

Supplemental Methods

1	Supplemental Methods	
2	Generation and merging of the sSNVs.....	2
3	Splice site consensus	2
4	Exonic splice regulatory elements.....	2
5	sSNVs from vertebrate species.....	3
6	Economic traits-related sSNV from cultivated plants and domesticated animals	4
7	Literature resource in SynMall.....	5
8	PMC/PubMed query statement.....	5
9	Criterion standard for human sSNVs.....	5
10	Fields and description of structured information extracted from the literature	5
11	Batch query performance in SynMall.....	6
12	Curation of datasets	6
13	Performance evaluation.....	7
14	Reference.....	8
15		
16		
17		

18 **Generation and merging of the sSNVs**

19 To obtain all possible point mutations in coding regions, we first retrieved all hg38 human
20 protein-coding transcripts from BioMart (Kinsella et al. 2011) (Version: Ensembl 110). After
21 removing transcripts containing unknown bases or lacking start/stop codons in their coding
22 sequences, 84,067 transcripts remained. Next, based on their coding sequences, we simulated
23 each base mutating into the other three bases, resulting in all possible point mutations in the
24 coding region. In the second step, we used the Variant Effect Predictor (McLaren et al. 2016)
25 (Version: Ensembl 110) tool to filter out mutations possessing synonymous consequences.
26 Additionally, to provide mutation information for the hg19 reference, we performed genome
27 coordinate conversion using the LiftOver (Kuhn et al. 2013) tool. Due to varying quality control
28 procedures, sequence variation annotation strategies, and reference genome versions, in case of
29 omitted data, we also merged sSNVs from CADD (Schubach et al. 2024), FAVOR (Zhou et al.
30 2023), and synVep (Zeng et al. 2021), which all contain synonymous mutations across the
31 human genome.

32 **Splice site consensus**

33 In higher eukaryotes, precise splicing is regulated by three weakly conserved cis-elements,
34 5' and 3' splice sites, and the branch site. According to RegSNPs-splicing (Zhang et al. 2017),
35 if an sSNV falls at the +1, +2, or +3 position of the 5' splice site or the -1 position of the 3'
36 splice site, we classified it as a variant on splice site consensus (VSS). Otherwise, we classified
37 it as variants in internal exons (VIE).

38 **Exonic splice regulatory elements**

39 Referring to sSNVs pathogenic prediction tools such as SliVA (Buske et al. 2013), DDIG-

40 SN (Livingstone et al. 2017), and regSNPs-splicing (Zhang et al. 2017), exonic splicing
41 regulatory (ESR) sequences are considered to be important features when constructing models.
42 The potential of sSNVs to result in a gain/loss of an ESR may be correlated with its
43 pathogenicity. Hence, a comprehensive ESR motifs set is curated from RESCUE-ESE
44 (Fairbrother et al. 2002), FAS-HEX3 (Wang et al. 2004), SpliceAID, RegRNA2 (These two
45 datasets are retrieved from SynMICdb (Sharma et al. 2019)), Composite-ESR (Ke et al. 2008),
46 NI-ESR (Stadler et al. 2006), and Ast-ESR (Goren et al. 2006). After merging and removing
47 duplicates, we have 701 exonic splicing silencer motifs, 1,048 exonic splicing enhancer motifs,
48 and 285 ESR motifs remaining. The detailed table is available in Supplemental Table S5.

