

# Revisiting Adaptive Introgression at the HLA Genes in Lithuanian Genomes with Machine Learning

Josef Hackl<sup>1</sup>, Xin Huang<sup>1,2,\*</sup>

<sup>1</sup> Department of Evolutionary Anthropology, University of Vienna, Vienna, Austria

<sup>2</sup> Human Evolution and Archaeological Sciences (HEAS), University of Vienna, Vienna, Austria

\* Corresponding author: [xin.huang@univie.ac.at](mailto:xin.huang@univie.ac.at)

**Keywords:** Adaptive introgression; Balancing selection; HLA; Population genetics; Machine learning

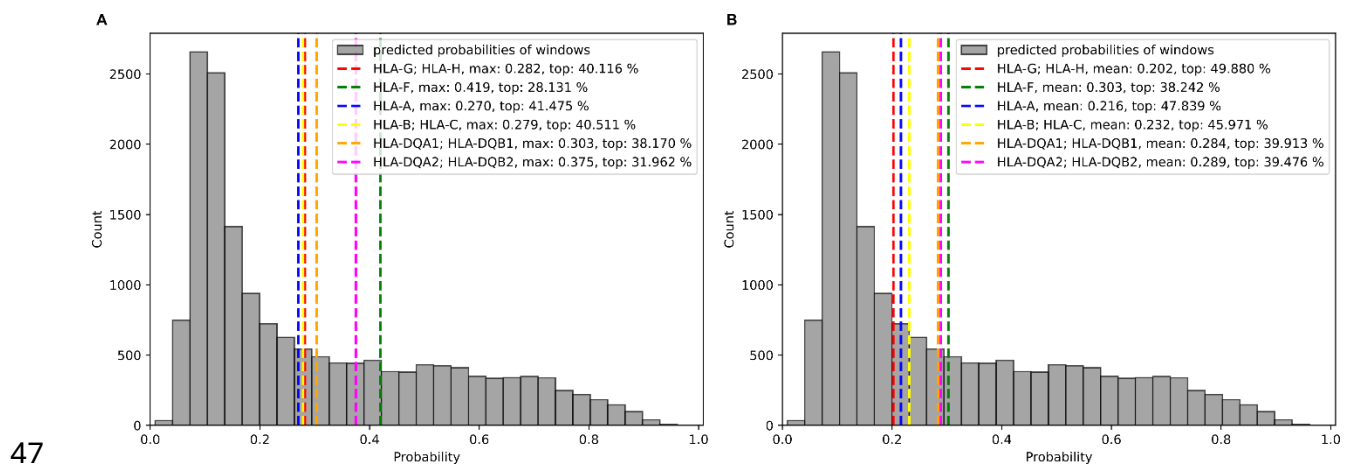
Dear Editor,

We are writing to discuss the article titled ‘Disentangling Archaic Introgression and Genomic Signatures of Selection at Human Immunity Genes,’ published by Urniyete et al (2023). This study employed an *ad-hoc* approach, first applying the machine learning tool, ArchIE (Durvasula and Sankararaman 2019), to detect introgression candidates, followed by the use of the iHS statistic (Voight et al. 2006) to identify candidates under positive selection. According to the authors, the *HLA-C* gene displays both introgression and positive selection signals, suggesting it as a candidate for adaptive introgression in Lithuanians.

However, this approach is problematic due to the varying effectiveness of the methods employed (Zhang et al. 2023) and the confounding effects of introgression on methods used to detect selection (Racimo et al. 2015). More specifically, adaptive introgression can be confounded by balancing selection (Fijarczyk and Babik 2015), and the human leukocyte antigen (HLA) genes are well known examples for long-term balancing selection (Andrés et al. 2009; Gelabert et al. 2024). Considering this, we reanalyzed the Lithuanian genomic data using a recently developed machine learning approach, MaLAdapt (Zhang et al. 2023), which is specifically designed to detect adaptive introgression through supervised learning. Our results suggest that the HLA genes are not candidates for adaptive introgression.

