

1 **The *Arabidopsis thaliana* pan-NLRome**

2 **Authors**

3 Anna-Lena Van de Weyer¹, Freddy Monteiro^{2,3}, Oliver J. Furzer^{2,5}, Marc T. Nishimura⁴,
4 Volkan Cevik^{5,6}, Kamil Witek⁵, Jonathan D.G. Jones^{*5}, Jeffery L. Dangl^{*2}, Detlef Weigel^{*1},
5 Felix Bemm¹

6 **Affiliations**

7 1. Department of Molecular Biology, Max Planck Institute for Developmental Biology, 72076
8 Tübingen, Germany.

9 2. Howard Hughes Medical Institute and Department of Biology, University of North Carolina,
10 Chapel Hill, North Carolina 27599-3280, USA.

11 3. Center for Research in Agricultural Genomics (CRAG), CSIC-IRTA-UAB-UB, 08193
12 Barcelona, Spain.

13 4. Department of Biology, Colorado State University, Fort Collins, CO 80523, USA.

14 5. The Sainsbury Laboratory, University of East Anglia, Norwich Research Park, Norwich,
15 NR4 7UH, United Kingdom.

16 6. The Milner Centre for Evolution, Department of Biology and Biochemistry, University of
17 Bath, Bath, BA2 7AY, United Kingdom.

18

19 Anna-Lena Van de Weyer, Freddy Monteiro, Oliver J. Furzer and Felix Bemm contributed
20 equally to this work.

21

22 *co-corresponding authors, e-mail: jonathan.jones@tsl.ac.uk; dangl@email.unc.edu;
23 weigel@weigelworld.org

24 Abstract

25 Disease is both among the most important selection pressures in nature and among the
26 main causes of yield loss in agriculture. In plants, resistance to disease is often conferred by
27 Nucleotide-binding Leucine-rich Repeat (NLR) proteins. These proteins function as
28 intracellular immune receptors that recognize pathogen proteins and their effects on the
29 plant. Consistent with evolutionarily dynamic interactions between plants and pathogens,
30 NLRs are known to be encoded by one of the most variable gene families in plants, but the
31 true extent of intraspecific NLR diversity has been unclear. Here, we define the majority of
32 the *Arabidopsis thaliana* species-wide “NLRome”. From NLR sequence enrichment and
33 long-read sequencing of 65 diverse *A. thaliana* accessions, we infer that the pan-NLRome
34 saturates with approximately 40 accessions. Despite the high diversity of NLRs, half of the
35 pan-NLRome is present in most accessions. We chart the architectural diversity of NLR
36 proteins, identify novel architectures, and quantify the selective forces that act on specific
37 NLRs, domains, and positions. Our study provides a blueprint for defining the pan-NLRome
38 of plant species.

39 Introduction

40 Plant immune receptor repertoires have been shaped by millennia of plant-microbe
41 coevolution^{1,2}. Immunity is activated either by cell surface receptors that recognize microbe-
42 associated molecular patterns (PAMPs), or by intracellular receptors that detect pathogen
43 effectors¹. These intracellular receptors are typically encoded by highly polymorphic genes.
44 About two thirds of disease resistance genes encode nucleotide-binding leucine-rich repeat
45 receptors (NLRs)³, and most plant genomes carry hundreds of NLR genes⁴. The majority of
46 plant NLRs contain a central nucleotide binding domain shared between Apaf-1, Resistance
47 proteins and CED4 (NB-ARC, hereafter NB for simplicity)⁵. Most contain also leucine-rich
48 repeats (LRRs)^{6,7}, and either a Toll/Interleukin-1 receptor (TIR) or coiled-coil (CC) domain at
49 the N-terminus⁸⁻¹⁰. Proteins with similar arrangements of functional domains are also
50 involved in host defense in animals and fungi¹¹⁻¹³.

51
52 Recognition by NLRs generally involves one of three main mechanisms¹⁴. NLRs can directly
53 detect pathogen effectors through interaction with the canonical NLR domains¹⁵⁻¹⁷, or with
54 an NLR-incorporated integrated domain (ID) that resembles known domains of pathogen
55 effector targets¹⁸⁻²². Alternatively, NLRs detect effector activity indirectly by monitoring a
56 host virulence target (“guardee”)²³⁻²⁵, or detect effectors via direct interactions. Importantly,

57 these mechanisms have been directly demonstrated only for a very small number of NLRs,
58 and additional mechanisms might await discovery.

59

60 To date, NLR complements, or NLRomes, have been defined from available genome
61 annotations for single cultivars of plants or for multiple species across different taxonomic
62 levels, respectively ^{2,4,26–29}. The most striking findings were the repetitive modular
63 arrangement of NLRs and the discovery of head-to-head paired NLR genes, of which one
64 member included an ID ^{2,4,22,26–28}. The potential use of those IDs as modular building blocks
65 has opened new possibilities for the engineering of novel resistances to pathogens ^{30–33}. The
66 existing list of IDs, however, likely represents only a glimpse of the true diversity across
67 plants.

68

69 The definition of pan-NLRomes, or repertoires of NLR genes, across different species, or
70 higher taxonomic groups, has provided estimates of the variation in size of the NLR family ^{34–}
71 ³⁶, presence/absence relations ^{35,36}, categorical distribution into structural classes across the
72 phylogeny, and diversity of IDs ^{26,27}. Publicly available plant genome annotations have been
73 the foundation of most NLRome studies and, their systematic integration has allowed
74 ancestry reconstruction of key NLR lineages and illuminated ancient and recent expansion-
75 contraction events ⁴. In contrast, knowledge of the true diversity of within species pan-
76 NLRomes is scarce and has so far been derived from only a limited number of individuals,
77 and thus covers a narrow diversity within the population ^{34–37}. Across individuals of the same
78 species, which often has only a single reference genome annotation, the remarkable
79 differences in NLR family size between rice, tomato, and *Arabidopsis thaliana* might be due
80 to low coverage of available genomes, or the difficulty of accurate assembly of tandem
81 paralogous genes often found in NLR clusters when short-read sequencing is used under
82 conditions of insufficient depth ^{34,35}.

83

84 Despite these potential shortcomings, early intraspecific pan-NLRome studies revealed
85 patterns of allelic and structural variation consistent with adaptive evolution and balancing
86 selection for subsets of NLR-encoding genes ³⁷, fitting a model of co-evolution of host and
87 pathogens. Allelic variation seems to be reflected in many different haplotypes that are found
88 across NLR loci ^{38,39}. These can include recombination “hotspots” generating NLR clusters
89 ^{36,40}, and true allelic series ^{15,41}. The patterns of presence/absence polymorphisms as well as
90 copy number variation at loci with multiple NLR genes imply that reference genomes may
91 not include representatives of all distinct NLR clades within a species ^{36,38,42–44}. *Resistance*

92 gene enrichment sequencing (RenSeq) facilitates discovery of “missing” NLR genes in a
93 species, especially when hybridization-based capture of genomic fragments with sequence
94 similarity to known NLR-coding genes is combined with Single-Molecule Real Time
95 sequencing (SMRT RenSeq)⁴⁵.

96

97 Our objective was to define the full NLR repertoire and its variability in the reference species
98 *A. thaliana*, by analyzing a panel of 65 diverse accessions using SMRT RenSeq. We show
99 that we reach saturation of the pan-NLRome with this well-chosen set of accessions; we
100 define the core NLR complement of the species and detail novel domain architectures; and
101 we describe presence-absence polymorphisms in non-core NLRs. Together, our work
102 provides a foundation for the identification and cloning of disease resistance genes in more
103 complex species of agronomic importance.

104 Results

105 The Samples

106 A set of 65 *A. thaliana* accessions was selected to explore the diversity of the pan-NLRome
107 (Fig. 1a, Supplementary Table 1). The selection included 46 accessions from the 1001
108 Genomes Project, of which 21 belonged to previously identified relict populations
109 characterized by an unusually high amount of genetic diversity⁴⁶. Additionally, the 19
110 founder accessions of MAGIC lines, a resource to dissect the genetics of complex traits^{47,48}
111 and the widely-studied accession Ws-2, were included.

112 NLR Complements

113 Baits were designed to hybridize with 736 NLR-coding genes from multiple Brassicaceae,
114 including *Arabidopsis thaliana*, *Arabidopsis lyrata*, *Brassica rapa*, *Aethionema arabicum* and
115 *Eutrema parvulum* (see Online Methods for details of bait design, sequencing, assembly,
116 annotation and quality control approaches). A combination of NLR sequence capture
117 (RenSeq) and single-molecule real-time sequencing (SMRT) was used to reconstruct full
118 NLR complements. In total, we produced 13,167 NLR-coding gene annotations, with a range
119 of 167 to 251 genes per accession (Fig. 1b). Individual accessions had between 47% and
120 71% physically clustered NLR genes (more than one NLR in 200 kb of genomic sequence;
121 adapted from⁴⁹. A particularly interesting subset of NLR-coding genes are those in head-to-
122 head orientation^{50,51}, and we found 10 to 34 NLRs per accession in such an orientation, or
123 with high sequence similarity to known functional pairs (see Online Methods). NLRs were

124 grouped into four classes (TNL, NL, CNL, and RNL; see Glossary) based on canonical
125 protein domains (TIR, NB, CC, RPW8 and LRR). Across all accessions, TNLs formed the
126 largest and most size-variable class, followed by NLS, CNLs, and RNLs (Fig. 1c,
127 Supplementary Fig. 1). Of the 13,167 NLR genes, 663 contained at least one additional
128 integrated protein domain (ID), in which we found 36 distinct Pfam domains (Fig. 2a,b,
129 Supplementary Table 2 and Supplementary Table 3). Individual accessions had 5 to 17 IDs
130 distributed across 4 to 16 NLR genes, in line with reports for specific accessions^{4,36}. This
131 result reveals an unprecedented incidence of previously unreported *A. thaliana* IDs.
132 Annotated RenSeq Col-0 identifiers and sequences are provided as supplementary material,
133 but the Col-0 TAIR10/Araport11 sequences and identifiers were used in all downstream
134 analysis (see Online Methods).