49 **sSNVs from vertebrate species**

50 First, we targeted the species included in the UCSC 100-way vertebrate multiple sequence
51 alignment. Since Ensembl focuses on vertebrate genomes, and Ensembl Variation performs
52 quality control on mutations while providing evidence status and functional consequence
53 annotations, its data quality and reliability are relatively high. Therefore, we obtained sSNV
54 information for all 17 non-human vertebrate species from Ensembl Variation (Hunt et al. 2018)
55 (<http://www.ensembl.org/info/genome/variation/index.html>, downloaded on 2024-05-09).
56 Additionally, the European Variation Archive (EVA) (Cezard et al. 2022)
57 (<https://www.ebi.ac.uk/eva/>, downloaded on 2025-01-17), as the most comprehensive platform
58 for genetic mutations across all species, offers extensive information on non-human species. To
59 ensure that the sequence consequences of mutations could be annotated, we filtered vertebrate
60 species mutations that could be annotated using the VEP with available cache files
61 (https://ftp.ensembl.org/pub/release-113/variation/indexed_vep_cache/) and added sSNV

62 information for 7 additional vertebrate species. Since most annotation resources are designed
63 primarily for humans, we aim to map non-human sSNVs to human reference genomes to enable
64 shared annotation. However, LiftOver may introduce artifacts when mapping genomic
65 coordinates across species. We referred to the method used in PrimateAI (Sundaram et al. 2018).
66 Based on the multiple sequence alignment (MSA,
67 <https://hgdownload.soe.ucsc.edu/goldenPath/hg38/multiz100way/>

68), we mapped non-human mutations onto the human genome. Additionally, common variants
69 in other primates are largely benign in humans (Gao et al. 2023; Cheng et al. 2023). Therefore,
70 we selected mutations from five primate species, including those from the Great Ape
71 (<https://eichlerlab.gs.washington.edu/greatape/data>) and Han
72 (https://figshare.com/articles/dataset/Han_et.al_Data_tsv_gz/7855850). Since these mutation
73 datasets were already mapped to hg18 or hg19, we used the LiftOver method to map them to
74 hg38. Detailed statistics can be found in Supplemental Table S4.

75 **Economic traits-related sSNV from cultivated plants and domesticated animals**

76 To understand the relationship between sSNVs and economic traits in major livestock and
77 crops, we integrated 11,963 GWAS associations of sSNVs from 17 crops and 4 animals in the
78 GWAS Atlas (Liu et al. 2023) (<https://ngdc.cncb.ac.cn/gwas/>, downloaded on 2025-01-06) and
79 CropGS-Hub (Chen et al. 2024) (<https://iagr.genomics.cn/CropGS/#>, downloaded on 2024-11-
80 25) databases. These databases are built based on literature-extracted GWAS information. In
81 CropGS-Hub, the literature related to Rice and Sorghum is already included in the GWAS Atlas,
82 so we focused on the literature not covered by it. Detailed statistical information about the data
83 can be found in Supplemental Table S4.

84 **Literature resource in SynMall**

85 *PMC/PubMed query statement*

86 SynMall automatically collects literature from PMC and PubMed using the following
87 search query: ("synonymous"[Title/Abstract]) AND ("mutation"[Title/Abstract] OR
88 "variation"[Title/Abstract] OR "variant"[Title/Abstract] OR "mutant"[Title/Abstract]) NOT
89 ("non-synonymous") NOT ("nonsynonymous").

90 *Criterion standard for human sSNVs*

91 For sSNVs in humans, we aim to extract evidence-supporting associations, categorizing
92 them as either benign or pathogenic. Therefore, we referred to the criteria of the ACMG
93 (Richards 2015).

94 *Fields and description of structured information extracted from the literature*

95 In total, 21 fields across three domains are considered when curating each paper. The
96 literature-central domain provides basic publication details and key supporting evidence
97 sentences. The variant-central section offers detailed information about the variant, including
98 allele change, genomic position, strand, coding sequence position, reference single-nucleotide
99 polymorphism ID, and codon change, as well as additional information related to the gene and
100 species. The phenotype-central part contains manually annotated data inferring the phenotypic
101 effects of sSNVs, where the Mechanism field describes how the sSNVs induce the disorder
102 (e.g., through splicing regulation, mRNA structure stability, protein synthesis, etc.). The Trait
103 field is designed for non-human species to capture traits associated with sSNVs, while the Trait
104 Impact field describes the effect of the mutation on the trait, such as promoting or inhibiting, if
105 applicable.

