

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

Defining the functional significance of intergenic transcribed regions

John P. Lloyd¹, Zing Tsung-Yeh Tsai², Rosalie P. Sowers³, Nicholas L. Panchy⁴, Shin-Han Shiu^{1,4,5*}

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

² Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI 48109, USA

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

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

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

* To whom correspondence should be addressed.

Correspondence to:

Shin-Han Shiu

Michigan State University

E-mail: shius@msu.edu

Telephone number: +1-517-353-7196

Running title: Functionality of intergenic transcripts

Keywords: Intergenic transcription, ncRNAs, definition of function, molecular evolution, machine learning, data integration

Author contributions: 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. performed the research. J.P.L., Z.T.-Y.T., R.P.S., N.L.P., and S.-H.S. wrote the article.

30 **ABSTRACT**

31 With advances in transcript profiling, the presence of transcriptional activities in intergenic
32 regions has been well established. However, whether intergenic expression reflects
33 transcriptional noise or activity of novel genes remains unclear. We identified intergenic
34 transcribed regions (ITRs) in 15 diverse flowering plant species and found that the amount of
35 intergenic expression correlates with genome size, a pattern that could be expected if intergenic
36 expression is largely nonfunctional. To further assess the functionality of ITRs, we first built
37 machine learning classifiers using *Arabidopsis thaliana* as a model that accurately distinguish
38 functional sequences (phenotype genes) and likely nonfunctional ones (pseudogenes and
39 unexpressed intergenic regions) by integrating 93 biochemical, evolutionary, and sequence-
40 structure features. Next, by applying the models genome-wide, we found that 4,427 ITRs (38%)
41 and 796 annotated ncRNAs (44%) had features significantly similar to benchmark protein-
42 coding or RNA genes and thus were likely parts of functional genes. Approximately 60% of
43 ITRs and ncRNAs were more similar to nonfunctional sequences and were likely transcriptional
44 noise. The predictive framework established here provides not only a comprehensive look at how
45 functional, genic sequences are distinct from likely nonfunctional ones, but also a new way to
46 differentiate novel genes from genomic regions with noisy transcriptional activities.
47

48 INTRODUCTION

49 Advances in sequencing technology have helped to identify pervasive transcription in
50 intergenic regions with no annotated genes. These intergenic transcripts have been found in
51 metazoa and fungi, including human (ENCODE Project Consortium 2012), *Drosophila*
52 *melanogaster* (Brown et al. 2014), *Caenorhabditis elegans* (Boeck et al. 2016), and
53 *Saccharomyces cerevisiae* (Nagalakshmi et al. 2008). In plants, 7,000 to 15,000 intergenic
54 transcripts have also been reported in *Arabidopsis thaliana* (Yamada et al. 2003; Stolc et al.
55 2005; Moghe et al. 2013; Krishnakumar et al. 2015) and *Oryza sativa* (Nobuta et al. 2007). The
56 presence of intergenic transcripts indicates that there may be additional genes in genomes that
57 have escaped gene finding efforts thus far, including those that function as RNA genes (Simon
58 and Meyers 2011; Guil and Esteller 2012; Fei et al. 2013; Tan et al. 2015). Meanwhile, it is also
59 possible that some of these intergenic transcripts are products of un-regulated noise (Struhl
60 2007). Given the functional significance of most intergenic transcripts remains unclear, the
61 identification of functional intergenic transcribed regions (ITRs) represents a fundamental task
62 that is critical to our understanding of genome evolution.

63 Loss-of-function study represents the gold standard by which the functional significance
64 of genomic regions, including ITRs, can be confirmed (Ponting and Belgard 2010; 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). In human, 162 long intergenic non-
67 coding RNAs harbor phenotype-associated SNPs (Ning et al. 2013). In addition to intergenic
68 expression, most model organisms feature an abundance of annotated non-coding RNA (ncRNA)
69 sequences (Zhao et al. 2016), which are mostly identified through the presence of transcriptional
70 evidence occurring outside of annotated protein-coding genes. Thus, the only difference between
71 ITRs and most ncRNA sequences is whether or not they have been annotated. Similar to the ITR
72 examples above, a small number of ncRNAs have been confirmed as functional through loss-of-
73 function experiments including *Xist* in mouse (Penny et al. 1996; Marahrens et al. 1997), *Malat1*
74 in human (Bernard et al. 2010), *bereft* in *D. melanogaster* (Hardiman et al. 2002), and *At4* in *A.*
75 *thaliana* (Shin et al. 2006).

76 However, the number of ITRs and ncRNAs with well-established functions is dwarfed by
77 those without functional evidence. While some ITRs and ncRNAs can be novel genes, intergenic
78 transcription may also be the byproduct of noisy transcription that can occur due to nonspecific

79 landing of RNA Polymerase II (RNA Pol II) or spurious regulatory signals that drive expression
80 in random genomic regions (Struhl 2007). In the ENCODE project (ENCODE Project
81 Consortium 2012), ~80% of the human genome was defined as biochemically functional as
82 reproducible biochemical activities, e.g. transcription, could be detected. This has drawn
83 considerable critique because the existence of a biochemical activity is not an indication of
84 selection (Eddy 2013; Graur et al. 2013; Niu and Jiang 2013). Instead, it is advocated that a
85 genomic region with discernible activity is only functional if it is under selection (Amundson and
86 Lauder 1994; Graur et al. 2013; Doolittle et al. 2014). Under this “selected effect” functionality
87 definition, ITRs and most annotated ncRNA genes remain functionally ambiguous.

88 Due to the debate on the definitions of function post-ENCODE, Kellis et al. (Kellis et al.
89 2014) suggested that evolutionary, biochemical, and genetic evidences provide complementary
90 information to define functional genomic regions. Consistent with this, integration of
91 biochemical and conservation evidence was successful in identifying regions in the human
92 genome that are under selection (Gulko et al. 2014) and classification of human disease genes
93 and pseudogenes (Tsai et al. 2017). In this study, we adopt a similar framework to investigate the
94 functionality of intergenic transcription in plants. We first identified ITRs in 15 flowering plant
95 species with 17-fold genome size differences and evaluated the relationship between the
96 prevalence of intergenic expression and genome size. Next, we established machine learning
97 models using *A. thaliana* data to predict likely-functional ITRs and ncRNAs based on 93
98 evolutionary, biochemical, and sequence-structure features. Finally, we applied the models to
99 ITRs and annotated ncRNAs to determine whether these functionally ambiguous sequences are
100 more similar to benchmark functional or likely nonfunctional sequences.

101 **RESULTS & DISCUSSION**

102 **Genome size versus prevalence of intergenic transcripts indicates ITRs may generally be** 103 **nonfunctional**

104 Transcription of an unannotated, intergenic region could be due to nonfunctional
105 transcriptional noise or the activity of a novel gene. If noisy transcription occurs due to random
106 landing of RNA Pol II or spurious regulatory signals, a naïve expectation is that, as genome size
107 increases, the total nucleotides covered by ITRs would increase accordingly. By contrast, we
108 expect that the extent of expression for genic sequences will not be significantly correlated with

109 genome size because larger plant genomes do not necessarily have more genes ($r^2=0.01$;
110 $p=0.56$).

111 To gauge if ITRs generally behave more like what we expect of noisy or genic
112 transcription, we first identified genic and intergenic transcribed regions using leaf transcriptome
113 data from 15 flowering plants with 17-fold differences in genome size (Supplementary Table 1).
114 As expected, the coverage of expression originating from annotated genic regions had no
115 significant correlation with genome size ($r^2=0.03$; $p=0.53$; **Fig. 1A**). By contrast, the coverage of
116 intergenic expression was significantly and positively correlated ($r^2=0.30$; $p=0.04$; **Fig. 1B**),
117 consistent with the interpretation that a significant proportion of intergenic expression represents
118 transcriptional noise. However, the correlation between genome size and intergenic expression
119 explained ~30% of the variation (**Fig. 1B**), suggesting that other factors also affect ITR content,
120 including the possibility that some ITRs are truly functional, novel genes. To further evaluate the
121 functionality of intergenic transcripts, we next identified the biochemical and evolutionary
122 features of functional genic regions and tested whether intergenic transcripts in *A. thaliana* were
123 more similar to functional or nonfunctional sequences.

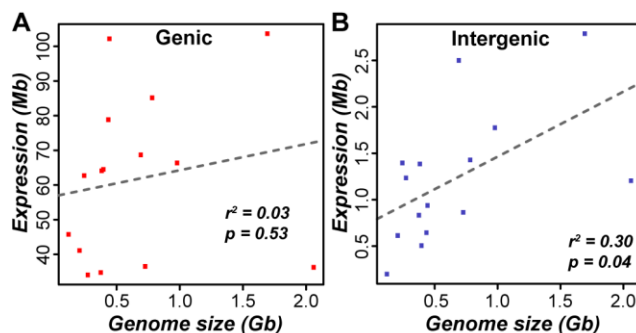


Figure 1. Relationship between genome size and number of nucleotides covered by RNA-seq reads (expression) in 15 flowering plant species. (A) annotated genic regions. (B) intergenic regions excluding any annotated features. Mb: megabase. Gb: gigabase. Dotted lines: linear model fits. r^2 : square of Pearson's correlation coefficient.

124

125 **Benchmark functional protein-coding and nonfunctional genomic sequences are** 126 **significantly distinct in multiple features**

127 To determine whether intergenic transcripts resemble functional sequences, we first
128 asked what features allow benchmark functional protein-coding and nonfunctional genomic
129 regions to be distinguished in the model plant *Arabidopsis thaliana*. For benchmark functional

130 sequences, we used protein-coding genes with visible loss-of-function phenotypes when mutated
131 (referred to as phenotype genes, $n=1,876$; see Materials and Methods). Because their mutations
132 have significant growth and/or developmental impact and likely contribute to reduced fitness,
133 these phenotype genes can be considered functional under the selected effect definition (Neander
134 1991). For benchmark nonfunctional genomic regions, we utilized pseudogene sequences
135 ($n=761$; see Materials and Methods). Considering that only 2% of pseudogenes are maintained
136 over 90 million years of divergence between human and mouse (Svensson et al. 2006), it is
137 expected that the majority of pseudogenes are no longer under selection (Li et al. 1981).

138 We evaluated 93 gene or gene product features for their ability to distinguish between
139 phenotype genes and pseudogenes. These features were grouped into seven categories, including
140 chromatin accessibility, DNA methylation, histone 3 (H3) marks, sequence conservation,
141 sequence-structure, transcription factor (TF) binding, and transcription activity (Supplementary
142 Table 2). We emphasize that no features were exclusive to protein-coding sequences. We used
143 Area Under the Curve - Receiver Operating Characteristic (AUC-ROC) as a metric to measure
144 how well a feature distinguished between phenotype genes and pseudogenes, which ranges
145 between 0.5 (random guessing) and 1 (perfect separation of functional and nonfunctional
146 sequences). Among the seven feature categories, transcription activity features were highly
147 informative (median AUC-ROC=0.88; **Fig. 2A**). Despite the strong performance of transcription
148 activity-related features, the presence of expression (i.e. transcript evidence) was a poor predictor
149 of functionality (AUC-ROC=0.58; **Fig. 2A**). This is because 80% of pseudogenes were
150 considered expressed in ≥ 1 of 51 RNA-seq datasets, demonstrating that presence of transcripts
151 should not be used by itself as evidence of functionality. Sequence conservation, DNA
152 methylation, TF binding, and H3 mark features were also fairly distinct between phenotype
153 genes and pseudogenes (median AUC-ROC ~ 0.7 for each category; **Fig. 2B-E**). By contrast,
154 chromatin accessibility and sequence-structure features were largely uninformative (median
155 AUC-ROC=0.51 and 0.55, respectively; **Fig. 2F,G**). We also observed high performance
156 variability within feature categories (see Supplementary Information). While many features are
157 distinct between phenotype genes and pseudogenes, functional predictions based on single
158 features yield high error rates (Supplementary Table 3; Supplementary Information), indicating a
159 need to jointly consider multiple features for distinguishing phenotype genes and pseudogenes.

160

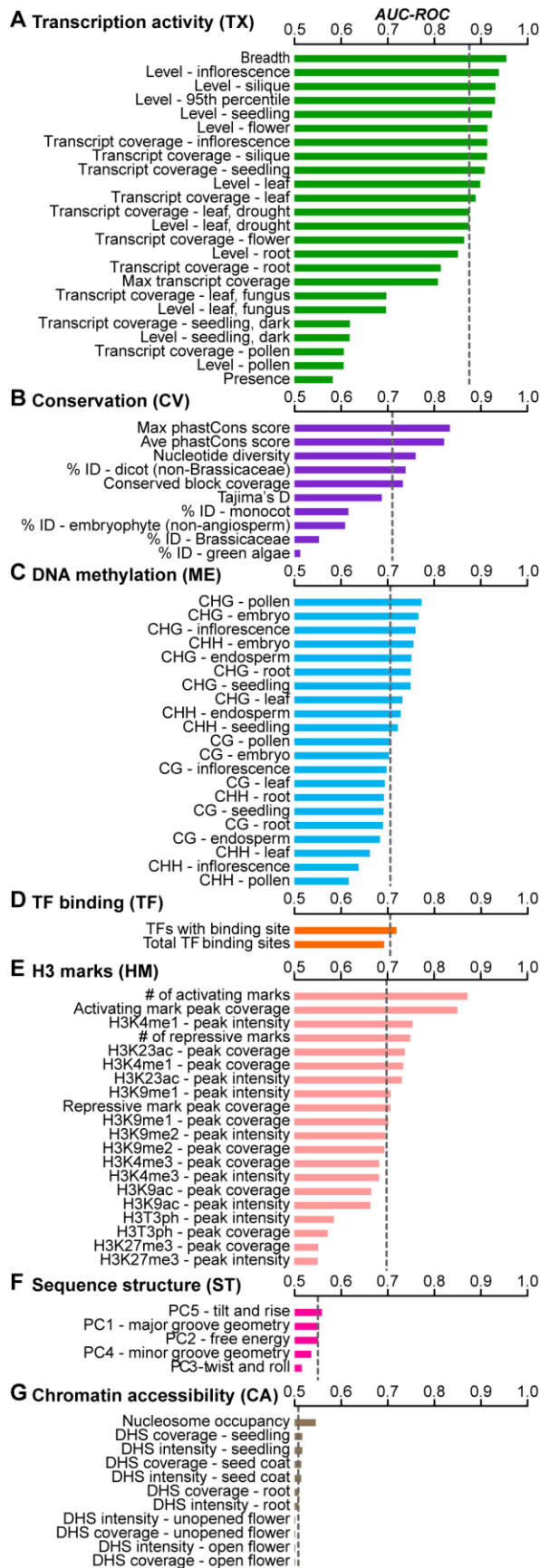


Figure 2. Predictions of functional (phenotype gene) and non-functional (pseudogene) sequences based on each individual feature. Prediction performance is measured using Area Under the Curve - Receiver Operating Characteristic (AUC-ROC). Features include those in the categories of (A) transcription activity, (B) sequence conservation, (C) DNA methylation, (D) transcription factor (TF) binding, (E) histone 3 (H3) marks, (F) sequence structure, and (G) chromatin accessibility. AUC-ROC ranges in value from 0.5 (equivalent to random guessing) to 1 (perfect predictions). Dotted lines: median AUC-ROC of features in a category.

