

1 **Cross-species modeling of plant genomes at single** 2 **nucleotide resolution using a pre-trained DNA** 3 **language model**

4 Jingjing Zhai^{1*}, Aaron Gokaslan^{2*}, Yair Schiff², Ana Berthel¹, Zong-Yan Liu³, Wei-Yun Lai¹,
5 Zachary R. Miller¹, Armin Scheben⁵, Michelle C. Stitzer¹, M. Cinta Romay¹, Edward S.
6 Buckler^{1,3,4+}, Volodymyr Kuleshov²⁺

7 1 Institute for Genomic Diversity, Cornell University, Ithaca, NY USA 14853

8 2 Department of Computer Science, Cornell University, Ithaca, NY, USA 14853

9 3 Section of Plant Breeding and Genetics, Cornell University, Ithaca, NY USA 14853

10 4 USDA-ARS; Ithaca, NY, USA 14853

11 5 Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, 1 Bungtown Road,
12 Cold Spring Harbor, NY USA 11724

13 * These authors contributed equally to this work

14 + To whom correspondence may be addressed. Email: jz963@cornell.edu, ed.buckler@usda.gov
15 and yk379@cornell.edu

16

17 **Abstract**

18 Interpreting function and fitness effects in diverse plant genomes requires transferable models.
19 Language models (LMs) pre-trained on large-scale biological sequences can learn evolutionary
20 conservation and offer cross-species prediction better than supervised models through fine-tuning
21 limited labeled data. We introduce PlantCaduceus, a plant DNA LM based on the Caduceus and
22 Mamba architectures, pre-trained on a curated dataset of 16 Angiosperm genomes. Fine-tuning
23 PlantCaduceus on limited labeled Arabidopsis data for four tasks, including predicting translation
24 initiation/termination sites and splice donor and acceptor sites, demonstrated high transferability
25 to 160 million year diverged maize, outperforming the best existing DNA LM by 1.45 to 7.23-
26 fold. PlantCaduceus is competitive to state-of-the-art protein LMs in terms of deleterious mutation
27 identification, and is threefold better than PhyloP. Additionally, PlantCaduceus successfully
28 identifies well-known causal variants in both Arabidopsis and maize. Overall, PlantCaduceus is a
29 versatile DNA LM that can accelerate plant genomics and crop breeding applications.

30 **Main**

31 Over 1,000 plant genomes have been published during the past 20 years, and this number will
32 continue to increase significantly in the coming decades¹⁻³. Understanding the functional elements
33 and fitness effects of these genomes at both transcriptional and translational levels is crucial for
34 advancing plant genomics and crop breeding. Unlike biomedical applications that primarily focus
35 on a few key species, plant genomics must account for the vast diversity of hundreds of crop
36 species, each with unique variations in size, composition, and complexity⁴. Extensive genomic
37 resources have been generated for model plants, such as *Arabidopsis*⁵, rice⁶ and maize⁶,
38 significantly advancing plant genomics research. However, generating analogous genomic
39 resources experimentally for all plant genomes is time-consuming, costly, and impractical. This
40 highlights the need for developing cross-species models capable of capturing evolutionary
41 conservation across diverse plant species.

42
43 Supervised deep learning (DL) sequence models are successful in understanding DNA sequence
44 functions such as transcription initiation⁷, alternative splicing⁸ and gene expression⁹. However,
45 supervised DL models typically require large-scale labeled data, such as ENCODE-scale datasets
46^{10,11}, to achieve robust performance. Such extensive labeled data is often scarce in plant genomics.
47 Moreover, training supervised models on model species, such as *Arabidopsis*, presents challenges
48 when transferring to other plant species. However, the success of self-supervised language models
49 (LMs) offers a promising alternative. In this paradigm, a foundation model is pre-trained on vast
50 amounts of unlabeled biological sequences to learn evolutionary conservation. Pre-trained models
51 are then fine-tuned on limited labeled data, enabling better performance on downstream tasks and
52 enhancing generalizability across species relative to existing methods. For example, protein LMs,
53 pre-trained on diverse protein sequences spanning the evolutionary tree, have shown successful
54 applications in predicting atomic-level protein structure¹² and disease-causing variants¹³ as well
55 as in engineering protein design¹⁴. These models provide valuable tools for understanding protein
56 function and facilitating innovative solutions in biotechnology and medicine¹⁵.

57
58 Unlike protein LMs that are limited to coding regions, DNA LMs enable a comprehensive
59 understanding of the entire genome, offering deeper insights into gene regulation and evolution.
60 Protein LMs have shown success in identifying pathogenic missense mutations in human genetics

61 ^{13,16}, but increasing evidence shows that mutations in noncoding regions, including both intergenic
62 and intronic regions, contribute significantly to both agronomic traits ¹⁷ and human diseases ^{18,19}.
63 Additionally, training multi-species DNA LMs can capture evolutionary conservation at the DNA
64 level, enhancing our understanding of genetic variation across different species.

65
66 However, DNA LMs face significant challenges compared to protein LMs. Firstly, eukaryotes,
67 especially plants ²⁰, contain varied percentages of repetitive sequences, complicating the pre-
68 training task. Given that LMs are pre-trained to either predict the next token or tokens are masked
69 arbitrarily in a sequence, repetitive sequences that are easier to predict but do not necessarily
70 improve downstream applications can reduce overall model quality ²¹. Additionally, noncoding
71 regions are less conserved than coding regions, leading to potential biases if entire genomes are
72 included in pre-training. Lastly, unlike protein sequences, modeling double-stranded DNA
73 requires consideration of reverse complementary base pairing ²² and a bi-directional model that
74 accounts for both upstream and downstream sequences.

75
76 To tackle these challenges, we introduce PlantCaduceus, a DNA language model pre-trained on a
77 curated dataset consisting of 16 angiosperm genomes (**Fig. 1A-1B**). PlantCaduceus employs
78 single-nucleotide tokenization, enabling precise modeling at the base-pair-resolution across
79 diverse plant genomes. By down-sampling noncoding regions and down-weighting repetitive
80 sequences, we generated an unbiased genomic dataset for pre-training. In contrast, other publicly
81 available DNA LMs, such as AgroNT ²³ and Nucleotide Transformer ²⁴, use entire genomes for
82 pre-training, potentially introducing biases toward certain genomes and repetitive sequences.
83 Additionally, both models use non-overlapping kmer tokenizers that disrupt the genome into
84 arbitrary segments. Unlike the unidirectional HyenaDNA ²⁵ or Evo ²⁶, PlantCaduceus offers bi-
85 directional context, providing a more comprehensive understanding of DNA interactions.
86 Furthermore, to handle double-stranded DNA, we used the Caduceus architecture ²⁷, which builds
87 on the Mamba architecture ²⁸ and supports reverse complement equivariance, unlike GPN ²¹, which
88 uses convolutional neural network and manually augments reverse complement sequences. By
89 evaluating the pre-trained PlantCaduceus model on five cross-species tasks, including translation
90 initiation/termination sites, splice donor and acceptor sites, and evolutionary conservation
91 prediction. We found that our model demonstrated the best performance compared to baseline

92 models for all five tasks. Notably, downstream classifiers fine-tuned on PlantCaduceus with
93 limited labeled data in Arabidopsis maintained the best performance on other crop species such as
94 maize, improving the PRAUC from 1.45-fold to 7.23-fold as compared to the best existing DNA
95 LM, indicating that PlantCaduceus effectively captures broad evolutionary conservation.
96 Additionally, deleterious mutations identified with the zero-shot strategy of PlantCaduceus
97 showed a three-fold enrichment of rare alleles when compared to the most commonly used
98 evolutionary-based methods such as phyloP and phastCons²⁹. For missense mutations,
99 PlantCaduceus matches the performance of state-of-the-art protein LMs, suggesting that
100 PlantCaduceus can be effectively used for genome-wide deleterious mutation identification.
101 Furthermore, PlantCaduceus successfully identifies well-known causal variants in both
102 Arabidopsis and maize. These results indicate that PlantCaduceus can serve as a foundational
103 model to accelerate plant genomics and crop breeding applications.

104 **Results**

105 **PlantCaduceus: a pre-trained DNA language model with 16 Angiosperm genomes**

106 Caduceus²⁷ is a DNA LM architecture that builds upon the recently introduced Mamba²⁸
107 architecture, a selective state space sequence model that has demonstrated competitive
108 performance to transformers³⁰ in various natural language processing tasks, with more efficient
109 scaling for longer range sequences. Unlike Mamba, Caduceus is specifically designed for DNA
110 sequences, taking into account the bi-directional nature of DNA and introducing reverse
111 complement (RC) equivariance. Here, we trained PlantCaduceus using the Caduceus architecture
112 on 16 Angiosperm genomes (**Fig. 1A-1B; Supplemental Table 1**), spanning 160 million years of
113 evolutionary history (**METHODS**). PlantCaduceus takes 512 base pair (bp) windows of input
114 sequences, tokenizing them into single nucleotides, and is pre-trained using a masked language
115 modeling objective (**Fig. 1B; METHODS**). To address the substantial variation in genome sizes
116 and the high proportion of repetitive sequences in these genomes, we emphasized non-repetitive
117 sequences by down-weighting and down-sampling repetitive sequences during pre-training
118 (**METHODS**). To scale Caduceus, we trained a series of PlantCaduceus models with parameter
119 sizes ranging from 20 million to 225 million (**Table 1**). The training and validation losses for each
120 model are detailed in **Supplemental Table 2**. After pre-training, we conducted a preliminary

121 assessment to verify the model's learning capabilities. Taking the sorghum genome as an example,
122 we employed Uniform Manifold Approximation and Projection (UMAP) ³¹ to visualize the
123 embeddings generated by the four pre-trained PlantCaduceus models. By segmenting the genome
124 into 512 bp windows, we observed distinct clustering in the UMAP visualization, corresponding
125 to different genomic regions (**Fig. 1C**). Due to the high proportion of repetitive intergenic
126 sequences in the sorghum genome, the embedding spaces appeared dispersed in the UMAP
127 visualization (**Fig. 1D; Supplemental Fig. 1**). Even without any supervision, PlantCaduceus was
128 able to differentiate between coding and noncoding regions with high clarity.