106 **Batch query performance in SynMall**

107 We evaluated the response speed of batch retrieval in SynMall using Apache JMeter (with
108 10 concurrent threads simulating multiple users), as summarized in Supplemental Table S3.
109 When querying 1,000 records, the average response times were 7.35 s for Genomic Coordinates,
110 0.99 s for Gene Names, and 2.99 s for RS IDs. Please note that the first request or prolonged
111 inactivity may trigger reinitialization of the database connection pool, leading to slower
112 response times than those shown in the table. Currently, batch queries support up to 1,000
113 records per request, as larger queries may result in timeout errors. For datasets exceeding this
114 limit (1,000–50,000 records), users are advised to use the Annotation module.

115 **Curation of datasets**

116 We compile a benchmark dataset for machine learning using sSNVs curated from multiple
117 external databases and literature. The dataset includes a balanced training set of 2,362 sSNVs
118 and a balanced test set of 238 sSNVs. First, we retrieve initial data from ClinVar (downloaded
119 April 2025) (Landrum et al. 2020), HGMD (Professional 2023.3) (Stenson et al. 2020) , dbDSM
120 (Wen et al. 2016), and manually reviewed sSNVs from SynMall. In ClinVar, we select variants
121 labeled as "Benign", "Likely Benign", "Likely Benign/Benign", "Likely
122 Pathogenic/Pathogenic", "Likely Pathogenic", and "Pathogenic". Only records with a review
123 status of "criteria provided, multiple submitters, no conflicts", "criteria provided, single
124 submitter", or "reviewed by expert panel" are included. In HGMD, we include DM-classified
125 pathogenic synonymous mutations, and in dbDSM, we select variants from manually curated
126 sources. We remove any variants that appear in both benign and pathogenic categories. To
127 evaluate VEP performance on rare variants, we filter out common variants with AF > 1e-3,

128 retaining only rare variants. We control sequence similarity using CD-HIT (Fu et al. 2012),
129 ensuring that protein sequences in the training and test sets share less than 40% identity. Finally,
130 we apply a "close-by" strategy (Cheng et al. 2020) to balance pathogenic and benign samples.

131 For each minority pathogenic sample, we select a benign sample with the closest genomic
132 position, creating a dataset with balanced positive and negative samples. The full dataset is
133 available for download on the "Download" page.

134 **Performance evaluation**

135 To evaluate the performance of VEP tools on the synonymous variant test set, we use the
136 Area Under the Receiver Operating Characteristic Curve (AUC) and the Area Under the
137 Precision-Recall Curve (AUPR). The ROC curve plots the true positive rate (TPR) against the
138 false positive rate (FPR) across different classification thresholds, while the Precision-Recall
139 (PR) curve plots Precision against Recall. Both metrics provide threshold-independent
140 measures commonly used for assessing binary classification performance.

141 Notably, some tools produce missing values and fail to provide predictions for certain
142 variants in the independent test set. To address this, we apply both the "subset" and "pairwise"
143 evaluation strategies. Specifically, the subset approach extracts the portion of the test set for
144 which all VEP tools provide prediction scores, and evaluates the performance of all tools
145 simultaneously on this subset. In contrast, the pairwise approach compares synScore against
146 each target VEP individually by selecting the subset of variants without missing values for that
147 tool, and evaluates their relative performance within this set.

148

149 **Reference**

150 Buske OJ, Manickaraj A, Mital S, Ray PN, Brudno M. 2013. Identification of deleterious
151 synonymous variants in human genomes. *Bioinformatics* **29**: 1843–50.
152 doi:10.1093/bioinformatics/btt308

153 Cezard T, Cunningham F, Hunt SE, Koylass B, Kumar N, Saunders G, Shen A, Silva AF, Tsukanov
154 K, Venkataraman S, et al. 2022. The European Variation Archive: a FAIR resource of
155 genomic variation for all species. *Nucleic Acids Res* **50**: D1216–D1220.
156 doi:10.1093/nar/gkab960