162 **Consideration of multiple features produces accurate predictions of functional genomic**
163 **regions**

164 To consider multiple features in combination, we instead integrated all 93 features to
165 establish a machine learning model distinguishing phenotype gene and pseudogenes (referred to
166 as the full model; see Materials and Methods). The full model provided more accurate
167 predictions (AUC-ROC=0.98; False Negative Rate (FNR)=4%; False Positive Rate (FPR)=10%;
168 **Fig. 3A**) than any individual feature (**Fig. 2**; Supplementary Fig. 1, Supplementary Table 3). An
169 alternative measure of performance based on the precision (proportion of predicted functional
170 sequences that are truly functional) and recall (proportion of truly functional sequences predicted
171 correctly) values also indicated that the model was performing well (**Fig. 3B**). When compared
172 to the best-performing single feature (expression breadth), the full model had a similar FNR but
173 half the FPR (10% compared to 21%). Thus, the full model is highly capable of distinguishing
174 between phenotype genes and pseudogenes.

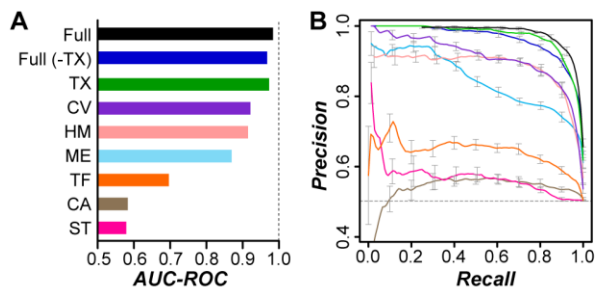
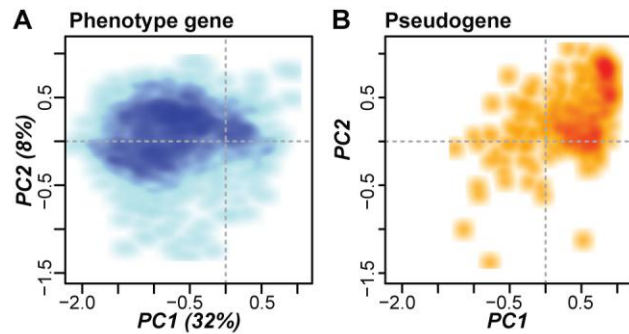


Figure 3. Predictions of functional and nonfunctional sequences based on multiple features. (*A*) AUC-ROC values of function prediction models built when considering all features (Full), all except transcription activity (TX)-related features (Full (-TX)), and all features from each category. The category abbreviations follow those in **Fig. 2**. (*B*) Precision-recall curves of the models with matching colors from (*A*). The models were built using feature values calculated from 500 bp sequence windows. We also conducted principle component (PC) analysis to investigate how well phenotype genes and pseudogenes could be separated and found that phenotype genes (Supplementary Fig. 1A) and pseudogenes (Supplementary Fig. 1B) were distributed in largely distinct space. However, there remained substantial overlap, indicating that standard parametric approaches are not well suited to distinguishing between benchmark functional and nonfunctional sequences.

175
176 We next determined the relative contributions of different feature categories in predicting
177 phenotype genes and pseudogenes and established seven prediction models each using only the
178 subset of features from a single category (**Fig. 2**). Although none of these category-specific
179 models had performance as high as the full model (**Fig. 3C**), the transcription activity feature



Supplementary Figure 1. Smoothed scatterplots of the first two principle components (PCs) of (A) phenotype gene and (B) pseudogene features. The percentages on the axes in (A) indicate the feature value variation explained by the associated PC.

180
181 category model performed almost as well as the full model (AUC-ROC=0.97, FNR=6%,
182 FPR=12%). Instead of the presence of expression evidence, the breadth and level of transcription
183 are the causes of the strong performance of the transcription activity-only model. We also found
184 that a model excluding transcription activity features (full (-TX), **Fig. 3C,D**) performed almost as
185 well as the full model and similarly to the transcription activity-feature-only model, but with an
186 increased FPR (AUC-ROC=0.96; FNR=3%; FPR=20%). These findings indicate that a diverse
187 array of features can be considered jointly to make highly accurate predictions of the
188 functionality of a genomic sequence. Meanwhile, our finding of the high performance of the
189 transcription activity-only model highlights the possibility of establishing an accurate model for
190 functional prediction in species with only a modest amount of transcriptome data.

191 **The functional likelihood measure can be used to classify functional and non-functional** 192 **sequences**

193 To provide a measure of the potential functionality of any sequence in the *A. thaliana*
194 genome, including ITRs and ncRNAs, we utilized the confidence score from the full model as a
195 “functional likelihood” value (Tsai et al. 2017). The functional likelihood (FL) score ranges
196 between 0 and 1, with high values indicating that a sequence is more similar to phenotype genes
197 (functional) and low values indicating a sequence more closely resembles pseudogenes
198 (nonfunctional). FL values for all genomic regions examined in this study are available in
199 Supplementary Table 4. As expected, phenotype genes had high FL values (median=0.97; **Fig.**
200 **4A**) and pseudogenes had low values (median=0.01; **Fig. 4B**). To call sequences as functional or
201 not, we defined a threshold FL value (0.35) by maximizing the F-measure (see Materials and

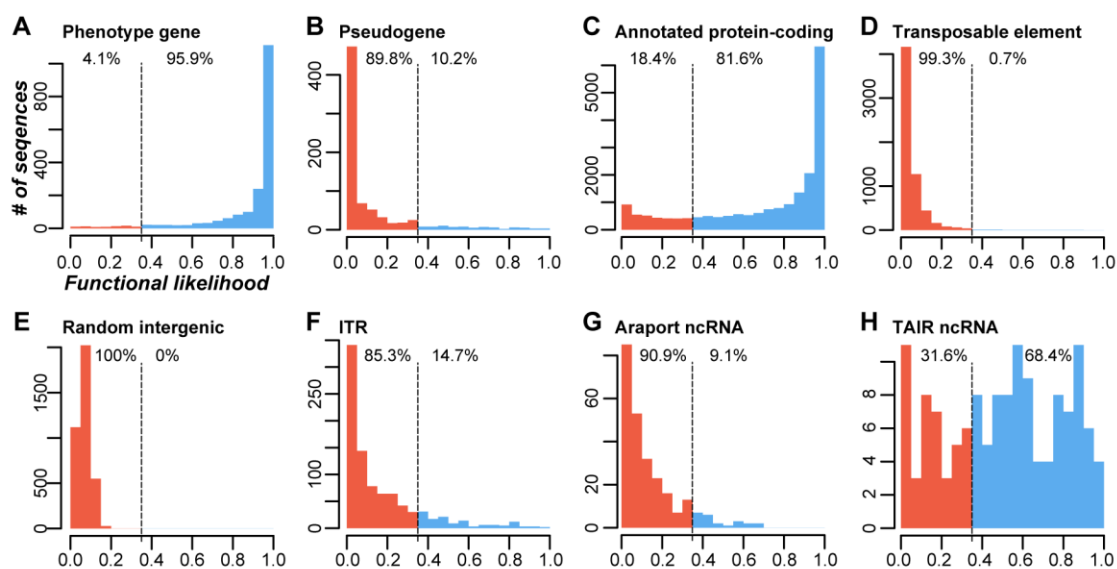


Figure 4. Functional likelihood distributions of various sequence classes based on the full model. (A) Phenotype genes. (B) Pseudogenes. (C) Annotated protein-coding genes. (D) Transposable elements. (E) Random unexpressed intergenic sequences. (F) Intergenic transcribed regions (ITR). (G) Araport11 ncRNAs. (H) TAIR10 ncRNAs. The full model was established using 500 bp sequence windows. Higher and lower functional likelihood values indicate greater similarity to phenotype genes and pseudogenes, respectively. Vertical dashed lines indicate the threshold for calling a sequence as functional or nonfunctional. The percentages to the left and right of the dashed line indicate the percent of sequences predicted as functional or nonfunctional, respectively.

202

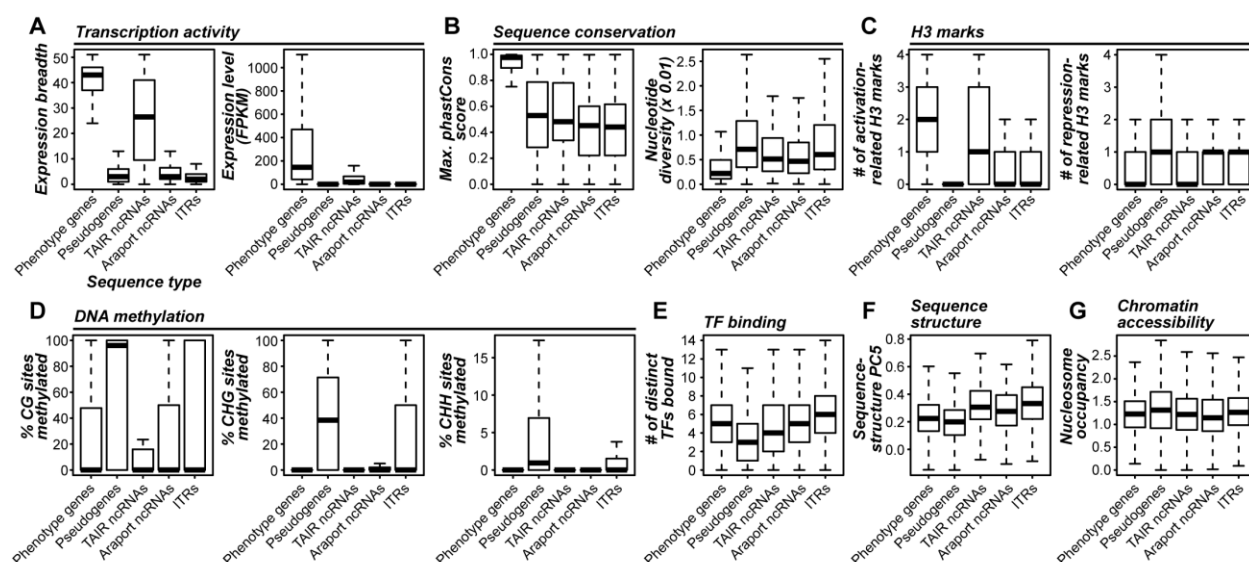
203 Methods). Using this threshold, 96% of phenotype genes (**Fig. 4A**) and 90% of pseudogenes
 204 (**Fig. 4B**) are correctly classified as functional and nonfunctional, respectively, demonstrating
 205 that the full model is highly capable of distinguishing functional and nonfunctional sequences.

206 We next applied our model to predict the functionality of annotated protein-coding genes,
 207 transposable elements (TEs), and unexpressed intergenic regions. Most annotated protein-coding
 208 genes not included in the phenotype gene dataset had high FL scores (median=0.86; **Fig. 4C**) and
 209 80% were predicted as functional. The features exhibited by low-scoring protein-coding genes
 210 and high-scoring pseudogenes are discussed in the Supplementary Information. By contrast, the
 211 FLs were low for both TEs (median=0.03, **Fig. 4D**) and unexpressed intergenic regions
 212 (median=0.07; **Fig. 4E**), and 99% of TEs and all unexpressed intergenic sequences were
 213 predicted as nonfunctional. Overall, the FL measure provides a useful metric to distinguish
 214 between phenotype genes and pseudogenes. In addition, the FLs of annotated protein-coding
 215

216 genes, TEs, and unexpressed intergenic sequences agree with *a priori* expectations regarding the
 217 functionality of these sequences.

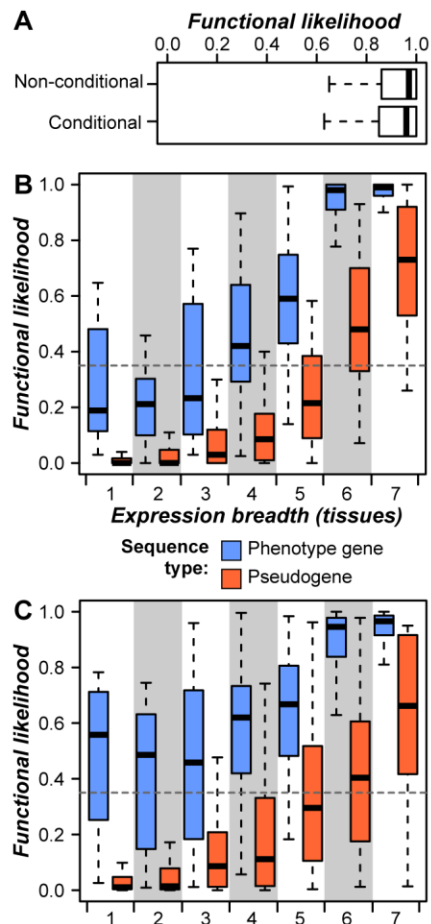
218 Most ITRs and annotated ncRNAs do not resemble benchmark phenotype genes

219 We next evaluated functional predictions of ITRs and ncRNAs. Consistent with previous
 220 studies (Moghe et al. 2013), ITRs and ncRNAs in our dataset were more narrowly and weakly
 221 expressed and less conserved compared to phenotype genes (Supplementary Fig. 2A,B). In
 222 addition, ITRs in particular had biochemical characteristics that were generally more similar to
 223 pseudogenes (Supplementary Fig. 2C-F). Given the association between transcription activity



224 **Supplementary Figure 2.** Distributions of 12 example features from (A) transcription activity, (B)
 225 sequence conservation, (C) histone 3 (H3) marks, (D) DNA methylation, (E) transcription factor (TF)
 226 binding, (F) sequence structure, and (G) chromatin accessibility feature categories (Fig. 2). Feature
 227 distributions are shown for phenotype genes, pseudogenes, TAIR- and Araport-annotated ncRNAs and
 228 intergenic transcribed regions (ITR).

229 features and functional predictions (**Fig 2A; Fig. 3A**), we investigated how functional prediction
 230 models performed for conditionally-functional and narrowly-expressed sequences before
 231 applying them to ITRs and ncRNAs. We found that genes with conditional phenotypes had no
 significant differences in FLs (median=0.96) as those with phenotypes under standard growth
 conditions (median=0.97; U test, $p=0.38$, Supplementary Fig. 3A), indicating the full model can
 capture conditionally functional sequences. However, the full model is biased against narrowly-
 expressed (≤ 3 tissues) phenotype genes as 65% of them were predicted as nonfunctional



Supplementary Figure 3. Impacts of conditional phenotypes and expression breadth on the function prediction model. (A) Functional likelihood distributions of phenotype genes with mutant phenotypes under standard growth conditions (non-conditional) and non-standard growth conditions such as stressful environments (conditional) based on the 500 bp full model. Feature values were calculated from a random 500 bp region from within the sequence body. Higher and lower functional likelihood values indicate a greater similarity to phenotype genes and pseudogenes, respectively. (B,C) Distributions of functional likelihood scores for phenotype genes (blue) and pseudogenes (red) for sequences with various breadths of expression for (B) the 500 bp full model and (C) the 500 bp tissue-agnostic model generated by excluding the expression breadth and features available from multiple tissues. The tissue-agnostic model is aimed toward minimizing the effects of biochemical activity occurring across multiple tissues and predicts a greater proportion of narrowly-expressed phenotype genes as functional compared to the full model.

