Tensor Decomposition-based Feature Extraction and Classification to Detect Natural Selection from Genomic Data

Md Ruhul Amin,* Mahmudul Hasan, Sandipan Paul Arnab (D), and Michael DeGiorgio (D)*

Department of Electrical Engineering and Computer Science, Florida Atlantic University, Boca Raton, FL 33431, USA

*Corresponding authors: E-mails: aminm2021@fau.edu; mdegiorg@fau.edu

Abstract

Inferences of adaptive events are important for learning about traits, such as human digestion of lactose after infancy and the rapid spread of viral variants. Early efforts toward identifying footprints of natural selection from genomic data involved development of summary statistic and likelihood methods. However, such techniques are grounded in simple patterns or theoretical models that limit the complexity of settings they can explore. Due to the renaissance in artificial intelligence, machine learning methods have taken center stage in recent efforts to detect natural selection, with strategies such as convolutional neural networks applied to images of haplotypes. Yet, limitations of such techniques include estimation of large numbers of model parameters under nonconvex settings and feature identification without regard to location within an image. An alternative approach is to use tensor decomposition to extract features from multidimensional data although preserving the latent structure of the data, and to feed these features to machine learning models. Here, we adopt this framework and present a novel approach termed T-REx, which extracts features from images of haplotypes across sampled individuals using tensor decomposition, and then makes predictions from these features using classical machine learning methods. As a proof of concept, we explore the performance of T-REx on simulated neutral and selective sweep scenarios and find that it has high power and accuracy to discriminate sweeps from neutrality, robustness to common technical hurdles, and easy visualization of feature importance. Therefore, T-REx is a powerful addition to the toolkit for detecting adaptive processes from genomic data.

Key words: candecomp/parafac, tensor decomposition, positive natural selection, dimensionality reduction.

Introduction

Natural selection refers to the evolutionary processes that differentially affect the number of offspring organisms may leave in the next generation based on the fitness of particular traits in an environment (Gillespie 2004). As traits will typically have some genetic basis, changes in the frequencies of traits in the population will also influence frequencies of genetic variants, or alleles, that contribute to these traits. Specifically, positive natural selection is the process by which beneficial traits increase in frequency in a population, leading to increases in the frequencies of alleles coding for the traits they contribute to, and ultimately a decrease in genetic variation at the locus under selection (Gillespie 2004). Because positive selection may cause particular alleles to rapidly rise in frequency in a population, through the process of genetic hitchhiking neutral genetic variants at sites nearby the selected locus will also rise to high frequency with it (Maynard Smith and Haigh 1974; Przeworski 2002; Kim and Nielsen 2004; Hermisson and Pennings 2017). This indirect influence of positive selection on neighboring sites causes a loss of neutral genetic variation, resulting in the phenomenon coined as selective sweep (Hermisson and Pennings 2005; Pennings and Hermisson 2006a, 2006b).

Inferences of such selective sweep events have been important for learning about a number of traits, such as how some human populations have evolved to digest lactose after infancy due to the advent of agriculture (Tishkoff et al. 2007; Field et al. 2016; Ségurel and Bon 2017; Taliun et al. 2021), the ability of organisms to survive at extreme environments such as high altitudes (Beall et al. 2010; Bigham et al. 2010; Simonson et al. 2010; Yi et al. 2010; Peng et al. 2011; Wang et al. 2011; Xu et al. 2011; Huerta-Sánchez et al. 2013, 2014; Zhang et al. 2014; Wei et al. 2016; Lindo et al. 2018; Graham and McCracken 2019; Liu et al. 2019; Szpiech et al. 2021; Zhang et al. 2021), and the rapid spread of certain viral variants that require societies to regularly generate new drugs and vaccines (Rambaut et al. 2008; Bedford et al. 2011; Feder et al. 2016, 2021; Kim and Kim 2016; Kang et al. 2021). These important applications to human and other study systems have fueled significant interest in detecting sweeps among evolutionary, ecological, anthropological, and epidemiological researchers over the last several decades. Initial efforts toward identifying signatures of selective sweeps from genetic data were with summary statistics, which classically explored deviations from expected genetic variation under simple models of neutrality. Such

Received: March 02, 2023. Revised: August 10, 2023. Accepted: September 14, 2023 © The Author(s) 2023. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution. This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (https:// creativecommons.org/licenses/by-nc/4.0/), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

Open Access

approaches have been expanded in recent years, to employ diverse forms of variation, such as haplotype diversity within and among populations to increase both power to detect sweeps and robustness against confounding factors (Sabeti et al. 2002, 2007; Voight et al. 2006; Ferrer-Admetlla et al. 2014; Garud et al. 2015; Harris et al. 2018; Torres et al. 2018; Harris and DeGiorgio 2020b; Szpiech et al. 2021). However, with the growth in computational power and theoretical advances for modeling sweeps, complementary model-based approaches have become ever more common, as they provide a probabilistic approach for detecting sweeps and typically exhibit greater power than summary statistic approaches, provided assumptions of the underlying model fits observed data well enough (Kim and Stephan 2002; Nielsen et al. 2005; Chen et al. 2010; Vy and Kim 2015; DeGiorgio et al. 2016; Huber et al. 2016; Racimo 2016; Lee and Coop 2017; Harris and DeGiorgio 2020a; Setter et al. 2020; DeGiorgio and Szpiech 2022). Yet, these approaches still suffer in that the complexity of scenarios they can model are limited, as they are typically grounded in simple theoretical models for expected genomic variation.

