

1 **Defining functional intergenic transcribed regions based on heterogeneous**
2 **features of phenotype genes and pseudogenes**

3

4 John P. Lloyd¹, Zing Tsung-Yeh Tsai^{1,2,a}, Rosalie P. Sowers³, Nicholas L. Panchy⁴, Shin-Han
5 Shiu^{1,4,5*}

6 ¹ Department of Plant Biology, Michigan State University, East Lansing, MI 48824, USA

7 ² Institute of Information Science, Academia Sinica, Taipei, 115, Taiwan

8 ³ Department of Biochemistry and Molecular Biology, Pennsylvania State University, University
9 Park, PA 16802, USA

10 ⁴ Genetics Program, Michigan State University, East Lansing, MI 48824, USA

11 ⁵ Ecology, Evolutionary Biology, and Behavior Program, Michigan State University, East
12 Lansing, MI 48824, USA

13 * To whom correspondence should be addressed.

14

15 *Correspondence to:*

16 Shin-Han Shiu

17 Michigan State University

18 E-mail: shius@msu.edu

19 Telephone number: 517-353-7196

20

21 ^a Present address: Department of Computational Medicine and Bioinformatics, University of
22 Michigan, Ann Arbor, MI 48109, USA

23

24 Running title: Functionality of intergenic transcription

25

26 Keywords: Intergenic transcription, ncRNAs, definition of function, machine learning, molecular
27 evolution, chromatin state

28

29

30 **ABSTRACT**

31 With advances in transcript profiling, the presence of transcriptional activities in intergenic
32 regions has been well established in multiple model systems. However, whether intergenic
33 expression reflects transcriptional noise or the activity of novel genes remains unclear. We
34 identified intergenic transcribed regions (ITRs) in 15 diverse flowering plant species and found
35 that the amount of intergenic expression correlates with genome size, a pattern that could be
36 expected if intergenic expression is largely non-functional. To further assess the functionality of
37 ITRs, we first built machine learning classifiers using *Arabidopsis thaliana* as a model that can
38 accurately distinguish functional sequences (phenotype genes) and non-functional ones
39 (pseudogenes and random unexpressed intergenic regions) by integrating 93 biochemical,
40 evolutionary, and sequence-structure features. Next, by applying the models to ITRs, we found
41 that 2,453 (21%) had features significantly similar to phenotype genes and thus were likely parts
42 of functional genes, while an additional 17% resembled benchmark RNA genes. However, ~60%
43 of ITRs were more similar to nonfunctional sequences and should be considered transcriptional
44 noise unless falsified with experiments. The predictive framework establish here provides not
45 only a comprehensive look at how functional, genic sequences are distinct from likely non-
46 functional ones, but also a new way to differentiate novel genes from genomic regions with noisy
47 transcriptional activities.

48

49 INTRODUCTION

50 Advances in sequencing technology have helped to identify pervasive transcription in intergenic
51 regions with no annotated genes. These intergenic transcripts have been found in metazoa and
52 fungi, including *Homo sapiens* (human; ENCODE Project Consortium 2012), *Drosophila*
53 *melanogaster* (Brown et al. 2014), *Caenorhabditis elegans* (Boeck et al. 2016), and
54 *Saccharomyces cerevisiae* (Nagalakshmi et al. 2008). In plants, ~7,000 and ~15,000 intergenic
55 transcripts have also been reported in *Arabidopsis thaliana* (Yamada et al. 2003; Stolc et al.
56 2005; Moghe et al. 2013; Krishnakumar et al. 2015) and *Oryza sativa* (Nobuta et al. 2007),
57 respectively. The presence of intergenic transcripts indicates that there may be additional genes
58 in genomes that have escaped gene finding efforts thus far. Considering that knowledge of the
59 complete suite of functional elements present in a genome is an important goal for large-scale
60 functional genomics efforts and the quest to connect genotype to phenotype, identifying
61 functional intergenic transcribed regions (ITRs) represents a fundamental task that is critical to
62 our understanding of the gene space in a genome.

63 Loss-of-function phenotyping analysis represents the gold standard by which the
64 functional significance of genomic regions, including ITRs, can be confirmed (Niu and Jiang
65 2013). In *Mus musculus* (mouse), at least 25 ITRs with loss-of-function mutant phenotypes have
66 been identified (Sauvageau et al. 2013; Lai et al. 2015), indicating that they are *bona fide* genes.
67 In addition, loss-of-function mutants have been used to confirm ITR functionality in mouse
68 embryonic stem cell proliferation (Ivanova et al. 2006; Guttman et al. 2009) and male
69 reproductive development (Heinen et al. 2009), as well as brain and eye development in *Danio*
70 *rario* (Ulitsky et al. 2011). In human, 162 long intergenic non-coding RNAs (lincRNAs) harbor
71 phenotype-associated SNPs, suggesting that these expressed intergenic regions may be
72 functional (Ning et al. 2013). In addition to intergenic expression, most model organisms feature
73 an abundance of annotated non-coding RNA (ncRNA) sequences (Zhao et al. 2016), which are
74 mostly identified through the presence of expression occurring outside of annotated genes. Thus,
75 the only difference between ITRs and most ncRNA sequences is whether or not they have been
76 annotated. Similar to the ITR examples above, a small number of ncRNAs have been confirmed
77 as functional through loss-of-function experimental characterization, including *Xist* in mouse
78 (Penny et al. 1996; Marahrens et al. 1997), *Malat1* in human (Bernard et al. 2010), *bereft* in *D.*
79 *melanogaster* (Hardiman et al. 2002), and *At4* in *A. thaliana* (Shin et al. 2006). However, despite

80 the presence of a few notable examples, the number of ITRs and ncRNAs with well-established
81 functions is dwarfed by those with no known function.

82 While some ITRs and ncRNAs are likely novel genes, intergenic transcription can also be
83 the byproduct of noisy expression that can occur due to nonspecific landing of RNA Polymerase
84 II (RNA Pol II) or spurious regulatory signals that drive expression in random genomic regions
85 (Struhl 2007). Thus, whether an intergenic transcript is functional cannot be depend on solely the
86 fact that it is expressed. In addition to the biochemical activity, the genomic region with the
87 activity must be under selection. This line of logic has revived the classical idea on how function
88 can be defined based on “causal role” or “selected effect” functionality (Doolittle et al. 2014). A
89 “causal role” definition requires a definable activity to consider a genomic region as functional
90 (Cummins 1975; Amundson and Lauder 1994), which is adopted by the ENCODE Consortium
91 (2012) to classify ~80% of the human genome as having biochemical functions. This finding has
92 been used as evidence disproving the presence of junk DNA that are not under natural selection
93 (see Eddy 2013). This has drawn considerable critique because biochemical activity itself is not
94 an indication of selection (Graur et al., 2013; Niu and Jiang, 2013). Instead, if we are interested
95 in if a genomic region with discernible activity is under selection, selected effect functionality is
96 advocated to be a more suitable definition for function (Amundson and Lauder 1994; Graur et
97 al., 2013; Doolittle et al. 2014). Under the selected effect functionality definition, ITRs and most
98 annotated ncRNA genes remain functionally ambiguous.

99 If an ITR is functional, it would represent a genic sequence that is not identified with
100 conventional gene finding programs. Gene finding programs incorporate sequence
101 characteristics, transcriptional evidence, and conservation information to define genic regions
102 that are expected to be functional. Thus, genes that lack the features typically associated with
103 genic regions will remain unidentified. Due to the debate on the definitions of function post
104 ENCODE, Kellis et al. (2014) has suggested that evolutionary, biochemical, and genetic
105 evidences provide complementary information to define functional genomic regions. Integrating
106 chromatin accessibility, transcriptome, and conservation evidence was shown to be successful in
107 identifying regions in the human genome that are under selection (Gulko et al. 2014). Moreover,
108 a comprehensive integration of biochemical, evolutionary, and genetic evidence resulted in
109 highly-accurate identification of human disease genes and pseudogenes (Tsai et al. 2017).
110 However, it is not known if such predictions are possible outside of animal systems or if the

111 features that define functional genomic regions in animals are applicable in other biological
112 kingdoms. In plant species, despite the fact that many biochemical signatures are known to be
113 associated with genic regions, these signatures have not been incorporated to assist in identifying
114 the functional genomic regions.

115 To investigate the prevalence of intergenic transcription across species with a wide range
116 of genome sizes, we identified ITRs in 15 flowering plant species with 17-fold genome size
117 differences. To assess the functionality of plant intergenic transcripts, we first determined
118 whether 93 evolutionary, biochemical, and sequence-structure features could distinguish
119 functional sequences (phenotype genes) and non-functional ones (pseudogenes and random
120 unexpressed intergenic regions) using *A. thaliana* as a model. Next, we jointly considered all
121 features to establish functional gene prediction models using machine learning methods. Finally,
122 we applied the models to ITRs and putative ncRNAs to determine whether these functionally
123 ambiguous sequences are more similar to known functional or likely non-functional sequences.

124 **RESULTS & DISCUSSION**

125 **Relationship between genome size and intergenic expression indicates that intergenic** 126 **transcripts may generally be non-functional**

127 Transcription of unannotated, intergenic regions can be due to either activities of novel genes or
128 non-functional transcriptional noise. If noisy transcription occurs due to random landing of RNA
129 Pol II or spurious regulatory signals, a naïve expectation is that, as genome size increases, the
130 amount of intergenic expression would increase accordingly. By contrast, we expect that the
131 extent of genic sequence expression will not be significantly correlated with genome sizes
132 because larger plant genomes do not necessarily have more genes ($r^2=0.01$; $p=0.56$). Thus, to
133 gauge if intergenic transcribed regions (ITRs) generally behave more like what we expect of
134 noisy or genic transcription, we assessed the correlation between genome size and the amount of
135 intergenic expression occurring within a species.

136 We first identified genic and intergenic transcribed regions using leaf transcriptome data
137 from 15 flowering plant species with 17-fold differences in genome size (Supplemental Table 1).
138 Identical numbers of RNA-sequencing (RNA-seq) reads (30 million) and the same mapping
139 procedures were used in all species to facilitate cross-species comparisons (see Methods).
140 Transcribed regions were considered as ITRs if they did not overlap with any protein-coding or

141 RNA gene annotation and had no significant translated sequence similarity to plant protein
142 sequences (see Methods). As expected, the amount of expression originating from annotated
143 genic regions has no significant correlation with genomes size ($r^2=0.03$; $p=0.53$; **Fig. 1A**). In
144 contrast, the amount of intergenic expression occurring is significantly and positively correlated
145 ($r^2=0.30$; $p=0.04$; **Fig. 1B**). Because more intergenic expression is occurring in species with
146 more genome space, this is consistent with the interpretation that a significant proportion of
147 intergenic expression represents transcriptional noise. However, the correlation between genome
148 size and intergenic expression explained ~30% of the variation, suggesting that other factors also
149 affect ITR content, including the possibility that some ITRs are truly functional, novel genes. To
150 further evaluate the functionality of intergenic transcripts, we next identified the biochemical and
151 evolutionary features of functional genic regions and tested whether intergenic transcripts in *A.*
152 *thaliana* were more similar to functional or non-functional sequences.

153 **Expression, conservation, and epigenetic features are significantly distinct between** 154 **benchmark functional and non-functional genomic sequences**

155 To determine whether intergenic transcripts resemble functional sequences, we first asked what
156 features may allow benchmark functional and non-functional genomic regions to be
157 distinguished. For benchmark functional sequences, we used genes with visible loss-of-function
158 phenotypes when mutated (referred to as phenotype genes, $n=1,876$; see Methods). These
159 phenotype genes were considered functional based on the selected effect functionality criterion
160 (Neander 1991) because their mutations have significant growth and/or developmental impact
161 and likely contributes to reduced fitness. For benchmark non-functional genomic regions, we
162 utilized pseudogene sequences ($n=761$; see Methods). These pseudogenes exhibit sequence
163 similarity to known genes, but harbor disabling mutations including frame shifts and/or in-frame
164 stop codons, that result in the production of presumably non-functional protein products.
165 Considering that only 2% of pseudogenes are maintained over 90 million years of divergence
166 between human and mouse (Svensson et al. 2006), it is expected that the majority of
167 pseudogenes are no longer under selection (Li et al. 1981).

168 We evaluated 93 gene or gene product features for their ability to distinguish between
169 phenotype genes and pseudogenes. These features were grouped into seven categories, including
170 chromatin accessibility, DNA methylation, histone 3 (H3) marks, sequence conservation,

171 sequence-structure characteristics, transcription factor (TF) binding, and transcription activity.
172 Feature values (Supplemental Table 2) were calculated for a randomly-selected 500 base pair
173 (bp) window inside a phenotype gene or pseudogene. We used Area Under the Curve - Receiver
174 Operating Characteristic (AUC-ROC) as a metric to measure how well a feature distinguishes
175 between phenotype genes and pseudogenes. AUC-ROC values range between 0.5 (random
176 guessing) and 1 (perfect separation of functional and non-functional sequences), with AUC-ROC
177 values of 0.7, 0.8, and 0.9 considered fair, good, and excellent performance, respectively. Among
178 the seven feature categories, transcription activity features were highly informative (median
179 AUC-ROC=0.88; **Fig. 2A**). Sequence conservation, DNA methylation, TF binding, and H3
180 mark features were also fairly distinct between phenotype genes and pseudogenes (median AUC-
181 ROC ~ 0.7 for each category; **Fig. 2B-E**). By contrast, chromatin accessibility and sequence-
182 structure features were largely uninformative (median AUC-ROC=0.51 and 0.55, respectively;
183 **Fig. 2F-G**). The poor performance of chromatin accessibility features is likely because the
184 DNase I hypersensitivity (HS) datasets are sparse, as only 2-6% of phenotype gene and
185 pseudogene sequences overlap a DNase I HS site. Further, median nucleosome occupancy
186 nucleosome occupancy of phenotype genes (median normalized nucleosome occupancy = 1.22)
187 is only slightly higher than that of pseudogenes (median = 1.31; Mann Whitney U test, $p < 2e-4$).
188 For sequence-structure features based on dinucleotide structures (see Methods), we found that
189 poor performance was likely due to phenotype genes and pseudogenes sharing similar
190 dinucleotide sequence compositions (Supplemental Fig. 1).

