1	Review		
2	Decoding Genomic Landscapes of Introgression		
3	Xin Huang (黄欣) ^{1, 2, *} , Josef Hackl ¹ , Martin Kuhlwilm ^{1, 2, *}		
4	1 Department of Evolutionary Anthropology, University of Vienna, Vienna, Austria		
5	2 Human Evolution and Archaeological Sciences (HEAS), University of Vienna, Vienna, Austria		
6	* Correspondence: xin.huang@univie.ac.at (X. Huang) and martin.kuhlwilm@univie.ac.at (M.		
7	Kuhlwilm)		
8			
9	Genomic landscapes of introgression provide valuable information for how different		
10	evolutionary processes interact and leave signatures in genomes. The recent expansion of		
11	genomic datasets across diverse taxa, together with advances in methodological		
12	development, has created new opportunities to investigate the impact of introgression along		
13	individual genomes in various clades, making the precise identification of introgressed loci		
14	a rapidly evolving area of research. In this review, we summarize recent methodological		
15	progress within three major categories: summary statistics, probabilistic modeling, and		
16	supervised learning. We examine how these approaches have been applied to data beyond		
17	humans and discuss the challenges associated with their application. Finally, we outline		
18	future directions for each category, including accessible implementation, transparent		
19	analysis, and systematic benchmarking.		
20			
21	Highlights		
22	• Recent advances in methods and tools have enabled the study of genomic landscapes of		
23	introgression across diverse and complex evolutionary scenarios, including adaptive and		
24	ghost introgression.		
25	• Despite their long history, summary statistics-based methods continue to evolve, with		
26	new implementations broadening their applicability across taxa.		

- Probabilistic modeling is a major approach that provides a powerful framework to
 explicitly incorporate evolutionary processes and has yielded fine-scale insights across
 diverse species.
- Supervised learning is an emerging approach with great potential, particularly when the
 detection of introgressed loci is framed as a semantic segmentation task.
- Various methods have been applied across clades, revealing introgressed loci linked to
 immunity, reproduction, and environmental adaptation, especially in cases of adaptive
 and ghost introgression.
- 35

36 Genomic Landscapes of Introgression

37 Introgression (see Glossary) plays an important role in evolution. Beyond merely studying 38 introgression events through phylogenetic approaches [1], understanding their genomic 39 footprints-how introgressed loci are retained, eliminated, or distributed within genomes-is 40 essential (Figure 1A), because these patterns are shaped by demographic histories, selective 41 pressures, and genomic architectures. Although many methods for detecting introgressed loci 42 have provided crucial insights into past gene flow from archaic hominins to contemporary 43 human populations, they have largely been developed with modern and archaic human genomes [2]. 44

45 With the increasing availability of genomic data from diverse taxa, such as hickories, peafowls,

46 and corn earworms [3–5], studies on the detection of introgressed loci are becoming more

47 prevalent in the field. Consequently, characterizing genomic landscapes of introgression offers

48 deeper insights into the evolutionary forces driving hybridization outcomes and the functional

49 roles of introgressed loci in different species (Figure 1B). Although developed in the pre-

50 genomic era, summary statistics-based approaches remain widely used. In recent years,

51 substantial methodological advances have emerged from probabilistic modeling and supervised

52 learning (Figure 2). These developments motivate a critical assessment of existing approaches

and a comprehensive review of emerging strategies for decoding genomic landscapes of

54 introgression.



56 Figure 1. Introgression and its genomic landscapes. (A) Individual genomes in the target population 57 carry short introgressed segments because recombination breaks down long haplotypes. (B) Due to genetic drift, the frequency of these segments in the target population varies along chromosomes. 58 59 (C) Population setting in the context of introgression: deriving from a common ancestor, lineage-60 specific genetic variation arises over time in the diverging populations (color gradients). 61 Introgression transfers this variation from the source population (blue) into the target population 62 (orange). A reference population (green), more closely related to the target than the source 63 population, is often used to determine non-introgressed variation in a lineage. Some methods use an 64 outgroup (grey) to infer whether an allele is ancestral or derived within this topology. If no data is 65 available from the source population, the scenario is referred to as ghost introgression. (D) Adaptive 66 introgression represents a special case where introgressed ancestry surrounding an adaptive locus rises in frequency beyond the expectation under neutrality. When multiple source populations are 67 68 involved, fragments from divergent origins may co-occur in the same genome, potentially 69 confounding the detection of introgressed loci. Genomic landscapes of introgression might have 70 different distribution patterns and dynamics in different clades (illustrated by different color 71 gradients).

72 Summary Statistics-based Methods

55

Summary statistics are simple yet effective approaches for exploratory data analysis, commonly
used to distill complex genomic data into simple numeric representations, including frequency-

- 75 based, linkage-based, or topology-based measures (Figure 2A). The D statistic is a widely used
- 76 method for detecting genome-wide evidence of introgression [6], based on asymmetries in
- derived allele sharing between the target and source populations relative to a reference
- 78 population and an outgroup (Figure 1C). However, it is not well suited for pinpointing
- introgressed loci, as it can be biased in regions of low genetic diversity [7-9].
- 80 To address this, alternative statistics that are calculated in windows along the genome have been
- 81 developed. The dynamic estimator of the proportion of introgression (f_d) and the distance fraction
- 82 (d_f) both scale the observed excess of shared derived alleles but differ in their normalization
- approaches: f_d uses the maximum possible level of derived allele sharing due to introgression,
- 84 whereas d_f normalizes against the expected derived allele sharing under the species tree topology
- 85 [7,8]. Both reduce the bias of *D* by avoiding direct dependence on the genetic distance between
- the reference and target populations [8]. The D^+ statistic further extends D by incorporating both
- 87 shared derived and ancestral alleles, thereby increasing the number of informative sites and
- reducing variance [9]. Additionally, f_d has a bounded variant, f_{dM} , which ranges from -1 to 1, and
- is symmetrically distributed around zero under no introgression [10].
- 90 To detect loci under adaptive introgression (Figure 1D), additional methods have been
- 91 introduced that also leverage allele sharing patterns between the target and source populations.
- 92 These include the number of uniquely shared sites (U) and the quantile of derived allele
- 93 frequency distributions (Q) [11]. Such methods retain variants shared between target and source
- 94 populations that are rare or absent in the reference population, thereby enriching for candidates
- 95 likely introduced via introgression. While adaptive variants in the target population often reach
- 96 high frequency due to **positive selection**, these methods are primarily sensitive to such cases and
- 97 may miss beneficial alleles that are at low or intermediate frequencies.
- In cases of ghost introgression, where source samples are unavailable, inference relies on alleles
 present in the target but absent in the reference population. S* was initially developed to detect
- 100 archaic introgression in human populations without a source genome, by identifying clusters of
- 101 private mutations in strong linkage within the target population [12]. However, S^* does not
- 102 account for local mutation or recombination rate variation. To handle this, one approach involves
- 103 simulating data under varying local rates, fitting a **generalized additive model** to the resulting
- 104 S* scores, and using this model to estimate expected values and assess significance in real data

105 [13,14]. Another approach, *S'*, modifies the calculation of *S** by incorporating local mutation and 106 recombination rates [15].

