

1 **Influence of neighboring small sequence variants on functional impact prediction**

2 Jan-Simon Baasner¹, Dakota Howard^{1,2} and Boas Pucker^{1,3}

3

4 1 Genetics and Genomics of Plants, Center for Biotechnology, Bielefeld University, Bielefeld,
5 Germany

6 2 Biology and Computer Science Department, Furman University, Greenville, South Carolina,
7 USA

8 3 Evolution and Diversity, Plant Sciences, University of Cambridge, Cambridge, United Kingdom

9

10 JSB: janbaas@cebitec.uni-bielefeld.de

11 DH: dhoward@cebitec.uni-bielefeld.de

12 BP: bpucker@cebitec.uni-bielefeld.de

13

14 Corresponding author: Boas Pucker, bpucker@cebitec.uni-bielefeld.de

15

16 ORCIDs:

17 JSB: 0000-0003-3996-6817

18 DH: 0000-0002-7674-0385

19 BP: 0000-0002-3321-7471

20

21

22

23

24 **Abstract**

25 Once a suitable reference sequence is generated, genomic differences within a species are
26 often assessed by re-sequencing. Variant calling processes can reveal all differences between
27 two strains, accessions, genotypes, or individuals. These variants can be enriched with
28 predictions about their functional implications based on available structural annotations i.e. gene
29 models. Although these functional impact predictions on a per variant basis are often accurate,
30 some challenging cases require the simultaneous incorporation of multiple adjacent variants into
31 this prediction process. Examples are neighboring variants which modify each others' functional
32 impact. Neighborhood-Aware Variant Impact Predictor (NAVIP) considers all variants within a
33 given protein coding sequence when predicting the functional consequences. As a proof of
34 concept, variants between the *Arabidopsis thaliana* accessions Columbia-0 and Niederzenz-1
35 were annotated. NAVIP is freely available on github: <https://github.com/bpucker/NAVIP>.

36

37 **Introduction**

38 Re-sequencing projects e.g. investigating many individuals or accessions of one species [1–3]
39 are gaining relevance in plant research. Approaches similar to genome-wide association studies
40 which are based on mapping-by-sequencing (MBS) were frequently applied [4–6]. They are
41 boosted by an increasing availability of high quality reference genome sequences [7–12] and
42 dropping sequencing costs [13, 14]. *De novo* assemblies are still beneficial for the detection of
43 large structural variants [8, 11, 12, 15–17] and especially to reveal novel sequences [8, 11, 12,
44 18], but the reliable detection of modifying single nucleotide variants (SNVs) can be achieved
45 based on (short) read mappings.

46 Once identified, the annotation of sequence variants in most species is performed by predicting
47 their functional implications based on the available annotation of genes. Leading tools like
48 ANNOVAR [19] and SnpEff [20] are currently performing this prediction by focusing on a single
49 variant at a time. An impact prediction facilitates the identification of targets for post-GWAS
50 analyses [21, 22]. Although the effect prediction for single variants is very efficient and usually
51 correct, there is a minority of challenging cases in which predictions cannot be accurate based
52 on a single variant alone. Multiple InDels could either lead to frameshifts or they compensate for
53 each others' effect leaving the sequence with minimal modifications [23–25]. Two SNVs

54 occurring in the same codon could lead to a different amino acid substitution compared to the
55 apparent effect resulting from an isolated analysis of each of these SNVs.

56 Here we present a new tool to accurately predict the combined effect of phased variants on
57 annotated coding sequences. Neighborhood-Aware Variant Impact Predictor (NAVIP) was
58 developed to investigate large variant data sets of plant re-sequencing projects, but is not
59 limited to the annotation of variants in plants. As a proof of concept, NAVIP was deployed to
60 identify cases between the *A. thaliana* accessions Columbia-0 (Col-0) and Niederzenz-1 (Nd-1)
61 where an accurate impact prediction needs to consider multiple variants at a time [15].

62

63

64 **Materials & Methods**

65 Variant detection

66 Sequencing reads of Nd-1 [15] were mapped to the Col-0 reference genome sequence (TAIR9)
67 [26] via BWA MEM v.0.7.13 [27] using the `-m` option to avoid spurious hits. Variant calling was
68 performed via GATK v3.8 [28] based on the developers' recommendation. All processes were
69 wrapped into custom Python scripts (https://github.com/bpucker/variant_calling) to facilitate
70 automatic execution on a high performance compute cluster. An initial variant set was
71 generated based on hard filtering criteria recommended by the GATK developers. The two
72 following variant calling runs considered the set of surviving variants of the previous round as
73 gold standard to avoid the need for hard filtering.

74

75 Variant validation

76 Since a high quality genome sequence assembly of Nd-1 was recently generated [12], we
77 harnessed this sequence to validate all variants identified by short read mapping. Starting at the
78 north end of each chromosome sequence, sorted variants were tested one after the other by
79 taking the upstream sequence from Col-0, modifying it according to all upstream *bona fide*
80 variants, and searching for it in the Nd-1 assembly (AdditionalFile1). Variants were admitted to
81 the following analysis if the assembly supports them. This consecutive inspection of all variants
82 enabled a reliable removal of false positives.