We downloaded the Lithuanian genomes from Urniyete et al. (2023), as well as the chromosome 6 variants of the Altai Neanderthal (Prüfer et al. 2014) from <http://cdna.eva.mpg.de/neandertal/Vindija/> and the chromosome 6 variants of modern humans identified by the 1000 Genomes Project (The 1000 Genomes Project Consortium 2015) from <https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>. Since the HLA genes are located on human chromosome 6, our analysis focused exclusively on this

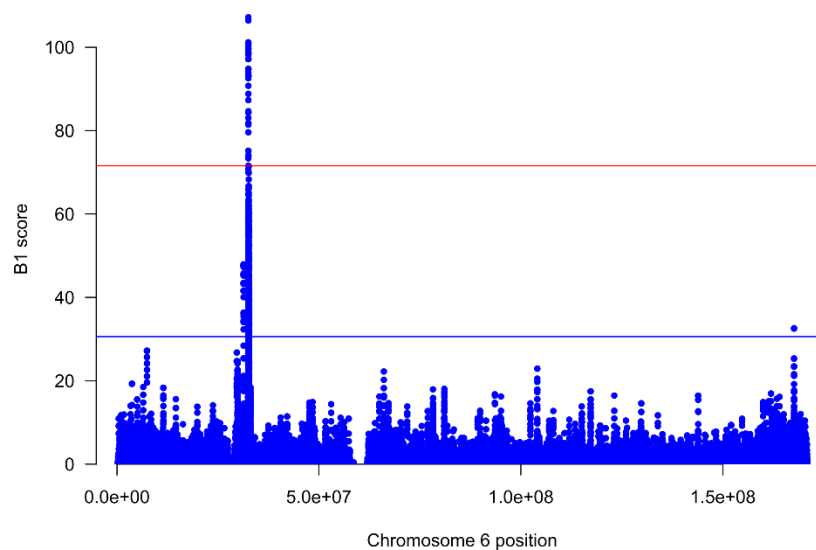
31 chromosome. We used BEAGLE 5.4 (version: 06Aug24.a91) to phase and impute missing genotypes in  
 32 these genomes, utilizing all European populations (CEU, FIN, GBR, IBS, TSI) from the 1000 Genomes  
 33 Project as the reference panel. In BEAGLE 5.4 (Browning et al. 2018; Browning et al. 2021), we set the  
 34 argument  $ne = 10,000$  for effective population size, with other parameters left as default. Next, we  
 35 extracted biallelic single nucleotide polymorphisms (SNPs) from the imputed data and merged them with  
 36 biallelic SNPs from the YRI population in the 1000 Genomes Project, as the YRI population is considered  
 37 a reference population that did not experience introgression (Huang et al. 2022). We applied MaLAdapt to  
 38 extract various input features using sliding windows of 50,000 base pairs with a step size of 10,000 base  
 39 pairs from this merged dataset. Finally, we used the pretrained model provided by MaLAdapt to predict  
 40 the probability of adaptive introgression for each 50 kb window from these input features. We estimated  
 41 the probability of adaptive introgression for each HLA gene by either averaging or taking the maximum  
 42 of the probabilities from the 50 kb windows that overlap with the gene. Our results (Figure 1) show that  
 43 the probabilities of adaptive introgression for the HLA genes fall within the middle range of the empirical  
 44 distribution, which indicates that they are not outliers. Since adaptive introgression is rare and typically  
 45 identified through outliers, these findings suggest that the HLA genes are not candidates for adaptive  
 46 introgression.



48 **Figure 1. Distribution of adaptive introgression probabilities on chromosome 6 in Lithuanian**  
 49 **genomes predicted by MaLAdapt using the pretrained model MaLAdapt\_25\_-sweep-all\_model. A,**  
 50 **the adaptive introgression probabilities for the HLA genes by taking the maximum of the probabilities**  
 51 **from the 50 kb windows that overlap with the genes; B, the adaptive introgression probabilities for the**  
 52 **HLA genes by averaging the probabilities from the 50 kb windows that overlap with the genes. The grey**  
 53 **bars represent the distribution of adaptive introgression probabilities for all 50 kb windows across**  
 54 **chromosome 6 in the Lithuanian genomes. The dashed lines indicate the adaptive introgression**  
 55 **probabilities for the HLA genes. We downloaded the pretrained model from**

56 [https://drive.google.com/drive/folders/10r8e5WbhcgAljC0DVMle4saVYODRgFCO?usp=share\\_link](https://drive.google.com/drive/folders/10r8e5WbhcgAljC0DVMle4saVYODRgFCO?usp=share_link) on  
57 June 5, 2024, and used the ranges of the HLA genes as outlined in Table 1 of Urnikyte et al. (2023). The  
58 hg19 coordinates are as follows: *HLA-G*; *HLA-H*: 29798610–29897944; *HLA-F*: 29698426–29746527;  
59 *HLA-A*: 29898105–29947740; *HLA-B*; *HLA-C*: 31198205–31348022; *HLA-DQA1*; *HLA-DQB1*:  
60 32598042–32644388; *HLA-DQA2*; *HLA-DQB2*: 32698044–32748039; *HLA-DOB*: 32748045–32797488.