232

233 (Supplementary Fig. 3B). Further, pseudogenes that were more highly and broadly expressed

234 were disproportionately predicted as functional (**Fig. 5**; Supplementary Fig. 3B). To tailor

235 functional predictions to narrowly-expressed sequences, particularly ITRs and ncRNAs, we

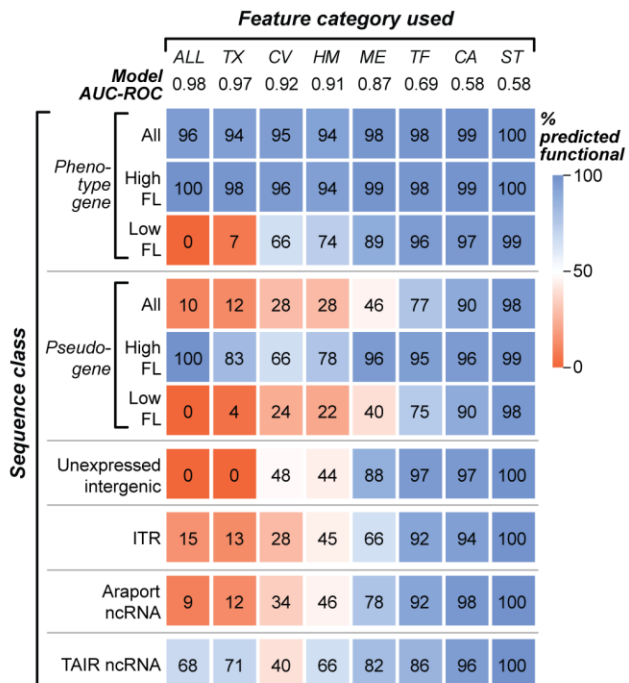


Figure 5. Proportion of phenotype genes, pseudogenes, ITRs, and ncRNAs predicted as functional in the full and single-category models. Percentages of sequence classes that are predicted as functional in models based on all features and the single category models, each using all features from a category (abbreviated according to **Fig. 2**). The models are sorted from left to right based on performance (AUC-ROC). The colors of and numbers within the blocks indicate the proportion sequences predicted as functional by a given model. Phenotype gene and pseudogene sequences are shown in three sub-groups: all sequences (All), and those predicted as functional (high functional likelihood (FL)) and nonfunctional (low FL) in the full model. ITR: intergenic transcribed regions. A greater proportion of ITRs and Araport ncRNAs are predicted as functional when considering only DNA methylation or H3 mark features compared to the full (**Fig. 3**) or tissue-agnostic (Supplementary Fig. 5) models. However, these two category-specific models also had higher false positive rates (unexpressed intergenic sequences and pseudogenes). Thus, these single feature-category models do not provide additional support for the functionality of most Araport ncRNAs and ITRs.

236

237

238

239

240

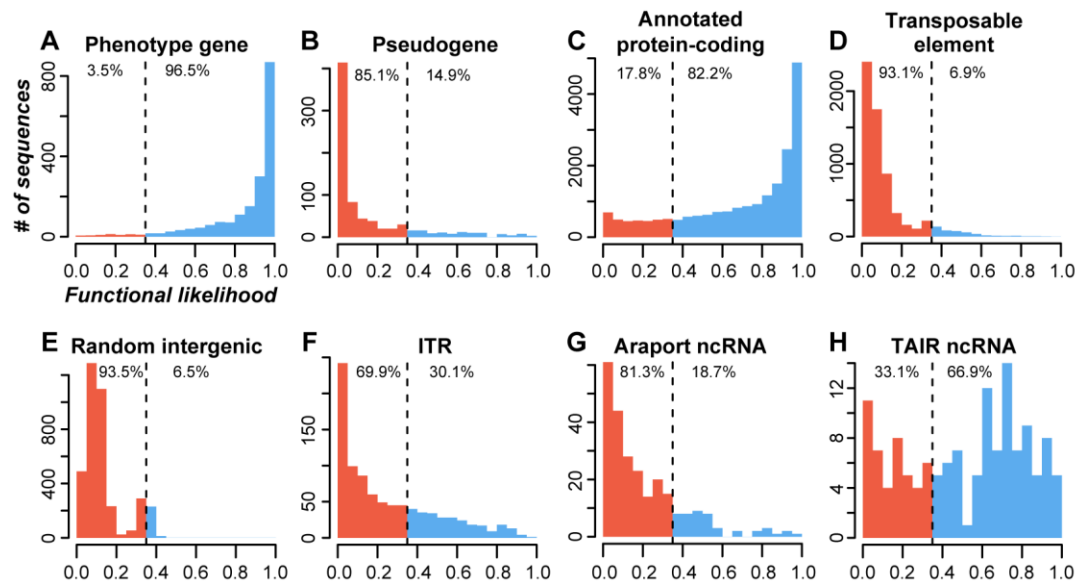
241

242

243

generated a “tissue-agnostic” model by excluding expression breadth and features available across multiple tissues (see Materials and Methods). This tissue-agnostic model performed similarly to the full model (AUC-ROC=0.97; FNR=4%; FPR=15%; Supplementary Fig. 4; Supplementary Table 4), although there was a 5% increase in FPR (from 10% to 15%).

Importantly, the proportion of phenotype genes expressed in ≤ 3 tissues predicted as functional increased by 23% (35% in the full model to 58% in the tissue-agnostic model; Supplementary



Supplementary Figure 4. Distributions of functional likelihood scores based on the 500 bp tissue-agnostic model. (A) Phenotype genes. (B) Pseudogenes. (C) Annotated protein-coding genes. (D) Transposable elements. (E) Random unexpressed intergenic sequences. (F) Intergenic transcribed regions (ITR). (G) Araport11 ncRNAs. (H) TAIR10 ncRNAs. Vertical dashed lines display the threshold to define a sequence as functional or nonfunctional. The numbers to the left and right of the dashed line show the percentage of sequences predicted as functional or nonfunctional, respectively.

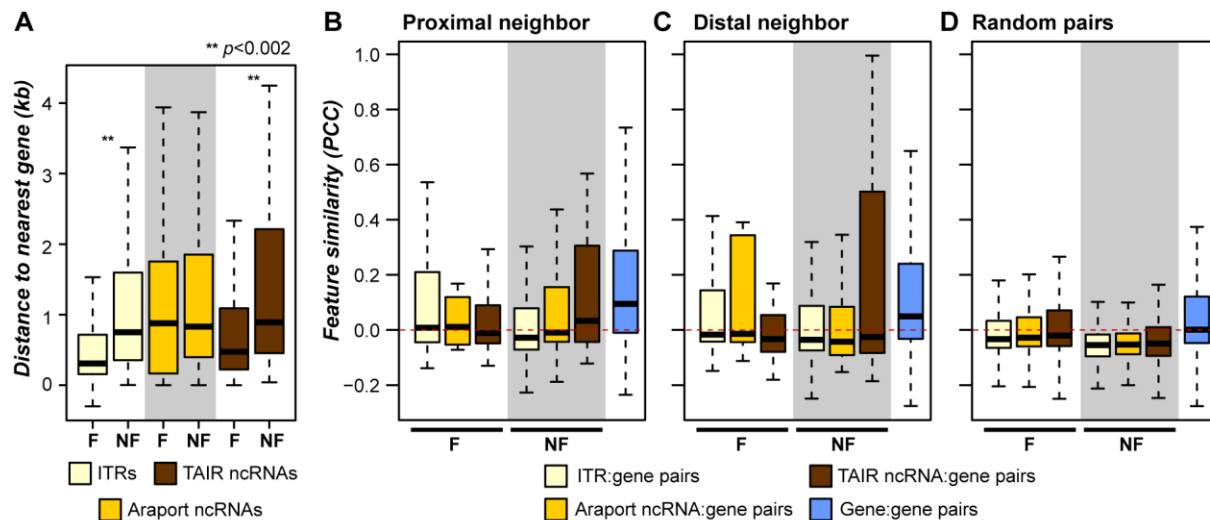
244

245 Fig. 3C), indicating that the tissue-agnostic model is more suitable for predicting the

246 functionality of narrowly-expressed sequences than the full model.

247

248 We next applied both the full and tissue-agnostic models to 895 ITRs, 136 ncRNAs
 249 annotated by The Arabidopsis Information Resource (TAIR), and 252 long ncRNAs annotated
 250 by the Araport database that do not overlap with any other annotated genome features. The
 251 median FLs based on the full model were low (0.09) for both ITRs (**Fig. 4F**) and Araport
 252 ncRNAs (**Fig. 4G**), and only 15% and 9% of these sequences were predicted as functional,
 253 respectively. By contrast, TAIR ncRNAs had a significantly higher median FL value (0.53; U
 254 tests, both $p < 5e-31$; **Fig. 4H**) and 68% were predicted as functional, which is best explained by
 255 differences in features from the transcription activity category (**Fig. 5**). We also note that ITRs
 256 and annotated ncRNAs that were close to genes were more frequently predicted as functional,
 257 suggesting a subset may represent unannotated exons of known genes (Supplementary Fig. 5;
 258 Supplementary Information). Next we applied the tissue-agnostic model to ITRs and
 TAIR/Araport ncRNAs. Compared to the full model, around twice as many ITRs (30%) and



Supplementary Figure 5. Distance of ITRs and annotated ncRNA regions to and feature similarity with neighboring genes. (A) Distance from intergenic transcribed regions (ITRs) and annotated ncRNAs to the closest neighboring gene. ITR and ncRNA sequences are separated by whether they are predicted as functional (F) or nonfunctional (NF) by the 500 bp full model. (B) Feature similarity between proximal neighbors (within 95th percentile (456 bp) of intron lengths), and (C) distal neighbors (>456 bp). Pairs involving ITRs and annotated ncRNAs were divided by whether the ITR or ncRNA sequence was predicted as functional (F) or nonfunctional (NF) by the full model. Feature values were quantile normalized prior to calculating correlations. (D) Feature similarity based on Pearson's Correlation Coefficients (PCC) between random pairs of ITRs, Araport11 ncRNAs, TAIR10 ncRNAs, or annotated genes.

259

260 Araport ncRNAs (19%) but a similar number of TAIR ncRNA (67%) were predicted as
 261 functional. Considering the union of the full and tissue-agnostic model predictions, 268 ITRs
 262 (32%), 57 Araport ncRNAs (23%), and 105 TAIR ncRNAs (77%) were likely functional. Thus,
 263 the majority of ITRs and Araport ncRNAs are more similar to pseudogenes than to phenotype
 264 genes that are predominantly protein coding.

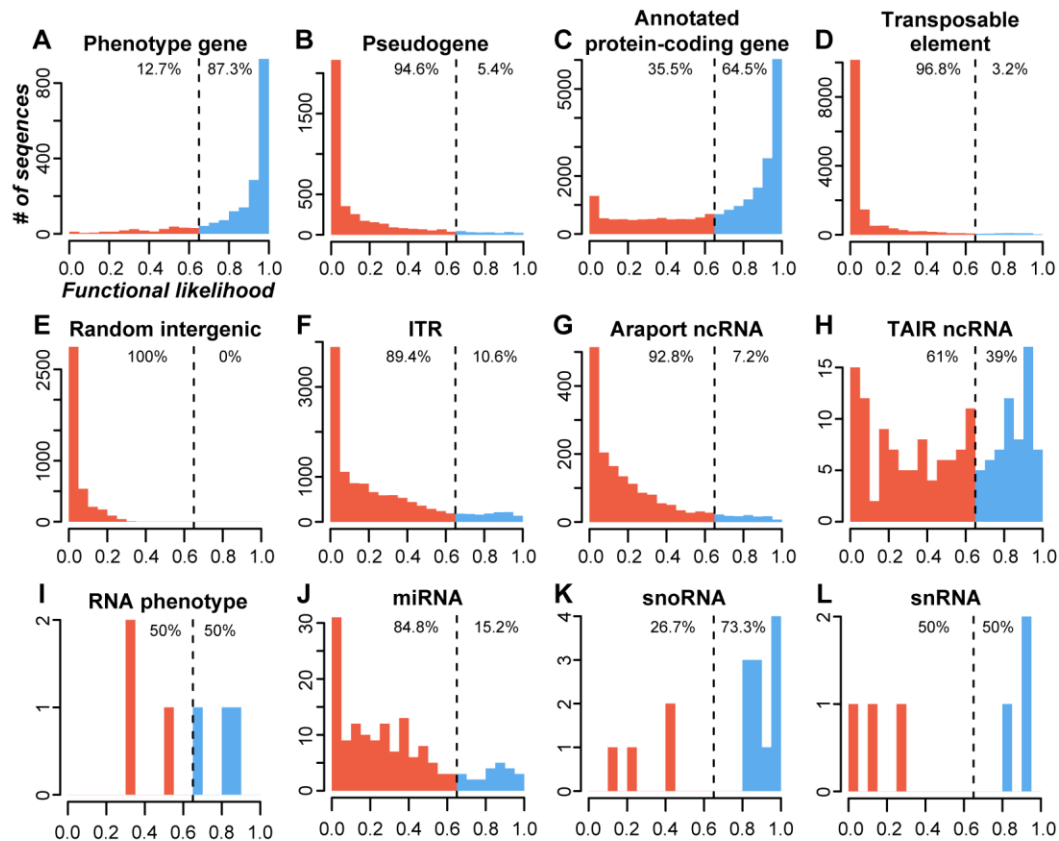
265 **Benchmark protein-coding and RNA genes exhibit distinct characteristics**

266 We demonstrated that the majority of ITR and annotated ncRNA sequences do not
 267 exhibit characteristics of benchmark phenotype genes. Note that the phenotype genes are
 268 predominantly protein coding. Although the features utilized to generate functional predictions
 269 were not exclusive to protein-coding sequences, RNA genes may exhibit a distinct feature profile
 270 from protein-coding genes. The full and tissue agnostic models described above were established

271 with 500 bp windows and most known RNA genes are too short to be considered by these
272 models. Thus, to evaluate functional predictions among annotated RNA genes, we generated a
273 new tissue-agnostic model using 100 bp sequences (for features, see Supplementary Table 6) that
274 performed similarly to the full 500 bp model, except with 9% higher FNR (AUC-ROC=0.97;
275 FNR=13%; FPR=5%; Supplementary Fig. 6). With this new tissue agnostic model, 50% (three
276 out of six) of RNA genes with documented mutant phenotypes (phenotype RNA genes) were
277 predicted as functional (Supplementary Fig. 6I). We also applied this model to other RNA Pol II-
278 transcribed RNA genes (without documented phenotypes) and found that 15% of microRNA
279 (miRNA) primary transcripts (Supplementary Fig. 6J), 73% of small nucleolar RNAs (snoRNAs;
280 Supplementary Fig. 6K), and 50% of small nuclear RNAs (snRNAs; Supplementary Fig. 6L)
281 were predicted as functional. Although the proportion of phenotype RNA genes predicted as
282 functional (50%) is significantly higher than the proportion of pseudogenes predicted as
283 functional (5%, FET, $p < 0.004$), this finding suggests that a model trained using protein-coding
284 genes has a substantial FNR for detecting RNA genes.