83

84

85 Variant impact prediction

86 Our Neighborhood-Aware Impact Predictor (NAVIP, <https://github.com/bpucker/NAVIP>) takes a
87 VCF file containing sequence variants, a FASTA file containing the reference sequence, and a
88 GFF3 file containing the annotation as input. Provided variants need to be homozygous or in a
89 phased state to allow an accurate impact prediction per allele. Effects on all annotated
90 transcripts are assessed per gene by taking the presence of all given variants into account.
91 NAVIP generates a new VCF file with an additional annotation field and additional report files
92 including FASTA files with the resulting sequences (see manual for details:
93 <https://github.com/bpucker/NAVIP/wiki>).

94

95 Assessing predicted premature stop codons and frameshifts

96 SnpEff [20] was applied to the validated variant data set to predict the effects of single variants.
97 To assess the influence of the underlying annotation, this prediction was performed based on
98 TAIR10 [26] and Araport11 [29]. Predicted premature stop codons with two variants within the
99 same codon were selected for comparison to the NAVIP prediction, because these cases have
100 the potential to show different results.

101 Transcripts with predicted frameshifts were analyzed to identify downstream insertions/deletions
102 which are compensating each others' effect i.e. the second frameshift is reverting an upstream
103 frameshift. The distance between these events was analyzed by the third module of NAVIP.

104

105 Experimental validation of variants

106 *A. thaliana* Nd-1 plants were grown as previously described [15]. DNA for PCR experiments was
107 extracted from leaf tissue using a cetyltrimethylammonium bromide (CTAB)-based method as
108 previously described [30]. Oligonucleotides flanking regions with variants of interest were
109 designed manually (AdditionalFile2) and purchased from Metabion (<http://www.metabion.com/>).
110 Amplification via PCR, analysis of PCR products, purification of PCR products, Sanger
111 sequencing, and evaluation of results was performed as previously described [31].