61 We also applied the B1 statistic from BetaScan (Siewert and Voight 2017) to identify candidates for long-  
62 term balancing selection on chromosome 6 in the Lithuanian genomes. For the scan, only SNPs with  
63 minor allele frequencies greater than 0.05 in the imputed data were used. SNPs located in regions defined  
64 by the RepeatMasker table, simple repeats table, and segmental duplication table from the UCSC Table  
65 Browser (hg19 coordinates, last accessed in October 2024) were removed. Additionally, SNPs with *p*-  
66 values less than  $10^{-3}$  from exact Hardy-Weinberg equilibrium tests in each population, performed using  
67 PLINK 1.9 (Chang et al. 2015), were excluded. As per Siewert and Voight (2017), only SNPs with folded  
68 allele frequencies greater than 0.15 were used as cores for calculating the B1 scores. All other parameters  
69 in BetaScan were kept at their default values. Our results show a peak in B1 scores within the HLA genes,  
70 particularly in *HLA-B*; *HLA-C* (highest B1 score of 36.350388), *HLA-DQA1*; *HLA-DQB1* (highest B1  
71 score of 71.446036), and *HLA-DQA2*; *HLA-DQB2* (highest B1 score of 43.007245), suggesting they are  
72 candidates for long-term balancing selection. This is consistent with previous studies (DeGiorgio et al.  
73 2014; Bitarello et al. 2018) and was also noted by Urnikyte et al. (2023).



74  
75 **Figure 2. Manhattan plot of B1 scores on chromosome 6 in Lithuanian genomes.** The red horizontal  
76 line (B1 = 71.572528) represents the top 0.05%, and the blue horizontal line (B1 = 30.589559) represents  
77 the top 1%. This plot was created using the qqman package (Turner 2018).

78 It remains controversial whether the HLA genes are under balancing selection or if they experienced  
79 adaptive introgression from archaic humans (Ding et al. 2014; Yasukochi and Ohashi 2017). Recently  
80 developed machine learning-based methods for detecting adaptive introgression, such as genomatnn and  
81 MaLAdapt, have not reported signals of adaptive introgression in the HLA genes using populations from  
82 the 1000 Genomes Project (Gower et al. 2021; Zhang et al. 2023). However, the HLA regions have been  
83 consistently identified as candidates for balancing selection in both modern and ancient human  
84 populations across various recent studies employing different approaches (Siewert and Voight 2017;  
85 Bitarello et al. 2018; Gelabert et al. 2024). Moreover, a recent method, based on the ancestral  
86 recombination graph, strongly supports balancing selection in the HLA regions with trans-species  
87 polymorphism (Deng et al. 2024). Considering that polymorphisms maintained by balancing selection are  
88 typically shared across populations or species (Hedrick 2007; Bitarello et al. 2023), it is likely that the  
89 HLA genes in Lithuanian genomes are maintained by balancing selection shared among various human  
90 populations. Although a study by Abi-Rached et al. (2011) suggested that the *HLA-B* locus was under  
91 adaptive introgression, they used a simulator that assumed neutrality at this locus, even though the  
92 classical class I *HLA* loci are well-known examples of balancing selection (Yasukochi and Ohashi 2017).

93 We would like to point out a similar issue with the machine learning tool, ArchIE, used by Urnikyte et al.  
94 (2023), which relies on the ms simulator based on the Wright–Fisher neutral model (Hudson 2002) to  
95 generate training data. Since ms cannot simulate data under natural selection, this raises concerns about  
96 how ArchIE performs when analyzing data that includes natural selection. Simulation misspecification  
97 can impact the performance of supervised learning tools (Mo and Siepel 2023). Therefore, it is critical to  
98 document the specific demographic model used for generating the simulated data, which was not reported  
99 in Urnikyte et al. (2023). If Urnikyte et al. (2023) used the demographic model hard-coded in ArchIE, an  
100 additional issue arises: the ArchIE code trains on data simulated from a four-population model, whereas a  
101 three-population model was reported (Huang 2024). This discrepancy may also affect the performance of  
102 ArchIE. Furthermore, a recent study (Ray et al. 2024) highlighted the importance of balancing the training  
103 data to achieve results similar to those originally reported by ArchIE, as introgression is rare and only a  
104 small proportion of the training data contains introgressed fragments. It remains unclear whether Urnikyte  
105 et al. (2023) balanced their training data—ensuring a similar amount of non-introgressed and introgressed  
106 fragments—before applying ArchIE to detect archaic introgressed fragments in Lithuanian genomes.  
107 Hence, it is important to thoroughly document the details when applying machine learning approaches  
108 (Walsh et al. 2021), as factors like data preprocessing, training data, and hyperparameters can  
109 significantly impact the final performance of these models. Using version control tools with code hosting  
110 platforms like GitHub, model hosting platforms like Hugging Face, and reproducible workflow

