

1 **Accurate identification of de novo genes in plant genomes using** 2 **machine learning algorithms**

3

4 Claudio Casola^{1,2}, Adekola Owoyemi¹, Alan E. Pepper^{2,3} and Thomas R. Ioerger⁴

5

6 ¹ Department of Ecology and Conservation Biology, Texas A&M University, College Station,
7 USA

8 ² Interdisciplinary Doctoral Degree Program in Ecology and Evolutionary Biology, Texas A&M
9 University, College Station, USA

10 ³ Department of Biology, Texas A&M University, College Station, USA

11 ⁴ Department of Computer Science and Engineering, Texas A&M University, College Station,
12 USA

13

14

15

16 *Author for Correspondence: Claudio Casola, Department of Ecology and Conservation

17 Biology, Texas A&M University, College Station, TX, 77843, 979-845-8803, ccasola@tamu.edu

18

19

20

21

22

23

24 **Abstract**

25 De novo gene birth—the evolution of new protein-coding genes from ancestrally noncoding
26 DNA—is increasingly appreciated as an important source of genetic and phenotypic innovation.
27 However, the frequency and overall biological impact of de novo genes (DNGs) remain
28 controversial. Large-scale surveys of de novo genes are critical to address these issues, but DNG
29 identification represents a persistent challenge due to the lack of standardized protocols and the
30 laborious analyses traditionally used to detect DNGs. Here, we introduced novel approaches to
31 identify de novo genes that rely on Machine Learning Algorithms (MLAs) and are poised to
32 accelerate DNG discovery. We specifically investigated if MLAs developed in one species using
33 known DNGs can accurately predict de novo genes in other genomes. To maximize the
34 applicability of these methods across species, we relied only on DNA and protein sequence
35 features that can be easily obtained from annotation data. Using hundreds of published and
36 newly annotated DNGs from three angiosperms, we trained and tested both Decision Tree (DT)
37 and Neural Network (NN) algorithms. Both MLAs showed high levels of accuracy and recall
38 within-genomes. Although accuracies and recall decreased in cross-species analyses, they
39 remained elevated between evolutionary closely related species. A few training features,
40 including presence of a protein domain and coding probability, held most of the MLAs
41 predictive power. In analyses of all genes from a genome, recall was still elevated. Although
42 false positive rates were relatively high, MLA screenings of whole-genome datasets reduced by
43 up to ten-fold the number of genes to be examined by conventional comparative genomic
44 methods. Thus, a combination of MLAs and traditional strategies can significantly accelerate the
45 accurate discovery of DNG and the annotation in angiosperm genomes.

46

47 **Introduction**

48 Novel genes are major drivers of adaptation and evolutionary innovation. A large body of work
49 suggests that new protein-coding genes form at a high rate and represent major contributors to
50 genome evolution and phenotypic variation in plants (1-8). As a result of this rapid evolutionary
51 gene turnover, all species contain hundreds to thousands young, lineage-specific protein-coding
52 sequences that are absent from other taxa (9-12). Both small-scale and whole-genome
53 duplications are responsible for the formation of many novel genes, which tend to share the
54 biological function of their parent genes (1, 3, 4). Occasionally, plants acquired novel genes
55 throughout non-duplicative mechanisms, including horizontal DNA transfer from other species
56 (13, 14) and via the “exaptation” of the coding regions of transposable elements (15-17).
57 Although fundamentally different in nature, all these processes generate new protein-coding
58 sequences from pre-existing genes.

59
60 Evolutionary genomic analyses have revealed that new genes can also emerge from ancestrally
61 noncoding DNA sequences wherein an open reading frame (ORF) and new regulatory sequences
62 originate through ‘enabler substitutions’, typically nucleotide substitutions (18-20). Long
63 considered unlikely evolutionary accidents (Jacob 1977), these so-called de novo genes
64 (hereafter, DNGs) encode for evolutionary novel protein sequences that share no homology with
65 genes from other species. Some DNGs have been shown to bear significant phenotypic impacts,
66 from regulating the mating pathway in budding yeast (21, 22), to producing antifreeze proteins in
67 some fish (23, 24) and affecting human health (25-27).

68

69 First unequivocally discovered in *Drosophila* (28, 29), de novo genes have been identified across
70 several other groups of animals (23, 30-33) and are well characterized in some model organisms,
71 particularly *Saccharomyces cerevisiae* (21, 34-38). In plants, thousands of potential DNGs have
72 been retrieved through computational surveys in *Arabidopsis thaliana*, *Brassica rapa*, poplar,
73 rice, sweet orange and *Triticeae* (8, 10-12, 39-47).

74

75 Despite the growing number of species and genomes analyzed, the current understanding of de
76 novo gene birth and evolution in plants remain severely limited for several reasons. For instance,
77 the identification of DNGs has traditionally relied on comparative genomic strategies that can be
78 computationally demanding, remain difficult to implement on a large scale, and tend to produce
79 many false positives (48). DNG surveys typically include an initial step wherein all genes (or
80 proteins) in a focal species are queried against genes from other species through sequence
81 homology searches, for example using Blast algorithms (49). The pattern of presence/absence of
82 homologous coding sequences of a given gene along a phylogeny of species allows to estimate
83 its approximate time of origin, a procedure known as ‘phylostratigraphy’ (50). Genes that share
84 no homology outside of a given genome according to the phylostratigraphic method should thus
85 be considered species-specific (taxonomically restricted). However, there are significant caveats
86 associated with phylostratigraphy.

87

88 First, homology searches can only generate catalogs of all genes that lack homology, also known
89 as ‘orphan genes’, of which DNGs represent only a subset. Rapidly evolving ancestral genes (51,
90 52), genes derived from exapted transposable elements (11, 53), horizontally transferred genes
91 (53) and genes with alternative coding frame (11, 53) also contribute to the pool of orphan genes.

92 Although some of these processes are thought to occur at much lower rates than de novo gene
93 birth, the proportion of DNGs among orphan genes might be low (54, 55). Discriminating DNGs
94 from other types of orphan genes requires accurate investigation of synteny conservation across
95 species to identify the enabler substitutions that are uniquely associated with de novo gene birth.
96 Substitutions that enable longer ORFs are especially useful but can be observed only by
97 comparing the DNG coding region with the syntenic genomic regions from several other species
98 (18, 19). While critical to the correct identification of DNGs, the search for enabler substitutions
99 is rarely implemented.

100 Second, it has been shown that homology searches can underestimate the age of some types of
101 genes, i.e. rapidly evolving genes (56, 57), which can directly affect estimates of DNG rate of
102 formation (48).

103 Third, homology searches against large datasets require extensive computing resources.
104 Although strategies exist to accelerate this analysis (58), a thorough search of the known
105 sequence space remains challenging and is not reproducible over time, given that sequences
106 databases are expanding exponentially.

107

108 Additionally, a wide spectrum of strategies and bioinformatic protocols have been applied to the
109 search of de novo genes in plants and in other organisms. The combination of these issues can
110 produce significant discrepancies in estimates of DNGs even within the same species. For
111 instance, the number of DNGs discovered in *A. thaliana* ranges from 364 (11) to 782 (41). The
112 lack of standardized accurate approaches to assess the number of DNGs represents a major
113 challenge to estimates the rate of de novo gene birth and is diminishing our ability to characterize
114 the biological impact of DNG across plant species.

115

116 Machine learning algorithms (MLAs) offer a set of approaches with the potential to mitigate the
117 limitations in DNG detection outlined above. MLAs have proven to be powerful methods for
118 learning models from non-linear datasets in a variety of domains, including many applications in
119 genomics and bioinformatics. In fact, MLAs have been developed to annotate genomic features
120 that include protein-coding genes, RNA genes, enhancers, transcription start sites, splice sites
121 and gene function (59-61). To the best of our knowledge, MLAs applied to the detection of de
122 novo genes have not been developed yet. Interestingly, a few studies have explored the ability of
123 MLAs to identify the broader category of orphan genes in plants (62, 63). A variety of MLAs
124 trained and tested on 1,784 orphan genes from *A. thaliana* showed up to 92% accuracy and 95%
125 sensitivity (62), whereas hybrid deep-learning algorithms applied to 1,544 moso bamboo
126 (*Phyllostachys edulis*) orphan genes reached up to 87% balanced accuracy (63). These results
127 suggest that MLAs can achieve high levels of accuracy for orphan gene prediction. However, as
128 discussed above, DNGs represent only a fraction of orphan genes, and are evolutionarily distinct
129 from other types of genes that lack homology across species. Thus, the ability of MLAs to
130 accurately predict DNGs remains untested.

131

132 One of the challenges in applying MLAs to learning classifiers for DNGs is the (typically) small
133 number of positive examples compared to the size of the rest of the genome. While current
134 annotations of plant genomes usually contain tens of thousands of genes, only a few hundred
135 genes can be confidently identified as de novo genes for training. Some MLAs can be sensitive
136 to this asymmetry, outputting models with low information content that appear accurate only
137 because the majority of genes are ancestral genes (hereafter, AGs), while being very inaccurate

138 for DNGs. There are various methods that have been proposed for handling this significant class
139 imbalance (64). We show that sub-sampling of AGs as negative examples can be effective in
140 training accurate models for DNGs. However, although the accuracy on balanced testing sets
141 can be high, even a moderate false positive rate (FPR) can lead to many false positive predictions
142 when the classifier is applied to tens-of-thousands of genes in a whole genome. We show that
143 the FPR can be reduced somewhat by adjusting the selection of examples during training, though
144 at the cost of increasing the false negative rate (FNR). However, false positives can also be
145 removed through traditional comparative genomic analyses that allows to detect signatures of de
146 novo gene birth, i.e. lack of homology in other species and presence of enabler substitutions.

147

148 The ability to detect de novo genes through machine learning classifiers depends on the presence
149 of features that show different distributions of values between DNGs and AGs, defined here as
150 genes with no recent de novo origin. Studies across eukaryotes indicate that DNGs and AGs
151 exhibit different distributions in multiple features associated with gene and protein sequences.
152 For instance, DNGs tend to be shorter and with fewer exons than most AGs (8, 11, 35, 36).
153 Proteins encoded by DNGs typically have fewer annotated domains and possess more
154 structurally disordered regions than ancestral proteins in some eukaryotes (18, 20, 51, 65).

155

156 Leveraging these observations, we sought to develop and test MLAs aimed at discriminating
157 DNGs from ancestral genes using sequence-derived information. MLAs were trained using DNA
158 and protein sequence attributes from DNGs and AGs obtained from three plant species. The de
159 novo gene catalogs consisted of 331 putative DNGs from *Arabidopsis thaliana* (41), 175 and 343
160 DNGs from *Oryza sativa* (8) and 754 novel DNGs from *Brassica rapa*. These species represent

161 evolutionary lineages with different levels of divergence, as *A. thaliana* and *B. rapa* are
162 relatively closely related species belonging to the Brassicaceae family within the dicotyledon
163 (dicots) clade, whereas rice is a much more evolutionary distant species in the monocotyledon
164 (monocots) clade.

165

166 Our investigation had the following goals: (1) Assessing the accuracy and recall of different
167 MLA approaches, including decision trees and neural networks, in discriminating DNGs and
168 AGs; (2) Identifying DNA and protein features with high predictive power for detecting DNGs;
169 (3) Assessing the accuracy and recall of MLAs based exclusively on sequence features compared
170 to those incorporating both sequence features and functional genomic data (gene expression
171 levels, translation level, protein-protein interactions, etc.); (4) Determining the predictive ability
172 of MLAs built on data from one species across other taxa; (5) Determining the predictive ability
173 of MLA approaches in detecting DNGs using whole-genome sequence data that include all genes
174 in a species/accession.