107 Instead of frequency-based or linkage-based information, Twisst and its upgraded version

108 Twisst2 summarize subtree topologies, each formed by selecting one sample per species, to

assess how often gene trees support alternative species tree relationships [16,17]. This topology

110 weighting approach indirectly captures discrepancies between gene trees and the species tree,

111 which are explicitly tested by *D*[9]. While *D* is limited to four taxa, Twisst can in principle

112 handle any number, but becomes increasingly impractical with more than six due to the large

113 number of possible species tree topologies [16].

114 Probabilistic Model-based Methods

115 Probabilistic modeling enables model-based inference of introgressed loci by defining the

relationship between observed genetic variation and underlying evolutionary processes through

probability distributions and performing inference under likelihood-based or Bayesian

118 frameworks (Figure 2B). For example, IBDmix estimates the probability of identity-by-descent

119 between the target and source populations at each locus, thereby eliminating the need for a

120 reference population [19]. It surpasses S^* on simulated data, especially in scenarios where

121 introgressed fragments are also present in the reference population [19].

Probabilistic methods can also address compound scenarios such as adaptive introgression from
a ghost population. VolcanoFinder is a likelihood-based method designed for this setting [20]. It

identifies genomic regions exhibiting a characteristic "volcano" pattern of genetic diversity,

125 which is marked by reduced diversity at the selected site flanked by elevated diversity, reflecting

a selective sweep on an introgressed, highly divergent haplotype. However, it often fails to

127 distinguish adaptive introgression from classic selective sweeps in both real and simulated data

128 [21,22]. This limitation may stem from its reliance on a specific demographic model that

assumes a beneficial allele first undergoes fixation via a selective sweep in the source population

130 before introgressing into the target population, where it sweeps again. A recent study suggests

131 that VolcanoFinder performs well only under conditions of strong selection and high divergence

132 between source and target populations [22].

133 Probabilistic models can further capture gene tree discordance with fine resolution. For example,

134 ancestral recombination graphs (ARGs) provide a more detailed representation of ancestry and

- 135 can be interpreted as a sequential model of local genealogies along the genome, which is also
- 136 known as a tree sequence [23,24]. ARGweaver-D is a recent method that applies ARGs to detect
- 137 introgressed loci. In addition to identifying Neanderthal and Denisovan introgressed loci in
- 138 modern human genomes, it also detects putative introgressed loci from a super-archaic ghost
- 139 population into Denisovans [25]. However, its Bayesian framework results in high computational
- 140 cost, limiting its scalability to larger sample sizes and more complex demographic models.
- 141 Although recent methods have improved the scalability of ARG inference [23], most are
- 142 designed to reconstruct genealogies rather than directly detecting introgressed loci, which
- 143 requires additional implementation to extract such signals.
- 144 While ARGs explicitly represent the full genealogical history, hidden Markov models
- 145 (HMMs)—a classical probabilistic model for sequential data in machine learning [26]—treat 146 ancestry as a latent state to be inferred from observed sequences. HMMs have long been used for 147 local ancestry inference (LAI) [27]; early applications for detecting introgressed loci adapted 148 these methods, which rely on a reference panel (i.e., the source population in this review) to label 149 ancestry [28]. However, they are unsuitable for detecting ghost introgression where the source 150 population is unavailable and their performance may be also questionable when the source 151 population is represented by only a few samples. More recent work has designed HMMs 152 specifically for introgression detection, including applications without source populations, in 153 multi-source scenarios, and using low-coverage, contaminated sequencing data, as well as for 154 jointly inferring adaptive introgression and the strength of natural selection acting on the 155 introgressed loci [29-33].

156 Supervised Learning-based Methods

157 The emergence of machine learning, particularly deep learning, reflects a broader trend across 158 disciplines, including genetics [34]. As large-scale genomic data become available for an 159 expanding range of populations and taxa, traditional model-based methods increasingly struggle 160 with both processing feasibility and the challenge of constructing detailed models for each group. 161 Moreover, as sample sizes and variant densities grow, the curse of dimensionality may further 162 limit the effectiveness of traditional methods [26]. As these approaches do not scale well, data-163 driven alternatives that forgo explicit mechanistic models of the evolutionary process are gaining 164 importance for their scalability and flexibility in analyzing high-dimensional genomic data.

- 165 Supervised learning, a machine learning paradigm that maps inputs to labeled outputs, has
- 166 attracted growing interest in population genetics method development [35]. For detecting
- 167 introgressed loci, current approaches comprise two main groups: those that use predefined
- 168 summary statistics, such as the ones discussed above, as input features, and those that
- 169 automatically extract features from raw data. These methods frame the detection of introgressed
- 170 loci as a classification problem, aiming to determine whether a given genomic region, variant, or
- allele is introgressed. Methods such as ArchIE, FILET, and MaLAdapt belong to the first group,
- 172 employing logistic regression or Extra-Trees classifiers [36–38]. In contrast, genomatnn,
- 173 ERICA, and IntroUNET represent the second group, using convolutional neural networks to
- 174 directly learn from genotype matrices [39–41].

A major challenge for applying supervised learning in evolutionary biology is the lack of labeled data, that is, ground truth indicating whether a locus is introgressed. To overcome this limitation, simulated datasets are used to train machine learning models, which are then applied to empirical data for prediction (Figure 2C). Although simulated data may not perfectly reflect reality, several supervised methods have demonstrated promising results. For example, ArchIE has been shown to surpass *S** and *S'* by incorporating genetic distance between genomes from the reference and target populations [36]. Similarly, MaLAdapt outperforms the *U* statistic, the *Q* statistic, and

182 VolcanoFinder under a Neanderthal introgression model, and have revealed novel candidates for

adaptive introgression in modern human populations [38].

