

UC Berkeley

UC Berkeley Previously Published Works

Title

Analysis of intraspecies diversity reveals a subset of highly variable plant immune receptors and predicts their binding sites

Permalink

<https://escholarship.org/uc/item/6j80q7xk>

Journal

The Plant Cell, 33(4)

ISSN

1040-4651

Authors

Prigozhin, Daniil M
Krasileva, Ksenia V

Publication Date

2021-05-31

DOI

10.1093/plcell/koab013

Peer reviewed

RESEARCH ARTICLE

Analysis of intraspecies diversity reveals a subset of highly variable plant immune receptors and predicts their binding sites

Daniil M Prigozhin¹ and Ksenia V Krasileva²

¹Berkeley Center for Structural Biology, Molecular Biophysics and Integrated Bioimaging Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720

²Department of Plant and Microbial Biology, University of California, Berkeley, CA 94720
Corresponding Authors: kseniak@berkeley.edu and daniilprigozhin@lbl.gov

Short title: Highly variable NLR immune receptors

One-sentence summary: NLR immune receptor complements of 62 ecotypes of *A. thaliana* and 54 lines of *B. distachyon* help identify highly variable NLR subfamilies responsible for the generation of new receptor specificities.

The author responsible for distribution of materials integral to the findings presented in this article in accordance with the policy described in the Instructions for Authors (www.plantcell.org) is: Ksenia V. Krasileva (kseniak@berkeley.edu).

ABSTRACT

The evolution of recognition specificities by the immune system depends on the generation of receptor diversity and on connecting the binding of new antigens with the initiation of downstream signaling. In plant immunity, the innate Nucleotide-Binding Leucine Rich Repeat (NLR) receptor family enables antigen binding and immune signaling. In this study, we surveyed the NLR complements of 62 ecotypes of *Arabidopsis thaliana* and 54 lines of *Brachypodium distachyon* and identified a limited number of NLR subfamilies that show high allelic diversity. We show that the predicted specificity-determining residues cluster on the surfaces of Leucine Rich Repeat domains, but the locations of the clusters vary among NLR subfamilies. By comparing NLR phylogeny, allelic diversity, and known functions of the Arabidopsis NLRs, we formulate a hypothesis for the emergence of direct and indirect pathogen-sensing receptors and of the autoimmune NLRs. These findings reveal the recurring patterns of evolution of innate immunity and can inform NLR engineering efforts.

1 INTRODUCTION

2 Plants lack the adaptive immunity of vertebrates. With their immune receptor specificities
3 encoded in the germline, plants can achieve remarkable receptor diversity at the population level
4 (Bakker et al., 2006). The mechanisms that generate this diversity and select for useful (and
5 against deleterious) receptor variants are thus of great importance to both basic science and crop
6 improvement (Dangl et al., 2013). Ongoing efforts at pan-genome sequencing of both model and
7 crop species reveal the intraspecies diversity of plant immune receptors, their natural history,

8 mechanisms of action, and the evolutionary forces that shape plant immunity (Van de Weyer et
9 al., 2019; Stam et al., 2019a, 2019b; Seong et al., 2020; Gordon et al., 2017).

10
11 Two types of plant immune receptors form the basis of pathogen recognition: extracellular
12 receptors, including receptor-like kinases (RLK) and receptor-like proteins (RLP); and
13 intracellular Nucleotide-binding Leucine Rich Repeat (NLR) proteins (Dangl et al., 2013). While
14 RLKs and RLPs monitor the extracellular environments of plants, NLRs are cell death-executing
15 receptors that are shared across the plant and animal kingdoms (Jones et al., 2016). Plant NLRs
16 are typically composed of three domains, including a central Nucleotide Binding (NB-ARC)
17 domain that mediates receptor oligomerization upon activation, the C-terminal Leucine Rich
18 Repeat (LRR) domain that defines receptor specificity, and one of three N-terminal domains:
19 Resistance To Powdery Mildew 8 (RPW8), Coiled-Coil (CC), or Toll/Interleukin-1 Receptor
20 homology (TIR) domains, which mediate the immune effector function. NLRs are divided into
21 three monophyletic classes based on the N-terminal domains and their evolutionary origin: RNL,
22 CNL, and TNL (Shao et al., 2016).

23
24 NLRs can function as sensors or signal transducers (helpers) (Wu et al., 2017). For example, all
25 RNL genes are thought to be helpers (Jubic et al., 2019), while TNLs and CNLs can fulfill either
26 function. Sensor NLRs recognize pathogens using three main modes: i) direct binding to the
27 pathogen-derived effector molecules, ii) indirect recognition of effector activities on other plant
28 proteins, and iii) recognition of modifications to a non-canonical integrated domain of the NLR,
29 which acts as a bait for the effector (Cesari, 2018). The recognition mode of a given sensor NLR
30 is likely to have a large effect on the evolutionary pressure it experiences. Indirect recognition
31 NLRs likely undergo balancing or purifying selection based on the monitoring of conserved
32 effector activity. By contrast, effector recognition upon direct binding likely requires NLRs to
33 adapt rapidly to keep track of easy-to-mutate effector surface residues. Among the best studied
34 NLRs that directly bind pathogen-derived effectors are the flax (*Linum usitatissimum*) L genes
35 (Ellis et al., 2007; Catanzariti et al., 2010), the MLA/Sr50 locus in barley (*Hordeum vulgare*)
36 and wheat (*Triticum* spp.) (Chen et al., 2017; Saur et al., 2019), and the RPP1 genes in
37 Arabidopsis (Krasileva et al., 2010; Goritschnig et al., 2016). Their effector targets are

38 structurally diverse, suggesting that the current recognition specificities of individual alleles are
39 recently derived, rather than ancestral.

40
41 The continuous generation of diversity in sensor NLRs is required to provide protection from
42 diverse pathogens and is thought to result from divergent (diversifying) selection and a birth-
43 and-death process acting on NLR gene clusters (Michelmore and Meyers, 1998). NLRs diversify
44 through copy number variation, recombination, gene conversion, gene fusion, and point
45 mutations (Baggs et al., 2017). In a subset of NLRs, these mechanisms combine to produce an
46 astounding array of alleles (Bakker et al., 2006; Ding et al., 2007). Not unexpectedly, such
47 diversity comes at a price. Hybrid necrosis has been observed widely in inbreeding and
48 outcrossing plants in both cultivated and wild populations and can be considered a plant version
49 of autoimmunity (Bomblies, 2009). Hybrid necrosis occurs due to a mismatch between immune
50 receptor variants and other plant genes, leading to autoimmune recognition, as exemplified by
51 Dangerous Mix genes in *A. thaliana* (Bomblies et al., 2007; Chae et al., 2014; Atanasov et al.,
52 2018) and *Ne* genes in wheat (Zhang et al., 2016). Tomato (*Solanum lycopersicum*) *Cf-2* is an
53 example of a non-NLR immune receptor that shows this type of phenotype (Kruger, 2002;
54 Santangelo et al., 2003). These negative interactions revealed in crosses are likely only a small
55 fraction of the cost of derivation of new immune specificities in the presence of the whole
56 intracellular plant proteome.

57
58 Cross-species phylogenetic analyses of the NLR gene family have provided important insights
59 into NLR evolution. A combined phylogeny of maize (*Zea mays*), sorghum (*Sorghum bicolor*),
60 brachypodium, and rice (*Oryza sativa*) NLRs was used to identify recently derived NLR immune
61 specificities against rice blast disease (Yang et al., 2013). An expansion of a network of helper
62 and sensor NLRs was identified across asterids in which a set of diverse sensors signal through a
63 redundant set of helpers that show reduced diversity (Wu et al., 2017). Phylogenetic analyses in
64 grasses identified major integration clades of NLRs that incorporate additional domains that
65 serve as baits for pathogens (Bailey et al., 2018). In view of the recent progress in elucidating the
66 intra-species NLR complements of both model and non-model plants (Van de Weyer et al., 2019;
67 Stam et al., 2019a, 2019b; Seong et al., 2020; Gordon et al., 2017), a systematic analysis is

68 needed to uncover the relationships between NLR phylogeny, mode of recognition, and the
69 amount of allelic diversity.

70

71 The recent elucidation of both the pre-activation monomeric and activated resistosome-forming
72 conformations of ZAR1, an indirect recognition CNL, dramatically improved our understanding
73 of both target binding and the receptor activation mechanisms of NLRs (Wang et al., 2019b,
74 2019a). The structures of Roq1 and RPP1, both direct binders, in complex with their targets,
75 were recently revealed (Martin et al., 2020; Ma et al., 2020), further shedding light on LRR and
76 post-LRR domain-dependent target recognition and downstream TIR domain activation. While
77 more NLR structures are likely to be revealed in the future, structure determination efforts will
78 likely lag behind the pan-genome sequencing due to the cost and difficulty of the experiments
79 involved. Therefore, the prediction of the mode of recognition and specificity-determining
80 residues of NLRs based on sequence data is an attractive direction that is yet to be fully explored.
81 The idea that highly variable residues determine immune receptor specificity predates the
82 elucidation of the first antibody structure by three years (Kabat, 1970). In the subsequent
83 decades, several measures of amino-acid diversity were advanced. Shannon entropy, which
84 originated in information theory, is given by the formula:

$$H = - \sum_{i=1}^{20} p_i \log_2 p_i$$

85 where p_i is the fraction of one of the twenty amino acids in a column of a protein sequence
86 alignment. This measure was first applied to study residues that determine antibody and T-cell
87 receptor specificity (Shenkin et al., 1991; Stewart et al., 1997). High entropy values correlate
88 strongly with surface exposure and hydrophilic character (Liao et al., 2005) and can be used to
89 predict rapidly evolving ligand-binding sites (Magliery and Regan, 2005). In addition to B- and
90 T-cell receptors, entropy-based measures have been applied to identify binding sites in TRP
91 repeat proteins, ankyrin repeat proteins, Zn-finger transcription factors, and G protein-coupled
92 receptors (Sanders et al., 2011; Magliery and Regan, 2005).

93

94 In the current study, we used phylogenetic analyses to group *Arabidopsis* and *Brachypodium*
95 NLRs into near allelic series and applied Shannon entropy analyses of protein alignments to
96 define highly variable NLRs and their candidate specificity-determining residues. Our results

97 show that, depending on the ecotype, 15 to 35 Arabidopsis NLRs belong to rapidly diversifying
98 families. These families are distributed in the NLR phylogeny among both CC- and TIR-
99 containing NLRs and encompass the known Dangerous Mix NLRs. We further show that in the
100 highly variable NLRs (hvNLRs), the highly variable residues identified by Shannon entropy
101 cluster on the surface of the LRR domain and contain surface-exposed hydrophobic residues,
102 thus identifying likely binding sites. The exact location of the putative binding sites on the LRR
103 surface is not conserved across different NLRs. Based on the phylogenetic distribution of
104 hvNLRs, we formulate a hypothesis regarding the origin of indirect recognition sensor NLRs.
105 When applied to *Brachypodium distachyon* pan-genome, our methods reveal a similarly
106 dispersed phylogenetic distribution of highly variable NLRs in this model grass species.
107 Collectively, our results reveal the origins of novel recognition specificities in NLR innate
108 immune receptors and the common patterns in the evolution of innate immunity.

109

110 **RESULTS**

111 **Arabidopsis NLRome shows variable rates of NLR diversification**

112 The recent elucidation of the NLR complements of over 60 accessions of the model plant *A.*
113 *thaliana* (Van de Weyer et al., 2019) provided a unique opportunity to examine rapidly evolving
114 clades of Arabidopsis NLRs. The unique advantage of the Arabidopsis dataset is the ability to
115 correlate observed diversity to known functional classes of the extensively characterized NLRs.
116 Previous NLRome analyses of this dataset were performed using OrthoMCL followed by
117 orthogroup refinement. While these analyses provided a valuable basis for global analyses of
118 selection pressures, they did not produce robust allelic series for each gene. This is likely due to
119 the divergent rates of diversification across NLRs, which complicate orthogroup assignment. To
120 circumvent this challenge, we adopted a phylogeny-based approach. To group NLRs into near
121 allelic series, we first built a unified phylogeny of all NLRs based on their shared nucleotide-
122 binding domain (Figure 1A). This tree contained 7,818 NB-ARC sequences that had >70%
123 coverage across the NB-ARC domain and represented 7,716 NLR genes, including 168 NB-ARC
124 sequences of NLRs from the reference Arabidopsis Col-0 assembly. Even though the N-terminal
125 domains were not included in the analysis, this phylogeny clearly split into clades corresponding
126 to the three canonical architectures: RPW8, Coiled-coil, and TIR domain-containing NLRs (Shao

127 et al., 2016; Tamborski and Krasileva, 2020). We arbitrarily placed the root of the tree between
128 TNL and non-TNL NLRs to simplify downstream analyses.