285 To determine whether the suboptimal predictions by the phenotype protein-coding gene-
286 based models are because RNA genes belong to a class of their own, we next built a multi-class
287 function prediction model (as opposed to the binary, two-class models described above) aimed at
288 distinguishing four classes of sequences: benchmark RNA genes (n=46, Supplementary Table 6),
289 , phenotype protein-coding genes (1,882), pseudogenes (3,916), and randomly-selected,
290 unexpressed intergenic regions (4,000). In the four-class model, 87% of benchmark RNA genes,
291 including all six phenotype RNA genes, were predicted as functional sequences (65% RNA
292 gene-like and 22% phenotype protein-coding gene-like; **Fig. 6A**). In addition, 95% of phenotype
293 protein-coding genes were predicted as functional (**Fig. 6B**), including 80% of narrowly
294 expressed genes, an increase of 22% over the 500 bp tissue-agnostic model (Supplementary Fig.
295 3C). For benchmark non-functional sequences, 70% of pseudogenes (**Fig. 6C**) and 100% of
296 unexpressed intergenic regions (**Fig. 6D**) were predicted as nonfunctional (either as pseudogenes
297 or unexpressed intergenic sequences). Overall, the four-class model improves prediction
298 accuracy of RNA genes and narrowly expressed genes. In addition, given the 30% FPR among
299 pseudogenes, the four-class model provides a liberal estimate of sequence functionality and high
300 confidence estimate of non-functionality.

301



Supplementary Figure 6. Distributions of functional likelihood scores based on the 100 bp tissue-agnostic model. (A) Phenotype genes. (B) Pseudogenes. (C) Protein-coding gene. (D) transposable elements. (E) Random unexpressed intergenic sequences. (F) Intergenic transcribed regions (ITR). (G) Araport11 ncRNAs. (H) TAIR10 ncRNAs. (I) RNA genes with loss-of-function mutant phenotypes. (J) MicroRNAs, (K) Small nucleolar RNAs, (L) Small nuclear RNAs. The tissue-agnostic model was built with 100 bp features and while excluding the expression breadth and tissue-specific features as annotated RNA genes tend to be more narrowly expressed than phenotype genes (U tests, all $p < 2e-05$; Supplementary Fig. 2A). Higher functional likelihood values indicate greater similarity to phenotype genes while lower values indicate similarity to pseudogenes. Vertical dashed lines display the threshold to define a sequence as functional or nonfunctional. The numbers to the left and right of the dashed line show the percentage of sequences predicted as functional or nonfunctional, respectively.

302

303 **Most intergenic transcribed regions and annotated ncRNAs do not resemble benchmark** 304 **RNA genes**

305 By applying the four-class model on ITRs and annotated ncRNAs, we found that 34% of
306 ITRs, 38% of Araport ncRNAs, and of 65% TAIR ncRNAs were predicted as functional
307 sequences (**Fig. 6E-G**). Specifically, $\leq 20\%$ of ITR and annotated ncRNA sequences were

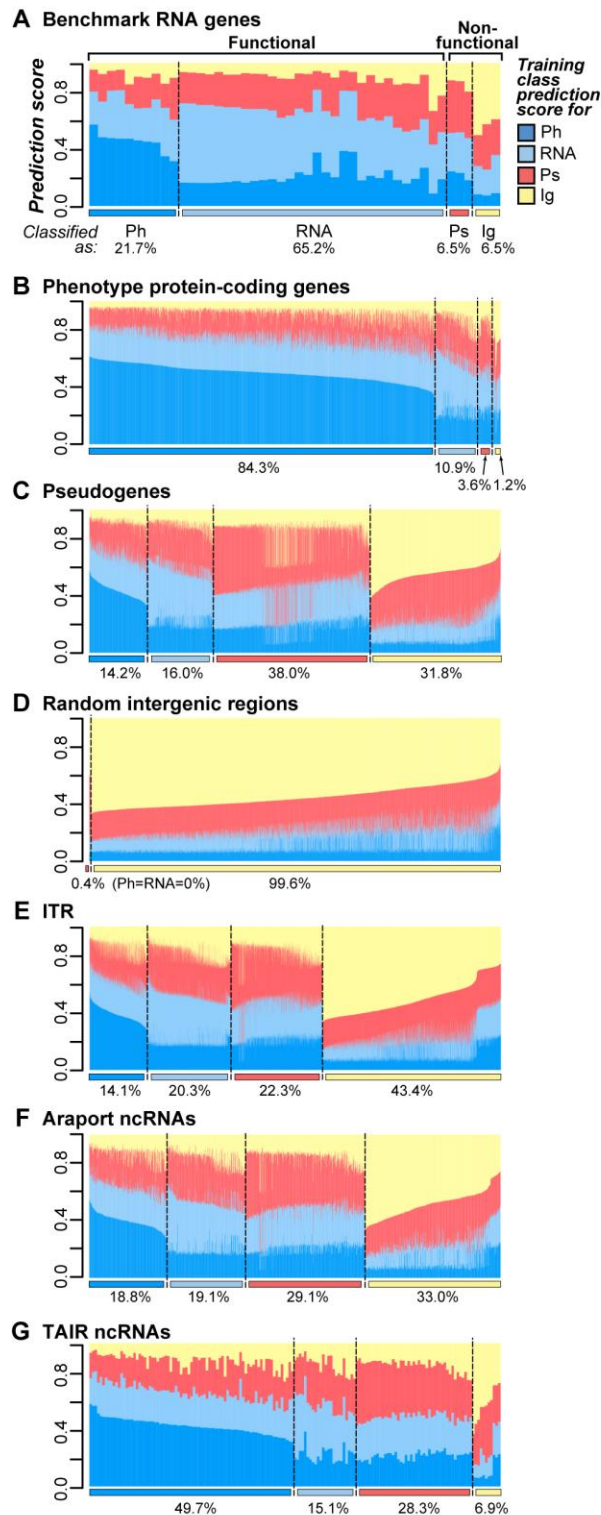


Figure 6. Function predictions based on a four-class prediction model. **(A)** Stacked bar plots indicate the prediction scores of benchmark RNA genes for each of the four classes: dark blue - phenotype protein-coding gene (Ph), cyan - RNA gene (RNA), red - pseudogene (Ps), yellow – random intergenic sequence (Ig). Ig were included to provide another set of likely nonfunctional sequences distinct from pseudogenes. Expression breadth and tissue-specific features were excluded and 100 bp sequences were used. A benchmark RNA gene is classified as one of the four classes according to the highest prediction score. The color bars below the chart indicate the predicted class, with the same color scheme as the prediction score. Sequences classified as Ph or RNA were considered functional, while those classified as Ps or Ig were considered nonfunctional. Percentages below a classification region indicate the proportion of sequences classified as that class. **(B)** Phenotype protein-coding gene prediction scores. **(C)** Pseudogene prediction scores. **(D)** Random unexpressed intergenic region prediction scores. Note that no sequence was predicted as functional. **(E)** Intergenic transcribed region (ITR), **(F)** Araport11 ncRNA regions. **(G)** TAIR10 ncRNA regions. Note that the 100 bp model used here allowed us to evaluate an additional 10,938 ITRs and 1,406 annotated ncRNAs compared to the 500 bp full and tissue-agnostic models.

308

309 classified as RNA genes (**Fig. 6E-G**). Although miRNAs dominate the benchmark RNA

310 sequences, we should emphasize that the four-class prediction model also increased the

311 proportion of snoRNAs and snRNAs that were predicted as functional (91%) compared to the

312 100bp tissue agnostic model (67%). Thus a lack of similarity to benchmark miRNAs provides
313 evidence that most ITRs and Araport ncRNAs are not functioning as RNA genes.

314 To provide an overall estimate of the proportion of likely functional and nonfunctional
315 ITRs and annotated ncRNAs, we considered the predictions from the four-class model (**Fig. 6**),
316 the full model (**Fig. 3,4**), and the tissue-agnostic models (Supplementary Fig. 4, 6), which cover
317 both protein-coding and RNA gene functions. Based on support from ≥ 1 of the four models, we
318 classified 4,437 ITRs (38%) and 796 annotated ncRNAs (44%) as likely functional, as they
319 resembled either phenotype protein-coding or RNA genes. Although our findings lend support
320 that they are likely parts of novel or annotated genes, we should stress that, given the relatively
321 high FPR in the tissue-agnostic and four-class models, the estimate of functional ITR/ncRNA is
322 a liberal one. Most importantly, we find that a substantial number of ITRs (62%) and annotated
323 ncRNAs (56%) are predicted as nonfunctional. Moreover, at least a third of ITRs (**Fig. 6E**) and
324 Araport ncRNAs (**Fig. 6F**) most closely resemble unexpressed intergenic regions. Based on these
325 findings, we conclude that the majority of ITRs and annotated ncRNA regions resemble
326 nonfunctional genomic regions, and therefore are derived from noisy transcription.

327 **CONCLUSION**

328 Discerning the location of functional regions within a genome represents a key goal in
329 genomic biology and is fundamental to molecular evolutionary studies. Despite advances in
330 computational gene finding, it remains challenging to determine whether ITRs represent
331 functional or noisy biochemical activity. We established robust function prediction models based
332 on the evolutionary, biochemical, and structural characteristics of phenotype genes and
333 pseudogenes in *A. thaliana*. The prediction models accurately define functional and
334 nonfunctional regions and are applicable genome-wide and echo recent findings using human
335 data to evaluate RNA gene functionality (Tsai et al. 2017). We utilized prediction models to
336 assess the functionality of both protein-coding and annotated RNA genes. As benchmark
337 examples of more recently identified RNA gene classes become available in *A. thaliana*, such as
338 *cis*-acting regulatory (Guil and Esteller 2012) or competitive endogenous (Tan et al. 2015)
339 RNAs, it will be interesting to see if sequences that encode RNA products with these roles can be
340 predicted as functional based on a similar predictive framework. Given that function predictions
341 were successful in both plants and metazoans, integrating the evolutionary and biochemical

342 features of known genes for functional genomic region prediction will likely be applicable to any
343 species. The next step will be to test whether function prediction models can be applied across
344 species, which could ultimately allow the phenotype data and omics resources available in model
345 systems to effectively guide the identification of functional regions in non-models.

346 Expression data was highly informative to functional predictions. We found that the
347 prediction model based on only 24 transcription activity-related features performs nearly as well
348 as the full model that integrates additional information including conservation, H3 mark,
349 methylation, and TF binding data. In human, use of transcription data from cell lines also
350 produced highly accurate predictions of functional genomic regions (Tsai et al. 2017).
351 Importantly, our findings suggest that function prediction models can be established in any
352 species, model or not, with a modest number of transcriptome datasets (e.g. 51 in this study and
353 19 in human). Despite the importance of transcription data, we emphasize that the presence of
354 expression evidence is an extremely poor predictor of functionality. With the effectiveness of our
355 model noted, one major caveat of our models is that narrowly-expressed phenotype genes are
356 frequently predicted as pseudogene and broadly-expressed pseudogenes tend to be called
357 functional. To improve the function prediction model, it will be important to explore additional
358 features unrelated to transcription, particularly those relevant to broadly expressed pseudogenes
359 that are mostly likely recently pseudogenized. In addition, because few phenotype genes are
360 narrowly-expressed (5%) in the *A. thaliana* training data, more phenotyping data for narrowly
361 expressed genes will be crucial as well.

362 Upon application of the function prediction models genome-wide, 4,427 ITRs and 796
363 annotated ncRNAs in *A. thaliana* are predicted as functional sequences. However, considering
364 the high false positive rates (e.g. 10% for the full and 31% for the four-class model), this is most
365 likely an overestimate of the functional sequences contributed by ITRs and annotated ncRNAs.
366 While we err on the side of calling non-functional sequences as functional, we reduce the error
367 rate for calling a functional sequence as non-functional. Despite this conservative approach to
368 classifying sequences as non-functional, the majority of ITRs and ncRNAs resemble
369 pseudogenes and random unexpressed intergenic regions. Similarly, most human ncRNAs are
370 more similar to nonfunctional sequences than they are to protein coding and RNA genes (Tsai et
371 al. 2017). Together with our finding of a significant relationship between the amount of
372 intergenic expression and genome size, we conclude that a significant proportion of intergenic

373 transcripts are nonfunctional noise. Thus, instead of assuming any expressed sequence must be
374 functionally significant, we advocate that the null hypothesis should be that it is not, particularly
375 considering that most ITRs and annotated ncRNAs have not been experimentally characterized.

376 **MATERIALS AND METHODS**

377 **Identification of transcribed regions in leaf tissue of 15 flowering plants**

378 RNA-sequencing (RNA-seq) datasets were retrieved from the Sequence Read Archive
379 (SRA) at the National Center for Biotechnology Information (NCBI;
380 www.ncbi.nlm.nih.gov/sra/) for 15 flowering plant species (Supplementary Table 1). All datasets
381 were generated from leaf tissue and sequenced on Illumina HiSeq 2000 or 2500 platforms.
382 Genome sequences and gene annotation files were downloaded from Phytozome v.11
383 (www.phytozome.net) (Goodstein et al. 2012) or Oropetium Base v.01 (www.sviridis.org)
384 (VanBuren et al. 2015). Genome sequences were repeat masked using RepeatMasker v4.0.5
385 (www.repeatmasker.org) if a repeat-masked version was not available. Only one end from
386 paired-end read datasets were utilized in downstream processing. Reads were trimmed to be rid
387 of low scoring ends and residual adaptor sequences using Trimmomatic v0.33 (LEADING:3
388 TRAILING:3 SLIDINGWINDOW:4:20 MINLEN:20) (Bolger et al. 2014) and mapped to
389 genome sequences using TopHat v2.0.13 (default parameters except as noted below) (Kim et al.
390 2013). Reads ≥ 20 nucleotides in length that mapped uniquely within a genome were used in
391 further analysis.

392 For each species, thirty million mapped reads were randomly selected from among all
393 datasets and assembled into transcript fragments using Cufflinks v2.2.1 (default parameters
394 except as noted below) (Trapnell et al. 2010), while correcting for sequence-specific biases
395 during the sequencing process by providing an associated genome sequence with the -b flag. The
396 expected mean fragment length for assembled transcript fragments in Cufflinks was set to 150
397 from the default of 200 so that expression levels in short fragments would not be overestimated.
398 The 1st and 99th percentile of intron lengths for each species were used as the minimum and
399 maximum intron lengths, respectively, for both the TopHat2 and Cufflinks steps. Intergenic
400 transcribed regions (ITRs) were defined by transcript fragments that did not overlap with gene
401 annotation and did not have significant six-frame translated similarity to plant protein sequences
402 in Phytozome v.10 (BLASTX E-value < 1E-05). The correlation between assembled genome

403 size and gene counts was determined with data from the first 50 published plant genomes
404 (Michael and Jackson 2013).