- 184 By intersecting introgressed regions predicted by summary statistic-based outlier detection,
- 185 MaLAdapt, and genomatnn, circadian loci have been implicated as adaptively introgressed from
- archaic hominins [42]. Moreover, genomatnn, MaLAdapt, and the U and Q statistics all support
- 187 *BNC2*, a gene associated with human pigmentation and previously identified as a target of
- 188 positive selection in modern Europeans [43,44], whereas VolcanoFinder does not detect such a
- signal [20]. While ERICA employs deep learning, it shares a core principle with Twisst by
- 190 predicting the proportions of gene tree topologies within a genomic region to identify
- 191 introgressed loci through gene tree discordance [40].
- 192 Among these approaches, IntroUNET is particularly interesting, as it frames the identification of
- 193 introgressed alleles from the ghost population as a **semantic segmentation** task, which is a
- 194 fundamental problem in modern compute vision [45], and thus, in principle, enables high-

- resolution predictions at the allele level (Figure 3A), which may be difficult to achieve using
- 196 other approaches. This capability may be valuable for precisely delineating the boundaries of
- 197 introgressed segments.



198

199 Figure 2. Conceptual overview of computational approaches for detecting introgressed loci. (A) 200 Underlying genomic data from different individuals representing different populations, where 201 variants can be observed within sliding windows along the genome. The star denotes a private 202 variant that is observed in genomes of the target population and absent in the reference population. 203 (B) Summary statistics-based methods summarize genomic information into statistic values (S) from 204 the reference and target genomes, and optionally from a source genome. They typically apply outlier 205 detection to identify putative introgressed loci based on a threshold (T). An outlier is highlighted 206 with three asterisks. (C) Probabilistic model-based methods describe how the data are generated 207 under a probabilistic framework, based on various strategies to determine different patterns. For 208 example, an HMM represents transitions between hidden states, where S₁ denotes the non-209 introgressed state and S_2 denotes the introgressed state. The model defines how these states emit 210 observations, such as number of private variants, enabling likelihood estimation and model fitting to

- 211 observed data. (D) Supervised learning-based methods rely on labeled training data to learn models
- 212 that predict the introgression status of regions or variants in target genomes. For example, in
- 213 genotype matrices of ancestral (grey) and derived (colored) alleles, an artificial neural network
- 214 predicts whether alleles are introgressed or non-introgressed. Labels are typically generated from
- 215 computer simulations.

216 Recent Applications beyond Humans

- 217 Recent methodological innovations have provided a variety of tools for decoding the genomic
- 218 landscapes of introgression beyond humans (Figure 3B). A wide range of taxa have been
- 219 investigated with summary statistics-based methods [46–55]. These approaches are now widely
- used thanks to recent implementations like Dsuite [56], and a comprehensive survey is beyond
- scope. Still, recent studies offer notable examples of adaptive introgression: f_{dM} plus selection
- scans identified flowering-time genes in *Brassica napus*; and the U and Q statistics detected
- sperm function genes in sticklebacks and high-altitude candidates in Tibetan cattle [46,49,54].
- 224 These findings highlight recurrent targets among genes involved in key biological functions.
- 225 Probabilistic model-based approaches are also extensively utilized in non-human species,
- 226 demonstrating their versatility across diverse introgression scenarios. For example, IBDmix has
- inferred introgressed fragments in baboons and bears [58–60]. hmmix has identified loci from
- 228 ghost introgression in orcas, canids, and the extinct Columbian mammoth, and in great apes
- 229 when combined with the S* statistic, underscoring the benefits of integrative analyses
- 230 [3,50,53,61–63]. In Tibetan canids, the high-altitude adaptation gene *EPAS1*, previously linked to
- 231 Denisovan introgression in humans, may also derive from a ghost lineage [62,64]. admixfrog,
- leveraging both ancient and modern genomes, has estimated ancestry proportions in ancient
- bears, detected immunity-related introgression in Alpine ibex, and resolved fine-scale
- 234 introgression patterns in chimpanzees using non-invasive fecal samples resembling degraded
- ancient DNA [65–67]. For adaptive introgression, VolcanoFinder has identified candidate loci
- from ghost lineages in gorillas, hickories, and pigs, associated with bitter taste perception,
- defense response, and commercial traits, respectively [3,53,68]. AHMM-S has detected
- 238 insecticide-resistance loci in fruit flies, though its multi-locus extension, AHMM-MLS, suggests
- AHMM-S may overestimate selection coefficients [32,33].
- 240 To date, supervised learning-based methods have been applied to a limited number of non-human
- taxa, such as fruit flies, butterflies, and rice [37,40,41]. In these applications, the primary goal

242 has been to assess whether model predictions align with results from previous approaches, rather 243 than to investigate new biological questions. Nonetheless, ERICA has identified multiple 244 domestication-related loci in rice as candidates for adaptive introgression. This limited scope is 245 likely due to the technical complexity of machine learning, as most implementations are tailored 246 to specific datasets and are not easily applied beyond their original training context, reflecting broader challenges in software development within evolutionary biology [69,70]. As population-247 248 scale genomic datasets and machine learning algorithms continue to develop, broader application 249 across diverse lineages is expected. Such efforts will help refine our understanding of 250 introgression landscapes and population interactions throughout evolutionary history. 251 Beyond analyzing empirical datasets, some studies have explored how different approaches 252 perform in non-human scenarios using simulated data. For example, both S* and S' perform well 253 under a Neanderthal introgression model, but only S* remains effective in a bonobo ghost 254 introgression scenario [18]. This difference may be due to its recent implementation, sstar, which 255 provides a flexible computational framework applicable to diverse demographic scenarios, 256 whereas SPrime, the implementation of S', is hard-coded with parameters specific to an out-of-257 Africa Neanderthal-admixture model [18,71]. Furthermore, a recent study suggests that the Q 258 statistic performs comparably or better than genomatnn, MaLAdapt, and VolcanoFinder under non-human demographic models, indicating that summary statistics continue to be valuable even 259

when more advanced methods are available [22,39].