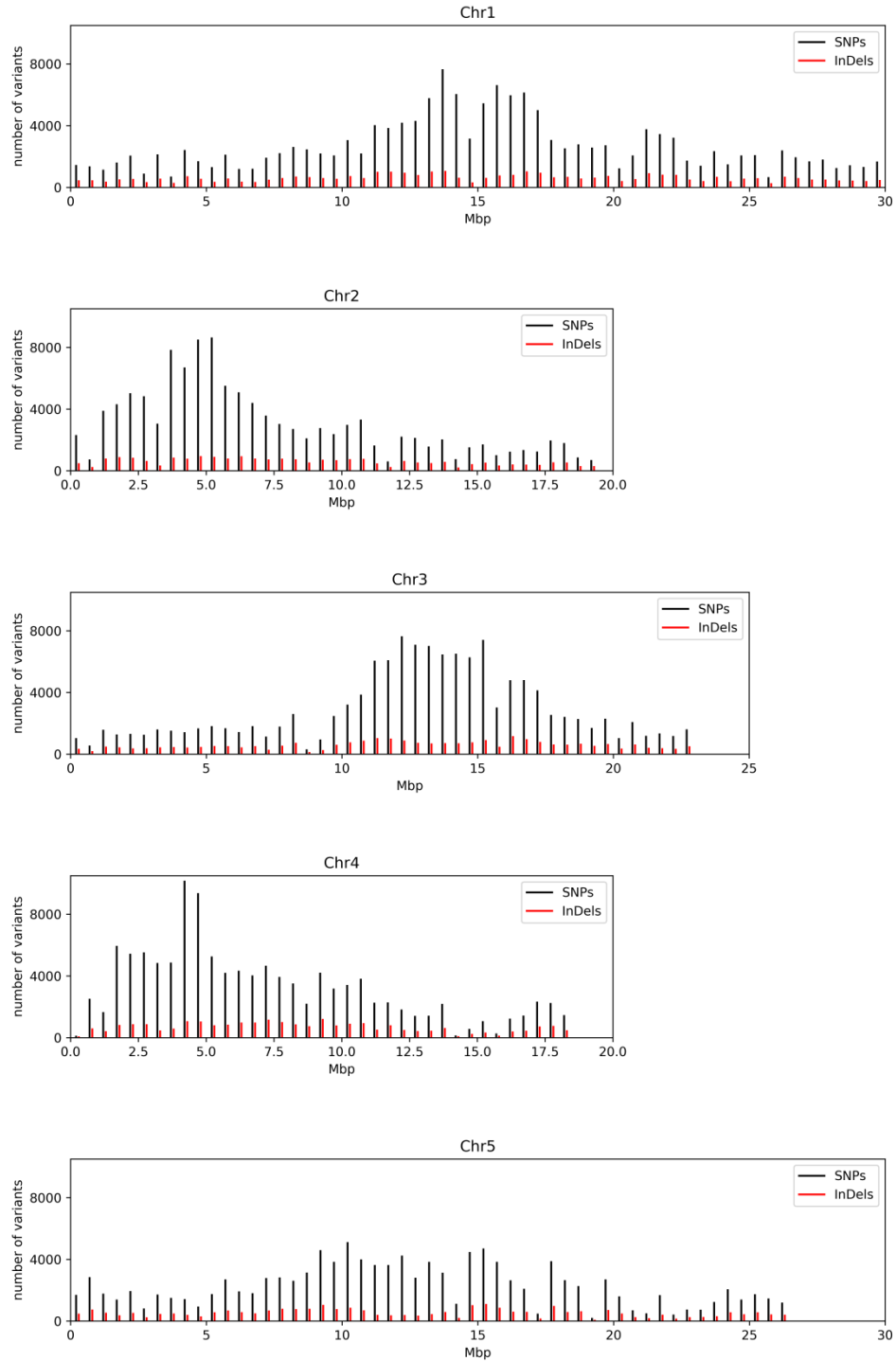
112

113

114 **Results**

115 Variant detection and validation

116 Nd-1 reads were mapped against the Col-0 reference genome sequence (TAIR9). Based on
117 124,662,140 mapped paired-end reads, 384,622 variants were detected in the first variant
118 calling round of this study. This initial set was extended over three additional rounds of variant
119 calling leading to over one million of variants. The variant calling was stopped, because no
120 substantial increase in the number of novel variants was observed during the last rounds. An
121 assembly based on independent Single Molecule Real Time (SMRT) sequencing reads
122 supported 772,644 (76.6%) of all variants detected during the last iteration (AdditionalFile3, Fig.
123 1). On average, one variant was observed every 154 bp between Col-0 and Nd-1. SNV
124 frequencies ranged from one event in 225 bp on Chr5 to one event in 158 bp on Chr4. InDel
125 frequencies ranged from one event in 1,051 bp on Chr5 to one event in 809 bp on Chr4.



126

127 **Fig. 1: Genome-wide distribution of sequence variants between Col-0 and Nd-1.**

128 Distributions of SNVs and InDels over the chromosome sequences of Col-0 were visualized as previously

129 described [15].

130

131 Although the repeated variant calling processes were intended to increase the sensitivity, we did
132 not observe a substantial improvement between the second and third round. This saturation
133 indicates that no additional variants would be detected in further variant calling rounds. The
134 number of detected variants as well as the validation rate was almost constant (Table 1).

135

136 **Table 1: Total and validated number of variants.**

Variant data set	Total variants	Validated variants
Initial set based on hard filtering	384,617	350,005 (90.1%)
Soft filtering round 1	1,006,920	771,449 (76.6%)
Soft filtering round 2	1,008,610	772,612 (76.6%)
Soft filtering round 3	1,008,629	772,643 (76.6%)

137

138 Experimental validation

139 Randomly selected loci with two SNVs within one codon were experimentally validation via PCR
140 and amplicon sequencing (Table 2). Successful sequencing reactions show a validation rate of
141 >95%.

142

143 **Table 2: Neighboring SNVs validated in Nd-1 via PCR and amplicon sequencing.** Sequences of
144 oligonucleotides used for the amplicon generation are listed in AdditionalFile2.

AGI	Fw primer	Rv primer	Status
At1g30545	N400	N401	Validated
At3g55500	N402	N403	Validated
At3g26770	N406	N407	Validated
At4g30570	N408	N409	Validated
At1g28150	N410	N411	Validated
At1g35430	N412	N413	Validated
At4g27230	N414	N415	One SNV failed
At5g60230	N424	N425	2 validated
At1g31820	N426	N427	4 validated / 1 failed