191 **Error rates for functional region predictions are high when only single features are** 192 **considered**

193 Within each feature category, there is often a wide range of performance between features (**Fig.**
194 **2**, Supplemental Table 3). There are often clear biological or technical explanations for features
195 that perform poorly. For the transcription activity category, 17 features have an AUC-ROC
196 performance >0.8 , including the best-performing feature, expression breadth (AUC-ROC=0.95;
197 **Fig. 2A**). However, five transcription activity-related features perform poorly, including the
198 presence of expression (transcript) evidence (AUC-ROC=0.58; **Fig. 2A**). This is because 80% of
199 pseudogenes are considered expressed in ≥ 1 of 51 RNA-seq datasets, demonstrating that
200 presence of transcripts should not be used by itself as evidence of functionality. For the sequence

201 conservation category, maximum and average phastCons conservation scores were highly
202 distinct between phenotype genes and pseudogenes (AUC-ROC=0.83 and 0.82, respectively;
203 **Fig. 2B**). On the other hand, identity to best matching nucleotide sequences found in the
204 *Brassicaceae* and algal species were not informative (AUC-ROC=0.55 and 0.51, respectively;
205 **Fig. 2B**). This is because 99.8% and 95% of phenotype genes and pseudogenes, respectively, had
206 a potentially homologous sequence within the *Brassicaceae* family compared to only 3% and
207 1%, respectively, in algal species. Thus *Brassicaceae* genomes were too similar and algal
208 genomes were too dissimilar to *A. thaliana* to provide meaningful information. H3 mark features
209 also display high variability. The most informative H3 mark features are based on the number
210 and coverage of activation-related marks (AUC-ROC=0.87 and 0.85, respectively; **Fig. 2E**),
211 consistent with the notion that histone marks are often jointly associated with active genomic
212 sequences to potentially provide a robust regulatory signal (Schreiber and Bernstein 2002; Wang
213 et al. 2008). By comparison, the coverage and intensity of H3 lysine 27 trimethylation
214 (H3K27me3) and H3 threonine 3 phosphorylation (H3T3ph) are largely indistinct between
215 phenotype genes and pseudogenes (AUC-ROC range: 0.55-0.59; **Fig. 2E**).

216 Despite this high variability in performance, some features and feature categories have
217 high AUC-ROCs suggesting that these features may individually provide sufficient information
218 for distinguishing between functional and non-functional genomic regions. To assess this
219 possibility, we next evaluated the error rates of functional predictions based on single features.
220 We first considered expression breadth of a sequence, the best predicting feature of functionality.
221 Despite high AUC-ROC (0.95), the false positive rate (FPR; % of pseudogenes predicted as
222 phenotype genes) is 21% when only expression breadth is used, while the false negative rate
223 (FNR; % of phenotype genes predicted as pseudogenes) is 4%. Similarly, the best-performing
224 H3 mark- and sequence conservation-related features have FPRs of 26% and 32%, respectively,
225 while also incorrectly classifying at least 10% of phenotype genes as pseudogenes. Thus, even
226 when considering well-performing single features, error rates remain high indicating the need to
227 jointly consider multiple features for distinguishing phenotype genes and pseudogenes.

228 **Consideration of multiple features in combination produces accurate predictions of**
229 **functional genomic regions**

230 To consider multiple features in combination, we first conducted principle component (PC)
231 analysis to investigate how well phenotype genes and pseudogenes could be separated. Between
232 the first two PCs, which jointly explain 40% of the variance in the feature dataset, phenotype
233 genes (**Fig. 3A**) and pseudogenes (**Fig. 3B**) are distributed in largely distinct space. However,
234 there remains substantial overlap, indicating that standard parametric approaches are not well
235 suited to distinguishing between benchmark functional and non-functional sequences. Thus, we
236 instead considered all 93 features in combination using random forest (see Methods), which
237 generated a binary classifier that can be used to predict whether a sequence resembled phenotype
238 genes or pseudogenes. This classifier is referred to as the full model. The phenotype gene and
239 pseudogene sequences and associated conservation, biochemical, and sequence-structure features
240 were separated into distinct training and testing sets such that the full model was generated and
241 validated using independent data subsets (cross-validation). The resulting full prediction model
242 provided much more accurate predictions (AUC-ROC=0.98; FNR=4%; FPR=10%; **Fig. 3C**)
243 compared to any individual feature (**Fig. 2**). An additional measure of performance based on the
244 precision (proportion of predicted functional sequences that are truly functional) and recall
245 (proportion of functional sequences predicted as functional) values among predictions generated
246 by the full model (**Fig. 3D**) also indicate that the model is performing well. When compared to
247 the best-performing single feature (expression breadth), the full model has a similar FNR but
248 only half the FPR (10% compared to 21%). Thus, the full model is more capable of
249 distinguishing between phenotype genes and pseudogenes.

250 We next determined what the relative contributions of different feature categories were in
251 predicting phenotype genes and pseudogenes and whether models based on a subset of features
252 would perform similarly as the full model. Seven prediction models were established, each using
253 only the subset of features from a single category (**Fig. 2**). Although none of these category-
254 specific models had performance as high as the full model, the models based on transcription
255 activity, sequence conservation, and H3 mark features scored highly (AUC-ROC=0.97, 0.92, and
256 0.91, respectively; **Fig. 3C**). Particularly, the transcription activity feature category model
257 performed almost as well as the full model (FNR=6%, FPR=12%). We should emphasize that,
258 instead of the presence of expression evidence, other transcription activity-related features are
259 significantly distinct between functional and non-functional regions that produce useful
260 predictions.

261 Considering that investigating the functionality of ITRs is a primary goal of this study
262 and that ITRs are defined based on the presence of expression evidence, we also built a model
263 did not consider any transcription activity features (full w/o TX, **Fig. 3C-D**). We found that the
264 model excluding transcription activity features performed almost as well as the full model and
265 similarly to the transcription activity-feature-only model although with an increased FPR (AUC-
266 ROC=0.96; FNR=3%; FPR=20%). This indicates that predictions of functional regions are not
267 reliant solely on transcription data, but instead a diverse array of features can be considered to
268 make highly accurate predictions of the functionality of a genomic sequence. Meanwhile, our
269 finding of the high performance of the transcription activity-only model highlights the possibility
270 of establishing an accurate model for distinguishing functional genic and non-functional genomic
271 sequences in plant species with only a modest amount of transcriptome data.

272 **Functional likelihood allows the prediction of functional and non-functional genomic** 273 **regions**

274 To provide a measure of the potential functionality of any sequence, including ITRs and
275 ncRNAs, in the *A. thaliana* genome, we utilized the confidence score from the full model as a
276 “functional likelihood” value (Tsai et al. 2017; see Methods). The functional likelihood score
277 ranges between 0 and 1, with high values indicating that a sequence is more similar to phenotype
278 genes (functional) and low values indicating a sequence more closely resemble pseudogenes
279 (non-functional). Functional likelihood values for all genomic regions examined in this study are
280 available in Supplemental Table 4. As expected, phenotype genes have high functional
281 likelihood values (median=0.97; **Fig. 4A**) and pseudogenes have low values (median=0.01; **Fig.**
282 **4B**). To call sequences as functional or not, we defined a threshold functional likelihood value of
283 0.35 (see Methods). Using this threshold, 96% of phenotype genes (**Fig. 4A**) and 90% of
284 pseudogenes (**Fig. 4B**) are correctly classified as functional and non-functional, respectively,
285 demonstrating that the full model is highly capable of distinguishing functional and non-
286 functional sequences.

287 We next applied our model to predict the functionality of annotated protein-coding genes,
288 transposable elements, and random unexpressed intergenic regions. Most annotated protein-
289 coding genes not included in the phenotype gene dataset have high functional likelihood scores
290 (median=0.86; **Fig. 4C**) and 80% are predicted as functional. Of the 20% of protein-coding

291 genes that were predicted as non-functional, we expect that at least 4% represent false negatives
292 based on the FNR of the full model. The actual FNR among protein-coding genes may be higher,
293 however, as phenotype genes represent a highly active and well conserved subset of all genes.
294 However, a subset of the low-scoring protein-coding genes may also represent gene sequences
295 undergoing functional decay and *en route* to pseudogene status. To assess this possibility, we
296 examined 1,940 *A. thaliana* "decaying" genes that may be experiencing pseudogenization due to
297 promoter disablement (Yang et al. 2011) and found that while they represent only 7% of all *A.*
298 *thaliana* annotated protein-coding genes, they make up 45% of protein-coding genes predicted as
299 non-functional (Fisher's Exact Test (FET), $p < 1E-11$). In addition to protein-coding genes, we
300 evaluated the functional likelihoods of transposable elements (TEs) and randomly-selected,
301 unexpressed intergenic regions that are most likely non-functional. As expected, the functional
302 likelihoods were low for both TEs (median=0.03, **Fig. 4D**) and unexpressed intergenic regions
303 (median=0.07; **Fig. 4E**), and 99% of TEs and all unexpressed intergenic sequences were
304 predicted as non-functional, further demonstrating the utility of the function prediction model.
305 Overall, the functional likelihood measure provides a useful metric to distinguish between
306 phenotype genes and pseudogenes. In addition, the functional likelihoods of annotated protein-
307 coding genes, TEs, and unexpressed intergenic sequences agree with *a priori* expectations
308 regarding the functionality of these sequences.

309 **Exclusion of features from multiple tissues increases prediction performance for narrowly-** 310 **expressed sequences**

311 Although the full model performs exceedingly well, there remain false predictions. There are 76
312 phenotype genes (4%) predicted as non-functional (referred to as low-FL phenotype genes). We
313 assessed why these phenotype genes were not correctly identified by first asking what category
314 of features were particularly distinct between low-FL and the remaining phenotype genes. We
315 found that the major category that led to the misclassification of phenotype genes was
316 transcription activity, as only 7% of low-scoring phenotype genes were predicted as functional in
317 the transcription activity-only model, compared to 98% of high FL phenotype genes (**Fig. 5A**).
318 By contrast, >65% of low-FL phenotype genes were predicted as functional when sequence
319 conservation, H3 mark, or DNA methylation features were used. This could suggest that the full
320 model is less effective in predicting functional sequences that are weakly or narrowly expressed.

321 While sequence conservation features are distinct between functional and non-functional
322 sequences when considered in combination, a significantly higher proportion of low-FL
323 phenotype genes were specific to the *Brassicaceae* family, with only 33% present in
324 dicotyledonous species outside of the *Brassicaceae*, compared to 78% of high-scoring phenotype
325 genes (FET, $p < 4e-12$), thus our model likely has reduced power in detecting lineage-specific
326 genes.

327 Given the association between transcription activity features and functional predictions,
328 we next investigated how functional predictions performed for conditionally-functional and
329 narrowly-expressed sequences. We found that genes with conditional phenotypes (see Methods)
330 had no significant differences in functional likelihoods (median=0.96) as those with phenotypes
331 under standard growth conditions (median=0.97; U test, $p=0.38$), indicating that our model can
332 capture conditionally functional sequences. Next, we evaluated functional likelihood
333 distributions among sequences with different breadths of gene expression. For this comparison,
334 we focused on non-stress, single-tissue expression datasets (Supplemental Table 5), which was
335 distinct from the expression breadth feature in the prediction model that considered all datasets.
336 While phenotype genes are better predicted than pseudogenes among sequences with the same
337 number of tissues with expression evidence (U tests, all $p < 1.7E-06$; Supplemental Fig. 2A),
338 65% of the 62 phenotype genes expressed in ≤ 3 tissues are predicted as non-functional. Further,
339 there is a significant correlation between the number of tissues with expression evidence and
340 functional likelihood values of all sequences in our analysis ($r^2=0.77$; $p < 2E-16$). Thus, the
341 functional prediction model is biased against narrowly-expressed phenotype genes.

342 We also found that 80 pseudogenes (10%) were defined as functional (high-FL
343 pseudogenes). Consistent with misclassifications among phenotype genes, a key difference
344 between high-FL pseudogenes and those that were correctly predicted as non-functional was that
345 high-FL pseudogenes tend to be highly and broadly expressed (**Fig. 5A**). A significantly higher
346 proportion of high-FL pseudogenes come from existing genome annotation as 19% of annotated
347 pseudogenes were classified as functional, compared to 4% of pseudogenes identified through a
348 computational pipeline (Zou et al 2009) (FET, $p < 1.5E-10$). We found that high-FL pseudogenes
349 might be more recently pseudogenized and thus have not yet lost many genic signatures, as the
350 mean number of disabling mutations (premature stop or frameshift) per kb in high-scoring
351 pseudogenes (1.9) was significantly lower than that of low-scoring pseudogenes (4.0; U test, $p <$

352 0.02). Lastly, we cannot rule out the possibility that a small subset of high-scoring pseudogenes
353 represent truly functional sequences, rather than false positives (e.g. Karreth et al. 2015; Poliseno
354 et al. 2010). Overall, the misclassification of both narrowly-expressed phenotype genes and
355 broadly-expressed pseudogenes highlights the need for an updated prediction model that is less
356 influenced by expression breadth.

357 To tailor functional predictions to narrowly-expressed sequences, we generated a “tissue-
358 agnostic” model that attempts to minimize the contribution of biochemical activities occurring in
359 many tissues by excluding expression breadth and features that were available across multiple
360 tissues (see Methods). The tissue-agnostic model performed similarly to the full model (AUC-
361 ROC=0.97; FNR=4%; FPR=15%; Supplemental Fig. 3; Supplemental Table 4). Importantly, the
362 proportion of phenotype genes expressed in ≤ 3 tissues predicted as functional increased by 23%
363 (35% in the full model to 58% in the tissue-agnostic model, Supplemental Fig. 2B), indicating
364 that the tissue-agnostic model is more suitable for predicting the functionality of narrowly-
365 expressed sequences than the full model, although there is an increase in FPR (from 10% to
366 15%). We next sought to evaluate the functional likelihood of ITR and annotated ncRNA
367 sequences utilizing both the full model and the tissue-agnostic model, in case that these
368 sequences are narrowly-expressed.

369 **Intergenic transcribed regions and annotated ncRNAs are mostly predicted as non-** 370 **functional**

371 ITRs and ncRNAs represent functionally ambiguous sequences, as they are usually identified by
372 the presence of expression evidence and few have been functionally characterized. Nevertheless,
373 a subset of ITRs likely represent novel genes and may also represent unannotated exon
374 extensions of known genes (Johnson et al. 2005). To evaluate the functionality of ITRs and
375 ncRNAs, we next applied both the full and tissue-agnostic models to these sequences.

376 Additionally, we investigated whether likely-functional ITRs and ncRNAs are close to annotated
377 genes, and if so, if they may be extensions of the gene neighbors. We assessed functional
378 likelihood values for 895 ITRs from three sources: Araport 11 annotation, Moghe et al. (2013),
379 and an additional set identified in this study from 206 RNA-seq datasets. We also analyzed the
380 functional likelihood of TAIR ncRNAs (n=136), and Araport long ncRNAs (referred to as
381 Araport ncRNAs, n=252) TAIR and Araport ncRNAs are collectively referred to as annotated

382 ncRNAs. The functional likelihoods based on the full model were low (median=0.09) for both
383 ITRs (**Fig. 4F**) and Araport ncRNAs (**Fig. 4G**), and only 15% and 9% of these sequences are
384 predicted as functional, respectively. By contrast, TAIR ncRNAs have higher functional
385 likelihood values (median=0.53; **Fig. 4H**) and 68% are predicted as functional. We next asked
386 what features were distinct among TAIR ncRNAs compared to ITRs and Araport ncRNAs that
387 led to a greater proportion of these sequences predicted as functional and found that transcription
388 activity features of TAIR ncRNAs are more similar to phenotype genes when compared to ITRs
389 and Araport ncRNAs (**Fig. 5B**). By contrast, only 40% of TAIR ncRNAs are predicted as
390 functional if sequence conservation features are considered, potentially because RNA genes
391 experience less selective constraint at the primary sequence level compared to protein-coding
392 genes (Pang et al. 2006). When looking at the performance of single-category predictions, we
393 also find that a greater proportion of ITRs and Araport ncRNAs are predicted as functional when
394 considering only DNA methylation or H3 mark features (**Fig. 5B**). However, these two category-
395 specific models are also marked by increased false positive rates and predict a substantial
396 proportion of unexpressed intergenic sequences as functional (**Fig. 5B**). Notably, 88% of
397 unexpressed intergenic sequences are predicted as functional based on the DNA methylation-
398 only model. Thus, while single-category models are useful for determining features that are
399 similar or dissimilar across sequences types, they may not be useful as a basis for predicting
400 sequences as functional or non-functional.