- 262 Figure 3. Methodological features of different open-source implementations. (A) Different
- **263** prediction levels across genomes refer to the resolution at which an introgressed label is assigned.
- 264 Window level: fixed-length genomic windows, defined by base pairs or by number of variants, are
- 265 typically summarized across all variants and samples without individual-level resolution. Segment
- 266 level: continuous introgressed haplotypes of varying lengths are inferred for each individual genome.
- 267 Variant level: individual genomic variants, such as single nucleotide polymorphisms, are classified as
- 268 introgressed or not. Allele level: the status of each allele at a segregating site is individually
- 269 determined. (B) Feature assessment of implementations representing different methodological
- 270 approaches. Implementations with summary statistics-based methods include Dsuite
- 271 (https://github.com/millanek/Dsuite) for the d_{f_s} f_d, and f_{dM} statistics; sai (https://github.com/xin-
- 272 <u>huang/sai</u>) for the U and Q statistics; sstar (<u>https://github.com/xin-huang/sstar</u>) for the S^* statistic;
- 273 SPrime (<u>https://github.com/browning-lab/sprime</u>) for the S' statistic; and Twisst2
- 274 (https://github.com/simonhmartin/twisst2) for topology weighting. Implementations with
- 275 probabilistic model-based methods comprise IBDmix
- 276 (https://github.com/PrincetonUniversity/IBDmix), VolcanoFinder

- 277 (<u>https://degiorgiogroup.fau.edu/vf.html</u>), ARGweaver-D
- 278 (http://compgen.cshl.edu/ARGweaver/doc/argweaver-d-manual.html), admixfrog
- 279 (https://github.com/BenjaminPeter/admixfrog), ArchaicSeeker 2.0 (https://github.com/Shuhua-
- 280 Group/ArchaicSeeker2.0), AHMM-S (https://github.com/jesvedberg/Ancestry HMM-S), AHMM-
- 281 MLS (<u>https://github.com/genicos/ahmm_mls</u>), hmmix
- 282 (https://github.com/LauritsSkov/Introgression-detection). Implementations with supervised
- 283 learning-based methods encompass ArchIE (<u>https://github.com/sriramlab/ArchIE</u>), FILET
- 284 (https://github.com/kr-colab/FILET), MaLAdapt (https://github.com/xzhang-popgen/maladapt),
- 285 ERICA (<u>https://github.com/YuboZhangPKU/ERICA</u>), genomatnn
- 286 (https://github.com/grahamgower/genomatnn), IntroUNET
- 287 (<u>https://github.com/SchriderLab/introNets</u>). For methods with multiple implementations such as
- 288 the d_{f_a} f_a, and S^{*} statistics, whether by the same or different authors, only the most recent version is
- assessed. "Assumes a demographic model" refers to using a specific model to tune parameters
- 290 (SPrime), condition inference (ARGweaver-D and VolcanoFinder), or simulate training data (sstar
- and supervised learning-based methods).

292 Challenges

- 293 Despite substantial methodological progress, decoding genomic landscapes of introgression is 294 still fraught with challenges due to confounding factors, model misspecification, and analysis 295 opacity that can bias or obscure inference. One major source of confounding arises from 296 evolutionary processes that mimic the genomic signatures of introgression. For example, incomplete lineage sorting (ILS) can produce gene tree discordance similar to that expected 297 298 under introgression. Although the D statistic is expected to distinguish ILS and introgression at 299 the whole-genome level, this is not the case at the locus level [9]. Also, population structure in 300 unsampled or ancestral lineages can generate spurious signals resembling ghost introgression, 301 even in the absence of gene flow [72]. A long-standing debate in human evolution concerns 302 whether ghost introgression occurred in African populations, with different conclusions from 303 different demographic inference approaches, although one study nonetheless applied ArchIE to
- and examined putatively introgressed fragments [73–75]. Furthermore, **long-term balancing**
- **selection** may give rise to patterns that resemble adaptive introgression [76,77].
- Another challenge is model misspecification, which can arise in several forms. First, somemethods embed rigid assumptions in their model design. For instance, hmmix assumes the

308 presence of both introgressed and non-introgressed hidden states in the data [29]. If introgression 309 is absent, the model may nonetheless infer false signals simply by fitting its predefined state 310 structure. Second, methods like SPrime or VolcanoFinder often assume simplified or idealized 311 demographic models. Such assumptions may not hold in more complex or less ideal evolutionary 312 scenarios, potentially limiting applicability. Third, in supervised learning, performance may degrade if the demographic model used to simulate training data differs substantially from the 313 314 test scenario [78]. Finally, some deep learning architectures, such as convolutional neural 315 networks, have architectural constraints such as requiring fixed input shapes and being sensitive 316 to the order of input samples, which may limit their flexibility across various datasets [34]. 317 In practice, a third challenge is analysis opacity. While current critiques of machine learning 318 often focus on their interpretability [79], a lack of transparency can also arise from the analysis 319 procedure, including undocumented preprocessing steps, hard-coded parameters, or 320 discrepancies between published methods and their actual implementations [57,70,77]. These 321 issues, not only relevant for machine learning [57], frequently force researchers to inspect source 322 code directly to verify correctness, thereby impeding reproducibility and slowing scientific 323 progress. For example, the performance of IntroUNET may be affected by training data that 324 inadvertently retained information of polymorphic sites from an unintended fourth population, as it reused the demographic model from ArchIE for training, which was described as a three-325 326 population model but in fact included a fourth, and by repeated training datasets caused by 327 unexpected behavior in its modified simulator [70,77]. This challenge can be addressed through 328 transparent reporting, reproducible workflows, and community standards. Furthermore, the lack 329 of accessible and robust implementations has hindered consistent benchmarking, which is 330 essential for method evaluation. For example, the performance of the Q statistic differs between 331 two studies [22,38]. The lack of standardized Q statistic implementations and transparent 332 documentation makes it difficult to determine whether discrepancies result from implementation, 333 data processing, or demographic models. Moreover, studies have used different approaches to 334 detect adaptive introgression, either by combining introgression signals with selection scans or 335 by applying dedicated methods [21,38]. However, it is still uncertain which approach performs 336 best, or under what circumstances each should be applied. The recent emergence of machine 337 learning benchmarks provides a valuable reference and highlights the importance of structured

comparison and shared standards when evaluating the performance of different methods
 (https://iclr.cc/virtual/2024/invited-talk/21799).

340 **Outlook**

341 As introgression detection continues to mature as a methodological field, future progress will 342 rely on extending current approaches-including statistical, probabilistic, and supervised 343 learning-based methods-to accommodate increasingly complex evolutionary scenarios and data 344 types. Owing to their simplicity, interpretability, and computational efficiency, summary 345 statistics-based methods remain attractive. In particular, those incorporating linkage 346 disequilibrium (LD) information may prove useful in more complex cases. However, no 347 existing summary statistic can currently distinguish loci resulting from multi-source 348 introgression (Figure 3B). Extending such approaches toward locus-level detection, especially 349 under selection or multiple pulses of gene flow, remains an open and important challenge.