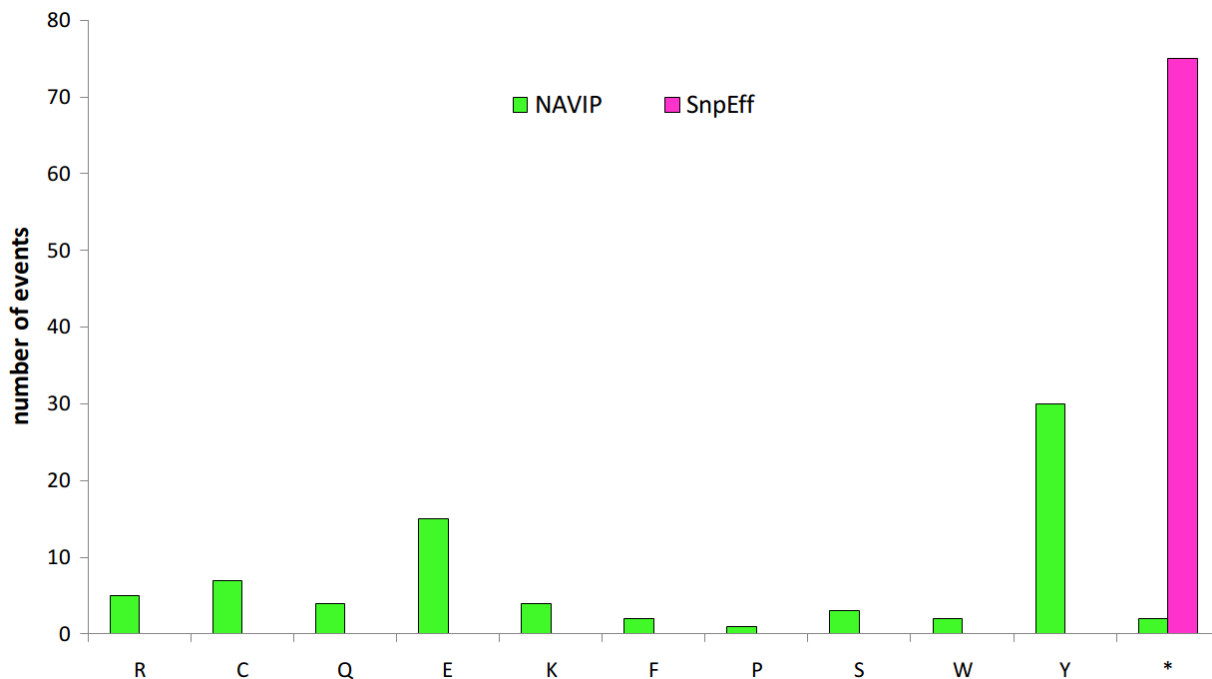
145

146

147 Relevance of NAVIP

148 Running NAVIP on this *A. thaliana* data set (AdditionalFile4) took about 5 minutes with a single
149 core and a peak memory usage of about 3 GB RAM. Since SnpEff is one of the most frequently
150 applied tools for the annotation of variants, the NAVIP output was compared with SnpEff
151 predictions. SnpEff was applied to the same data set based on the Araport11 annotation.
152 Interesting cases for comparison are codons containing at least two SNVs. Of 75 premature
153 stop codons predicted in such codons by SnpEff, 73 were predicted as amino acid substitutions
154 by NAVIP (Fig. 2). While a single SNV would cause a premature stop codon, the simultaneous
155 presence of two SNVs results in an amino acid encoding codon. In total, 702 premature stop
156 codons were predicted by SnpEff thus 9.6 % of them were false positives. NAVIP revealed that
157 tyrosine occurs frequently instead of a premature stop codon, because the tyrosine codons are
158 very similar to two of the three stop codons.

159



160

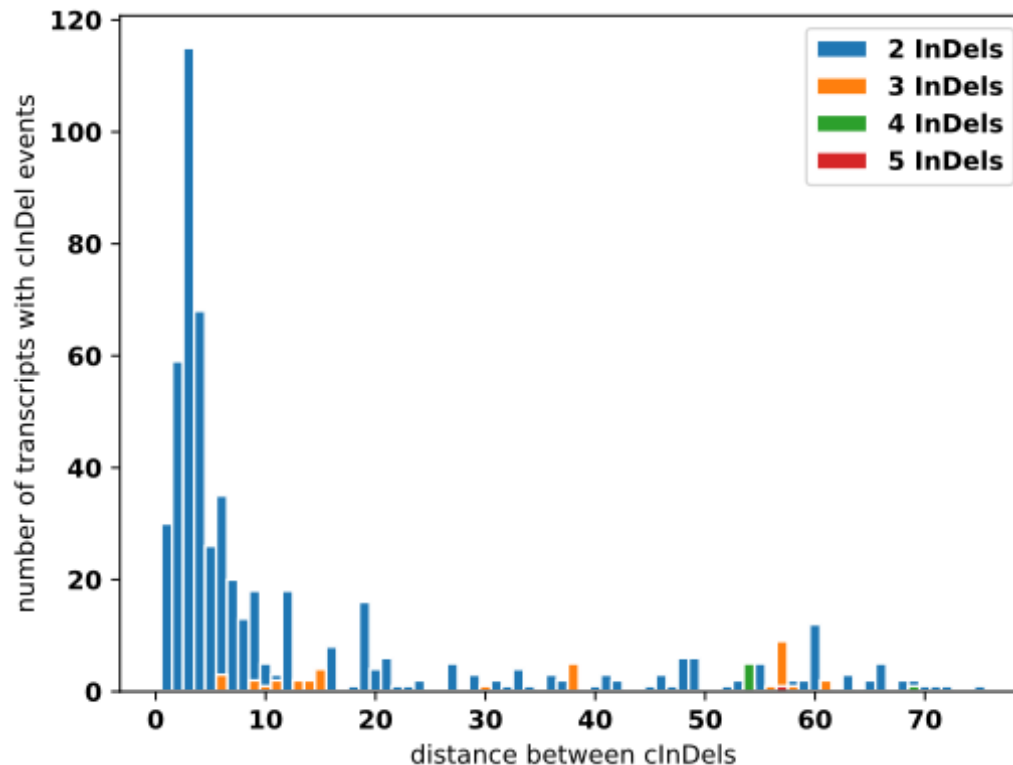
161 **Fig. 2: Second site variants turn predicted premature stop codons into amino acid substitutions.**

162 Premature stop codons predicted by SnpEff (pink) are frequently amino acid substitutions if a second
163 variant is located within the same codon. NAVIP revealed 73 false positive predictions of premature stop
164 codons by SnpEff which are in fact amino acid substitutions (green).

165

166 InDels can compensate each others' frameshift when occurring together. Since premature stop
167 codons can emerge by chance following a frameshift, the distance between such InDels was
168 analyzed. This length distribution revealed that most compensating InDels (cInDels) occur within
169 a short distance of 2-8 bp (Fig. 3). Multiples of three are more frequent than other distances of a
170 similar size.