129 **Improving the accuracy and cross-species transferability of modeling transcription and** 130 **translation through fine-tuning PlantCaduceus**

131 Transcription and translation are two key processes in the central dogma of molecular biology, and
132 the precise identification of junction sites during these processes is essential for comprehensive
133 gene annotation. To assess PlantCaduceus's performance in modeling these processes, we
134 designed four gene annotation tasks: predicting the translation initiation site (TIS), translation
135 termination site (TTS), and splice donor and acceptor sites (**METHODS**). We employed a feature-
136 based approach to fine-tune PlantCaduceus by keeping the pre-trained model weights frozen while
137 training XGBoost models using embeddings extracted from the last hidden state of PlantCaduceus
138 (**Fig. 2A**). Compared to full fine-tuning, this approach allows us to leverage the rich
139 representations learned by PlantCaduceus while minimizing the usage of computational resources.
140 Previous LMs focus on evaluation within the same species ^{24,25,32-34}. However, given that the DNA
141 LM model is pre-trained on multiple species, we wanted to investigate whether a model fine-tuned
142 with limited labeled data in Arabidopsis could be used for prediction in other species. Therefore,
143 we trained and validated XGBoost models in Arabidopsis and tested their performance on both
144 species included (*Oryza sativa* and *Sorghum bicolor*) and not included (*Gossypium hirsutum*,
145 *Glycine max* and *Zea mays*) in the pre-training (**Fig. 2B; Supplemental Table 3**). We
146 benchmarked the performance of PlantCaduceus against three DNA LMs: GPN ²¹, AgroNT ²³, and
147 Nucleotide Transformer ²⁴, as well as a supervised hybrid model comprising a convolutional neural
148 network (CNN) and a long short-term memory (LSTM) network ³⁵, hereafter referred to as
149 CNN+LSTM. For DNA LMs, we used the same feature-based approach as PlantCaduceus (**Fig.**

150 **2A)** to train XGBoost models using embeddings extracted from the last hidden state of each DNA
151 LM (**Fig. 2C**). The CNN+LSTM model was trained from scratch in a supervised manner.

152
153 First, focusing on within species evaluation on Arabidopsis hold-out test set, PlantCaduceus (32
154 layers) showed consistently superior performance across the four tasks of predicting TIS (**Fig. 2C**),
155 TTS (**Fig. 2D**), splice donor site (**Fig. 2E**), and splice acceptor site (**Fig. 2F**). Other DNA LMs
156 like GPN and AgroNT also performed well, particularly in predicting splice donor and acceptor
157 sites. Additionally, for splice donor and acceptor site prediction, even the supervised CNN+LSTM
158 model achieved near perfect PRAUC values, indicating that within-species prediction is a
159 relatively straightforward task.

160
161 We then assessed the cross-species generalization ability of these models by testing them on *O.*
162 *sativa* and *S. bicolor*, which were included in pre-training, as well as *G. hirsutum*, *G. max*, and *Z.*
163 *mays*, which were not (**Fig. 2B; Fig. 1A**). When tested across these five species, all models except
164 PlantCaduceus exhibited a significant drop in average PRAUC, decreasing from 0.789 in *A.*
165 *thaliana* to 0.237 in these species (**Fig. 2C-2F**). For instance, transferring the supervised
166 CNN+LSTM model to *Z.mays*—which diverged 160 million years ago—resulted in a PRAUC
167 drop from 0.713 to nearly zero for the TIS task. This significant drop was expected, as the
168 supervised model had never seen sequences from these species, making cross-species
169 generalization challenging. Although GPN maintained decent cross-species predictions, it still
170 showed significant performance drops, with the average PRAUC decreasing from 0.944 in *A.*
171 *thaliana* to 0.509 in other species (**Fig. 2C-2F; Supplemental Table 4**). As expected, the non-
172 plant DNA NT-v2 model performed poorly on these tasks due to the significant divergence
173 between plant and animal genomes. Even though AgroNT was pre-trained on 48 plant genomes,
174 its performance fell short of expectations in cross-species evaluations. In contrast, PlantCaduceus
175 consistently maintained high PRAUC values across all species, with an average PRAUC of 0.764,
176 regardless of whether the species were included in pre-training, demonstrating its superior
177 generalization ability across diverse plant species (**Fig. 2C-2F; Supplemental Table 4**).

178
179 GPN, as the second-best DNA LM, was not pre-trained on any of the five testing species. To ensure
180 a fairer comparison with GPN and to understand why PlantCaduceus achieved superior

181 performance on these cross-species tasks, we conducted an ablation test by re-training a custom
182 GPN model (**METHODS**) using the same datasets as PlantCaduceus and scaling it to 130 million
183 parameters, on the same order of magnitude as PlantCaduceus. We observed that including more
184 genomes in the pre-training and scaling model size significantly improved GPN's cross-species
185 predictability (**Supplemental Fig. 2; Supplemental Table 4**), especially for TIS and TTS tasks.
186 This indicates that when more genomes are included during pre-training, the embeddings learned
187 by DNA LMs are more general across species. However, PlantCaduceus still exhibited the best
188 performance, indicating that its architecture is superior to that of GPN. Moreover, even with a
189 parameter size of 20 million—6.5 times smaller than the custom 130 million GPN and 3.25x times
190 smaller than the original GPN—PlantCaduceus still outperformed all models in predicting TIS,
191 TTS, splice donor, and splice acceptor sites. These results demonstrate that PlantCaduceus not
192 only captures broader evolutionary conservation features but also is more parameter-efficient than
193 other DNA LMs.

194 **Cross-species evolutionary constraint prediction through fine-tuning PlantCaduceus**

195 Genome-wide association studies (GWAS) have identified thousands of variants associated with
196 complex traits³⁶. However, identifying causal variants is complicated by linkage disequilibrium
197 (LD), as significant SNPs identified by GWAS are usually in LD with causal variants³⁷. In contrast,
198 evolutionary constraint, as evidenced by DNA conservation across species, can identify potential
199 causal mutations by revealing their fitness effects³⁸. Given that PlantCaduceus is pre-trained on
200 16 Angiosperm genomes, we hypothesize that it can be fine-tuned to predict evolutionary
201 constraint using DNA sequences alone. Maize and sorghum are both members of the
202 Andropogoneae clade, descended from a common ancestor approximately 18 million years ago³⁹.
203 To generate evolutionary constraints in the sorghum genome, we aligned 34 genomes from the
204 Andropogoneae clade, with rice as an outgroup (**Supplemental Table 5**), to the *Sorghum bicolor*
205 reference genome (**Supplemental Fig. 3**). We focused on the 277 million sites with nearly
206 complete coverage and defined those sites with an identity threshold of 15 as conserved versus
207 neutral with an identity threshold of 15 (**Fig. 3A**). We used sites chromosomes 1 to 9 to train an
208 XGBoost model and evaluated it on sorghum chromosome 10. As mentioned above, we
209 benchmarked this task against GPN, AgroNT, NT-v2, and the supervised CNN+LSTM model. On
210 the validation set, PlantCaduceus achieved the best performance, with an AUC of 0.896 (**Fig. 3B**)

211 and a PR-AUC of 0.876 (**Fig. 3C**). In comparison, the best AUC and PR-AUC for other DNA
212 LMs were 0.778 and 0.790, respectively. As expected, the supervised CNN+LSTM model
213 performed the worst, with an AUC of 0.638, as it had only seen sequences from sorghum (**Fig.**
214 **3B-3C**). This demonstrates that PlantCaduceus enables predicting evolutionary constraint without
215 multiple sequence alignment.

216

217 To further explore the cross-species predictive power of the model fine-tuned on sorghum
218 evolutionary constraint data, we generated an analogous testing dataset for maize (**METHODS**).
219 Remarkably, when our PlantCaduceus model, originally fine-tuned on sorghum, was applied to
220 the maize dataset, it demonstrated strong cross-species prediction performance, achieving an AUC
221 of 0.829 (**Fig. 4D**) and a PR-AUC of 0.797 (**Fig. 4E**). In contrast, all other models consistently
222 showed poor performance on maize (**Fig. 4D-4E**). We also evaluated the performance of our
223 custom GPN model which was trained on the same dataset as PlantCaduceus. While the custom
224 GPN model showed improved performance with an AUC of 0.833 and a PR-AUC of 0.814,
225 PlantCaduceus, with only 20 million parameters, outperformed both the original GPN and the
226 custom GPN models (**Supplemental Fig. 4**). These results highlight the robustness and
227 effectiveness of our DNA LM for cross-species predictions of evolutionary constraints using only
228 sequence data as input. The transferability of our model across different species within the
229 Andropogoneae clade suggests that it captures fundamental evolutionary patterns and can be
230 readily adapted to predict evolutionary constraint in related species with limited additional training
231 data.

232 **Zero-shot variant effect prediction identifies deleterious mutations in different species**

233 The training objective of PlantCaduceus is to predict masked nucleotides based on sequence
234 context; if a pre-trained multi-species DNA LM can accurately predict masked tokens, it suggests
235 that similar sequence patterns, conserved across different species, were frequently observed during
236 pre-training. We hypothesize that the predicted likelihood of the reference allele versus the
237 alternate allele can identify deleterious mutations, as mutations in conserved regions across species
238 are likely deleterious⁴⁰⁻⁴³. To test this hypothesis, we employed the same zero-shot strategy as
239 GPN²¹ to estimate the effect of each mutation (**Fig. 4A**). Specifically, for each mutation, we
240 calculated the log-likelihood difference between the reference and alternate alleles, where a more

241 negative value indicates higher conservation. We generated 1.1 million sites from sorghum
242 chromosome 8 (included in pre-training) and 1.3 million sites from maize chromosome 8 (not
243 included) through in silico mutagenesis of SNPs. We then calculated zero-shot scores for these
244 mutations to assess how PlantCaduceus performs on both seen and unseen genomes. As expected,
245 mutations in highly conserved functional regions—such as stop-gained, splice acceptor, splice
246 donor, and start-lost sites—exhibited the most negative zero-shot scores, underscoring their
247 potential deleterious effects (**Fig. 4B; Supplemental Fig. 5A**). Missense mutations also showed
248 notably negative zero-shot scores. In contrast, intergenic regions and introns displayed scores
249 closer to zero, indicating lower evolutionary constraint and a reduced likelihood of deleterious
250 effects (**Fig. 4B; Supplemental Fig. 5A**). However, we observed that a subset of mutations in
251 repetitive regions still received very low zero-shot scores, suggesting that repetitive regions may
252 be too easy for the model to predict the masked tokens. Overall, the zero-shot score of
253 PlantCaduceus aligns with established concepts of deleteriousness^{44,45}.