350 By modeling introgression under explicitly defined evolutionary scenarios, probabilistic methods 351 allow key processes such as mutation, recombination, and natural selection to be incorporated 352 into a unified framework. This capacity enables not just detection of introgressed loci but also 353 quantitative characterization of their properties, such as estimating the age of introgressed 354 variants, the length distribution of introgressed fragments, or the strength of selection acting on 355 them. Extending these probabilistic models to decode genomic landscapes of introgression 356 shaped by multiple evolutionary forces continues to be a key area of development [58,80,81]. 357 Further efforts may focus on improving scalability to large datasets and enhancing robustness to 358 model misspecification.

359 Supervised learning, and machine learning more broadly, holds great potential. In principle, these 360 approaches can integrate diverse data types and achieve high-resolution predictions. However, 361 current applications require high-quality input data and do not support polyploid datasets (Figure 362 3B). While other approaches are also limited in this regard, machine learning approaches, being 363 data-driven, may offer greater flexibility for accommodating such complexities in the future. 364 Extending these implementations to support low-coverage data and to unify the analysis of 365 neutral, adaptive, ghost, and multi-source introgression would be highly desirable. Most 366 importantly, such implementations should be accessible to the community and not tailored to a 367 specific species. Beyond model performance, software engineering is critical for ensuring that

368 machine learning methods are reproducible, maintainable, and broadly usable across datasets and369 research groups.

370 Another open question is how well machine learning models can generalize across demographic 371 conditions. Although supervised learning-based methods typically employ simulations to 372 generate labeled training data but indeed make no demographic assumptions during inference. 373 Performance degradation under mismatched test scenarios is likely caused by current 374 implementations relying on training data simulated under a specific demographic model, which 375 may lead to overfitting. To mitigate this, training on a diverse set of demographic models may 376 improve generalizability, allowing models to integrate signals from multiple evolutionary 377 processes. For example, ERICA was trained on data with a range of ILS and gene flow settings, 378 which makes it adaptable to diverse scenarios, although it may underperform compared to 379 models trained under the exact test scenario [40].

Alongside supervised learning, **unsupervised learning** and **self-supervised learning** also show promise, as these paradigms do not use labeled data and therefore avoid the requirement to generate simulated training data. For instance, outlier detection, which is central to summary statistics-based methods, can be naturally extended using **deep generative models** such as **variational autoencoders** [82]. Similarly, recent extensions of LAI have incorporated deep learning architectures [27,34], raising the possibility that these approaches could be repurposed for detecting introgressed loci. Additionally, recent trends in genetics involve developing

genomic language models using self-supervised learning [83,84]. This paradigm has

demonstrated strong generalization ability across diverse scenarios and could help improve the
robustness and transferability of machine learning models for detecting introgressed loci [85].

390 The integration of different methodological paradigms also presents an important opportunity. 391 For example, summary statistics-based and probabilistic model-based approaches, which are 392 typically interpretable and grounded in explicit evolutionary assumptions, can serve as valuable 393 baselines for assessing the performance and reliability of emerging machine learning-based 394 methods, while also helping to improve interpretability and robustness. Furthermore, developing 395 standardized benchmark datasets that span a range of demographic scenarios and evolutionary 396 processes will be crucial for systematic comparisons across methods. Such integrative efforts 397 will not only support methodological advancement but also accelerate biological discovery.

398 Concluding Remarks

399 The detection of introgressed loci has evolved into a diverse methodological field encompassing 400 summary statistics, probabilistic modeling, and supervised learning. Each class of methods offers 401 distinct advantages: summary statistics remain efficient and interpretable for initial scans, 402 probabilistic modeling enables principled inference under explicit evolutionary assumptions, and 403 machine learning offer scalability and potential for discovering patterns beyond predefined 404 models. Rather than converging on a single best method, future progress will likely depend on 405 leveraging the complementarity between approaches, improving transparency and benchmarking, 406 and developing tools that are robust to real-world complexity (see Outstanding questions). As 407 evolutionary datasets continue to expand in scale and scope, refining the approaches for 408 decoding genomic landscapes of introgression is essential for understanding how gene flow has 409 shaped genomes across the tree of life.

410

411 **Outstanding Questions**

410		Ano there should not terms of intrographical landscenes compare spacing that could inform
412	•	Are there shared patterns of introgression landscapes across species that could inform
413		general evolutionary principles?
414	•	How can the biological interpretability and analytical transparency of introgressed loci
415		inferred by complex models, especially those using machine learning, be improved?
416	•	How can the potential of machine learning, including genomic language models, be fully
417		leveraged to decode introgression landscapes, and under what conditions do these
418		approaches outperform traditional methods?
419	•	How can computational tools be developed to be accessible, generalizable across species,
420		and robust under variation in data quality, confounding factors, and model
421		misspecification?
422	•	Can different methods be systematically evaluated under diverse demographic models
423		and confounding factors to clarify performance discrepancies across studies, establish
424		consistent benchmarks, and identify which methods are best suited to specific scenarios
425		for accurate and reliable inference?

427	Glossary
428	Adaptive introgression: the transfer of genetic variants from one distant lineage into another via
429	gene flow, followed by natural selection favoring those variants in the recipient population.
430	Ancestral recombination graph: a graph-based data structure that exhaustively represents the
431	evolutionary relationships between a set of genomes, accounting for both coalescent and
432	recombination events.
433	Convolutional neural network: a type of artificial neural network that is particularly effective
434	for grid-like data.
435	Deep generative model: generative models that use deep neural networks to learn and sample
436	from complex data distributions.
437	Deep learning: a machine learning approach that utilizes deep neural networks to learn
438	hierarchical representation from data.
439	Extra-Trees classifier: an ensemble machine learning algorithm that builds multiple decision
440	trees using random splits and the entire training dataset to reduce variance and enhance
441	generalization.
442	Generalized additive model: a statistical model that captures non-linear relationships between
443	inputs and outputs by combining smooth functions additively, while preserving interpretability.
444	Genomic language model: a machine learning model adapted from natural language processing
445	to investigate genomic problems.
446	Ghost introgression: gene flow from a population which is not directly represented in genomic
447	data, as either unsampled or extinct but inferred from recipient populations.
448	Identity-by-descent: identical genomic segments shared among individuals that are inherited
449	from a common ancestor without being broken by recombination.
450	Incomplete lineage sorting (ILS): the phenomenon where gene trees fail to match the species
451	tree because ancestral polymorphisms are retained and randomly sorted through rapid speciation
452	events.
453	Introgression: the phenomenon of transferring genetic material across genetically divergent

454 populations.