171



172

173 **Fig. 3: Distance between compensating InDels (cInDels).**

174 An InDel can compensate the frameshift caused by an upstream InDel. Distances between such cInDels
175 are short and frequently multiples of three. In total, 484 genes contain cInDels.

176

177

178

179

180 **Discussion**

181 Variant validation, frequency, and distribution

182 Although differentiation between *bona fide* variants (true positives) and false positives based on
183 a high quality genome sequence assembly worked very well, false negatives were not taken into
184 account and might even bias this classification approach by preventing the validation of
185 neighboring variants (AdditionalFile1). If a variant is missed by the initial variant calling, its
186 presence in the flanking sequence used during the validation process will prevent a proper
187 match. Therefore, the number of variants could be slightly higher than reported here.
188 Nevertheless, this conservative approach was selected to minimize the risk of keeping false
189 positive variants. There is always a trade-off between sensitivity and specificity in the variant
190 calling process [32] and our approach is in strong favor of specificity. However, the number of
191 identified and validated variants exceeds previous reports of 485,887 variants between Col-0
192 and Nd-1 [15]. Instead the observed variant frequency is closer to the results of a comparison
193 between Bur-0 and Col-0 [33]. Despite the difference in total numbers, the distribution on the
194 chromosome scale is similar to the previous comparison of Col-0 and Nd-1 [15]. It seems that
195 Chr4 is the most variable one, while Chr5 is the least variable one between both compared
196 accessions.

197 Successful validation via PCR and amplicon sequencing supported the presence of two SNVs
198 within one codon. Although these variants are perceived as two SNVs, the underlying
199 mechanism could be a multiple nucleotide polymorphism (MNP). It would be interesting to see if
200 these SNVs occur independently in other accessions in the *A. thaliana* population.

201

202 Functional implications of variants

203 We developed NAVIP to assess the impact of neighboring variants on protein coding
204 sequences. The presence of the 557 cases described here for the comparison of two *A. thaliana*
205 accessions demonstrates the necessity to have such a tool at hand. NAVIP revealed the

206 presence of second site mutations that compensate other variants e.g. turning a premature stop
207 codon into an amino acid substitution or compensation of a frameshift. The purpose of NAVIP is
208 not to replace existing tools, but to add novel functionalities to established tools like SnpEff [20].
209 This could boost the power of re-sequencing studies by opening up the field of compensating or
210 in general mutually influencing variants. Such variants have the potential to reveal new insights
211 into patterns of molecular evolution and especially co-evolution of sites. Although the number of
212 cases is probably small, the consideration of multiple variants during the effect prediction could
213 reveal novel targets in GWAS-like approaches. The remaining challenge is now the reliable
214 detection of sequence variants prior to the application of NAVIP. For heterozygous species
215 phasing of these variants is another task that needs to be addressed. The correct prediction of
216 functional implications relies on the correct assignment of variants to respective haplophases. If
217 provided with accurately phased variants, NAVIP can perform predictions for highly
218 heterozygous and even polyploid species.

219

220 **Availability of data**

221 The data sets supporting the results of this article are included within the article and its
222 additional files. Python scripts developed and applied for this study are available on github:
223 <https://github.com/bpucker/NAVIP> (<https://doi.org/10.5281/zenodo.2620396>)
224 https://github.com/bpucker/variant_calling (<https://doi.org/10.5281/zenodo.2616418>).

225

226 **Authors' contribution**

227 BP designed research. JSB wrote the NAVIP code. JSB, DH, and BP conducted bioinformatic
228 analyses. DH and BP performed experimental validation. BP wrote the manuscript. All authors
229 read and approved the final version.

230

231 **Acknowledgements**

232 We acknowledge support by members of Genetics and Genomics of Plants, Bioinformatics
233 Resource Facility, and Sequencing Core Facility at the Center of Biotechnology. We thank
234 Hanna Schilbert for critical reading of the manuscript.

235

236

237

238 **References**

239 1. Alonso-Blanco C, Andrade J, Becker C, Bemm F, Bergelson J, Borgwardt KM, et al. 1,135 Genomes
240 Reveal the Global Pattern of Polymorphism in *Arabidopsis thaliana*. *Cell*. 2016;166:481–91.

241 2. Duan N, Bai Y, Sun H, Wang N, Ma Y, Li M, et al. Genome re-sequencing reveals the history of apple
242 and supports a two-stage model for fruit enlargement. *Nat Commun*. 2017;8:249.

243 3. Lobaton JD, Miller T, Gil J, Ariza D, de la Hoz JF, Soler A, et al. Resequencing of Common Bean
244 Identifies Regions of Inter-Gene Pool Introgression and Provides Comprehensive Resources for
245 Molecular Breeding. *Plant Genome*. 2018;11. doi:10.3835/plantgenome2017.08.0068.

246 4. James GV, Patel V, Nordström KJ, Klasen JR, Salomé PA, Weigel D, et al. User guide for mapping-by-
247 sequencing in *Arabidopsis*. *Genome Biol*. 2013;14:R61.

248 5. Mascher M, Jost M, Kuon J-E, Himmelbach A, Aßfalg A, Beier S, et al. Mapping-by-sequencing
249 accelerates forward genetics in barley. *Genome Biol*. 2014;15:R78.