175

176

177 **Results and Discussion**

178

179 **A set of high confidence *A. thaliana*-specific DNGs**

180 A dataset 782 of putative *A. thaliana*-specific DNGs was recently generated by Li et al. (41).
181 These genes were identified using sequence homology searches on a limited set of databases and
182 without validating the de novo status of each gene through synteny conservation with closely
183 related species. To produce a set of high confidence DNGs from the Li et al. (2016) catalogue,

184 we conducted additional homology searches and retained only genes that passed a series of
185 stringent criteria, including conserved synteny with other Brassicaceae and lack of homology
186 with genomes and transcriptomes deposited on NCBI, as described in **Methods**. This approach
187 follows a robust computational framework developed to identify DNGs (58). A total of 298
188 putative DNGs were excluded as they shared homology with genes in other species (**S1 Table**),
189 leaving 331 high confidence *A. thaliana*-specific DNGs, similarly to the number of DNGs
190 reported by Donoghue et al. (11); however, we could not directly compare our list of DNGs to
191 those identified in this paper, as de novo gene IDs were not provided by Donoghue et al. (11).
192 The remaining 153 putative *A. thaliana* DNGs shared no conserved synteny with other species
193 and were removed from the catalog as the de novo birth pathway could not be determined.

194

195 **Identification of high confidence DNGs in *Brassica***

196 We first analyzed all available *Brassica* genome assemblies to determine gene annotation
197 completeness. Most assemblies showed a high level of completeness according to BUSCO
198 scores (**S3 Table**). We selected *B. rapa* as the primary focal species for a de novo gene survey
199 due to its agricultural importance and due to the extensive genomic and functional data available
200 for this species (66-71). As the main focal genome to identify de novo genes, we selected the
201 recently improved *B. rapa* v3.0 genome assembly from the Chiifu-401-42 genotype (71), which
202 contained 46,248 protein coding annotated genes and showed a slightly higher annotation
203 completeness than other available assemblies. We performed extensive sequence similarity
204 searches to identify putative *Brassica*-specific genes (see **Methods**) and identified 754 candidate
205 *B. rapa* DNGs.

206

207 **Sequence and structural features of de novo genes, ancestral genes and their proteins**

208 On average, de novo gene and protein sequences are known to diverge from ancestral genes in
 209 several features. To determine if these qualities can be used to discriminate between DNGs and
 210 AGs, we compiled 22 sequence features from the 331 DNGs from *A. thaliana*, the 754 DNGs
 211 from *B. rapa*, and the recently published catalogs of 175 and 343 DNGs from rice, *Oryza sativa*
 212 (8) (**S1-3 Datasets**). These sequence features are straightforward to calculate for any annotated
 213 genome, and do not require any additional data collection (such as gene expression by RNAseq
 214 or ribosome profiling). The rice DNGs were obtained integrating sequence homology searches
 215 with synteny analyses and were therefore considered well-curated datasets comparable to those
 216 of *A. thaliana* and *B. rapa*. To the best of our knowledge, these datasets represent the largest
 217 comparative catalog of plant de novo genes to date. The same features were also retrieved by all
 218 AGs in these three species (**Table 1**).

219

220

221 **Table 1. Primary sequence features and mean values of de novo and ancestral genes and**
 222 **proteins.**

	AT DNGs	AT AGs	BR DNGs	BR AGs	OS DNGs-1	OS DNGs-2	OS AGs
# Genes	331	26,423	754	34,354	175	343	38,405
Gene length (bp)	376	2,269	630	1,988	4,612	3,788	3,729
CDS length (bp)	218	1,251	350	1,169	418	365	1157
CDS #exons	1.4	5.3	1.9	5.4	3.0	2.7	4.7
GC-content	43	44	49	47	58	58	58
SSRs in CDS (bp)	22	18	16	11	16	10	16
%SSRs in CDS	4.1	1.3	4.2	0.9	2.8	1.9	1.8
Gene overlap	0.0	0.1	0.5	0.1	13.1	8.7	7.9
CAI	0.76	0.77	0.77	0.79	0.79	0.79	0.81
%Coding	31	98	39	92	44	38	84

Coding probability	0.20	0.93	0.49	0.91	0.50	0.43	0.85
Protein domains	0.02	0.85	0.02	1.15	0	0.01	0.95
%ISD	0.28	0.31	0.33	0.30	0.46	0.44	0.34
pI	8.7	7.5	8.6	7.5	9.7	9.6	8.0
#TM helices	0.15	0.78	0.09	0.60	0.06	0.05	0.51
%w/TM helices	13	25	7	15	5	4	14
%w/SP	15	15	7	11	5	5	12
%w/mTP	5	4	3	3	1	1	2
%w/cTP	1	6	2	5	2	2	5
%w/luTP	0	0.5	0	0	0	0	0
%w/NoTP	78	74	88	80	92	92	81
Kozak score	-5.44	-4.99	-5.40	-5.06	-5.15	-5.37	-4.93

223 AT: *A. thaliana*.

224 BR: *B. rapa*.

225 OS: *O. sativa*.

226 CDS: coding sequence.

227 SSRs: simple sequence repeats (microsatellites).

228 CAI: codon adaptation index.

229 Coding probability according to CPC2 (72).

230 %Coding: proportion of genes predicted by CPC2 to be coding.

231 %ISD: proportion of protein sequence with intrinsic structural disorder.

232 pI: Isoelectric point.

233 #TM helices: average of transmembrane helices per gene.

234 %w/TM helices: proportion of genes with transmembrane helices.

235 %w/SP: proportion of genes with signal peptide.

236 %w/mTP: proportion of genes with Mitochondrial transit peptide.

237 %w/cTP: proportion of genes with Chloroplast transit peptide.

238 %w/luTP: proportion of genes with Luminal transit peptide.

239 %w/NoTP: proportion of genes with No-targeting peptide.

240

241

242 Along the lines of previous studies (8, 41), all DNGs components tend to be shorter (except for
 243 introns in the *O. sativa* DNG sets), especially at the level of the coding region, and contain fewer
 244 exons compared to AGs (**Table 1; S2A-C Fig**). The GC-content varied more significantly across
 245 species than between DNGs and AGs (**Table 1; S2D Fig**). Interestingly, the GC-content
 246 distribution peaks at higher values for AGs in *A. thaliana* and for DNGs in *B. rapa* (**S2D Fig**). In
 247 rice, the GC-content distribution is bimodal in AG, as previously described (73), with DNG GC
 248 values peaking in between (**S2D Fig**). This pattern in the rice genome is mirrored at the level of

249 the codon adaptation index (**S2G Fig**). Additionally, DNG coding regions contain a higher
250 proportion of DNA derived from microsatellites identified by Tandem Repeat Finder.
251 Interestingly, the microsatellite content is lower in the large rice DNG dataset, which includes
252 older DNGs, suggesting that simple repeat content may decrease with time due to substitutions.
253 As expected, the predicted coding potential of DNGs is much lower than in AGs, with much
254 fewer DNGs being labeled ‘coding’ according to the coding potential calculator 2. In agreement
255 with this, the coding adaptation index of DNGs is on average significantly lower than in AGs
256 (**Table 1; S2E-G Fig**). Furthermore, Kozak scores were lower in DNGs compared to AGs,
257 although their distributions largely overlapped between the two types of genes (**Table 1; S2H**
258 **Fig**).

259
260 Protein features included the presence of conserved domains, predicted proportion of intrinsic
261 structural disorder (ISD) residues, isoelectric point, predicted number of transmembrane helices
262 and subcellular localization peptides (**Table 1; S2I-K Fig**). As expected, given the small size of
263 de novo proteins and their recent origin, very few of them contained conserved domains.
264 Similarly, we found a much lower number of transmembrane motifs (TMs) and genes with TMs
265 in de novo proteins than ancestral proteins. This is in contrast with the observation that *S.*
266 *cerevisiae* de novo ORFs with adaptive potential are enriched for TM domains (74).

267
268 Across all species, DNG proteins showed consistently higher isoelectric point (pI), indicating a
269 higher proportion of basic residues, and very few transmembrane domains compared to AG
270 proteins (**Table 1**). Elevated pIs due to a depletion of acid residues have been also found in
271 mammalian orphan proteins (75) and in *S. cerevisiae*-specific translated ORFs (37), but no

272 explanations for these trends have been put forward. We also found that the distribution of pI
273 values is very different between proteins encoded by DNGs and AGs. While in AG proteins the
274 pI values follow an approximate trimodal distribution, isoelectric points in DNG proteins cluster
275 around two peaks around low (~4) and high (~11) values (**S2I Fig**).

276

277 Additionally, we observed higher ISD levels in DNGs than AGs only in *B. rapa* and rice (**Table**
278 **1; S2K Fig**). This is in agreement with what observed in *Drosophila melanogaster* DNGs (76)
279 and in orphan genes in *D. melanogaster* (77) and *Leishmania* (78). A similar trend was also
280 reported in rodents young genes (65), although follow up studies suggest that this pattern is
281 likely to be an artifact (51, 79). Some authors have shown that high ISD levels are associated
282 with elevated GC content in orphan genes or young genes (80), possibly a result of some young
283 genes overlapping with ancestral genes (51). It is unclear if the modest difference in GC content
284 between DNGs and AGs in *B. rapa* and rice drives their elevated ISD levels.

285

286 Overall, DNG proteins also contained fewer localization peptides compared to AG proteins, with
287 the notable exception in *A. thaliana* of both a higher proportion of mitochondrial transit peptides
288 in DNGs vs. AGs, and a comparable number of DNG and AG proteins with a signal peptide
289 (**Table 1**).

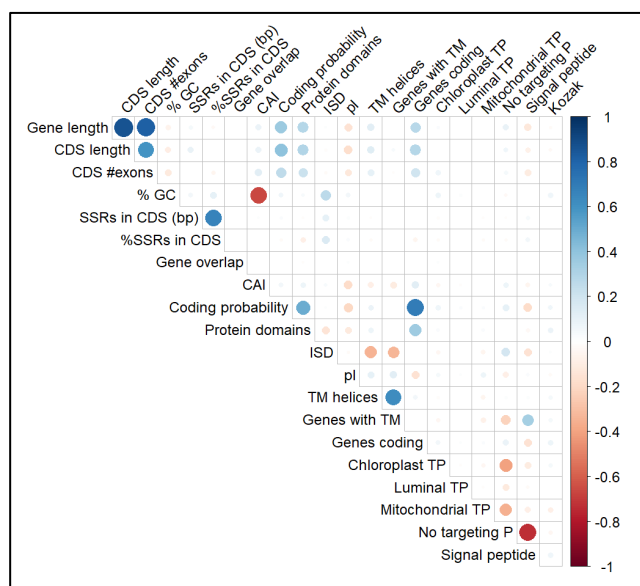
290

291 We further investigated possible correlation among features in each species (**Fig 1**). As expected,
292 gene length, CDS length and CDS #exons (the number of coding exons) were positively
293 correlated. Similarly, the coding probability and the presence of protein domains increased with
294 gene length and CDS length and, to a smaller extent, with CDS #exons. Longer coding regions

295 are more likely to be predicted as coding by the coding potential calculator (72). Longer genes
 296 are also more likely to encode protein domains. In agreement with previous studies (51, 80), we
 297 found that GC-content and ISD were positively correlated. Another expected pattern is the
 298 negative correlation between ISD and the presence of transmembrane helices, which cannot form
 299 in disordered protein regions. The anticorrelation between GC-content and codon adaption index
 300 may be due to the paucity of GC-rich codons among the most-used codons in angiosperms.
 301 Similar correlations were observed in *B. rapa* and rice (S3 Fig).

302

303



312 **Fig 1.** Correlation matrix of *A. thaliana* sequence features used to train ML classifiers. Feature
 313 names follow the nomenclature of Table 1.