254
255 Besides in silico mutagenesis, we also evaluated if zero-shot score can be used to identify
256 deleterious mutations in natural populations. Deleterious mutations tend to have lower frequencies
257 within a population due to selective constraints³⁸, we therefore used minor allele frequency (MAF)
258 to quantify the deleteriousness of mutations predicted by different methods. Despite the potential
259 for low MAF in neutral/beneficial alleles, we believe this approach provides useful signals for
260 assessing deleterious mutations³⁸. We benchmarked PlantCaduceus against two evolutionary-
261 informed methods, phyloP and phastCons²⁹, as well as GPN²¹. Both phyloP and phastCons assess
262 evolutionary constraint using multiple sequence alignments and phylogenetic models
263 (**METHODS**), assigning higher scores to conserved regions. We analyzed 4.6 million SNPs in the
264 sorghum TERRA population⁴² and 9.4 million SNPs from maize Hapmap 3.2.1 population⁴⁶ and
265 observed that most of the SNPs had neutral zero-shot score, while there was still a heavy tail with
266 negative zero-shot scores (**Fig. 4C; Supplemental Fig. 5B**). By defining the top 0.1% as the most
267 deleterious mutations, we observed a significant enrichment in coding regions, as reflected by the
268 high odds ratios in both sorghum (40.70) and maize (42.42) with p-values less than 2.2e-16
269 (**Supplemental Fig. 6**). We then categorized SNPs into four percentiles based on zero-shot scores:
270 the top 50%, 10%, 1%, and 0.1% most deleterious mutations and observed that all models showed
271 a decreasing average MAF of SNPs in higher percentiles for missense, nonsynonymous, and

272 noncoding SNPs in both sorghum (**Supplemental Fig. 7**) and maize (**Fig. 4D**). Notably, the
273 putative deleterious mutations identified by PlantCaduceus exhibited the lowest average MAF
274 across all percentiles, outperforming GPN and significantly surpassing phyloP and phastCons
275 (**Supplemental Fig. 7; Fig. 4D**). Given the success of protein LMs in predicting deleterious
276 missense mutations^{13,16}, we also incorporated ESM¹² as a benchmark. For missense mutations,
277 we found that PlantCaduceus matches the performance of the state-of-the-art protein LM ESM¹².
278 At the top 50%, 10%, and 1% percentiles, PlantCaduceus even slightly outperforms ESM in
279 sorghum (**Supplemental Fig. 7**).

280

281 However, since GPN is only pre-trained with genomes from eight Brassicales species and
282 specifically designed for mutation effect prediction in Arabidopsis, we further validated
283 PlantCaduceus by analyzing over 10 million mutations from the Arabidopsis 1001 Genomes
284 Project⁴⁷. Being pre-trained with a broader range of evolutionarily distant genomes,
285 PlantCaduceus effectively captured deleterious mutations in Arabidopsis and slightly
286 outperformed GPN (**Supplemental Fig. 8**). For missense mutations, PlantCaduceus consistently
287 matched the performance of the state-of-the-art protein language model ESM and was nearly
288 competitive with GPN for noncoding mutations.

289

290 We further verified if PlantCaduceus could pinpoint known causal deleterious mutations. We
291 collected 19 candidate phenotype-impacting and potentially deleterious mutations identified in
292 homozygous EMS mutants in Arabidopsis⁴⁸. Among these, 15 mutations were ranked in the top
293 1% or top 10% by the zero-shot score (**Table 2**), highlighting the zero-shot score of PlantCaduceus
294 can be used for pinpointing causal deleterious mutations. Additionally, PlantCaduceus
295 successfully identified a well-studied causal sweet corn mutation, which derives its characteristic
296 sweetness from the W578R mutation at the *sugary1* (*Su1*) locus⁴⁹. This mutation disrupts starch
297 metabolism, leading to the accumulation of phytoglycogen, which lowers seedling vigor and
298 reduces germination, ultimately decreasing fitness⁵⁰. Although GWAS revealed numerous
299 significantly sweet-trait-associated variants on chromosome 4, identifying the exact causal
300 mutations is challenging due to high LD in this low recombination region (**Fig. 5A**). By integrating
301 zero-shot scores from PlantCaduceus with GWAS data (**Fig. 5B-5C**), we successfully identified
302 the W578R mutation as the sole causal variant in this QTL region (**Fig. 5D**). Taken together, these

303 results demonstrate that the PlantCaduceus model effectively pinpoints known causal deleterious
304 mutations, highlighting its potential as a powerful tool for identifying causal variants underlying
305 important agronomic traits.

306 **Discussion**

307 Functional annotation of plant genomes is crucial for plant genomics and crop breeding but
308 remains limited by the lack of functional genomic data and accurate predictive models. Here, we
309 introduced PlantCaduceus, a multi-species plant DNA LM pretrained on a curated set of 16
310 evolutionarily distant Angiosperm genomes, enabling cross-species prediction of functional
311 annotations with limited data. PlantCaduceus leverages Mamba²⁸ and Caduceus²⁷ architectures
312 to support bi-directional, reverse complement equivariant sequence modeling. We demonstrated
313 the superior cross-species performance of PlantCaduceus on five tasks involving transcription,
314 translation, and evolutionary constraint modeling. These results highlight the potential of
315 PlantCaduceus to serve as a foundational model for comprehensively understanding plant genomes.

316
317 PlantCaduceus has the potential to accurately annotate any newly sequenced Angiosperm
318 genomes. Unlike supervised deep learning models that easily overfit on limited labeled data,
319 PlantCaduceus demonstrates robust cross-species performance in modeling transcription,
320 translation, and evolutionary constraints, even for species not included in pre-training (**Fig. 2;**
321 **Supplemental Fig. 2**). This indicates that through self-supervised pre-training on large-scale
322 genomic datasets, PlantCaduceus has captured broad evolutionary conservation and DNA
323 sequence grammar. The cross-species prediction ability of PlantCaduceus can significantly
324 accelerate plant genomics research, aiding initiatives such as the 1000 Plant Genomes Project¹ by
325 providing accurate annotations and insights across diverse plant species.

326
327 PlantCaduceus offers a more effective approach to estimate deleterious mutations without relying
328 on multiple sequence alignments (MSAs). Deleterious mutations are considered as the genetic
329 basis of heterosis, where hybrids yield more due to the suppression of deleterious recessives from
330 one parent by dominant alleles from the other⁵¹. Historically, deleterious mutations have been
331 estimated by generating MSAs^{38,52,53} and using evolutionary methods such as phyloP and
332 phastCons²⁹. However, the prevalence of transposable elements and polyploidy in plant genomes

333 complicates the genome-wide MSA generation ^{54,55}. PlantCaduceus overcomes these challenges
334 by using a masked language modeling strategy to learn conservation from large scale genomic
335 datasets of diverse species. Promisingly, the deleterious mutations prioritized by PlantCaduceus
336 with the zero-shot strategy showed three-fold rare allele enrichment compared to phyloP and
337 phastCons, and our approach is competitive with state-of-the-art protein LM for missense
338 mutations. Furthermore, PlantCaduceus enables pinpointing causal variants from significant
339 GWAS signals, which are usually confounded by LD. These results suggest that PlantCaduceus
340 can be utilized as a powerful tool in crop breeding, enhancing genome-wide deleterious mutation
341 identification, optimizing parental line selection, and promoting hybrid vigor ⁵¹.

342

343 In future work, we plan to incorporate additional plant genomes from diverse lineages, such as
344 gymnosperms, to capture broader evolutionary conservation. Additionally, we plan to pre-train
345 PlantCaduceus with longer context windows, enabling it to capture long-range DNA interactions
346 and better handle tasks benefiting from long-range cis-effects, such as allele-specific expression,
347 chromatin state prediction, and chromatin interaction mapping. Furthermore, it would also be
348 interesting to explore how to better tokenize repetitive sequences in plant genomes. We envision
349 that these approaches will allow us to push the boundaries of what PlantCaduceus can achieve,
350 establishing it as an even more powerful and versatile foundation model for advancing genomic
351 research and facilitating crop improvement.

352 **Methods**

353 **Pre-training dataset**

354 The pre-training dataset comprises 16 genomes from two distinct clades: eight genomes from the
355 family Poaceae and eight genomes from the order Brassicales (**Supplemental Table 1**). To
356 visualize their relatedness, we subset these taxa from a large phylogeny of seed plants⁵⁶. The
357 Poaceae species displayed substantial variation in genome size and repetitive sequence content,
358 with the hexaploid wheat genome exhibiting a size of 15 Gbp. For each Poaceae genome, except
359 for *Tripsacum*, we obtained the genome and corresponding genome annotation and repeat-masked
360 annotation from the Joint Genome Institute (JGI). For the *Tripsacum* genome, the genome FASTA
361 and annotation files were downloaded from MaizeGDB
362 (https://maizegdb.org/genome/assembly/Td-FL_9056069_6-DRAFT-PanAnd-1.0), and the
363 EDTA tool⁵⁷ was used to identify repetitive sequences within the genome. Based on the repeat-
364 masked annotation, each genome was softmasked with bedtools⁵⁸ and subsequently divided into
365 genomic windows of 512 bp with a step size of 256 bp. Each window was assigned to a unique
366 class based on the genome annotation, and all coding sequence regions were selected for pre-
367 training. The remaining genomic regions were then down-sampled to ensure an equal number of
368 CDS regions and noncoding regions. It is important to note that for the hexaploid wheat genome,
369 only subgenome A was utilized to avoid species bias. The Brassicales genomes datasets were
370 acquired from a Hugging Face repository ([https://huggingface.co/datasets/songlab/genomes-
371 brassicales-balanced-v1](https://huggingface.co/datasets/songlab/genomes-brassicales-balanced-v1)). The validation and testing datasets were randomly selected and
372 constituted 5% of the total dataset.