135 NLR Domain Architecture Diversity

136 We investigated the repertoire of the 36 IDs, since these might function as pathogen effector
137 binding platforms^{19–21,32}. 29 of the 36 IDs were already known from other Brassicaceae
138 including *A. thaliana* Col-0 (Fig. 2a, b; Supplementary Table 2, Supplementary Table 4,
139 Supplementary Table 5, Supplementary Table 6). Nine of the 36 IDs were reported
140 concordantly in the two major NLR-ID censuses, namely WRKY, PP2, Pkinase, PAH,
141 DUF640, B3, Pkinase_Tyr, PPR_2 and Alliinase_C^{26,27}. Five of those nine occur in
142 genetically linked paired NLRs in the pan-NLRome (pair ratio > 0.5 in Fig. 2b, see Online
143 Methods and Glossary; Supplementary Table 2). Rediscovery of these nine IDs is of
144 relevance, since these are enriched for domains similar to known effector targets^{26,27,52,53}.
145 Our sequencing and annotation effort expands the *A. thaliana* ID repertoire beyond the ten
146 IDs found in the Col-0 reference accession. IDs found in only one gene model did not
147 receive particular attention, as they are conceivably an artefact of our annotation pipeline.

148
149 A hallmark of NLRome variation across species is the variation in the relative fraction of
150 different domain architectures^{4,11}. Examining the arrangement of NLR domains in the *A.*
151 *thaliana* pan-NLRome we identified 97 distinct architectures (Supplementary Fig. 2). Whilst
152 27 canonical architectures (without IDs) account for the vast majority of the identified NLRs
153 (95% of the pan-NLRome), the remaining 5% contain at least one of 36 different IDs (Fig.
154 2c). The 97 architectures greatly augment the 22 architectures found in the reference Col-0
155 genome (Fig. 2d, Supplementary Table 7), with most of the new *A. thaliana* architectures
156 containing at least one ID (Supplementary Fig. 2). Half of the new *A. thaliana* architectures
157 contain more than one gene (38/75) (Fig. 2e), of which, 17 predominantly comprise paired

158 NLRs (pair ratio > 0.5, see Online Methods and Glossary) and contain at least one ID (Fig.
159 2e). About half of the architectures have not been previously described in the Brassicaceae
160 (including *A. thaliana* Col-0) (49/97) (Fig. 2d, Supplementary Table 7). These novel
161 architectures account only for 1.3% of the pan-NLRome (175 NLRs), with all but one
162 containing an ID (Fig. 2d, e, Supplementary Table 7, Supplementary Table 8 and
163 Supplementary Table 9). Finally, 12 IDs are repeatedly recruited into different novel
164 architectures (labeled “novel > known” in Fig. 2f, Supplementary Table 8 and Supplementary
165 Table 9), reflecting the recycling of a limited set of IDs into new domain arrangements. It is
166 likely that these IDs are derived from proteins repeatedly targeted by pathogen virulence
167 effectors.

168 The pan-NLRome

169 To begin to understand the diversity of both NLR content and alleles, we grouped sets of
170 homology-related NLRs from different ecotypes. The resulting clusters were termed
171 orthogroups. We clustered 11,497 NLRs into 464 high confidence orthogroups (OGs) (Fig.
172 3a), plus 1,663 singletons. Ninety-five percent of the OGs could be discovered within 38
173 randomly chosen accessions (Fig. 3b). Additional sampling only recovered OGs with three or
174 fewer members, indicating that the pan-NLRome we describe is largely, if not completely,
175 saturated. OGs were classified according to size, domain architecture and structural
176 features. We define the core NLRome as the 106 OGs found in at least 52 accessions
177 (6,080 genes), 143 OGs found in at least 13, but fewer than 52 accessions as shell (3,932
178 genes), and the 215 OGs found in 12 or fewer accessions as cloud (1,485 genes) (Fig. 3a).
179 The majority of OGs, 58%, were TNLs, in concordance with TNLs being the prevalent NLR
180 class in the *Brassicaceae*⁵⁴, 22% were CNLs, 7% RNLs, and 13% NLs (Fig. 3c). TNLs
181 showed a strong tendency towards larger shell and core OGs compared to CNLs
182 (Supplementary Fig. 3). Sixty-four OGs included genetically paired NLRs (see Online
183 Methods), and 28 contained members with an ID, with almost none being present in the
184 cloud NLRome (Fig. 3d). Shell and core OGs contained most paired NLRs (98% in 55 OGs,
185 Supplementary Fig. 3). This shows that conserved NLR pairs are widely distributed in the
186 population and that incorporation of IDs into NLRs is widespread in *A. thaliana*.

187 Placement of non-reference OGs

188 We discovered 296 high confidence OGs without a reference Col-0 allele, with six belonging
189 to the core, 205 to the cloud, and 85 to the shell NLRome. In order to anchor these OGs to
190 the reference genome, we asked how often orthogroups co-occurred, using OGs with known

191 location (NLR and non-NLR OGs with a Col-0 reference allele) to anchor contigs with OGs
192 lacking a reference allele. Anchoring efficiency of newly discovered OGs was calculated for
193 co-occurrences in 2 to 39 accessions (Supplementary Table 10). With a minimum threshold
194 of 10 accessions, we derived 42 co-occurrence subnetworks (Supplementary Fig. 4),
195 anchoring 24 out of 132 OGs present in at least 10 accessions, but missing from the Col-0
196 reference. Most were anchored to other NLRs (Supplementary Table 10 and Supplementary
197 Fig. 4). Newly anchored OGs include one CNL pair and three TNL pairs (Fig. 4,
198 Supplementary Fig. 4, and Supplementary Fig. 5), with one ID-containing sensor-type OG
199 (205.1) arranged in head-to-head orientation to the executor-type OG 204.1 (Supplementary
200 Fig. 5). The use of annotated non-NLR genes in the assembled contigs allowed us to
201 properly place these novel OGs.

202 Pan-NLRome Diversity

203 In an orthogonal approach to classifying NLR genes by their architectures, we assessed
204 sequence diversity, an indication of the evolutionary forces shaping the pan-NLRome.
205 Average nucleotide diversity was similar for CNLs, NLs and TNLs (Fig. 5a). It was lowest in
206 RNLs, but because this is the smallest group, this difference is not statistically significant.
207 The same trend was true for haplotype diversity (Fig. 5c). Nucleotide diversity was lowest in
208 core and higher in shell and cloud OGs across TNLs, CNLs and NLs (Fig. 5a), suggesting
209 that selection is relaxed in OGs with larger presence/absence variation. The pattern was
210 reversed for haplotype diversity (Fig. 5c) and core OGs generally showed high values. We
211 noticed that a few shell TNL and NL OGs stand out because of their ultra-low haplotype
212 diversity, suggesting a conserved but rarely encountered selective pressure (without any
213 correlation between geographic location and the accessions carrying these orthogroups).
214 The average nucleotide diversity saturated with 32 randomly selected accessions, and the
215 haplotype diversity with 49 accessions, suggesting a prevalence of low frequency haplotypes
216 (Supplementary Fig. 6). Compared to non-clustered OGs, physically clustered OGs had
217 significantly higher nucleotide diversity (Supplementary Fig. 3). This finding may indicate
218 relaxed selection after gene duplication in these clusters⁵⁵. When considering different NLR
219 protein domains, the highest diversity was found in LRRs across all major classes and
220 subclasses (Fig. 5b). Combining population genetics statistics for a Principal Component
221 Analysis (PCA) revealed that more than 60% of the variance can be explained by as little as
222 two components (Supplementary Fig. 7). However, none of the collected metadata, such as
223 orthogroup size, type, class or the presence of IDs or a potential partner, explained the
224 clustering of the first two principal components (Supplementary Fig. 7). This suggests a

225 complex interplay of the different factors driving NLR evolution. Tajima's D values, which can
226 indicate balancing and purifying selection⁵⁶, were similarly distributed across different NLR
227 classes, with all classes containing extremes in both directions (Fig. 5d). Low Tajima's D
228 values were most common in TNLs, largely driven by core- and shell-type OGs.

229

230 Site-specific selection analyses revealed core and shell OGs that have likely experienced
231 constant (248), pervasive (165) or episodic (130) positive selection (Fig. 6a, b;
232 Supplementary Fig. 8). Codons completely invariable, indicating constant positive selection,
233 can be found across all types (e.g., core, shell) and classes (e.g., TNLs, CNLs). Pervasive
234 positive selection seemed more likely in core-like OGs (71%) than in shell-like OGs (63%)
235 while episodic positive selection patterns showed at a similar rate (52%). Subclasses
236 showed a more uneven pattern of positive selection (Fig. 6e). Sites under constant positive
237 selection were mostly found in NB, TIR and LRR regions when comparing all annotated
238 protein domains (Fig. 6f). Pervasive and episodic positive selection patterns appeared
239 predominantly in NB and TIR domains (Fig. 6g, h). A few OGs stood out because of the
240 large fraction of codons of annotated protein domains under positive selection, including
241 *RPP13* which confers race-specific downy mildew resistance⁵⁷ (Supplementary Fig. 8).
242 Sites under positive selection were also found in 11 IDs, including WRKY, TCP, B3 and
243 DA1-like domains (Fig. 6c). Notably, invariant sites were detected in the WRKY domains of
244 all three OGs containing a WRKY, and in a surprisingly high proportion of sites in the BRX
245 domains of the RLM3-containing OG (Supplementary Table 11). We conclude that positive
246 selection is widespread in the core-NLRome, being most prevalent in canonical NLR
247 domains.