314

315

316

317 **Both Decision Tree and Neural Network classifiers detect the majority of de novo genes in**
318 **test gene sets**

319 As a preliminary attempt to develop a predictive model for DNGs, we trained both a Decision
320 Tree (DT) classifier and a Neural Network (NN) classifier on the *A. thaliana* gene datasets. The
321 set of positive examples was formed by the 331 DNGs we identified in this species. Because
322 using all 26,423 AGs as negative examples led to degenerate tree where every gene was
323 classified as negative due to class imbalance (64), the AGs were sub-sampled (81) by choosing
324 an equal number of 331 genes at random for the negative examples. The NN classifier consisted
325 of a fully-connected network with a hidden layer of 20 units (see **Methods**). We also evaluated
326 the effect of a different number of hidden units (from 10 to 100 in intervals of 10 units) and 2
327 hidden layers, but these did not significantly improve the accuracy of the NN. When the 662
328 selected examples were divided randomly into 70% for training and 30% for testing, the DT and
329 NN models were found to have 91-92.0% accuracy in predicting DNGs in the test set. Thus,
330 even though Decision Trees and Neural Networks represent completely different methods for
331 capturing patterns in training data, they are both able to learn how to discriminate DNGs from
332 AGs using sequenced-based features. The confusion matrix for the 100 genes from the 30% test
333 sets shows that the errors are evenly distributed between false positives and false negatives, with
334 recall values above 90% (**Table 2**).

335

336

337

338

339 **Table 2. Confusion matrix of DT and NN classifiers, trained on *A. thaliana* genes, applied**
340 **to an independent balanced test set. Values represent counts of test examples.**

		predicted class labels			
		DNG	AG	Recall	
DT classifier	actual class labels	DNG	91	9	91%
		AG	7	93	
NN classifier	actual class labels	DNG	92	8	92%
		AG	8	92	

341

342

343

344 In order to determine a more general estimate of the predictive accuracy (since the single tree
345 above is dependent on the specific AGs chosen as negative examples), 10-fold cross-validation
346 was carried out (where a decision tree was generated based on 90% of the data and tested on the
347 remaining 10%, repeated 10 times in a rotated manner), which resulted in a performance estimate
348 of 89.7%, with a 95%-confidence interval (CI95) of 87.3-92.1% (**Table 3**). The same protocol
349 was used to train classifiers on species-specific sequence features retrieved in DNGs and AGs of
350 *B. rapa* and *O. sativa* (using the larger dataset of 343 de novo genes in the latter, see **Table 1**). In
351 all species, NN models performed slightly better than DT models, significantly so in *A. thaliana*
352 and *B. rapa* but not in rice (**Table 3**). DT and NN classifiers showed significantly lower
353 accuracies in *B. rapa* and *O. sativa* compared to *A. thaliana* (**Table 3**; DT *A. thaliana*-*B. rapa*
354 $P=0.0038$; DT *B. rapa*-*O. sativa* $P=0.0003$; NN *A. thaliana*-*B. rapa* $P=0.0005$; NN *B. rapa*-*O.*
355 *sativa* $P=0.0001$, unpaired T-test). Overall, these results indicate that MLAs trained on species-
356 specific datasets can successfully retrieve the vast majority of DNGs. Confusion matrices for
357 both classifiers indicate that NN models achieve substantially higher recall than DT models in all

358 species (**S4 Table**). Recall is comparably high (~92%) in NN models of Brassicaceae, but
359 decreased to ~83% in rice, where DT models achieved only ~76% of recall (**S4 Table**).

360

361 **Table 3. Accuracy of DT and NN classifiers in the three angiosperms.**

Species	Decision Tree balanced accuracy [†] (95% C.I.)	Neural Network balanced accuracy [†] (95% C.I.)
<i>A. thaliana</i>	89.7% (87.3-92.1%)	93.2% (91.5-94.8)*
<i>B. rapa</i>	85.1% (84.3-86.0%)	88.9% (88.1-89.8)*
<i>O. sativa</i>	76.5% (73.1-80.0%)	80.2% (77.0-83.3)

362 [†]Averaged over 10-fold cross-validation

363 *NN model vs. DT model within species, *P*-value<0.05, unpaired T-test

364

365

366 The variation in accuracy and recall across species may be due to several factors. A higher
367 quality of gene annotation in *A. thaliana* may explain the increased accuracy of MLAs in this
368 species. The lower accuracy in rice could in part depend on the slightly older age of the larger
369 dataset of 343 DNGs used for this species (8), as DNGs should acquire features that are more of
370 typical genes through time (35). Differences in age of DNGs could also explain the significant
371 overlap in the distribution of continuous features between DNGs and AGs (**S2 Fig**). We also
372 observed that some features show a varying degree of predictive importance between
373 Brassicaceae and rice, which could further contribute to differences in accuracy (see below).

374

375

376 **Training DT classifiers with functional genomic features**

377 We evaluated whether adding functional genomic data could improve the accuracy of the
378 decision tree classifier using datasets available for *A. thaliana*. The classifier was re-trained
379 using 28 additional features which are not systematically available in all genomes, including

380 transcription data (RNAseq), translation levels estimated through Ribosomal profiling
381 (RiboSeq), proximity to transposable elements, selective constraint, and phenotype data for gene
382 knockout mutants (**S5 Table**). The 10-fold cross-validated accuracy of models extended with
383 these functional features was 91.4%, (89.7-93.1%), which is not significantly greater than
384 models without these functional features ($P>0.05$, unpaired T-test). Some of the functional
385 features were occasionally used as decision criteria in lower branches of some of the decision
386 trees; the functional feature with highest importance (0.04) was “AVG RiboP RPKM 25
387 samples”, which suggests that lack of expression evidence can be an important discriminator for
388 DNGs.

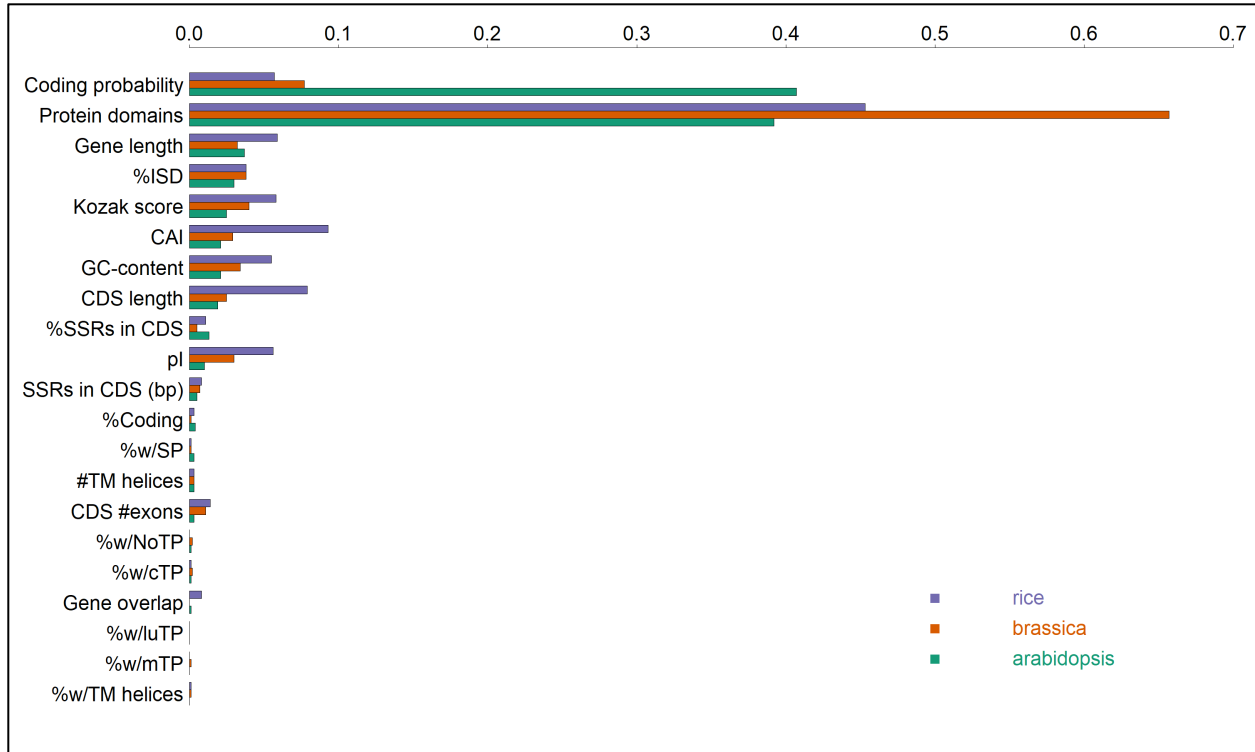
389

390

391 **Feature Importance in Decision Trees**

392 To better assess the contribution of each feature to DT classifiers in the three species, we
393 calculated feature importances, which are based on the relative contribution of each feature to
394 splitting the data in the tree and range between 0 and 1 (see **Methods** and **Fig 2**). Feature
395 importances averaged over 10 runs for each species are shown in **S6 Table**. The DT classifier
396 developed in *A. thaliana* consisted of 70 nodes (including 24 leaves) with a depth of 11. The
397 attribute tested for splitting at the root of the tree, considered the most important based on
398 reduction of Gini impurity, was “Protein domains”. This turns out to be a highly discriminating
399 feature: most AGs (22458/26423=85.0%) have at least one domain (recognized as a known fold
400 family by Pfam based on amino acid sequence), whereas most DNGs do not (only 7/331=2.1%
401 had a recognized domain) (**Fig 2**).

402



403

404 **Fig 2.** Top ten features in the DT classifier ranked by importance.

405

406

407 Overall, the top ten features by importance are largely the same across all species. “Protein
408 domains” is the most prominent feature across species (**Fig 2; S6 Table**). Several other top
409 features that are known to be significantly different between DNGs and AGs, including “Coding
410 probability”, “Gene length”, “CDS length”, “Codon adaptation index (CAI)”, “%GC” and the
411 proportion of “intrinsic structurally disordered (ISD)” regions in proteins, show high importance.
412 Although the presence of a conserved Kozak motif has not been investigated before in DNGs,
413 the “Kozak score” feature showed a relatively high importance, in agreement with other findings
414 suggesting that de novo genes might acquire more ‘gene-like’ regulatory sequences by natural
415 selection after their emergence (31).

416

417 Furthermore, some features exhibited substantially higher importance in rice than Brassicaceae.
418 For instance, “Codon adaptation index (CAI)” represents the second most important feature in
419 rice while ranking sixth and eight in *A. thaliana* and *B. rapa*, respectively (**Fig 2; S6 Table**).
420 This is interesting as CAI is only slightly higher in AGs than DNGs in rice (**Table 1; S2G Fig**).
421 “Gene length”, “CDS length”, “Kozak score” are also more prominent features in the monocot.
422
423 Interestingly, the top 2 features alone in *A. thaliana* (“Coding probability” and “Protein
424 domains”) can be used to construct decision trees with nearly equivalent accuracy of the ones
425 trained on all features. For example, in *A. thaliana*, the balanced accuracy with such trees
426 (averaged over 10-fold CV) is 91.4% (95%CI: 89.1-93.7%). The decision trees still have
427 multiple nodes in them, typically around 30; they just include splits on multiple threshold values
428 (i.e. sub-ranges) of “Coding probability”. The average feature importances are 0.407 for
429 “Coding probability” and 0.392 for “Protein domains”. In fact, the neural network performs
430 even better, probably due to the reduction in parameters (weights in the network) with just two
431 inputs: 94.0% (95%CI: 92.2-95.7%). Further analyses will be necessary to determine if this
432 applies to other plant genomes.