129
130 We split the overall phylogeny into 65 clades based on clade size (40-500 sequences) and
131 bootstrap support. Of these, 43 clades had bootstrap scores of 100, 12 additional clades had
132 bootstrap scores >70, and only 10 clades had low bootstrap values, grouping sequences that
133 could not be confidently assigned elsewhere (Figure 1B). To gain insight into the relative ages of
134 the initial clades, we used the Evolutionary Placement Algorithm to place *Arabidopsis lyrata* and
135 *Capsella rubella* NLRs in the *A. thaliana* pan-NLRome (Supplemental Figure 1). Of the 65
136 initial clades, 53 had representative sequences from either or both outgroups (Supplemental Data
137 Set 1). In the initial partition, the largest clade contained 431 sequences, allowing us to construct
138 *de novo* full-length alignments and clade phylogenies for all clades. A tree of one of the initial
139 clades, Int14015, containing the resistance gene *RPP8*, is representative of observed evolutionary
140 dynamics and is shown in Figure 1C. This tree contains five well supported subclades that differ
141 in size and internal diversity, as reflected by the very short internal branch lengths in four out of
142 five subclades. The observation that closely related sequences evolve at very different rates is
143 true not only for *RPP8*, but throughout the NLR family. *RPP1*, a well characterized NLR that
144 directly interacts with its target *ATR1*, also has closely related sequences that are largely
145 identical in different ecotypes (Figure 1D). In fact, all clades with longer branches, i.e. higher
146 amino acid divergence, have closely related clades with paralogous genes that show very little
147 variation between ecotypes. These observations are consistent with the notion that closely related
148 NLR genes are experiencing different selection pressures (Ding et al., 2007).

149
150 We iteratively refined the initial clades by splitting them into two or more subclades and
151 repeating the alignment and phylogeny generation steps. We prioritized cutting long, well
152 supported internal branches and therefore tended to preserve both rapidly evolving and low
153 variability subclades (see Methods). After two iterations, the NLRs fell into 223 non-singleton
154 and 14 singleton clades. The distribution of clade representatives across all ecotypes is
155 summarized in Supplemental Data Set 2. This NLRome partition is somewhat more conservative
156 than the OrthoMCL-based analysis, which produced 464 orthogroups and 1663 singletons (Van
157 de Weyer et al., 2019). In our final clade assignments, 83% of all clades contained no more than

158 one gene for all represented ecotypes, thus approximating allelic series. Over 90% of all NLRs
159 fell into clades of 20 or more genes, allowing sampling for sequence diversity analysis. Only six
160 large clades that ranged in size from 73 to 323 sequences contained multiple genes for ten or
161 more ecotypes and could not be split further due to the lack of long internal branches with strong
162 support (Supplemental Data Set 2). The large clades contained RPP1, RPP4/5, RPP39, and
163 RPP8, suggesting that interallelic exchange complicated the phylogeny and prevented separation
164 into allelic series. Taken together, our analyses suggest that pan-genomic NLR repertoires can be
165 clustered into near-allelic series using phylogenetic approaches.

166

167 **Sequence analysis of the NLRome clades identifies highly variable NLRs**

168 NLR genes encode immune receptors that provide protection during pathogen infection. Their
169 highly variable regions are expected to contain the specificity-determining residues. We used
170 Shannon entropy as a sensitive and robust measure of amino acid diversity. Entropy is zero at
171 positions that are invariant, and it reaches a theoretical maximum of $\log_2 20$ or ~ 4.32 when all 20
172 amino acids are present in equal ratios; a position with two variant amino acids present at equal
173 ratios produces a value of 1 bit. A Shannon entropy plot thus represents a fingerprint of sequence
174 diversity encoded in the alignment (Figure 2A).

175

176 Several functional classes of NLRs produced entropy plots with limited diversity. The ancient
177 helper RNL NRG1.1, the indirect recognition CNL RPS2, and the integrated-domain TNL
178 RRS1B produced entropy plots in which entropy never exceeded 1 bit. The low sequence
179 variability in these clades is consistent with their conserved functions. By contrast, 30 NLR
180 genes in the reference ecotype Col-0, including 14 CNL genes and 16 TNL genes, belonged to
181 clades whose alignments repeatedly scored above 1.5 bits and revealed a series of periodic spikes
182 in the LRR region. Among these genes were the known direct recognition proteins from the
183 RPP13 and RPP1 clades. Using Shannon entropy as a metric, we defined highly variable NLRs
184 (hvNLRs) as those with 10 or more positions exceeding 1.5 bit cutoff (see Supplemental Figure 2
185 for the relevant distribution). No protein known to indirectly recognize pathogen effector was
186 found among hvNLRs, and all known direct binders were detected among hvNLRs (Figure 2B).
187 When we ran Shannon entropy analyses on the previously identified NLR orthogroups (Van de
188 Weyer et al., 2019), we only detected 15 hvNLRs, 5 of which did not overlap with our

189 phylogeny-based analyses (3 slightly below 1.5 bits cutoff and 2 not supported as true
190 orthogroups by phylogeny). This suggests that phylogeny-based orthogroup assignment is a
191 better option for preserving and detecting hvNLR clades. We predict that phylogeny-based NLR
192 clade analysis combined with Shannon entropy can be applied to non-model plants to
193 computationally separate candidate direct binders from other NLRs based on their sequence
194 diversity.

195

196 **Highly variable NLRs are distributed throughout the TNL and CNL clades**

197 We observed that hvNLRs were distributed over the NLR tree of the reference accession Col-0
198 with representatives in both TNL and CNL major clades. Within both major clades, there were
199 multiple hvNLR genes right next to conserved paralogs that did not show excess diversity. This
200 is consistent with our prior observation that NLR subclades with long branches have close
201 paralogs with limited subclade diversity. Recent duplications of hvNLRs have produced local
202 hvNLR clusters such as those near *RPP7*, *RPP39*, *RPP4/5*, and *RPP1*. NLRs found in
203 phylogenetic proximity often also cluster physically on the Arabidopsis chromosomes
204 (Supplemental Figure 3). Nonetheless, genomic clustering with close paralogs is not required for
205 an NLR to become highly variable, as shown by *RPP9*, *RPP13*, and *RPP28*. Also, presence in a
206 physical cluster does not force a gene to become an hvNLR, as shown by *RLM3* in the *RPP4/5*
207 genomic cluster and *CW9* in the *RPP7* genomic cluster. Thus, it appears that the copy number
208 variation observed in the clusters is an independent process that helps create material for NLR
209 evolution, but the generation of highly variable NLRs can proceed outside of genomic clusters.

210

211 The physical proximity and phylogenetic relationships of hvNLRs and their closely related low
212 variability paralogs suggest that rapid switches in selective pressure were involved in generating
213 the apparent diversity. Since the selection of an NLR is likely to correlate with its function, we
214 located the known guardian NLRs within the phylogeny. Since these NLRs are expected to
215 maintain binding sites for conserved plant proteins, we expected them to show low entropy
216 scores. As we have already seen for *RPS2*, other known guardian NLRs including *RPM1*, *RPS5*,
217 and *ZAR1* all showed low variability. However, they did not form a separate clade within the
218 phylogeny; instead, they were interspersed by hvNLRs. This phylogenetic arrangement, together
219 with the excess of both copy number variation and amino acid diversity in the hvNLRs, argue for

220 a mechanism where hvNLRs mostly act in direct recognition mode but are infrequently able to
221 generate indirect recognition alleles that are preserved due to their competitive advantage.

222

223 **Highly variable NLRs contain the known NLR autoimmune loci**

224 Generating diverse receptors in the immune system carries with it a cost of autoimmune
225 recognition. In the known Dangerous Mix gene pairs, at least one and sometimes both causative
226 alleles are NLRs (Chae et al., 2014). If our prediction that highly variable NLRs are sources of
227 novel direct binding is correct, we would expect to find a strong overlap between hvNLRs and
228 Dangerous Mix NLRs. Indeed, hvNLR clades contain all the known NLR Dangerous Mix genes
229 including *RPP7*, *RPP8*, *RPP4/5*, and *RPP1*. We suspect that in the future, more Dangerous Mix
230 NLRs will be found that will map to other hvNLR loci. This finding also suggests that targeted
231 resequencing of NLRs in crop species could help identify loci responsible for hybrid necrosis
232 phenotypes, which are a frequent impediment to breeding.

233

234 **Highly variable residues cluster on the surfaces of LRR domains of hvNLRs**

235 The LRR domains are known to encode the recognition specificities of plant NLRs. First, we
236 wanted to know whether highly variable residues occur predominantly in the LRR domain. This
237 was indeed the case for all 30 hvNLRs examined (Table 1). We noticed, however, that regions in
238 the NB-ARC domain also had high entropy scores in multiple NLRs (*RPP1* and *RPP8* in Figure
239 2A). This suggests that a limited number of residues in the NB-ARC domain could participate in
240 target binding in these receptors. Alternatively, these could compensate for changes in the LRR
241 in order to maintain the off state in the absence of the ligand. Many TNLs have post-LRR
242 domains that lack the characteristic LRR pattern of residues yet are predicted to be folded and
243 form a contiguous structure with the preceding repeats (Van Ghelder and Esmenjaud, 2016). We
244 observed that the post-LRR domains also often contained residues with high entropy scores
245 (*RPP1* in Figure 2A). Together, these data suggest that the LRR carries the majority of binding
246 residues, while NB-ARC and post-LRR domains can also participate in ligand binding.

247

248 If the high entropy residues do indeed make up the target binding sites, we would expect to find
249 them in one or two clusters on the receptor surfaces and to include exposed hydrophobic
250 residues. LRR domains fold in a predictable manner that buries the conserved leucines and

251 exposes the variable residues on the protein surface; this allows us to skip structure prediction
252 and to approximate LRR surfaces based on repeat annotation. The concave side of LRR domains
253 contains a beta-sheet with a regular array of surface-exposed residues, and it can be represented
254 as a table with one line per repeat unit and the columns corresponding to variable positions in the
255 canonical $LX_2X_3LX_5LX_6X_7$ repeat. In the case of ZAR1, the first plant NLR whose structure was
256 elucidated, such matrix representation based on repeat annotation perfectly matches the one that
257 is based on the experimental structure (Figure 3A).

258
259 In order to test whether entropy analysis can predict NLR binding sites, we annotated LRRs for
260 each hvNLR gene in Col-0 and mapped entropy scores onto this representation. This analysis
261 revealed that in all the hvNLRs, the periodic spikes in entropy signal over the LRR likely
262 correspond to one or two surface clusters in the NLR protein (Figure 3B for three representative
263 examples, Supplemental Data Set 3 for all Col-0 hvNLRs). In the first example, AT5G43740, the
264 strongest variability signal is found in LRRs 8 through 12 and positions 3, 5, 7, and 8 of the
265 repeat. Additional high entropy signal comes from LRR1 through LRR5 positions 8 and 10. In
266 RPP13, the positions C-terminal to the predicted beta sheet appear to play an important role in
267 determining binding specificity. Unlike AT5G43740, highly variable residues in positions 8, 9,
268 and 10 of the repeats appear throughout the annotated LRR region, while all residues in position
269 2 and 3 are conserved. We therefore predict that in RPP13, loops that follow the beta strands
270 play a key role in determining substrate specificity. Our prediction that specificity determinants
271 of RPP13 stretch between LRR1 and LRR12 are in agreement with the large experimentally
272 identified specificity-determining region in RPP13 (Rentel et al., 2008a).

273
274 RPP1 is a well-studied example of a direct recognition NLR where multiple alleles have different
275 recognition profiles of the effector ATR1 of the downy mildew pathogen *Hyaloperonospora*
276 *arabidopsidis* (Rehmany et al., 2005). In RPP1, we observed a large number of contiguous
277 residues that likely contribute to binding specificity stretching from LRR1 to LRR15. Highly
278 variable residues are concentrated in positions 5, 7 and 8 at the beginning of the domain but shift
279 towards the start of the beta strands in the later repeat units, with residues 2, 3 and 5 lighting up
280 uniformly in LRR7 - LRR15. Rather unusually, we also observed some variable residues in the -
281 1 and -2 positions. We conclude that in RPP1 (and in AT5G43740), the targets likely bind

282 through the middle of the horseshoe LRR shape rather than on one side of it, as in the case of
283 RPP13. The high-entropy residues in RPP1 contain the amino acids previously shown to extend
284 recognition specificity of the RPP1 allele NdA towards ATR1-Maks9 (Krasileva, 2011) and
285 those that directly interact with ATR1 in the cryo-EM structure (Ma et al., 2020) (see below for
286 discussion).