111 management systems like Snakemake (Mölder et al. 2021) helps document the details and ensures that  
112 computational steps can be easily reproduced or modified by others (Huang 2024).

113 As interest in applying machine learning, particularly deep learning, to population genetics and  
114 evolutionary biology continues to grow (Huang et al. 2024), it is crucial for researchers to understand the  
115 underlying principles. For instance, since ArchIE is trained using simulated data that does not account for  
116 natural selection, its performance on data with natural selection should be carefully examined when  
117 applied in such contexts. Additionally, it is essential to develop robust machine learning applications that  
118 allow users to easily comprehend and adapt them to their own data (Huang 2024). One limitation of our  
119 study is that we used the pretrained model provided by MaLAdapt, which was trained on a specific  
120 human demographic model (Zhang et al. 2023). Currently, retraining MaLAdapt for a specific dataset is  
121 challenging due to its implementation. Reduced performance of MaLAdapt has been observed in other  
122 species (Romieu et al. 2024), likely due to demographic model misspecification. However, since the  
123 pretrained model was trained on a human population setting that includes Eurasians, to which the  
124 Lithuanians belongs, we expect the reduction in model performance in our analysis to be minimal.

## 125 **Acknowledgements**

126 J.H. and X.H. thank Martin Kuhlwilm for discussions and comments on the manuscript; and the Life  
127 Science Compute Cluster at the University of Vienna for providing computing resources.

## 128 **Competing interests**

129 J.H. and X.H. declare no conflict of interests.

## 130 **Author contributions**

131 X.H. designed the study. J.H. and X.H. analyzed the data and wrote the manuscript.

## 132 **Data availability**

133 The Snakemake workflow for reproducing the analysis can be found in [https://github.com/xin-](https://github.com/xin-huang/Lithuanian-archaic-introgression)  
134 [huang/Lithuanian-archaic-introgression](https://github.com/xin-huang/Lithuanian-archaic-introgression), last accessed October 10, 2024.

## 135 **References**

136 Abi-Rached, L., Jobin, M.J., Kulkarni, S., McWhinnie, A., Dalva, K., Gragert, L., Babrzadeh, F.,  
137 Gharizadeh, B., Luo, M., Plummer, F.A., Kimani, J., Carrington, M., Middleton, D., Rajalingam, R.,  
138 Beksac, M., Marsh, S.G.E., Maisers, M., Guethlein, L.A., Tavoularis, S., Little, A., Green, R.E.,

139 Norman, P.J., Parham, P., 2011. The shaping of modern human immune systems by multiregional  
140 admixture with archaic humans. *Science* **334**, 89–94.

141 Andrés, A.M., Hubisz, M.J., Indap, A., Torgerson, D.G., Degenhardt, J.D., Boyko, A.R., Gutenkunst,  
142 R.N., White, T.J., Green, E.D., Bustanmante, C.D., Clark, A.G., Nielsen, R., 2009. Targets of  
143 balancing selection in the human genome. *Mol Biol Evol* **12**, 2755–2764.

144 Bitarello, B.D., de Filippo, C., Teixeira, J.C., Schmidt, J.M., Kleinert, P., Meyer, D., Andrés, A.M., 2018.  
145 Signatures of long-term balancing selection in human genomes. *Genome Biol Evol* **10**, 939–955.

146 Bitarello, B.D., Brandt, D.Y.C., Meyer, D., Andrés, A.M., 2023. Inferring balancing selection from  
147 genome-scale data. *Genome Biol Evol* **15**, evad032.

148 Browning, B.L., Zhou, Y., Browning, S.R., 2018. A one-penny imputed genome from next generation  
149 reference panels. *Am J Hum Genet* **103**, 338–348.