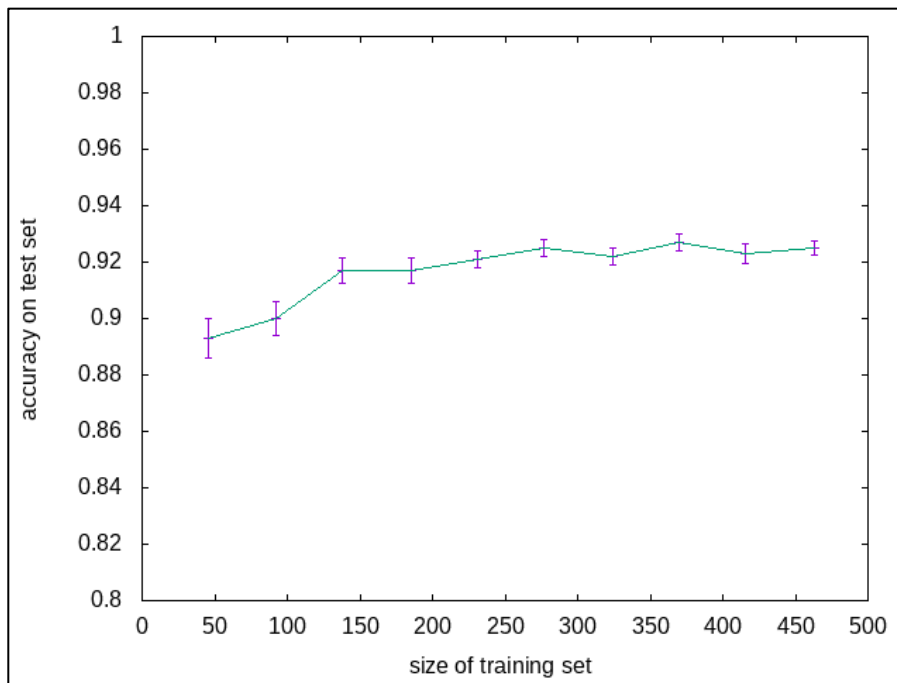
433

434

435 **Accuracy is not limited by small training set size**

436 The size of the training dataset can affect accuracy in MLA predictions. We tested if this is the
437 case for the DT classifiers using randomly selected subsets of *A. thaliana* de novo genes and
438 ancestral genes equal to 10-100% of the original training set of 464 genes. This was repeated 30
439 times for each set size to obtain a DT classifier learning curve. The testing dataset for each

440 analysis was carried out on a group of DNGs that did not overlap with the training set.
441 Additionally, during training and testing, the number of negative examples (AGs) was always
442 balanced with an equal number of positive examples. We found that the accuracy of the DT
443 remains elevated (>85%) even for training with one-tenth of the full training-set size of 464
444 genes and is nearly equal to the highest accuracy with only 30% of the full training size,
445 corresponding to 138 genes (**Fig 3**). The observed trend also suggests that the accuracy would
446 not be significantly improved by using a larger training set with more known DNGs.
447
448



449
450 **Fig 3.** Learning curve of the *A. thaliana* DT classifier on 10-100% of the original training
451 dataset. Standard errors are shown.

452
453
454

455 **Cross-species models for de novo gene prediction are nearly as accurate as species-specific**
456 **models**

457 An important question is whether the patterns extracted by the MLAs for discriminating DNGs
458 are species-specific, or whether the MLAs are capturing general properties of DNGs that extend
459 across multiple plant species. To assess this, we trained a DT and NN classifiers using data on
460 genes from *A. thaliana*, and then applied these models to the *B. rapa* and *O. sativa* datasets. In
461 10-fold cross-validation analyses, the *A. thaliana* DT model achieved 83.1% and 67.2% accuracy
462 in *B. rapa* and rice, respectively (**Table 4**). The *A. thaliana* NN model reached slightly higher
463 accuracy in both species (**Table 4**). Overall, species-specific models (**Table 3**) achieved
464 significantly higher accuracy than cross-species models (P -value <0.05 , unpaired T-test), except
465 for the *A. thaliana* DT model applied to *B. rapa* datasets (P -value=0.0855, unpaired T-test).
466 Taken together, these results indicate that models trained on the *A. thaliana* genome, which has
467 been more carefully and thoroughly annotated, can be applied with nearly equal accuracy on *B.*
468 *rapa*, and does not require an MLA to be re-trained on each new gene set. Conversely, *A.*
469 *thaliana* classifiers achieved substantially lower accuracies in *O. sativa* compared to the rice-
470 specific models. It appears that the features associated with DNGs in *A. thaliana*, such as coding
471 probability, lacking recognizable protein domains, and having lower Kozak score, generalize
472 across species, and are associated with DNGs in other plant genomes.

473

474

475

476

477 **Table 4. Comparison of performance of DT and NN models using a common model trained**
478 **on *A. thaliana* versus species-specific models.**

Species	Decision Tree balanced accuracy [†] (95% C.I.)	Neural Network balanced accuracy [†] (95% C.I.)
<i>B. rapa</i>	83.1% (81.2-85.0)	85.1% (83.5-86.6)
<i>O. sativa</i>	67.2% (63.9-70.5)	69.4% (66.5-72.2)

479 [†]Averaged over 10-fold cross-validation

480

481

482 The higher predictive ability of *A. thaliana* MLA classifiers in *B. rapa* compared to rice suggests
483 that cross-species DNG identification with MLAs tend to be more accurate in closely related
484 genomes. Features that differ significantly between Brassicaceae and rice genes, including gene
485 length, %GC and simple repeat content in the coding region, may drive the lower sensitivity of
486 the *A. thaliana* DT classifier in rice. A broader taxonomic sampling at varying phylogenetic
487 distances from *A. thaliana* will be required to test this hypothesis more thoroughly. We noticed
488 that the recall decreased in cross-species prediction, with a limited difference in *B. rapa* (from
489 ~84-92% to ~82-84%) and a significant drop in in rice (from ~76-83% to ~55-59%) (**S7 Table**).
490 This indicates that cross-species MLA predictions of DNGs may achieve acceptable levels of
491 sensitivity within taxonomic families but fail to detect a substantial fraction of DNGs in more
492 distantly related species. Thus, broad MLA-based de novo gene surveys in plants may require the
493 training of classifiers using DNGs detected with comparative genomic approaches in at least one
494 species per family.

495

496

497

498 **Whole-genome predictions of DNGs using DT and NN classifiers**

499 We next assessed the accuracy of species-specific DT and NN classifiers to predict DNGs in
500 whole-genome gene sets. We calculated the balanced accuracy, which is better suited to assess
501 the performance of classifiers when classes are imbalanced. The overall balanced accuracy
502 ranged from 92.3% in *A. thaliana* to 76.9% in *O. sativa*, with slightly higher accuracy in DT vs.
503 NN models (**Table 5**). A high recall of ~94-99% was found across species and classifiers,
504 although DT models also achieved lower FPRs compared to NN models (**S8 Table**). Given the
505 much higher number of AGs than DNGs, the total number of false positives reached ~2,055 in *A.*
506 *thaliana* and a maximum of ~8,948 genes in *O. sativa* (**S8 Table**). Given the high FPRs, MLAs
507 alone may achieve the level of accuracy required to entirely replace traditional comparative
508 genomic analyses in DNG surveys; however, the application of MLAs as a first step would
509 decrease by up to 10-fold the number of genes that need to be investigated with homology
510 searches and other time-consuming approaches in order to remove false positives.

511

512

513 **Table 5. Comparison of performance of DT and NN species-specific**
514 **models applied to whole genomes.**

species	Decision Tree genome-wide accuracy [†] (95% C.I.)	Neural Network genome-wide accuracy [†] (95% C.I.)
<i>A. thaliana</i>	92.3% (91.9-92.7)	91.1% (90.4-91.9)
<i>B. rapa</i>	87.4% (87.1-87.7)	86.0% (85.6-86.4)
<i>O. sativa</i>	79.3% (78.7-79.9)	76.9% (74.9-78.8)

515 [†]Averaged over 10-fold cross-validation

516

517

518 We further examined if the performance of models trained on *A. thaliana* data but applied to the
519 two other species (in this analysis, there are no confidence intervals because a single input model
520 trained on one species was tested for accuracy on the whole genome of another). The *A. thaliana*
521 classifiers, particularly NN models, showed relatively high balanced accuracy in *B. rapa* and in
522 rice (**Table 6**). However, recall values dropped significantly in both species, from ~97-98% to
523 ~82-84% in *B. rapa* and from ~94-98% to ~55-59% in rice (**S9 Table**). Interestingly, the *A.*
524 *thaliana* NN model resulted in lower FPRs than the within-species models while the opposite
525 was true for the *A. thaliana* DT model (**S9 Table**).

526

527

528 **Table 6. Comparison of performance of DT and NN models using a common model trained**
529 **on *A. thaliana* versus species-specific models applied to whole genomes.**

species	Decision Tree genome-wide accuracy	Neural Network genome-wide accuracy
<i>B. rapa</i>	84.9%	86.2%
<i>O. sativa</i>	74.6%	81.6%

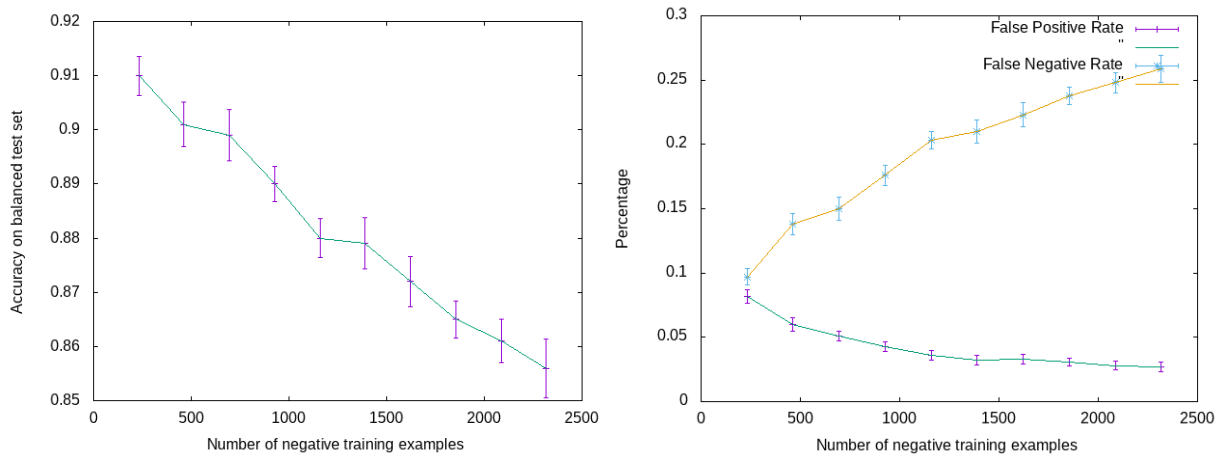
530

531

532 **The DT classifier specificity can be adjusted by increasing the proportion of negative**
533 **examples during training**

534 Given the high false positive rates in classifiers, we investigated if including negative examples
535 during training could increase the model specificity using *A. thaliana* datasets and the DT model.
536 We scaled up the number of negative examples up to 10 times the original set of 232 AGs while
537 maintaining a balanced test set with equal numbers of DNGs and AGs. Each iteration at different
538 training sizes was repeated 30 times. We found that for increasing numbers of negative
539 examples, the balanced accuracy steadily decreases from ~91% to ~85%, primarily due to the

540 increasing false negative rate from <10% up to 25% (**Fig 4**). Concomitantly, the false positive
541 rate decreased from ~8% to <4% (**Fig 4**). As the main goal of the MLA approach is to detect
542 DNGs, the loss of sensitivity associated with the higher number of negative examples might be
543 not worthwhile. However, we noticed that this tradeoff between false positive and false negative
544 rates might be acceptable for limited increases in the size of negative examples.
545



546 **Fig 4.** Accuracy (left) and false positive and false negative rates (right) for increasing number of
547 negative examples used in the training of the DT classifier in *A. thaliana*. Standard errors are
548 shown.
549

550

551

552

553 **Conclusions**

554 In this study, we have developed and assessed the first machine learning framework to identify
555 de novo genes. In order to make these approaches readily applicable in species with limited
556 functional genomic data, we have specifically selected basic sequence features that can be
557 obtained from DNA and protein sequences available in annotated genomes. Using DNG datasets

558 from three plant species, including an updated gene set from *A. thaliana* and the first group of de
559 novo genes in *B. rapa*, we have found that both decision tree (DT) and neural network (NN)
560 classifiers achieve high levels of accuracy and recall in predicting DNGs. Using DT algorithms
561 applied to sub-sampled sets of DNGs and AGs, we identified a few features with significant
562 predictive power for DNGs. This is in line with performance ability of MLAs to discover orphan
563 genes in *A. thaliana* based on six DNA sequence features (62). Importantly, orphan genes are not
564 equivalent to de novo genes, as the former appear to be mostly constituted by rapidly evolving
565 genes (54). Training MLA models with additional features derived from functional genomic data
566 (transcription and translation data) and information from phenotypic assays that are not readily
567 available in most sequenced genomes does not lead to a substantial increase in accuracy or
568 reduction in the number false positives.