401 As ITRs and annotated ncRNAs are generally narrowly-expressed, it is likely that we are
402 underestimating the proportion that is functional. We next applied the tissue-agnostic model to
403 ITRs and annotated ncRNAs, as this model is less biased against narrowly-expressed sequences
404 (Supplemental Fig. 2B). Compared to the full model, twice as many ITRs (30% compared to
405 15% in the full model; FET, $p < 4E-15$) and Araport ncRNAs (19% compared to 9%; FET, $p <$
406 0.003) are predicted as functional. A similar proportion of TAIR ncRNAs are predicted as
407 functional (67% compared to 68%; FET, $p=0.80$), which is likely a result of TAIR ncRNAs
408 being more broadly expressed than ITRs and Araport ncRNAs (Supplemental Fig. 4A).
409 Considering both the full and tissue-agnostic models, we predict a total of 268 ITRs (32%), 57
410 Araport ncRNAs (23%), and 105 TAIR ncRNAs (77%) as functional.

411 Intergenic transcripts can represent evidence for unannotated extensions or alternative
412 splicing variants of known genes (Johnson et al. 2005). Thus, we next evaluated whether ITRs

413 and annotated ncRNAs that are predicted as functional are close to annotated genes and if these
414 sequences share features with neighboring genes. We found that ITRs and annotated ncRNAs
415 closer to annotated genes tend to be predicted as functional (Supplemental Fig. 5A). Using the
416 95th percentile of intron lengths for all genes as a threshold to call ITRs and annotated ncRNAs
417 as proximal or distant to neighboring genes, 57% of functional ITRs and annotated ncRNAs are
418 considered proximal, compared to 35% for non-functional ITRs and annotated ncRNAs (FET, p
419 $< 2E-09$), suggesting that a subset these likely-functional sequences may be unannotated exons
420 of known genes. If ITRs and annotated ncRNAs represent unannotated extensions of known
421 genes, they may share features with their gene neighbors. However, functional ITRs/ncRNAs
422 have features that bear little similarity to neighboring genes, regardless of if they are proximal or
423 distant to neighboring genes (Supplemental Fig. 5B-C). In contrast, genes are generally more
424 similar to their neighbors, regardless of proximity, than ITRs or annotated ncRNAs are to their
425 nearest neighbor (Supplemental Fig. 5B-C). This is also true compared to random gene pairs
426 (Supplemental Fig. 5D). Thus, despite their proximity to annotated genes, we expect that few
427 ITRs or annotated ncRNAs represent unannotated exon extensions of known genes. For proximal
428 functional ITRs/annotated ncRNAs, we cannot rule out the possibility that they represent false-
429 positive functional predictions due to the accessible and active chromatin states of nearby genes
430 that serve as a confounding factor. For the 116 functional ITRs and annotated ncRNAs that are
431 distal, they may represent fragments of novel genes.

432 Overall, we find that ITRs and annotated ncRNAs are generally predicted as non-
433 functional. Furthermore, tissue-specific or conditional functionality does not fully explain these
434 non-functional predictions and few predicted-functional ITRs and ncRNAs are likely
435 unannotated extensions of neighboring genes. In addition to the ITRs and ncRNAs investigated
436 thus far, there are 12,344 ITR and ncRNA sequences that are shorter than 500 bp and were
437 unable to be investigated by the full model. We next evaluated methods to assess the
438 functionality of these shorter sequences.

439 **Short RNA genes have mixed predictions based on a binary classification model**

440 The functional predictions performed thus far require 500 bp of sequence. However, there are an
441 additional 10,938 ITRs and 1,406 annotated ncRNAs (12,344 in total) that are shorter than 500
442 bp. To evaluate the functionality of short ITRs and ncRNAs, we generated a new binary

443 classification model using features calculated from a randomly-selected 100 bp sequence within
444 a gene or pseudogene body (for features, see Supplemental Table 6). ITRs and annotated
445 ncRNAs tend to be more narrowly expressed than phenotype genes (U tests, all $p < 6e-15$;
446 Supplemental Fig. 4B) and the tissue-agnostic model was shown to improve false negative rates
447 among low-FL phenotype genes. Therefore, we generated this model while excluding expression
448 breadth and tissue-specific features (referred to as 100 bp tissue-agnostic model). The 100bp
449 tissue-agnostic model performed similarly to the full 500 bp model in distinguishing between
450 phenotype genes and pseudogenes (AUC-ROC=0.97; FNR=13%; FPR=5%; Supplemental Fig.
451 6). Most importantly, focusing on entries <500 bp in length, this 100 bp model led to the
452 prediction of an additional 366 ITRs (11%), 109 Araport ncRNAs (8%), and 10 TAIR ncRNAs
453 (44%) as functional (Supplemental Fig. 6F-H).

454 In addition to allowing the evaluation of 12,344 short ITRs and annotated ncRNAs, the
455 100 bp tissue-agnostic model can be applied to annotated short RNA genes. Thus, we next
456 sought to evaluate functional likelihood scores for Pol II-transcribed RNA genes that have been
457 annotated in TAIR10, including the primary transcripts of microRNAs (miRNAs; n=151), small
458 nucleolar RNAs (snoRNAs; n=15), and small nuclear RNAs (snRNAs; n=6). We found that 15%
459 of miRNAs (Supplemental Fig. 6I), 73% of snoRNAs (Supplemental Fig. 5J), and 50% of
460 snRNAs (Supplemental Fig. 6K) were predicted as functional. Because most TAIR10 annotated
461 RNA genes are computationally predicted and have not been experimentally validated, it is
462 possible that some may represent false positive gene annotations, particularly among miRNA
463 entries. Meanwhile, we cannot rule out the possibility that the 100bp tissue-agnostic model
464 performs sub-optimally for RNA genes. To further assess these possibilities, we identified six
465 RNA genes (four miRNAs, one lncRNA, and one *trans*-acting small interfering RNA) with loss-
466 of-function mutant phenotypes (referred to as RNA phenotype genes; Supplemental Table 7). Of
467 these six genes, we correctly identify three as functional (Supplemental Fig. 6L). Although this is
468 significantly higher than the proportion of pseudogenes (FET, $p < 0.004$) and miRNAs ($p = 0.05$)
469 predicted as functional, this finding suggests that the 100 bp tissue-agnostic model has a
470 substantial false negative rate for detecting functional RNA genes. One immediate question is
471 whether the suboptimal prediction is because RNA genes belong to a class of their own. To
472 further evaluate functional predictions of RNA gene sequences, TAIR ncRNAs, Araport

473 ncRNAs, and ITRs, we next built multi-class functional prediction models for distinguishing
474 RNA genes from other types of functional and non-functional sequences.

475 **Intergenic transcribed regions and annotated ncRNAs do not resemble benchmark RNA** 476 **genes**

477 To build a model that considers genomic sequences that are likely functional at the RNA level as
478 a distinct class, we generated a four-class function prediction model aimed at distinguishing four
479 classes of sequences: benchmark RNA genes, phenotype protein-coding genes (same as
480 phenotype genes from previous sections), pseudogenes, and randomly-selected, unexpressed
481 intergenic regions. Here, unexpressed intergenic sequences were included to provide another set
482 of likely non-functional sequences distinct from pseudogenes. The benchmark RNA gene
483 training set was composed of six RNA phenotype genes discussed in the previous section and 40
484 high-confidence primary miRNA sequences from miRBase (Kozomara and Griffiths-Jones
485 2014). The model provides four scores, one for each sequence class (for scores, see
486 Supplemental Table 4), and the maximum score was used to classify sequences. We excluded
487 expression breadth and tissue-specific features when generating the four-class model.

488 Based on predictions from the four-class model, the RNA gene training set was well-
489 classified, with 87% predicted as either RNA gene-like (65%) or phenotype protein-coding gene-
490 like (22%; **Fig. 6A**). Notably, all six RNA phenotype genes were predicted as functional (four
491 and two predicted as RNA genes and phenotype protein-coding genes, respectively). To assess
492 whether sequences predicted as RNA gene-like had evidence of translation, we identified
493 genomic regions with translation evidence based on two shotgun proteomics datasets
494 (Baerenfaller et al. 2008; Castellana et al. 2008). We find that phenotype protein-coding genes
495 and other protein-coding genes predicted as benchmark RNA gene-like are less likely to have
496 evidence of translation compared to those predicted as phenotype protein-coding gene-like (FET,
497 both $p < 6e-5$, Supplemental Fig. 7). Taken together with the predictions of benchmark RNA
498 genes, these results suggest that the benchmark RNA gene prediction score allows sequences that
499 function at the RNA level to be distinguished from other sequence types. For the remaining three
500 classes in the four-class model, 95% of phenotype genes were predicted as either phenotype
501 protein-coding gene-like or benchmark RNA gene-like (**Fig. 6B**), while 70% of pseudogenes
502 (**Fig. 6C**) and 100% of unexpressed intergenic regions (**Fig. 6D**) resembled either pseudogenes

503 or unexpressed intergenic sequences. Importantly, among phenotype genes expressed in ≤ 3
504 tissues, 80% were correctly predicted as phenotype protein-coding or benchmark RNA gene-like
505 in the four-class model, an increase of 22% over the 500 bp tissue-agnostic model.

506 Since the four-class model was generally able to distinguish benchmark RNA genes from
507 other sequence classes, regardless of breadth of expression, we next evaluated whether ITRs and
508 annotated ncRNAs resemble benchmark RNA genes. We find that 20%, 19%, and 15% of ITRs,
509 Araport ncRNAs, and TAIR ncRNAs, respectively, are predicted as RNA genes (**Fig. 6E-G**). We
510 also considered that ITRs and annotated ncRNAs that were predicted as phenotype protein-
511 coding gene-like may also be functioning at the RNA level. Consistent with this notion, fewer
512 than 5% of phenotype protein-coding gene-like ITRs and annotated ncRNAs have evidence of
513 translation, compared to 37% of phenotype genes and 27% of protein-coding genes
514 (Supplemental Fig. 7). This suggests that the majority of ITRs and annotated ncRNAs predicted
515 as benchmark RNA gene-like or phenotype protein-coding gene-like are likely functional RNA
516 genes.

517 To provide an overall estimate the proportion of likely-functional ITRs and annotated
518 ncRNAs, we considered the outcome of all four models presented in this study (full 500 bp, 500
519 bp and 100 bp tissue-agnostic, and four-class models) in combination. We classify 2,453 ITRs
520 (21%) and 506 annotated ncRNAs (28%) as functional, as they resemble phenotype protein-
521 coding genes in at least one of the four models. An additional 1,984 ITRs (17%) and 290
522 ncRNAs (16%) resemble benchmark RNA genes and therefore could be functional at the RNA
523 level. Ultimately, we find that the majority of ITRs (62%) and annotated ncRNAs (56%) are
524 predicted as non-functional, suggesting that these sequences do not primarily represent novel
525 protein-coding or RNA genes. Moreover, at least a third of ITRs (**Fig. 6E**) and Araport ncRNAs
526 (**Fig. 6F**) are most similar to unexpressed intergenic regions. Given that these sequences have not
527 been functionally characterized, it is possible that many represent regions of noisy transcription
528 and, in the cases of annotated ncRNAs, false positive gene annotations.

529 CONCLUSION

530 We identify a collection of evolutionary, biochemical, and sequence-structure signatures that
531 represent defining features of functional genic regions in a plant genome. Considering these
532 features jointly via machine learning methods produces highly accurate predictions that can

533 distinguish between functional and non-functional genomic regions with low false positive and
534 false negative rates. Expression evidence is particularly distinct between phenotype genes and
535 pseudogenes. However, it is the level and breadth of expression that is important for predictions
536 as most pseudogenes have evidence of expression. In addition, predictions performed without
537 expression evidence also performed well, indicating that functional regions are not defined solely
538 by expression features. We also identified ITRs occurring across 15 diverse land plant species
539 with a wide range of genome sizes and find that the amount of intergenic expression occurring in
540 a species increases with genome size while the amount of genic expression does not. Considering
541 that noisy expression should be expected to increase with additional genome space, this hints that
542 much of the intergenic transcription occurring in a species may be non-functional.

543 Among the 11,833 ITRs analyzed in this study, we predict 2,453 (21%) are likely
544 functional as they exhibit the biochemical, evolutionary, and sequence-structure characteristics
545 of known functional genomic regions. For annotated ncRNA regions, we classify 506 of 1,794
546 (28%) as likely-functional. An additional 1,984 ITRs (17%) and 290 ncRNAs (16%) resemble
547 benchmark RNA genes and therefore could be functional at the RNA level. However, the false
548 positive rate among RNA gene predictions could be quite high, as 15% of pseudogenes were
549 predicted as RNA genes. More robust and reliable predictions would be possible if additional
550 benchmark RNA genes with loss-of-function phenotype information were available. Ultimately,
551 the ITRs and annotated ncRNAs that are predicted as functional are likely-genic regions that
552 could be responsible for biological novelties and represent an important component of the
553 functional gene set in *A. thaliana*. Therefore, they should be considered high priority targets in
554 future experimental studies. However, the remaining 7,396 ITRs (63%) and 998 annotated
555 ncRNAs (56%) are most similar to pseudogenes or unexpressed intergenic sequences, suggesting
556 these sequences are likely non-functional and byproducts of transcriptional noise. Given that the
557 majority of ITRs and annotated ncRNAs are predicted as non-functional, we recommend that the
558 null hypothesis for the functionality of expressed intergenic sequences is that they represent
559 transcriptional noise. We do not suggest that all novel intergenic transcription represents non-
560 functional activity, but instead that ITRs should be generally regarded as non-functional until
561 convincing experimental evidence is provided that a transcribed genomic region is functional.