287
288 To further investigate whether the identified highly variable surfaces indeed represent target-
289 binding sites, we surveyed these regions of high diversity for the presence of exposed
290 hydrophobic residues, which are commonly found at the centers of protein-protein binding sites
291 (Figure 3C). Indeed, in every case, the highly variable residues included exposed hydrophobic
292 amino acids, often including bulky aromatics such as tryptophan and phenylalanine. We also
293 tested whether the entropy-based predictions agree with the results of positive selection analyses
294 that have been used in the past to identify functionally important residues in NLRs (Kuang et al.,
295 2004). In RPP13, 66% of all high-entropy residues (>1.5 bits) were under positive selection
296 according to Phylogenetic Analysis by Maximum Likelihood (PAML) Model 8 (Supplemental
297 Figure 4). All of the remaining high-entropy residues fell into regions that contained gaps in the
298 alignment and could not be analyzed by PAML. Thus, the results of the entropy analyses of
299 hvNLR surfaces are consistent with the results of the widely accepted molecular evolution
300 analyses performed on the underlying nucleotide sequences.

301
302 **NLR binding sites are largely similar across the NLRome**

303 We next examined how the placement of the highly variable residues and the predicted ligand
304 binding site evolved across the NLR phylogeny (Figure 4). Overall, closely related paralogs
305 shared a similar binding site location, and most variation was apparent between CNLs and TNLs.
306 We observed that the clustering of highly variable residues was largely similar across CNLs,
307 with most sites clustering together in C-terminal repeats and most variability introduced by the
308 repeat number variation. In TNLs, highly variable sites were more dispersed across the LRRs,
309 and the predicted binding site was stretched across NLRs with a larger number of repeats. Across
310 both TNLs and CNLs, the N-termini of LRRs 1-4 were invariable: this region is in contact with
311 the invariable part of the NB-ARC domain and might be important for regulating NLR
312 activation.

313

314 **The ZAR1-RKS1 binding site overlaps with the binding site of RPP13 predicted by entropy-**
315 **based analysis**

316 Arabidopsis ZAR1 is an indirect-recognition NLR and the first one with an elucidated structure.
317 In our phylogeny, its closest hvNLR is RPP13 (Figure 2B). While the ZAR1 entropy plot lacked
318 high-entropy residues, we wanted to compare the known footprint of RKS1, the ZAR1 binding
319 partner, with the positions of highly variable residues in RPP13. Unusually for hvNLRs, highly
320 variable residues of RPP13 cluster on the C-terminal side of the repeats, with positions 7-10 of
321 the repeat units showing the highest diversity (Figure 4). Surprisingly, the similarly positioned
322 residues in ZAR1 are used to bind its stable complex partner, RKS1 (Figure 5). This finding is
323 consistent with the notion that ZAR1 and RPP13 emerged from an hvNLR common ancestor that
324 had a binding site similar to that observed in ZAR1 and predicted in RPP13.

325

326 **High-entropy residues in RPP13 are required for recognition of ATR13**

327 To experimentally test our prediction, we created synthetic RPP13 constructs and transiently
328 expressed them in *Nicotiana benthamiana* together with the ATR13 d49 Emco5 allele, which is
329 recognized by RPP13-Nd but not RPP13-Col. We used another effector, ATR1 d51 Emoy2,
330 which is not recognized by either RPP13 variant, as a negative control. RPP13-Col containing
331 the 509-729 amino acid region from the Nd allele showed a gain of ATR13 recognition, which is
332 consistent with our prediction (Figure 5C). Similarly, swapping 21 amino acids with Shannon
333 entropy >1.5 bit from Nd to Col created a loss-of-function allele, despite stable protein
334 expression, confirming the functional requirement for highly variable residues (Figure 5C,
335 Supplemental File 3). However, the same 21 amino acids transferred from RPP13-Nd to RPP13-
336 Col were not sufficient for a gain of recognition, suggesting that residues with lower entropy
337 scores also participate in target binding. (Neither functional nor non-functional RPP13-Col
338 variants could be observed by immunoblotting, as reported previously (Rentel et al. 2008).)

339

340 **The majority of RPP1 target-binding site residues show high sequence variability**

341 While this manuscript was in review, the cryo-EM structure of RPP1 bound to ATR1 was
342 published (Ma et al. 2020), allowing us to directly evaluate the accuracy of our binding site
343 predictions. The majority of binding residues had entropy values above one bit (Figure 6A). Both

344 precision (fraction of positives among all predictions) and recall (fraction of positives recovered)
345 varied with the entropy cutoff chosen. Maximal recall was achieved at a cutoff of 0.8 bit, and
346 precision improved up to a cutoff of 1.8 bits. Thus, cutoff values in this range are likely to be
347 useful, with higher cutoffs achieving greater accuracy at the cost of missing a greater number of
348 true positives (Figure 6B). Our empirical 1.5 bit cutoff used to define hvNLR clades is therefore
349 a conservative one. It is also important to note that sequence-based analyses predicted a number
350 of RPP1 binding residues past the LRR domain (Table 1); the structure revealed that these
351 residues form a contiguous surface on the post-LRR domain that is characteristic of a number of
352 TNL receptors.

353

354 **hvNLRs show a similar phylogenetic distribution in *Brachypodium distachyon***

355 To test whether our methods and findings are applicable beyond *A. thaliana*, we performed a
356 similar analysis on 54 lines of *Brachypodium distachyon*, a model grass species. The automatic
357 short-read assembly and annotation pipeline used to generate the *Brachypodium* data is less
358 reliable than the targeted resequencing approach used to generate *Arabidopsis* pan-NLRome.
359 Specifically, only 45% of hvNLRs present in reference strain Bd-21 were recovered in the
360 assembly control. Nonetheless, the overall picture that emerged from the analysis of
361 *Brachypodium* NLR clades is similar to that of *Arabidopsis*. After splitting the overall
362 *Brachypodium* NLR tree into 91 initial clades, we performed four rounds of clade refinement to
363 arrive at a final clade partition with 433 subclades. Of these, 28 produced alignments that
364 fulfilled the hvNLR criteria. Altogether, 40 hvNLRs in the reference accession Bd21 were
365 identified as hvNLRs.

366

367 Similar to *A. thaliana*, *Brachypodium* hvNLRs were distributed throughout the phylogeny,
368 including in the highly expanded monocot-specific CNL clade. Here too, hvNLRs had sister
369 clades that showed little amino-acid diversity. Importantly, when we constructed the joint tree
370 for Col-0 and Bd21 reference NLRomes, the only hvNLRs from the two species that appeared
371 close together belonged to the RPP13-like clades (Figure 7). This highlights the importance of
372 sequencing the pan-NLRomes of plants of interest, as the identification of hvNLRs is unlikely to
373 be transferable except for closely related species.

374

375 **DISCUSSION**

376 Even before the first NLR structure or the extensive sequence datasets were available,
377 Michelmore and Meyers predicted that hypervariable amino acid positions in the NLRs would
378 map to the concave surface of the LRR domain based on the signatures of positive selection in a
379 number of selected examples (Michelmore and Meyers, 1998). They generalized that this might
380 be true for all NLRs. This model was challenged by the discovery of indirect recognition and of
381 strongly conserved NLRs. Our analysis proposes a powerful methodology to study NLR-omes,
382 predicts NLR mode of action through sequence analysis, and reconciles the evolution of direct
383 recognition NLRs (under diversifying selection) and indirect recognition NLRs (under purifying
384 or balancing selection).

385
386 In this study, we observed that hvNLRs account for the known direct recognition NLRs and for
387 autoimmune NLRs. We also observed that the hvNLRs have close paralogs with little allelic
388 diversity that include the known indirect recognition NLRs. Based on this observation, we
389 propose that indirect recognition NLRs are a functional byproduct of hvNLR evolution,
390 providing an important update of the birth-and-death model (Michelmore and Meyers, 1998).
391 Our analyses suggest that in a given species, diversity generation occurs in a limited subset of
392 NLR genes, creating a wide recognition potential, including binding to endogenous plant
393 proteins. When recognition of endogenous proteins is beneficial, such as under perturbations by
394 the pathogen, the NLR evolves into indirect recognition and begins to experience different
395 selective forces.

396
397 The resolution and sensitivity of our analyses became possible when we adopted two key
398 approaches: identifying orthologous groups of NLR receptors by phylogeny in place of
399 commonly used distance metrics; and using simpler Shannon entropy measure of diversity in
400 place of more complex evolutionary models. Separating rapidly evolving protein families into
401 meaningful clades or groups for downstream analyses is a common challenge. In the NLR family
402 of plant immune receptors, this process is further complicated by ongoing information flow
403 between close paralogs through recombination and gene conversion (Kuang et al., 2004).
404 Phylogeny-based analyses are considered to be more accurate than distance-based methods for
405 similar problems such as classifying Human Immunodeficiency Virus isolates (Pineda-Peña et

406 al., 2013). Our phylogeny-based partition of NLR immune receptors into clades improved on the
407 published OrthoMCL-based partition by producing more encompassing clades and (in particular)
408 fewer singletons. OrthoMCL is a distance-based algorithm that was originally developed to
409 separate members of different protein families rapidly; it uses a single parameter to determine the
410 rate of convergence (Li et al., 2003). This makes its use to partition the pan-NLRome
411 problematic, because closely related NLRs are known to experience vastly different selection
412 pressures and thus are expected to contain very different amounts of allelic diversity (Bakker et
413 al., 2006; Kuang et al., 2004). The specific danger for hvNLR identification is that highly
414 variable clades will be split, losing the relevant signal. This is indeed what we observed, as the
415 OrthoMCL-based analysis identified only one out of three hvNLRs and missed key sources of
416 new NLR specificity such as the RPP1 cluster, which was split into small orthogroups. The
417 drawback of the phylogeny-based approach is that it is not yet fully automated; however, we are
418 hopeful that phylogeny-aware algorithms will emerge that will fill this gap. One alternate
419 approach that would simplify the analysis would be to replace the initial clade assignment with
420 iterative matching of NLR sequences against a set of inferred ancestral NLR models (Shao et al.,
421 2016).

422

423 It is well established that closely related NLRs experience different modes of selection (Ding et
424 al., 2007; Wang et al., 2011; Kuang et al., 2004). By expanding this observation to the pan-
425 NLRome and combining it with the wealth of characterized NLRs in Arabidopsis, we were able
426 to decipher a larger evolutionary pattern where hvNLRs act as sources of new specificities and
427 encompass the known direct-recognition NLRs. Their diversification, while advantageous to the
428 plant, comes at a cost. All known Dangerous Mix NLR genes that can trigger autoimmune
429 recognition belong to hvNLR clades. Thus, the generation of novel specificities goes hand in
430 hand with the potential for self-recognition and auto-immunity. We also propose that during their
431 continuous evolution, hvNLRs can generate indirect-recognition NLRs at a low frequency.
432 Because indirect recognition usually tracks a conserved effector activity, it is more robust than
433 direct recognition of the effector surface. Duplication of such successful variants might then be
434 favored due to the increased fitness of the progeny where one copy could eventually be
435 preserved while the other could continue to generate novel specificities (Kondrashov et al.,
436 2002). The latter inference is consistent with our observation that ZAR1, an indirect-recognition

437 NLR, binds to its stable complex partner RKS1 through the same surface on the LRR that
438 contains highly variable residues in RPP13, its closest hvNLR.

439

440 When we applied Shannon entropy analysis to the NLR clades, only a subset of clades gave
441 strong signals; these clades included known direct recognition NLRs and autoimmune NLRs.
442 When we looked at the distribution of high-entropy amino acids in the 30 hvNLRs of
443 Arabidopsis reference strain Col-0, we found that these residues commonly clustered on the
444 predicted surfaces of LRR domains. This observation is consistent with the finding that binding
445 specificities are largely encoded in the LRR domains, as supported by multiple genetic and
446 biochemical studies (Ellis et al., 2007; Krasileva et al., 2010), as well as the prediction (by
447 evolutionary studies) that amino acid residues under positive selection are located within LRRs
448 (Kuang et al., 2004; Rose et al., 2004; Wang et al., 2011). When we carried out a positive
449 selection analysis on the RPP13 clade, we found that the majority of residues with entropy >1.5
450 bits were under positive selection. The only exceptions were residues that could not be analyzed
451 for positive selection due to the presence of gaps in the relevant alignment columns. Shannon
452 entropy calculation does not count gap characters. Instead, it works without making complex
453 assumptions about the data and is therefore much faster computationally.