248 Linking Diversity to Function

249 Because NLRs that had been experimentally implicated in resistance to biotrophic
250 pathogens showed enhanced diversity, we sorted OGs by resistance to adapted biotrophs
251 (*Hyaloperonospora arabidopsidis*), non-adapted biotrophs (*Brassica*-infecting races of
252 *Albugo candida*)⁵⁸ and hemibiotrophs (mostly *Pseudomonas* spp.). OGs that provide
253 resistance against adapted biotrophs are significantly more diverse than other categories
254 (Fig. 7a; ANOVA and Tukey's HSD $p < 0.01$), suggesting that host-adapted biotrophic
255 pathogens are driving diversification of NLRs more than other pathogens. That RNL helper
256 NLRs have low diversity is consistent with their requirement to function with several sensor
257 NLRs⁵⁹⁻⁶¹.

258

259 Among the OGs with the lowest Tajima's D values, a prominent example was *RPM1*, which
260 confers resistance to a hemibiotrophic bacterial pathogen, and for which an ancient, stably
261 balanced presence/absence polymorphism across *A. thaliana* is well established⁶². OGs
262 that provide resistance to adapted biotrophs tend to have higher Tajima's D values,
263 indicating that they experience not only diversifying, but also balancing selection. Tajima's D
264 values within sensor-executor pairs encoded in head-to-head orientation were correlated
265 whereas other closely linked NLR genes or random pairs were not (Fig. 7b, Supplementary
266 Table 12). As an example, two OGs with high Tajima's D values are the paired NLRs *CSA1*
267 (OG91) and *CHS3* (OG130). *CHS3* featured two very different groups of alleles
268 distinguished by the presence of LIM and DA1-like IDs⁶³. This pattern was perfectly mirrored
269 by the one for *CSA1*, the paired "executor" partner NLR of *CHS3* (Fig. 7c).

270 Discussion

271 We defined the full species repertoire of the gene family that encodes NLR immune
272 receptors in the model plant *A. thaliana*. Importantly, the pan-NLRome inventory became
273 saturated with ~40 accessions randomly selected from the 65 accessions we analyzed.
274 Before our work, it was known that there was excessive variation at some NLR loci, such
275 that in the small number of accessions in which the relevant genomic region was analyzed in
276 detail, every accession was very different, including significant presence/absence variation
277^{41,64}. The fact that our pan-NLRome saturates with ~40 accessions indicates that the number
278 of divergent loci is not unlimited. It also provides some guidance for future efforts in other
279 species. It will be fascinating to compare the allelic and diversity saturation of self-fertilising
280 *A. thaliana* with obligate out-crossers and with domesticated species. Among functionally
281 annotated genes, we found the highest sequence diversity in NLR-coding genes whose
282 products recognize evolutionarily adapted biotrophic pathogens.

283
284 We have also found an astonishing diversity of IDs, which allow hosts to rapidly accrue the
285 ability to recognize the biochemical action of pathogen effector proteins. ID-containing NLRs
286 that have been functionally characterized are all found in paired orientation. In these pairs,
287 the ID member functions as pathogen sensor, and the other member as signaling executor
288^{19-21,50,51,63,65}, with both members contributing to repression and activation of NLR signaling
289⁶⁶. The correlation between Tajima's D values of such paired NLRs support a co-evolutionary
290 scenario whereby mutations into the sensor component lead to compensatory changes in
291 the executor, or vice versa.

292

293 However, half of the 22 most commonly found IDs did not occur in an arrangement indicative
294 of sensor/executor pairs. An open question is whether these function with unlinked executor
295 partners, or whether they can function as dual sensor/executor proteins. Within the *A.*
296 *thaliana* pan-NLRome, we identified three key families of defense-related TCP, WRKY and
297 CBP60 transcription factors, represented as IDs in sensors of the class defined by RRS1.
298 TCP domains are particularly interesting, as TCP transcription factors are preferentially
299 targeted by pathogen effectors from divergently evolved pathogens^{52,53,67,68}. The TCP
300 domain may open a new avenue to engineering of NLR specificity, through TCP swap or
301 inclusion of known effector-interacting platforms from TCP14⁶⁵, as recently demonstrated
302 with protease cleavage site swaps^{30,67}. Furthermore, since many of the novel IDs were
303 found at intermediate frequency in the population, and were novel compared to the Col-0
304 reference genome, we predict that this will apply to other plant species, suggesting that the
305 number and diversity of NLR-IDs greatly exceeds that which has been so far reported.

306 Figure Legends

307 **Figure 1. Overview of NLR complements in 65 accessions.** a) Map of accession
308 provenance. 1001 Genomes (relicts, blue, non-relicts, purple), MAGIC founders (yellow). b)
309 Number of total, clustered and paired NLRs in accessions. Means, solid black lines;
310 Bayesian 95% Highest Density Intervals (HDIs), solid bands. Individual data, open circles;
311 full densities shown as bean plots. c) Number of different structural classes in accessions.
312 Mean, HDI, and individual data as in (b).

313

314 **Figure 2. Diversity of IDs and domain architectures.** a) UpSet intersection of IDs in the
315 pan-NLRome, the Col-0 reference accession and 19 other Brassicaceae. The number of IDs
316 in each set is indicated between parentheses at the lower left. Set intersections shown on
317 the bottom. b) ID analysis. IDs known for *A. thaliana* in black and IDs newly described for *A.*
318 *thaliana* in light grey. Asterisks indicate IDs not found in other Brassicaceae. Numbers next
319 to y-axis indicate ratio of paired NLRs among ID containing NLRs. c) Cumulative size
320 contribution to the pan-NLRome (y-axis) of each of the 97 size sorted domain architectures
321 (x-axis). d) UpSet intersection plot of architectures shared between pan-NLRome, *A.*
322 *thaliana* reference- and 19 other Brassicaceae. e) 38 architectures with at least two
323 representatives and not found in the Col-0 reference. Asterisks indicate 20 architectures not
324 found in 19 Brassicaceae family genomes or the *A. thaliana* Col-0 reference. Numbers next
325 to y-axis shows fraction of paired NLRs. f) Number of known and novel architectures
326 containing the 27 overlapping Brassicaceae IDs (see a). “a” and “b” to the right of the bars
327 indicate putative IDs^{26,27} (see also Supplementary Table 2).

328

329 **Figure 3. OG sizes, saturation, and distribution of core-, shell- and cloud-NLRs.** a) OG
330 size distribution for TNLs (yellow), NLs (green), CNLs (blue), RNLs (purple), and all NLRs

331 (grey dashed line). The vertical lines at $x=13$ and $x=51$ differentiate cloud, shell and core. b)
332 Saturation of the pan-NLRome. The blue boxes show the percentage of the pan-NLRome
333 that can be recovered when randomly drawing a fixed number of accessions (1000x
334 bootstrapping). The horizontal dashed line indicates where 90% of the pan-NLRome is
335 found. The green boxes shown for each subset of drawn accessions indicate the average
336 size of undiscovered orthogroups. The vertical dashed line indicates that 95% of the pan-
337 NLRome can be recovered with 38 accessions. c) OG-type specific distribution of NLR
338 classes in the cloud (dark blue), the shell (grey), and the core (olive green), and the relative
339 fraction (numbers in the bars). d) OG-type specific distribution of paired and unpaired NLRs,
340 and NLRs with and without IDs in the cloud (dark blue), the shell (grey), and the core (olive
341 green), and the percentage (numbers in the bars).

342

343 **Figure 4. Genomic location of NLR genes in the reference assembly.** The five *A.*
344 *thaliana* chromosomes are shown as horizontal bars with centromeres in red. Reference
345 NLRs are shown as black line segments. Text labels are shown only for functionally defined
346 Col-0 NLRs. Anchored OGs found in at least 10 accessions are shown below each
347 chromosome. Anchored OGs with paired NLR members are shown in orange, while
348 remaining anchored OGs are shown in blue.

349

350 **Figure 5. Population genetic statistics across the Pan-NLRome.** a) Nucleotide diversity
351 distribution grouped by orthogroup type and NLR class. Nucleotide diversity was defined as
352 average pairwise nucleotide difference normalized to the number of sites in the respective
353 orthogroup. b) Nucleotide diversity distribution grouped by domain type and NLR class.
354 Sites of each orthogroup annotated with the same domain were aggregated. Type C within
355 class NL occurs where a minority of OG members had an identifiable CC-domain. c)
356 Haplotype diversity distribution grouped by orthogroup type and NLR class. Haplotype
357 diversity was defined as average pairwise haplotype differences. A value of 1 indicates a
358 high chance of finding two different haplotypes for two randomly chosen sequences of a
359 given orthogroup. d) Tajima's D (a measure of genetic selection) distribution grouped by
360 orthogroup type and NLR class. RNL orthogroups are not shown because of the low number
361 of orthogroups showing this class.

362

363 **Figure 6. Positive selection landscape of the Pan-NLRome.** (a-e) Ratio of different
364 positive selection classes grouped by NLR class (a), orthogroup type (b), presence of a non-
365 canonical domain (c), presence of a paired NLR (d) or NLR subclasses (e). An orthogroup
366 was considered if at least one positive selected site of a given class was detectable. (f-h)
367 Ratio between orthogroups showing constant (f), pervasive (g) or episodic (h) selection or no
368 selection grouped by annotated protein domains.

369

370 **Figure 7. Population genetics of different OG classes grouped by known resistances
371 and pairs.** a) Nucleotide diversity distributions by functional class according to pathogen
372 type to which they provide resistance. b) Correlation of Tajima's D values in sensor/executor
373 and control pairs. c) Maximum-likelihood phylogenetic trees of two OGs 90.1 and 130.1,
374 which form a sensor/executor pair⁶³ (100 bootstrap support indicated at major nodes). Scale
375 bar refers to substitutions per site. Lines connecting the trees denote same accession.