405 **Phenotype data sources**

406 Mutant phenotype data for *A. thaliana* protein-coding genes was collected from a
407 published dataset (Lloyd and Meinke 2012), the Chloroplast 2010 database (Ajjawi et al. 2010;
408 Savage et al. 2013), and the RIKEN phenome database (Kuromori et al. 2006) as described by
409 Lloyd et al. (2015). Phenotype genes used in our analyses were those whose disruption resulted
410 in lethal or visible defects under standard laboratory growth conditions. Genes with documented
411 mutant phenotypes under standard conditions were considered as a distinct and non-overlapping
412 category from other annotated protein-coding genes. We identified six RNA genes with
413 documented loss-of-function phenotypes through literature searches (Supplementary Table 7):
414 *At4* (AT5G03545) (Shin et al. 2006), *MIR164A* and *MIR164D* (AT2G47585 and AT5G01747,
415 respectively) (Guo et al. 2005), *MIR168A* (AT4G19395) (W. Li et al. 2012), and *MIR828A* and
416 *TAS4* (AT4G27765 and AT3G25795, respectively) (Hsieh et al. 2009). Conditional phenotype
417 genes were those belonging to the Conditional phenotype group as described by Lloyd and
418 Meinke (2012). Loss-of-function mutants of these genes exhibited phenotype only under stress
419 conditions.

420 **Arabidopsis thaliana genome annotation**

421 *A. thaliana* protein-coding gene, miRNA gene, snoRNA gene, snRNA gene, ncRNA
422 region, pseudogene, and transposable element annotations were retrieved from The Arabidopsis
423 Information Resource v.10 (TAIR10; www.arabidopsis.org) (Berardini et al. 2015). Additional
424 miRNA gene and lncRNA region annotations were retrieved from Araport v.11
425 (www.araport.org). A primary difference between the TAIR ncRNAs and Araport lncRNAs
426 (referred to as Araport ncRNAs in the Results & Discussion section) is the date in which they
427 were annotated. For example, 221 ncRNAs were present in the v.7 release of TAIR, which dates
428 back to 2007 (TAIR10 contains 394 ncRNA annotations) (Swarbreck et al. 2008; Lamesch et al.
429 2012; Berardini et al. 2015). However, Araport lncRNAs were annotated in the past five years
430 (Krishnakumar et al. 2015). Thus, that TAIR ncRNAs are generally more highly and broadly
431 expressed is likely a result of the less sensitive transcript identification methods available for
432 early TAIR releases. A pseudogene-finding pipeline (Zou et al. 2009) was used to identify

433 additional pseudogene fragments and count the number of disabling mutations (premature stop or
434 frameshift mutations). Genes, pseudogenes, and transposons with overlapping annotation were
435 excluded from further analysis. Overlapping lncRNA annotations were merged for further
436 analysis. When pseudogenes from TAIR10 and the pseudogene-finding pipeline overlapped, the
437 longer pseudogene annotation was used.

438 *A. thaliana* ITRs analyzed include: (1) the Set 2 ITRs in Moghe et al. (Moghe et al.
439 2013), (2) the novel transcribed regions from Araport v.11, and (3) additional ITRs from 206
440 RNA-seq datasets (Supplementary Table 5). Reads were trimmed, mapped, and assembled into
441 transcript fragments as described above, except that overlapping transcript fragments from across
442 datasets were merged. ITRs analyzed did not overlap with any TAIR10, Araport11, or
443 pseudogene annotation. Overlapping ITRs from different annotated subsets were kept based on a
444 priority system: Araport11 > Set 2 ITRs from Moghe et al. (Moghe et al. 2013) > ITRs identified
445 in this study. For each sequence entry (gene, ncRNA, pseudogene, transposable element, or
446 ITR), a 100 and 500 base pair (bp) window was randomly chosen for calculating feature values
447 and subsequent model building steps. Feature descriptions are provided in the following sections.
448 The feature values for randomly selected 500 and 100 bp windows are provided in
449 Supplementary Tables 2 and 6, respectively. Additionally, non-expressed intergenic sequences
450 were randomly-sampled from genome regions that did not overlap with annotated genes,
451 pseudogenes, transposable elements, or regions with genic or intergenic transcript fragments
452 (100 bp, n=4,000; 500 bp, n=3,716). All 100 and 500 bp windows described above are referred
453 to as sequence windows throughout the Methods section.

454 **Sequence conservation and structure features**

455 There were 10 sequence conservation features examined. The first two were derived from
456 comparisons between *A. thaliana* accessions including nucleotide diversity and Tajima's D
457 among 81 accessions (Cao et al. 2011) using a genome matrix file from the 1,001 genomes
458 database (www.1001genomes.org). The python scripts are available through GitHub
459 (<https://github.com/ShiuLab/GenomeMatrixProcessing>). The remaining eight features were
460 derived from cross-species comparisons, three based on multiple sequence and five based on
461 pairwise alignments. Three multiple sequence alignment-based features were established using
462 aligned genomic regions between *A. thaliana* and six other plant species (*Glycine max*,
463 *Medicago truncatula*, *Populus trichocarpa*, *Vitis vinifera*, *Sorghum bicolor*, and *Oryza sativa*)

464 (F. Li et al. 2012), which are referred to as conserved blocks. For each conserved block, the first
465 feature was the proportion of a sequence window that overlapped a conserved block (referred to
466 as coverage), and the two other features were the maximum and average phastCons scores within
467 each sequence window. The phastCons score was determined for each nucleotide within
468 conserved blocks (F. Li et al. 2012). Nucleotides in a sequence window that did not overlap with
469 a conserved block were assigned a phastCons score of 0. For each sequence window, five
470 pairwise alignment-based cross-species conservation features were the percent identities to the
471 most significant BLASTN match (if $E\text{-value} < 1E\text{-}05$) in each of five taxonomic groups. The five
472 taxonomic groups included the *Brassicaceae* family ($n_{\text{species}}=7$), other dicotyledonous plants (22),
473 monocotyledonous plants (7), other embryophytes (3), and green algae (5). If no sequence with
474 significant similarity was present, percent identity was scored as zero.

475 For sequence-structure features, we used 125 conformational and thermodynamic
476 dinucleotide properties collected from DiProDB database (Friedel et al. 2009). Because the
477 number of dinucleotide properties was high and dependent, we reduced the dimensionality by
478 utilizing principal component (PC) analysis as described previously (Tsai et al. 2015). Sequence-
479 structure values corresponding to the first five PCs were calculated for all dinucleotides in and
480 averaged across the length of a sequence window and used as features when building function
481 prediction models.

482 **Transcription activity features**

483 We generated four multi-dataset and 20 individual dataset transcription activity features.
484 To identify a set of RNA-seq datasets to calculate multi-dataset features, we focused on the 72 of
485 206 RNA-seq datasets each with ≥ 20 million reads (see above; Supplementary Table 5).
486 Transcribed regions were identified with TopHat2 and Cufflinks as described in the RNA-seq
487 analysis section except that the 72 *A. thaliana* RNA-seq datasets were used. Following transcript
488 assembly, we excluded 21 RNA-seq datasets because they had unusually high RPKM (Reads Per
489 Kilobase of transcript per Million mapped reads) values (median RPKM value
490 $\text{range}=272\sim 2,504,294$) compared to the rest ($2\sim 252$). The remaining 51 RNA-seq datasets were
491 used to generate four multi-dataset transcription activity features including: expression breadth,
492 95th percentile expression level, maximum transcript coverage, and presence of expression
493 evidence (for values see Supplementary Table 2, 6). Expression breadth was the number of
494 RNA-seq datasets that have ≥ 1 transcribed region that overlapped with a sequence window. The

495 95th percentile expression level was the 95th percentile of RPKM values across 51 RNA-seq
496 datasets where RPKM values were set to 0 if there was no transcribed region for a sequence
497 window. Maximum transcript coverage was the maximum proportion of a sequence window that
498 overlapped with a transcribed region across 51 RNA-seq datasets. Presence of expression
499 evidence was determined by overlap between a sequence window and any transcribed region in
500 the 51 RNA-seq datasets.

501 In addition to features based on multiple datasets, 20 individual dataset features were
502 derived from 10 datasets: seven tissue/organ-specific RNA-seq datasets including pollen
503 (SRR847501), seedling (SRR1020621), leaf (SRR953400), root (SRR578947), inflorescence
504 (SRR953399), flower, (SRR505745) and silique (SRR953401), and three datasets from non-
505 standard growth conditions, including dark-grown seedlings (SRR974751) and leaf tissue under
506 drought (SRR921316) and fungal infection (SRR391052). For each of these 10 RNA-seq
507 datasets, we defined two features for each sequence window: the maximum transcript coverage
508 (as described above) and the maximum RPKM value of overlapping transcribed regions (referred
509 to as Level in **Fig. 2**). If no transcribed regions overlapped a sequence window, the maximum
510 RPKM value was set as 0. For the analysis of narrowly- and broadly-expressed phenotype genes
511 and pseudogenes (Supplementary Fig. 3B,C), we used 28 out of 51 RNA-seq datasets generated
512 from a single tissue and in standard growth conditions to calculate the number of tissues with
513 evidence of expression (tissue expression breadth). In total, seven tissues were represented
514 among the 28 selected RNA-seq datasets (see above; Supplementary Table 5), and thus tissue
515 expression breadth ranges from 0 to 7 (note that only 1 through 7 are shown in Supplementary
516 Fig. 3B,C due to low sample size of phenotype genes in the 0 bin). The tissue breadth value is
517 distinct from the expression breadth feature used in model building that was generated using all
518 51 datasets and considered multiple RNA-seq datasets from the same tissue separately (range: 0-
519 51).

520 **Histone 3 mark features**

521 Twenty histone 3 (H3) mark features were calculated based on eight H3 chromatin
522 immunoprecipitation sequencing (ChIP-seq) datasets from SRA. The H3 marks examined
523 include four associated with activation (H3K4me1: SRR2001269, H3K4me3: SRR1964977,
524 H3K9ac: SRR1964985, and H3K23ac: SRR1005405) and four associated with repression
525 (H3K9me1: SRR1005422, H3K9me2: SRR493052, H3K27me3: SRR3087685, and H3T3ph:

526 SRR2001289). Reads were trimmed as described in the RNA-seq section and mapped to the
527 TAIR10 genome with Bowtie v2.2.5 (default parameters) (Langmead et al. 2009). Spatial
528 Clustering for Identification of ChIP-Enriched Regions v.1.1 (Xu et al. 2014) was used to
529 identify ChIP-seq peaks with a false discover rate ≤ 0.05 with a non-overlapping window size of
530 200, a gap parameter of 600, and an effective genome size of 0.92 (Koehler et al. 2011). For each
531 H3 mark, two features were calculated for each sequence window: the maximum intensity
532 among overlapping peaks and peak coverage (proportion of overlap with the peak that overlaps
533 maximally with the sequence window). In addition, four multi-mark features were generated.
534 Two of the multi-mark features were the number of activating marks (0-4) overlapping a
535 sequence window and the proportion of a sequence window overlapping any peak from any of
536 the four activating marks (activating mark peak coverage). The remaining two multi-mark
537 features were the same as the two activating multi-mark features except focused on the four
538 repressive marks.

539 **DNA methylation features**

540 Twenty-one DNA methylation features were calculated from bisulfite-sequencing (BS-
541 seq) datasets from seven tissues (pollen: SRR516176, embryo: SRR1039895, endosperm:
542 SRR1039896, seedling: SRR520367, leaf: SRR1264996, root: SRR1188584, and inflorescence:
543 SRR2155684). BS-seq reads were trimmed as described above and processed with Bismark v.3
544 (default parameters) (Krueger and Andrews 2011) to identify methylated and unmethylated
545 cytosines in CG, CHH, and CHG (H = A, C, or T) contexts. Methylated cytosines were defined
546 as those with ≥ 5 mapped reads and with $>50\%$ of mapped reads indicating that the position was
547 methylated. For each BS-seq dataset, the percentage of methylated cytosines in each sequence
548 window for CG, CHG, and CHH contexts were calculated if the sequence window had ≥ 5
549 cytosines with ≥ 5 reads mapping to the position. To determine whether the above parameters
550 were reasonable, we assessed the false positive rate of DNA methylation calls by evaluating the
551 proportion of cytosines in the chloroplast genome that are called as methylated, as the
552 chloroplast genome has few DNA methylation events (Ngernprasirtsiri et al. 1988; Zhang et al.
553 2006). Based on the above parameters, 0-1.5% of cytosines in CG, CHG, or CHH contexts in the
554 chloroplast genome were considered methylated in any of the seven BS-seq datasets. This
555 indicated that the false positive rates for DNA methylation calls were low and the parameters
556 were reasonable.

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

558 Chromatin accessibility features consisted of ten DHS-related features and one
559 micrococcal nuclease sequencing (MNase-seq)-derived feature. DHS peaks from five tissues
560 (seed coat, seedling, root, unopened flowers, and opened flowers) were retrieved from the Gene
561 Expression Omnibus (GSE53322 and GSE53324) (Sullivan et al. 2014). For each of the five
562 tissues, the maximum DHS peak intensity and DHS peak coverage were calculated for each
563 sequence window. Normalized nucleosome occupancy per bp based on MNase-seq was obtained
564 from Liu et al. (Liu et al. 2015). The average nucleosome occupancy value was calculated across
565 each sequence window. Transcription factor (TF) binding site features were based on *in vitro*
566 DNA affinity purification sequencing data of 529 TFs (O'Malley et al. 2016). Two features were
567 generated for each sequence window: the total number of TF binding sites and the number of
568 distinct TFs bound.

569 **Single-feature prediction performance**

570 The ability for each single feature to distinguish between functional and nonfunctional
571 regions was evaluated by calculating AUC-ROC value with the Python scikit-learn package
572 (Pedregosa et al. 2011). Thresholds to predict sequences as functional or nonfunctional using a
573 single feature were defined by the feature value that produced the highest F-measure, the
574 harmonic mean of precision (proportion of sequences predicted as functional that are truly
575 functional) and recall (proportion of truly functional sequences predicted as functional). The F-
576 measure allows consideration of both false positives and false negatives at a given threshold.
577 FPR were calculated as the percentage of negative (nonfunctional) cases with values above or
578 equal to the threshold and thus falsely predicted as functional. FNR were calculated as the
579 percentage of positive (functional) cases with values below the threshold and thus falsely
580 predicted as nonfunctional.