373 **Caduceus model architecture and pre-training**

374 We use the recently proposed Caduceus architecture²⁷, which is tailored to DNA sequence
375 modeling. Caduceus is based on the Mamba architecture²⁸, a model which scales to long sequences
376 more efficiently than attention-based methods while maintaining accuracy. Mamba stems from the
377 class of structured state space models (SSMs)⁵⁹, which are defined by a pair of linear differential
378 equations:

379

380

$$\dot{h}(t) = A_t h(t) + B_t x(t),$$

381
$$y(t) = C_t h(t) + D_t x(t),$$

382 where $x, y \in \mathbb{R}$ represent the input and output, respectively, $h \in \mathbb{R}^n$ is the state's hidden
383 representation, and the (potentially time dependent) parameters $A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times 1}, C \in$
384 $\mathbb{R}^{1 \times n}, D \in \mathbb{R}$ govern the system dynamics. For multi-dimensional inputs and outputs $x, y \in \mathbb{R}^d$, a
385 separate linear system is applied to each of the d channels. In practice, using some discretization
386 scheme that is a function of a discrete time parameter Δ , the system is discretized in time, yielding
387 the following:

388
$$h_{t+1} = \bar{A}_t h_t + \bar{B}_t x_t$$

389
$$y_{t+1} = C_t h_{t+1} + D_t x_t$$

390 Much of the SSM literature relies on parameters that are fixed in time, allowing for efficient
391 computation during training by means of the convolutional perspective of linear *time invariant*
392 systems⁶⁰. In contrast to previous SSMs, Mamba enables more expressive models that are *time*
393 *dependent*, by making the parameters functions of the inputs. This time dependence is crucial in
394 allowing Mamba to overcome the limitations of previous SSMs and rival Transformers³⁰ on
395 sequence modeling tasks across domains. For efficient computation, Mamba employs a parallel
396 algorithm to compute the recurrence relation defined above and an IO-aware implementation that
397 limits potentially bottlenecking memory transfer operations incurred on modern GPU hardware.

398
399 To account for upstream and downstream gene interactions, Caduceus employs weight sharing to
400 enable memory-efficient bi-directionality. Finally, Caduceus is designed to consider the reverse
401 complement (RC) symmetry of DNA sequences. This is accomplished by encoding RC
402 equivariance as an inductive bias: the Caduceus language model commutes with the RC operation.
403 Combining these three design decisions, Caduceus has shown promising results when applied to
404 human genome modeling²⁷.

405
406 The implementation of RC equivariance in Caduceus entails doubling the number of channels for
407 intermediate representations. At a high level, half the channels are used to encode information
408 about a sequence and the other half are used to encode information about its RC. For downstream
409 tasks in which we fine-tuned a classifier on top of learned embeddings, the labels were invariant
410 to the RC operation, since both DNA strands carry the same label. To account for this, we therefore
411 split embeddings of the Caduceus model along the channel dimension and averaged. This ensures

412 that both a sequence and its RC will have the same final embedding, i.e., we render the embeddings
413 invariant to the RC operation as well.

414

415 For the pre-training of PlantCaduceus, each model was trained for 240,000 steps using a
416 Decoupled AdamW optimizer⁶¹ with the global batch size of 2,048. The learning rate is 2E-4 with
417 a cosine decay scheduler, and 6% of the training duration was dedicated to warm up. The learning
418 rate decayed to 4E-6 by the end of training. The default BERT⁶² masking recipe was used with a
419 masking probability of 0.15. For each masked token: (i) there is an 80% probability it will be
420 replaced by a special token ([MASK]), (ii) a 10% probability it will be replaced by a random token,
421 and (iii) a 10% probability it will remain unchanged. Unless otherwise specified, all models were
422 trained using a sequence length of 512 base pairs. A weight decay of 1E-5 was applied throughout
423 the training process.

424 **TIS, TTS, splice donor and acceptor training, validation and testing dataset generation**

425 To generate high-quality training datasets for translation initiation sites (TIS), translation
426 termination sites (TTS), splice donor sites, and splice acceptor sites, we used the well-annotated
427 model plant genome of Arabidopsis with Araport 11 annotation⁶³. To accurately reflect the
428 inherent imbalance in junction sites prediction, all annotated junction sites were considered as
429 positive observations, while a randomly selected subset of sites (5%) that matched specific
430 appropriate motifs (e.g., ATG for TIS, UAA, UAG, and UGA for TTS, GT for donor splice sites,
431 and AG for acceptor splice sites) were used as negative observations. For each task, the pre-trained
432 model weights were frozen, and XGBoost models (`n_estimators=1000`, `max_depth=6`,
433 `learning_rate=0.1`) were trained using embeddings extracted from the last hidden state of the pre-
434 trained model. To ensure robust model training and validation, chromosome 5 was used for hold-
435 out testing, and the rest of the Arabidopsis genome was used for training.

436

437 Given the relatively poor annotation in other species compared to Arabidopsis, we used the
438 BUSCO tool⁶⁴ to identify 3,236 orthologous genes specific to monocotyledons in *O. sativa*, *S.*
439 *bicolor* and *Z. mays* and 2,326 orthologous genes specific to eudicotyledons in *G. hirsutum* and *G.*
440 *max* to generate reliable testing datasets in other species. This approach ensures that the selected
441 annotated genes are highly conserved and likely to be correctly annotated, mitigating the issue of

442 inaccurate performance evaluations. Specifically, BUSCO was utilized to scan the annotated
443 protein isoforms, and only complete BUSCO genes were considered as true positives. For those
444 BUSCO genes with multiple transcripts, we selected the longest transcript to avoid sequence
445 redundancy in the testing dataset. Subsequently, BUSCO gene/transcript-supported junction sites
446 were used as positive examples for their respective tasks. To generate negative sites, all sites within
447 the BUSCO genes that matched appropriate motifs (e.g., ATG for TIS, TAA, TAG, and TGA for
448 TTS, GT for donor splice sites, and AG for acceptor splice sites) but were not part of any annotated
449 gene models were used as true sites. Sites belonging to alternate transcripts were excluded to avoid
450 ambiguity. Furthermore, to expand the negative observations and capture a broader range of non-
451 junction sites, we included sites in the intergenic regions flanking the BUSCO genes that matched
452 the appropriate junction motifs. By incorporating both genic and intergenic sites from the BUSCO
453 gene set as negatives, we created an extremely imbalanced testing dataset to reflect the real-world
454 scenario of junction site prediction (**Supplemental Table 3**).

455 **Evolutionary constraint estimation**

456 The evolutionary constraint was estimated primarily within the Andropogoneae tribe, a large clade
457 of grasses comprising approximately 1,200 species that descended from a common ancestor
458 approximately 18 million years ago³⁹. In this analysis, 34 genomes from Andropogoneae and the
459 rice genome were used to estimate the evolutionary constraint. Due to the substantial transposable
460 element (TE) content in these genomes, AnchorWave, a sensitive genome-to-genome alignment
461 tool⁵⁴, was used to align the 35 genomes to the sorghum reference genome using the parameters
462 "-R 1 -Q 1". Following the alignments to the sorghum reference genome, we counted the number
463 of identities, SNPs, and coverages (**Supplemental Fig. 3**). Then the fine-tuned labels were
464 generated based on per-site identity and coverage (**Fig. 3A**). Conserved sites were defined as
465 having an identity greater than 34, while neutral sites were defined as having an identity of 15 or
466 less and coverage of at least 34. Sites with low coverage were excluded due to their potential
467 ambiguity. Given the large size of the training dataset, only 5% of conserved sites were randomly
468 selected for training, and an equivalent proportion of neutral sites was also randomly selected.
469 Sites from chromosomes 1 to 9 were used for training, while those from chromosome 10 were
470 used for validation. To generate the testing dataset in maize, the maize reference genome B73 was
471 used. Then, using the same approach, genome-wide evolutionary constraints were generated by

472 aligning 35 genomes to the maize reference genome with AnchorWave, using the parameters "-R
473 1 -Q 2," except for Tripsacum clades. For Tripsacum and maize, which share the most recent whole
474 genome duplication, we used "-R 1 -Q 1".

475 **phyloP and phastCons calculation**

476 With the same 34 genomes from Andropogoneae, we generated pairwise genome-to-genome
477 alignments using Cactus⁶⁵, a multiple genome alignment tool that uses a progressive alignment
478 strategy. The neutral model was calculated from fourfold degenerate coding sites across the entire
479 genome. The resulting alignments were then analyzed using PHAST²⁹ to quantify evolutionary
480 conservation with phyloP conservation scores – using the SPH scoring method (--method SPH)
481 and CONACC mode (--mode CONACC) – and phastCons scores.

482 **In silico mutagenesis.**

483 All potential mutations in the genic regions and 1 kb flanking regions of maize and sorghum
484 chromosome 8 were generated and annotated using the Ensembl Variant Effect Predictor (VEP)
485 local API⁴⁴, with the upstream/downstream parameter set to 1,000 to classify variants as either
486 upstream or downstream. For intergenic variants, we randomly sampled 100,000 SNPs from the
487 intergenic regions across chromosome 8 to ensure more even coverage of the entire chromosome.
488 For each variant type, we randomly sampled 100,000 mutations and calculated zero-shot scores.

489 **Genome-wide association study for sweet phenotype**

490 To perform a GWAS for the sweet phenotype, we used a subset of genotypes from the Hapmap
491 3.2.1 population ⁴⁶, where sweet phenotype data is available ⁶⁶. This subset consists of 272 diverse
492 inbred lines with recorded sweet phenotype data. We coded starchy corn as 0 and sweet corn as 1,
493 with 266 entries in the first category and 6 in the second. To map the sweet phenotype, we utilized
494 a model specifically designed to account for population structure: $y = X\beta + 5 \text{ global PCs} + e$. The
495 methods for GWAS followed those outlined in Kpaipho-Burch et al ⁶⁷. Briefly, the five global
496 principal components (PCs) were derived from 66,527 SNPs across 3,545 diverse inbred lines, and
497 the SNPs from 272 inbred lines were then rotated to such PCs. The selected SNPs had no missing
498 data across three maize populations, ensuring effective control for population structure and kinship.
499 This approach also reduced computational time compared to mixed linear models while
500 maintaining consistent trait mapping across populations.