455 Linkage disequilibrium (LD): the discrepancy between the probability to observe two alleles at
456 two loci together and the probability to observe them independently.

457 **Local ancestry inference (LAI):** the process that determines the ancestral origin of genomic

458 segments along the genome in admixed individuals.

459 **Logistic regression:** a statistical model that estimates the probability that a given input belongs

460 to one of two categories using a logistic function.

461 Long-term balancing selection: a form of natural selection that maintains ancestral genetic

462 variants over long evolutionary timescales, even across speciation events, resulting in trans-

463 species polymorphisms.

464 Machine learning: an algorithmic approach that automatically learns patterns or structures from465 data to make predictions or decisions.

466 Positive selection: a form of natural selection that increases the frequency of beneficial467 mutations.

468 Self-supervised learning: a machine learning paradigm that learns patterns or structures from469 data through automatically generated labels derived from the data itself.

470 Semantic segmentation: a machine learning task that determines the category of each individual471 pixel in an image.

472 Unsupervised learning: a machine learning paradigm that learns patterns or structures from data473 without using labeled examples.

474 Variational autoencoder: a deep generative model that learns a probabilistic latent

475 representation of data by combining artificial neural networks with variational inference.

476

477 Acknowledgements

478 X.H. thanks Jie Wang for helpful discussions during the preparation of this manuscript. This

- 479 project has been funded by the Vienna Science and Technology Fund (WWTF)
- 480 [10.47379/VRG20001] to M.K.
- 481

482 **Resources List**

- 483 ⁱhttps://github.com/millanek/Dsuite
- 484 ⁱⁱhttps://github.com/xin-huang/sai
- 485 ⁱⁱⁱhttps://github.com/xin-huang/sstar
- 486 ^{iv}<u>https://github.com/browning-lab/sprime</u>
- 487 vhttps://github.com/simonhmartin/twisst2
- 488 ^{vi}https://github.com/PrincetonUniversity/IBDmix
- 489 viihttps://degiorgiogroup.fau.edu/vf.html
- 490 viii<u>http://compgen.cshl.edu/ARGweaver/doc/argweaver-d-manual.html</u>
- 491 ^{ix}https://github.com/BenjaminPeter/admixfrog
- 492 *<u>https://github.com/Shuhua-Group/ArchaicSeeker2.0</u>
- 493 xihttps://github.com/jesvedberg/Ancestry_HMM-S
- 494 xii<u>https://github.com/genicos/ahmm_mls</u>
- 495 xiii<u>https://github.com/LauritsSkov/Introgression-detection</u>
- 496 xiv<u>https://github.com/sriramlab/ArchIE</u>
- 497 ^{xv}<u>https://github.com/kr-colab/FILET</u>
- 498 ^{xvi}https://github.com/xzhang-popgen/maladapt
- 499 ^{xvii}<u>https://github.com/YuboZhangPKU/ERICA</u>
- 500 ^{xviii}<u>https://github.com/grahamgower/genomatnn</u>
- 501 xix<u>https://github.com/SchriderLab/introNets</u>
- 502 xx<u>https://iclr.cc/virtual/2024/invited-talk/21799</u>
- 503
- 504 **Declaration of interests**
- 505 The authors declare no conflict of interests.
- 506
- 507 Declaration of generative AI in scientific writing

508 During the preparation of this work the authors used ChatGPT in order to edit the language. After

using this tool/service, the authors reviewed and edited the content as needed and take full

510 responsibility for the content of the published article.

511

512 **References**

- Hibbins, M.S. and Hahn, M.W. (2022) Phylogenomic approaches to detecting and
 characterizing introgression. *Genetics* 220, iyab173.
- 515 2. Ongaro, L. and Huerta-Sanchez, E. (2024) A history of multiple Denisovan introgression
 516 events in modern humans. *Nat. Genet.* 56, 2612–2622.
- 517 3. Zhang, W. *et al.* (2024) Uncovering ghost introgression through genomic analysis of a
 518 distinct eastern Asian hickory species. *Plant J.* 119, 1386–1399.
- 4. Wang, G. *et al.* (2025) Genomic evidence for hybridization and introgression between blue
 peafowl and endangered green peafowl and molecular foundation of leucistic plumage of
 blue peafowl. *GigaScience* 14, giae124.
- 5. North, H.L. *et al.* (2024) Rapid adaptation and interspecific introgression in the North
 American crop pest *Helicoverpa zea. Mol. Biol. Evol.* 41, msae129.
- 6. Patterson, N. *et al.* (2012) Ancient admixture in human history. *Genetics* 192, 1065–1093.
- 525 7. Martin, S.H. *et al.* (2015) Evaluating the use of ABBA-BABA statistics to locate introgressed
 526 loci. *Mol. Biol. Evol.* 32, 244–257.
- 527 8. Pfeifer, B. and Kapan, D.D. (2019) Estimates of introgression as a function of pairwise
 528 distances. *BMC Bioinform*. 20, 207.
- 529 9. Fang, L.L. *et al.* (2024) Leveraging shared ancestral variation to detect local introgression.
 530 *PLoS Genet.* 20, e1010155.
- 531 10. Malinsky, M. *et al.* (2015) Genomic islands of speciation separate cichlid ectomorphs in an
 532 East African crater lake. *Science* 350, 1493–1498.
- 533 11. Racimo, F. *et al.* (2017) Signatures of archaic adaptive introgression in present-day human
 534 populations. *Mol. Biol. Evol.* 34, 296–317.
- 535 12. Plagnol, V. and Wall, J.D. (2006) Possible ancestral structure in human populations. *PLoS*536 *Genet.* 2, e105.