569
570 A major advantage offered by MLAs is the significant decrease in computational time compared
571 to traditional genomic approaches to find DNGs—essentially, a time contraction from weeks or
572 months to minutes. For MLAs to be successfully applied in DNG surveys across hundreds to
573 thousands of species, it is critical to train models using datasets of known DNGs from a few key
574 species, and for these models to obtain high accuracy and recall in other species. We conducted
575 initial tests to explore this possibility using DT and NN models trained on *A. thaliana* to predict
576 DNGs in *B. rapa* and rice. We found that these cross-species predictions achieved comparable
577 accuracy of species-specific models in *B. rapa*, and somewhat lower accuracy in rice. This
578 suggests that MLAs trained in one species can likely be used to infer DNGs in closely related
579 species, as in the case of the two Brassicaceae, *A. thaliana* and *B. rapa*. Given that many
580 angiosperm families now contain sequenced genomes from multiple species, and considering

581 both the rapid increased of the number and quality of new genome assemblies, de novo gene
582 discovery based on MLAs could likely be applied to a large number of flowering plant taxa.
583 Future work on more taxa should help better determine how cross-species MLAs accuracy
584 decreases when the evolutionary distance between taxa increases, as this study indicates.
585
586 Genome-wide analyses showed that species-specific models predicted well above 90% of known
587 DNGs, although the much higher number of ancestral genes lead to several thousand false
588 positive cases. In cross-species genome-wide analyses, *A. thaliana* models identified 82-84% of
589 true DNGs in *B. rapa*, but achieved less than 60% recall in rice. Notably, FPRs were
590 substantially lower in cross-species NN models compared to species-specific models. The
591 combination of high recall, at least in some cross-species tests, and high FPRs suggests that a
592 three-step pathway can be employed to accelerate DNG discovery in angiosperms. First, NN
593 models are developed in one or two species with the best gene annotation quality in each
594 angiosperm family. Second, these NN models are applied to an array of target species in the
595 same family. Third, the candidate DNGs predicted from each target species, comprising only a
596 few thousand genes, are analyzed post-hoc with traditional comparative genomic approaches to
597 remove false positives.
598
599 As this represents the first systematic study to assess machine learning approaches in de novo
600 gene discovery, we expect that further developments of in this area could significantly increase
601 accuracy and recall while reducing the false positive rate in DNG detection. Along these lines,
602 alternative machine learning approaches, including deep neural networks, and methods to

603 address the class imbalance between DNGs and AGs different from sub-sampling, such as
604 synthetic minority over-sampling algorithms, or SMOTE (62, 82), warrant future investigations.

605

606

607

608 **Methods**

609 **Validation of *Arabidopsis thaliana* de novo genes**

610 The TAIRv10 (TAIR10) DNA and protein sequences of *A. thaliana* were obtained from the
611 folder “TAIR10 blastsets” in the TAIR repository (83). Data from the files
612 “TAIR10_pep_20101214”, “TAIR10_cds_20101214” and “TAIR10_exon_20101028” were
613 used in sequence similarity searches and sequence feature analyses. *A. thaliana* de novo genes
614 were retrieved from a set of 782 putative DNGs recently described by Li et al (2016) using
615 sequence homology searches. These genes were screened to identify high-confidence DNGs
616 supported by further comparative genomic data, particularly synteny information. Specifically,
617 we performed Blast v2.11.0 (49) searches of the corresponding protein sequences against several
618 NCBI databases. We used Blastp to search the “nr” and “tsa_nr” databases, and tBlastn to search
619 the “nr/nt”, “refseq_rna”, “est” and “TSA” databases with the following modified parameters:
620 num_descriptions 100 -num_alignments 100 -max_hsp 5 -evalue 0.001 -seg yes. The seg filter
621 was turned on in order to remove spurious hits due to nonhomologous stretches of similar amino
622 acids. These searches were carried out against increasingly broader taxonomic units that include
623 *Arabidopsis* species excluding *A. thaliana*, Brassicaceae (excluding *Arabidopsis*), rosidae
624 (excluding Brassicaceae), angiosperms (excluding rosidae), green plants (excluding
625 angiosperms). We also searched for homologous sequences of the 782 DNGs in fungi, bacteria

626 and archaea to identify and remove possible horizontal transfer cases (**S1 Table**). Subject
627 sequences were screened using unix scripts to remove truncated proteins. We excluded from the
628 catalog of DNGs all cases with any homology with sequences from a non-focal species.

629

630 To determine synteny conservation of DNG coding regions we searched 45 Brassicaceae genome
631 assemblies obtained from Phytozome (<https://phytozome-next.jgi.doe.gov>) using the translated
632 DNA sequence of each exon of the 782 putative DNGs, which allowed us to detect conserved
633 coding regions for genes with one or multiple exons, in tBlastx run with the following modified
634 parameters: num_descriptions 100 -num_alignments 100 -max_hsps 5 -evalue 0.001 -seg yes. A
635 list of the 45 species investigated is available in **S2 Table**.

636

637 First, we used these alignments to identify Brassicaceae genomic regions that were syntenic with
638 DNGs and with the potential to encode proteins. Although these regions did not include
639 annotated genes, they maintained long coding regions and were thus considered *bona fide*
640 homologs to DNGs. We based this selection on two criteria. We selected for hits with a
641 conserved methionine within the first five amino acids of the *A. thaliana* DNG protein, in order
642 to account for alternative first codon positions (84). Additionally, we included only hits wherein
643 the total alignment length from the first codon to the first stop codon equal to at least 75% of the
644 query protein. This threshold was selected to include hits that are likely to encode a protein,
645 given the lack of disabling mutations along most of the coding, while allowing for slightly
646 shorter loci, as stop codons may also vary slightly between orthologs. Using this strategy, we
647 identified 44 DNGs with putative homologous genes in non-*A. thaliana* Brassicaceae, which
648 were thus discarded.

649 Second, we used the same alignment data for the remaining DNGs to identify those that maintain
650 synteny conservation *in noncoding regions* with other Brassicaceae. To this end, we applied a
651 minimum threshold of 30% coverage between *A. thaliana* DNGs and Brassicaceae genomes as
652 corresponding to conserved synteny. This length threshold is lower than those used in previous
653 DNG analyses in animals (31, 85, 86), in order to account for the decreased overall synteny
654 conservation among Brassicaceae. Furthermore, the syntenic alignments were screened for the
655 presence of enabler substitutions, represented by a novel start codon, the removal of stop codons
656 and/or frameshifts in the DNG coding region compared to the syntenic DNA of the outgroup
657 species. Overall, we found 604 *A. thaliana* putative DNGs with apparent synteny conservation
658 with at least one other Brassicaceae genome.

659

660 **Identification of *Brassica rapa* de novo genes**

661 We retrieved all *Brassica rapa* genome assemblies and gene sets available as of October 2021
662 and screened each assembly for completeness using BUSCO v3.0 (87) (**S3 Table**). The *B. rapa*
663 coding regions, protein and genome assembly fasta files were downloaded from:

664 <https://ngdc.cncb.ac.cn/search/?dbId=gwh&q=GWHAAES000000000>.

665 *B. rapa* proteins containing stops (“Xs”) within their amino acid sequences were removed,
666 leaving 45,912 proteins. We further screened the remaining *B. rapa* proteins to identify and
667 remove sequences mostly formed by transposable elements (TEs), as they likely represent
668 misannotated TEs. To this aim, we first downloaded the sequences of 39,197 TE families from
669 31 Brassicaceae reported in PlantRep (88). Proteins containing TE sequences were retrieved by
670 performing a tBlastn search with the following modified parameters: -evalue 1e-10 -
671 max_target_seqs 10 -max_hsp 5.

672 A total of 9,791 *B. rapa* proteins shared sequence similarity with TEs over at least 50% of their
673 length and were removed from the dataset. Proteins with Blast matches uniquely with unknown
674 repeats were not removed as those repeats might represent microsatellites, which could
675 potentially form a portion of de novo gene coding regions.

676

677 To search for DNGs among the remaining 36,121 *B. rapa* proteins, we carried out a
678 multipronged homology search strategy to identify proteins with homology in non-*Brassica rapa*
679 genomes. First, the *B. rapa* proteins were searched against the plant NCBI refseq protein set
680 obtained on August 31, 2021 throughout a Blastp run with the following modified parameters: -
681 num_descriptions 5 -num_alignments 5 -evaluate 0.001 -seg yes. A total of 18,413 proteins
682 showed no sequence homology to NCBI refseq proteins with the exception of *Brassica*
683 sequences, thus representing candidate *Brassica* orphan proteins. A further Blastp search was
684 performed against the NCBI non-random, tsa_nr, refseq_rna and est databases of all
685 Brassicaceae proteins with the following modified parameters: -max_target_seqs 50 -max_hsp
686 5. All hits containing premature stop codons were removed as they could represent expressed
687 sequences of noncoding genes or truncated and thus non-functional proteins.

688

689 Similarly to the procedure applied to *A. thaliana* putative DNGs, we carried out Blast searches
690 against increasingly broad taxonomic units containing *B. rapa*, starting from *Brassica* but
691 excluding *B. rapa*, then other *Brassicaceae*, rosidae, angiosperms, green plants, excluding the
692 previous taxon at each step, and using the following modified parameters: -max_target_seqs 50 -
693 max_hsp 5 -evaluate 0.001 -seg yes. Fungal proteomes were also screened, whereas Archaea and
694 Bacteria sequences were not included as our analyses in *A. thaliana* DNGs showed that

695 prokaryotic databases contributed marginally to the detection of homologs. We found 2,089 *B.*
696 *rapa* proteins sharing no sequenced homology with two or more non-*Brassica* proteins, thus
697 representing *B. rapa* orphan genes. A cut-off of at least two non-*Brassica* proteins was
698 implemented to take into account possible contamination from *A. thaliana* into other
699 Brassicaceae genome datasets (89, 90). This number of orphan genes in *B. rapa* is similar to but
700 higher than the 1,540 *B. rapa* orphan genes recently described (45), probably because of
701 differences in the homology search criteria and in the gene annotation version. We further
702 removed from the list of orphan genes 35 genes with annotated protein domains in eggNOG
703 Mapper (91), as they likely represent fast evolving non-DNGs.

704

705 In order to identify enabler changes uniquely associated with de novo gene birth, we inspected
706 the alignments from the Blast searches between *B. rapa* orphan genes and the genome of 45
707 Brassicaceae (**S2 Table**). After applying the same approach described to filter *A. thaliana*
708 putative DNGs, we obtained 754 *B. rapa*-specific de novo genes.