157 Chen J, Tan C, Zhu M, Zhang C, Wang Z, Ni X, Liu Y, Wei T, Wei X, Fang X, et al. 2024. CropGS-
158 Hub: a comprehensive database of genotype and phenotype resources for genomic
159 prediction in major crops. *Nucleic Acids Res* **52**: D1519–D1529. doi:10.1093/nar/gkad1062

160 Cheng J, Novati G, Pan J, Bycroft C, Žemgulytė A, Applebaum T, Pritzel A, Wong LH, Zielinski M,
161 Sargeant T, et al. 2023. Accurate proteome-wide missense variant effect prediction with
162 AlphaMissense. *Science* **381**: eadg7492. doi:10.1126/science.adg7492

163 Cheng N, Li M, Zhao L, Zhang B, Yang Y, Zheng CH, Xia J. 2020. Comparison and integration of
164 computational methods for deleterious synonymous mutation prediction. *Brief Bioinform*
165 **21**: 970–981. doi:10.1093/bib/bbz047

166 Fairbrother WG, Yeh R-F, Sharp PA, Burge CB. 2002. Predictive Identification of Exonic Splicing
167 Enhancers in Human Genes. *Science* **297**: 1007–1013. doi:10.1126/science.1073774

168 Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: accelerated for clustering the next-generation
169 sequencing data. *Bioinformatics* **28**: 3150–3152. doi:10.1093/bioinformatics/bts565

170 Gao H, Hamp T, Ede J, Schraiber JG, McRae J, Singer-Berk M, Yang Y, Dietrich ASD, Fiziev PP,
171 Kuderna LFK, et al. 2023. The landscape of tolerated genetic variation in humans and
172 primates. *Science* **380**: eabn8153. doi:10.1126/science.abn8197

173 Goren A, Ram O, Amit M, Keren H, Lev-Maor G, Vig I, Pupko T, Ast G. 2006. Comparative
174 Analysis Identifies Exonic Splicing Regulatory Sequences—The Complex Definition of
175 Enhancers and Silencers. *Mol Cell* **22**: 769–781. doi:10.1016/j.molcel.2006.05.008

176 Hunt SE, McLaren W, Gil L, Thormann A, Schuilenburg H, Sheppard D, Parton A, Armean IM,
177 Trevanion SJ, Flück P, et al. 2018. Ensembl variation resources. *Database (Oxford)* **2018**:
178 bay119. doi:10.1093/database/bay119

179 Ke S, Zhang XH-F, Chasin LA. 2008. Positive selection acting on splicing motifs reflects
180 compensatory evolution. *Genome Res* **18**: 533–543. doi:10.1101/gr.070268.107

181 Kinsella RJ, Kähäri A, Haider S, Zamora J, Proctor G, Spudich G, Almeida-King J, Staines D,
182 Derwent P, Kerhornou A, et al. 2011. Ensembl BioMarts: a hub for data retrieval across
183 taxonomic space. *Database (Oxford)* **2011**: bar030. doi:10.1093/database/bar030

184 Kuhn RM, Haussler D, Kent WJ. 2013. The UCSC genome browser and associated tools. *Brief
185 Bioinform* **14**: 144–161. doi:10.1093/bib/bbs038

186 Landrum MJ, Chitipiralla S, Brown GR, Chen C, Gu B, Hart J, Hoffman D, Jang W, Kaur K, Liu C,
187 et al. 2020. ClinVar: improvements to accessing data. *Nucleic Acids Res* **48**: D835–D844.
188 doi:10.1093/nar/gkz972

189 Liu X, Tian D, Li C, Tang B, Wang Z, Zhang R, Pan Y, Wang Y, Zou D, Zhang Z, et al. 2023. GWAS
190 Atlas: an updated knowledgebase integrating more curated associations in plants and
191 animals. *Nucleic Acids Res* **51**: D969–D976. doi:10.1093/nar/gkac924