562 **METHODS**

563 **Identification of leaf intergenic transcribed regions**

564 RNA-sequencing (RNA-seq) datasets were retrieved from the Sequence Read Archive (SRA) at
565 the National Center for Biotechnology Information (NCBI) for 15 flowering plant species
566 (Supplemental Table 1). All datasets were generated from leaf tissue and sequenced on Illumina
567 HiSeq 2000 or 2500 platforms. Genome sequences and gene annotation files were downloaded
568 from Phytozome v11 (www.phytozome.net; Goodstein et al. 2011) or Oropetium Base v01
569 (www.sviridis.org; VanBuren et al. 2015). Genome sequences were repeat masked using
570 RepeatMasker v4.0.5 (www.repeatmasker.org) if a repeat-masked version of a genome assembly
571 was not available. Only one end from paired-end read datasets were utilized in downstream
572 processing. Reads were trimmed of low scoring ends and residual adaptor sequences using
573 Trimmomatic v0.33 (Bolger et al. 2014) and mapped to associated genome sequences using
574 Tophat v2.0.13 (Kim et al. 2013). Reads ≥ 20 nucleotides in length that mapped uniquely within a
575 genome at our mapping threshold were used in further analysis. Thirty million mapped reads
576 were randomly selected from among all datasets for a species and assembled into transcript
577 fragments using Cufflinks v2.2.1 (Trapnell et al. 2010). The expected mean fragment length for
578 assembled transcript fragments in Cufflinks was set to 150 from the default of 200 so that
579 expression levels in short fragments would not be overvalued. The 1st and 99th percentile of
580 intron lengths in a given gene annotation set were used as the minimum and maximum intron
581 lengths, respectively, for both the TopHat2 and Cufflinks steps. Intergenic transcribed regions
582 (ITRs) were defined by transcript fragments that did not overlap existing gene annotation and did
583 not have significant six-frame translated sequence similarity to annotated plant proteins in
584 Phytozome v10 (BLASTX E-value < 1E-05). To determine the relationship between genome
585 size and number of annotated genes, we calculated the correlation between assembled genome
586 size and gene counts from the first 50 published plant genomes as described by Michael and
587 Jackson (2013).

588 ***Arabidopsis thaliana* genome annotation**

589 *Arabidopsis thaliana* protein-coding gene, miRNA gene, snoRNA gene, snRNA gene, ncRNA
590 region, pseudogene, and transposable element annotations were retrieved from The Arabidopsis
591 Information Resource v10 (TAIR10; www.arabidopsis.org; Berardini et al. 2015). Additional
592 miRNA gene and lncRNA region annotations were retrieved from Araport v11

593 (www.araport.org; Krishnakumar et al. 2015). A pseudogene-finding pipeline similar to that
594 described by Chen et al. (Zou et al. 2009) was used to identify additional putative pseudogene
595 fragments and count the number of disabling mutations (early stop or frameshift mutations)
596 present in these sequences. To avoid potential confounding effects from overlapping gene
597 annotation, protein-coding and RNA gene annotation that overlapped other gene or pseudogene
598 annotation were excluded from further analysis, except for lncRNA annotation that overlapped
599 with other lncRNAs, which were merged. Pseudogenes and transposable elements that
600 overlapped genic regions were also removed. When pseudogenes from TAIR10 and the
601 pseudogene-finding pipeline overlapped, the longer pseudogene annotation was retained.

602 ITRs were defined by Moghe et al. (2013; “Set 2” ITRs; coordinates provided by the
603 authors) and Araport v11 (described as “novel transcribed regions”). Overlapping ITR
604 annotations from Araport were merged. Additional ITRs were identified from 206 RNA-seq
605 datasets generated using wild-type, Columbia-0 tissue on Illumina sequencing platforms
606 (Supplemental Table 5). Datasets were identified by querying NCBI-SRA for datasets from *A.*
607 *thaliana* with RNA as the source. Reads were trimmed, mapped, and assembled into transcript
608 fragments using the steps described in the previous section, except that reads from multiple
609 datasets were not merged and subsampled. Instead, overlapping assembled transcript fragments
610 from across datasets were merged. ITRs were identified by transcribed fragments that did not
611 overlap with any annotated feature from TAIR10 or Araport11 or any pseudogenes defined by
612 the pseudogene-finding pipeline. Overlaps among ITR annotations were resolved using a priority
613 system: Araport11 > Moghe et al. > ITRs identified in this study.

614 For each gene, ncRNA, pseudogene, transposable element, and intergenic transcribed
615 sequence, a randomly-selected 100 and 500 base pair (bp) window was chosen for feature
616 calculation (Supplemental Table 2; Supplemental Table 6; see below for feature descriptions).
617 Sequences that were not at least 100 or 500 bp in length were excluded. This controlled for
618 effects of sequence length and simplified gene structure considerations (e.g. exon/intron
619 boundaries). In addition, random 100 bp (n=4,000) and 500 bp (n=3,716) regions of intergenic
620 space (genome regions outside of gene, pseudogene, or transposable element annotation) that did
621 not overlap with any genic or intergenic transcript fragments were also selected for feature
622 calculation. These 100 and 500 bp windows in gene, ncRNA, pseudogene, transposable element,

623 and ITR annotation and unexpressed intergenic space are referred to as “feature regions”
624 throughout the Methods section.

625 **Single-feature prediction performance**

626 The ability for single features to distinguish between functional and non-functional regions was
627 tested using Area Under the Curve - Receiver Operating Characteristic (AUC-ROC) values
628 calculated using the scikit-learn package in Python. AUC-ROC values range between 0.5
629 (equivalent to random guessing) and 1 (perfect predictions) and values above 0.7, 0.8, and 0.9
630 are considered to be fair, good, and excellent, respectively. Thresholds to predict sequences as
631 functional or non-functional using a single feature were defined by the feature value that
632 produced the highest F-measure (harmonic mean of precision and recall), which gives
633 consideration to both false positives and false negatives at a given threshold. False positive rates
634 (FPR) were calculated as the percentage of negative cases with values above or equal to the
635 threshold and false negative rates (FNR) were calculated was the percentage of positive cases
636 with values below the threshold.

637 **Phenotype data sources**

638 Mutant phenotype data for *Arabidopsis thaliana* protein-coding genes was collected from a
639 published dataset (Lloyd and Meinke 2012), the Chloroplast 2010 database (Ajjawi et al. 2010;
640 Savage et al. 2013), and the RIKEN Phenome database (Kuromori et al. 2006) as described by
641 Lloyd et al. (2015). Phenotype genes used in our analyses were those whose disruption resulted
642 in lethal or visible defects under standard laboratory growth conditions (i.e. non-stress
643 conditions). Genes with documented mutant phenotypes under standard conditions were
644 considered as a distinct and non-overlapping category from other annotated protein-coding
645 genes. We identified six RNA genes with documented loss-of-function phenotypes through
646 literature searches: *At4* (AT5G03545; Shin et al. 2006), *MIR164A* and *MIR164D* (AT2G47585
647 and AT5G01747, respectively; Guo et al. 2005), *MIR168A* (AT4G19395; Li et al. 2012b), and
648 *MIR828A* and *TAS4* (AT4G27765 and AT3G25795, respectively; Hsieh et al. 2009). An
649 additional 23 RNA genes with documented overexpression mutant phenotypes were identified
650 from the literature links at miRBase. Conditional phenotype genes were those belonging to the
651 “Conditional” phenotype class as described by Lloyd and Meinke (2012). These genes had no

652 obvious mutant phenotype under standard growth conditions, but did exhibit a loss-of-function
653 phenotype under stress conditions. These were compared with phenotype genes belonging to the
654 “Morphological” phenotype class from the same study, which have visible growth defects under
655 standard growth conditions.

656 **Sequence conservation and structure features**

657 Nucleotide diversity and Tajima’s D were calculated among 81 *A. thaliana* accessions (Cao et al.
658 2011) for each feature region using custom Python scripts. The genome matrix file for the Cao et
659 al. study was retrieved from the 1,001 genomes database (www.1001genomes.org) and analyzed
660 with Python scripts available through GitHub (github.com/panchyni/GenomeMatrixProcessing).
661 The genomic regions that align between *A. thaliana* and six other plant species were retrieved
662 from Li et al. (2012a). The coverage of each feature region with these aligned blocks was
663 calculated. In addition, phastCons conservation scores were available for each nucleotide within
664 an aligned block. The maximum and average of phastCons scores were calculated for each
665 feature region. Nucleotides in a feature region that did not overlap with an aligned block were
666 assigned a phastCons score of 0. BLASTN searches were performed between feature region
667 nucleotide sequences and Phytozome v10 genome sequences. Five plant lineages were
668 considered: *Brassicaceae* ($n_{\text{species}}=7$), other dicotyledonous plants ($n=22$), monocots ($n=7$), other
669 embryophyte plants ($n=3$), and algae ($n=5$). The percent identity to the most significant match by
670 E-value (maximum E-value: $1E-05$) within a lineage group for each feature region was used as
671 the feature in functional predictions. DNA sequence-structure features consisted of the first five
672 principal components of the 125 conformational and thermodynamic dinucleotide properties
673 collected from DiProDB database (Friedel et al. 2009). The first five principal components (83%
674 of variation) correspond primarily to DNA major groove geometry, free energy, twist and roll,
675 DNA minor groove geometry, and tilt and rise, respectively (Tsai et al. 2015). Sequence-
676 structure values corresponding to principal components were calculated in dinucleotide windows
677 and averaged across the length of a feature region.

678 **Transcription activity features**

679 From the 206 *A. thaliana* RNA-seq datasets described above, we removed datasets with fewer
680 than 20 million reads ($n=134$) or abnormally high RPKM distributions among resulting transcript

681 fragments (n=21; median of median RPKM values among retained and removed datasets=10.2
682 and 4065.2, respectively; Supplemental Table 5), which indicated technical issues during the
683 read cleaning, read mapping, or transcript assembly processes. Transcript fragments assembled
684 from the remaining 51 RNA-seq datasets were used to calculate expression breadth across
685 datasets, 95th percentile RPKM expression levels, maximum transcription coverage in a single
686 dataset, and presence or absence of expression evidence. Ten datasets from diverse tissues and
687 conditions with a high number of reads were chosen to calculate max RPKM expression levels
688 and transcription coverage as single-dataset features. The tissues and conditions included: pollen
689 (SRR847501), light- and dark-grown seedlings (SRR1020621 and SRR974751, respectively),
690 leaf tissue under standard, drought, and fungal-infection conditions (SRR953400, SRR921316,
691 and SRR391052, respectively), root (SRR578947), inflorescence (SRR953399), flower
692 (SRR505745), and silique (SRR953401). RNA-seq datasets generated from a single tissue and in
693 standard growth conditions were used for tissue-specific expression analysis. The seven tissues
694 were pollen, seedling, leaf, root, inflorescence, flower, and silique (Supplemental Table 5). Two
695 additional datasets generated by sequencing RNA molecules associated with ribosomes
696 (SRR966480 and SRR966484) were retrieved from NCBI-SRA and processed using the same
697 steps as those used on other RNA-seq datasets.

698 **Histone 3 mark features**

699 Chromatin immunoprecipitation sequencing (ChIP-seq) datasets for four activation-associated
700 (H3K4me1: SRR2001269, H3K4me3: SRR1964977, H3K9ac: SRR1964985, and H3K23ac:
701 SRR1005405) and four repression-associated (H3K9me1: SRR1005422, H3K9me2:
702 SRR493052, H3K27me3: SRR3087685, and H3T3ph: SRR2001289) histone 3 (H3) marks were
703 retrieved from NCBI-SRA. Datasets were chosen due to high number of reads and presence of
704 histone 3 or total protein controls. Reads were trimmed with Trimmomatic v0.33 (Bolger et al.
705 2014) and mapped to the TAIR10 genome sequence with Bowtie v2.2.5 (Langmead et al. 2009).
706 H3 mark peaks were identified with the Spatial Clustering for Identification of ChIP-Enriched
707 Regions (SICER) software v1.1 (Xu et al. 2014). SICER requires an effective genome size input,
708 which was calculated according to Koehler et al. (2011). The maximum H3 mark peak intensity
709 and the coverage with each H3 mark peak were calculated for each feature region. The count and
710 coverage of all activating or repressing marks in a feature region were also calculated.

711 **DNA methylation features**

712 Bisulfite-sequencing (BS-seq) datasets from seven tissues (pollen: SRR516176, embryo:
713 SRR1039895, endosperm: SRR1039896, seedling: SRR520367, leaf: SRR1264996, root:
714 SRR1188584, and inflorescence: SRR2155684) were retrieved from NCBI-SRA. BS-seq reads
715 were trimmed with Trimmomatic v0.33 (Bolger et al. 2014) and processed with Bismark v3
716 (Krueger and Andrews 2011). A cytosine was considered to be methylated if at least five reads
717 mapped to the position and >50% of the reads indicated the position was methylated. For each
718 feature region, the percentage of methylated cytosines in CpG, CHG, and CHH contexts were
719 calculated if the feature region had ≥ 5 cytosines with ≥ 5 reads mapping to the position. To test
720 the false positive rate of DNA methylation calls, we evaluated the proportion of cytosines in the
721 chloroplast genome that are called as methylated, as the chloroplast genome is known to have
722 few DNA methylation events (Ngernprasirtsiri et al. 1988; Zhang et al. 2006). In any nucleotide
723 context for any BS-seq dataset, 0-1.5% (median=0) of cytosines in the chloroplast genome were
724 defined as methylated and only 0.1-2.4% of reads suggested that a cytosine position was
725 methylated. This indicated that the false positive rates for DNA methylation calls were low.

726 **Chromatin accessibility and transcription factor binding features**

727 Chromatin accessibility features consisted of DNase I hypersensitive peaks and micrococcal
728 nuclease sequencing (MNase-seq)-derived nucleosome occupancy. DNase I peaks from five
729 tissues (seed coat, seedling, root, unopened flowers, and opened flowers) were available from the
730 Gene Expression Omnibus (experiment identifiers: GSE53322 and GSE53324; Sullivan et al.
731 2014). The max DNase I peak intensity and coverage with DNase I peaks were calculated for
732 each feature region. MNase-seq nucleosome occupancy was produced by Liu et al. (2015). The
733 authors provided the normalized nucleosome occupancy for each nucleotide in the TAIR10
734 genome sequence. The average of nucleosome occupancy values across a feature region was
735 calculated. Transcription factor (TF) binding sites were identified from *in vitro* DNA affinity
736 purification sequencing data of 529 TFs (O'Malley et al. 2016). The total number of TF binding
737 sites and the number of distinct TFs bound were calculated for each feature region.

738 **Machine learning approach**

739 For two-class models using 500 bp sequences (positive class: phenotype genes, n=1,876;
740 negative class: pseudogenes, n=763), the random forest (RF) implementation in the Waikato
741 Environment for Knowledge Analysis software (WEKA; Hall et al. 2009) was utilized. We
742 generated 100 datasets with an equal proportion of phenotype genes and pseudogenes by
743 randomly selecting 763 phenotype genes and pairing them with all 763 pseudogene examples.
744 For each of these 100 datasets, 10-fold stratified cross-validation was utilized during model
745 building and testing. Therefore, training and testing of each model was performed on
746 independent datasets. The median score from the 100 prediction models was used as the final
747 functional prediction score (“functional likelihood”). Five hundred trees using 2, 4, 6, 9, 15, 20,
748 and 25 randomly-selected features were built using the RF algorithm. Fifteen features provided
749 the highest performance, as determined by AUC-ROC (calculated and visualized using the
750 ROCR package; Sing et al. 2005). The same methods were used to test two-class RF models
751 using 100 bp sequences (phenotype genes, n=1,882; pseudogenes, n=3,916), except that 100
752 datasets with equal proportions of phenotype genes and pseudogenes were generated by
753 randomly-selecting 1,882 pseudogenes to pair with all 1,882 phenotype gene examples. In
754 single-category predictions, fewer features were considered in parameter searches. For the H3
755 mark, DNA methylation, and transcription activity categories 2, 4, 7, and 10 features were tested.
756 For the chromatin accessibility and sequence conservation categories 2, 4, and 6 features were
757 tested. For the sequence-structure category 2, 3, and 4 features were tested. For the transcription
758 factor binding category 1 and 2 features were tested. A tissue-agnostic model was generated by
759 excluding the expression breadth feature and all features from tissue-specific RNA-seq, BS-seq,
760 and DNase I hypersensitivity datasets. Tissue-specific features were replaced with the maximum
761 FPKM and coverage from RNA-seq datasets, minimum DNA methylation proportion from any
762 one tissue in CpG, CHG, and CHH contexts, and maximum intensity and coverage with DNaseI
763 peaks in a single tissue.