250 6. Wu Y, Zheng Z, Visscher PM, Yang J. Quantifying the mapping precision of genome-wide association
251 studies using whole-genome sequencing data. *Genome Biol*. 2017;18:86.

252 7. Dohm JC, Minoche AE, Holtgräwe D, Capella-Gutiérrez S, Zakrzewski F, Tafer H, et al. The genome of
253 the recently domesticated crop plant sugar beet (*Beta vulgaris*). *Nature*. 2014;505:546–9.

254 8. Zapata L, Ding J, Willing E-M, Hartwig B, Bezdán D, Jiao W-B, et al. Chromosome-level assembly of
255 *Arabidopsis thaliana* Ler reveals the extent of translocation and inversion polymorphisms. *Proc Natl*
256 *Acad Sci*. 2016;113:E4052–60.

257 9. Jarvis DE, Ho YS, Lightfoot DJ, Schmöckel SM, Li B, Borm TJA, et al. The genome of *Chenopodium*
258 *quinoa*. *Nature*. 2017;542:307–12.

- 259 10. Lightfoot DJ, Jarvis DE, Ramaraj T, Lee R, Jellen EN, Maughan PJ. Single-molecule sequencing and Hi-
260 C-based proximity-guided assembly of amaranth (*Amaranthus hypochondriacus*) chromosomes provide
261 insights into genome evolution. *BMC Biol.* 2017;15:74.
- 262 11. Michael TP, Jupe F, Bemm F, Motley ST, Sandoval JP, Lanz C, et al. High contiguity *Arabidopsis*
263 *thaliana* genome assembly with a single nanopore flow cell. *Nat Commun.* 2018;9:541.
- 264 12. Pucker B, Holtgraewe D, Stadermann KB, Frey K, Huettel B, Reinhardt R, et al. A Chromosome-level
265 Sequence Assembly Reveals the Structure of the *Arabidopsis thaliana* Nd-1 Genome and its Gene Set.
266 *PLOS ONE.* 2019: e0216233. doi:10.1371/journal.pone.0216233.
- 267 13. Stein LD. The case for cloud computing in genome informatics. *Genome Biol.* 2010;11:207.
- 268 14. Christensen KD, Dukhovny D, Siebert U, Green RC. Assessing the Costs and Cost-Effectiveness of
269 Genomic Sequencing. *J Pers Med.* 2015;5:470–86.
- 270 15. Pucker B, Holtgräwe D, Sörensen TR, Stracke R, Viehöver P, Weisshaar B. A *De Novo* Genome
271 Sequence Assembly of the *Arabidopsis thaliana* Accession Niederzenz-1 Displays Presence/Absence
272 Variation and Strong Synteny. *PLOS ONE.* 2016;11:e0164321.
- 273 16. Fan X, Chaisson M, Nakhleh L, Chen K. HySA: a Hybrid Structural variant Assembly approach using
274 next-generation and single-molecule sequencing technologies. *Genome Res.* 2017;27:793–800.
- 275 17. Wala JA, Bandopadhyay P, Greenwald NF, O'Rourke R, Sharpe T, Stewart C, et al. SvABA: genome-
276 wide detection of structural variants and indels by local assembly. *Genome Res.* 2018;28:581–91.
- 277 18. Zhou Y, Massonnet M, Sanjak JS, Cantu D, Gaut BS. Evolutionary genomics of grape (*Vitis vinifera* ssp.
278 *vinifera*) domestication. *Proc Natl Acad Sci.* 2017;114:11715–20.
- 279 19. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-
280 throughput sequencing data. *Nucleic Acids Res.* 2010;38:e164.
- 281 20. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for annotating and
282 predicting the effects of single nucleotide polymorphisms, SnpEff. *Fly (Austin).* 2012;6:80–92.

- 283 21. Hou L, Zhao H. A review of post-GWAS prioritization approaches. *Front Genet.* 2013;4.
284 doi:10.3389/fgene.2013.00280.
- 285 22. Ries D, Holtgräwe D, Viehöver P, Weisshaar B. Rapid gene identification in sugar beet using deep
286 sequencing of DNA from phenotypic pools selected from breeding panels. *BMC Genomics.* 2016;17.
287 doi:10.1186/s12864-016-2566-9.
- 288 23. Cao J, Schneeberger K, Ossowski S, Günther T, Bender S, Fitz J, et al. Whole-genome sequencing of
289 multiple *Arabidopsis thaliana* populations. *Nat Genet.* 2011;43:956–63.
- 290 24. Gan X, Stegle O, Behr J, Steffen JG, Drewe P, Hildebrand KL, et al. Multiple reference genomes and
291 transcriptomes for *Arabidopsis thaliana*. *Nature.* 2011;477:419–23.
- 292 25. Schneeberger K, Ossowski S, Ott F, Klein JD, Wang X, Lanz C, et al. Reference-guided assembly of four
293 diverse *Arabidopsis thaliana* genomes. *Proc Natl Acad Sci U S A.* 2011;108:10249–54.
- 294 26. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The *Arabidopsis* Information
295 Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res.* 2012;40 Database
296 issue:D1202–10.
- 297 27. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
298 ArXiv13033997 Q-Bio. 2013. <http://arxiv.org/abs/1303.3997>. Accessed 16 Oct 2018.
- 299 28. Auwera GAV der, Carneiro MO, Hartl C, Poplin R, Angel G del, Levy-Moonshine A, et al. From FastQ
300 Data to High-Confidence Variant Calls: The Genome Analysis Toolkit Best Practices Pipeline. *Curr Protoc*
301 *Bioinforma.* 2013;43:11.10.1-11.10.33.
- 302 29. Cheng C-Y, Krishnakumar V, Chan AP, Thibaud-Nissen F, Schobel S, Town CD. Araport11: a complete
303 reannotation of the *Arabidopsis thaliana* reference genome. *Plant J.* 2017;89:789–804.
- 304 30. Rosso MG, Li Y, Strizhov N, Reiss B, Dekker K, Weisshaar B. An *Arabidopsis thaliana* T-DNA
305 mutagenized population (GABI-Kat) for flanking sequence tag-based reverse genetics. *Plant Mol Biol.*
306 2003;53:247–59.