501

502 **GPN, custom GPN, AgroNT and NT-v2 baselines**

503 To comprehensively evaluate our foundation model's performance, four foundation models
504 including GPN ²¹ (<https://huggingface.co/songlab/gpn-brassicales>), custom GPN, AgroNT ²³
505 (<https://huggingface.co/InstaDeepAI/agro-nucleotide-transformer-1b>) and NT-v2 ²⁴
506 (<https://huggingface.co/InstaDeepAI/nucleotide-transformer-v2-500m-multi-species>) were used
507 as baselines for various tasks. GPN is a convolutional DNA LM pre-trained on eight genomes of
508 Arabidopsis and seven other species from the Brassicales order. However, since GPN was pre-
509 trained with only eight evolutionarily close species and has only 65M parameters and most of the
510 tasks in this paper focus on evaluation in crops, we re-trained a custom GPN with 130M parameters
511 using 50 convolutional layers and the same dataset as PlantCaduceus for a fair comparison. The
512 other hyperparameters were kept identical to the original GPN (**Supplemental Table 6**). In
513 contrast, AgroNT ²³ is a transformer-based ³⁰ language model with 1 billion parameters, pre-trained
514 on 48 plant genomes. NT-v2 ²⁴, is a non-plant multi-species transformer model pre-trained on 850
515 genomes excluding plant species. These models employ different tokenization strategies: GPN
516 uses single-nucleotide tokenization, while AgroNT and NT-v2 use 6-mer tokenization. To ensure
517 a fair comparison, we extracted the middle token embeddings for GPN and the middle k-mer token
518 embeddings for AgroNT and NT-v2.

519 **Supervised CNN+LSTM baseline**

520 To establish a fair comparison between our DNA LM and existing supervised models, which are
521 primarily trained on human data, we used the DanQ model architecture ³⁵ as the supervised
522 baseline. DanQ is a hybrid convolutional and recurrent neural network specifically designed for
523 predicting the function of DNA sequences. It has demonstrated impressive performance in
524 predicting chromatin states in plant species, making it a suitable choice for our comparative
525 analysis ⁶⁸. For each task, the CNN+LSTM model was trained from scratch using one-hot encoded
526 DNA sequences as input. The Adam optimizer with a learning rate of 0.01 was employed for model
527 optimization. The batch size was set to 2,048. Early stopping with a patience of 20 steps was
528 implemented.

529 **Data availability**

530 The pre-training genomes are available at: [https://huggingface.co/datasets/kuleshov-](https://huggingface.co/datasets/kuleshov-group/Angiosperm_16_genomes)
531 [group/Angiosperm_16_genomes](https://huggingface.co/datasets/kuleshov-group/Angiosperm_16_genomes). All datasets used for fine-tuning are available at Hugging Face:
532 <https://huggingface.co/datasets/kuleshov-group/cross-species-single-nucleotide-annotation>
533

534 **Code availability**

535 The pre-trained models, along with documentation on how to use them, are available at Hugging
536 Face: <https://huggingface.co/collections/kuleshov-group/plantcaduceus-512bp-len-665a229ee098db706a55e44a>. The pre-training and fine-tuning codes are available at GitHub:
537 <https://github.com/kuleshov-group/PlantCaduceus>

539

540 **Acknowledgments**

541 This work is funded by the USDA-ARS, NSF PanAnd grant (#1822330), NSF CAREER grant
542 (#2145577) and NIH MIRA grant (#1R35GM151243-01). We thank Edgar Marroquin (Cornell
543 University) for discussing fine-tuning tasks, Travis Wrightsman (Cornell University) for providing
544 DanQ code, Arun S. Seetharam and Matthew B Hufford (Iowa State University) for sharing
545 Andropogoneae assemblies, Merritt Khaipho-Burch (Cornell University) for sharing the leftover
546 version HapMap3 VCF file, Sara Miller (Cornell University) for helpful comments and all
547 members of the E.S.B. laboratory (Cornell University) for helpful discussions. We would also like
548 to thank the SCINet project, the AI Center of Excellence of the USDA Agricultural Research
549 Service (0201-88888-003-000D and 0201-88888-002-000D) and MosaicML for providing
550 compute resources for pre-training and fine-tuning experiments.

551 **Author contributions**

552 J.Z., A.G., E.S.B., and V.K. designed the research; J.Z., A.B., Z.-Y.L., Z.R.M., A.S., M.C.S. and
553 C.R. curated the data; J.Z., A.G., Y.S., and Z.R.M. pre-trained models; J.Z., A.B., Z.-Y.L., and
554 Z.R.M. performed fine-tuning tasks; J.Z., A.B., Z.-Y.L., W.-Y.L., and Z.R.M. analyzed results;
555 J.Z. wrote the manuscript and all authors edited the manuscript.

556 **Competing interests**

557 The authors declare no competing interests.

558 **References**

- 559 1. One Thousand Plant Transcriptomes Initiative. One thousand plant transcriptomes and the
560 phylogenomics of green plants. *Nature* **574**, 679–685 (2019).
- 561 2. Marks, R. A., Hotaling, S., Frandsen, P. B. & VanBuren, R. Representation and participation across
562 20 years of plant genome sequencing. *Nat Plants* **7**, 1571–1578 (2021).
- 563 3. Sun, Y., Shang, L., Zhu, Q.-H., Fan, L. & Guo, L. Twenty years of plant genome sequencing:
564 achievements and challenges. *Trends Plant Sci.* **27**, 391–401 (2022).
- 565 4. Soltis, P. S. & Soltis, D. E. Plant genomes: Markers of evolutionary history and drivers of
566 evolutionary change. *Plants People Planet* **3**, 74–82 (2021).
- 567 5. Provar, N. J. *et al.* Anno genominis XX: 20 years of Arabidopsis genomics. *Plant Cell* **33**, 832–845
568 (2021).
- 569 6. Fu, L.-Y. *et al.* ChIP-Hub provides an integrative platform for exploring plant regulome. *Nat.*
570 *Commun.* **13**, 3413 (2022).
- 571 7. Dudnyk, K., Cai, D., Shi, C., Xu, J. & Zhou, J. Sequence basis of transcription initiation in the
572 human genome. *Science* **384**, eadj0116 (2024).
- 573 8. Jaganathan, K. *et al.* Predicting Splicing from Primary Sequence with Deep Learning. *Cell* **176**, 535–
574 548.e24 (2019).
- 575 9. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range
576 interactions. *Nat. Methods* **18**, 1196–1203 (2021).
- 577 10. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome.
578 *Nature* **489**, 57–74 (2012).
- 579 11. ENCODE Project Consortium *et al.* Expanded encyclopaedias of DNA elements in the human and
580 mouse genomes. *Nature* **583**, 699–710 (2020).
- 581 12. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model.
582 *Science* **379**, 1123–1130 (2023).

- 583 13. Brandes, N., Goldman, G., Wang, C. H., Ye, C. J. & Ntranos, V. Genome-wide prediction of disease
584 variant effects with a deep protein language model. *Nat. Genet.* **55**, 1512–1522 (2023).
- 585 14. Madani, A. *et al.* Large language models generate functional protein sequences across diverse
586 families. *Nat. Biotechnol.* **41**, 1099–1106 (2023).
- 587 15. Ruffolo, J. A. & Madani, A. Designing proteins with language models. *Nat. Biotechnol.* **42**, 200–202
588 (2024).
- 589 16. Cheng, J. *et al.* Accurate proteome-wide missense variant effect prediction with AlphaMissense.
590 *Science* **381**, eadg7492 (2023).
- 591 17. Engelhorn, J. *et al.* Genetic variation at transcription factor binding sites largely explains phenotypic
592 heritability in maize. *bioRxiv* 2023.08.08.551183 (2024) doi:10.1101/2023.08.08.551183.
- 593 18. Gaulton, K. J., Preissl, S. & Ren, B. Interpreting non-coding disease-associated human variants using
594 single-cell epigenomics. *Nat. Rev. Genet.* **24**, 516–534 (2023).
- 595 19. Leeman-Neill, R. J. *et al.* Noncoding mutations cause super-enhancer retargeting resulting in protein
596 synthesis dysregulation during B cell lymphoma progression. *Nat. Genet.* **55**, 2160–2174 (2023).
- 597 20. Novák, P. *et al.* Repeat-sequence turnover shifts fundamentally in species with large genomes. *Nat*
598 *Plants* **6**, 1325–1329 (2020).
- 599 21. Benegas, G., Batra, S. S. & Song, Y. S. DNA language models are powerful predictors of genome-
600 wide variant effects. *Proc. Natl. Acad. Sci. U. S. A.* **120**, e2311219120 (2023).
- 601 22. Zhou, H., Shrikumar, A. & Kundaje, A. Towards a Better Understanding of Reverse-Complement
602 Equivariance for Deep Learning Models in Genomics. in *Proceedings of the 16th Machine Learning*
603 *in Computational Biology meeting* (eds. Knowles, D. A., Mostafavi, S. & Lee, S.-I.) vol. 165 1–33
604 (PMLR, 22--23 Nov 2022).
- 605 23. Mendoza-Revilla, J. *et al.* A Foundational Large Language Model for Edible Plant Genomes.
606 *bioRxiv* 2023.10.24.563624 (2023) doi:10.1101/2023.10.24.563624.
- 607 24. Dalla-Torre, H. *et al.* The Nucleotide Transformer: Building and Evaluating Robust Foundation
608 Models for Human Genomics. *bioRxiv* 2023.01.11.523679 (2023) doi:10.1101/2023.01.11.523679.