376 Supplementary Figure Legends

377 **Supplementary Figure 1. NLR frequency for different subclasses.** For each subclass,
378 the corresponding class is color coded (TNLs: yellow, NLs: green, CNLs: blue, and RNLs:
379 purple), and classes are in addition divided by the vertical dashed lines. NLRs are grouped
380 into subclasses by their domains content: T (TIR), N (NB), C (CC), R (RPW8), and X (all
381 other integrated domains). Each domain must be present at least once, domains in brackets
382 may be present. Domain order is not considered. The mean is shown as a solid black
383 horizontal line and the 95% Highest density Intervals (HDI: points in the interval have a
384 higher probability than points outside) are shown as solid bands around the sample mean.
385 All raw data points are plotted as open circles and the full densities are shown as a bean
386 plot.

387

388 **Supplementary Figure 2. Schematic representation of NLR domain architecture**
389 **diversity and simplification of consecutively repeated domains.** a) Examples of NLR
390 domain architecture diversity. On top, a generic NLR, with an ID (Integrated Domain) is
391 shown at the C-terminus. IDs can also be found at the N-terminus, and more rarely between
392 the three canonical domain types. b) Reduction of domain combinations by collapsing
393 duplicated/repetitive domains. The number of NLRs grouped by each of the original
394 architectures is shown on the left, along with one example that can be visualized in the
395 genome browser. Ellipsis in the bottom left represent 19 other architectures containing 4,079
396 proteins exclusively composed of TIR, NB and LRR domains. The same strategy was
397 applied to all other architectures containing at least one duplicated domain in the RPW8, NB
398 and CC classes. c) Full set of the novel *A. thaliana* NLR architectures. Includes the
399 architectures contributed by only one gene. Domain architectures are shown in the y-axis.
400 The number of NLRs in each architecture is shown in the x-axis. Asterisks indicate the 49
401 architectures not yet detected in the Brassicaceae outside of *A. thaliana*, or in the reference
402 accession Col-0. Numbers next to y-axis show the ratio of paired NLRs divided by the total
403 number of NLRs in each architecture.

404

405 **Supplementary Figure 3. OG size distribution comparisons.** Vertical black lines divide
406 cloud (left section) from shell (middle section) and core (right section) NLRs. a) Comparison
407 of OG size distributions of TNL OGs (blue) and CNL OGs (green) b) Comparison of paired
408 (blue) and non-paired (green) OGs. c) Comparison of clustered (blue) and non-clustered
409 (green) OGs. d) Comparison of ID-containing (blue) and non-ID-containing OGs (green). e)
410 Distribution of Paired NLRs and NLRs with IDs across the Cloud- (dark blue), the Shell-
411 (grey), and the Core- (olive green) pan-NLRome.

412

413 **Supplementary Figure 4. Orthogroup co-occurrence network.** Annotated NLR (green
414 nodes) and non-NLR genes (white nodes) clustered into OGs were analyzed for co-
415 occurrence in the same contig. The number of co-occurrences is represented by grey lines
416 connecting nodes (edges). The minimal co-occurrence threshold imposed was 10
417 accessions, but similar networks can be derived for any number accessions. NLR OGs
418 without a Col-0 allele (green square nodes) are highlighted in blue boxes. Hypothetically
419 paired OGs not known in Col-0 are highlighted in orange boxes.

420

421 **Supplementary Figure 5. Quantitative co-occurrence of the novel hypothetical paired**
422 **NLRs in OG205.1 and OG204.1.** Abbreviations: OG, Orthogroup; H2H, Head-to-head; NA,
423 Not available.

424

425 **Supplementary Figure 6. Nucleotide and haplotype diversity saturation.** Saturation of
426 nucleotide (a) and haplotype (b) diversity after random subsetting of the complete Pan-
427 NLRome into bins of increasing sizes. For each size 100 bootstrap were carried out.

428

429 **Supplementary Figure 7. Population genetics statistics based PCA.** Principal
430 component analysis carried out on 10 population genetics statistics, namely Nucleotide
431 diversity / Pi, Haplotype diversity, Fu and Li's D, Fu and Li's F, Tajima's D, Rozas' R₂,
432 Strobeck's S and the number of segregating sites. Panels are colored according to
433 categorical variables.

434

435 **Supplementary Figure 8. Positive selection landscape of the Pan-NLRome.** (a-e)
436 Absolute number of orthogroups in positive selection classes grouped by NLR class (a),
437 orthogroup type (b), presence of a non-canonical domain (c), presence of a paired NLR (d)
438 or NLR subclasses (e). An orthogroup was considered if at least one positive selected site of
439 a given class was detectable. (f-i) Domain coverage with positively selected sites grouped by
440 NLR class and positive selection type across canonical domains (f-g, i) and all aggregated
441 non-canonical domains (h).

442

443 **Supplementary Figure 9. Read and assembly statistics.** a) Read lengths distribution
444 (Q20-filtered CCS reads) for all accessions (black). The mean is shown as a solid black
445 horizontal line. The full densities are shown as a bean plot. The total number of CCS reads
446 (blue circles) and the total number of bases (orange diamonds) are plotted in addition. b)
447 Contig lengths distribution (black). The mean is shown as a solid black horizontal line and
448 the 95% Highest density Intervals (HDI: points in the interval have a higher probability than
449 points outside) are shown as solid bands around the sample mean. The full densities are
450 shown as a bean plot. Raw data points are plotted using black dots. The total assembly
451 sizes (orange circles) are plotted in addition. c) Quality (black) and completeness values
452 (orange) for sub-sampled Col-0 datasets. The amount of input data for each sub-sampling
453 experiment is shown as a second x axis. d) Quality (black) and completeness values
454 (orange) for all RenSeq accessions. Unfilled circles indicate accessions with qualities higher
455 than any sub-sampled dataset. The vertical black line is drawn at 95% completeness. e)
456 Correlations between the Assembly Quality, the amount of Input Reads, the amount of Input
457 Bases [bp], the read length N50 [bp], and the similarity to Col-0 are shown for the RenSeq
458 datasets. Histograms and kernel densities (red line) are plotted for each variable. Scatter
459 plots for variable pairs are shown together with a fitted line (red) and the Pearson's
460 correlation coefficient (significance 0 '****', 0.001 '***', 1 ' ').

461

462 **Supplementary Figure 10. Phylogenetic tree of NB domain alignments for refining**
463 **sensor/executor pairs from all pairs.** RPS4/RRS1-like and SOC3/CHS1-like paired TNLs
464 fall into distinct subclades. These are indicated by color: RPS4-like (Silver, executors),
465 RRS1-like (Gold, sensors), SOC3-like (Pink, executors) and CHS1-like (Bronze, sensors).
466 This phylogeny was constructed by aligning the NB domain (~240 amino acids) of all TIR

467 and NB containing Col-0 proteins and selected additional representatives of pair flagged
468 orthogroups (OGs) from the pan-NLRome that are not represented in Col-0 (identified by
469 their OG and protein numbers). NB domains from APAF1 (Human) and AT1G58602.1 (A.
470 thaliana CNL) were also included. Amino acid sequences were aligned with MUSCLE
471 (Neighbor joining clustering), refined by manual trimming and the phylogeny produced with
472 the WAG maximum likelihood method allowing for 3 discrete Gamma categories.
473 AT4G36140 contains two distinct NB domains, both of which were included and the second
474 of which clusters with other RRS1-like NB domains. Number of 100 bootstraps supporting
475 topology shown at major node vertices. Scale bar represents amino acid substitutions per
476 site.

477

478 **Supplementary Figure 11. Orthogroup (OG) size frequencies before OG refinement.**

479 Data are shown separately for the different NLR classes (TNL: yellow, NL: green, CNL: blue,
480 RNL: purple) and for OGs with (solid lines) and without (dashed lines) complex paralogs
481 (duplications spread across the whole phylogeny).

482 Supplementary Tables

- 483 Supplementary Table 1. 65 accessions metadata. Column content explained in the github
484 repository.
- 485 Supplementary Table 2. Domains found in the *A. thaliana* pan-NLRome. Column content
486 explained in the github repository.
- 487 Supplementary Table 3. Gene models included in the *A. thaliana* pan-NLRome. Column
488 content explained in the github repository.
- 489 Supplementary Table 4. Universe of IDs detected in the *A. thaliana* reference Col-0
490 accession and/or a 19 other Brassicaceae NLRomes.
- 491 Supplementary Table 5. Gene models used to generate the Brassicaceae NLRome.
492 Supplementary Table 6 Sources of the 22 proteomes from 19 Brassicaceae species used to
493 generate the Brassicaceae NLRome. see Online Methods.
- 494 Supplementary Table 7. Domain architecture prevalence, number of paired NLRs and
495 presence in the *A. thaliana* reference Col-0 accession and/or 19 other Brassicaceae
496 NLRomes.
- 497 Supplementary Table 8. Number of architectures detected for each of the 27 IDs already
498 known from 22 Brassicaceae NLRomes.
- 499 Supplementary Table 9. Architecture metadata. Columns are explained in the github
500 repository.
- 501 Supplementary Table 10. Number of non-reference OGs that can be anchored to Col-0
502 genomic positions through same contig co-occurrence with reference OGs.
- 503 Supplementary Table 11. Frequency of positively selected sites in putative integrated
504 domains.
- 505 Supplementary Table 12. List of orthogroups categorized as sensor/executor or control
506 pairs, with functional metadata and Tajima's D.
- 507 Supplementary Table 13. Oligos used to introduce custom barcoded adapters in TSL
508 accessions.
- 509 Supplementary Table 14. Bait library (v2.4) used to capture NLR-coding genes.
- 510 Supplementary Table 15. First and second subdivisions used to define NLR classes.
- 511 Supplementary Table 16. Prediction of coiled-oil motifs in functionally validated Col-0 CNLs.
- 512 Supplementary Table 17. Homology-based clustering into orthogroups without refinement
- 513 Supplementary Table 18. Pervasive diversifying positive selection posterior probabilities
- 514 Supplementary Table 19. Positional pervasive diversifying positive selection posterior
515 probabilities
- 516 Supplementary Table 20. Episodic diversifying positive selection p-values
- 517 Supplementary Table 21. Positional episodic diversifying positive selection p-values
- 518 Supplementary Table 22. Refined final Orthogroups
- 519 Supplementary Table 23. Enrichment of curation flags in orthogroups calculated. Only q-
520 values below 0.1 are reported.
- 521 Supplementary Table 24. Misannotated genes manually removed from the NLRome.
- 522 Supplementary Table 25. Non-NLR Orthogroups used for non-reference OG placement.
- 523

524 Data Availability Statement

525 Raw data and assembled sequences were deposited at the European Nucleotide Archive
526 (ENA) under accession number PRJEB23122. Genome browser is available at [http://ann-](http://ann-nblrrrome.tuebingen.mpg.de/annotator/index)
527 [nblrrrome.tuebingen.mpg.de/annotator/index](http://ann-nblrrrome.tuebingen.mpg.de/annotator/index). Manually curated gene models(gff), domain
528 annotations, orthogroups, protein and transcript alignments, phylogenetic trees, scripts
529 necessary to produce figures and further metadata files containing information parsed and
530 restructured from the supplemental tables in this manuscript are available at github
531 (<https://github.com/weigelworld/pan-nlrome/>). Visualization of OG phylogenetic trees and
532 metadata is available at iTOL (https://itol.embl.de/shared/pan_NLRome).

