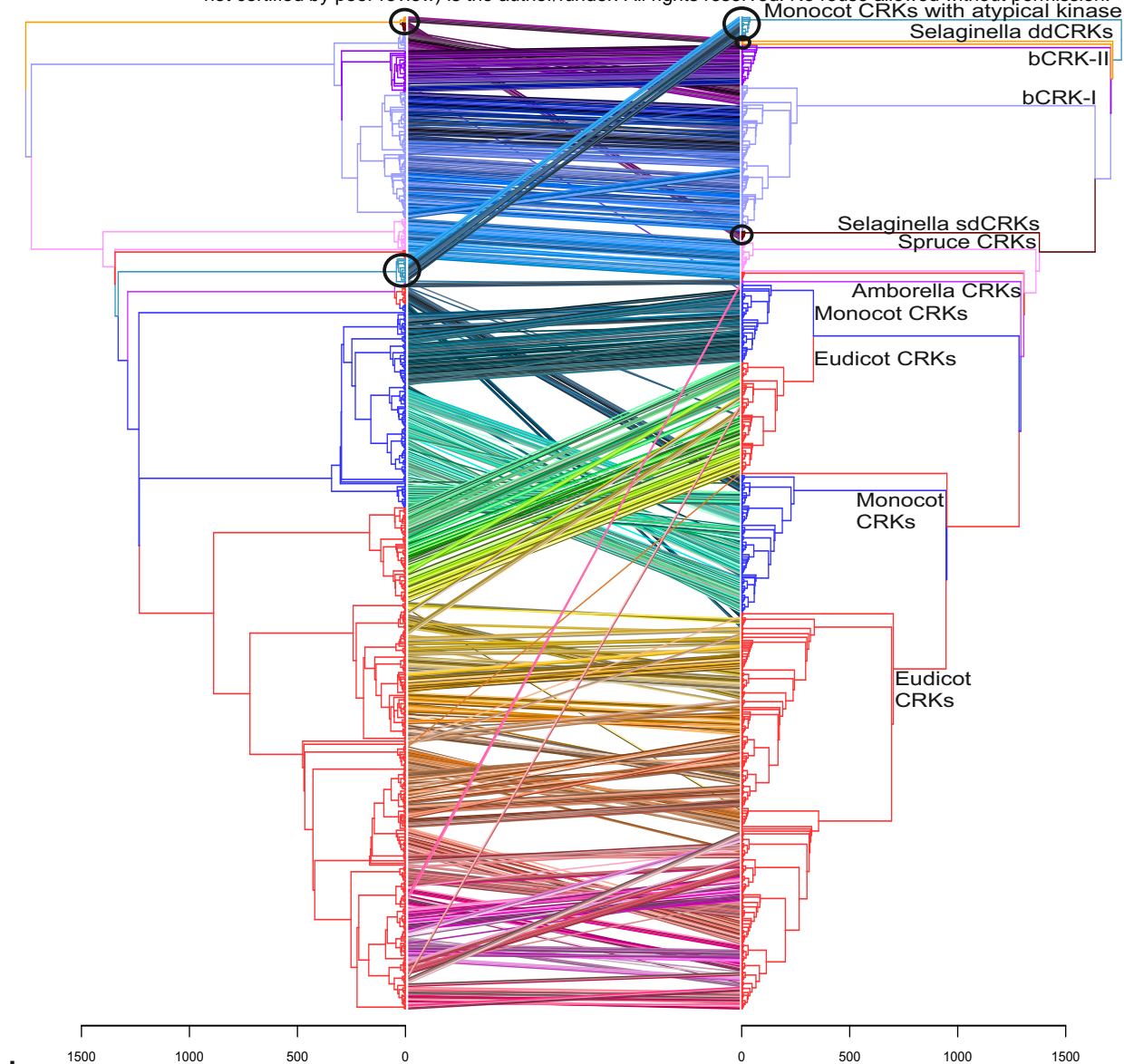


Figure 6. CRKs experienced domain rearrangements. **a)** Comparison of phylogenetic trees based on ectodomain region and kinase domain of 880 CRKs. Phylogenetic maximum-likelihood trees are presented as tanglegram where the tree of the CRK ectodomain region is plotted against the tree of the kinase domain. The kinase tree is rooted to group of monocot CRKs with a Concanavalin-A (ConcA) type kinase domain and the ectodomains tree is rooted to CRKs from *Selaginella moellendorffii*. The ectodomain tree was detangled based on the kinase domain tree. Lines connect the ectodomain and kinase domain belonging to same gene, and connection are drawn in different colors for better visibility. Juxtaposition of the trees shows rearrangements and domain swaps of ecto- and kinase domains. Black circles highlight the difference between the ectodomains and kinase domains of the *Selaginella* sdCRKs and ddCRKs and also the group of the atypical monocot CRKs which have exchanged the kinase domain. **b)** The exon-intron structure of the CRKs. Usually CRKs contain seven exons: one encoding DUF26 domains, one encoding transmembrane region (TMR) and five exons encoding the kinase domain. In atypical monocot CRKs with exchanged kinase domain, whole gene is encoded by one or two exons. The scale bar for each gene represents 100 bases. Regions encoding the DUF26-A are colored with blue, the DUF26-B with orange, the transmembrane region (TMR) with pink and the kinase domain with green.