- 609 25. Nguyen, E. *et al.* HyenaDNA: Long-Range Genomic Sequence Modeling at Single Nucleotide
610 Resolution. *arXiv [cs.LG]* (2023).
- 611 26. Nguyen, E. *et al.* Sequence modeling and design from molecular to genome scale with Evo. *bioRxiv*
612 2024.02.27.582234 (2024) doi:10.1101/2024.02.27.582234.
- 613 27. Schiff, Y. *et al.* Caduceus: Bi-Directional Equivariant Long-Range DNA Sequence Modeling. *arXiv*
614 *[q-bio.GN]* (2024).
- 615 28. Gu, A. & Dao, T. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. *arXiv*
616 *[cs.LG]* (2023).
- 617 29. Hubisz, M. J., Pollard, K. S. & Siepel, A. PHAST and RPHAST: phylogenetic analysis with
618 space/time models. *Brief. Bioinform.* **12**, 41–51 (2011).
- 619 30. Vaswani, A. *et al.* Attention Is All You Need. *arXiv [cs.CL]* (2017).
- 620 31. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for
621 Dimension Reduction. *arXiv [stat.ML]* (2018).
- 622 32. Ji, Y., Zhou, Z., Liu, H. & Davuluri, R. V. DNABERT: pre-trained Bidirectional Encoder
623 Representations from Transformers model for DNA-language in genome. *Bioinformatics* **37**, 2112–
624 2120 (2021).
- 625 33. Zhou, Z. *et al.* DNABERT-2: Efficient Foundation Model and Benchmark For Multi-Species
626 Genomes. (2023).
- 627 34. Zhang, D. *et al.* DNAGPT: A Generalized Pretrained Tool for Multiple DNA Sequence Analysis
628 Tasks. *bioRxiv* 2023.07.11.548628 (2023) doi:10.1101/2023.07.11.548628.
- 629 35. Quang, D. & Xie, X. DanQ: a hybrid convolutional and recurrent deep neural network for
630 quantifying the function of DNA sequences. *Nucleic Acids Res.* **44**, e107 (2016).
- 631 36. Loos, R. J. F. 15 years of genome-wide association studies and no signs of slowing down. *Nat.*
632 *Commun.* **11**, 1–3 (2020).
- 633 37. Tam, V. *et al.* Benefits and limitations of genome-wide association studies. *Nat. Rev. Genet.* **20**,
634 467–484 (2019).

- 635 38. Ramstein, G. P. & Buckler, E. S. Prediction of evolutionary constraint by genomic annotations
636 improves functional prioritization of genomic variants in maize. *Genome Biol.* **23**, 183 (2022).
- 637 39. Welker, C. A. D. *et al.* Phylogenomics enables biogeographic analysis and a new subtribal
638 classification of Andropogoneae (Poaceae—Panicoideae). *J. Syst. Evol.* **58**, 1003–1030 (2020).
- 639 40. Gossmann, T. I. *et al.* Genome Wide Analyses Reveal Little Evidence for Adaptive Evolution in
640 Many Plant Species. *Mol. Biol. Evol.* **27**, 1822–1832 (2010).
- 641 41. Cao, J. *et al.* Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. *Nat. Genet.*
642 **43**, 956–963 (2011).
- 643 42. Lozano, R. *et al.* Comparative evolutionary genetics of deleterious load in sorghum and maize.
644 *Nature Plants* **7**, 17–24 (2021).
- 645 43. Wu, Y. *et al.* Phylogenomic discovery of deleterious mutations facilitates hybrid potato breeding.
646 *Cell* **186**, 2313–2328.e15 (2023).
- 647 44. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).
- 648 45. Ng, P. C. & Henikoff, S. SIFT: Predicting amino acid changes that affect protein function. *Nucleic*
649 *Acids Res.* **31**, 3812–3814 (2003).
- 650 46. Bukowski, R. *et al.* Construction of the third-generation *Zea mays* haplotype map. *Gigascience* **7**, 1–
651 12 (2018).
- 652 47. 1001 Genomes Consortium. Electronic address: magnus.nordborg@gmi.oeaw.ac.at & 1001
653 Genomes Consortium. 1,135 Genomes Reveal the Global Pattern of Polymorphism in *Arabidopsis*
654 *thaliana*. *Cell* **166**, 481–491 (2016).
- 655 48. Capilla-Perez, L. *et al.* The HEM lines: A new library of homozygous *Arabidopsis thaliana* EMS
656 Mutants and its potential to detect meiotic phenotypes. *Front. Plant Sci.* **9**, 1339 (2018).
- 657 49. Tracy, W. F., Whitt, S. R. & Buckler, E. S. Recurrent mutation and genome evolution: Example of
658 *Sugary1* and the origin of sweet maize. *Crop Sci.* **46**, S–49–S–54 (2006).
- 659 50. Djemel, A. *et al.* Genomic regions affecting fitness of the sweet corn mutants *sugary1*. *J. Agric. Sci.*
660 **151**, 396–406 (2013).

- 661 51. Crow, J. F. 90 years ago: the beginning of hybrid maize. *Genetics* **148**, 923–928 (1998).
- 662 52. Mezmouk, S. & Ross-Ibarra, J. The pattern and distribution of deleterious mutations in maize. *G3* **4**,
663 163–171 (2014).
- 664 53. Lye, Z., Choi, J. Y. & Purugganan, M. D. Deleterious Mutations and the Rare Allele Burden on Rice
665 Gene Expression. *Mol. Biol. Evol.* **39**, (2022).
- 666 54. Song, B. *et al.* AnchorWave: Sensitive alignment of genomes with high sequence diversity,
667 extensive structural polymorphism, and whole-genome duplication. *Proc. Natl. Acad. Sci. U. S. A.*
668 **119**, (2022).
- 669 55. Song, B., Buckler, E. S. & Stitzer, M. C. New whole-genome alignment tools are needed for tapping
670 into plant diversity. *Trends Plant Sci.* **29**, 355–369 (2024).
- 671 56. Smith, S. A. & Brown, J. W. Constructing a broadly inclusive seed plant phylogeny. *Am. J. Bot.* **105**,
672 302–314 (2018).
- 673 57. Ou, S. *et al.* Benchmarking transposable element annotation methods for creation of a streamlined,
674 comprehensive pipeline. *Genome Biol.* **20**, 275 (2019).
- 675 58. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
676 *Bioinformatics* **26**, 841–842 (2010).
- 677 59. Gu, A., Goel, K. & Ré, C. Efficiently Modeling Long Sequences with Structured State Spaces. *arXiv*
678 [*cs.LG*] (2021).
- 679 60. Gu, A. *et al.* Combining Recurrent, Convolutional, and Continuous-time Models with Linear State-
680 Space Layers. *arXiv [cs.LG]* (2021).
- 681 61. Loshchilov, I. & Hutter, F. Decoupled Weight Decay Regularization. *arXiv [cs.LG]* (2017).
- 682 62. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional
683 Transformers for Language Understanding. *arXiv [cs.CL]* (2018).
- 684 63. Cheng, C.-Y. *et al.* Araport11: a complete reannotation of the Arabidopsis thaliana reference
685 genome. *Plant J.* **89**, 789–804 (2017).
- 686 64. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO:

- 687 assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*
688 **31**, 3210–3212 (2015).
- 689 65. Paten, B. *et al.* Cactus: Algorithms for genome multiple sequence alignment. *Genome Res.* **21**,
690 1512–1528 (2011).
- 691 66. Romay, M. C. *et al.* Comprehensive genotyping of the USA national maize inbred seed bank.
692 *Genome Biol.* **14**, R55 (2013).
- 693 67. Khaipho-Burch, M. *et al.* Elucidating the patterns of pleiotropy and its biological relevance in maize.
694 *PLoS Genet.* **19**, e1010664 (2023).
- 695 68. Wrightsman, T., Marand, A. P., Crisp, P. A., Springer, N. M. & Buckler, E. S. Modeling chromatin
696 state from sequence across angiosperms using recurrent convolutional neural networks. *Plant*
697 *Genome* **15**, e20249 (2022).
- 698

699 **Tables**

700 **Table 1.** PlantCaduceus model parameters

Models	# of layers	Hidden size	# of parameters (million)
PlantCaduceus_132	32	1024	225
PlantCaduceus_128	28	768	112
PlantCaduceus_124	24	512	40
PlantCaduceus_120	20	384	20

701

702 **Table 2.** The zero-shot score of deleterious mutations identified in homozygous EMS mutants

Chr	Position	Change	Mutation effect	Phenotype	Zero-shot score	Percentile
2	14297325	G>A	Splice change	Univalent chromosomes	-10.344	Top 1%
3	23443192	C>T	Splice change	Univalent chromosomes	-10.219	Top 1%
3	10277172	C>T	Splice change	Univalent chromosomes	-9.820	Top 1%
4	5820399	C>T	Splice change	Univalent chromosomes	-9.547	Top 1%
3	17827101	G>A	Splice change	Fragmentation	-9.531	Top 1%
3	3248339	C>T	Premature stop	Univalent chromosomes	-9.125	Top 1%
3	17823207	G>A	Splice change	Fragmentation	-9.000	Top 1%
1	1298121	C>T	Splice change	Univalent chromosomes	-8.859	Top 1%
3	17812658	G>A	Splice change	Fragmentation	-8.719	Top 1%
5	1625685	G>A	Splice change	Fragmentation	-8.203	Top 10%
3	3246364	G>A	Splice change	Univalent chromosomes	-7.660	Top 10%
3	3246274	G>A	Splice change	Univalent chromosomes	-7.406	Top 10%
5	23446256	G>A	Premature stop	All univalent chromosomes	-6.203	Top 10%
4	16868745	C>T	Premature stop	Univalent chromosomes	-6.008	Top 10%
3	17824467	G>A	Missense	Fragmentation	-5.570	Top 10%
4	16868001	C>T	Premature stop	Univalent chromosomes	-5.141	Top 50%
3	17807938	G>A	Premature stop	Fragmentation	-4.688	Top 50%
5	26302687	G>A	Premature stop	Univalent chromosomes	-3.805	Top 50%
1	19964116	G>A	Start gained	Univalent chromosomes	0.243	Top 50%