533 Acknowledgements

534 We thank Florian Jupe for contributing with methods before publication. Eunyong Chae for
535 contributing with alleles for bait design. Burkhard Steuernagel for assistance with
536 demultiplexing. Johannes Hofberger and Eric Schranz for providing sequences used in a
537 former version of the bait library. JLD is an Investigator of the Howard Hughes Medical
538 Institute (HHMI). This work was supported by a grant from the Gordon and Betty Moore
539 Foundation to the 2 Blades Foundation (GBMF4725), and the Howard Hughes Medical
540 Institute (HHMI). OF, JJ and KW acknowledge support from the Gatsby Charitable
541 Foundation.

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559 Glossary

560

pan-NLRome	NLR content of the 65 reseq'ed <i>A.thaliana</i> accessions.
NLR	Any gene encoding TIR, or NB-ARC, or RPW8 Pfam domains.
NLR classes	Classification scheme for NLRs. TNL (contains at least one TIR domain), RNL (>=1 RPW8 domain), CNL (>=1 CC domain), or NL (>=1 NB domain).
<i>R</i> gene	Functionally validated Disease Resistance gene (Pathogen/effector recognition).
RenSeq	<i>R</i> gene <u>en</u> richment and <u>seq</u> uencing.
SMRT	Single Molecule Real-Time.
Integrated domain (ID)	A predicted protein domain in addition to TIR, CC, NB, RPW8 and LRR within an NLR encoding gene model. May be gene prediction or domain prediction artefacts.
Cluster	In silico similarity-based clustering of genes/proteins.
Physically clustered NLRs	more than one NLR in 200 kb of genomic sequence.
Orthogroup (OG)	Homology-related NLRs from different ecotypes obtained using OrthAgogue and MCL clustering (see online methods for parameters used).
Pseudogene	Any gene with an exonerate minimal mapping score >= (50%) to a Col-0 gene annotated as pseudogene in Araport11.
Pseudo-genome	A genomic dataset which contains the TAIR10 genome without NLR genes (softmasked) and a RenSeq assembly as additional 'chromosomes'.
Paired NLRs	Two genetically linked NLRs that putatively function together to confer resistance.
Pair ratio	Number of manually flagged paired NLRs divided by the total number of NLRs in the same group (architecture, class, orthogroup, etc.).
Executor	NLR proteins that partner with sensor NLR(s). Presents typical domain arrangement and can trigger immune response ¹⁸ .
Sensor	NLR protein with a role in pathogen effector perception that may not be capable of triggering immunity without an executor partner ^{18,19} .
Helper	NLR protein with the domain architecture RPW8-NB-LRR. Have a role genetically defined as downstream of classic NLRs ^{59,61} .

Episodic (diversifying) positive selection	Selection at a specific site only affects a subset of lineages.
Pervasive (diversifying) positive selection	Selection at a specific site in all lineages.
Highest Density Interval (HDI)	points in the interval have a higher probability than points outside, analogous to 95% confidence intervals.

561 Author Contributions

562	AVDW	Anna-Lena Van de Weyer	orcid.org/0000-0002-5180-897X
563	DW	Detlef Weigel	orcid.org/0000-0002-2114-7963
564	FB	Felix Bemm	orcid.org/0000-0001-6557-4898
565	FM	Freddy Monteiro	orcid.org/0000-0002-9080-6715
566	JJ	Jonathan D. G. Jones	orcid.org/0000-0002-4953-261X
567	JLD	Jeffery L. Dangl	orcid.org/0000-0003-3199-8654
568	MN	Marc Nishimura	orcid.org/0000-0003-4666-6900
569	OF	Oliver J. Furzer	orcid.org/0000-0002-3536-9970
570	VC	Volkan Cevik	orcid.org/0000-0002-3545-3179
571	KW	Kamil Witek	orcid.org/0000-0003-0659-5562

572

573	Project Conception	JJ, JLD, DW
574	Project Management	FB, FM, AVDW, OF
575	Bait Design	OF, MN, VC
576	Pre-publication methods access	KW
577	Data Generation	AVDW, FM, MN, OF, VC
578	Data Preparation & Assembly	AVDW, FB, FM
579	Gene Annotation	AVDW, FB
580	Gene Curation	AVDW, FM, OF
581	Architecture Analysis	FM
582	Genomic Placement	FM
583	OG Co-occurrence network	FM
584	Pan-NLRome Generation	FB
585	Population Genetics Analysis	FB, OF
586	Positive Selection Analysis	FB
587	Metadata Collection	FB, AVDW, DW
588	Data Visualization (iTOL, JBrowse)	FB, AVDW
589	Initial draft manuscript	FB, AVDW, FM, OF
590	Manuscript Revision	FB, JLD, OF, JJ, FM, AVDW, DW, MN, VC

591 References

- 592 1. Jones, J. D. G. & Dangl, J. L. The plant immune system. *Nature* **444**, 323–329 (2006).
- 593 2. Gao, Y. *et al.* Out of Water: The Origin and Early Diversification of Plant -Genes. *Plant*
594 *Physiol.* **177**, 82–89 (2018).
- 595 3. Kourelis, J. & van der Hoorn, R. A. L. Defended to the Nines: 25 Years of Resistance
596 Gene Cloning Identifies Nine Mechanisms for R Protein Function. *Plant Cell* **30**, 285–
597 299 (2018).
- 598 4. Shao, Z.-Q. *et al.* Large-Scale Analyses of Angiosperm Nucleotide-Binding Site-
599 Leucine-Rich Repeat Genes Reveal Three Anciently Diverged Classes with Distinct
600 Evolutionary Patterns. *Plant Physiol.* **170**, 2095–2109 (2016).
- 601 5. van der Biezen, E. A. & Jones, J. D. The NB-ARC domain: a novel signalling motif
602 shared by plant resistance gene products and regulators of cell death in animals. *Curr.*
603 *Biol.* **8**, R226–7 (1998).
- 604 6. Takken, F. L. W. & Goverse, A. How to build a pathogen detector: structural basis of
605 NB-LRR function. *Curr. Opin. Plant Biol.* **15**, 375–384 (2012).
- 606 7. Maekawa, T., Kufer, T. A. & Schulze-Lefert, P. NLR functions in plant and animal
607 immune systems: so far and yet so close. *Nat. Immunol.* **12**, 817–826 (2011).
- 608 8. Nishimura, M. T. *et al.* TIR-only protein RBA1 recognizes a pathogen effector to
609 regulate cell death in. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E2053–E2062 (2017).
- 610 9. Bernoux, M. *et al.* Structural and Functional Analysis of a Plant Resistance Protein TIR
611 Domain Reveals Interfaces for Self-Association, Signaling, and Autoregulation. *Cell*
612 *Host Microbe* **9**, 200–211 (2011).
- 613 10. Qi, D., DeYoung, B. J. & Innes, R. W. Structure-function analysis of the coiled-coil and
614 leucine-rich repeat domains of the RPS5 disease resistance protein. *Plant Physiol.* **158**,
615 1819–1832 (2012).
- 616 11. Li, X., Kapos, P. & Zhang, Y. NLRs in plants. *Curr. Opin. Immunol.* **32**, 114–121 (2015).