454

455 In our analysis, we went a step further to predict binding sites in hvNLRs directly from pan-
456 NLRome sequence data. The identified binding sites are large. This is likely in due (in part) to
457 the concave shape of the LRR scaffold, which can place many of the beta strands in contact with
458 a relatively small target. Comparisons of antibody sequence-based predictions with experimental
459 structures showed that the predictions correctly recover ~80% of residues that do contact the
460 antigen, while also producing many false-positives (<50% precision) (Kunik et al., 2012). Unlike
461 the antibody, where the binding determinants are present on loops away from the core of the
462 structure, in the LRR, many predicted binding residues fall within the beta sheet located on the
463 concave side of the domain. This suggests that the accuracy of the prediction might be higher in
464 this system due to stronger structural constraints. Additional highly variable residues were
465 located in post-LRR domains and in specific sites within NB-ARC, suggesting their involvement
466 in substrate binding, or in case of NB-ARC of a compensatory mechanism to maintain self-
467 inhibition in the absence of the ligand. Further mutational and structural experiments in well-

468 established NLR-effector systems would be needed to test the accuracy of these predictions and
469 to help refine them.

470
471 Identification of the immense allelic diversity across hvNLRs argues that plant immunity is not
472 far in its allele-generating potential from the most well-known adaptive immune systems. Indeed,
473 LRRs are deployed in the adaptive immune systems of early-diverging vertebrates,
474 demonstrating that their modularity is sufficient for the generation of binding to any foreign
475 molecule (Das et al., 2013; Han et al., 2008). In the case of plants, enormous diversity is
476 generated at the population level rather than within a single organism, and therefore, defending
477 against new pathogens is a community effort. The identification of specific genes within crop
478 species capable of such diversity generation and their deployment in protein engineering efforts
479 could provide valuable material for plant health.

480
481 We conclude that phylogenetic analysis of pan-NLRomes combined with sequence diversity
482 analysis can rapidly classify NLRs into functional groups given sequencing information for at
483 least 40-60 diverse samples. We also believe that our method would be generally applicable to
484 the identification of highly variable receptor-like proteins, such as Cf-9 in tomato (Wulff et al.
485 2009), and the prediction of binding sites of highly variable extracellular immune receptors. Our
486 method can also predict incompatibility loci, which can be taken into account in breeding new
487 crop varieties. Similar allelic diversity analyses in other non-vertebrate eukaryotes with
488 expanded immune receptor families are needed to test whether the patterns of innate immune
489 receptor evolution we observed are shared across the eukaryotic kingdoms of life.

490

491 **METHODS**

492

493 **Phylogenetic analysis**

494 Phylogenetic tree construction for the *A. thaliana* and *B. distachyon* NLRomes and the
495 NLRomes of reference accessions was performed as previously described (Bailey et al., 2018).
496 Briefly, amino acid sequences were searched for the presence of the NB-ARC domain using
497 hmmsearch (Mistry et al., 2013) and the extended NB-ARC Hidden Markov Model (HMM)
498 13059_2018_1392_MOESM16_ESM.hmm (Bailey et al., 2018), and initial alignment was made

499 on this HMM using the -A option. The resulting alignment was processed with Easel tools
500 (<https://github.com/EddyRivasLab/easel>) to remove insertions and retain aligned sequences that
501 matched at least 70% of the HMM model. This alignment was used to construct maximum
502 likelihood phylogenetic trees using RAxML software version 8.2.12 (Kozlov et al., 2019) (raxml
503 -T 8 -n Raxml.out -f a -x 12345 -p 12345 -# 100 -m PROTCATJTT). The sequences of outgroup
504 species were aligned to the same NB-ARC HMM and placed in the pan-NLRome tree using
505 RAxML Evolutionary Placement Algorithm (Kozlov et al., 2019). The trees were visualized in
506 iTOL (Letunic and Bork, 2019).

507

508 **Initial clade assignments**

509 The phylogeny was used to separate protein sequences into clades using R scripts
510 prefix_Initial_Assignment.R (hereafter the prefix is either *Atha_NLRome* or *Brachy_NLRome*
511 for the two species under analysis). This and other scripts referenced below are available at
512 (<https://github.com/krasileva-group/hvNLR>). First, for each NB-ARC sequence, a clade 40 to
513 500 in size with the strongest bootstrap support was chosen. For sequences that did not belong to
514 clades in this size range, smaller clades were allowed. Second, the resulting set of clades was
515 made non-redundant by excluding all nesting clades. The resulting partitions uniquely assigned
516 the 7,818 *A. thaliana* NLR sequences to 65 clades and 11,488 *B. distachyon* NLR sequences to
517 91 clades.

518

519 **Iterative clade refinement**

520 For each identified clade, full-length protein sequences were aligned using the PRANK
521 algorithm (Löytynoja, 2014), and phylogenetic trees based on full-length alignments were
522 constructed as described above using RAxML (Kozlov et al., 2019). Trees were visualized in
523 iTOL, along with subclade statistics calculated in R, and R scripts were used to produce subclade
524 lists based on the trimmed branches (prefix_Refinement.R). For the first iteration, gappy
525 columns in the full-length alignments were masked (90% cutoff), and later iterations were
526 analyzed without masking gappy columns. Clade refinement was performed as follows: all tree
527 branches longer than 0.3 were cut to form two or more subclades. All branches 0.1 and shorter
528 were retained in the first iteration, and for the branches between 0.1 and 0.3, the decision to cut
529 was made by visually inspecting the tree in iTOL and considering bootstrap support and overlap

530 in ecotypes on either side of a branch. The sequences belonging to the refined subclades were
531 realigned using PRANK, and tree construction repeated. In the following iterations, some
532 branches shorter than 0.1 were cut via tree inspection in iTOL based on bootstrap support and
533 ecotype overlap. The refinement process converged to produce the final assignment of all genes
534 into 237 final clades for *A. thaliana* and 433 clades for *B. distachyon*.

535

536 **Identification of hvNLR clades and prediction of binding sites in hvNLRs**

537 We used R scripts (prefix_CladeAnalysis.R) to calculate alignment Shannon entropy scores
538 using the package “entropy”. Alignments that contained 10 or more positions with at least 1.5
539 bits were considered highly variable. All highly variable clades were examined for the presence
540 of Arabidopsis Col-0 alleles. For these Col-0 alleles, we predicted the LRR coordinates manually
541 and cross-checked these predictions with an LRRpredictor online server (Martin et al., 2020). R
542 script was used to map entropy scores to the predicted concave surface of the LRR domain
543 (Atha_NLRome_GeneEntropy.R). The entropy scores for the individual strands of LRRs
544 (LxxLxLxx) were exported in tabular format. The hydrophobicity scores of these residues were
545 calculated as the percent of hydrophobic residues at a given amino acid position and exported as
546 a second table. The resulting 2D representations of entropy and hydrophobicity of the concave
547 sides were visually examined for clustering of residues that showed both high entropy scores and
548 the presence of hydrophobic residues. Positive selection analysis of the RPP13 clade alignment
549 was carried out in PAML (Yang, 2007).

550

551 **Structural analysis of RPP13 homology model, ZAR1 structure, and RPP1 structure**

552 In order to compare the 3D spatial distribution of highly variable residues in RPP13 with the
553 ZAR1-RKS1 binding site, we used phyre2 in one-to-one threading mode to produce a model for
554 RPP13 (Kelley et al., 2015) based on the ZAR1 experimental structure. The alignment had 24%
555 identity over the complete sequences, with 31% identity before and 15% over the LRR domain.
556 Important for the model accuracy, there were only two gaps of 7 residues and two gaps of 3
557 residues, with several more single-residue gaps in the LRR domain. Thus, it is unlikely that
558 whole repeat units are missing from the model. R script (Atha_NLRome_GeneEntropy.R) was
559 used to produce a Chimera-formatted attribute files to color the model surfaces by entropy
560 scores, and figures were generated in Chimera (Pettersen et al., 2004). The dependence of

561 binding residue prediction recall and precision on the entropy cutoff was determined using a
562 custom R script (RPP1_Precision_Recall.R).

563

564 **Constructs**

565 RPP13-Nd and RPP13-Col cDNA without a stop codon fused to C-terminal HA tag in
566 pENTRY/TOPO-D were obtained from the Staskawicz laboratory (Rentel et al., 2008b) and
567 were used to generate chimeric and synthetic RPP13 variants. RPP13 501-729 synthetic
568 constructs with highly variable residues (Shannon entropy cutoff >1.5) swapped between Nd and
569 Col (Supplemental Data Set 4) were designed in SnapGene and synthesized as gene fragments by
570 IDT. The clones were digested with uniquely cutting restriction enzymes SacI (New England
571 BioLabs) and MspI (New England BioLabs). The chimeric constructs were ligated for 2 h at
572 room temperature with T4-DNA ligase (New England BioLabs) and transformed into
573 electrocompetent *E. coli* Top 10b (Invitrogen). The resulting constructs were introduced into
574 binary vector pMD:npRPP13 (Rentel et al., 2008b) using LR clonase II (Invitrogen) and
575 transformed into *Agrobacterium tumefaciens* GV3101(pMP90RK). ATR1 d51 Emoy2 tagged
576 with C-terminal citrine in pEarleyGate103 (Krasileva et al., 2010) and ATR13 d41 Emco5 in
577 p1776 (Rentel et al., 2008b) were used for transient transformation.

578

579 **Transient expression**

580 *A. tumefaciens* strains were grown for 24-48 hours at 28°C in Luria–Bertani broth (100 µg/mL,
581 gentamicin 50 µg/mL, kanamycin 25 µg/mL) with constant shaking. After pelleting, the cells
582 were resuspended in induction medium (10 mM MgCl₂, 10 mM MES, and 150 µM
583 acetosyringone, adjusted to pH 5.6 with KOH), adjusted to a final OD₆₀₀ of 0.6, and induced for
584 3 h at room temperature. Co-infiltrations were done at a final OD₆₀₀ of 0.6 and contained
585 constructs mixed in a 1:1 ratio. Fully expanded leaves of 4-5-week-old *Nicotiana benthamiana*
586 plants grown in Supersoil mix #4 supplemented with Miracle Gro Plant Food fertilizer at 24°C
587 under a 16 h light (fluorescent lamps)/8 h dark cycle were infiltrated using a blunt end syringe.
588 After infiltration, the plants were kept at constant light (fluorescent lamps, GE Cat #F405941-
589 ECO) and room temperature. The hypersensitive response reaction was monitored for 4 days,
590 with pictures taken 3 days post infiltration. Two leaf disks (1.5 cm² in diameter) were collected

591 from RPP13/ATR1 co-infiltrations for protein extraction 2 days post infiltration, frozen in liquid
592 nitrogen, and stored at -80°C.

593

594 **Protein extraction and immunoblotting**

595 Tissue in a 1.5 mL Eppendorf tube was frozen in liquid nitrogen a ground with a manual drill
596 using a pre-chilled plastic pestle. Total protein was extracted by re-suspending the ground tissue
597 in 2x Laemmli buffer (Bio-Rad, Cat. #1610737) supplemented with fresh β -mercaptoethanol to a
598 final concentration of 5% (by volume), boiling for 5 minutes, and pelleting the debris for 10
599 minutes at 14,000 rpm. 15 μ L of each protein sample was separated on a 4–15% Mini-
600 PROTEAN gel (BioRad) for 1 hour at 100V and transferred onto a nitrocellulose membrane
601 using wet transfer for 1.5 h at 300 mA. The membranes were blocked overnight in 5% milk in
602 TBST-T, incubated for 1 h in rat α -HA-horseradish peroxidase antibody (clone 3F10; Roche, Cat
603 #12013819001) at 1:1000 dilution in TBST-T, washed once for 15 minutes and twice for 5
604 minutes in TBST-T, and imaged using SuperSignal West Pico PLUS Luminol substrate (Thermo
605 Scientific) inside a Gel Imager (BioRad). Total protein loading was confirmed by staining the
606 membrane in Ponceau S and destaining in 5% acetic acid.