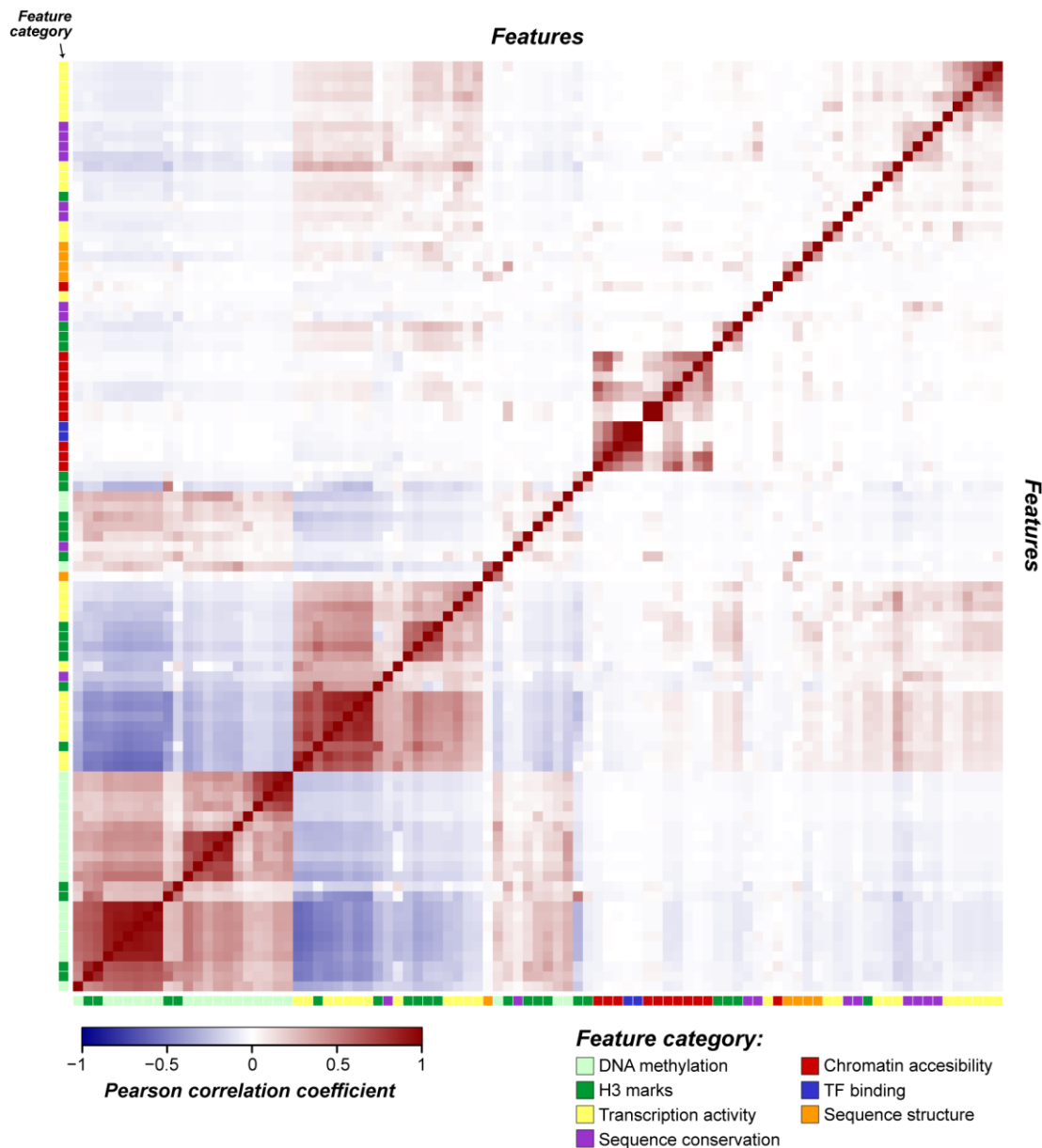
581 **Binary classification with machine learning**

582 For binary classification (two-class) models that contrasted phenotype genes and
583 pseudogenes, the random forest (RF) implementation in the Waikato Environment for
584 Knowledge Analysis software (WEKA) (Hall et al. 2009) was utilized. Three types of two-class
585 models were established, including the full model (500 bp sequence window, **Fig. 3A,B** and **Fig.**
586 **4**), tissue-agnostic models (500 bp, Supplementary Fig. 4; 100 bp, Supplementary Fig. 6), and

587 single feature category models (**Fig. 3A,B**). For each model type, we first generated 100
588 balanced datasets by randomly selecting equal numbers of phenotype genes (positive examples)
589 and pseudogenes (negative examples). For each of these 100 datasets, 10-fold stratified cross-
590 validation was utilized, where the model was trained using 90% of sequences and tested on the
591 remaining 10%. Thus, for each model type, a sequence window had 100 prediction scores, where
592 each score was the proportion of 500 random forest trees that predicted a sequence as a
593 phenotype gene in a balanced dataset. The median of 100 prediction scores was used as the
594 functional likelihood (FL) value (Supplementary Table 4). The FL threshold to predict a
595 sequence as functional or nonfunctional was defined based on maximum F-measure as described
596 in the previous section.

597 We tested multiple -K parameters (2 to 25) in the WEKA-RF implementation, which
598 alters the number of randomly-selected features included in each RF tree (Supplementary Table
599 8), and found that 15 randomly-selected features provided the highest performance based on
600 AUC-ROC (calculated and visualized using the ROCR package) (Sing et al. 2005). Feature
601 importance was assessed by excluding one feature at a time to determine the associated reduction
602 in prediction performance (Supplementary Table 9). All leave-one-out models performed well
603 (AUC-ROC >0.97), indicating that no single feature was dominating the function predictions
604 and/or many features are correlated (Supplementary Fig. 7). Binary classification models were
605 also built using all features from 500 bp sequences (equivalent to the full model) with the
606 Sequential Minimal Optimization - Support Vector Machine (SMO-SVM) implementation in
607 WEKA (Hall et al. 2009). The results of SMO-SVM models were highly similar to the full RF
608 results: *PCC* between the FL values generated by RF and SMO-SVM=0.97; AUC-ROC of
609 SMO-SVM=0.97; FPR=12%; FNR=3%. By comparison, the full RF model had AUC-
610 ROC=0.98, FPR=10%, FNR=4%.

611 Tissue-agnostic models were generated by excluding the expression breadth feature and
612 95th percentile expression level and replacing all features from RNA-seq, BS-seq, and DHS
613 datasets that were available in multiple tissues. For multiple-tissue RNA-seq data, the maximum
614 expression level across 51 RNA-seq datasets (in RPKM) and maximum coverage (as described
615 in the transcription activity section) of a sequence window in any of 51 RNA-seq datasets were
616 used. For multi-tissue DNA methylation features, minimum proportions of methylated cytosines
617 in any tissue in CG, CHG, and CHH contexts were used. For DHS data, the maximum peak



Supplementary Figure 7. Correlation between features used in functional predictions. Colors within the heatmap indicate pairwise correlation between two features. Colors on the left-most and bottom-most edges indicate the associated feature category (Fig. 2). Feature values were quantile normalized prior to calculating correlation.

618

619 intensity and peak coverage was used instead. In single feature category predictions, fewer total
620 features were used and therefore lower $-K$ values (i.e. the number of random features selected
621 when building random forests) were considered in parameter searches (Supplementary Table 8).

622 **Multi-class machine learning model**

623 For the four-class model, benchmark RNA gene, phenotype protein-coding gene,
624 pseudogene, and random unexpressed intergenic sequences were used as the four training
625 classes. Benchmark RNA genes consisted of six RNA genes with documented loss-of-function
626 phenotypes and 40 high-confidence miRNA genes from miRBase (www.mirbase.org)
627 (Kozomara and Griffiths-Jones 2014). We considered that the decreased numbers of benchmark
628 RNA genes would not allow us to effectively distinguish between sequence classes. However,
629 binary predictions generated using 35 phenotype gene and pseudogene instances and the 100 bp
630 tissue-agnostic feature set resulted in an AUC-ROC performance of 0.96. We generated 250
631 datasets with equal proportions (larger classes randomly sampled) of training sequences. Two-
632 fold stratified cross-validation was utilized due to the low number of benchmark RNA genes.
633 The features included those described for the tissue-agnostic model and focused on 100 bp
634 sequence windows. The RF implementation, *cforest*, in the *party* package of R (Strobl et al.
635 2008) was used to build the classifiers. The four-class predictions provide prediction scores for
636 each sequence type: an RNA gene, phenotype protein-coding gene, pseudogene, and unexpressed
637 intergenic score (Supplementary Table 4). The prediction scores indicate the proportion of
638 random forest trees that classify a sequence as a particular class. Median prediction scores from
639 across 100 balanced runs were used as final prediction scores. Scores from a single balanced
640 dataset models sum to 1, but not the median from 100 balanced runs. Thus, the median scores
641 were scaled to sum to 1. For each sequence window, the maximum prediction score among the
642 four classes was used to classify a sequence as phenotype gene, pseudogene, unexpressed
643 intergenic region, or RNA gene.

644 **AVAILABILITY**

645 All relevant data are within the article and supplementary data files.

646 **SUPPLEMENTARY MATERIALS**

647 Supplementary Data are available online.

648 **ACKNOWLEDGEMENTS**

649 The authors wish to thank Christina Azodi, Ming-Jung Liu, Gaurav Moghe, Bethany Moore, and
650 Sahra Uygun and for providing processed data and discussion. This work was supported by the
651 National Science Foundation (grant numbers IOS-1126998, IOS-1546617, and DEB-1655386 to

652 S.H.S), and Research Experience for Undergraduates support [to R.P.S]; and the Michigan State
653 University Dissertation Continuation Fellowship [to J.P.L].

654 **CONFLICT OF INTEREST**

655 The authors have no conflicts of interest to disclose.

656 **REFERENCES**

- 657 Ajjawi I, Lu Y, Savage LJ, Bell SM, Last RL. 2010. Large-scale reverse genetics in *Arabidopsis*:
658 case studies from the Chloroplast 2010 Project. *Plant Physiol.* 152:529–540.
- 659 Amundson R, Lauder G V. 1994. Function without purpose. *Biol. Philos.* 9:443–469.
- 660 Berardini TZ, Reiser L, Li D, Mezheritsky Y, Muller R, Strait E, Huala E. 2015. The
661 *Arabidopsis* information resource: Making and mining the “gold standard” annotated
662 reference plant genome. *Genesis* 53:474–485.
- 663 Bernard D, Prasanth K V, Tripathi V, Colasse S, Nakamura T, Xuan Z, Zhang MQ, Sedel F,
664 Jourden L, Couplier F, et al. 2010. A long nuclear-retained non-coding RNA regulates
665 synaptogenesis by modulating gene expression. *EMBO J.* 29:3082–3093.
- 666 Boeck ME, Huynh C, Gevirtzman L, Thompson OA, Wang G, Kasper DM, Reinke V, Hillier
667 LW, Waterston RH. 2016. The time-resolved transcriptome of *C. elegans*. *Genome Res.*
668 26:1441–1450.
- 669 Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence
670 data. *Bioinformatics* 30:2114–2120.
- 671 Brown JB, Boley N, Eisman R, May GE, Stoiber MH, Duff MO, Booth BW, Wen J, Park S,
672 Suzuki AM, et al. 2014. Diversity and dynamics of the *Drosophila* transcriptome. *Nature*
673 512:393–399.
- 674 Cao J, Schneeberger K, Ossowski S, Günther T, Bender S, Fitz J, Koenig D, Lanz C, Stegle O,
675 Lippert C, et al. 2011. Whole-genome sequencing of multiple *Arabidopsis thaliana*
676 populations. *Nat. Genet.* 43:956–963.
- 677 Doolittle WF, Brunet TDP, Linquist S, Gregory TR. 2014. Distinguishing between “function”
678 and “effect” in genome biology. *Genome Biol. Evol.* 6:1234–1237.
- 679 Eddy SR. 2013. The ENCODE project: missteps overshadowing a success. *Curr. Biol.* 23:R259--
680 61.
- 681 ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the human

- 682 genome. *Nature* 489:57–74.
- 683 Fei Q, Xia R, Meyers BC. 2013. Phased, Secondary, Small Interfering RNAs in
684 Posttranscriptional Regulatory Networks. *Plant Cell* 25:2400–2415.
- 685 Friedel M, Nikolajewa S, Sühnel J, Wilhelm T. 2009. DiProDB: a database for dinucleotide
686 properties. *Nucleic Acids Res.* 37:D37--40.
- 687 Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W, Hellsten
688 U, Putnam N, et al. 2012. Phytozome: a comparative platform for green plant genomics.
689 *Nucleic Acids Res.* 40:D1178-86.
- 690 Graur D, Zheng Y, Price N, Azevedo RBR, Zufall R a, Elhaik E. 2013. On the immortality of
691 television sets: “function” in the human genome according to the evolution-free gospel of
692 ENCODE. *Genome Biol. Evol.* 5:578–590.
- 693 Guil S, Esteller M. 2012. Cis-acting noncoding RNAs: friends and foes. *Nat. Struct. Mol. Biol.*
694 19:1068–1075.
- 695 Gulko B, Gronau I, Hubisz MJ, Siepel A. 2014. Probabilities of Fitness Consequences for Point
696 Mutations Across the Human Genome.
- 697 Guo H-S, Xie Q, Fei J-F, Chua N-H. 2005. MicroRNA directs mRNA cleavage of the
698 transcription factor NAC1 to downregulate auxin signals for arabidopsis lateral root
699 development. *Plant Cell* 17:1376–1386.
- 700 Hall M, Frank E, Holmes G, Pfahringer B, Reutemann P, Witten IH. 2009. The WEKA data
701 mining software. *ACM SIGKDD Explor. Newsl.* 11:10.
- 702 Hardiman KE, Brewster R, Khan SM, Deo M, Bodmer R. 2002. The bereft gene, a potential
703 target of the neural selector gene cut, contributes to bristle morphogenesis. *Genetics*
704 161:231–247.
- 705 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.
706 Uncovering small RNA-mediated responses to phosphate deficiency in Arabidopsis by deep
707 sequencing. *Plant Physiol.* 151:2120–2132.
- 708 Karreth FA, Reschke M, Ruocco A, Ng C, Chapuy B, Léopold V, Sjöberg M, Keane TM, Verma
709 A, Ala U, et al. 2015. The BRAF pseudogene functions as a competitive endogenous RNA
710 and induces lymphoma in vivo. *Cell* 161:319–332.
- 711 Kellis M, Wold B, Snyder MP, Bernstein BE, Kundaje A, Marinov GK, Ward LD, Birney E,
712 Crawford GE, Dekker J, et al. 2014. Defining functional DNA elements in the human

- 713 genome. *Proc. Natl. Acad. Sci. USA* 111:6131–6138.
- 714 Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. 2013. TopHat2: accurate
715 alignment of transcriptomes in the presence of insertions, deletions and gene fusions.
716 *Genome Biol.* 14:R36.
- 717 Koehler R, Issac H, Cloonan N, Grimmond SM. 2011. The uniqueome: a mappability resource
718 for short-tag sequencing. *Bioinformatics* 27:272–274.
- 719 Kozomara A, Griffiths-Jones S. 2014. miRBase: annotating high confidence microRNAs using
720 deep sequencing data. *Nucleic Acids Res.* 42:D68--73.
- 721 Krishnakumar V, Hanlon MR, Contrino S, Ferlanti ES, Karamycheva S, Kim M, Rosen BD,
722 Cheng C-Y, Moreira W, Mock SA, et al. 2015. Araport: the Arabidopsis information portal.
723 *Nucleic Acids Res.* 43:D1003--9.
- 724 Krueger F, Andrews SR. 2011. Bismark: a flexible aligner and methylation caller for Bisulfite-
725 Seq applications. *Bioinformatics* 27:1571–1572.
- 726 Kuromori T, Wada T, Kamiya A, Yuguchi M, Yokouchi T, Imura Y, Takabe H, Sakurai T,
727 Akiyama K, Hirayama T, et al. 2006. A trial of phenome analysis using 4000 *Ds*-insertional
728 mutants in gene-coding regions of Arabidopsis. *Plant J.* 47:640–651.
- 729 Lai K-MV, Gong G, Atanasio A, Rojas J, Quispe J, Posca J, White D, Huang M, Fedorova D,
730 Grant C, et al. 2015. Diverse Phenotypes and Specific Transcription Patterns in Twenty
731 Mouse Lines with Ablated LincRNAs. *PLoS One* 10:e0125522.
- 732 Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, Muller R, Dreher K,
733 Alexander DL, Garcia-Hernandez M, et al. 2012. The Arabidopsis Information Resource
734 (TAIR): improved gene annotation and new tools. *Nucleic Acids Res.* 40:D1202-10.
- 735 Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment
736 of short DNA sequences to the human genome. *Genome Biol.* 10:R25.
- 737 Li F, Zheng Q, Vandivier LE, Willmann MR, Chen Y, Gregory BD. 2012. Regulatory impact of
738 RNA secondary structure across the Arabidopsis transcriptome. *Plant Cell* 24:4346–4359.
- 739 Li W, Cui X, Meng Z, Huang X, Xie Q, Wu H, Jin H, Zhang D, Liang W. 2012. Transcriptional
740 regulation of Arabidopsis MIR168a and argonaute1 homeostasis in abscisic acid and abiotic
741 stress responses. *Plant Physiol.* 158:1279–1292.
- 742 Li W, Gojobori T, Nei M. 1981. Pseudogenes as a paradigm of neutral evolution. *Nature*.
- 743 Liu M-J, Seddon AE, Tsai ZT-Y, Major IT, Floer M, Howe GA, Shiu S-H. 2015. Determinants