703

704

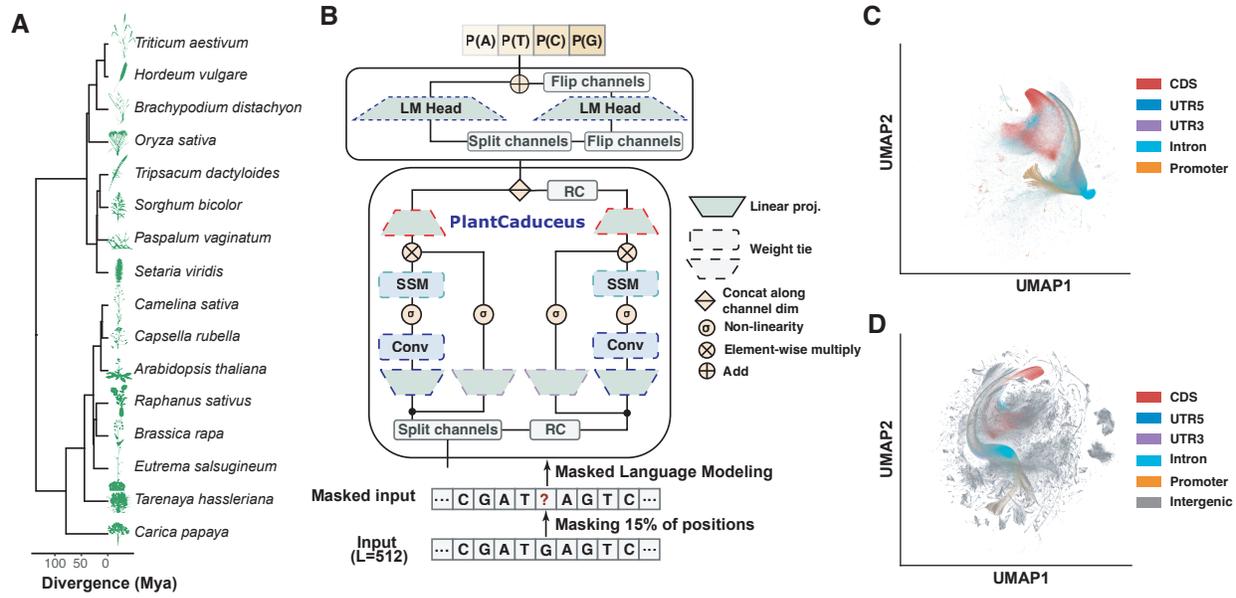


Fig 1. Overview of PlantCaduceus. (A) Phylogenetic tree of 16 Angiosperm species used for pre-training the PlantCaduceus model. (B) The input for PlantCaduceus consists of 512-bp DNA sequences with 15% of positions randomly masked. The pre-training objective is cross-entropy loss on the masked positions. The sequences are processed through the bi-directional Caduceus architecture, which is based on the Mamba sequence operator—a recently proposed structured state space model. Caduceus also contains a reverse complement equivariance inductive bias. (C) UMAP visualization of embeddings from PlantCaduceus (32 layers) averaged over non-overlapping 100-bp windows along the sorghum genome without intergenic regions. (D) The same UMAP visualization as in (C) but with intergenic regions.

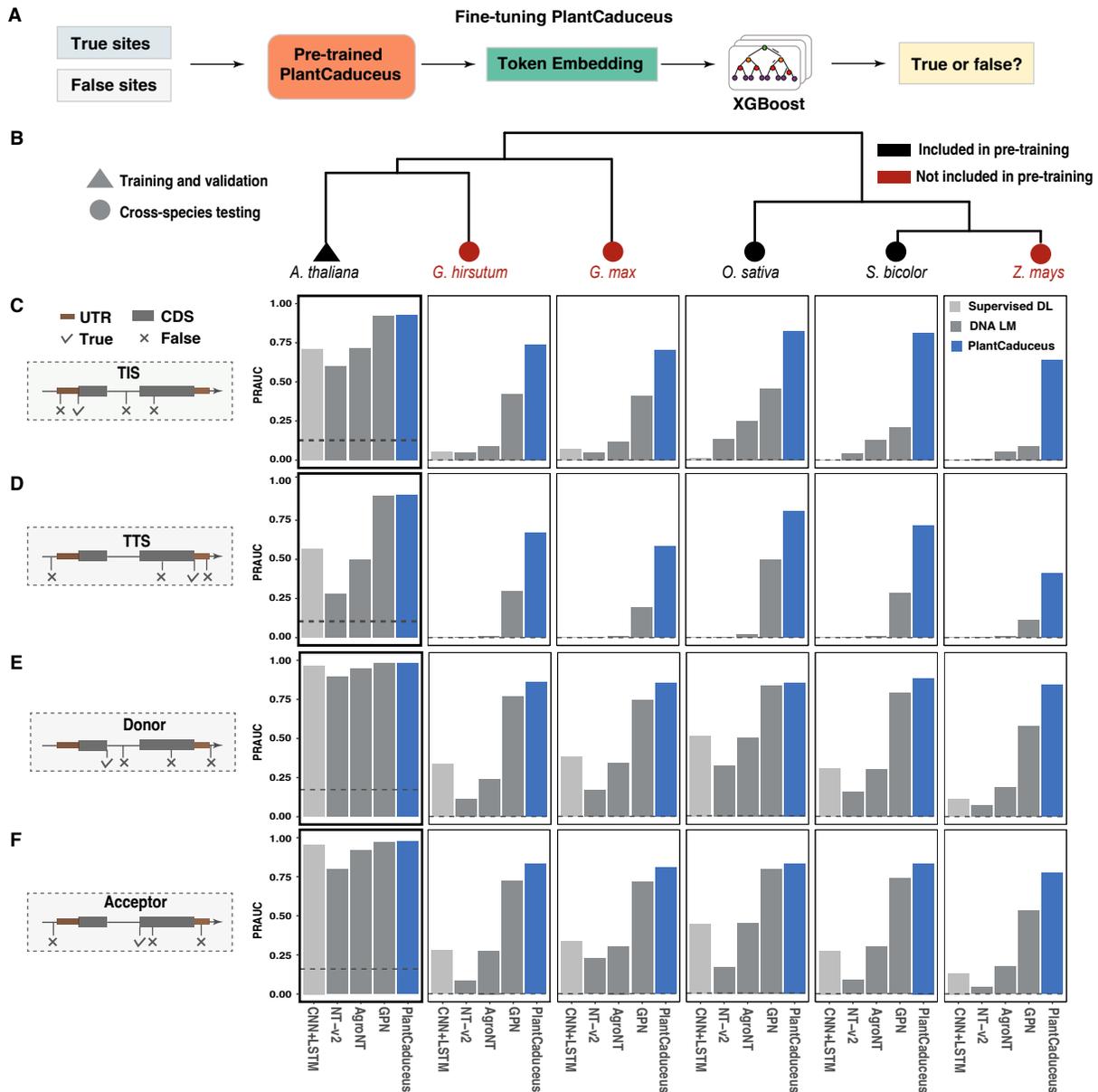


Fig 2. Modeling translation and transcription through fine-tuning PlantCaduceus. (A) Fine-tuning strategy for PlantCaduceus: The weights of the pre-trained PlantCaduceus model are kept frozen during pre-training. The last hidden state of PlantCaduceus is then used as features for the XGBoost model. (B) Phylogenetic tree of species used for training, validation, and testing during the fine-tuning of PlantCaduceus. (C-F) Bar plots displaying the PRAUC scores for six species across four tasks: TIS (C), TTS (D), splice donor (E), and splice acceptor (F). The gene structures on the left illustrate how positive and negative samples are obtained for each classification task. Blue bars represent the PlantCaduceus model with 32 layers. Gray bars denote three DNA language models: NT-v2, AgroNT, and GPN. Light gray bars represent a traditional supervised model, a hybrid of CNN and LSTM. The gray dashed line in each panel indicates the baseline for each dataset, corresponding to the negative sample ratio.