607

608 **Accession Numbers**

609 Arabidopsis pan-NLRome nucleotide assemblies were downloaded from the 2Blades foundation
610 (<http://2blades.org/resources/>). Gene annotations were downloaded from GitHub pan-NLRome
611 repository (<https://github.com/weigelworld/pan-nlrome/>). The gene models that matched
612 assemblies were available for 62 *A. thaliana* accessions (Van de Weyer et al., 2019), and these
613 were processed to extract the amino acid sequences of captured protein-coding genes using
614 bedtools getfasta program (Quinlan, 2014). The reference set of 168 NLR alleles (including
615 splice variants) of the Arabidopsis Col-0 genome was extracted as described before (Sarris et al.,
616 2016). The accession numbers of RPP13 used in the laboratory experiments are: RPP13-Nd
617 (AF209732.1) and RPP13-Col (AF209730.1). The PDB accession number of the RPP1 structure
618 used in this study is 7crb. Brachypodium proteomes for 54 lines were downloaded from
619 BrachyPan (<https://brachypan.jgi.doe.gov>) (Gordon et al., 2017). The R scripts used to analyze
620 project data are available via GitHub (<https://github.com/krasileva-group/hvNLR>), the complete
621 data set for the project including clade alignments and clade trees is available via Zenodo (DOI:

622 10.5281/zenodo.3951781), and the clade trees can be viewed in iTOL
623 (<http://itol.embl.de/shared/daniilprigozhin>).

624

625 **Supplemental Data**

626 **Supplemental Figure 1.** *A. thaliana* pan-NLRome tree showing initial clades and phylogenetic
627 placements of outgroup sequences from *A. lyrata* and *C. rubella*.

628 **Supplemental Figure 2.** Distribution of highly variable sites per final clade alignment.

629 **Supplemental Figure 3.** Comparison of phylogenetic versus physical clustering of Col-0 NLRs.

630 **Supplemental Figure 4.** Comparison of entropy-based and positive selection-based binding site
631 predictions.

632 **Supplemental Data Set 1.** Number of NLRs from *A. lyrata* and *C. rubella* in the initial NLR
633 clades.

634 **Supplemental Data Set 2.** Number of NLRs in the final NLR clades across the 62 *A. thaliana*
635 ecotypes.

636 **Supplemental Data Set 3.** 2D representations of LRR surfaces of 30 hvNLRs from ecotype Col-
637 0.

638 **Supplemental Data Set 4.** Nucleotide and amino acid fasta sequences of RPP13 501-729
639 synthetic constructs that have highly variable residues swapped between Col and Nd allele.

640

641 **ACKNOWLEDGEMENTS**

642 We thank Detlef Weigel and coauthors for making Arabidopsis NLRome data publicly available
643 in advance of publication. We thank members of the Krasileva lab and of the Berkeley Lab
644 Advanced Light Source structural biology community for helpful discussions. We thank Brian
645 Staskawicz, Raoul Martin, and Kyungyong Seong for their advice and for critical reading of the
646 manuscript. We are grateful to Douglas Dahlbeck for providing ATR13 and RPP13 clones
647 amidst the pandemic. We thank Marc Allaire and members of the Berkeley Center for Structural
648 Biology for support, encouragement, and the use of computational resources. KVK research on
649 plant NLRs is supported by the Gordon and Betty Moore Foundation (8802) and Two Blades
650 Foundation together Foundation for Food and Agriculture Research (CA19-SS-0000000046).
651 The Berkeley Center for Structural Biology is supported by the Howard Hughes Medical
652 Institute, the National Institutes of Health, and through participating research team partnerships.

653

654 **AUTHOR CONTRIBUTIONS**

655 DMP and KVK designed and performed the research and wrote the paper.

656

657 **REFERENCES**

- 658 **Atanasov, K.E., Liu, C., Erban, A., Kopka, J., Parker, J.E., and Alcázar, R.** (2018).
659 Mutations Suppressing Immune Hybrid Incompatibility and Their Effects on Disease
660 Resistance. *Plant Physiol.* **177**: 1152–1169.
- 661 **Baggs, E., Dagdas, G., and Krasileva, K.V.** (2017). NLR diversity, helpers and integrated
662 domains: making sense of the NLR IDentity. *Curr. Opin. Plant Biol.* **38**: 59–67.
- 663 **Bailey, P.C., Schudoma, C., Jackson, W., Baggs, E., Dagdas, G., Haerty, W., Moscou, M.,
664 and Krasileva, K.V.** (2018). Dominant integration locus drives continuous diversification
665 of plant immune receptors with exogenous domain fusions. *Genome Biol.* **19**: 23.
- 666 **Bakker, E.G., Toomajian, C., Kreitman, M., and Bergelson, J.** (2006). A genome-wide
667 survey of R gene polymorphisms in Arabidopsis. *Plant Cell* **18**: 1803–1818.
- 668 **Bombliès, K.** (2009). Too much of a good thing? Hybrid necrosis as a by-product of plant
669 immune system diversification. *Botany* **87**: 1013–1022.
- 670 **Bombliès, K., Lempe, J., Epple, P., Warthmann, N., Lanz, C., Dangl, J.L., and Weigel, D.**
671 (2007). Autoimmune response as a mechanism for a Dobzhansky-Muller-type
672 incompatibility syndrome in plants. *PLoS Biol.* **5**: e236.
- 673 **Catanzariti, A.-M., Dodds, P.N., Ve, T., Kobe, B., Ellis, J.G., and Staskawicz, B.J.** (2010).
674 The AvrM Effector from Flax Rust Has a Structured C-Terminal Domain and Interacts
675 Directly with the M Resistance Protein. *Molecular Plant-Microbe Interactions* **23**: 49–57.
- 676 **Cesari, S.** (2018). Multiple strategies for pathogen perception by plant immune receptors. *New
677 Phytol.* **219**: 17–24.
- 678 **Chae, E. et al.** (2014). Species-wide genetic incompatibility analysis identifies immune genes as
679 hot spots of deleterious epistasis. *Cell* **159**: 1341–1351.
- 680 **Chen, J. et al.** (2017). Loss of by somatic exchange in stem rust leads to virulence for resistance
681 in wheat. *Science* **358**: 1607–1610.
- 682 **Dangl, J.L., Horvath, D.M., and Staskawicz, B.J.** (2013). Pivoting the plant immune system
683 from dissection to deployment. *Science* **341**: 746–751.
- 684 **Das, S., Hirano, M., Aghaallaei, N., Bajoghli, B., Boehm, T., and Cooper, M.D.** (2013).
685 Organization of lamprey variable lymphocyte receptor C locus and repertoire development.
686 *Proc. Natl. Acad. Sci. U. S. A.* **110**: 6043–6048.

- 687 **Ding, J., Cheng, H., Jin, X., Araki, H., Yang, Y., and Tian, D.** (2007). Contrasting patterns of
688 evolution between allelic groups at a single locus in Arabidopsis. *Genetica* **129**: 235–242.
- 689 **Ellis, J.G., Dodds, P.N., and Lawrence, G.J.** (2007). Flax Rust Resistance Gene Specificity is
690 Based on Direct Resistance-Avirulence Protein Interactions. *Annual Review of*
691 *Phytopathology* **45**: 289–306.
- 692 **Gordon, S.P. et al.** (2017). Extensive gene content variation in the *Brachypodium distachyon*
693 pan-genome correlates with population structure. *Nat. Commun.* **8**: 2184.
- 694 **Goritschnig, S., Steinbrenner, A.D., Grunwald, D.J., and Staskawicz, B.J.** (2016).
695 Structurally distinct *Arabidopsis thaliana* NLR immune receptors recognize tandem WY
696 domains of an oomycete effector. *New Phytol.* **210**: 984–996.
- 697 **Han, B.W., Herrin, B.R., Cooper, M.D., and Wilson, I.A.** (2008). Antigen recognition by
698 variable lymphocyte receptors. *Science* **321**: 1834–1837.
- 699 **Jones, J.D.G., Vance, R.E., and Dangl, J.L.** (2016). Intracellular innate immune surveillance
700 devices in plants and animals. *Science* **354**.
- 701 **Jubic, L.M., Saile, S., Furzer, O.J., El Kasmi, F., and Dangl, J.L.** (2019). Help wanted: helper
702 NLRs and plant immune responses. *Curr. Opin. Plant Biol.* **50**: 82–94.
- 703 **Kabat, E.A.** (1970). Heterogeneity and structure of antibody-combining sites. *Ann. N. Y. Acad.*
704 *Sci.* **169**: 43–54.
- 705 **Kondrashov, F.A., Rogozin, I.B., Wolf, Y.I., and Koonin, E.V.** (2002). Selection in the
706 evolution of gene duplications. *Genome Biol.* **3**: RESEARCH0008.
- 707 **Krasileva, K.V., Dahlbeck, D., and Staskawicz, B.J.** (2010). Activation of an *Arabidopsis*
708 resistance protein is specified by the in planta association of its leucine-rich repeat domain
709 with the cognate oomycete effector. *Plant Cell* **22**: 2444–2458.
- 710 **Krasileva, K.V.** (2011). The Molecular Basis for Recognition of Oomycete Effectors in
711 *Arabidopsis*. Doctoral Dissertation. UC Berkeley. <https://escholarship.org/uc/item/57x979rf>
- 712 **Kruger, J.** (2002). A Tomato Cysteine Protease Required for Cf-2-Dependent Disease
713 Resistance and Suppression of Autonecrosis. *Science* **296**: 744–747.
- 714 **Kuang, H., Woo, S.-S., Meyers, B.C., Nevo, E., and Michelmore, R.W.** (2004). Multiple
715 genetic processes result in heterogeneous rates of evolution within the major cluster disease
716 resistance genes in lettuce. *Plant Cell* **16**: 2870–2894.
- 717 **Kunik, V., Peters, B., and Ofran, Y.** (2012). Structural consensus among antibodies defines the
718 antigen binding site. *PLoS Comput. Biol.* **8**: e1002388.
- 719 **Liao, H., Yeh, W., Chiang, D., Jernigan, R.L., and Lustig, B.** (2005). Protein sequence
720 entropy is closely related to packing density and hydrophobicity. *Protein Eng. Des. Sel.* **18**:

721 59–64.

722 **Li, L., Stoeckert, C.J., Jr, and Roos, D.S.** (2003). OrthoMCL: identification of ortholog groups
723 for eukaryotic genomes. *Genome Res.* **13**: 2178–2189.

724 **Magliery, T.J. and Regan, L.** (2005). 10.1186/1471-2105-6-240. *BMC Bioinformatics* **6**: 240.

725 **Martin, R., Qi, T., Zhang, H., Liu, F., King, M., Toth, C., Nogales, E., and Staskawicz, B.J.**
726 (2020). Structure of the activated ROQ1 resistosome directly recognizing the pathogen
727 effector XopQ. *Science* **370**.

728 **Ma, S. et al.** (2020). Direct pathogen-induced assembly of an NLR immune receptor complex to
729 form a holoenzyme. *Science* **370**.

730 **Michelmore, R.W. and Meyers, B.C.** (1998). Clusters of resistance genes in plants evolve by
731 divergent selection and a birth-and-death process. *Genome Res.* **8**: 1113–1130.

732 **Pineda-Peña, A.-C., Faria, N.R., Imbrechts, S., Libin, P., Abecasis, A.B., Deforche, K.,**
733 **Gómez-López, A., Camacho, R.J., de Oliveira, T., and Vandamme, A.-M.** (2013).
734 Automated subtyping of HIV-1 genetic sequences for clinical and surveillance purposes:
735 performance evaluation of the new REGA version 3 and seven other tools. *Infect. Genet.*
736 *Evol.* **19**: 337–348.

737 **Rehmany, A.P., Gordon, A., Rose, L.E., Allen, R.L., Armstrong, M.R., Whisson, S.C.,**
738 **Kamoun, S., Tyler, B.M., Birch, P.R.J., and Beynon, J.L.** (2005). Differential
739 recognition of highly divergent downy mildew avirulence gene alleles by RPP1 resistance
740 genes from two Arabidopsis lines. *Plant Cell* **17**: 1839–1850.

741 **Rentel, M.C., Leonelli, L., Dahlbeck, D., Zhao, B., and Staskawicz, B.J.** (2008a). Recognition
742 of the *Hyaloperonospora parasitica* effector ATR13 triggers resistance against oomycete,
743 bacterial, and viral pathogens. *Proc. Natl. Acad. Sci. U. S. A.* **105**: 1091–1096.

744 **Rentel, M.C., Leonelli, L., Dahlbeck, D., Zhao, B., and Staskawicz, B.J.** (2008b).
745 Recognition of the *Hyaloperonospora parasitica* effector ATR13 triggers resistance against
746 oomycete, bacterial, and viral pathogens. *Proc. Natl. Acad. Sci. U. S. A.* **105**: 1091–1096.

747 **Rose, L.E., Bittner-Eddy, P.D., Langley, C.H., Holub, E.B., Michelmore, R.W., and**
748 **Beynon, J.L.** (2004). The maintenance of extreme amino acid diversity at the disease
749 resistance gene, RPP13, in *Arabidopsis thaliana*. *Genetics* **166**: 1517–1527.