764 The functional likelihood of a genomic sequence was calculated as the proportion of the
765 500 random forest trees that predicted a sequence as similar to a phenotype gene (Supplemental
766 Table 4). The functional likelihood threshold to predict a sequence as functional or non-
767 functional was defined based on the functional likelihood value that produced the maximum F-
768 measure among all possible thresholds. F-measure is the harmonic mean of precision (proportion
769 of predicted positive regions that are truly positive) and recall (proportion of truly positive

770 regions that are predicted as positive), which gives consideration to both false positives and false
771 negatives. FPR was calculated as the percentage of pseudogenes with functional likelihood
772 values above or equal to the functional threshold, while FNR was calculated as the percentage of
773 phenotype genes with functional likelihood values below the threshold. Functional prediction
774 models were also built using the Sequential Minimal Optimization - Support Vector Machine
775 (SMO-SVM) implementation in WEKA while considering a series of complexity constant
776 parameters: 0.01, 0.1, 0.5 (best by AUC-ROC), 1, 1.5, and 2.0. The results of SMO-SVM models
777 were highly similar to the RF results: PCC between RF and SMO-SVM=0.97; AUC-ROC of
778 SMO-SVM=0.97; FPR=12%; FNR=3%.

779 For the four-class model, phenotype gene, pseudogene, random unexpressed intergenic
780 sequences, and RNA training genes were used as training classes. RNA training genes consisted
781 of six RNA genes with documented loss-of-function phenotypes and high-confidence miRNA
782 genes from miRBase (www.mirbase.org; Kozomara and Griffiths-Jones 2014) Random-sampling
783 of the more populated classes in training cases was used to produce 250 datasets with equal
784 proportions of phenotype genes, pseudogenes, intergenic sequences, and RNA training genes.
785 Two-fold stratified cross-validation was utilized due to the low number of RNA training gene
786 examples. The features described from the tissue-agnostic model above were also used for the
787 four-class model. The random forest implementation in the *party* package of R with conditional
788 inference trees method utilized was used to build the random forest classifiers. The four-class
789 predictions provide prediction scores for each sequence type: a phenotype gene, pseudogene,
790 unexpressed intergenic, and RNA gene score (Supplemental Table 4). The scores indicate the
791 proportion of random forest trees that predict a given sequence as a phenotype gene, pseudogene,
792 unexpressed intergenic, or RNA gene sequence. The median prediction score from across 100
793 equal-proportion runs was used as the final prediction scores, which were then scaled to sum to
794 1. The maximum prediction score was used to classify a sequence as phenotype gene,
795 pseudogene, unexpressed intergenic, or RNA gene.

796 **FIGURE LEGENDS**

797 **Figure 1.** Relationship between genome size and extent of expression in 15 plant species. (A)
798 Amount of expression from annotated gene regions plotted against the size of assembled genome

799 for 15 diverse flowering plant species. The dotted gray line indicates the line of best fit. (B)
800 Amount of expression from intergenic regions plotted against the size of assembled genome.

801 **Figure 2.** Single feature predictions of functional and non-functional sequences. Area Under the
802 Curve - Receiver Operating Characteristic (AUC-ROC) prediction performances using single
803 features in the categories of transcription activity (A), sequence conservation (B), DNA
804 methylation (C), transcription factor binding (D), histone 3 (H3) marks (E), sequence structure
805 (F), and chromatin accessibility (G). AUC-ROC ranges in value from 0.5 (equivalent to random
806 guessing) to 1 (perfect predictions), with values greater than 0.7, 0.8, and 0.9 being considered
807 fair, good, and excellent, respectively. Dotted gray lines indicate the median AUC-ROC within a
808 feature category.

809 **Figure 3.** Multi-feature predictions of functional and non-functional sequences. Smoothed
810 scatterplots of the first two principle components (PCs) of phenotype gene (A) and pseudogene
811 (B) features. The percentages on the axes in (A) indicate the amount of total variation present in
812 the associated PC. (C) Receiver operating characteristic curves of machine learning integration
813 of all features (Full model), all non-transcription activity-related features (Full w/o TX), and
814 when using all features from a single feature category. Single categories are transcription activity
815 (TX), sequence conservation (CV), histone 3 marks (HM), DNA methylation (ME), transcription
816 factor binding (TF), chromatin accessibility (CA), and sequence structure (ST). (C) Precision-
817 recall curves of the models from (B).

818 **Figure 4.** Functional likelihood scores from the full, binary model. Functional likelihood
819 distributions for (A) phenotype genes, (B) pseudogenes, (C) protein-coding genes, (D)
820 transposable elements, (E) random unexpressed intergenic sequences, (F) intergenic transcribed
821 regions, (G) ncRNAs from Araport11, and (H) ncRNAs from TAIR10 from models built using
822 features calculated from 500 bp of sequence. Higher functional likelihood values indicate greater
823 similarity to phenotype genes while lower values indicate similarity to pseudogenes. Vertical
824 dashed lines display the threshold to predict a sequence as functional or non-functional. The
825 numbers to the left and right of the dashed line show the percentage of sequences predicted as
826 functional or non-functional, respectively.

827 **Figure 5.** Functional predictions from single-category predictions. (A) Percentages of phenotype
828 gene and pseudogene sequences predicted as functional (high FL) or non-functional (low FL) in
829 the full model (Full) that are predicted as functional in models based on a subset of features from
830 a single feature category. Single feature categories are transcription activity (TX), sequence
831 conservation (CV), histone 3 marks (HM), DNA methylation (ME), transcription factor binding
832 (TF), chromatin accessibility (CA), and sequence structure (ST). The single category models are
833 sorted from right to left on descending AUC-ROC and separated into informative (all AUC-ROC
834 ≥ 0.87) and uninformative (all AUC-ROC ≤ 0.70) groups. (B) Percentages of sequence classes
835 predicted as functional based on the same models in (A). ITR indicates intergenic transcribed
836 regions.

837 **Figure 6.** Phenotype gene, pseudogene, unexpressed intergenic, and RNA gene score
838 distributions from four-class predictions. Stacked bar plots indicate the phenotype protein-coding
839 gene (dark blue), RNA gene (light blue), pseudogene (red), intergenic (yellow) score for each (A)
840 RNA training set gene, (B) phenotype gene, (C) pseudogene, (D) random unexpressed intergenic
841 region, (E) intergenic transcribed region, (F) ncRNA from Araport11, and (G) ncRNA from
842 TAIR10. Black vertical lines indicate boundaries of classification regions, with sequences
843 classified according to highest prediction score. Numbers within or pointing toward a
844 classification regions within a chart indicate the percentage of sequences predicted as, in order
845 phenotype gene, RNA gene, pseudogene, or intergenic. The color bars at the bottom of the chart
846 indicate whether a region of the chart is considered phenotype protein-coding gene-like (dark
847 blue), RNA gene-like (light blue), pseudogene-like (red), or intergenic-like (yellow).

848 SUPPLEMENTAL FIGURE LEGENDS

849 **Supplemental Figure 1.** Relationship between dinucleotide frequencies in phenotype gene and
850 pseudogene sequences. Percentages of all 16 dinucleotides in phenotype genes (X-axis) and
851 pseudogenes (Y-axis). Gray dotted line indicates the line of best fit.

852 **Supplemental Figure 2.** Functional likelihood scores by expression breadth. Distributions of
853 functional likelihood scores for phenotype genes (blue) and pseudogenes (red) for sequences
854 expressed in one-to-seven tissues for (A) the full model and (B) a tissue-agnostic model

855 generated while excluding the expression breadth feature and merging tissue-specific features.
856 The tissue-agnostic model performs better for among narrowly-expressed phenotype genes.

857 **Supplemental Figure 3.** Functional likelihood scores from the 500 bp tissue-agnostic model.
858 Functional likelihood distributions for (A) phenotype genes, (B) pseudogenes, (C) protein-coding
859 genes, (D) transposable elements, (E) random unexpressed intergenic sequences, (F) intergenic
860 transcribed regions, (G) ncRNAs from Araport11, and (H) ncRNAs from TAIR10 from the
861 tissue-agnostic model built while excluding the expression breadth and tissue-specific features.
862 Features were calculated from a random 500 bp region from within the sequence body. Higher
863 functional likelihood values indicate greater similarity to phenotype genes while lower values
864 indicate similarity to pseudogenes. Vertical dashed lines display the threshold to predict a
865 sequence as functional or non-functional. The numbers to the left and right of the dashed line
866 show the percentage of sequences predicted as functional or non-functional, respectively.

867 **Supplemental Figure 4.** Expression breadth of sequence types. Expression breadth distributions
868 for sequence types from (A) 500 bp feature regions and (B) 100 bp feature regions.

869 **Supplemental Figure 5.** ITR and annotated ncRNA distance to and feature similarity with
870 neighboring genes. (A) Distance from intergenic transcribed regions (ITRs) and annotated
871 ncRNAs that are predicted as functional (F) or non-functional (NF) to the closest neighboring
872 gene. (B,C,D) Feature similarity based on Pearson's Correlation Coefficients between (A)
873 proximal neighbors (within 95th percentile of intron lengths; distance=456), (B) distal neighbors
874 (greater than 95th percentile of intron lengths), and (C) random pairs of ITRs, ncRNAs from
875 Araport11, and ncRNAs from TAIR10 and annotated genes, as well as pairs of annotated genes.
876 Pairs involving ITRs and annotated ncRNAs were further divided by whether the ITR or ncRNA
877 sequence was predicted as functional (F) or non-functional (NF). Features were quantile
878 normalized prior to calculating correlations.

879 **Supplemental Figure 6.** Functional likelihood scores from the 100 bp tissue-agnostic model.
880 Functional likelihood distributions for (A) phenotype genes, (B) pseudogenes, (C) protein-coding
881 genes, (D) transposable elements, (E) random unexpressed intergenic sequences, (F) intergenic
882 transcribed regions (ITR), (G) ncRNAs from Araport11, (H) ncRNAs from TAIR10, (I)
883 microRNAs, (J) small nucleolar RNAs, (K) small nuclear RNAs, and (L) RNA genes with

884 documented loss-of-function phenotypes from the tissue-agnostic model built while excluding
885 the expression breadth and tissue-specific features. Features were calculated from a random 100
886 bp region from within the sequence body. Higher functional likelihood values indicate greater
887 similarity to phenotype genes while lower values indicate similarity to pseudogenes. Vertical
888 dashed lines display the threshold to predict a sequence as functional or non-functional. The
889 numbers to the left and right of the dashed line show the percentage of sequences predicted as
890 functional or non-functional, respectively.

891 **Supplemental Figure 7.** Translation evidence for sequences predicted as phenotype protein-
892 coding gene-like and RNA gene-like. Translation evidence was based on sequence overlap in
893 two shotgun proteomics datasets.

894 **SUPPLEMENTAL TABLES**

895 **Supplemental Table 1. Leaf tissue RNA-sequencing datasets for 15 flowering plant species**

896

897 **Supplemental Table 2. Conservation, biochemical, and sequence-structure features**
898 **calculated from 500 bp sequences.**

899

900 **Supplemental Table 3. False positive and false negative rates for single feature**
901 **classifications.**

902

903 **Supplemental Table 4. Predictions for the full, tissue-agnostic, 100 bp, and four-class**
904 **models.**

905

906 **Supplemental Table 5. RNA-sequencing datasets for identifying intergenic transcribed**
907 **regions, calculating transcription activity features, and assessing tissue-specific predictions.**

908

909 **Supplemental Table 6. Conservation, biochemical, and sequence-structure features**
910 **calculated from 100 bp sequences.**

911

912 **Supplemental Table 7. RNA genes with documented loss-of-function phenotypes.**

913

914 **DATA ACCESS**

915 All data are available in the text of this article or in the supplemental material.

916 **ACKNOWLEDGEMENTS**

917 The authors wish to thank (in reverse alphabetical order - we do not wish our acknowledgees to
918 feel superior or inferior to one another) Sahra Uygun, Bethany Moore, Gaurav Moghe, Ming-
919 Jung Liu, and Christina Azodi for providing highly useful data used in this study. This work was
920 supported by the Michigan State University Dissertation Continuation Scholarship to J.P.L, the
921 National Science Foundation (NSF) Plant Genomics Research Experience for Undergraduates
922 support to R.P.S., the Taiwan Ministry of Science and Technology Postdoctoral Research
923 Abroad Program MOST-104-2917-I-564-070 to Z.T.-Y.T, and NSF grants (MCB-1119778,
924 IOS-1126998, and IOS-1546617) to S.-H.S.

925 **AUTHOR CONTRIBUTIONS**

926 J.P.L., Z.T.-Y.T., and S.-H.S. designed the research. J.P.L., Z.T.-Y.T., R.P.S., and N.L.P.
927 performed the research. J.P.L., Z.T.-Y.T., R.P.S., N.L.P., and S.-H.S. wrote the article.

928 **DISCLOSURE DECLARATION**

929 The authors have no conflicts of interest to disclose.

930 **REFERENCES**

- 931 Ajjawi I, Lu Y, Savage LJ, Bell SM, Last RL. 2010. Large-scale reverse genetics in Arabidopsis:
932 case studies from the Chloroplast 2010 Project. *Plant Physiol* **152**: 529–540.
- 933 Amundson R, Lauder GV. 1994. Function without purpose. *Biol Philos* **9**: 443–469
- 934 Bennetzen JL. 2005. Mechanisms of Recent Genome Size Variation in Flowering Plants. *Ann*
935 *Bot* **95**: 127–132.
- 936 Berardini TZ, Reiser L, Li D, Mezheritsky Y, Muller R, Strait E, Huala E. 2015. The
937 Arabidopsis information resource: Making and mining the “gold standard” annotated
938 reference plant genome. *Genesis* **53**: 474–485.