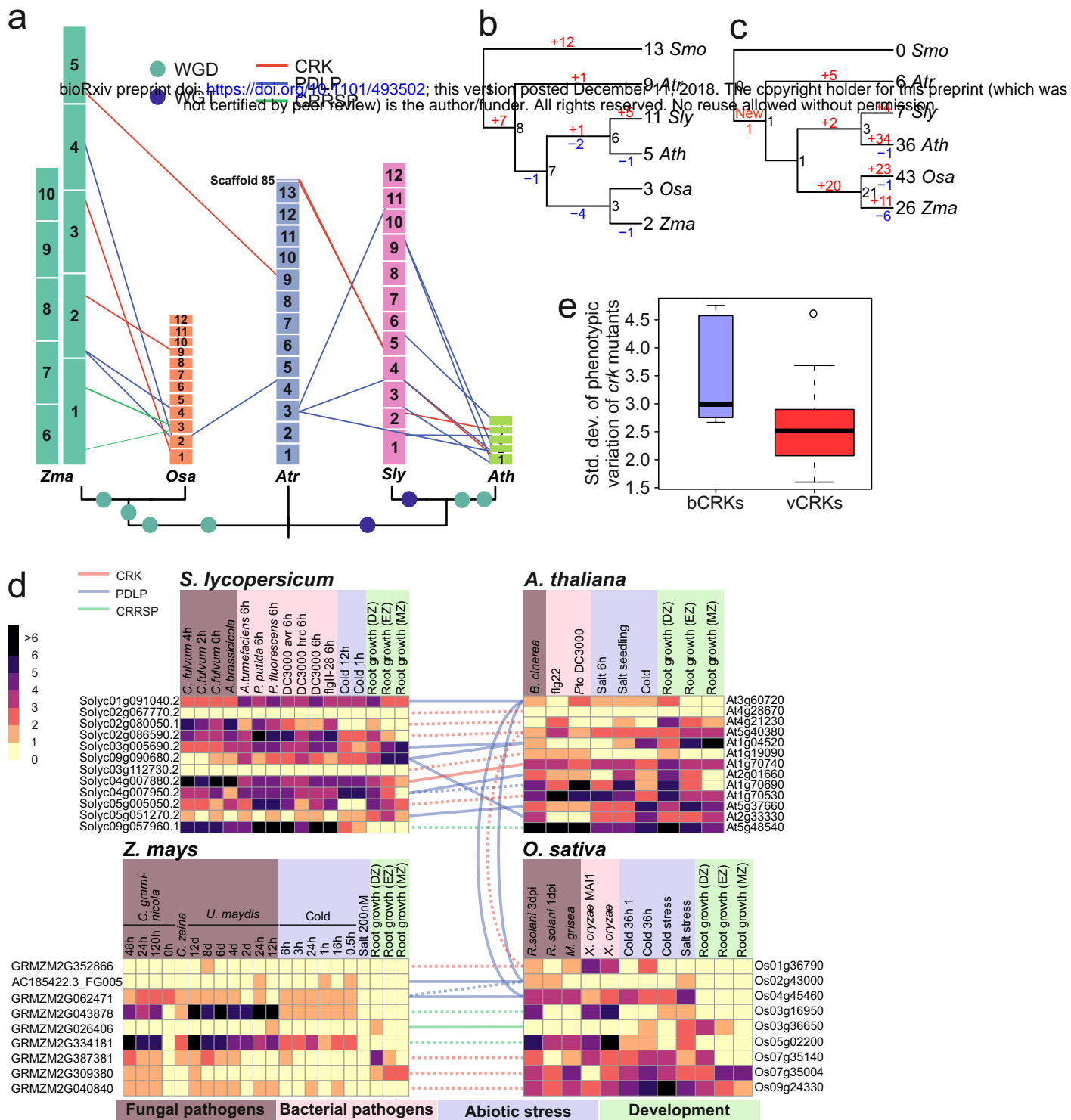


Figure 7. Identification of the mode of gene family evolution in DUF26-containing genes between *Arabidopsis thaliana*, tomato, rice, maize and *Amborella trichopoda*. Gene families that are preferentially retained after whole genome multiplications (WGMs) are typically identified by synteny analysis. **a**) Syntenic regions containing DUF26 genes from *Amborella trichopoda* to monocots *Oryza sativa* and *Zea mays* (to left from middle) and to eudicots *Solanum lycopersicum* and *Arabidopsis thaliana* (right from the middle). In the synteny analysis within monocots and dicots segments with at least 5 syntenic genes were included, whereas in comparisons to *Amborella* the minimum threshold was 3 syntenic genes. Analyses were carried out with Synmap software within CoGe. For *Amborella trichopoda* genomic locations of DUF26-containing genes are only known on chromosome/scaffold level based on physical mapping. Furthermore, gene families with a preferential retention pattern after WGMs show conserved gene counts over species. Phylogenetic tree of the five species shown in the panel was used to reconcile the gene trees and estimate gene counts in ancestral nodes for **b**) bCRKs and **c**) vCRKs, using *Selaginella moellendorffii* as outgroup. The gains are highlighted with red and losses with blue. Gene families with preferential retention pattern should also have many orthologs. **d**) Heatmaps of the normalized