750 **Sanders, M.P.A., Fleuren, W.W.M., Verhoeven, S., van den Beld, S., Alkema, W., de Vlieg,**
751 **J., and Klomp, J.P.G.** (2011). ss-TEA: Entropy based identification of receptor specific
752 ligand binding residues from a multiple sequence alignment of class A GPCRs. *BMC*
753 *Bioinformatics* **12**: 332.

754 **Santangelo, E., Fonzo, V., Astolfi, S., Zuchi, S., Caccia, R., Mosconi, P., Mazzucato, A., and**
755 **Soressi, G.P.** (2003). The Cf-2 / Rcr3esc gene interaction in tomato (*Lycopersicon*
756 *esculentum*) induces autonecrosis and triggers biochemical markers of oxidative burst at

757 cellular level. *Funct. Plant Biol.* **30**: 1117.

758 **Saur, I.M., Bauer, S., Kracher, B., Lu, X., Franzeskakis, L., Müller, M.C., Sabelleck, B.,**
759 **Kümmel, F., Panstruga, R., Maekawa, T., and Schulze-Lefert, P.** (2019). Multiple pairs
760 of allelic MLA immune receptor-powdery mildew AVR effectors argue for a direct
761 recognition mechanism. *Elife* **8**.

762 **Seong, K., Seo, E., Witek, K., Li, M., and Staskawicz, B.** (2020). Evolution of NLR resistance
763 genes with noncanonical N-terminal domains in wild tomato species. *New Phytol.*

764 **Shao, Z.-Q., Xue, J.-Y., Wu, P., Zhang, Y.-M., Wu, Y., Hang, Y.-Y., Wang, B., and Chen,**
765 **J.-Q.** (2016). Large-Scale Analyses of Angiosperm Nucleotide-Binding Site-Leucine-Rich
766 Repeat Genes Reveal Three Anciently Diverged Classes with Distinct Evolutionary
767 Patterns. *Plant Physiol.* **170**: 2095–2109.

768 **Shenkin, P.S., Erman, B., and Mastrandrea, L.D.** (1991). Information-theoretical entropy as a
769 measure of sequence variability. *Proteins* **11**: 297–313.

770 **Stam, R., Nosenko, T., Hörger, A.C., Stephan, W., Seidel, M., Kuhn, J.M.M., Haberer, G.,**
771 **and Tellier, A.** (2019a). The Reference Genome and Transcriptome Assemblies of the Wild
772 Tomato Species Highlights Birth and Death of NLR Genes Between Tomato Species. *G3* **9**:
773 3933–3941.

774 **Stam, R., Silva-Arias, G.A., and Tellier, A.** (2019b). Subsets of NLR genes show differential
775 signatures of adaptation during colonization of new habitats. *New Phytol.* **224**: 367–379.

776 **Stewart, J.J., Lee, C.Y., Ibrahim, S., Watts, P., Shlomchik, M., Weigert, M., and Litwin, S.**
777 (1997). A Shannon entropy analysis of immunoglobulin and T cell receptor. *Molecular*
778 *Immunology* **34**: 1067–1082.

779 **Tamborski, J. and Krasileva, K.V.** (2020). Evolution of Plant NLRs: From Natural History to
780 Precise Modifications. *Annu. Rev. Plant Biol.* **71**: 355–378.

781 **Van de Weyer, A.-L., Monteiro, F., Furzer, O.J., Nishimura, M.T., Cevik, V., Witek, K.,**
782 **Jones, J.D.G., Dangl, J.L., Weigel, D., and Bemm, F.** (2019). A Species-Wide Inventory
783 of NLR Genes and Alleles in *Arabidopsis thaliana*. *Cell* **178**: 1260–1272.e14.

784 **Van Ghelder, C. and Esmenjaud, D.** (2016). TNL genes in peach: insights into the post-LRR
785 domain. *BMC Genomics* **17**: 317.

786 **Wang, J., Hu, M., Wang, J., Qi, J., Han, Z., Wang, G., Qi, Y., Wang, H.-W., Zhou, J.-M.,**
787 **and Chai, J.** (2019a). Reconstitution and structure of a plant NLR resistosome conferring
788 immunity. *Science* **364**.

789 **Wang, J., Wang, J., Hu, M., Wu, S., Qi, J., Wang, G., Han, Z., Qi, Y., Gao, N., Wang, H.-**
790 **W., Zhou, J.-M., and Chai, J.** (2019b). Ligand-triggered allosteric ADP release primes a
791 plant NLR complex. *Science* **364**.

- 792 **Wang, J., Zhang, L., Li, J., Lawton-Rauh, A., and Tian, D.** (2011). Unusual signatures of
793 highly adaptable R-loci in closely-related Arabidopsis species. *Gene* **482**: 24–33.
- 794 **Wu, C.-H., Abd-El-Haliem, A., Bozkurt, T.O., Belhaj, K., Terauchi, R., Vossen, J.H., and**
795 **Kamoun, S.** (2017). NLR network mediates immunity to diverse plant pathogens. *Proc.*
796 *Natl. Acad. Sci. U. S. A.* **114**: 8113–8118.
- 797 **Yang, S., Li, J., Zhang, X., Zhang, Q., Huang, J., Chen, J.-Q., Hartl, D.L., and Tian, D.**
798 (2013). Rapidly evolving R genes in diverse grass species confer resistance to rice blast
799 disease. *Proc. Natl. Acad. Sci. U. S. A.* **110**: 18572–18577.
- 800 **Zhang, P., Hiebert, C.W., McIntosh, R.A., McCallum, B.D., Thomas, J.B., Hoxha, S.,**
801 **Singh, D., and Bansal, U.** (2016). The relationship of leaf rust resistance gene Lr13 and
802 hybrid necrosis gene Ne2m on wheat chromosome 2BS. *Theor. Appl. Genet.* **129**: 485–493.
- 803

804 **Table 1.** Number and locations of highly variable residues in hvNLR receptors. The number of
805 residues in clade alignment for each hvNLR with Shannon entropy values of at least 1.5 bits
806 (counted by domain) is shown. The majority of highly variable residues were found in the LRR
807 domain.

Gene Name	Type	preNB		NB-ARC		linker		LRR		postLRR	
		No hv aa	% total aa								
RPP9	TIR	0	0	0	0	0	0	23	5.8	11	5.3
RPP7	CC	0	0	0	0	1	1.5	34	6.1	0	0
AT1G58807.1	CC	1	0.6	0	0	0	0	29	6.7	1	3.4
AT1G58848.1	CC	1	0.6	0	0	0	0	37	7.2	0	0
AT1G59124.1	CC	1	0.6	0	0	0	0	17	5.6	0	0
AT1G59218.1	CC	1	0.6	0	0	0	0	36	7.1	1	7.7
AT1G61180.1	CC	2	1.3	7	2.1	0	0	35	9.9	0	0
RPP39	CC	2	1.3	4	1.2	0	0	36	8.7	1	2.7
AT1G61300.1	CC	2	4.8	7	2.1	0	0	35	9.9	0	0
AT1G61310.1	CC	2	1.3	7	2	0	0	35	9.9	0	0
AT1G62630.1	CC	0	0	4	1.2	0	0	23	7	2	4
AT1G69550.1	TIR	0	0	2	0.6	2	2.4	58	9.8	1	0.6
RPP28	TIR	1	0.4	0	0	1	3	18	3.7	5	3.4
AT3G44400.1	TIR	2	0.9	4	1.3	3	5.1	22	8.1	18	11.6
RPP1	TIR	3	1.1	6	1.9	3	5.1	35	9.6	17	9.1
AT3G44630.1	TIR	3	1.1	6	1.8	3	5.1	35	9.5	15	8.3
AT3G44670.1	TIR	4	1.5	4	1.3	3	5.1	30	8.9	19	8
RPP13	CC	0	0	0	0	1	2.5	34	11.6	0	0
RPP4	TIR	3	1.6	5	1.7	5	8.3	45	8.4	1	1.6
SNC1	TIR	6	3.2	5	1.7	5	8.5	34	5.5	5	3.6
AT4G16920.1	TIR	7	3.8	5	1.7	5	8.5	51	8.3	3	2.1
RPP5	TIR	7	3.7	5	1.7	5	8.5	41	6.6	8	2.9
AT5G38350.1	TIR	0	0	3	0.9	1	1.7	13	4.6	6	4.1
SSI4-LIKE	TIR	0	0	3	1	0	0	21	6.6	4	2.2
AT5G41750.1	TIR	0	0	2	0.6	1	1.9	19	6	2	1
RPP8	CC	0	0	10	2.9	0	0	19	5	0	0
AT5G43740.1	CC	0	0	2	0.6	2	4.1	27	8.2	0	0
AT5G46510.1	TIR	1	0.5	1	0.3	1	0.9	7	2.3	6	1.4
VICTR/ACQOS	TIR	1	0.5	1	0.3	1	0.9	7	2.3	6	2.3
AT5G48620.1	CC	0	0	10	2.9	0	0	19	5.3	0	0

808

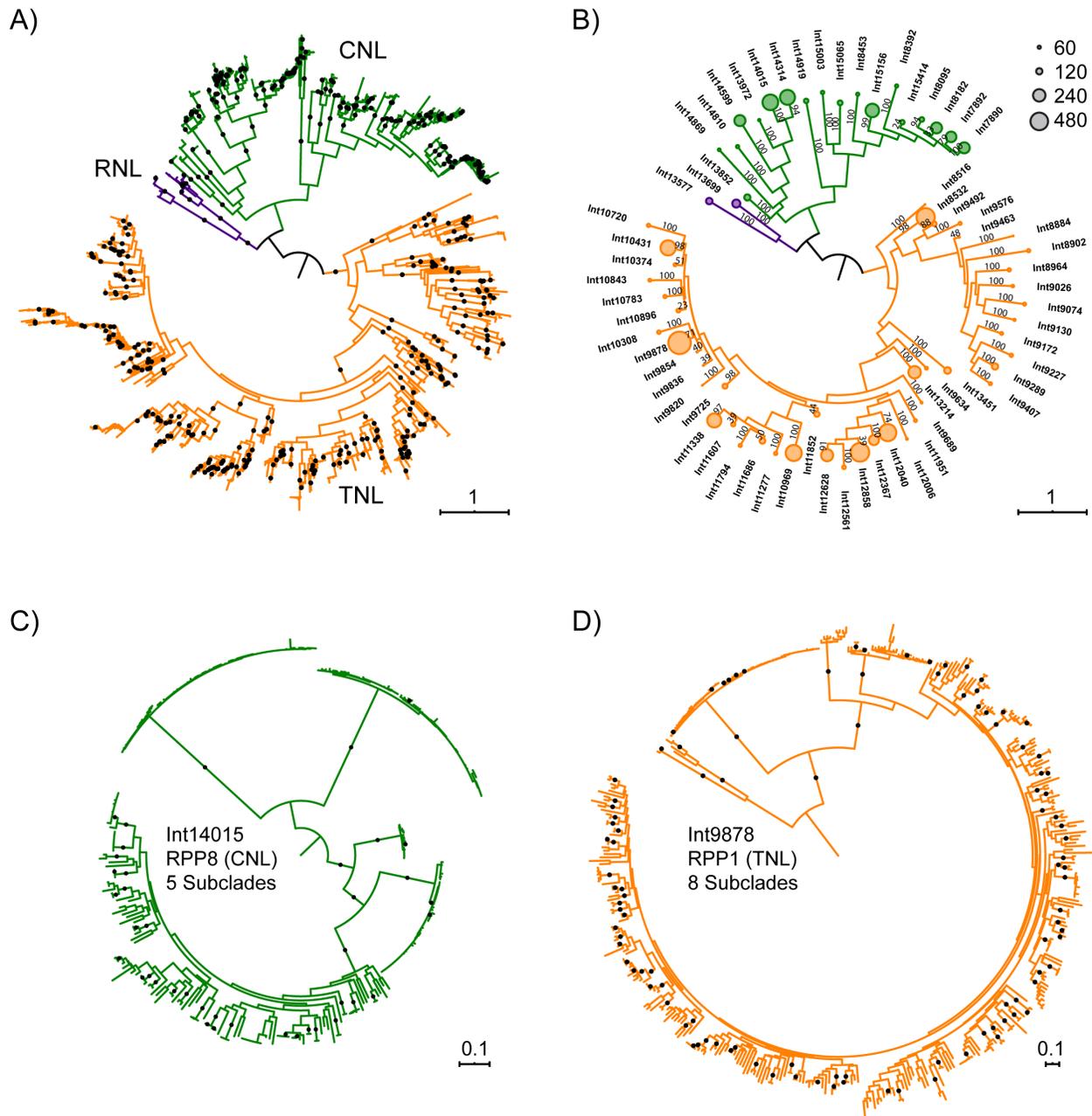


Figure 1. Phylogenetic analyses of Arabidopsis pan-NLRome. A) Maximum likelihood tree for 7,818 Arabidopsis NB-ARC sequences rooted on a branch connecting TNL and non-TNL clades. 99% or better bootstrap values are shown as dots. B) Same tree as in A) partitioned into 65 initial clades, with circle radius proportional to clade size, and indicating bootstrap support for each clade. C) Int14015 clade tree (rooted midpoint) based on a full-length alignment of the clade sequences. 99% or better bootstrap values are shown as dots. D) Int9878 clade ML tree (rooted midpoint) based on a full-length alignment of the clade sequences. 99% or better bootstrap values are shown as dots; branch length represents the number of substitutions per site.