307 31. Pucker B, Holtgräwe D, Weisshaar B. Consideration of non-canonical splice sites improves gene
308 prediction on the *Arabidopsis thaliana* Niederzenz-1 genome sequence. BMC Res Notes. 2017;10.
309 doi:<https://doi.org/10.1186/s13104-017-2985-y>.

310 32. Olson ND, Lund SP, Colman RE, Foster JT, Sahl JW, Schupp JM, et al. Best practices for evaluating
311 single nucleotide variant calling methods for microbial genomics. Front Genet. 2015;6.
312 doi:10.3389/fgene.2015.00235.

313 33. Ossowski S, Schneeberger K, Clark RM, Lanz C, Warthmann N, Weigel D. Sequencing of natural
314 strains of *Arabidopsis thaliana* with short reads. Genome Res. 2008;18:2024–33.

315

316

317 **Additional Files**

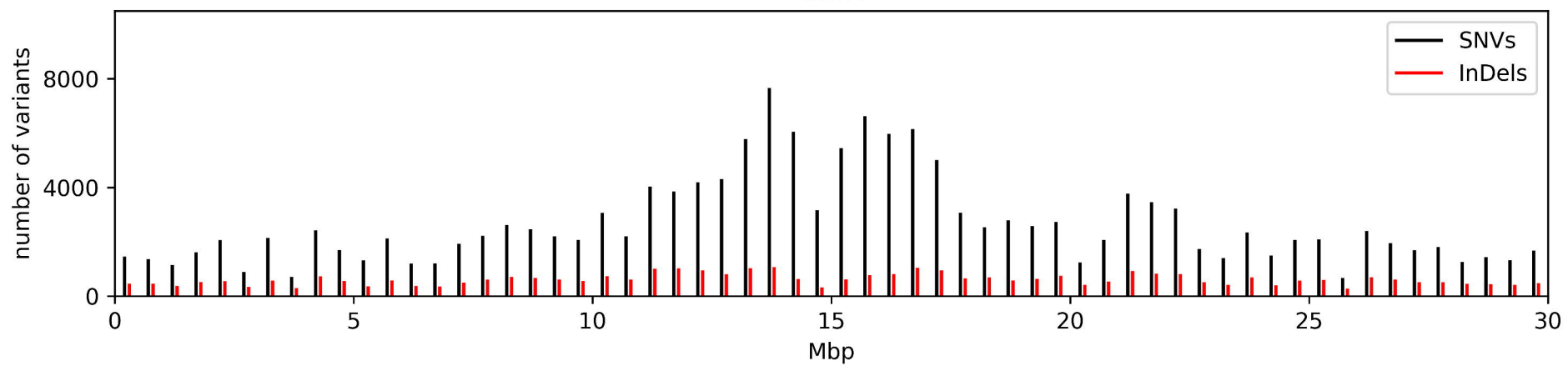
318 AdditionalFile1: Schematic illustration of the variant validation process.

319 AdditionalFile2: Oligonucleotide sequences used for the validation of randomly selected
320 variants.

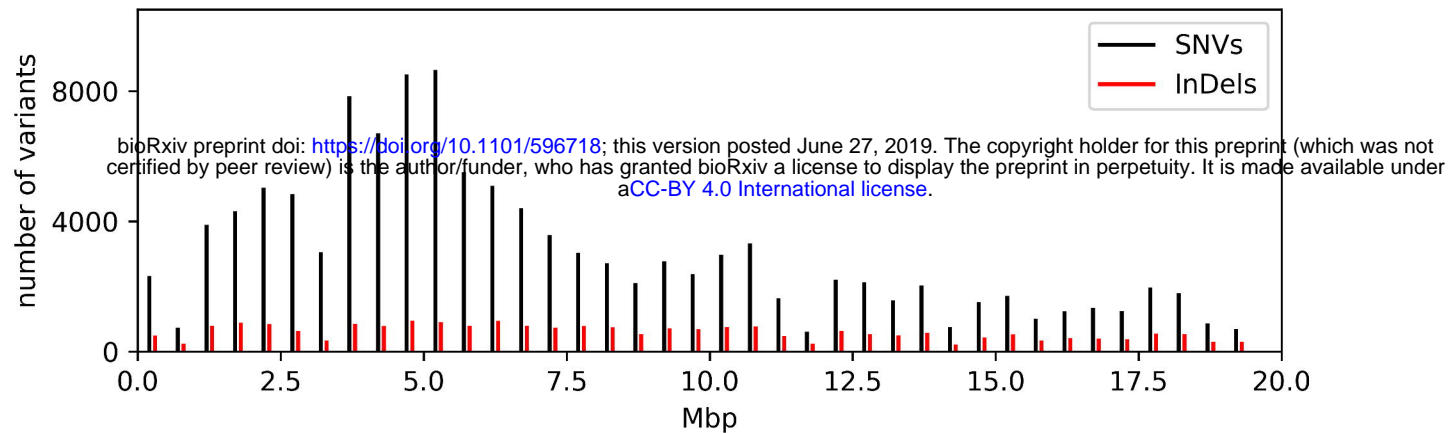
321 AdditionalFile3: Final set of validated variants.