709

710

711 **Examination of DNG coding and protein sequence features**

712 *DNA sequence features.* Length of genes, predicted coding sequences (CDS) and (where
713 available) UTRs were retrieved from gff files of the assemblies of each species. The GC-content
714 of each coding region was calculated using an in-house perl script. Transposable element (TE)
715 genome coordinates were obtained from the TAIR10 gff3 file in the *A. thaliana* TAIR repository
716 and from the *Brapa_genome_v3.0_TE.gff* file downloaded from
717 <http://brassicadb.cn/#/Download/> for *B. rapa*.

718 *O. sativa* de novo gene IDs were retrieved from supplementary information in Zhang et al. (8).

719 *O. sativa* release v3 coding region fasta sequences were downloaded from the Gramene

720 repository (http://ftp.gramene.org/oge/release-3/fasta/oryza_sativa/dna/). Rice protein fasta

721 sequences were obtained translating the CDS sequences using the ORF finder program in the

722 SMS suite (92). The initial set of 50,556 genes was parsed to retain only the longest isoform of

723 each locus and to remove genes with premature stop codons, leaving 38,748 genes.

724

725 Microsatellites were retrieved from the coding regions of each gene using Tandem Repeat Finder

726 v4.09 (93) with default options. Overlap of TEs and microsatellites with the coding region and

727 gene distance from TEs were obtained using bedtools (94). The coding potential was estimated

728 using the Coding Potential Calculator (72). The codon adaptation index was calculated using the

729 ‘cai’ tool in the EMBOSS suite (95).

730

731 Kozak scores were computed as the sum of the logs of nucleotide probabilities within a window

732 of 12 bp around the ATG start codon (96), based on probabilities extracted from all genes in each

733 species genome:

734
$$s = \sum_{i=-6}^5 \log(p(n_i))$$

735

736 where n_i is the observed nucleotide at position i relative to the start of the ATG. The scores are

737 generally negative with a mean around -5.1, but the closer to zero, the more like the consensus

738 sequence (AAAAAATGGCG) they are. This is similar, but not identical, to the

739 acACAATGGC consensus sequence for terrestrial plants (97), reflecting biases that promote

740 translation initiation by the ribosome. The nucleotide probability profiles surrounding start
741 codons in *A. thaliana*, *B. rapa*, and *O. sativa*, as well as other diverse plant genomes (*Petunia*
742 *inflata*: 36,489 genes; *Quercus robur*: 25,808 genes) are highly similar, although genes in *O.*
743 *sativa* appear to have a relaxed constraint in the nucleotide following the ATG, whereas it is
744 guanine over 50% of the time in the other two species (see **S1 Fig**), which could be related to
745 fact that rice is a monocot and thus distantly related to the dicot family Brassicaceae. The
746 preference for adenines upstream of the ATG in plants is much less pronounced than in
747 nucleotide profiles of Kozak sequences in other eukaryotes (e.g. human, *Drosophila*) (see **S1**
748 **Fig**).

749

750 *Protein sequence features*. Protein domains were obtained from the NCBI Conserved Domain
751 Database (98). Protein structural disorder was calculated using IUPred2 (99) after removing
752 cysteines from the protein sequences in order to account for the possible presence of the disulfide
753 bonds, which can strongly affect ISD estimates (100). Transmembrane helices were estimated
754 using TMHMM Server v. 2.0 (101). The identified helices may represent either transmembrane
755 structures or signal peptides, which tend to occur within the first 60 amino acid in the N-terminal
756 of the protein. Proteins were conservatively assigned putative transmembrane helices if they
757 contained one or more predicted helices and at least 18 amino acids in a helix past the first 60
758 amino acids. Signal peptides were predicted using TargetP2.0 (102). The isoelectric point of each
759 protein was calculated using the Sequence Manipulation Suite Protein Isoelectric Point tool (92).
760 All programs above were run using default settings.

761

762 *Functional genomic features in A. thaliana*. The *A. thaliana* TAIR10 gff3 file
763 “TAIR10_GFF3_genes.gff” in the TAIR repository was used to obtain the length of “5’UTR
764 length” and “3’UTR length” and the number of coding exons (“#Exons”). The proportion of the
765 coding region overlapping with transposable elements (TEs; “TEs in CDS (bp)” and “%TEs in
766 CDS”) and the gene distance to the nearest TE (“TEdist”) were calculated using bedtools (94)
767 and the genome coordinate of TAIR10 coding exons and TEs. Genome coordinate of TEs were
768 obtained from the gff3 file “TAIR10_GFF3_genes_transposons.gff” available in the TAIR
769 repository. Possible regulatory motifs (“#Motifs promoter”) in the promoter regions of DNGs
770 and AGs were identified using the MEME suite (103). The DNA sequences corresponding to the
771 300bp upstream of the transcription start site of each TAIR10 gene were retrieved using bedtools
772 and screened using the MEME Streme tools (104).
773
774 Transcription factor binding site (TFBS) information was retrieved from the Plant *cis*-Map
775 genome browser (<http://ucsc.gao-lab.org/index.html>). The Conserved TFBS dataset included
776 binding sites deposited in the PlantRegMap database (105, 106). Conservation of TFBSs was
777 assessed using multiple genome alignments across species from the Plant *cis*-Map genome
778 browser conservation track. Binding sites with at least 50% of their sequence falling within
779 conserved elements were considered conserved (“#Conserved TFBSs”). The pipeline to identify
780 putative functional transcription factor binding sites (“#FUN TFBSs”) is described in Tian et al.
781 (105). The AtRegNet confirmed TFBSs dataset (“#AtRegNet TFBSs”) was downloaded from the
782 AGRIS database (107). The bedtools suite was used to extract the DNA sequences of the 200bp
783 upstream of transcription start site of TAIR10 genes, corresponding to the putative promoter
784 regions, and intersected with TFBS genome coordinates.

785
786 Transcription quantification features were obtained from a study of 18 natural *A. thaliana*
787 accessions (108) study. The average (“RNAseq AVG”) and maximum (“RNAseq MAX”)
788 expression across 48 samples, reported as log(rpkm) values, were calculated for each gene across
789 45 samples. Expression in only 1 out of 45 samples was also added (“RNAseq <2 samples”).
790 Average (“RP AVG”) and maximum (“RP MAX”) ribosome profiling (RP) expression data,
791 reported as log(rpkm) values, were also calculated for root tissues including control and deficient
792 phosphorous nutrition conditions in a total of 25 samples (109). The maximum RP expression
793 was calculated only for genes expressed in at least two samples.

794
795 The “Missense variant” feature represents the number of missense (nonsynonymous)
796 substitutions divided by the length of the coding in each gene (“Missense variation”). Missense
797 substitutions were obtained from the 1001 *A. thaliana* Genomes portal
798 (<https://1001genomes.org/index.html>).

799
800 Protein domains and protein-protein interactions (PPIs) were obtained from the files
801 TAIR10_all.domains (“#HMMPfam Domain”) and TairProteinInteraction 20090527.txt
802 (“#PPIs”, “#PPIs (w/predicted)”), respectively, from the TAIR repository. The PPI data contain
803 interaction annotations extracted from the literature by TAIR and BIOGRID (110). The
804 frequency of three different categories of amino acids, “Tiny”, “Aromatic” and “Acidic”, were
805 obtained using the pepstats program in the EMBOSS suite (95).

806

807 We gathered phenotypic information using data deposited in the TAIR repository containing
808 phenotypic data extracted from the literature by TAIR. Data files names: TAIR_Phenotypes_9-
809 2019.txt (“TAIR_Phenotypes_9-2019.txt”), Locus_Germplasm_Phenotype_20190630.txt
810 (“LGP_20190630.txt”), Locus_Germplasm_Phenotype_20130122 (“LGP_20130122”).
811 Additionally, data from manually curated meta-analysis of loss-of-function phenotypes (111)
812 (“Lloyd and Meinke 2012”) and from a high-throughput phenotype screening of annotated genes
813 with unknown function (112) (“Luhua, et al. 2013”) were included.

814

815

816 **Machine learning training and testing**

817 Balanced sets of DNGs and AGs were selected for ML training and testing within each species.
818 For instance, in *A. thaliana* 331 AGs were randomly chosen among the 26,423 available AGs.
819 The Decision Tree (DT) classifier was trained using the *scikit-learn* package in Python (113).
820 The feed-forward fully-connected neural network (NN) with a single hidden layer with 20
821 hidden nodes was also trained with the MLPClassifier implementation in *scikit-learn*, using *tanh*
822 activation functions and the ‘adam’ solver. Log transformations were applied to length-based
823 features (Gene length, CDS length, Distance from TEs) and to functional genomic features with
824 RPKM in *A. thaliana*.

825

826 Feature importance (normalized decrease in Gini impurity index at each node where a feature is
827 used, weighted by the fraction of training examples represented at those nodes) was calculated
828 using the ‘feature_importance’ attribute of the DecisionTreeClassifier generated by *scikit-learn*.
829 In some runs, “At-least-one-domain” was the feature at the root of the tree, and in other runs,

830 “Coding potential (cpc2)” represented the splitting feature at the root. Thus, we estimated the
831 importance of features by averaging them over multiple runs.

832

833 Scripts and data from this study are available at <https://github.com/ioerger2/DNG>. This
834 repository contains instructions on how to run Decision Tree and Neural Network training and
835 testing. A python script generates accuracies, confusion matrices, decision trees and feature
836 importances for each dataset. The *A. thaliana* features are available in the repository.

837

838

839 **Statistical analyses**

840 All statistical analyses were performed in R (cit). In MLA testing, accuracy corresponds to the
841 sum of true positives (TPs) and true negatives (TNs), divided by all genes:
842 $(TP+TN)/(TP+FP+TN+FN)$. Sensitivity and specificity are represented by true positive rate
843 $(TPR=TP/(TP+FN))$ and true negative rate $(TNR=TN/(TN+FP))$. Test sets and training sets are
844 always kept disjoint (with no overlap of genes) in all runs.

845

846

847

848

849 **Acknowledgements**

850 The authors would like to thank Michael Dickens and the Texas A&M University High
851 Performance Research Computing facility for assistance with performing homology search
852 analyses on HPC clusters.