- 939 Bernard D, Prasanth KV, Tripathi V, Colasse S, Nakamura T, Xuan Z, Zhang MQ, Sedel F,
940 Jourden L, Couplier F, et al. 2010. A long nuclear-retained non-coding RNA regulates
941 synaptogenesis by modulating gene expression. *EMBO J* **29**: 3082–3093.
- 942 Boeck ME, Huynh C, Gevirtzman L, Thompson OA, Wang G, Kasper DM, Reinke V, Hillier
943 LW, Waterston RH. 2016. The time-resolved transcriptome of *C. elegans*. *Genome Res* **26**:
944 1441–1450.
- 945 Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence
946 data. *Bioinformatics* **30**: 2114–2120.
- 947 Brenchley R, Spannagl M, Pfeifer M, Barker GLA, D’Amore R, Allen AM, McKenzie N,
948 Kramer M, Kerhornou A, Bolser D, et al. 2012. Analysis of the bread wheat genome using
949 whole-genome shotgun sequencing. *Nature* **491**: 705–710.
- 950 Brown JB, Boley N, Eisman R, May GE, Stoiber MH, Duff MO, Booth BW, Wen J, Park S,
951 Suzuki AM, et al. 2014. Diversity and dynamics of the *Drosophila* transcriptome. *Nature*
952 **512**: 393–399.
- 953 Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, Rinn JL. 2011. Integrative
954 annotation of human large intergenic noncoding RNAs reveals global properties and
955 specific subclasses. *Genes Dev* **25**: 1915–1927.
- 956 Cao J, Schneeberger K, Ossowski S, Günther T, Bender S, Fitz J, Koenig D, Lanz C, Stegle O,
957 Lippert C, et al. 2011. Whole-genome sequencing of multiple *Arabidopsis thaliana*
958 populations. *Nat Genet* **43**: 956–963.
- 959 Chen Z-X, Sturgill D, Qu J, Jiang H, Park S, Boley N, Suzuki AM, Fletcher AR, Plachetzki DC,
960 FitzGerald PC, et al. 2014. Comparative validation of the *D. melanogaster* modENCODE
961 transcriptome annotation. *Genome Res* **24**: 1209–1223.
- 962 Cummins R. 1975. Functional Analysis. *J Philos* **72**: 741.
- 963 Doolittle WF, Brunet TDP, Linquist S, Gregory TR. 2014. Distinguishing between “function”
964 and “effect” in genome biology. *Genome Biol Evol* **6**: 1234–1237.

- 965 Eddy SR. 2013. The ENCODE project: missteps overshadowing a success. *Curr Biol* **23**: R259–
966 61
- 967 ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the human
968 genome. *Nature* **489**: 57–74.
- 969 Friedel M, Nikolajewa S, Sühnel J, Wilhelm T. 2009. DiProDB: a database for dinucleotide
970 properties. *Nucleic Acids Res* **37**: D37–40.
- 971 Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W, Hellsten
972 U, Putnam N, et al. 2011. Phytozome: a comparative platform for green plant genomics.
973 *Nucleic Acids Res* **40**: D1178–D1186.
- 974 Graur D, Zheng Y, Price N, Azevedo RBR, Zufall RA, Elhaik E. 2013. On the immortality of
975 television sets: “function” in the human genome according to the evolution-free gospel of
976 ENCODE. *Genome Biol Evol* **5**: 578–590.
- 977 Gulko B, Gronau I, Hubisz MJ, Siepel A. 2014. *Probabilities of Fitness Consequences for Point*
978 *Mutations Across the Human Genome*. <http://dx.doi.org/10.1101/006825>.
- 979 Guo H-S, Xie Q, Fei J-F, Chua N-H. 2005. MicroRNA directs mRNA cleavage of the
980 transcription factor NAC1 to downregulate auxin signals for arabidopsis lateral root
981 development. *Plant Cell* **17**: 1376–1386.
- 982 Guttman M, Amit I, Garber M, French C, Lin MF, Feldser D, Huarte M, Zuk O, Carey BW,
983 Cassady JP, et al. 2009. Chromatin signature reveals over a thousand highly conserved large
984 non-coding RNAs in mammals. *Nature* **458**: 223–227.
- 985 Hall M, Frank E, Holmes G, Pfahringer B, Reutemann P, Witten IH. 2009. The WEKA data
986 mining software. *ACM SIGKDD Explorations Newsletter* **11**: 10.
- 987 Hanada K, Zhang X, Borevitz JO, Li W-H, Shiu S-H. 2007. A large number of novel coding
988 small open reading frames in the intergenic regions of the *Arabidopsis thaliana* genome are
989 transcribed and/or under purifying selection. *Genome Res* **17**: 632–640.

- 990 Hanada K, Higuchi-Takeuchi M, Okamoto M, Yoshizumi T, Shimizu M, Nakaminami K, Nishi
991 R, Ohashi C, Iida K, Tanaka M, et al. 2013. Small open reading frames associated with
992 morphogenesis are hidden in plant genomes. *Proc Natl Acad Sci U S A* **110**: 2395–2400.
- 993 Hardiman KE, Brewster R, Khan SM, Deo M, Bodmer R. 2002. The bereft gene, a potential target of the
994 neural selector gene cut, contributes to bristle morphogenesis. *Genetics* **161**: 231–247.
- 995 Heinen TAJ, Staubach F, Häming D, Tautz D. 2009. Emergence of a New Gene from an Intergenic
996 Region. *Curr Biol* **19**: 1527–1531.
- 997 Hsieh L-C, Lin S-I, Shih AC-C, Chen J-W, Lin W-Y, Tseng C-Y, Li W-H, Chiou T-J. 2009.
998 Uncovering small RNA-mediated responses to phosphate deficiency in Arabidopsis by deep
999 sequencing. *Plant Physiol* **151**: 2120–2132.
- 1000 Ibarra-Laclette E, Lyons E, Hernández-Guzmán G, Pérez-Torres CA, Carretero-Paulet L, Chang T-H, Lan
1001 T, Welch AJ, Juárez MJA, Simpson J, et al. 2013. Architecture and evolution of a minute plant
1002 genome. *Nature* **498**: 94–98.
- 1003 Ivanova N, Dobrin R, Lu R, Kotenko I, Levorse J, DeCoste C, Schafer X, Lun Y, Lemischka IR.
1004 2006. Dissecting self-renewal in stem cells with RNA interference. *Nature* **442**: 533–538.
- 1005 Johnson JM, Edwards S, Shoemaker D, Schadt EE. 2005. Dark matter in the genome: evidence
1006 of widespread transcription detected by microarray tiling experiments. *Trends Genet* **21**: 93–
1007 102.
- 1008 Karreth FA, Reschke M, Ruocco A, Ng C, Chapuy B, Léopold V, Sjoberg M, Keane TM, Verma
1009 A, Ala U, et al. 2015. The BRAF pseudogene functions as a competitive endogenous RNA
1010 and induces lymphoma in vivo. *Cell* **161**: 319–332.
- 1011 Kellis M, Wold B, Snyder MP, Bernstein BE, Kundaje A, Marinov GK, Ward LD, Birney E,
1012 Crawford GE, Dekker J, et al. 2014. Defining functional DNA elements in the human
1013 genome. *Proc Natl Acad Sci U S A* **111**: 6131–6138.
- 1014 Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. 2013. TopHat2: accurate
1015 alignment of transcriptomes in the presence of insertions, deletions and gene fusions.
1016 *Genome Biol* **14**: R36.

- 1017 Koehler R, Issac H, Cloonan N, Grimmond SM. 2011. The uniqueome: a mappability resource
1018 for short-tag sequencing. *Bioinformatics* **27**: 272–274.
- 1019 Kozomara A, Griffiths-Jones S. 2014. miRBase: annotating high confidence microRNAs using
1020 deep sequencing data. *Nucleic Acids Res* **42**: D68–73.
- 1021 Krishnakumar V, Hanlon MR, Contrino S, Ferlanti ES, Karamycheva S, Kim M, Rosen BD,
1022 Cheng C-Y, Moreira W, Mock SA, et al. 2015. Araport: the Arabidopsis information portal.
1023 *Nucleic Acids Res* **43**: D1003–9.
- 1024 Krueger F, Andrews SR. 2011. Bismark: a flexible aligner and methylation caller for Bisulfite-
1025 Seq applications. *Bioinformatics* **27**: 1571–1572.
- 1026 Kuromori T, Wada T, Kamiya A, Yuguchi M, Yokouchi T, Imura Y, Takabe H, Sakurai T,
1027 Akiyama K, Hirayama T, et al. 2006. A trial of phenome analysis using 4000 Ds-insertional
1028 mutants in gene-coding regions of Arabidopsis. *Plant J* **47**: 640–651.
- 1029 Lai K-MV, Gong G, Atanasio A, Rojas J, Quispe J, Posca J, White D, Huang M, Fedorova D,
1030 Grant C, et al. 2015. Diverse Phenotypes and Specific Transcription Patterns in Twenty
1031 Mouse Lines with Ablated LincRNAs. *PLoS One* **10**: e0125522.
- 1032 Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment
1033 of short DNA sequences to the human genome. *Genome Biol* **10**: R25.
- 1034 Li F, Zheng Q, Vandivier LE, Willmann MR, Chen Y, Gregory BD. 2012a. Regulatory impact of
1035 RNA secondary structure across the Arabidopsis transcriptome. *Plant Cell* **24**: 4346–4359.
- 1036 Liu M-J, Seddon AE, Tsai ZT-Y, Major IT, Floer M, Howe GA, Shiu S-H. 2015. Determinants
1037 of nucleosome positioning and their influence on plant gene expression. *Genome Res* **25**:
1038 1182–1195.
- 1039 Li W, Cui X, Meng Z, Huang X, Xie Q, Wu H, Jin H, Zhang D, Liang W. 2012b. Transcriptional
1040 regulation of Arabidopsis MIR168a and argonaute1 homeostasis in abscisic acid and abiotic
1041 stress responses. *Plant Physiol* **158**: 1279–1292.

- 1042 Li WH, Gojobori T, Nei M. 1981. Pseudogenes as a paradigm of neutral evolution. *Nature* **292**:
1043 237–239.
- 1044 Lloyd J, Meinke D. 2012. A comprehensive dataset of genes with a loss-of-function mutant
1045 phenotype in Arabidopsis. *Plant Physiol* **158**: 1115–1129.
- 1046 Lloyd JP, Seddon AE, Moghe GD, Simenc MC, Shiu S-H. 2015. Characteristics of Plant
1047 Essential Genes Allow for within- and between-Species Prediction of Lethal Mutant
1048 Phenotypes. *Plant Cell* **27**: 2133–2147.
- 1049 Marahrens Y, Panning B, Dausman J, Strauss W, Jaenisch R. 1997. Xist-deficient mice are defective in
1050 dosage compensation but not spermatogenesis. *Genes Dev* **11**: 156–166.
- 1051 Mattick JS. 2009. The Genetic Signatures of Noncoding RNAs. *PLoS Genet* **5**: e1000459.
- 1052 Mercer TR, Dinger ME, Mattick JS. 2009. Long non-coding RNAs: insights into functions. *Nat*
1053 *Rev Genet* **10**: 155–159
- 1054 Michael TP, Jackson S. 2013. The First 50 Plant Genomes. *Plant Genome* **6**: 0.
- 1055 Moghe GD, Lehti-Shiu MD, Seddon AE, Yin S, Chen Y, Juntawong P, Brandizzi F, Bailey-
1056 Serres J, Shiu S-H. 2013. Characteristics and significance of intergenic polyadenylated RNA
1057 transcription in Arabidopsis. *Plant Physiol* **161**: 210–224.
- 1058 Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, Gerstein M, Snyder M. 2008. The
1059 Transcriptional Landscape of the Yeast Genome Defined by RNA Sequencing. *Science* **320**:
1060 1344–1349.
- 1061 Neander K. 1991. Functions as Selected Effects: The Conceptual Analyst's Defense. *Philos Sci*
1062 **58**: 168–184.
- 1063 Ngernprasirtsiri J, Kobayashi H, Akazawa T. 1988. DNA methylation as a mechanism of
1064 transcriptional regulation in nonphotosynthetic plastids in plant cells. *Proc Natl Acad Sci U*
1065 *S A* **85**: 4750–4754.

- 1066 Ning S, Wang P, Ye J, Li X, Li R, Zhao Z, Huo X, Wang L, Li F, Li X. 2013. A global map for dissecting
1067 phenotypic variants in human lincRNAs. *Eur J Hum Genet* **21**: 1128–1133.
- 1068 Niu D-K, Jiang L. 2013. Can ENCODE tell us how much junk DNA we carry in our genome?
1069 *Biochem Biophys Res Commun* **430**: 1340–1343.
- 1070 Nobuta K, Venu RC, Lu C, Beló A, Vemaraju K, Kulkarni K, Wang W, Pillay M, Green PJ,
1071 Wang G-L, et al. 2007. An expression atlas of rice mRNAs and small RNAs. *Nat Biotechnol*
1072 **25**: 473–477.
- 1073 O'Malley RC, Huang S-SC, Song L, Lewsey MG, Bartlett A, Nery JR, Galli M, Gallavotti A,
1074 Ecker JR. 2016. Cistrome and Epicistrome Features Shape the Regulatory DNA Landscape.
1075 *Cell* **166**: 1598.
- 1076 Pang KC, Frith MC, Mattick JS. 2006. Rapid evolution of noncoding RNAs: lack of
1077 conservation does not mean lack of function. *Trends Genet* **22**: 1–5.
- 1078 Penny GD, Kay GF, Sheardown SA, Rastan S, Brockdorff N. 1996. Requirement for Xist in X
1079 chromosome inactivation. *Nature* **379**: 131–137.
- 1080 Poliseno L, Salmena L, Zhang J, Carver B, Haveman WJ, Pandolfi PP. 2010. A coding-
1081 independent function of gene and pseudogene mRNAs regulates tumour biology. *Nature*
1082 **465**: 1033–1038.
- 1083 Ponting CP, Belgard TG. 2010. Transcribed dark matter: meaning or myth? *Hum Mol Genet* **19**:
1084 R162–R168.
- 1085 Sauvageau M, Goff LA, Lodato S, Bonev B, Groff AF, Gerhardinger C, Sanchez-Gomez DB,
1086 Haciosuleyman E, Li E, Spence M, et al. 2013. Multiple knockout mouse models reveal
1087 lincRNAs are required for life and brain development. *Elife* **2**: e01749.
- 1088 Savage LJ, Imre KM, Hall DA, Last RL. 2013. Analysis of essential Arabidopsis nuclear genes
1089 encoding plastid-targeted proteins. *PLoS One* **8**: e73291.
- 1090 Schnable JC, Pedersen BS, Subramaniam S, Freeling M. 2011. Dose–Sensitivity, Conserved
1091 Non-Coding Sequences, and Duplicate Gene Retention Through Multiple Tetraploidies in