- 617 12. Jones, J. D. G., Vance, R. E. & Dangl, J. L. Intracellular innate immune surveillance
618 devices in plants and animals. *Science* **354**, (2016).
- 619 13. Uehling, J., Deveau, A. & Paoletti, M. Do fungi have an innate immune response? An
620 NLR-based comparison to plant and animal immune systems. *PLoS Pathog.* **13**,
621 e1006578 (2017).
- 622 14. Dangl, J. L., Horvath, D. M. & Staskawicz, B. J. Pivoting the plant immune system from
623 dissection to deployment. *Science* **341**, 746–751 (2013).
- 624 15. Dodds, P. N. *et al.* Direct protein interaction underlies gene-for-gene specificity and
625 coevolution of the flax resistance genes and flax rust avirulence genes. *Proc. Natl.*
626 *Acad. Sci. U. S. A.* **103**, 8888–8893 (2006).
- 627 16. Catanzariti, A.-M. *et al.* The AvrM Effector from Flax Rust Has a Structured C-Terminal
628 Domain and Interacts Directly with the M Resistance Protein. *Mol. Plant. Microbe.*
629 *Interact.* **23**, 49–57 (2010).
- 630 17. Krasileva, K. V., Dahlbeck, D. & Staskawicz, B. J. Activation of an Arabidopsis
631 resistance protein is specified by the in planta association of its leucine-rich repeat
632 domain with the cognate oomycete effector. *Plant Cell* **22**, 2444–2458 (2010).
- 633 18. Wu, C.-H., Krasileva, K. V., Banfield, M. J., Terauchi, R. & Kamoun, S. The ‘sensor
634 domains’ of plant NLR proteins: more than decoys? *Front. Plant Sci.* **6**, 134 (2015).
- 635 19. Cesari, S., Bernoux, M., Moncuquet, P., Kroj, T. & Dodds, P. N. A novel conserved
636 mechanism for plant NLR protein pairs: the ‘integrated decoy’ hypothesis. *Front. Plant*
637 *Sci.* **5**, 606 (2014).
- 638 20. Le Roux, C. *et al.* A receptor pair with an integrated decoy converts pathogen disabling
639 of transcription factors to immunity. *Cell* **161**, 1074–1088 (2015).
- 640 21. Sarris, P. F. *et al.* A Plant Immune Receptor Detects Pathogen Effectors that Target
641 WRKY Transcription Factors. *Cell* **161**, 1089–1100 (2015).
- 642 22. Maqbool, A. *et al.* Structural basis of pathogen recognition by an integrated HMA

- 643 domain in a plant NLR immune receptor. *Elife* **4**, (2015).
- 644 23. Mackey, D., Holt, B. F., 3rd, Wiig, A. & Dangl, J. L. RIN4 interacts with *Pseudomonas*
645 *syringae* type III effector molecules and is required for RPM1-mediated resistance in
646 *Arabidopsis*. *Cell* **108**, 743–754 (2002).
- 647 24. Qi, D. *et al.* Recognition of the protein kinase AVRPPHB SUSCEPTIBLE1 by the
648 disease resistance protein RESISTANCE TO PSEUDOMONAS SYRINGAE5 is
649 dependent on s-acylation and an exposed loop in AVRPPHB SUSCEPTIBLE1. *Plant*
650 *Physiol.* **164**, 340–351 (2014).
- 651 25. Wang, G. *et al.* The Decoy Substrate of a Pathogen Effector and a Pseudokinase
652 Specify Pathogen-Induced Modified-Self Recognition and Immunity in Plants. *Cell Host*
653 *Microbe* **18**, 285–295 (2015).
- 654 26. Kroj, T., Chanclud, E., Michel-Romiti, C., Grand, X. & Morel, J.-B. Integration of decoy
655 domains derived from protein targets of pathogen effectors into plant immune receptors
656 is widespread. *New Phytol.* **210**, 618–626 (2016).
- 657 27. Sarris, P. F., Cevik, V., Dagdas, G., Jones, J. D. G. & Krasileva, K. V. Comparative
658 analysis of plant immune receptor architectures uncovers host proteins likely targeted
659 by pathogens. *BMC Biol.* **14**, 8 (2016).
- 660 28. Bailey, P. C. *et al.* Dominant integration locus drives continuous diversification of plant
661 immune receptors with exogenous domain fusions. *Genome Biol.* **19**, 23 (2018).
- 662 29. Stein, J. C. *et al.* Genomes of 13 domesticated and wild rice relatives highlight genetic
663 conservation, turnover and innovation across the genus *Oryza*. *Nat. Genet.* **50**, 285–296
664 (2018).
- 665 30. Kim, S. H., Qi, D., Ashfield, T., Helm, M. & Innes, R. W. Using decoys to expand the
666 recognition specificity of a plant disease resistance protein. *Science* **351**, 684–687
667 (2016).
- 668 31. Kourelis, J., van der Hoorn, R. A. L. & Sueldo, D. J. Decoy Engineering: The Next Step

- 669 in Resistance Breeding. *Trends Plant Sci.* **21**, 371–373 (2016).
- 670 32. Nishimura, M. T., Monteiro, F. & Dangl, J. L. Treasure your exceptions: unusual
671 domains in immune receptors reveal host virulence targets. *Cell* **161**, 957–960 (2015).
- 672 33. Monteiro, F. & Nishimura, M. T. Structural, Functional, and Genomic Diversity of Plant
673 NLR Proteins: An Evolved Resource for Rational Engineering of Plant Immunity. *Annu.*
674 *Rev. Phytopathol.* **56**, 243–267 (2018).
- 675 34. Zhang, M. *et al.* Numbers of genes in the NBS and RLK families vary by more than four-
676 fold within a plant species and are regulated by multiple factors. *Nucleic Acids Res.* **38**,
677 6513–6525 (2010).
- 678 35. Stam, R., Scheikl, D. & Tellier, A. Pooled Enrichment Sequencing Identifies Diversity
679 and Evolutionary Pressures at NLR Resistance Genes within a Wild Tomato Population.
680 *Genome Biol. Evol.* **8**, 1501–1515 (2016).
- 681 36. Guo, Y.-L. *et al.* Genome-wide comparison of nucleotide-binding site-leucine-rich
682 repeat-encoding genes in *Arabidopsis*. *Plant Physiol.* **157**, 757–769 (2011).
- 683 37. Bakker, E. G., Traw, M. B., Toomajian, C., Kreitman, M. & Bergelson, J. Low levels of
684 polymorphism in genes that control the activation of defense response in *Arabidopsis*
685 *thaliana*. *Genetics* **178**, 2031–2043 (2008).
- 686 38. Cao, J. *et al.* Whole-genome sequencing of multiple *Arabidopsis thaliana* populations.
687 *Nat. Genet.* **43**, 956–963 (2011).
- 688 39. Duan, N. *et al.* Genome re-sequencing reveals the history of apple and supports a two-
689 stage model for fruit enlargement. *Nat. Commun.* **8**, 249 (2017).
- 690 40. Choi, K. *et al.* Recombination Rate Heterogeneity within *Arabidopsis* Disease
691 Resistance Genes. *PLoS Genet.* **12**, e1006179 (2016).
- 692 41. Rose, L. E. *et al.* The maintenance of extreme amino acid diversity at the disease
693 resistance gene, *RPP13*, in *Arabidopsis thaliana*. *Genetics* **166**, 1517–1527 (2004).
- 694 42. Golicz, A. A. *et al.* The pangenome of an agronomically important crop plant *Brassica*

- 695 oleracea. *Nat. Commun.* **7**, 13390 (2016).
- 696 43. Kawakatsu, T. *et al.* Epigenomic Diversity in a Global Collection of *Arabidopsis thaliana*
697 Accessions. *Cell* **166**, 492–505 (2016).
- 698 44. Li, Y.-H. *et al.* De novo assembly of soybean wild relatives for pan-genome analysis of
699 diversity and agronomic traits. *Nat. Biotechnol.* **32**, (2014).
- 700 45. Witek, K. *et al.* Accelerated cloning of a potato late blight–resistance gene using
701 RenSeq and SMRT sequencing. *Nat. Biotechnol.* **34**, (2016).
- 702 46. 1001 Genomes Consortium. 1,135 Genomes Reveal the Global Pattern of
703 Polymorphism in *Arabidopsis thaliana*. *Cell* **166**, 481–491 (2016).
- 704 47. Kover, P. X. *et al.* A Multiparent Advanced Generation Inter-Cross to fine-map
705 quantitative traits in *Arabidopsis thaliana*. *PLoS Genet.* **5**, e1000551 (2009).
- 706 48. Scarcelli, N., Cheverud, J. M., Schaal, B. A. & Kover, P. X. Antagonistic pleiotropic
707 effects reduce the potential adaptive value of the FRIGIDA locus. *Proc. Natl. Acad. Sci.*
708 *U. S. A.* **104**, 16986–16991 (2007).
- 709 49. Holub, E. B. The arms race is ancient history in *Arabidopsis*, the wildflower. *Nat. Rev.*
710 *Genet.* **2**, 516–527 (2001).
- 711 50. Narusaka, M. *et al.* RRS1 and RPS4 provide a dual Resistance-gene system against
712 fungal and bacterial pathogens. *Plant J.* **60**, 218–226 (2009).
- 713 51. Saucet, S. B. *et al.* Two linked pairs of *Arabidopsis* TNL resistance genes independently
714 confer recognition of bacterial effector AvrRps4. *Nat. Commun.* **6**, 6338 (2015).
- 715 52. Mukhtar, M. S. *et al.* Independently evolved virulence effectors converge onto hubs in a
716 plant immune system network. *Science* **333**, 596–601 (2011).
- 717 53. Weßling, R. *et al.* Convergent targeting of a common host protein-network by pathogen
718 effectors from three kingdoms of life. *Cell Host Microbe* **16**, 364–375 (2014).
- 719 54. Peele, H. M., Guan, N., Fogelqvist, J. & Dixelius, C. Loss and retention of resistance
720 genes in five species of the Brassicaceae family. *BMC Plant Biol.* **14**, 1–11 (2014).