853

854

855 **Author Contributions**

856 CC, AEP and TRI conceived and designed the study. CC, AO and TRI analyzed the data. TRI.

857 CC, AEP and TRI wrote the paper.

858

859

860 **REFERENCES**

- 861 1. Cui L, Wall PK, Leebens-Mack JH, Lindsay BG, Soltis DE, Doyle JJ, et al. Widespread
862 genome duplications throughout the history of flowering plants. *Genome research*.
863 2006;16(6):738-49.
- 864 2. Lin R, Ding L, Casola C, Ripoll DR, Feschotte C, Wang H. Transposase-derived
865 transcription factors regulate light signaling in Arabidopsis. *Science (New York, NY)*.
866 2007;318(5854):1302-5.
- 867 3. Flagel LE, Wendel JF. Gene duplication and evolutionary novelty in plants. *New*
868 *Phytologist*. 2009;183(3):557-64.
- 869 4. Panchy N, Lehti-Shiu M, Shiu S-H. Evolution of gene duplication in plants. *Plant*
870 *physiology*. 2016;171(4):2294-316.
- 871 5. Wendel JF, Jackson SA, Meyers BC, Wing RA. Evolution of plant genome architecture.
872 *Genome biology*. 2016;17(1):37.
- 873 6. Qiao X, Li Q, Yin H, Qi K, Li L, Wang R, et al. Gene duplication and evolution in
874 recurring polyploidization–diploidization cycles in plants. *Genome biology*. 2019;20(1):1-23.
- 875 7. Li L, Wurtele ES. The QQS orphan gene of Arabidopsis modulates carbon and nitrogen
876 allocation in soybean. *Plant Biotechnol J*. 2015;13(2):177-87.
- 877 8. Zhang L, Ren Y, Yang T, Li G, Chen J, Gschwend AR, et al. Rapid evolution of protein
878 diversity by de novo origination in *Oryza*. *Nat Ecol Evol*. 2019;3(4):679-90.
- 879 9. Khalturin K, Hemmrich G, Fraune S, Augustin R, Bosch TCG. More than just orphans:
880 are taxonomically-restricted genes important in evolution? *Trends in Genetics*. 2009;25(9):404-
881 13.
- 882 10. Yang X, Jawdy S, Tschaplinski TJ, Tuskan GA. Genome-wide identification of lineage-
883 specific genes in Arabidopsis, *Oryza* and *Populus*. *Genomics*. 2009;93(5):473-80.
- 884 11. Donoghue MT, Keshavaiah C, Swamidatta SH, Spillane C. Evolutionary origins of
885 Brassicaceae specific genes in Arabidopsis thaliana. *BMC Evol Biol*. 2011;11:47.
- 886 12. Arendsee ZW, Li L, Wurtele ES. Coming of age: orphan genes in plants. *Trends Plant*
887 *Sci*. 2014;19(11):698-708.
- 888 13. Richardson AO, Palmer JD. Horizontal gene transfer in plants. *Journal of experimental*
889 *botany*. 2007;58(1):1-9.

- 890 14. Wickell DA, Li FW. On the evolutionary significance of horizontal gene transfers in
891 plants. *New Phytologist*. 2020;225(1):113-7.
- 892 15. Joly-Lopez Z, Forczek E, Hoen DR, Juretic N, Bureau TE. A gene family derived from
893 transposable elements during early angiosperm evolution has reproductive fitness benefits in
894 *Arabidopsis thaliana*. *PLoS Genet*. 2012;8(9):e1002931.
- 895 16. Lisch D. How important are transposons for plant evolution? *Nat Rev Genet*.
896 2013;14(1):49-61.
- 897 17. Gould SJ, Vrba ES. Exaptation - a Missing Term in the Science of Form. *Paleobiology*.
898 1982;8(1):4-15.
- 899 18. Van Oss SB, Carvunis AR. De novo gene birth. *PLoS Genet*. 2019;15(5):e1008160.
- 900 19. McLysaght A, Hurst LD. Open questions in the study of de novo genes: what, how and
901 why. *Nat Rev Genet*. 2016;17(9):567-78.
- 902 20. Bornberg-Bauer E, Hlouchova K, Lange A. Structure and function of naturally evolved
903 de novo proteins. *Curr Opin Struct Biol*. 2021;68:175-83.
- 904 21. Li D, Dong Y, Jiang Y, Jiang H, Cai J, Wang W. A de novo originated gene depresses
905 budding yeast mating pathway and is repressed by the protein encoded by its antisense strand.
906 *Cell Res*. 2010;20(4):408-20.
- 907 22. Bungard D, Copple JS, Yan J, Chhun JJ, Kumirov VK, Foy SG, et al. Foldability of a
908 Natural De Novo Evolved Protein. *Structure*. 2017;25(11):1687-96 e4.
- 909 23. Baalsrud HT, Torresen OK, Solbakken MH, Salzburger W, Hanel R, Jakobsen KS, et al.
910 De Novo Gene Evolution of Antifreeze Glycoproteins in Codfishes Revealed by Whole Genome
911 Sequence Data. *Mol Biol Evol*. 2018;35(3):593-606.
- 912 24. Zhuang X, Yang C, Murphy KR, Cheng CC. Molecular mechanism and history of non-
913 sense to sense evolution of antifreeze glycoprotein gene in northern gadids. *Proc Natl Acad Sci*
914 *U S A*. 2019;116(10):4400-5.
- 915 25. Li CY, Zhang Y, Wang Z, Zhang Y, Cao C, Zhang PW, et al. A human-specific de novo
916 protein-coding gene associated with human brain functions. *PLoS Comput Biol*.
917 2010;6(3):e1000734.
- 918 26. Samusik N, Krukovskaya L, Meln I, Shilov E, Kozlov AP. PBOV1 is a human de novo
919 gene with tumor-specific expression that is associated with a positive clinical outcome of cancer.
920 *PLoS One*. 2013;8(2):e56162.
- 921 27. Suenaga Y, Islam SM, Alagu J, Kaneko Y, Kato M, Tanaka Y, et al. NCYM, a Cis-
922 antisense gene of MYCN, encodes a de novo evolved protein that inhibits GSK3beta resulting in
923 the stabilization of MYCN in human neuroblastomas. *PLoS Genet*. 2014;10(1):e1003996.
- 924 28. Begun DJ, Lindfors HA, Thompson ME, Holloway AK. Recently evolved genes
925 identified from *Drosophila yakuba* and *D. erecta* accessory gland expressed sequence tags.
926 *Genetics*. 2006;172(3):1675-81.
- 927 29. Begun DJ, Lindfors HA, Kern AD, Jones CD. Evidence for de novo evolution of testis-
928 expressed genes in the *Drosophila yakuba*/*Drosophila erecta* clade. *Genetics*. 2007;176(2):1131-
929 7.
- 930 30. Knowles DG, McLysaght A. Recent de novo origin of human protein-coding genes.
931 *Genome Res*. 2009;19(10):1752-9.
- 932 31. Ruiz-Orera J, Hernandez-Rodriguez J, Chiva C, Sabido E, Kondova I, Bontrop R, et al.
933 Origins of De Novo Genes in Human and Chimpanzee. *PLoS Genet*. 2015;11(12):e1005721.
- 934 32. Murphy DN, McLysaght A. De novo origin of protein-coding genes in murine rodents.
935 *PLoS One*. 2012;7(11):e48650.

- 936 33. Prabh N, Rodelsperger C. De Novo, Divergence, and Mixed Origin Contribute to the
937 Emergence of Orphan Genes in *Pristionchus Nematodes*. *G3 (Bethesda)*. 2019;9(7):2277-86.
- 938 34. Cai J, Zhao R, Jiang H, Wang W. De novo origination of a new protein-coding gene in
939 *Saccharomyces cerevisiae*. *Genetics*. 2008;179(1):487-96.
- 940 35. Carvunis AR, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis N, et al.
941 Proto-genes and de novo gene birth. *Nature*. 2012;487(7407):370-4.
- 942 36. Vakirlis NN, Hebert AS, Opulente DA, Achaz G, Hittinger CT, Fischer G, et al. A
943 molecular portrait of de novo genes in yeasts. *Mol Biol Evol*. 2017;35(3):631-45.
- 944 37. Blevins WR, Ruiz-Orera J, Messegueur X, Blasco-Moreno B, Villanueva-Canas JL,
945 Espinar L, et al. Uncovering de novo gene birth in yeast using deep transcriptomics. *Nat*
946 *Commun*. 2021;12(1):604.
- 947 38. Durand E, Gagnon-Arsenault I, Hallin J, Hatin I, Dube AK, Nielly-Thibault L, et al.
948 Turnover of ribosome-associated transcripts from de novo ORFs produces gene-like
949 characteristics available for de novo gene emergence in wild yeast populations. *Genome Res*.
950 2019;29(6):932-43.
- 951 39. Poretti M, Praz CR, Sotiropoulos AG, Wicker T. A survey of lineage-specific genes in
952 Triticeae reveals de novo gene evolution from genomic raw 1 material. *bioRxiv*. 2022.
- 953 40. Campbell MA, Zhu W, Jiang N, Lin H, Ouyang S, Childs KL, et al. Identification and
954 characterization of lineage-specific genes within the Poaceae. *Plant Physiol*. 2007;145(4):1311-
955 22.
- 956 41. Li ZW, Chen X, Wu Q, Hagmann J, Han TS, Zou YP, et al. On the Origin of De Novo
957 Genes in *Arabidopsis thaliana* Populations. *Genome Biol Evol*. 2016;8(7):2190-202.
- 958 42. Lin H, Moghe G, Ouyang S, Iezzoni A, Shiu SH, Gu X, et al. Comparative analyses
959 reveal distinct sets of lineage-specific genes within *Arabidopsis thaliana*. *BMC Evol Biol*.
960 2010;10:41.
- 961 43. Wu DD, Wang X, Li Y, Zeng L, Irwin DM, Zhang YP. "Out of pollen" hypothesis for
962 origin of new genes in flowering plants: study from *Arabidopsis thaliana*. *Genome Biol Evol*.
963 2014;6(10):2822-9.
- 964 44. Xu Y, Wu G, Hao B, Chen L, Deng X, Xu Q. Identification, characterization and
965 expression analysis of lineage-specific genes within sweet orange (*Citrus sinensis*). *BMC*
966 *Genomics*. 2015;16:995.
- 967 45. Jiang M, Dong X, Lang H, Pang W, Zhan Z, Li X, et al. Mining of Brassica-Specific
968 Genes (BSGs) and Their Induction in Different Developmental Stages and under
969 *Plasmodiophora brassicae* Stress in *Brassica rapa*. *Int J Mol Sci*. 2018;19(7).
- 970 46. Jiang M, Zhan Z, Li H, Dong X, Cheng F, Piao Z. *Brassica rapa* orphan genes largely
971 affect soluble sugar metabolism. *Hortic Res*. 2020;7(1):181.
- 972 47. Cui X, Lv Y, Chen M, Nikoloski Z, Twell D, Zhang D. Young Genes out of the Male: An
973 Insight from Evolutionary Age Analysis of the Pollen Transcriptome. *Molecular plant*.
974 2015;8(6):935-45.
- 975 48. Moyers BA, Zhang J. Evaluating Phylostratigraphic Evidence for Widespread De Novo
976 Gene Birth in Genome Evolution. *Mol Biol Evol*. 2016;33(5):1245-56.
- 977 49. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
978 architecture and applications. *BMC Bioinformatics*. 2009;10:421.
- 979 50. Domazet-Loso T, Brajkovic J, Tautz D. A phylostratigraphy approach to uncover the
980 genomic history of major adaptations in metazoan lineages. *Trends in genetics : TIG*.
981 2007;23(11):533-9.