192 Livingstone M, Folkman L, Yang Y, Zhang P, Mort M, Cooper DN, Liu Y, Stantic B, Zhou Y. 2017.
193 Investigating DNA-, RNA-, and protein-based features as a means to discriminate
194 pathogenic synonymous variants. *Hum Mutat* **38**: 1336–1347. doi:10.1002/humu.23283

195 McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, Flliceck P, Cunningham F. 2016.
196 The Ensembl Variant Effect Predictor. *Genome Biol* **17**: 122. doi:10.1186/s13059-016-
197 0974-4

198 Richards S. 2015. Standards and guidelines for the interpretation of sequence variants: a joint
199 consensus recommendation of the American College of Medical Genetics and Genomics
200 and the Association for Molecular Pathology. *Genet Med* **17**. doi:10.1038/gim.2015.30

201 Schubach M, Maass T, Nazaretyan L, Röner S, Kircher M. 2024. CADD v1.7: using protein
202 language models, regulatory CNNs and other nucleotide-level scores to improve genome-
203 wide variant predictions. *Nucleic Acids Res* **52**: D1143–D1154. doi:10.1093/nar/gkad989

204 Sharma Y, Miladi M, Dukare S, Boulay K, Caudron-Herger M, Groß M, Backofen R, Diederichs S.
205 2019. A pan-cancer analysis of synonymous mutations. *Nat Commun* **10**: 2569.
206 doi:10.1038/s41467-019-10489-2

207 Stadler MB, Shomron N, Yeo GW, Schneider A, Xiao X, Burge CB. 2006. Inference of Splicing
208 Regulatory Activities by Sequence Neighborhood Analysis. *PLoS Genet* **2**: e191.
209 doi:10.1371/journal.pgen.0020191

210 Stenson PD, Mort M, Ball EV, Chapman M, Evans K, Azevedo L, Hayden M, Heywood S, Millar
211 DS, Phillips AD, et al. 2020. The Human Gene Mutation Database (HGMD®): optimizing
212 its use in a clinical diagnostic or research setting. *Hum Genet* **139**: 1197–1207.
213 doi:10.1007/s00439-020-02199-3

214 Sundaram L, Gao H, Padigepati SR, McRae JF, Li Y, Kosmicki JA, Fritzilas N, Hakenberg J, Dutta
215 A, Shon J, et al. 2018. Predicting the clinical impact of human mutation with deep neural
216 networks. *Nat Genet* **50**: 1161–1170. doi:10.1038/s41588-018-0167-z

217 Wang Z, Rolish ME, Yeo G, Tung V, Mawson M, Burge CB. 2004. Systematic Identification and
218 Analysis of Exonic Splicing Silencers. *Cell* **119**: 831–845. doi:10.1016/j.cell.2004.11.010

219 Wen P, Xiao P, Xia J. 2016. dbDSM: a manually curated database for deleterious synonymous
220 mutations. *Bioinformatics* **32**: 1914–1916. doi:10.1093/bioinformatics/btw086

221 Zeng Z, Aptekmann AA, Bromberg Y. 2021. Decoding the effects of synonymous variants. *Nucleic
222 Acids Res* **49**: 12673–12691. doi:10.1093/nar/gkab1159

223 Zhang X, Li M, Lin H, Rao X, Feng W, Yang Y, Mort M, Cooper DN, Wang Y, Wang Y, et al. 2017.
224 regSNPs-splicing: a tool for prioritizing synonymous single-nucleotide substitution. *Hum
225 Genet* **136**: 1279–1289. doi:10.1007/s00439-017-1783-x

226 Zhou H, Arapoglou T, Li X, Li Z, Zheng X, Moore J, Asok A, Kumar S, Blue EE, Buyske S, et al.
227 2023. FAVOR: functional annotation of variants online resource and annotator for variation
228 across the human genome. *Nucleic Acids Res* **51**: D1300–D1311. doi:10.1093/nar/gkac966

229