- 744 of nucleosome positioning and their influence on plant gene expression. *Genome Res.*
745 25:1182–1195.
- 746 Lloyd J, Meinke D. 2012. A comprehensive dataset of genes with a loss-of-function mutant
747 phenotype in *Arabidopsis*. *Plant Physiol.* 158:1115–1129.
- 748 Lloyd JP, Seddon AE, Moghe GD, Simenc MC, Shiu S-H. 2015. Characteristics of Plant
749 Essential Genes Allow for within- and between-Species Prediction of Lethal Mutant
750 Phenotypes. *Plant Cell* 27:2133–2147.
- 751 Marahrens Y, Panning B, Dausman J, Strauss W, Jaenisch R. 1997. Xist-deficient mice are
752 defective in dosage compensation but not spermatogenesis. *Genes Dev.* 11:156–166.
- 753 Michael TP, Jackson S. 2013. The First 50 Plant Genomes. *Plant Genome* 6:0.
- 754 Moghe GD, Lehti-Shiu MD, Seddon AE, Yin S, Chen Y, Juntawong P, Brandizzi F, Bailey-
755 Serres J, Shiu S-H. 2013. Characteristics and significance of intergenic polyadenylated
756 RNA transcription in *Arabidopsis*. *Plant Physiol.* 161:210–224.
- 757 Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, Gerstein M, Snyder M. 2008. The
758 Transcriptional Landscape of the Yeast Genome Defined by RNA Sequencing. *Science*
759 320:1344–1349.
- 760 Neander K. 1991. Functions as selected effects: The conceptual analyst's defense. *Philos. Sci.*
761 58:168–184.
- 762 Ngernprasirtsiri J, Kobayashi H, Akazawa T. 1988. DNA methylation as a mechanism of
763 transcriptional regulation in nonphotosynthetic plastids in plant cells. *Proc. Natl. Acad. Sci.*
764 U. S. A. 85:4750–4754.
- 765 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
766 dissecting phenotypic variants in human lincRNAs. *Eur. J. Hum. Genet.* 21:1128–1133.
- 767 Niu D-K, Jiang L. 2013. Can ENCODE tell us how much junk DNA we carry in our genome?
768 *Biochem. Biophys. Res. Commun.* 430:1340–1343.
- 769 Nobuta K, Venu RC, Lu C, Beló A, Vemaraju K, Kulkarni K, Wang W, Pillay M, Green PJ,
770 Wang G-L, et al. 2007. An expression atlas of rice mRNAs and small RNAs. *Nat.*
771 *Biotechnol.* 25:473–477.
- 772 O'Malley RC, Huang S-SC, Song L, Lewsey MG, Bartlett A, Nery JR, Galli M, Gallavotti A,
773 Ecker JR. 2016. Cistrome and Epicistrome Features Shape the Regulatory DNA Landscape.
774 *Cell* 166:1598.

- 775 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer
776 P, Weiss R, Dubourg V, et al. 2011. Scikit-learn: Machine Learning in Python. *J. Mach.*
777 *Learn. Res.* 12:2825–2830.
- 778 Penny GD, Kay GF, Sheardown SA, Rastan S, Brockdorff N. 1996. Requirement for Xist in X
779 chromosome inactivation. *Nature* 379:131–137.
- 780 Poliseno L, Salmena L, Zhang J, Carver B, Haveman WJ, Pandolfi PP. 2010. A coding-
781 independent function of gene and pseudogene mRNAs regulates tumour biology. *Nature*
782 465:1033–1038.
- 783 Ponting CP, Belgard TG. 2010. Transcribed dark matter: meaning or myth? *Hum. Mol. Genet.*
784 19:R162-8.
- 785 Sauvageau M, Goff LA, Lodato S, Bonev B, Groff AF, Gerhardinger C, Sanchez-Gomez DB,
786 Hacısuleyman E, Li E, Spence M, et al. 2013. Multiple knockout mouse models reveal
787 lincRNAs are required for life and brain development. *Elife* 2:e01749.
- 788 Savage LJ, Imre KM, Hall DA, Last RL. 2013. Analysis of essential Arabidopsis nuclear genes
789 encoding plastid-targeted proteins. *PLoS One* 8:e73291.
- 790 Schreiber SL, Bernstein BE. 2002. Signaling Network Model of Chromatin. *Cell* 111:771–778.
- 791 Shin H, Shin H-S, Chen R, Harrison MJ. 2006. Loss of At4 function impacts phosphate
792 distribution between the roots and the shoots during phosphate starvation. *Plant J.* 45:712–
793 726.
- 794 Simon SA, Meyers BC. 2011. Small RNA-mediated epigenetic modifications in plants. *Curr.*
795 *Opin. Plant Biol.* 14:148–155.
- 796 Sing T, Sander O, Beerenwinkel N, Lengauer T. 2005. ROCr: visualizing classifier performance
797 in R. *Bioinformatics* 21:3940–3941.
- 798 Stolc V, Samanta MP, Tongprasit W, Sethi H, Liang S, Nelson DC, Hegeman A, Nelson C,
799 Rancour D, Bednarek S, et al. 2005. Identification of transcribed sequences in Arabidopsis
800 thaliana by using high-resolution genome tiling arrays. *Proc. Natl. Acad. Sci. U. S. A.*
801 102:4453–4458.
- 802 Strobl C, Boulesteix A-L, Kneib T, Augustin T, Zeileis A. 2008. Conditional variable
803 importance for random forests. *BMC Bioinformatics* 9:307.
- 804 Struhl K. 2007. Transcriptional noise and the fidelity of initiation by RNA polymerase II. *Nat.*
805 *Struct. Mol. Biol.* 14:103–105.

- 806 Sullivan AM, Arsovski AA, Lempe J, Bubb KL, Weirauch MT, Sabo PJ, Sandstrom R, Thurman
807 RE, Neph S, Reynolds AP, et al. 2014. Mapping and dynamics of regulatory DNA and
808 transcription factor networks in *A. thaliana*. *Cell Rep.* 8:2015–2030.
- 809 Svensson O, Arvestad L, Lagergren J. 2006. Genome-wide survey for biologically functional
810 pseudogenes. *PLoS Comput. Biol.* 2:e46.
- 811 Swarbreck D, Wilks C, Lamesch P, Berardini TZ, Garcia-hernandez M, Foerster H, Li D, Meyer
812 T, Muller R, Ploetz L, et al. 2008. The Arabidopsis Information Resource (TAIR): gene
813 structure and function annotation. *Nucleic Acids Res.* 36:1009–1014.
- 814 Tan JY, Sirey T, Honti F, Graham B, Piovesan A, Merkschlager M, Webber C, Ponting CP,
815 Marques AC. 2015. Extensive microRNA-mediated crosstalk between lncRNAs and
816 mRNAs in mouse embryonic stem cells. *Genome Res.* 25:655–666.
- 817 Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold
818 BJ, Pachter L. 2010. Transcript assembly and quantification by RNA-Seq reveals
819 unannotated transcripts and isoform switching during cell differentiation. *Nat. Biotechnol.*
820 28:511–515.
- 821 Tsai ZT-Y, Lloyd JP, Shiu S-H. 2017. Defining Functional Genic Regions in the Human
822 Genome through Integration of Biochemical, Evolutionary, and Genetic Evidence. *Mol.*
823 *Biol. Evol.*
- 824 Tsai ZT-Y, Shiu S-H, Tsai H-K. 2015. Contribution of Sequence Motif, Chromatin State, and
825 DNA Structure Features to Predictive Models of Transcription Factor Binding in Yeast.
826 *PLoS Comput. Biol.* 11:e1004418.
- 827 VanBuren R, Bryant D, Edger PP, Tang H, Burgess D, Challabathula D, Spittle K, Hall R, Gu J,
828 Lyons E, et al. 2015. Single-molecule sequencing of the desiccation-tolerant grass
829 *Oropetium thomaeum*. *Nature* 527:508–511.
- 830 Wang Z, Zang C, Rosenfeld JA, Schones DE, Barski A, Cuddapah S, Cui K, Roh T-Y, Peng W,
831 Zhang MQ, et al. 2008. Combinatorial patterns of histone acetylations and methylations in
832 the human genome. *Nat. Genet.* 40:897–903.
- 833 Xu S, Grullon S, Ge K, Peng W. 2014. Spatial clustering for identification of ChIP-enriched
834 regions (SICER) to map regions of histone methylation patterns in embryonic stem cells.
835 *Methods Mol. Biol.* 1150:97–111.
- 836 Yamada K, Lim J, Dale JM, Chen H, Shinn P, Palm CJ, Southwick AM, Wu HC, Kim C,

837 Nguyen M, et al. 2003. Empirical analysis of transcriptional activity in the Arabidopsis
838 genome. *Science* 302:842–846.

839 Yang L, Takuno S, Waters ER, Gaut BS. 2011. Lowly expressed genes in *Arabidopsis thaliana*
840 bear the signature of possible pseudogenization by promoter degradation. *Mol. Biol. Evol.*
841 28:1193–1203.

842 Zhang X, Yazaki J, Sundaresan A, Cokus S, Chan SW-L, Chen H, Henderson IR, Shinn P,
843 Pellegrini M, Jacobsen SE, et al. 2006. Genome-wide high-resolution mapping and
844 functional analysis of DNA methylation in *Arabidopsis*. *Cell* 126:1189–1201.

845 Zhao Y, Li H, Fang S, Kang Y, Wu W, Hao Y, Li Z, Bu D, Sun N, Zhang MQ, et al. 2016.
846 NONCODE 2016: an informative and valuable data source of long non-coding RNAs.
847 *Nucleic Acids Res.* 44:D203--8.

848 Zou C, Lehti-Shiu MD, Thibaud-Nissen F, Prakash T, Buell CR, Shiu S-H. 2009. Evolutionary
849 and expression signatures of pseudogenes in *Arabidopsis* and rice. *Plant Physiol.* 151:3–15.
850

851 **FIGURE LEGENDS**

852 **Figure 1.** Relationship between genome size and number of nucleotides covered by RNA-seq
853 reads (expression) in 15 flowering plant species. **(A)** annotated genic regions. **(B)** intergenic
854 regions. Transcribed regions were considered as intergenic if they did not overlap with any gene
855 annotation and had no significant translated sequence similarity to plant protein sequences.
856 Identical numbers of RNA-sequencing reads (30 million) and the same mapping procedures were
857 used in all species. Mb: megabase. Gb: gigabase. Dotted lines: linear model fits. r^2 : square of
858 Pearson's correlation coefficient.

859 **Figure 2.** Predictions of functional (phenotype gene) and non-functional (pseudogene) sequences
860 based on each individual feature. Prediction performance is measured using Area Under the
861 Curve - Receiver Operating Characteristic (AUC-ROC). AUC-ROC values range between 0.5
862 (random) and 1 (perfect separation), with AUC-ROC values of 0.7, 0.8, and 0.9 considered fair,
863 good, and excellent performance, respectively. Features include those in the categories of **(A)**
864 transcription activity, **(B)** sequence conservation, **(C)** DNA methylation, **(D)** transcription factor
865 (TF) binding, **(E)** histone 3 (H3) marks, **(F)** sequence structure, and **(G)** chromatin accessibility.

866 AUC-ROC ranges in value from 0.5 (equivalent to random guessing) to 1 (perfect predictions).
867 Dotted lines: median AUC-ROC of features in a category.

868 **Figure 3.** Predictions of functional and nonfunctional sequences based on multiple features. **(A)**
869 AUC-ROC values of function prediction models built when considering all features (Full), all
870 except transcription activity (TX)-related features (Full (-TX)), and all features from each
871 category. The category abbreviations follow those in **Fig. 2.** **(B)** Precision-recall curves of the
872 models with matching colors from **(A)**. The models were built using feature values calculated
873 from 500 bp sequence windows. We also conducted principle component (PC) analysis to
874 investigate how well phenotype genes and pseudogenes could be separated and found that
875 phenotype genes (Supplementary Fig. 1A) and pseudogenes (Supplementary Fig. 1B) were
876 distributed in largely distinct space. However, there remained substantial overlap, indicating that
877 standard parametric approaches are not well suited to distinguishing between benchmark
878 functional and nonfunctional sequences.

879 **Figure 4.** Functional likelihood distributions of various sequence classes based on the full
880 model. **(A)** Phenotype genes. **(B)** Pseudogenes. **(C)** Annotated protein-coding genes. **(D)**
881 Transposable elements. **(E)** Random unexpressed intergenic sequences. **(F)** Intergenic
882 transcribed regions (ITR). **(G)** Araport11 ncRNAs. **(H)** TAIR10 ncRNAs. The full model was
883 established using 500 bp sequence windows. Higher and lower functional likelihood values
884 indicate greater similarity to phenotype genes and pseudogenes, respectively. Vertical dashed
885 lines indicate the threshold for calling a sequence as functional or nonfunctional. The
886 percentages to the left and right of the dashed line indicate the percent of sequences predicted as
887 functional or nonfunctional, respectively.

888 **Figure 5.** Proportion of phenotype genes, pseudogenes, ITRs, and ncRNAs predicted as
889 functional in the full and single-category models. Percentages of sequence classes that are
890 predicted as functional in models based on all features and the single category models, each
891 using all features from a category (abbreviated according to **Fig. 2**). The models are sorted from
892 left to right based on performance (AUC-ROC). The colors of and numbers within the blocks
893 indicate the proportion sequences predicted as functional by a given model. Phenotype gene and
894 pseudogene sequences are shown in three sub-groups: all sequences (All), and those predicted as

895 functional (high functional likelihood (FL)) and nonfunctional (low FL) in the full model. ITR:
896 intergenic transcribed regions. A greater proportion of ITRs and Araport ncRNAs are predicted
897 as functional when considering only DNA methylation or H3 mark features compared to the full
898 (**Fig. 3**) or tissue-agnostic (Supplementary Fig. 4) models. However, these two category-specific
899 models also had higher false positive rates (unexpressed intergenic sequences and pseudogenes).
900 Thus, these single feature-category models do not provide additional support for the functionality
901 of most Araport ncRNAs and ITRs.