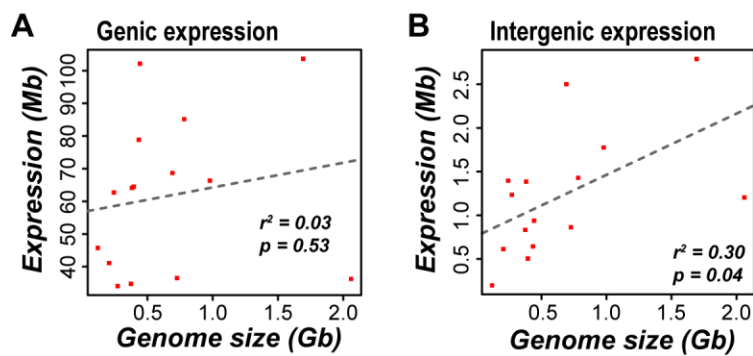
- 1092 the Grasses. *Front Plant Sci* **2**. <http://dx.doi.org/10.3389/fpls.2011.00002>.
- 1093 Schnable JC, Wang X, Pires JC, Freeling M. 2012. Escape from preferential retention following
1094 repeated whole genome duplications in plants. *Front Plant Sci* **3**: 94.
- 1095 Schreiber SL, Bernstein BE. 2002. Signaling Network Model of Chromatin. *Cell* **111**: 771–778
- 1096 Shin H, Shin H-S, Chen R, Harrison MJ. 2006. Loss of At4 function impacts phosphate
1097 distribution between the roots and the shoots during phosphate starvation. *Plant J* **45**: 712–
1098 726.
- 1099 Sing T, Sander O, Beerenwinkel N, Lengauer T. 2005. ROCr: visualizing classifier performance
1100 in R. *Bioinformatics* **21**: 3940–3941.
- 1101 Slotkin RK, Keith Slotkin R, Martienssen R. 2007. Transposable elements and the epigenetic
1102 regulation of the genome. *Nat Rev Genet* **8**: 272–285.
- 1103 Soltis DE, Albert VA, Leebens-Mack J, Bell CD, Paterson AH, Zheng C, Sankoff D,
1104 Depamphilis CW, Wall PK, Soltis PS. 2009. Polyploidy and angiosperm diversification. *Am*
1105 *J Bot* **96**: 336–348.
- 1106 Stergachis AB, Haugen E, Shafer A, Fu W, Vernot B, Reynolds A, Raubitschek A, Ziegler S,
1107 LeProust EM, Akey JM, et al. 2013. Exonic Transcription Factor Binding Directs Codon
1108 Choice and Affects Protein Evolution. *Science* **342**: 1367–1372.
- 1109 Stolc V, Samanta MP, Tongprasit W, Sethi H, Liang S, Nelson DC, Hegeman A, Nelson C,
1110 Rancour D, Bednarek S, et al. 2005. Identification of transcribed sequences in *Arabidopsis*
1111 *thaliana* by using high-resolution genome tiling arrays. *Proc Natl Acad Sci U S A* **102**:
1112 4453–4458.
- 1113 Struhl K. 2007. Transcriptional noise and the fidelity of initiation by RNA polymerase II. *Nat*
1114 *Struct Mol Biol* **14**: 103–105.
- 1115 Sullivan AM, Arsovski AA, Lempe J, Bubb KL, Weirauch MT, Sabo PJ, Sandstrom R, Thurman
1116 RE, Neph S, Reynolds AP, et al. 2014. Mapping and dynamics of regulatory DNA and

- 1117 transcription factor networks in *A. thaliana*. *Cell Rep* **8**: 2015–2030.
- 1118 Svensson O, Arvestad L, Lagergren J. 2006. Genome-wide survey for biologically functional
1119 pseudogenes. *PLoS Comput Biol* **2**: e46.
- 1120 Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold
1121 BJ, Pachter L. 2010. Transcript assembly and quantification by RNA-Seq reveals
1122 unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* **28**:
1123 511–515.
- 1124 Tsai ZT-Y, Shiu S-H, Tsai H-K. 2015. Contribution of Sequence Motif, Chromatin State, and
1125 DNA Structure Features to Predictive Models of Transcription Factor Binding in Yeast.
1126 *PLoS Comput Biol* **11**: e1004418.
- 1127 Tsai ZT-Y, Lloyd JP, Shiu S-H. 2017. Defining Functional Genic Regions in the Human
1128 Genome through Integration of Biochemical, Evolutionary, and Genetic Evidence. *Mol Biol*
1129 *Evol.* <http://dx.doi.org/10.1093/molbev/msx101>
- 1130 Ulitsky I, Shkumatava A, Jan CH, Sive H, Bartel DP. 2011. Conserved function of lincRNAs in
1131 vertebrate embryonic development despite rapid sequence evolution. *Cell* **147**: 1537–1550.
- 1132 VanBuren R, Bryant D, Edger PP, Tang H, Burgess D, Challabathula D, Spittle K, Hall R, Gu J,
1133 Lyons E, et al. 2015. Single-molecule sequencing of the desiccation-tolerant grass
1134 *Oropetium thomaeum*. *Nature* **527**: 508–511.
- 1135 Wang Z, Zang C, Rosenfeld JA, Schones DE, Barski A, Cuddapah S, Cui K, Roh T-Y, Peng W,
1136 Zhang MQ, et al. 2008. Combinatorial patterns of histone acetylations and methylations in
1137 the human genome. *Nat Genet* **40**: 897–903.
- 1138 Woodhouse MR, Schnable JC, Pedersen BS, Lyons E, Lisch D, Subramaniam S, Freeling M.
1139 2010. Following Tetraploidy in Maize, a Short Deletion Mechanism Removed Genes
1140 Preferentially from One of the Two Homeologs. *PLoS Biol* **8**: e1000409.
- 1141 Xu S, Grullon S, Ge K, Peng W. 2014. Spatial clustering for identification of ChIP-enriched
1142 regions (SICER) to map regions of histone methylation patterns in embryonic stem cells.

- 1143 *Methods Mol Biol* **1150**: 97–111.
- 1144 Yamada K, Lim J, Dale JM, Chen H, Shinn P, Palm CJ, Southwick AM, Wu HC, Kim C,
1145 Nguyen M, et al. 2003. Empirical analysis of transcriptional activity in the Arabidopsis
1146 genome. *Science* **302**: 842–846.
- 1147 Yang L, Takuno S, Waters ER, Gaut BS. 2011. Lowly expressed genes in Arabidopsis thaliana
1148 bear the signature of possible pseudogenization by promoter degradation. *Mol Biol Evol* **28**:
1149 1193–1203.
- 1150 Zhang X, Yazaki J, Sundaresan A, Cokus S, Chan SW-L, Chen H, Henderson IR, Shinn P,
1151 Pellegrini M, Jacobsen SE, et al. 2006. Genome-wide high-resolution mapping and
1152 functional analysis of DNA methylation in arabidopsis. *Cell* **126**: 1189–1201.
- 1153 Zhao J, Sun BK, Erwin JA, Song J-J, Lee JT. 2008. Polycomb proteins targeted by a short repeat
1154 RNA to the mouse X chromosome. *Science* **322**: 750–756.
- 1155 Zhao Y, Li H, Fang S, Kang Y, Wu W, Hao Y, Li Z, Bu D, Sun N, Zhang MQ, et al. 2016.
1156 NONCODE 2016: an informative and valuable data source of long non-coding RNAs.
1157 *Nucleic Acids Res* **44**: D203–8.
- 1158 Zou C, Lehti-Shiu MD, Thibaud-Nissen F, Prakash T, Buell CR, Shiu S-H. 2009. Evolutionary
1159 and expression signatures of pseudogenes in Arabidopsis and rice. *Plant Physiol* **151**: 3–15.
- 1160