- 537 13. Vernot, B. and Akey, J.M. (2014) Resurrecting surviving Neanderthal lineages from modern
 538 human genomes. *Science* 343, 1017–1021.
- 539 14. Vernot, B. *et al.* (2016) Excavating Neanderthal and Denisovan DNA from the genomes of
 540 Melanesian individuals. *Science* 352, 235–239.
- 541 15. Browning, S.R. *et al.* (2018) Analysis of human sequence data reveals two pulses of archaic
 542 Denisovan admixture. *Cell* 173, 53–61.
- 543 16. Martin, S.H. and van Belleghem, S.M. Exploring evolutionary relationships across the
 544 genome using topology weighting. *Genetics* 206, 429–438.
- 545 17. De-Kayne, R. et al. (2024) Incomplete recombination suppression fuels extensive haplotype
- 546 diversity in a butterfly color pattern supergene. *bioRixv*, Published online July 26, 2024.
 547 https://doi.org/10.1101/2024.07.26.605145
- 548 18. Huang, X. *et al.* (2022) sstar: a Python package for detecting archaic introgression from
 549 population genetic data with *S**. *Mol. Biol. Evol.* 39, msac212.
- 550 19. Chen, L. *et al.* (2020) Identifying and interpreting apparent Neanderthal ancestry in African
 551 individuals. *Cell* 180, 677–687.e16.
- 552 20. Setter, D. *et al.* (2020) VolcanoFinder: genomic scans for adaptive introgression. *PLoS Genet.*553 16, e1008867.
- 554 21. Moest, M. *et al.* (2020) Selective sweeps on novel and introgressed variation shape mimicry
 555 loci in a butterfly adaptive radiation. *PLoS Biol.* 18, e3000597.
- 22. Romieu, J. *et al.* (2024) Performance evaluation of adaptive introgression classification
 methods. *bioRixv*, Published online June 14, 2024.

558 https://doi.org/10.1101/2024.06.12.598278

- 559 23. Nielsen, R. *et al.* (2025) Inference and applications of ancestral recombination graphs. *Nat.*560 *Rev. Genet.* 26, 47–58.
- 561 24. Wong, Y. *et al.* (2024) A general and efficient representation of ancestral recombination
 562 graphs. *Genetics* 228, iyae100.
- 563 25. Hubisz, M.J. *et al.* (2020) Mapping gene flow between ancient hominis through demographyaware inference of the ancestral recombination graph. *PLoS Genet.* 16, e1008895.
- 565 26. Bishop, C.M. (2006) *Pattern Recognition and Machine Learning*. Springer.
- 566 27. Sun, Q. *et al.* (2025) Opportunities and challenges of local ancestry in genetic association
- **567** analyses. *Am. J. Hum. Genet.* 112, 727–740.

- 568 28. Sankararaman, S. (2020) Methods for detecting introgressed archaic sequences. *Curr. Opin.*569 *Genet. Dev.* 62, 85–90.
- 570 29. Skov, L. *et al.* (2018) Detecting archaic introgression using an unadmixed outgroup. *PLoS*571 *Genet.* 14, e1007641.
- 30. Peter, B.M. (2020) 100,000 years of gene flow between Neandertals and Denisovans in the
 Altai mountains. *bioRixv*, Published online, March 15, 2020.
- 574 https://doi.org/10.1101/2020.03.13.990523
- 575 31. Yuan, K. *et al.* (2021) Refining models of archaic admixture in Eurasia with ArchaicSeeker
 576 2.0. *Nat. Commun.* 12, 6232.
- 577 32. Svedberg, J. *et al.* (2021) Inferring adaptive introgression using hidden Markov model. *Mol.*578 *Biol. Evol.* 38, 2152–2165.
- 33. Ayala, N.M. *et al.* (2023) Inferring multi-locus selection in admixed populations. *PLoS Genet.* 19, e1011062.
- 34. Huang, X. *et al.* (2024) Harnessing deep learning for population genetic inference. *Nat. Rev. Genet.* 25, 61–78.
- 583 35. Schrider, D.R. and Kern, A.D. (2018) Supervised machine learning for population genetics: a
 new paradigm. *Trends. Genet.* 34, 301–312.
- 585 36. Durvasula, A. and Sankararaman, S. (2019) A statistical model for reference-free inference of
 archaic local ancestry. *PLoS Genet.* 15, e1008175.
- 587 37. Schrider, D.R. *et al.* (2018) Supervised machine learning reveals introgressed loci in the
 588 genomes of *Drosophila simulans* and *D. sechellia. PLoS Genet.* 14, e1007341.
- 38. Zhang, X. *et al.* (2023) *MaLAdapt* reveals novel targets of adaptive introgression from
 Neanderthals and Denisovans in worldwide human populations. *Mol. Biol. Evol.* 40,
 msad001.
- 592 39. Gower, G. *et al.* (2021) Detecting adaptive introgression in human evolution using
 593 convolutional neural networks. *eLife* 10, e64669.
- 40. Zhang, Y. *et al.* (2023) Inferring historical introgression with deep learning. *Syst. Biol.* 72, 1013–1038.
- 41. Ray, D.D. *et al.* (2024) IntroUNET: identifying introgressed alleles via semantic
 segmentation. *PLoS Genet.* 20, e1010657.

- 598 42. Velazquez-Archelay, K. *et al.* (2023) Archaic introgression shaped human circadian traits.
 599 *Genome Biol. Evol.* 15, evad203.
- 43. Cheng, J.Y. *et al.* (2021) Detecting selection in multiple populations by modeling ancestral
 admixture components. *Mol. Biol. Evol.* 39, msab294.
- 44. Huang, X. *et al.* (2021) Dissecting dynamics and differences of selective pressures in the
 evolution of human pigmentation. *Biol. Open* 10, bio056523.
- 45. Hao, S. *et al.* (2020) A brief survey on semantic segmentation with deep learning. *Neurocomputing* 406, 302–321.
- 46. Wang, T. *et al.* (2023) Interploidy introgression shaped adaptation during the origin and
 domestication history of *Brassica napus*. *Mol. Biol. Evol.* 40, msad199.
- 47. Stone, B.W. and Wessinger, C.A. (2024) Ecological diversification in an adaptive radiation of
 plants: the role of *de novo* mutation and introgression. *Mol. Biol. Evol.* 41, msae007.
- 610 48. Coughlan, J.M. *et al.* (2022) Patterns of population structure and introgression among
- 611 recently differentiated *Drosophila melanogaster* populations. *Mol. Biol. Evol.* 39, msac223.
- 612 49. Feng, X. *et al.* (2024) Secondary contact, introgressive hybridization, and genome
 613 stabilization in sticklebacks. *Mol. Biol. Evol.* 41, msae031.
- 614 50. Kuhlwilm, M. *et al.* (2019) Ancient admixture from an extinct ape lineage into bonobos. *Nat.*615 *Ecol. Evol.* 3, 957–965.
- 616 51. Evans, B.J. *et al.* (2021) Mitonuclear interactions and introgression genomics of macaque
 617 monkeys (*Macaca*) highlight the influence of behaviour on genome evolution. *Proc. R. Soc.*618 *B Biol. Sci.* 288, 20211756.
- 52. Jensen, A. *et al.* (2023) Complex evolutionary history with extensive ancestral gene flow in
 an African primate radiation. *Mol. Biol. Evol.* 40, msad247.
- 621 53. Pawar, H. *et al.* (2023) Ghost admixture in eastern gorillas. *Nat. Ecol. Evol.* 7, 1503–1514.
- 54. Lyu, Y. *et al.* (2024) Recent selection and introgression facilitated high-altitude adaptation in
 cattle. *Sci. Bull.* 69, 3415–3424.
- 55. van der Valk, T. *et al.* (2024) Comparative genomic analyses provide new insights into
 evolutionary history and conservation genomics of gorillas. *BMC Ecol. Evol.* 24, 14.
- 56. Malinsky, M. *et al.* (2021) Dsuite Fast D-statistics and related admixture evidence from
 VCF files. *Mol. Ecol. Resour.* 21, 584–595.