- 721 55. Ohno, S. *Evolution by Gene Duplication*. (1970).
- 722 56. Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA
723 polymorphism. *Genetics* **123**, 585–595 (1989).
- 724 57. Bittner-Eddy, P. D., Crute, I. R., Holub, E. B. & Beynon, J. L. RPP13 is a simple locus in
725 *Arabidopsis thaliana* for alleles that specify downy mildew resistance to different
726 avirulence determinants in *Peronospora parasitica*. *Plant J.* **21**, 177–188 (2000).
- 727 58. Cevik, V. *et al.* Transgressive segregation reveals mechanisms of Arabidopsis immunity
728 to Brassica-infecting races of white rust (*Albugo candida*). *Proceedings of the National*
729 *Academy of Sciences* **adv. online**, 10.1073/pnas.1812911116 (2019).
- 730 59. Bonardi, V. *et al.* Expanded functions for a family of plant intracellular immune receptors
731 beyond specific recognition of pathogen effectors. *Proceedings of the National Academy*
732 *of Sciences* **108**, 16463–16468 (2011).
- 733 60. Wu, C.-H. *et al.* NLR network mediates immunity to diverse plant pathogens.
734 *Proceedings of the National Academy of Sciences* 201702041 (2017).
- 735 61. Castel, B. *et al.* Diverse NLR immune receptors activate defence via the RPW8-NLR
736 NRG1. *New Phytol.* (2018). doi:10.1111/nph.15659
- 737 62. Stahl, E. A., Dwyer, G., Mauricio, R., Kreitman, M. & Bergelson, J. Dynamics of disease
738 resistance polymorphism at the Rpm1 locus of Arabidopsis. *Nature* **400**, 667–671
739 (1999).
- 740 63. Xu, F. *et al.* Autoimmunity conferred by chs3-2D relies on CSA1, its adjacent TNL-
741 encoding neighbour. *Sci. Rep.* **5**, 8792 (2015).
- 742 64. Noël, L. *et al.* Pronounced Intraspecific Haplotype Divergence at the RPP5 Complex
743 Disease Resistance Locus of Arabidopsis. *Plant Cell* **11**, 2099–2111 (1999).
- 744 65. Zhang, Y. *et al.* Temperature-dependent autoimmunity mediated by chs1 requires its
745 neighboring TNL gene SOC3. *New Phytol.* **213**, 1330–1345 (2017).
- 746 66. Ma, Y. *et al.* Distinct modes of derepression of an Arabidopsis immune receptor

- 747 complex by two different bacterial effectors. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 10218–
748 10227 (2018).
- 749 67. Yang, L. *et al.* *Pseudomonas syringae* Type III Effector HopBB1 Promotes Host
750 Transcriptional Repressor Degradation to Regulate Phytohormone Responses and
751 Virulence. *Cell Host Microbe* **21**, 156–168 (2017).
- 752 68. Sugio, A., MacLean, A. M. & Hogenhout, S. A. The small phytoplasma virulence effector
753 SAP11 contains distinct domains required for nuclear targeting and CIN-TCP binding
754 and destabilization. *New Phytol.* **202**, 838–848 (2014).
- 755 69. Helm, M. *et al.* Engineering a decoy substrate in soybean to enable recognition of the
756 Soybean Mosaic Virus NIa protease. *Mol. Plant. Microbe. Interact.* **adv. online**
757 doi:10.1094/MPMI-12-18-0324-R (2019).

Figure 1.

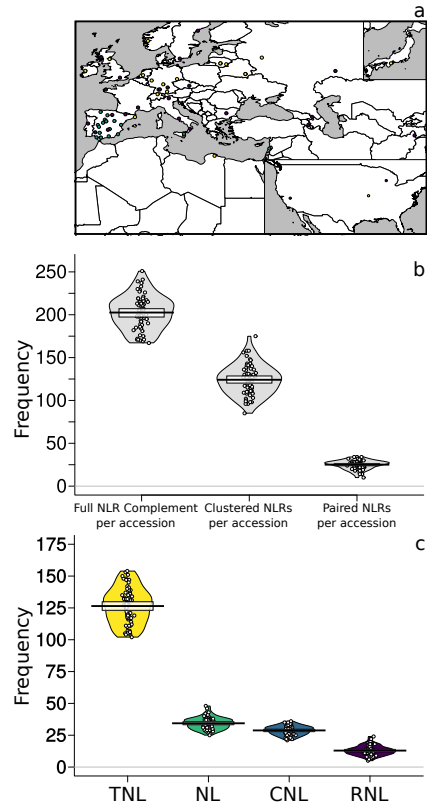


Figure 1. Overview of NLR complements in 65 accessions. a) Map of accession provenance. 1001 Genomes (relicts, blue, non-relicts, purple), MAGIC founders (yellow). b) Number of total, clustered and paired NLRs in accessions. Means, solid black lines; Bayesian 95% Highest Density Intervals (HDIs), solid bands. Individual data, open circles; full densities shown as bean plots. c) Number of different structural classes in accessions. Mean, HDI, and individual data as in (b).

Figure 2.

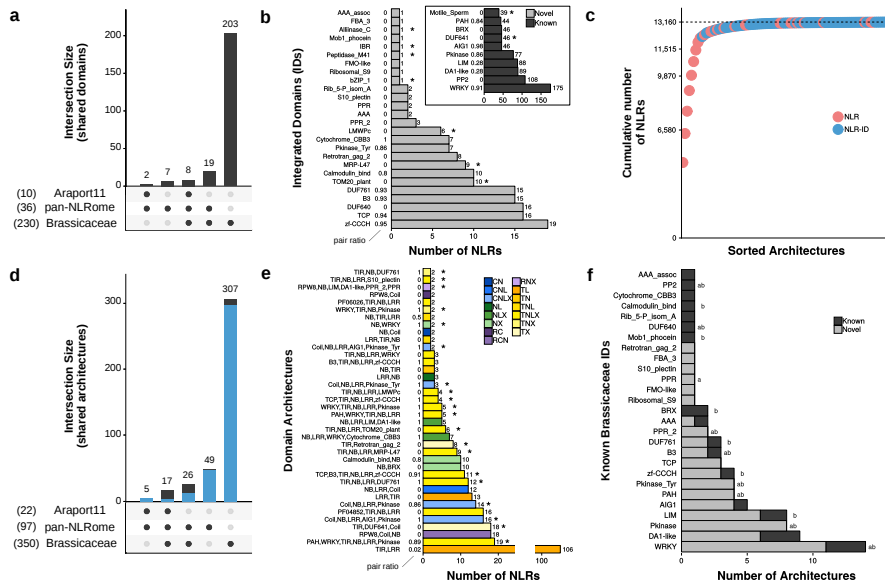


Figure 2. Diversity of IDs and domain architectures. a) UpSet intersection of IDs in the pan-NLRome, the Col-0 reference accession and 19 other Brassicaceae. The number of IDs in each set is indicated between parentheses at the lower left. Set intersections shown on the bottom. b) ID analysis. IDs known for *A. thaliana* in black and IDs newly described for *A. thaliana* in light grey. Asterisks indicate IDs not found in other Brassicaceae. Numbers next to y-axis indicate ratio of paired NLRs among ID containing NLRs. c) Cumulative size contribution to the pan-NLRome (y-axis) of each of the 97 size sorted domain architectures (x-axis). d) UpSet intersection plot of architectures shared between pan-NLRome, *A. thaliana* reference- and 19 other Brassicaceae. e) 38 architectures with at least two representatives and not found in the Col-0 reference. Asterisks indicate 20 architectures not found in 19 Brassicaceae family genomes or the *A. thaliana* Col-0 reference. Numbers next to y-axis shows fraction of paired NLRs. f) Number of known and novel architectures containing the 27 overlapping Brassicaceae IDs (see a). "a" and "b" to the right of the bars indicate putative IDs ^{26,27} (see also Supplementary Table 2).

Figure 3.

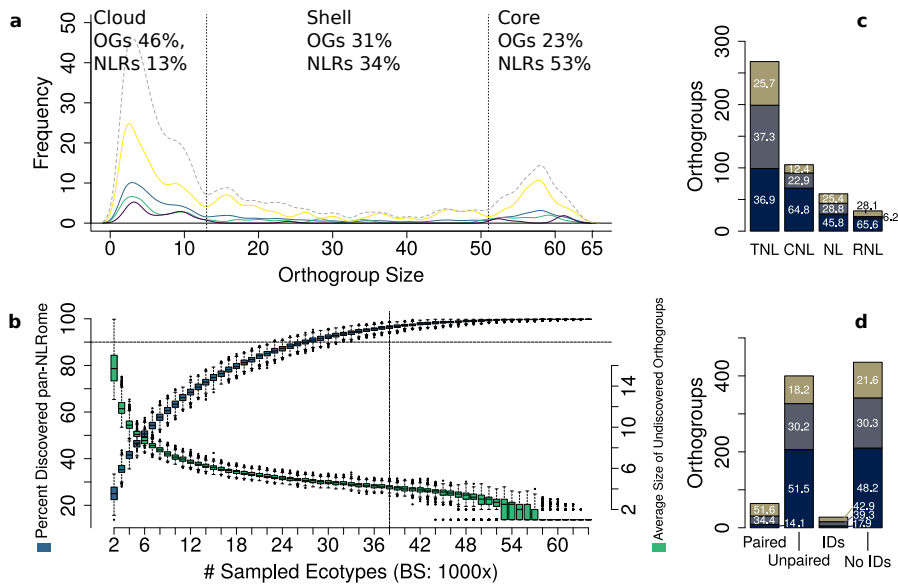


Figure 3. OG sizes, saturation, and distribution of core-, shell- and cloud-NLRs. a) OG size distribution for TNLs (yellow), NLs (green), CNLs (blue), RNLs (purple), and all NLRs (grey dashed line). The vertical lines at $x=13$ and $x=51$ differentiate cloud, shell and core. b) Saturation of the pan-NLRome. The blue boxes show the percentage of the pan-NLRome that can be recovered when randomly drawing a fixed number of accessions (1000x bootstrapping). The horizontal dashed line indicates where 90% of the pan-NLRome is found. The green boxes shown for each subset of drawn accessions indicate the average size of undiscovered orthogroups. The vertical dashed line indicates that 95% of the pan-NLRome can be recovered with 38 accessions. c) OG-type specific distribution of NLR classes in the cloud (dark blue), the shell (grey), and the core (olive green), and the relative fraction (numbers in the bars). d) OG-type specific distribution of paired and unpaired NLRs, and NLRs with and without IDs in the cloud (dark blue), the shell (grey), and the core (olive green), and the percentage (numbers in the bars).

Figure 4.