transcriptional expression counts (Transcript per million [TPM]) of candidate DUF26 orthologs from four of the species: *Solanum lycopersicum*, *Arabidopsis thaliana*, *Zea mays*, and *Oryza sativa*. Coloring in heatmaps is proportional to \log_2 (TPM) value that represents the gene expression level. The corresponding \log_2 (TPM) value is displayed next to the color key. The rows represent gene models and the columns show the experiments, collected from publicly available Sequence Read Archive (SRA) database. SRA accessions are annotated to relevant stress conditions (descriptions are presented in Table S4). Solid lines connect putative orthologs based on evidence from phylogenetic and synteny analyses; dashed lines connect putative orthologs based on evidence from either phylogenetic or synteny analyses. **e**) Final prediction of gene families evolving under dosage balance is that their knockouts demonstrate a high variance in phenotype. This can be seen by reanalysis of phenotype data from (Bourdais et al); the bCRK knockouts have a significantly higher standard deviation (Y-axis) over the different phenotyping experiments than vCRKs.

Pathogens: *Agrobacterium tumefaciens*, *Alternaria brassicicola*, *Botrytis cinerea*, *Cercospora zeina*, *Cladosporium fulvum*, *Colleotrichum graminicola*, *Magnaporthe grisei*, *Pseudomonas putida*, *Pseudomonas fluorescens*, *Pseudomonas syringae* pv. tomato DC3000, *Rizoctonia solani*, *Ustilago maydis*, *Xanthomonas oryzae*.