902 **Figure 6.** Function predictions based on a four-class prediction model. (**A**) Stacked bar plots
903 indicate the prediction scores of benchmark RNA genes for each of the four classes: dark blue -
904 phenotype protein-coding gene (Ph), cyan - RNA gene (RNA), red - pseudogene (Ps), yellow -
905 random intergenic sequence (Ig). Ig were included to provide another set of likely nonfunctional
906 sequences distinct from pseudogenes. Expression breadth and tissue-specific features were
907 excluded and 100 bp sequences were used. A benchmark RNA gene is classified as one of the
908 four classes according to the highest prediction score. The color bars below the chart indicate the
909 predicted class, with the same color scheme as the prediction score. Sequences classified as Ph or
910 RNA were considered functional, while those classified as Ps or Ig were considered
911 nonfunctional. Percentages below a classification region indicate the proportion of sequences
912 classified as that class. (**B**) Phenotype protein-coding gene prediction scores. (**C**) Pseudogene
913 prediction scores. (**D**) Random unexpressed intergenic region prediction scores. Note that no
914 sequence was predicted as functional. (**E**) Intergenic transcribed region (ITR), (**F**) Araport11
915 ncRNA regions. (**G**) TAIR10 ncRNA regions. Note that the 100 bp model used here allowed us
916 to evaluate an additional 10,938 ITRs and 1,406 annotated ncRNAs compared to the 500 bp full
917 and tissue-agnostic models.

918 SUPPLEMENTARY FIGURE LEGENDS

919 **Supplementary Figure 1.** Smoothed scatterplots of the first two principle components (PCs) of
920 (**A**) phenotype gene and (**B**) pseudogene features. The percentages on the axes in (**A**) indicate the
921 feature value variation explained by the associated PC.

922 **Supplementary Figure 2.** Distributions of 12 example features from (**A**) transcription activity,
923 (**B**) sequence conservation, (**C**) histone 3 (H3) marks, (**D**) DNA methylation, (**E**) transcription
924 factor (TF) binding, (**F**) sequence structure, and (**G**) chromatin accessibility feature categories

925 (Fig. 2). Feature distributions are shown for phenotype genes, pseudogenes, TAIR- and Araport-
926 annotated ncRNAs and intergenic transcribed regions (ITR).

927 **Supplementary Figure 3.** Impacts of conditional phenotypes and expression breadth on the
928 function prediction model. (A) Functional likelihood distributions of phenotype genes with
929 mutant phenotypes under standard growth conditions (non-conditional) and non-standard growth
930 conditions such as stressful environments (conditional) based on the 500 bp full model. Feature
931 values were calculated from a random 500 bp region from within the sequence body. Higher and
932 lower functional likelihood values indicate a greater similarity to phenotype genes and
933 pseudogenes, respectively. (B,C) Distributions of functional likelihood scores for phenotype
934 genes (blue) and pseudogenes (red) for sequences with various breadths of expression for (B) the
935 500 bp full model and (C) the 500 bp tissue-agnostic model generated by excluding the
936 expression breadth and features available from multiple tissues. The tissue-agnostic model is
937 aimed toward minimizing the effects of biochemical activity occurring across multiple tissues
938 and predicts a greater proportion of narrowly-expressed phenotype genes as functional compared
939 to the full model.

940 **Supplementary Figure 4.** Distributions of functional likelihood scores based on the 500 bp
941 tissue-agnostic model. (A) Phenotype genes. (B) Pseudogenes. (C) Annotated protein-coding
942 genes. (D) Transposable elements. (E) Random unexpressed intergenic sequences. (F) Intergenic
943 transcribed regions (ITR). (G) Araport11 ncRNAs. (H) TAIR10 ncRNAs. Vertical dashed lines
944 display the threshold to define a sequence as functional or nonfunctional. The numbers to the left
945 and right of the dashed line show the percentage of sequences predicted as functional or
946 nonfunctional, respectively.

947 **Supplementary Figure 5.** Distance of ITRs and annotated ncRNA regions to and feature
948 similarity with neighboring genes. (A) Distance from intergenic transcribed regions (ITRs) and
949 annotated ncRNAs to the closest neighboring gene. ITR and ncRNA sequences are separated by
950 whether they are predicted as functional (F) or nonfunctional (NF) by the 500 bp full model. (B)
951 Feature similarity between proximal neighbors (within 95th percentile (456 bp) of intron
952 lengths), and (C) distal neighbors (>456 bp). Pairs involving ITRs and annotated ncRNAs were
953 divided by whether the ITR or ncRNA sequence was predicted as functional (F) or nonfunctional

954 (NF) by the full model. Feature values were quantile normalized prior to calculating correlations.
955 (D) Feature similarity based on Pearson's Correlation Coefficients (PCC) between random pairs
956 of ITRs, Araport11 ncRNAs, TAIR10 ncRNAs, or annotated genes.

957 **Supplementary Figure 6.** Distributions of functional likelihood scores based on the 100 bp
958 tissue-agnostic model. (A) Phenotype genes. (B) Pseudogenes. (C) Protein-coding gene. (D)
959 transposable elements. (E) Random unexpressed intergenic sequences. (F) Intergenic transcribed
960 regions (ITR). (G) Araport11 ncRNAs. (H) TAIR10 ncRNAs. (I) RNA genes with loss-of-
961 function mutant phenotypes. (J) MicroRNAs, (K) Small nucleolar RNAs, (L) Small nuclear
962 RNAs. The tissue-agnostic model was built with 100 bp features and while excluding the
963 expression breadth and tissue-specific features as annotated RNA genes tend to be more
964 narrowly expressed than phenotype genes (U tests, all $p < 2e-05$; Supplementary Fig. 2A).
965 Higher functional likelihood values indicate greater similarity to phenotype genes while lower
966 values indicate similarity to pseudogenes. Vertical dashed lines display the threshold to define a
967 sequence as functional or nonfunctional. The numbers to the left and right of the dashed line
968 show the percentage of sequences predicted as functional or nonfunctional, respectively.

969 **Supplementary Figure 7.** Correlation between features used in functional predictions. Colors
970 within the heatmap indicate pairwise correlation between two features. Colors on the left-most
971 and bottom-most edges indicate the associated feature category (Fig. 2). Feature values were
972 quantile normalized prior to calculating correlation.

973 **SUPPLEMENTARY TABLES**

974 **Supplementary Table 1.** Leaf tissue RNA-sequencing datasets for 15 flowering plant species

975 **Supplementary Table 2.** Conservation, biochemical, and sequence-structure feature values
976 calculated from 500 bp sequences.

977 **Supplementary Table 3.** False positive and false negative rates for single feature classifications.

978 **Supplementary Table 4.** Function predictions for all models generated in this study.

979 **Supplementary Table 5.** RNA-sequencing datasets for identifying intergenic transcribed
980 regions, calculating transcription activity features, and assessing tissue-specific predictions.

981 **Supplementary Table 6.** Conservation, biochemical, and sequence-structure feature values
982 calculated from 100 bp sequences.

983 **Supplementary Table 7.** RNA genes with documented loss-of-function phenotypes.

984 **Supplementary Table 8.** K parameters tested for random forest runs.

985 **Supplementary Table 9.** AUC-ROC scores from leave-one-feature-out machine learning runs.

986 **SUPPLEMENTARY INFORMATION**

987 **Variability in single feature performance within feature categories**

988 Within each feature category, there was a wide range of performance between features (**Fig. 2,**
989 Supplementary Table 3) and there were clear biological or technical explanations for features
990 that perform poorly. For the transcription activity category, 17 out of 24 features had an AUC-
991 ROC performance >0.8 , including the best-performing feature, expression breadth (AUC-
992 ROC=0.95; **Fig. 2A**). However, five transcription activity-related features performed poorly
993 (AUC-ROC <0.65), including the presence of expression (transcript) evidence (AUC-ROC=0.58;
994 **Fig. 2A**). For the sequence conservation category, maximum and average phastCons
995 conservation scores were highly distinct between phenotype genes and pseudogenes (AUC-
996 ROC=0.83 and 0.82, respectively; **Fig. 2B**). On the other hand, identity to best matching
997 nucleotide sequences found in *Brassicaceae* and algal species were not informative (AUC-
998 ROC=0.55 and 0.51, respectively; **Fig. 2B**). This was because 99.8% and 95% of phenotype
999 genes and pseudogenes, respectively, had a potentially homologous sequence within the
1000 *Brassicaceae* family, and only 3% and 1%, respectively, in algal species. Thus, *Brassicaceae*
1001 genomes were too similar and algal genomes too dissimilar to *A. thaliana* to provide meaningful
1002 information. H3 mark features also displayed high variability. The most informative H3 mark
1003 features were based on the number and coverage of activation-related marks (AUC-ROC=0.87
1004 and 0.85, respectively; **Fig. 2E**), consistent with the notion that histone marks are often jointly
1005 associated with active genomic sequences to provide a robust regulatory signal (Schreiber and
1006 Bernstein 2002; Wang et al. 2008). By comparison, the coverage and intensity of H3 lysine 27
1007 trimethylation (H3K27me3) and H3 threonine 3 phosphorylation (H3T3ph) were largely
1008 indistinct between phenotype genes and pseudogenes (AUC-ROC range: 0.55-0.59; **Fig. 2E**).

1009 **Error rates from single-feature functional predictions**

1010 The differences between genes and pseudogenes in transcription, conservation, and epigenetic
1011 features and functional genomic regions suggested that these features may individually provide
1012 sufficient information for distinguishing between functional and nonfunctional genomic regions.
1013 To assess this possibility, we next evaluated the error rates of function predictions based on
1014 single features. We first considered expression breadth of a sequence, the best predicting single
1015 feature of functionality. Despite high AUC-ROC (0.95; **Fig. 2A**), the false positive rate (FPR; %
1016 of pseudogenes predicted as phenotype genes) was 21% when only expression breadth was used,
1017 while the false negative rate (FNR; % of phenotype genes predicted as pseudogenes) was 4%.
1018 Similarly, the best-performing H3 mark- and sequence conservation-related features (**Fig. 2B,E**)
1019 had FPRs of 26% and 32%, respectively, and also incorrectly classified at least 10% of
1020 phenotype genes as pseudogenes. Thus, error rates are high even when considering well-
1021 performing single features, indicating the need to jointly consider multiple features for
1022 distinguishing phenotype genes and pseudogenes.

1023 **Features of misclassified sequences**

1024 Although the full model performs exceedingly well, there remain false predictions. There are 76
1025 phenotype genes (4%) predicted as nonfunctional (referred to as low-functional likelihood (FL)
1026 phenotype genes). We assessed why these phenotype genes were not correctly identified by first
1027 asking what category of features were particularly distinct between low-FL and the remaining
1028 phenotype genes. We found that the major category that led to the misclassification of phenotype
1029 genes was transcription activity, as only 7% of low-scoring phenotype genes were predicted as
1030 functional in the transcription activity-only model, compared to 98% of high FL phenotype genes
1031 (**Fig. 5**). By contrast, >65% of low-FL phenotype genes were predicted as functional when
1032 sequence conservation, H3 mark, or DNA methylation features were used. This could suggest
1033 that the full model is less effective in predicting functional sequences that are weakly or
1034 narrowly expressed. While sequence conservation features are distinct between functional and
1035 nonfunctional sequences when considered in combination, a significantly higher proportion of
1036 low-FL phenotype genes were specific to the *Brassicaceae* family, with only 33% present in
1037 dicotyledonous species outside of the *Brassicaceae*, compared to 78% of high-scoring phenotype

1038 genes (FET, $p < 4e-12$), thus our model likely has reduced power in detecting lineage-specific
1039 sequences.

1040 We also predict 80 pseudogenes (10%) to be functional (high-FL pseudogenes). A
1041 significantly higher proportion of high-FL pseudogenes came from existing genome annotation
1042 as 19% of annotated pseudogenes were classified as functional, compared to 4% of pseudogenes
1043 identified through a computational pipeline (FET, $p < 1.5E-10$) (Zou et al. 2009). We found that
1044 high-FL pseudogenes might be more recently pseudogenized and thus have not yet lost many
1045 genic signatures, as the mean number of disabling mutations (premature stop or frameshift) per
1046 kb in high-scoring pseudogenes (1.9) were significantly lower than that of low-scoring
1047 pseudogenes (4.0; U test, $p < 0.02$). Lastly, we cannot rule out the possibility that a small subset
1048 of high-scoring pseudogenes represent truly functional sequences, rather than false positives
1049 (Poliseno et al. 2010; Karreth et al. 2015). Overall, the misclassification of both narrowly-
1050 expressed phenotype genes and broadly-expressed pseudogenes highlights the need for an
1051 updated prediction model that is less influenced by expression breadth.

1052 Among protein-coding genes without phenotype information, we predict 20% as
1053 nonfunctional. We expect that at least 4% represent false negatives based on the FNR of the full
1054 model. The actual FNR among protein-coding genes may be higher, however, as phenotype
1055 genes represent a highly active and well conserved subset of all genes. However, a subset of the
1056 low-scoring protein-coding genes may also represent gene sequences undergoing functional
1057 decay and *en route* to pseudogene status. To assess this possibility, we examined 1,940 *A.*
1058 *thaliana* "decaying" genes that may be experiencing pseudogenization due to promoter
1059 disablement (Yang et al. 2011) and found that, while these decaying genes represented only 7%
1060 of all *A. thaliana* annotated protein-coding genes, they made up 45% of protein-coding genes
1061 predicted as nonfunctional (Fisher's Exact Test (FET), $p < 1E-11$).

1062 **Feature correlation between likely-functional ITRs and ncRNAs and their neighboring** 1063 **genes**

1064 ITRs and annotated ncRNAs closer to annotated genes tended to be predicted as functional
1065 (Supplementary Fig. 5A), as 57% of likely functional and 35% of likely nonfunctional ITRs and
1066 ncRNAs were proximal to neighboring genes (within the 95th percentile of intron lengths for all
1067 genes) (FET, $p < 2E-09$). Likely functional ITRs and annotated ncRNAs that are proximal to
1068 genes may frequently represent unannotated exon extensions. If this is the case, it may be

1069 expected that such sequences would exhibit similar features as gene neighbors. However, no
1070 clear pattern of increased feature correlation was observed between likely functional ITRs /
1071 ncRNAs and neighboring genes when compared to likely non-functional ITRs / ncRNAs or
1072 random sequence pairs, regardless of proximity (Supplementary Fig. 5B-D). Thus, despite their
1073 proximity to annotated genes, it remains unclear if some ITRs or annotated ncRNAs represent
1074 unannotated exon extensions of known genes or not. In addition, for proximal functional
1075 ITRs/annotated ncRNAs, we cannot rule out the possibility that they represent false-positive
1076 functional predictions due to the accessible and active chromatin states of nearby genes. Given
1077 the challenge in ascertaining the origin of likely functional, proximal ITRs/ncRNAs, we instead
1078 conservatively estimate that 187 distal, functional ITRs and annotated ncRNAs may represent
1079 fragments of novel genes.

1080