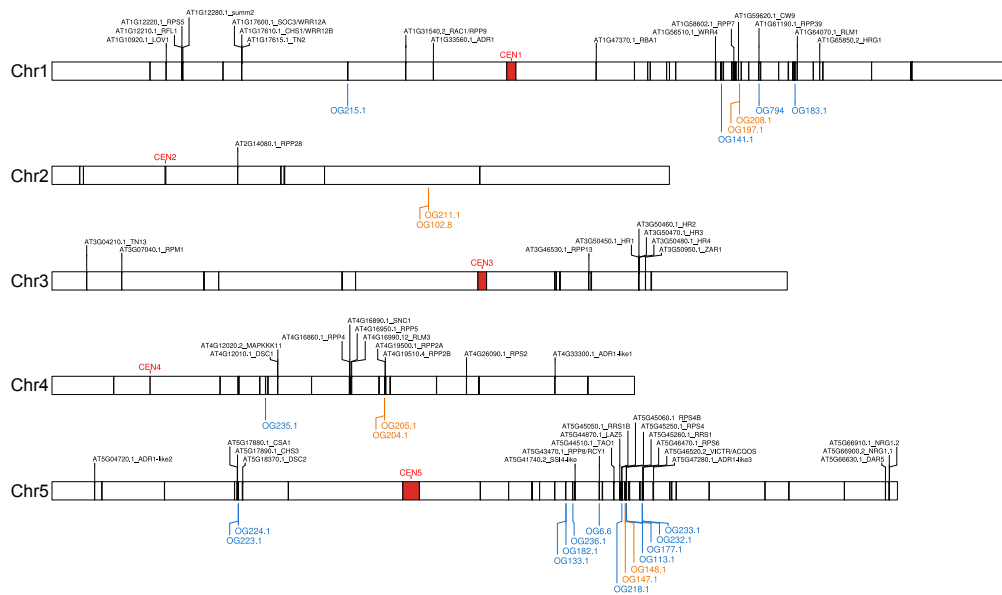


Figure 4. Genomic location of *NLR* genes in the reference assembly. The five *A. thaliana* chromosomes are shown as horizontal bars with centromeres in red. Reference NLRs are shown as black line segments. Text labels are shown only for functionally defined Col-0 NLRs. Anchored OGs found in at least 10 accessions are shown below each chromosome. Anchored OGs with paired NLR members are shown in orange, while remaining anchored OGs are shown in blue.

Figure 5.

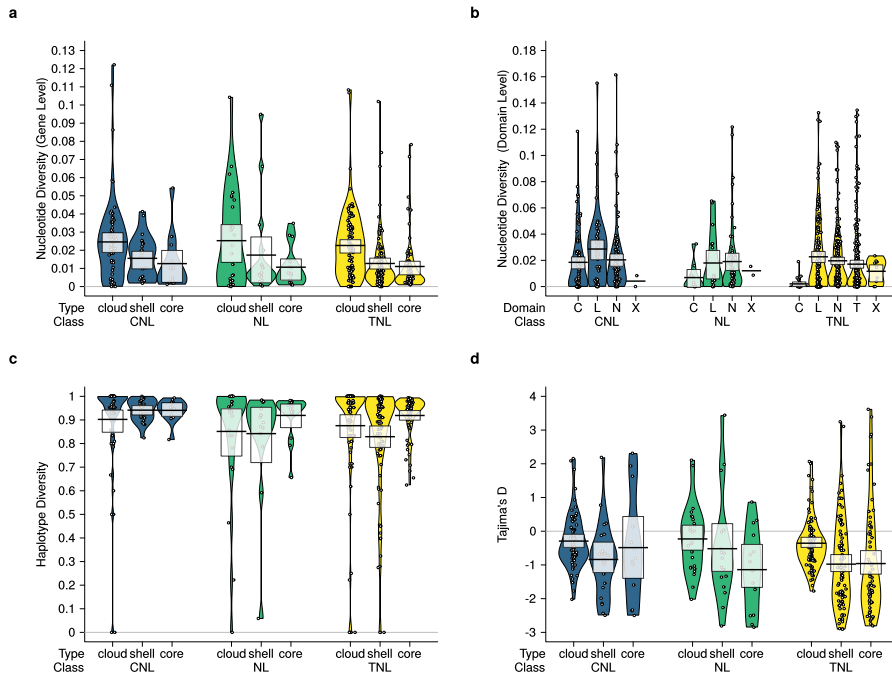


Figure 5. Population genetic statistics across the Pan-NLRome. a) Nucleotide diversity distribution grouped by orthogroup type and NLR class. Nucleotide diversity was defined as average pairwise nucleotide difference normalized to the number of sites in the respective orthogroup. b) Nucleotide diversity distribution grouped by domain type and NLR class. Sites of each orthogroup annotated with the same domain were aggregated. Type C within class NL occurs where a minority of OG members had an identifiable CC-domain. c) Haplotype diversity distribution grouped by orthogroup type and NLR class. Haplotype diversity was defined as average pairwise haplotype differences. A value of 1 indicates a high chance of finding two different haplotypes for two randomly chosen sequences of a given orthogroup. d) Tajima's D (a measure of genetic selection) distribution grouped by orthogroup type and NLR class. RNL orthogroups are not shown because of the low number of orthogroups showing this class.

Figure 6.

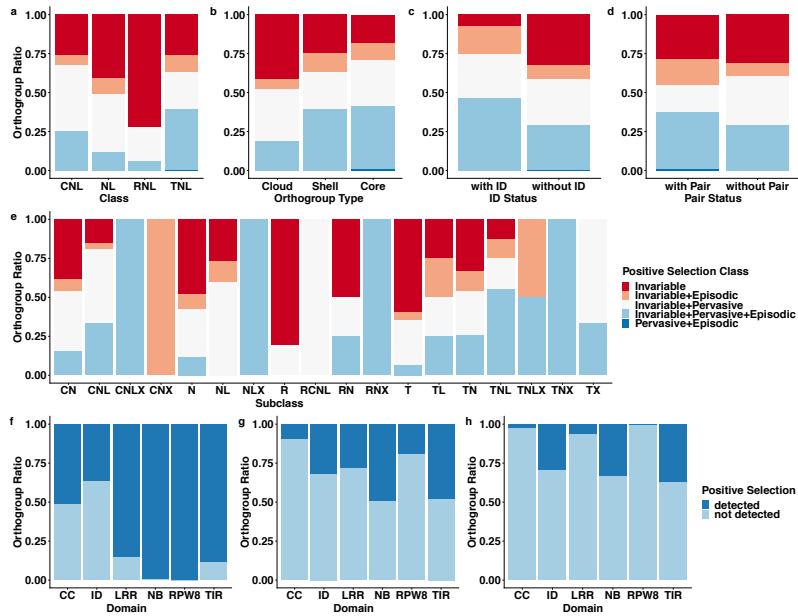


Figure 6. Positive selection landscape of the Pan-NLRome. (a-e) Ratio of different positive selection classes grouped by NLR class (a), orthogroup type (b), presence of a non-canonical domain (c), presence of a paired NLR (d) or NLR subclasses (e). An orthogroup was considered if at least one positive selected site of a given class was detectable. (f-h) Ratio between orthogroups showing constant (f), pervasive (g) or episodic (h) selection or no selection grouped by annotated protein domains.

Figure 7.

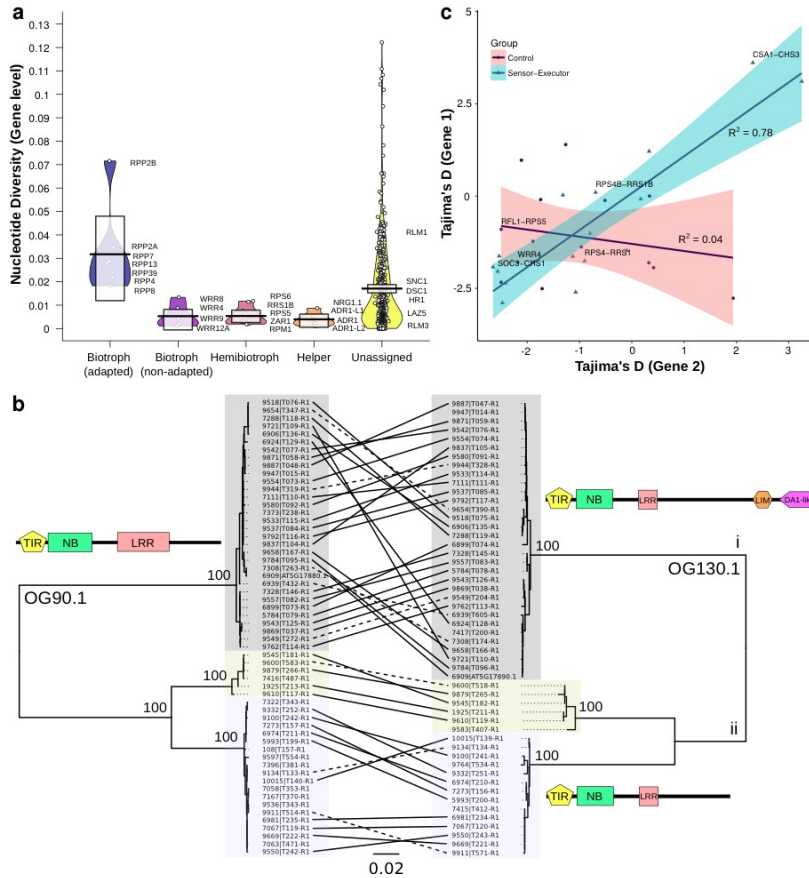


Figure 7. Population genetics of different OG classes grouped by known resistances and pairs. a) Nucleotide diversity distributions by functional class according to pathogen type to which they provide resistance. b) Correlation of Tajima's D values in sensor/executor and control pairs. c) Maximum-likelihood phylogenetic trees of two OGs 90.1 and 130.1, which form a sensor/executor pair⁶² (100 bootstrap support indicated at major nodes). Scale bar refers to substitutions per site. Lines connecting the trees denote same accession.