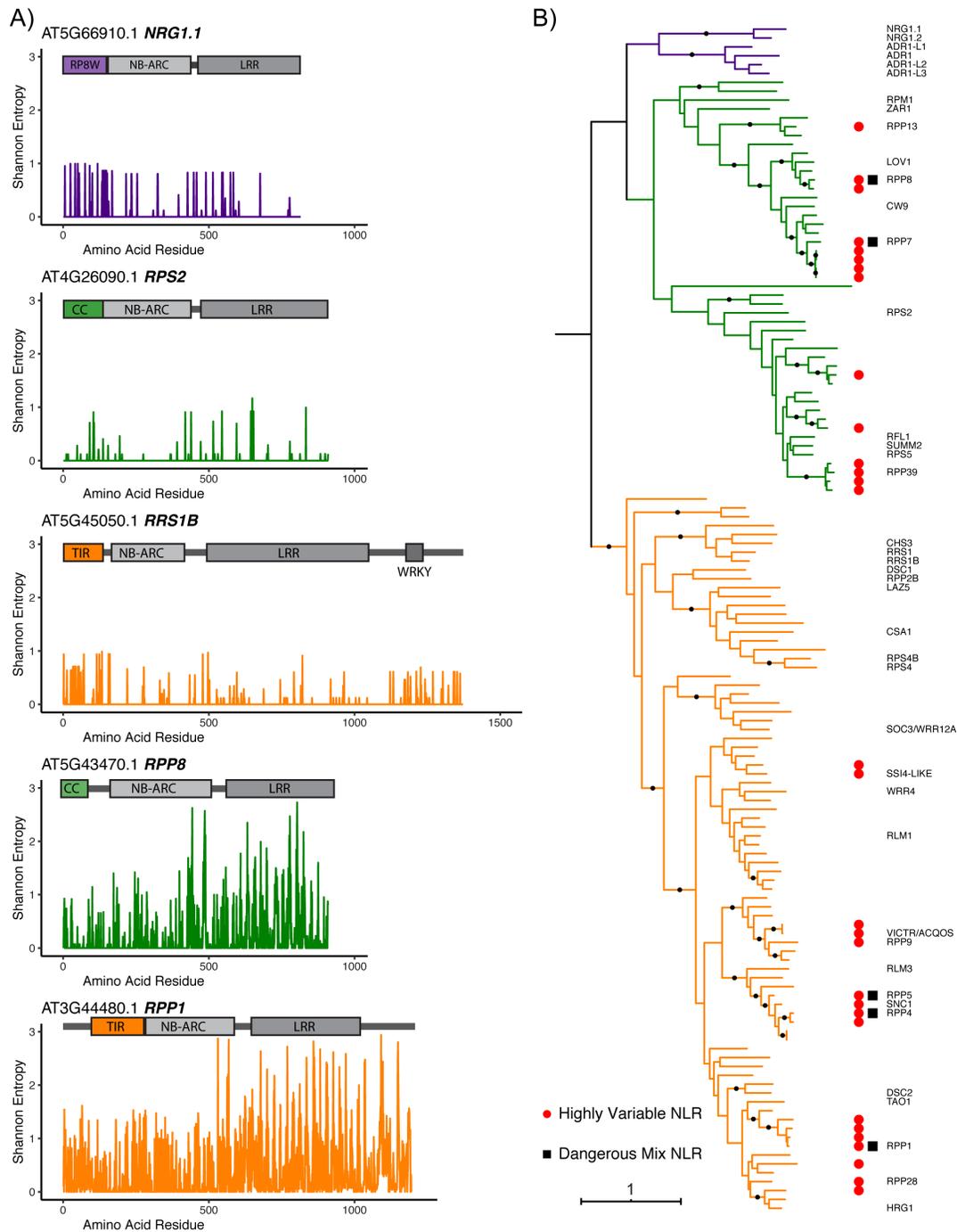


Figure 2. Identification and phylogenetic distribution of highly variable NLRs. A) Domain diagrams and Shannon entropy plots of clade alignments containing known NLRs from ancient helper (*NRG1.1*), guard (*RPS2*), integrated decoy (*RRS1B*), and direct recognition (*RPP1*) functional groups. It is not presently known whether *RPP8* is a direct recognition NLR. B) Phylogenetic distributions of NLRs of the reference ecotype, *Col-0*, indicating positions of known genes and showing the locations of hvNLRs and autoimmune Dangerous Mix NLRs. 99% or better bootstrap values are shown as dots; branch length represents the number of substitutions per site.

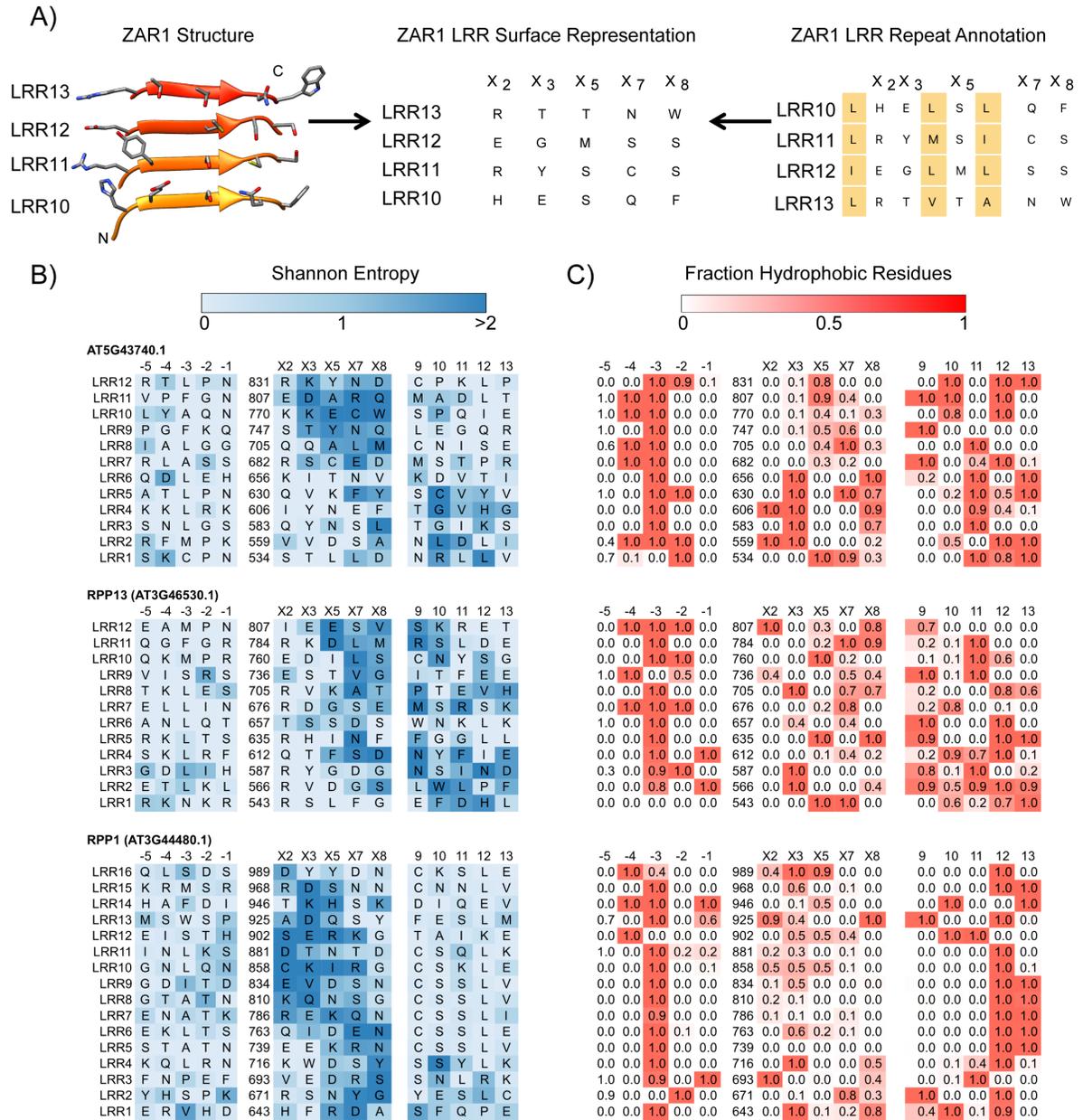


Figure 3. 2D representations of LRR surfaces allow comparisons of predicted NLR binding sites to be made in the absence of experimental structures. A) Beta-sheet on the concave side of ZAR1 LRR domain shows a regular array of surface-exposed residues that correspond to the variable positions in the LxxLxLxx LRR motif (left). Single-letter amino acid representation of the observed array (center). Identical representation is obtained from LRR repeat annotation by arranging the rows from bottom to top and hiding the columns containing conserved leucines (right). B) Shannon entropy scores and amino acid residues of three representative Col-0 hvNLRs mapped onto the 2D surface representation, including five additional amino acids on either side of the core repeat unit. C) Percentages of hydrophobic residues in the alignments of the same three proteins.

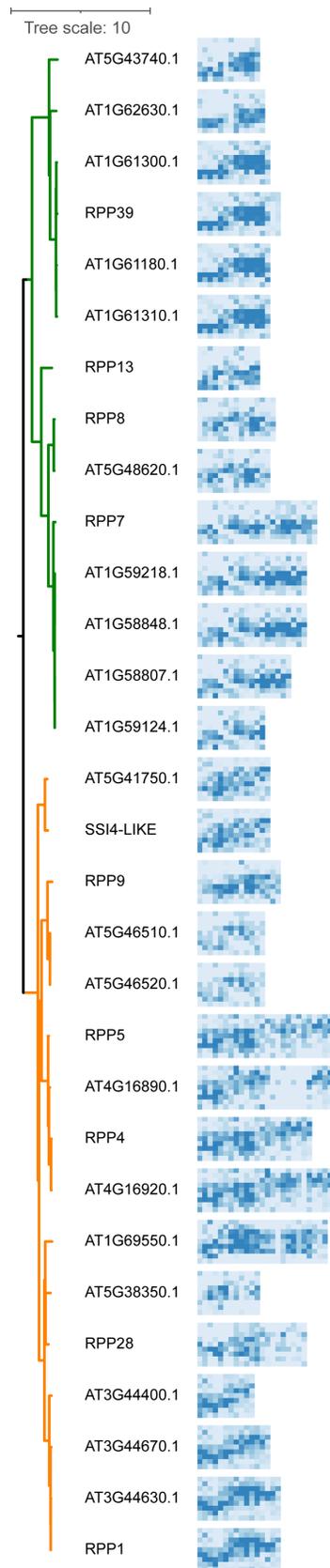


Figure 4. 2D representations of Col-0 hvNLR LRR surfaces in the context of the Col-0 NLR tree. The 2D binding site representations are those in Figure 3 and Supplemental Data Set 3 situated horizontally and trimmed to include positions -2, -1, 2, 3, 5, 7, 8, 10, and 11 of each repeat unit. For each cartoon the -2 position of the LRR1 is in the top left corner and the position 11 of the last LRR is in the bottom right. The tree on the left is a subset of the Col-0 NLR tree from Figure 2B with only the hvNLR leaves shown.