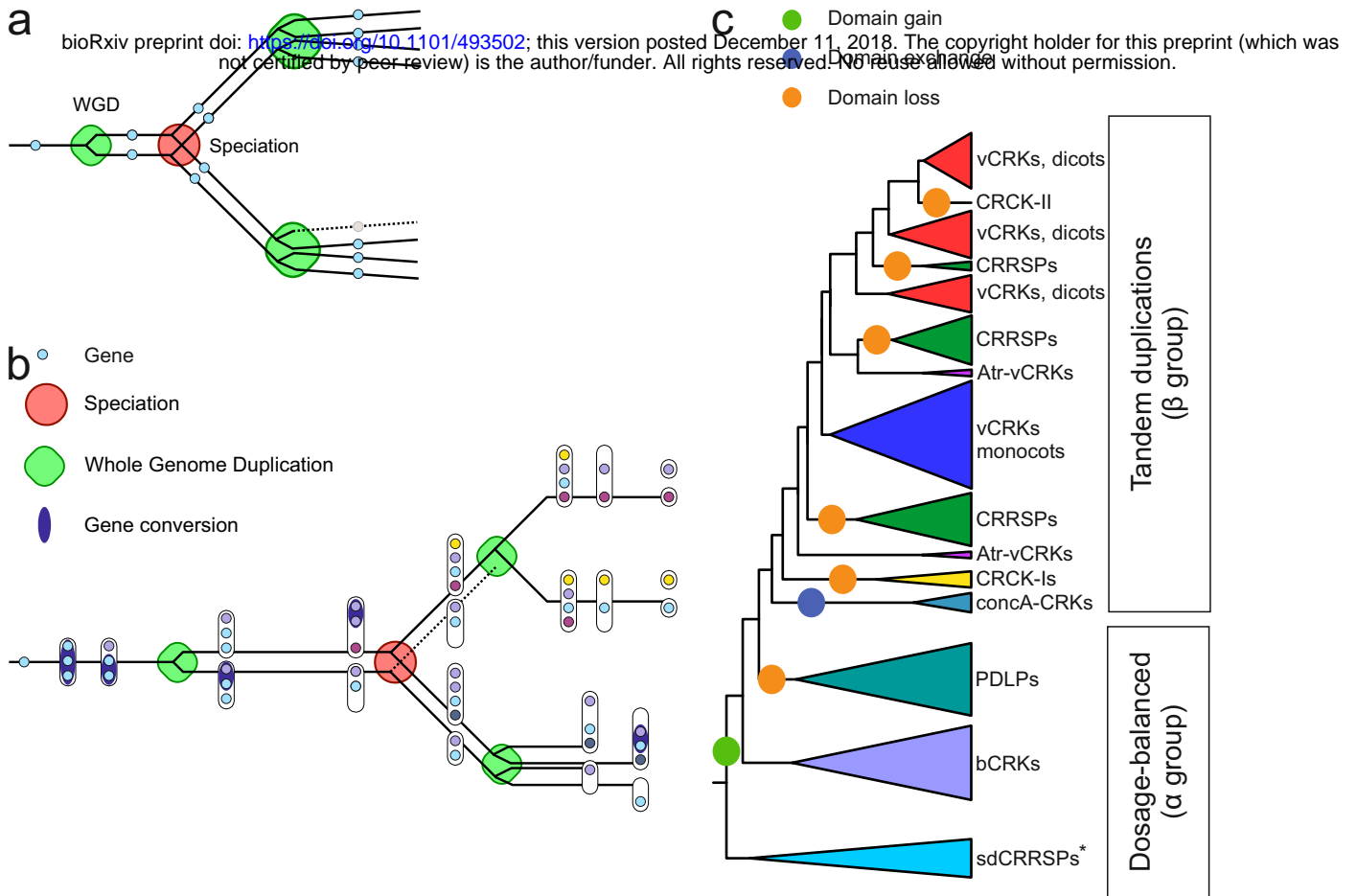


Figure 8. Model of mixed-type gene family evolution. Gene families evolve through two major events, whole genome multiplications (WGM) and small-scale duplications (SSD). Genes related to environmental responses and secondary metabolism experience SSDs in the form of tandems, whereas highly connected genes associated with transcriptional and developmental regulation or signal transduction functions are preferentially retained after WGMs. **a)** Prevailing hypothesis for the retention pattern is dosage-balance; for highly connected genes the stoichiometric balance needs to be maintained, and therefore selection acts against gene losses after WGMs and against duplications by SSDs. **b)** On the other hand, gene family evolving through tandem duplications (evolution before the speciation node) has a high birth rate and therefore the number of duplicates between species can vary. After duplications the homogeneity of the duplicates is maintained through gene conversion events, which has a high probability with homologous sequences that are near-by. This can be maintained for long periods, but eventually over time the sequences diverge by drift and selection based on dosage. Our data suggests that a tandemly expanding gene family may evolve into a dosage balance mode as a result of WGMs (evolution after speciation node). Following WGMs, the duplicated tandems may experience extensive fractionation due to drift and selection by dosage which breaks the tandem structure. At the same time, the connectivity of the gene family has been accumulating through sub- and neofunctionalization, and these phenomena together may result into a dosage balance model of evolution (top branch after speciation node). This does not necessarily occur across all WGM events, as was observed for bCRKs in *Solanaceae* (bottom branch), where there exist both single copies and tandem copies in the genome. Different subfamilies can be in different states of this process. **c)** We observed CRRSPs and PDLPs to follow dosage balance mode after the paleohexaploid event, whereas bCRKs have assumed the mode in later WGM events. The overall numbers of the bCRKs are preserved but identification of orthologs between species that have experienced independent WGMs is difficult, suggesting that convergent functionality of the members is recent. Gene families expanding through tandem duplications such as vCRKs and CRRSPs have high birthrate and demonstrate several lineage-specific expansions.

*Lost in Brassicaceae species and rice