150 Browning, B.L., Tian, X., Zhou, Y., Browning, S.R., 2021. Fast two-stage phasing of large-scale sequence  
151 data. *Am J Hum Genet* **108**, 1880–1890.

152 Chang, C.C., Chow, C.C., Tellier, L.C.A.M., Vattikuti, S., Purcell, S.M., Lee, J.J., 2015. Second-  
153 generation PLINK: Rising to the challenge of larger and richer datasets. *GigaScience* **4**, s13742-015-  
154 0047-8.

155 DeGiorgio, M., Lohmueller, K.E., Nielsen, R., 2014. A model-based approach for identifying signatures  
156 of ancient balancing selection in genetic data. *PLoS Genet* **10**, e1004561.

157 Ding, Q., Hu, Y., Jin, L., 2014. Non-Neanderthal origin of the HLA-DPB1\*0401. *J Biol Chem* **289**,  
158 10252.

159 Deng, Y., Nielsen, R., Song, Y.S., 2024. Robust and accurate Bayesian inference of genome-wide  
160 genealogies for large samples. bioRxiv. <https://doi.org/10.1101/2024.03.16.585351>, last accessed  
161 October 10, 2024.

162 Durvasula, A., Sankararaman, S., 2019. A statistical model for reference-free inference of archaic local  
163 ancestry. *PLoS Genet* **15**, e1008175.

164 Fijarczyk, A., Babik, W., 2015. Detecting balancing selection in genomes: Limits and prospects. *Mol Ecol*  
165 **14**, 3529–3545.

166 Gelabert, P., Bickle, P., Hofmann, D., Teschler-Nicola, M., Anders, A., Huang, X., Hämmerle, M., Olalde,  
167 I., Fournier, R., Ringbauer, H., Akbari, A., Cheronet, O., Lazaridis, I., Broomandkhoshbacht, N.,  
168 Fernandes, D.M., Buttinger, K., Callan, K., Candilio, F., Bravo, G., Curtis, E., Ferry, M., Keating, D.,  
169 Freilich, S., Kearns, A., Harney, É., Lawson, A.M., Mandl, K., Michel, M., Oberreiter, V., Zagorc, B.,  
170 Oppenheimer, J., Sawyer, S., Schattke, C., Ozdogan, K.T., Qiu, L., Workman, J.N., Zalzala, F.,  
171 Mallick, S., Mah, M., Micco, A., Pieler, F., Pavuk, J., Šefčáková, A., Lazar, C., Vasic, R., Starovic, A.,  
172 Djuric, M., Škrivanko, M.K., Šlaus, M., Bedić, Ž., Novotny, F., Szabó, L.D., Cserpák-Laczi, O.,

173 Hága, T., Hajdú, Z., Mirea, P., Nagy, E.G., Virág, Z.M., Horváth, A.M., Horváth, L.A., Biró, K.T.,  
174 Domboróczki, L., Szeniczey, T., Jakucs, J., Szelekovszky, M., Zoltán, F., Sztáncsuj, S., Tóth, K.,  
175 Csengeri, P., Pap, I., Patay, R., Putica, A., Vasov, B., Havasi, B., Sebők, K., Raczky, P., Lovász, G.,  
176 Tvrđý, Z., Rohland, N., Novak, M., Ruttkay, M., Krošláková, M., Batora, J., Cheben, I., Boric, D.,  
177 Dani, J., Kuhlwilm, M., Palamara, P.F., Hajdu, T., Pinhasi, R., Reich, D., 2024. Social and genetic  
178 diversity in the first farmers of Central Europe. *Nature Hum Behav.*

179 Gower, G., Picazo, P.I., Fumagalli, M., Racimo, F., 2021. Detecting adaptive introgression in human  
180 evolution using convolutional neural networks. *eLife* **10**, e64669.

181 Hedrick, P.W., 2007. Balancing selection. *Curr Biol* **17**, R230–R231.

182 Huang, X., Kruisz, P., Kuhlwilm, M., 2022. sstar: A Python package for detecting archaic introgression  
183 from population genetic data with  $S^*$ . *Mol Biol Evol* **39**, msac212.

184 Huang, X., 2024. Developing machine learning applications for population genetic inference: Ensuring  
185 precise terminology and robust implementation. *EcoEvoRxiv*. <https://doi.org/10.32942/X2N90M>, last  
186 accessed October 10, 2024.

187 Huang, X., Rymbekova, A., Dolgova, O., Lao, O., Kuhlwilm, M., 2024. Harnessing deep learning for  
188 population genetic inference. *Nat Rev Genet* **25**, 61–78.