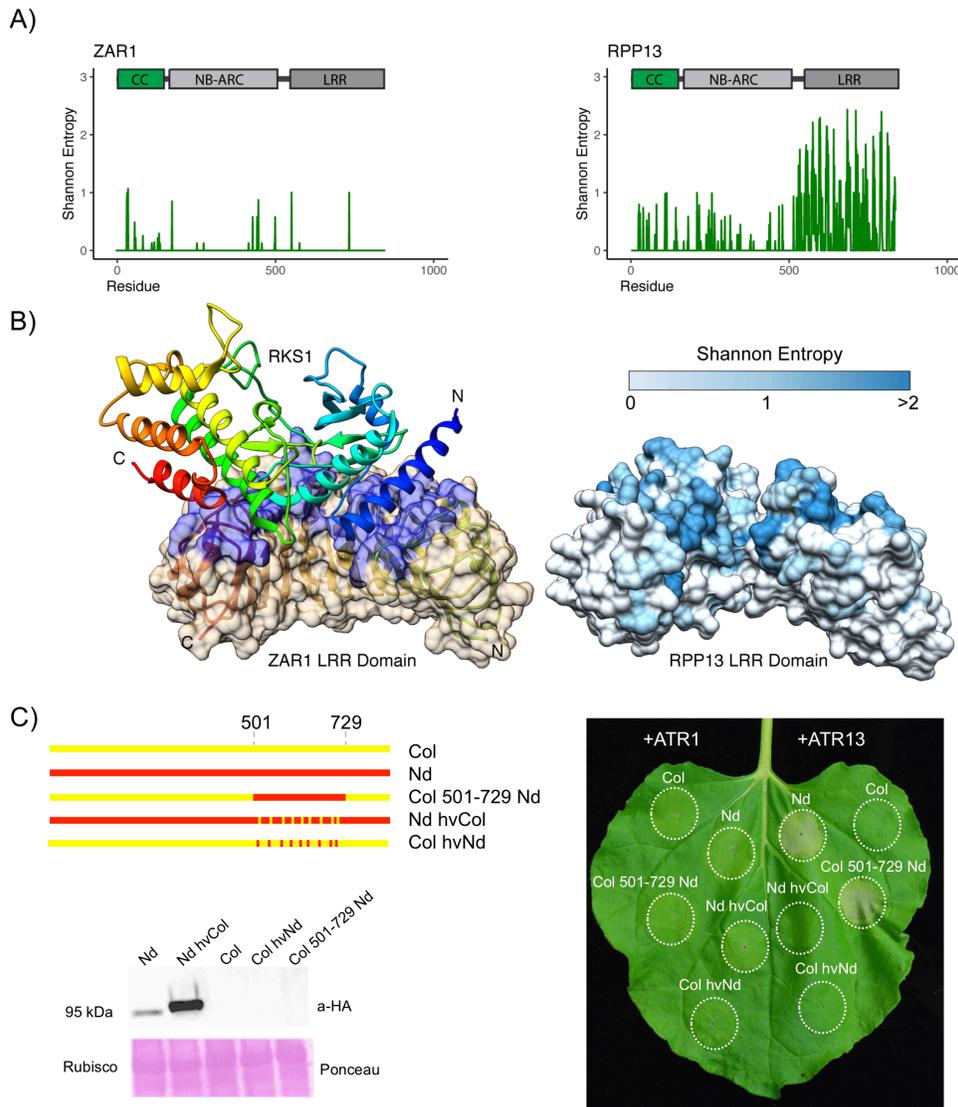


Figure 5. Highly variable residues in RPP13 overlap with the observed ZAR1-RKS1 binding site and are required for ATR13 recognition. A) Shannon entropy plots and domain diagrams for ZAR1, an indirect recognition CNL, and RPP13, a related hvNLR. B) Cryo-EM structure of RKS1 bound to ZAR1 (CC and NB-ARC domains omitted for clarity) (PDB ID: 6J5W). RKS1 shown as a secondary structure diagram with rainbow coloring from blue (N-terminus) to red (C-terminus), ZAR1 LRR as a secondary structure diagram and transparent surface with RKS1 contact residues colored blue. RPP13 LRR domain homology model with surface oriented as in ZAR1 and colored by Shannon Entropy of the RPP13 clade alignment from low (light blue) to high entropy (dark blue). C) Chimeric constructs of RPP13 region 501-729 containing highly variable LRR repeats. The constructs were designed by targeting amino acids with Shannon entropy >1.5 bits and functionally tested by *Agrobacterium*-mediated transient expression assays in *Nicotiana benthamiana* together with cognate ATR13d41-Emco5 effector or ATR1d51-Emoy2 negative control at the final OD600 of 0.6 with constructs mixed in equal ratio. The image was taken at 3 days post infiltration. Each construct was tested on 14 leaves and showed consistent presence/absence of HR on all leaves. Immunoblotting showed stable expression of both functional and mutated RPP13-Nd variants. No RPP13-Col variants could be detected despite having an intact HA tag similar to what has been reported previously (Rentel et al. 2008).

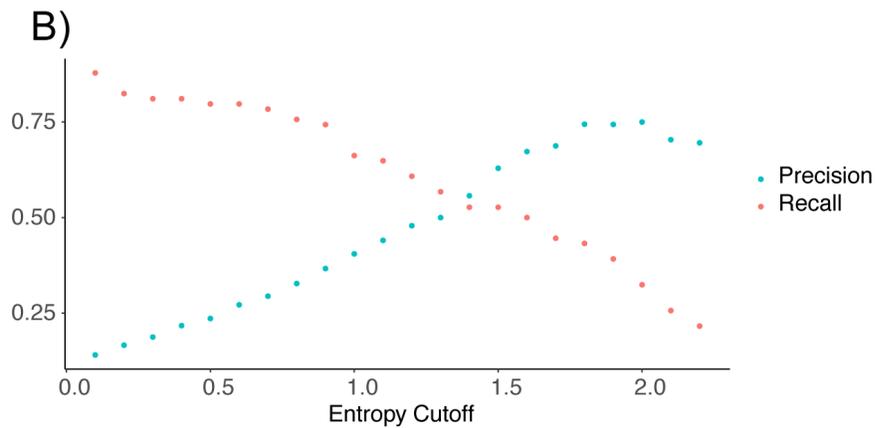
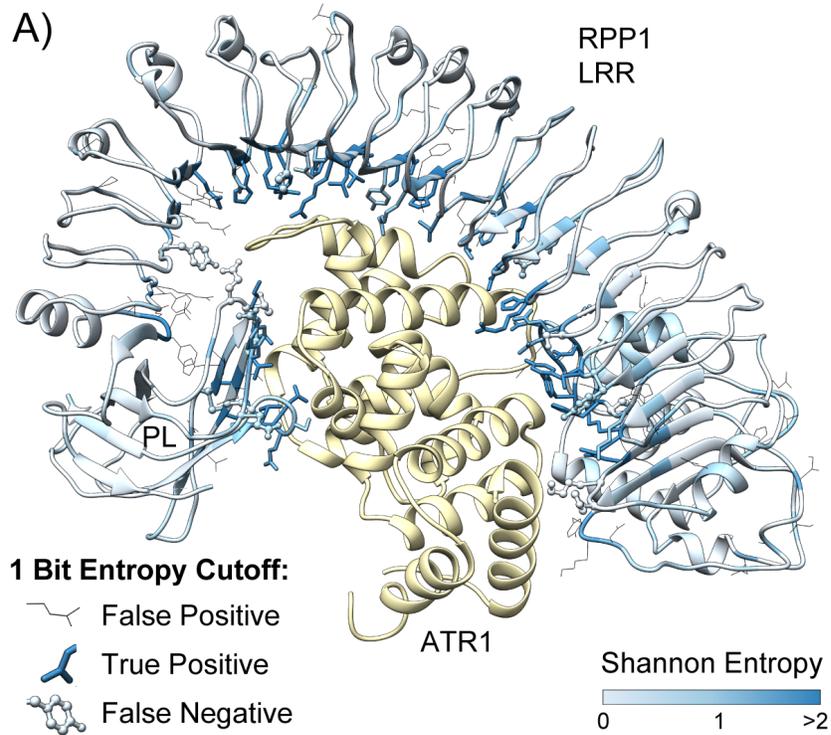


Figure 6. RPP1 contact residues show high sequence diversity. A) Structure of RPP1 LRR-ATR1 complex (PDB ID: 7CRB) colored by entropy scores, with contact residues shown as sticks for predicted true positives and ball and stick for false negatives using a 1 bit entropy cutoff. False positive predictions at the same cutoff are represented as wire. B) Precision and recall for the prediction of RPP1-ATR1 binding site residues based on the choice of entropy cutoff.

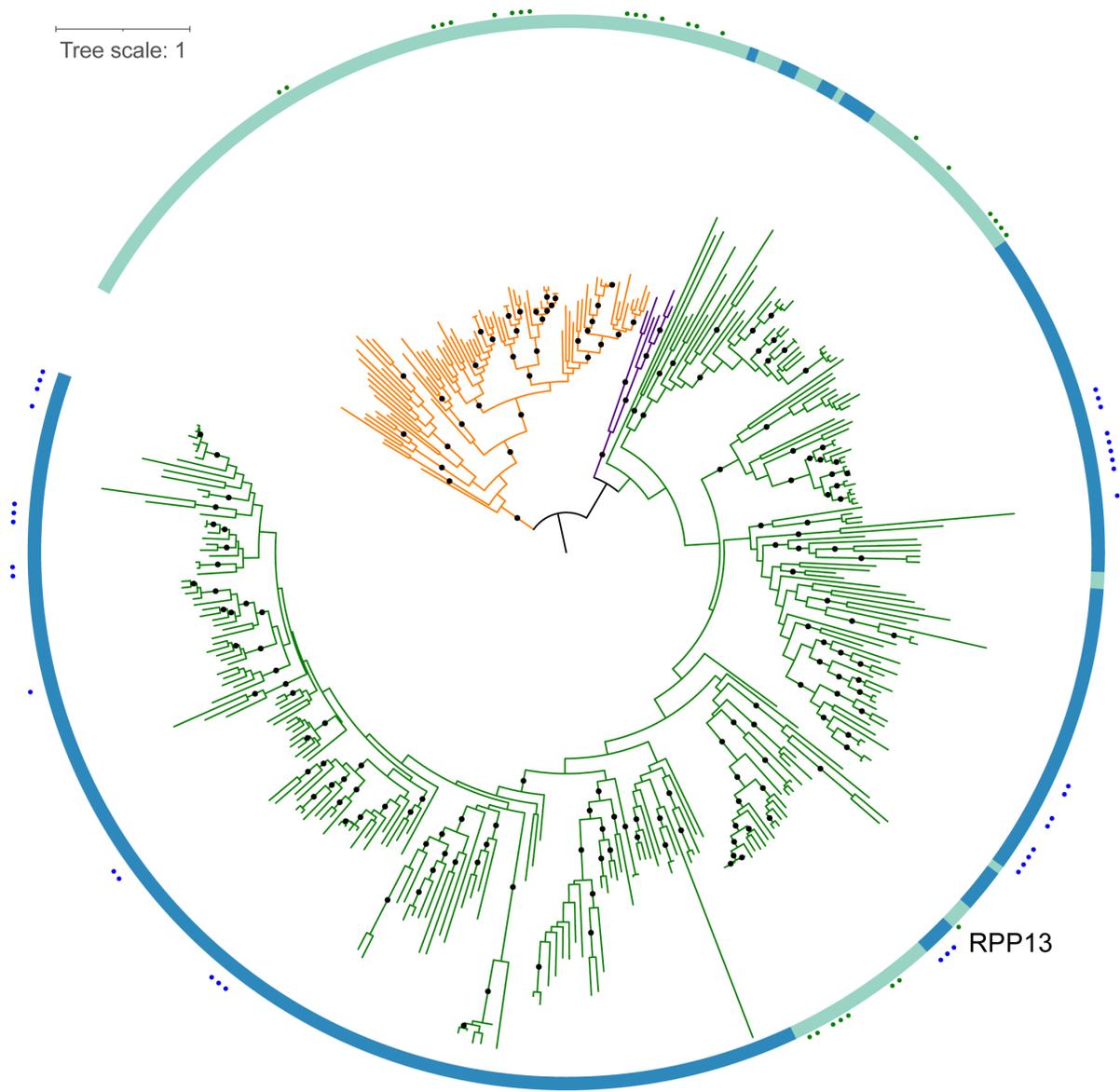


Figure 7. Dispersed distribution of hvNLRs in a joint phylogeny of *Brachypodium* Bd21 (blue ribbon) and *Arabidopsis* Col-0 (green ribbon). The *Arabidopsis* hvNLR clades (green dots) and *Brachypodium* hvNLRs (blue dots) do not cluster except for the RPP13 CNL clades. The tree is rooted arbitrarily on a branch connecting TNL clade (orange branches) and non-TNL clades (RNL branches are shown in purple, and CNL branches are shown in green). 99% or better bootstrap values are shown as dots; branch length represents the number of substitutions per site.