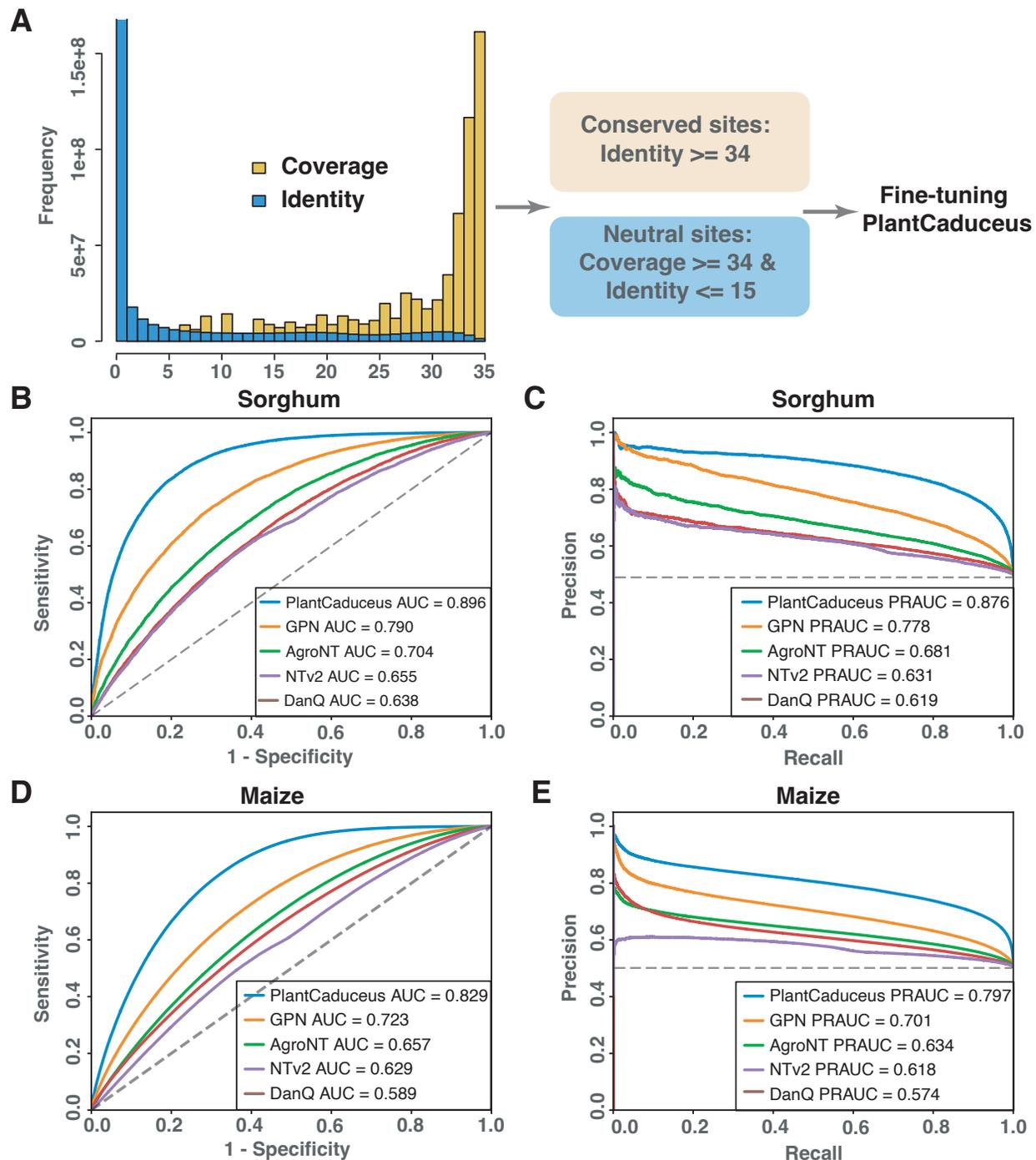


Fig 3. Evolutionary constraint prediction. (A) Illustration of the evolutionary conservation data curation. (B) Receiver operating characteristic (ROC) and (C) precision-recall (PR) curves of different models in sorghum. (D) ROC and (E) PR curves of transferring different models trained in sorghum to unseen maize data.

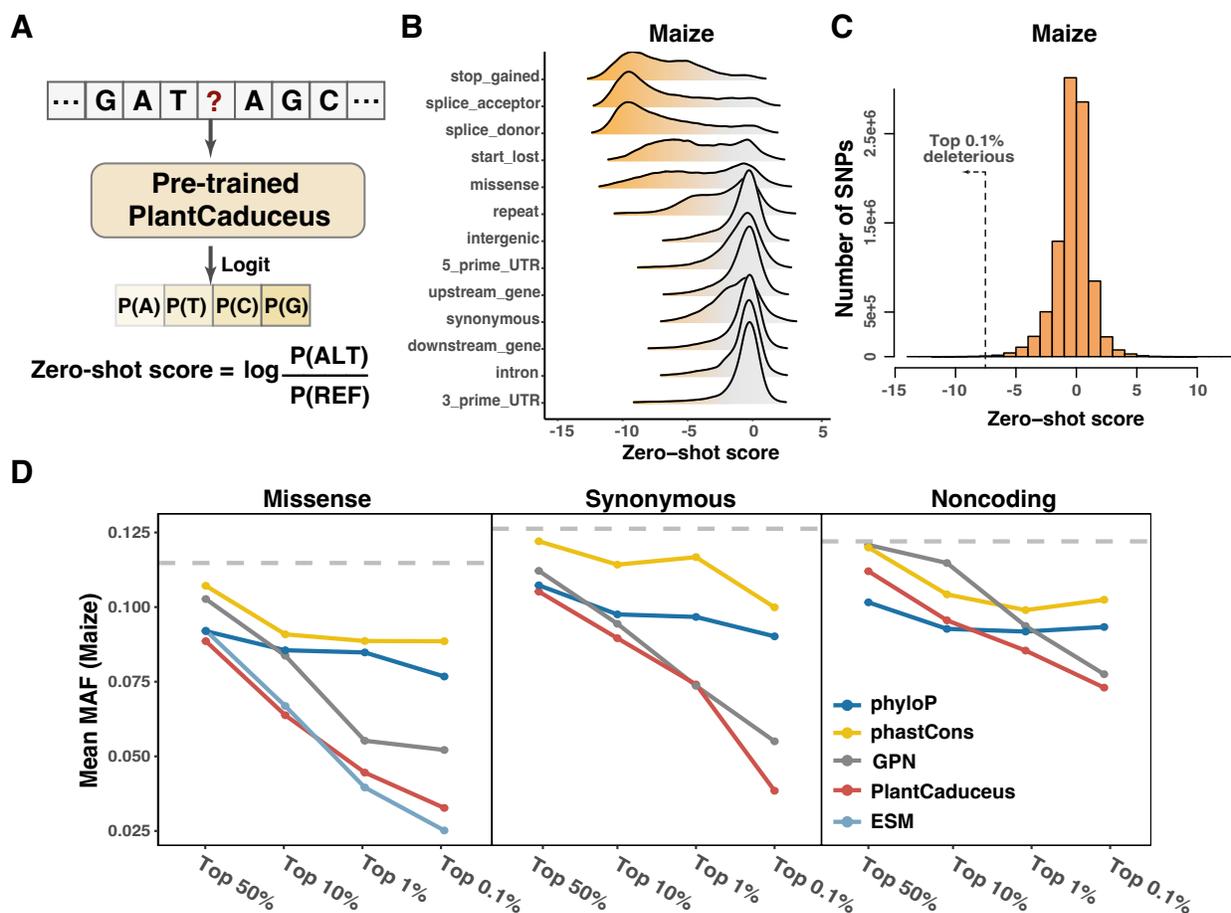


Fig 4. Deleterious mutations identification in maize. (A) The zero-shot strategy of PlantCaduceus for identifying deleterious mutations. (B) The zero-shot score distribution of different types of variants generated by in silico mutagenesis in maize chromosome 8. (C) The zero-shot score distribution of 9.4M SNPs in the maize Hapmap3 population. (D) The MAF of putative deleterious mutations prioritized by different models in maize.

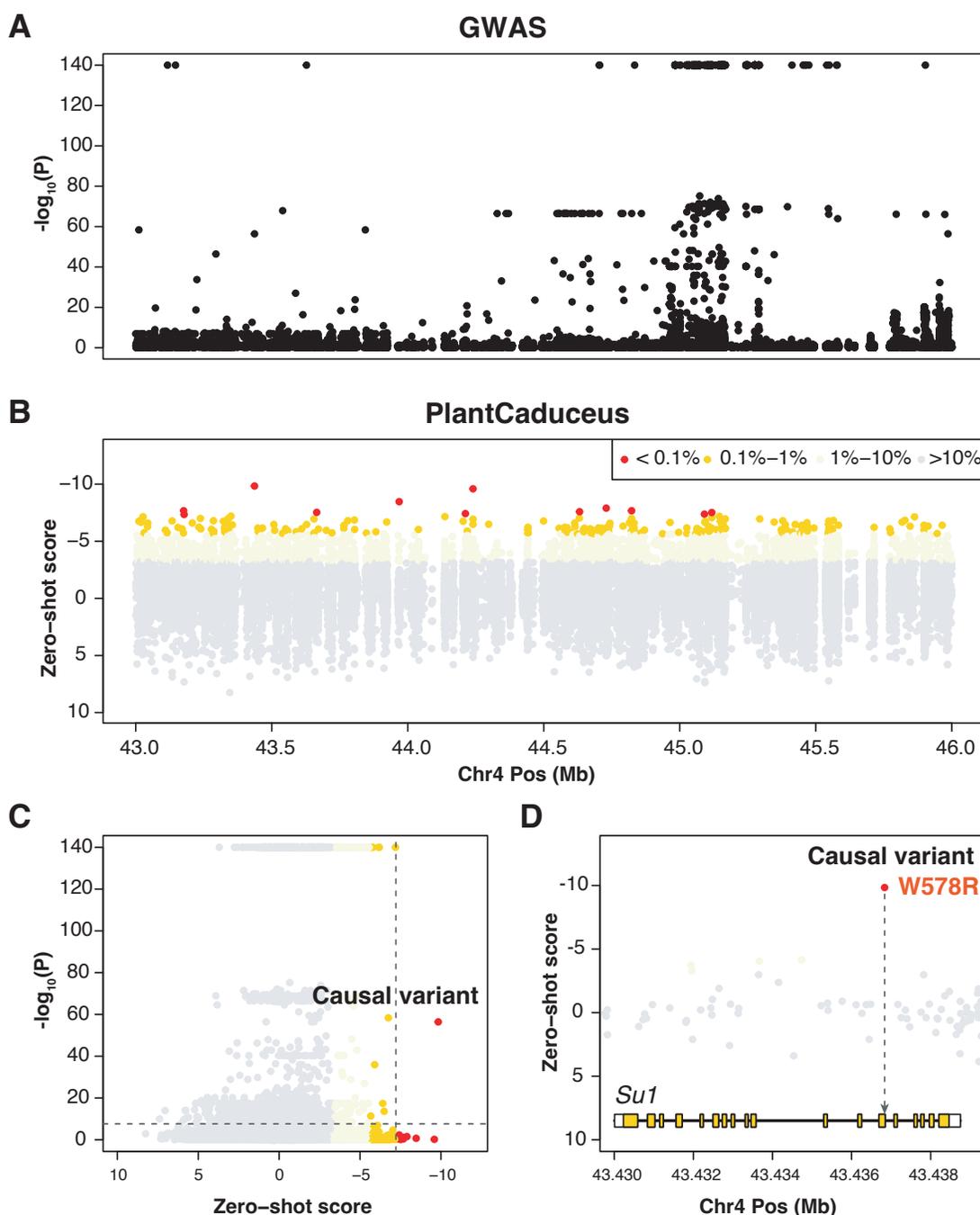


Fig 5. The causal mutation in *Su1* locus. (A) Manhattan plot of the sweet corn trait in the region from 43.0 to 46.0 Mb on chromosome 4. (B) The zero-shot scores of SNPs in 43.0 to 46.0 Mb in chromosome 4, corresponding to the same region as in (A). (C) Scatter plot of zero-shot scores from PlantCaduceus versus $-\log_{10}(P)$ values from GWAS result. The horizontal dashed line indicates the GWAS significance threshold (Bonferroni's threshold: $0.05/N$; $N=2,072,522$), and the vertical dashed line marks the top 0.1% percentile of zero-shot scores. (D) Zoomed-in view of the causal variant region and the *Su1* gene structure.