1161 **FIGURES**

1162 **Figure 1.**

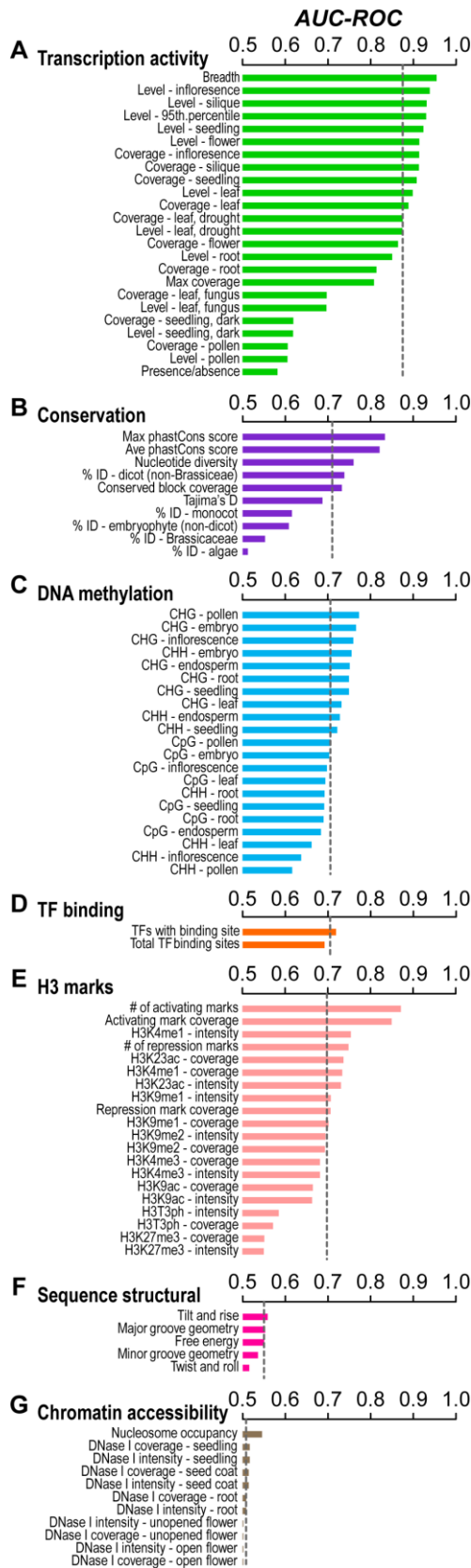


1163

1164

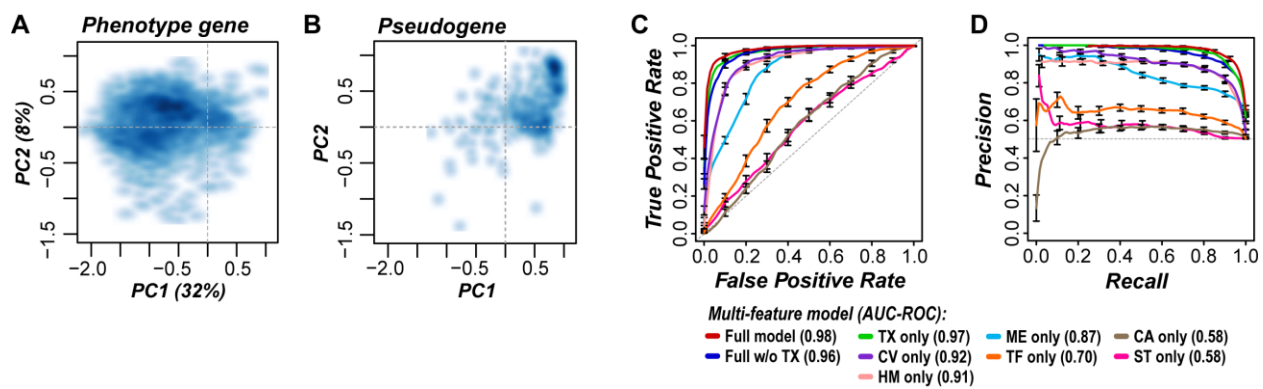
1165

1166 **Figure 2.**



1167

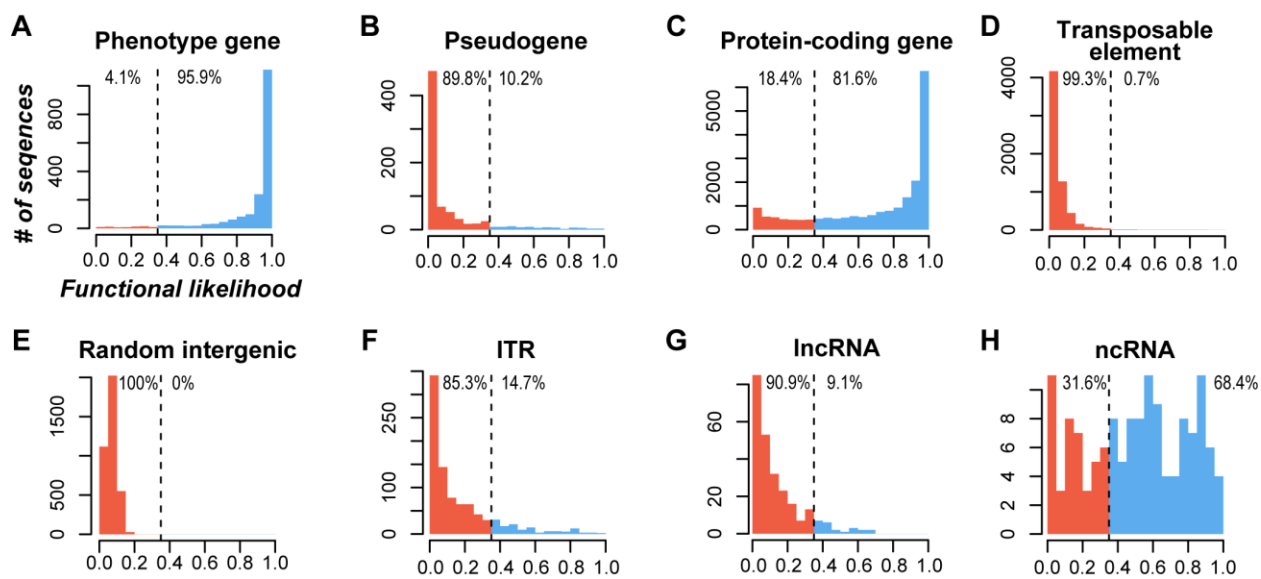
1168 **Figure 3.**



1169

1170

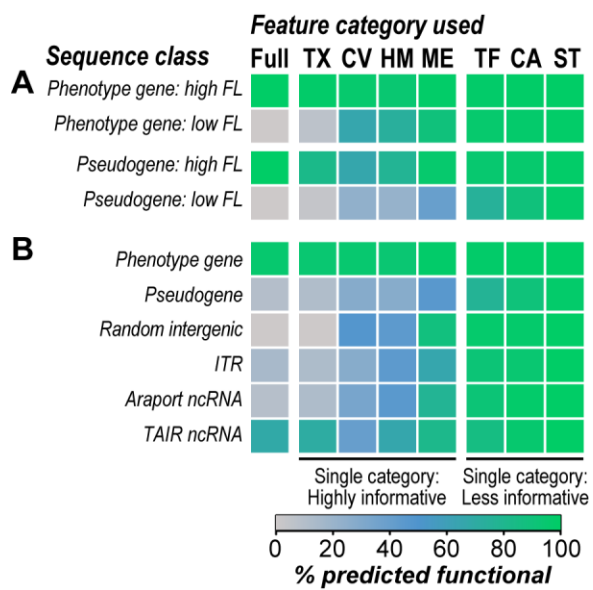
1171 **Figure 4.**



1172

1173

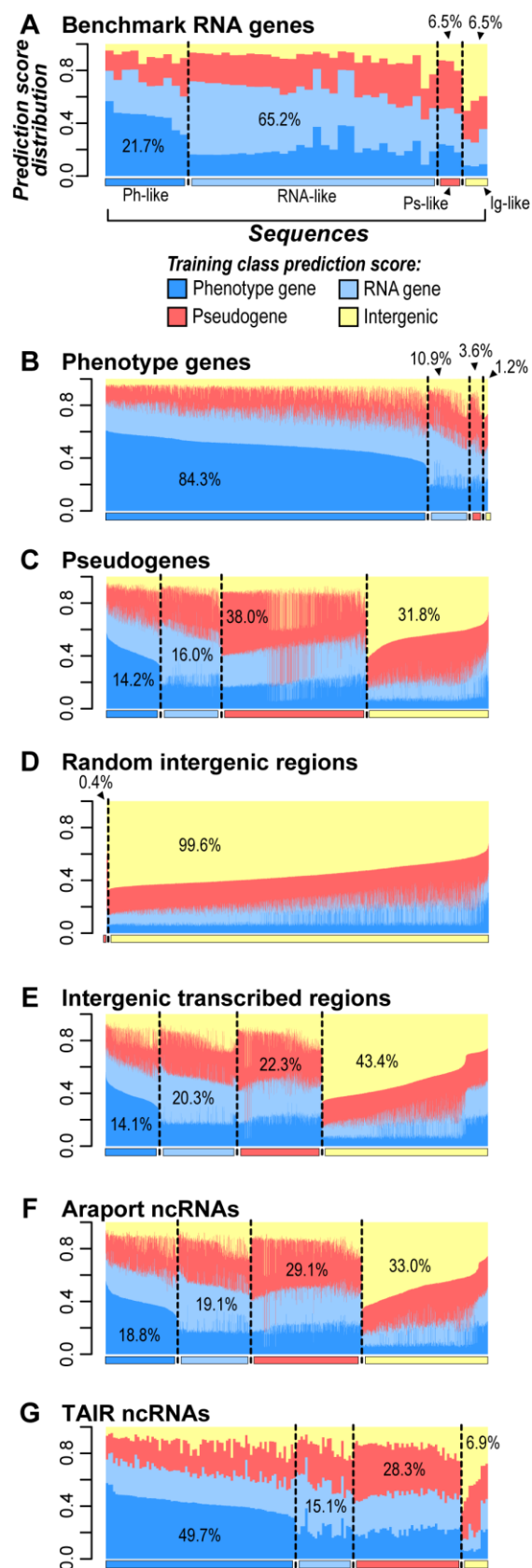
1174 **Figure 5.**



1175

1176

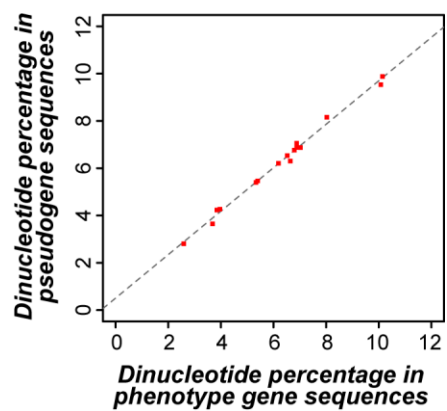
1177 **Figure 6.**



1178

1179 **SUPPLEMENTAL FIGURES**

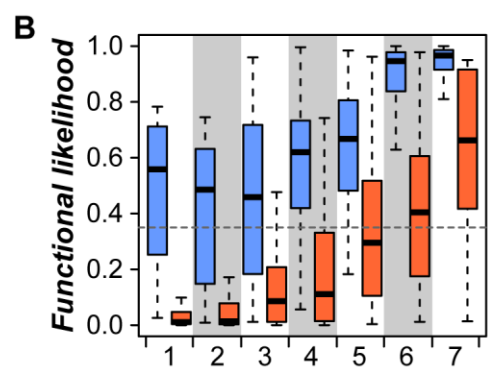
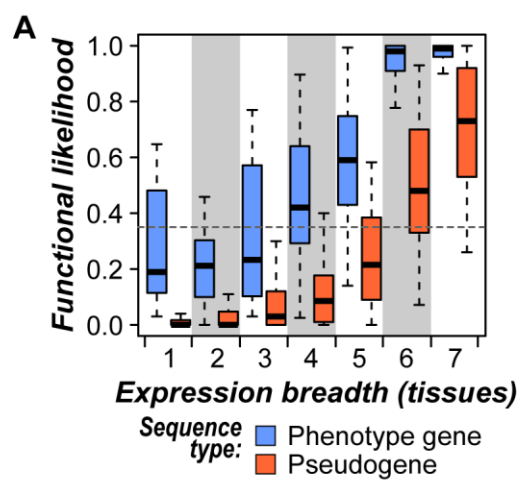
1180 **Supplemental Figure 1.**



1181

1182

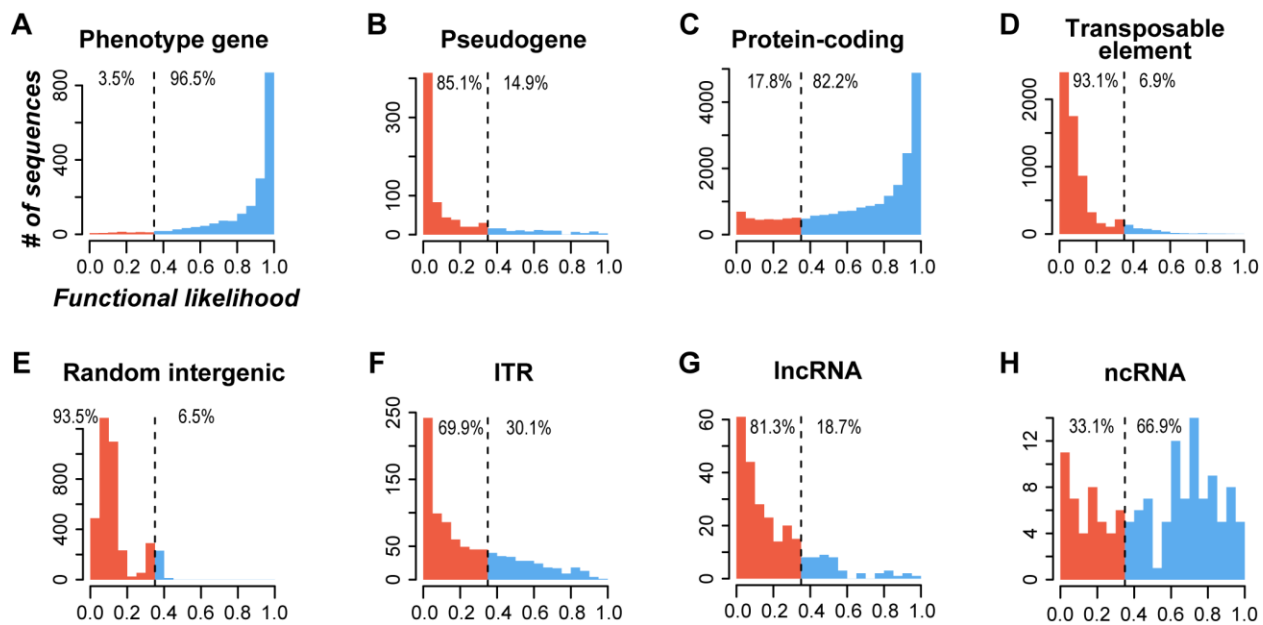
1183 **Supplemental Figure 2.**



1184

1185

1186 **Supplemental Figure 3.**

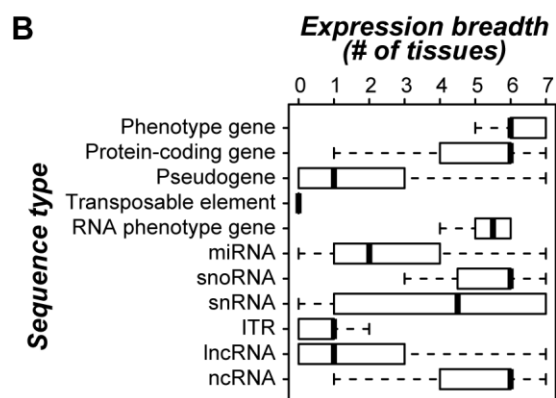
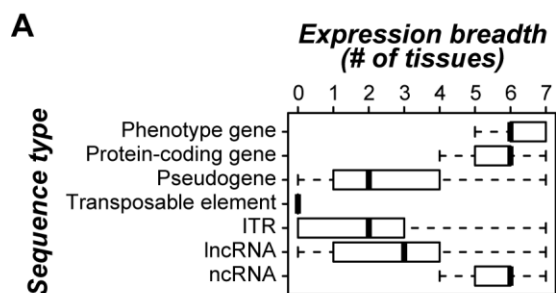


1187

1188

1189

1190 **Supplemental Figure 4.**

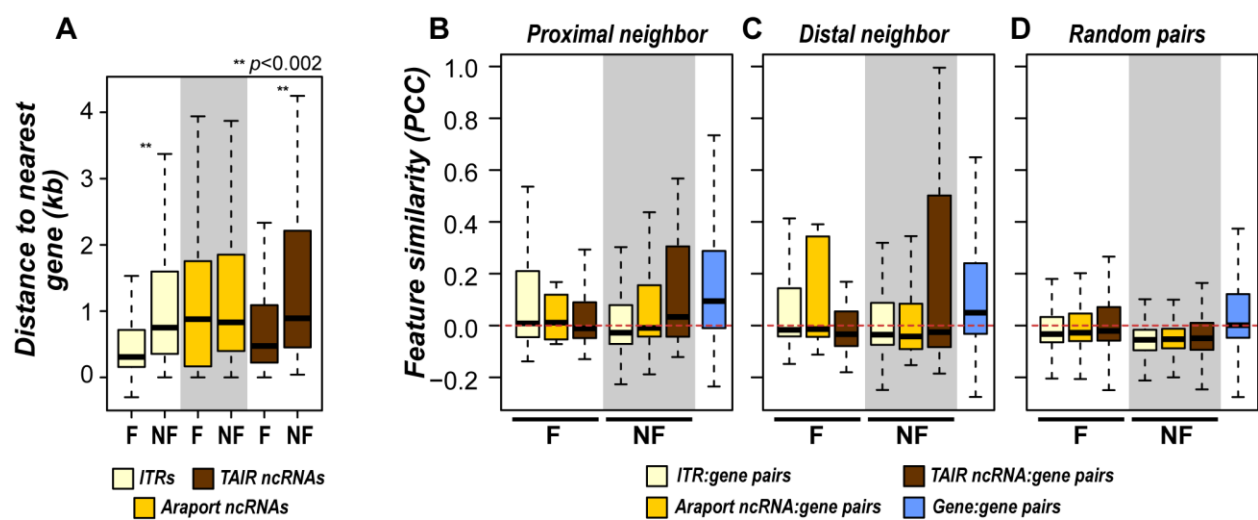


1191

1192

1193

1194 **Supplemental Figure 5.**

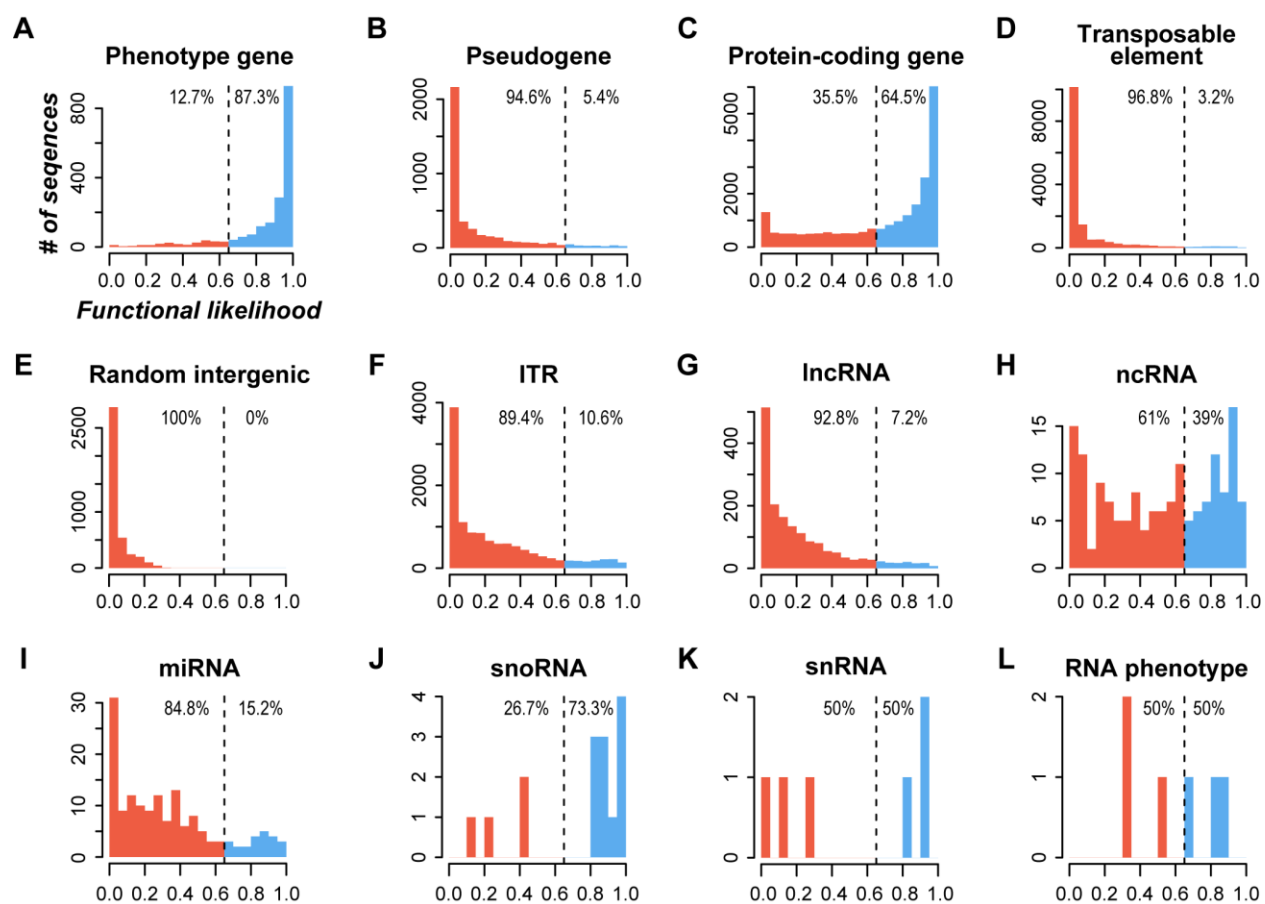


1195

1196

1197

1198 **Supplemental Figure 6.**

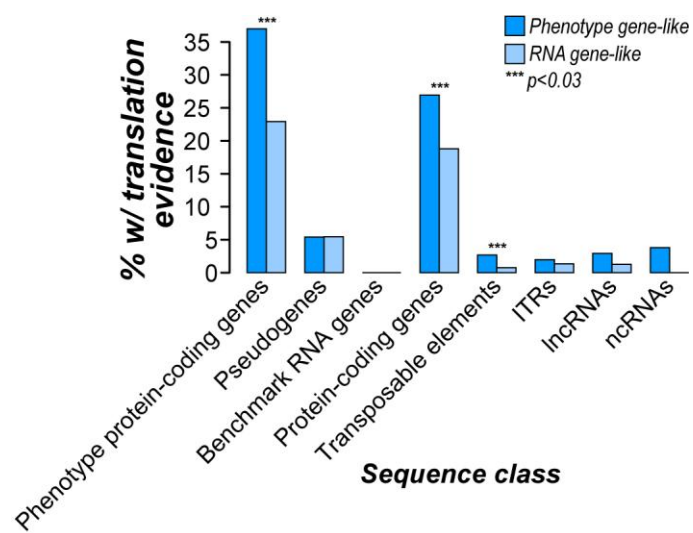


1199

1200

1201

1202 **Supplemental Figure 7.**



1203

1204