- 57. Huang, X. *et al.* (2025) SAI: a Python package for statistics for adaptive introgression. *bioRixy*, Published online April 22, 2025. https://doi.org/10.1101/2025.04.19.649497
- 58. Vilgalys, T.P. *et al.* (2022) Selection against admixture and gene regulatory divergence in a
 long-term primate field study. *Science* 377, 635–641.
- 59. Wang, M.S. *et al.* (2022) A polar bear paleogenome reveals extensive ancient gene flow from
 polar bears into brown bears. *Nat. Ecol. Evol.* 6, 936–944.
- 634 60. Puckett, E.E. (2025) Phylogeography of introgression: spatial and temporal analyses identify
 635 two introgression events between brown and American black bears. *Heredity*, Published
 636 online April 19, 2025. https://doi.org/10.1038/s41437-025-00762-0
- 637 61. Foote, A.D. *et al.* (2019) Killer whale genomes reveal a complex history of recurrent
 638 admixture and vicariance. *Mol. Ecol.* 28, 3427–3444.
- 639 62. Wang, M.S. *et al.* (2020) Ancient hybridization with an unknown population facilitated high640 altitude adaptation of canids. *Mol. Biol. Evol.* 37, 2616–2629.
- 641 63. van der Valk, T. *et al.* (2021) Million-year-old DNA sheds light on the genomic history of
 642 mammoths. *Nature* 591, 265–269.
- 643 64. Huerta-Sánchez, E. *et al.* (2014) Altitude adaptation in Tibetans caused by introgression of
 644 Denisovan-like DNA. *Nature* 512, 194–197.
- 645 65. Segawa, T. *et al.* (2024) The origins and diversification of Holarctic brown bear populations
 646 inferred from genomes of past and present populations. *Proc. R. Soc. B* 291, 20232411.
- 647 66. Münger, X. *et al.* (2024) Facilitated introgression from domestic goat into Alpine ibex at
 648 immune loci. *Mol. Ecol.* 33, e17429.
- 649 67. Fontsere, C. *et al.* (2022) Population dynamics and genetic connectivity in recent chimpanzee
 650 history. *Cell Genom.* 2, 100133.
- 68. Peng, Y. *et al.* (2023) Distinct traces of mixed ancestry in western commercial pig genomes
 following gene flow from Chinese indigenous breeds. *Front. Genet.* 13, 1070783.
- 653 69. Nekrutenko, A. *et al.* (2018) Biology needs evolutionary software tools: let's build them
 654 right. *Mol. Biol. Evol.* 35, 1372–1375.
- 655 70. Huang, X. (2024) Developing machine learning applications for population genetic
- 656 inference: Ensuring precise terminology and robust implementation. *EcoEvoRixv*, Pulished
- 657 online September 20, 2024. https://doi.org/10.32942/X2N90M

- 658 71. Zhou, Y. and Browning, S.R. (2021) Protocol for detecting introgressed archaic variants with
 659 SPrime. *STAR Protoc.* 2, 100550.
- 660 72. Tournebize, R. and Chikhi, L. (2025) Ignoring population structure in hominin evolutionary
 661 models can lead to the inference of spurious admixture events. *Nat. Ecol. Evol.* 9, 225–236.
- 662 73. Durvasula, A. and Sankararaman, S. (2020) Recovering signals of ghost archaic introgression
 663 in African populations. *Sci. Adv.* 6, eaax5097.
- 664 74. Fan, S. *et al.* (2023) Whole-genome sequencing reveals a complex African population
 665 demographic history and signatures of local adaptation. *Cell* 186, 923–939.e14.
- 666 75. Ragsdale, A.P. *et al.* (2023) A weakly structured stem for human origins in Africa. *Nature*667 617, 755–763.
- 668 76. Hedrick, P.W. (2013) Adaptive introgression in animals: examples and comparison to new
 669 mutation and standing variation as sources of adaptive variation. *Mol. Ecol.* 22, 4606–4618.
- 670 77. Hackl, J. and Huang, X. (2025) Revisiting adaptive introgression at the HLA genes in
 671 Lithuanian genomes with machine learning. *Infect. Genet. Evol.* 127, 105708.
- 672 78. Mo, Z. and Siepel, A. (2023) Domain-adaptive neural networks improve supervised machine
 673 learning based on simulated population genetic data. *PLoS Genet.* 19, e1011032.
- 674 79. Azodi, C.B. *et al.* (2020) Opening the black box: interpretable machine learning for
 675 geneticists. *Trends Genet.* 36, 442–455.
- 80. Duranton, M. and Pool, J.E. (2022) Interactions between natural selection and recombination
 shape the genomic landscape of introgression. *Mol. Biol. Evol.* 39, msac122.
- 81. Glasenapp, M.R. and Pogson, G.H. (2024) Selection shapes the genomic landscape of
 introgressed ancestry in a pair of sympatric sea urchin species. *Genome Biol. Evol.* 16,
 evae124.
- 82. Neloy, A.A. and Turgeon, M. (2024) A comprehensive study of auto-encoders for anomaly
 detection: efficiency and trade-offs. *Mach. Learn. Appl.* 17, 100572.
- 83. Benegas, G. *et al.* (2025) Genomic language models: opportunities and challenges. *Trends Genet.* 41, 286–302.
- 84. Consens, M.E. *et al.* (2025) Transformers and genome language models. *Nat. Mach. Intell.* 7,
 346–362.

- 687 85. Zhou, C. *et al.* (2024) A comprehensive survey on pretrained foundation models: a history
- from BERT to ChatGPT. Int. J. Mach. Learn. & Cyber., Published online November 24,
- 689 2024. https://doi.org/10.1007/s13042-024-02443-6