322 AdditionalFile4: NAVIP annotation of variants between Nd-1 and Col-0.

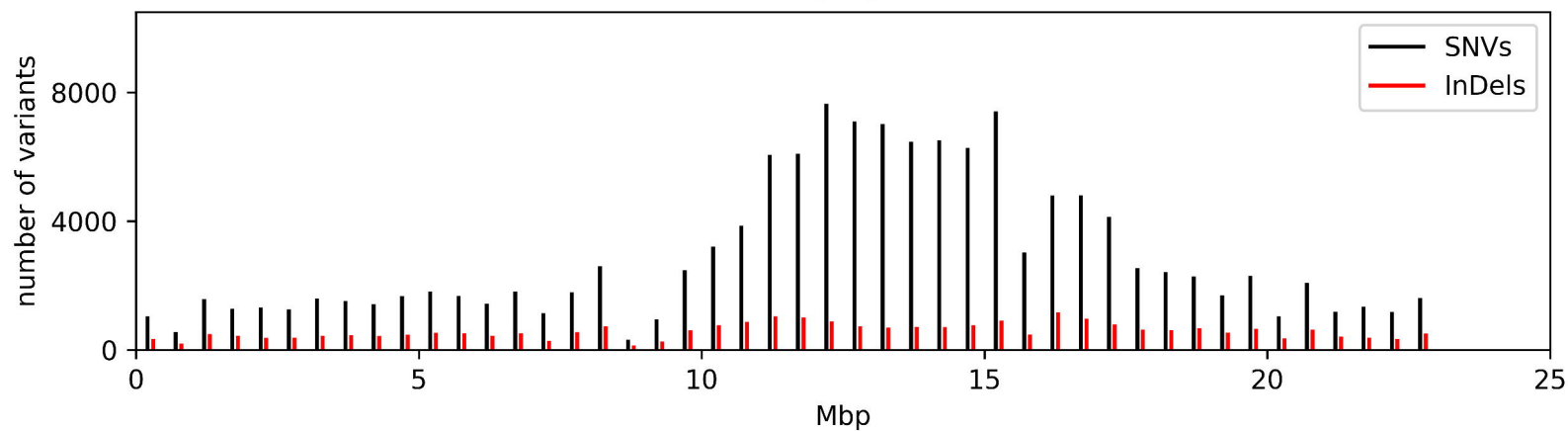
Chr1



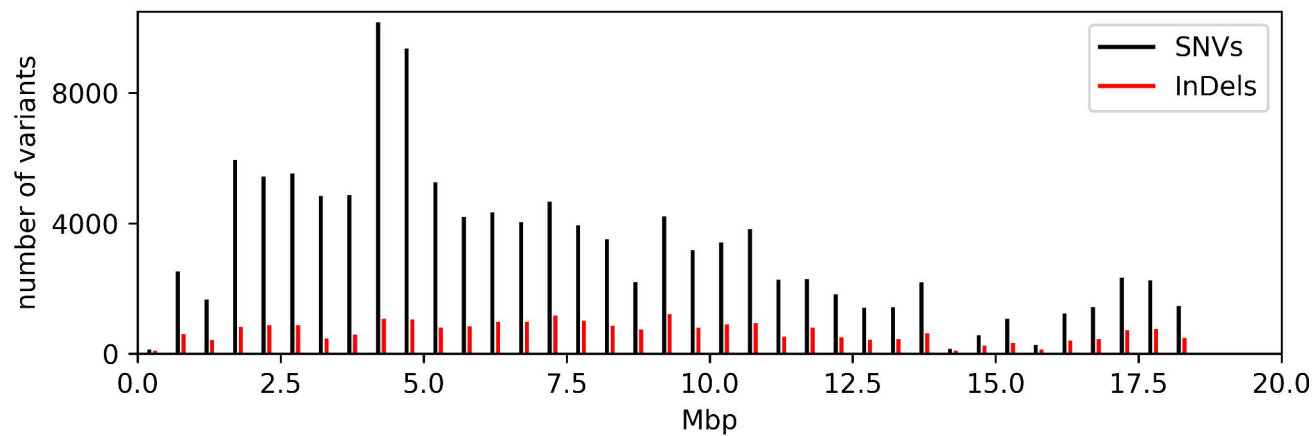
Chr2



Chr3



Chr4



Chr5