189 Hudson, R.R., 2002. Generating samples under a Wright–Fisher neutral model of genetic variation.  
190 *Bioinformatics* **18**, 337–338.

191 Mo, Z., Siepel, A., 2023. Domain-adaptive neural networks improve supervised machine learning based  
192 on simulated population genetic data. *PLoS Genet* **19**, e1011032.

193 Mölder, F., Jablonski, K.P., Letcher, B., Hall, M.B., Tomkins-Tinch, C.H., Sochat, V., Forster, J., Lee, S.,  
194 Twardziok, S.O., Kanitz, A., Wilm, A., Holtgrewe, M., Rahmann, S., Nahnsen, S., Köster, J., 2021.  
195 Sustainable data analysis with Snakemake. *FI000 Res* **10**, 33.

196 Prüfer, K., Racimo, F., Patterson, N., Jay, F., Sankararaman, S., Sawyer, S., Heinze, A., Renaud, G.,  
197 Sudmant, P.H., de Filippo, C., Li, H., Mallick, S., Dannemann, M., Fu, Q., Kircher, M., Kuhlwilm,  
198 M., Lachmann, M., Meyer, M., Ongyerth, M., Siebauer, M., Theunert, C., Tandon, A., Moorjani, P.,  
199 Pickrell, J., Mullikin, J.C., Vohr, S.H., Green, R.E., Hellmann, I., Johnson, P.L.F., Blanche, H., Cann,  
200 H., Kitzman, J.O., Shendure, J., Eichler, E.E., Lein, E.S., Bakken, T.E., Golovanova, L.V.,  
201 Doronichev, V.B., Shunkov, M.V., Derevianko, A.P., Viola, B., Slatkin, M., Reich, D., Kelso, J.,  
202 Pääbo, S., 2014. The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature*  
203 **505**, 43–49.

204 Racimo, F., Sankararaman, S., Nielsen, R., Huerta-Sánchez, E., 2015. Evidence for archaic adaptive  
205 introgression in humans. *Nat Rev Genet* **16**, 359–371.

206 Ray, D.D., Flagel, L., Schrider, D.R., 2024. IntroUNET: Identifying introgressed alleles via sematic  
207 segmentation. *PLoS Genet* **20**, e1010657.

208 Romieu, J., Camarata, G., Crochet, P.A., de Navascués, M., Leblois, R., Rousset, F., 2024. Performance  
209 evaluation of adaptive introgression classification methods. bioRxiv.  
210 <https://doi.org/10.1101/2024.06.12.598278>, last accessed October 10, 2024.

211 Siewert, K.M., Voight, B.F., 2017. Detecting long-term balancing selection using allele frequency  
212 correlation. *Mol Biol Evol* **34**, 2996–3005.

213 The 1000 Genomes Project Consortium. 2015. A global reference for human genetic variation. *Nature*  
214 **526**, 68–74.

215 Turner, S.D., 2018. qqman: an R package for visualizing GWAS results using Q-Q and manhattan plots. *J*  
216 *Open Source Softw* **3**, 731.

217 Urnikyte, A., Masiulyte, A., Pranckeniene, L., Kučinskas, V., 2023. Disentangling archaic introgression  
218 and genomic signatures of selection at human immunity genes. *Infect Genet Evol* **116**, 105528.

219 Voight, B.F., Kudravalli, S., Wen, X., Pritchard, J.K., 2006. A map of recent positive selection in the  
220 human genome. *PLoS Biol* **4**, e72.

221 Walsh, I., Fishman, D., Garcia-Gasulla, D., Titma, T., Pollastri, G., ELIXIR Machine Learning Focus  
222 Group, Harrow, J., Psomopoulos, F.E., Tosatto, S.C.E., 2021. DOME: Recommendations for  
223 supervised machine learning validation in biology. *Nat Methods* **18**, 1122–1127.

224 Yasukochi, Y., Ohashi, J., 2017. Elucidating the origin of *HLA-B\*73* allelic lineage: Did modern humans  
225 benefit by archaic introgression? *Immunogenetics* **69**, 63–67.

226 Zhang, X., Kim, B., Singh, A., Sankararaman, S., Durvasula, A., Lohmueller, K.E., 2023. *MaLAdapt*  
227 reveals novel targets of adaptive introgression from Neanderthals and Denisovans in worldwide  
228 human populations. *Mol Biol Evol* **40**, msad001.