- 982 51. Casola C. From De Novo to "De Nono": The Majority of Novel Protein-Coding Genes
983 Identified with Phylostratigraphy Are Old Genes or Recent Duplicates. *Genome Biol Evol.*
984 2018;10(11):2906-18.
- 985 52. Tautz D, Domazet-Loso T. The evolutionary origin of orphan genes. *Nat Rev Genet.*
986 2011;12(10):692-702.
- 987 53. Wissler L, Gadau J, Simola DF, Helmkampf M, Bornberg-Bauer E. Mechanisms and
988 dynamics of orphan gene emergence in insect genomes. *Genome Biol Evol.* 2013;5(2):439-55.
- 989 54. Weisman CM, Murray AW, Eddy SR. Many, but not all, lineage-specific genes can be
990 explained by homology detection failure. *PLoS Biol.* 2020;18(11):e3000862.
- 991 55. Vakirlis N, Carvunis AR, McLysaght A. Synteny-based analyses indicate that sequence
992 divergence is not the main source of orphan genes. *Elife.* 2020;9.
- 993 56. Moyers BA, Zhang J. Phylostratigraphic bias creates spurious patterns of genome
994 evolution. *Mol Biol Evol.* 2015;32(1):258-67.
- 995 57. Moyers BA, Zhang JZ. Further Simulations and Analyses Demonstrate Open Problems of
996 Phylostratigraphy. *Genome Biology and Evolution.* 2017;9(6):1519-27.
- 997 58. Vakirlis N, McLysaght A. Computational Prediction of De Novo Emerged Protein-
998 Coding Genes. *Methods Mol Biol.* 2019;1851:63-81.
- 999 59. Washburn JD, Mejia-Guerra MK, Ramstein G, Kremling KA, Valluru R, Buckler ES, et
1000 al. Evolutionarily informed deep learning methods for predicting relative transcript abundance
1001 from DNA sequence. *Proc Natl Acad Sci U S A.* 2019;116(12):5542-9.
- 1002 60. Lee D, Zhang J, Liu J, Gerstein M. Epigenome-based splicing prediction using a
1003 recurrent neural network. *PLoS Comput Biol.* 2020;16(6):e1008006.
- 1004 61. Wang H, Cimen E, Singh N, Buckler E. Deep learning for plant genomics and crop
1005 improvement. *Curr Opin Plant Biol.* 2020;54:34-41.
- 1006 62. Gao Q, Jin X, Xia E, Wu X, Gu L, Yan H, et al. Identification of Orphan Genes in
1007 Unbalanced Datasets Based on Ensemble Learning. *Front Genet.* 2020;11:820.
- 1008 63. Zhang X, Xuan J, Yao C, Gao Q, Wang L, Jin X, et al. A deep learning approach for
1009 orphan gene identification in moso bamboo (*Phyllostachys edulis*) based on the CNN +
1010 Transformer model. *BMC Bioinformatics.* 2022;23(1):162.
- 1011 64. Japkowicz N, Stephen S. The class imbalance problem: A systematic study. *Intelligent*
1012 *Data Analysis.* 2002;6(5):429-49.
- 1013 65. Wilson BA, Foy SG, Neme R, Masel J. Young Genes are Highly Disordered as Predicted
1014 by the Preadaptation Hypothesis of De Novo Gene Birth. *Nat Ecol Evol.* 2017;1(6):0146-146.
- 1015 66. Wang X, Wang H, Wang J, Sun R, Wu J, Liu S, et al. The genome of the mesopolyploid
1016 crop species *Brassica rapa*. *Nat Genet.* 2011;43(10):1035-9.
- 1017 67. Cheng F, Wu J, Fang L, Wang X. Syntenic gene analysis between *Brassica rapa* and
1018 other Brassicaceae species. *Front Plant Sci.* 2012;3:198.
- 1019 68. Tong C, Wang X, Yu J, Wu J, Li W, Huang J, et al. Comprehensive analysis of RNA-seq
1020 data reveals the complexity of the transcriptome in *Brassica rapa*. *BMC genomics.* 2013;14(1):1-
1021 10.
- 1022 69. Cheng F, Sun R, Hou X, Zheng H, Zhang F, Zhang Y, et al. Subgenome parallel selection
1023 is associated with morphotype diversification and convergent crop domestication in *Brassica*
1024 *rapa* and *Brassica oleracea*. *Nat Genet.* 2016;48(10):1218-24.
- 1025 70. Cheng F, Wu J, Cai C, Fu L, Liang J, Borm T, et al. Genome resequencing and
1026 comparative variome analysis in a *Brassica rapa* and *Brassica oleracea* collection. *Sci Data.*
1027 2016;3:160119.

- 1028 71. Zhang L, Cai X, Wu J, Liu M, Grob S, Cheng F, et al. Improved Brassica rapa reference
1029 genome by single-molecule sequencing and chromosome conformation capture technologies.
1030 Horticulture Res. 2018;5:50.
- 1031 72. Kang YJ, Yang DC, Kong L, Hou M, Meng YQ, Wei L, et al. CPC2: a fast and accurate
1032 coding potential calculator based on sequence intrinsic features. Nucleic Acids Res.
1033 2017;45(W1):W12-W6.
- 1034 73. Wang HC, Hickey DA. Rapid divergence of codon usage patterns within the rice
1035 genome. BMC Evol Biol. 2007;7 Suppl 1:S6.
- 1036 74. Vakirlis N, Acar O, Hsu B, Castilho Coelho N, Van Oss SB, Wacholder A, et al. De novo
1037 emergence of adaptive membrane proteins from thymine-rich genomic sequences. Nat Commun.
1038 2020;11(1):781.
- 1039 75. Luis Villanueva-Canas J, Ruiz-Orera J, Agea MI, Gallo M, Andreu D, Alba MM. New
1040 Genes and Functional Innovation in Mammals. Genome Biol Evol. 2017;9(7):1886-900.
- 1041 76. Zhao L, Saelao P, Jones CD, Begun DJ. Origin and spread of de novo genes in
1042 Drosophila melanogaster populations. Science (New York, NY). 2014;343(6172):769-72.
- 1043 77. Bitard-Feildel T, Heberlein M, Bornberg-Bauer E, Callebaut I. Detection of orphan
1044 domains in Drosophila using "hydrophobic cluster analysis". Biochimie. 2015;119:244-53.
- 1045 78. Mukherjee S, Panda A, Ghosh TC. Elucidating evolutionary features and functional
1046 implications of orphan genes in Leishmania major. Infect Genet Evol. 2015;32:330-7.
- 1047 79. Schmitz JF, Ullrich KK, Bornberg-Bauer E. Incipient de novo genes can evolve from
1048 frozen accidents that escaped rapid transcript turnover. Nat Ecol Evol. 2018;Epub.
- 1049 80. Basile W, Sachenkova O, Light S, Elofsson A. High GC content causes orphan proteins
1050 to be intrinsically disordered. PLoS Comput Biol. 2017;13(3):e1005375.
- 1051 81. Charfe F, Rivera AJ, del Jesus MJ, Herrera F. Addressing imbalance in multilabel
1052 classification: Measures and random resampling algorithms. Neurocomputing. 2015;163:3-16.
- 1053 82. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: synthetic minority over-
1054 sampling technique. Journal of artificial intelligence research. 2002;16:321-57.
- 1055 83. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The
1056 Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. Nucleic
1057 Acids Res. 2012;40(Database issue):D1202-10.
- 1058 84. Li YR, Liu MJ. Prevalence of alternative AUG and non-AUG translation initiators and
1059 their regulatory effects across plants. Genome Res. 2020;30(10):1418-33.
- 1060 85. Heames B, Schmitz J, Bornberg-Bauer E. A Continuum of Evolving De Novo Genes
1061 Drives Protein-Coding Novelty in Drosophila. J Mol Evol. 2020;88(4):382-98.
- 1062 86. Vakirlis N, Duggan KM, McLysaght A. De novo birth of functional, human-specific
1063 microproteins. bioRxiv. 2021.
- 1064 87. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
1065 assessing genome assembly and annotation completeness with single-copy orthologs.
1066 Bioinformatics. 2015;31(19):3210-2.
- 1067 88. Luo X, Chen S, Zhang Y. PlantRep: a database of plant repetitive elements. Plant Cell
1068 Rep. 2022;41(4):1163-6.
- 1069 89. Laurin-Lemay S, Brinkmann H, Philippe H. Origin of land plants revisited in the light of
1070 sequence contamination and missing data. Curr Biol. 2012;22(15):R593-4.
- 1071 90. Steinegger M, Salzberg SL. Terminating contamination: large-scale search identifies
1072 more than 2,000,000 contaminated entries in GenBank. Genome Biol. 2020;21(1):115.

- 1073 91. Huerta-Cepas J, Szklarczyk D, Heller D, Hernandez-Plaza A, Forslund SK, Cook H, et al.
1074 eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource
1075 based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.* 2019;47(D1):D309-D14.
1076 92. Stothard P. The sequence manipulation suite: JavaScript programs for analyzing and
1077 formatting protein and DNA sequences. *Biotechniques.* 2000;28(6):1102, 4.
1078 93. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids*
1079 *Res.* 1999;27(2):573-80.
1080 94. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
1081 features. *Bioinformatics.* 2010;26(6):841-2.
1082 95. Rice P, Longden I, Bleasby A. EMBOSS: the European Molecular Biology Open
1083 Software Suite. *Trends in genetics : TIG.* 2000;16(6):276-7.
1084 96. Kozak M. An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs.
1085 *Nucleic Acids Res.* 1987;15(20):8125-48.
1086 97. Lutcke HA, Chow KC, Mickel FS, Moss KA, Kern HF, Scheele GA. Selection of AUG
1087 initiation codons differs in plants and animals. *Embo J.* 1987;6(1):43-8.
1088 98. Marchler-Bauer A, Derbyshire MK, Gonzales NR, Lu S, Chitsaz F, Geer LY, et al. CDD:
1089 NCBI's conserved domain database. *Nucleic Acids Res.* 2015;43(Database issue):D222-6.
1090 99. Dosztanyi Z. Prediction of protein disorder based on IUPred. *Protein Sci.*
1091 2018;27(1):331-40.
1092 100. Uversky VN, Dunker AK. Understanding protein non-folding. *Biochim Biophys Acta.*
1093 2010;1804(6):1231-64.
1094 101. Moller S, Croning MD, Apweiler R. Evaluation of methods for the prediction of
1095 membrane spanning regions. *Bioinformatics.* 2001;17(7):646-53.
1096 102. Almagro Armenteros JJ, Salvatore M, Emanuelsson O, Winther O, von Heijne G,
1097 Elofsson A, et al. Detecting sequence signals in targeting peptides using deep learning. *Life Sci*
1098 *Alliance.* 2019;2(5).
1099 103. Bailey TL, Johnson J, Grant CE, Noble WS. The MEME Suite. *Nucleic Acids Res.*
1100 2015;43(W1):W39-49.
1101 104. Bailey TL. STREME: Accurate and versatile sequence motif discovery. *Bioinformatics.*
1102 2021.
1103 105. Tian F, Yang DC, Meng YQ, Jin J, Gao G. PlantRegMap: charting functional regulatory
1104 maps in plants. *Nucleic Acids Res.* 2020;48(D1):D1104-D13.
1105 106. Jin J, Tian F, Yang DC, Meng YQ, Kong L, Luo J, et al. PlantTFDB 4.0: toward a central
1106 hub for transcription factors and regulatory interactions in plants. *Nucleic Acids Res.*
1107 2017;45(D1):D1040-D5.
1108 107. Yilmaz A, Mejia-Guerra MK, Kurz K, Liang X, Welch L, Grotewold E. AGRIS: the
1109 Arabidopsis Gene Regulatory Information Server, an update. *Nucleic Acids Res.*
1110 2011;39(Database issue):D1118-22.
1111 108. Gan X, Stegle O, Behr J, Steffen JG, Drewe P, Hildebrand KL, et al. Multiple reference
1112 genomes and transcriptomes for Arabidopsis thaliana. *Nature.* 2011;477(7365):419-23.
1113 109. Bazin J, Baerenfaller K, Gosai SJ, Gregory BD, Crespi M, Bailey-Serres J. Global
1114 analysis of ribosome-associated noncoding RNAs unveils new modes of translational regulation.
1115 *Proc Natl Acad Sci U S A.* 2017;114(46):E10018-E27.
1116 110. Oughtred R, Rust J, Chang C, Breitkreutz BJ, Stark C, Willems A, et al. The BioGRID
1117 database: A comprehensive biomedical resource of curated protein, genetic, and chemical
1118 interactions. *Protein Sci.* 2021;30(1):187-200.

- 1119 111. Lloyd J, Meinke D. A comprehensive dataset of genes with a loss-of-function mutant
1120 phenotype in Arabidopsis. *Plant Physiol.* 2012;158(3):1115-29.
- 1121 112. Luhua S, Hegie A, Suzuki N, Shulaev E, Luo X, Cenariu D, et al. Linking genes of
1122 unknown function with abiotic stress responses by high-throughput phenotype screening. *Physiol*
1123 *Plant.* 2013;148(3):322-33.
- 1124 113. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-
1125 learn: Machine learning in Python. *the Journal of machine Learning research.* 2011;12:2825-30.
1126