Instead, due to a renaissance in artificial intelligence, machine learning methods have been at the forefront of recent efforts for detecting natural selection events from patterns in genomic variation (Schrider and Kern 2018). A number of approaches employ multiple summary statistics as input features, and differ in the types of summary statistics and the way at which input features are modeled (Lin et al. 2011; Schrider and Kern 2016; Sheehan and Song 2016; Kern and Schrider 2018; Sugden et al. 2018; Mughal and DeGiorgio 2019; Mughal et al. 2020; Lauterbur et al. 2022; Arnab et al. 2023). Because the summary statistics target different patterns of genetic variation, the ensemble of such statistics can be used to provide cumulative evidence for, or against, the probability of a selective sweep producing the set of summary statistic values. Importantly though, these machine learning approaches require that hand-engineered summary statistics are chosen in advance, when they may not necessarily be the best features for discriminating among diverse evolutionary events. As a complementary strategy concurrent with the rise of deep learning (LeCun et al. 2015), convolutional neural networks (CNNs; LeCun et al. 1998) have been recently employed as a mechanism to automatically extract features and detect sweeps from raw genotypic variation (Chan et al. 2018; Flagel et al. 2019; Torada et al. 2019; Gower et al. 2021; Isildak et al. 2021). To use CNNs as a way to extract features and detect selective sweeps, the genomic region has to be represented as images, and such approaches have matched or outperformed other statistical frameworks (Kern and Schrider 2018; Flagel et al. 2019; Torada et al. 2019; Isildak et al. 2021).

CNNs are powerful tools that have proven useful in image classification and deep learning tasks (LeCun et al. 1998; Gu et al. 2018). Despite their robustness, they may suffer some limitations for detecting sweeps. Because the majority of CNN architectures have at least one fully-

connected dense hidden layer prior to the output layer, such models often have an enormous number of parameters (Goodfellow et al. 2016). The increased number of parameters generally requires larger training sets to learn their parameters, and the computational complexity of finding the optimal parameters is often high. Moreover, CNN architectures are typically agnostic with respect to where in an input image an object to detect is located, thereby ignoring important information when detecting selective sweeps, as haplotype diversity should be altered nearby a selected locus (e.g. Hermisson and Pennings 2005; Pennings and Hermisson 2006a, 2006b) and support for a sweep centered on a particular genomic location should change depending on whether the altered diversity is at the center or periphery of the image. Instead, it may be useful to employ techniques that automatically extract features from images whereas retaining the spatial location within the image of important features, and to then use these features as input to the many powerful linear and nonlinear machine learning methods that have been developed (Hastie et al. 2009). One such approach for extracting features from image data is tensor decomposition (Kolda and Bader 2009).

Tensor decomposition is a class of dimensionality reduction techniques that can be applied to extract important features from data that has higher-order structure (Kolda and Bader 2009). Data with higher-order structure differs from typical data that is collected as a vector of feature values, as the feature values are organized in a specific manner. For example, image data has higher-order structure, as pixel (feature) values are organized into rows and columns, with pixels tending to have similar values if they have similar row-column coordinates. Traditional data analysis methods need higher-order data to be flattened into a vector for each observation before it can be analyzed. Moreover, this flattening procedure runs the risk of erasing information that might be encoded within the higher-order structure of the data. In situations where it is important to maintain the integrity of the structure of such higher-order structured data, tensor decomposition can be a useful tool for embedding this higher-order structured data in a low-dimensional space although retaining the information encoded in the original data. Tensor decomposition when applied to higher-order data can extract features, which in turn can be used for prediction tasks such as classification.

Additionally, working with high-dimensional data containing enormous numbers of features comes with an increased computational cost for a predictive model, which sometimes is referred to as "curse of dimensionality" (Bellman 1966). Most nonlinear methods suffer more from this curse of dimensionality than linear methods, as nonlinear methods involve a large number of parameters (Verleysen and François 2005). To circumvent this curse of dimensionality issue, dimensionality reduction-based (Salem and Hussein 2019) and ensemble-based methods (Sun et al. 2020) have been developed that operate on vector representations of data, whereas tensor decomposition-based dimensionality

reduction techniques are able to also retain the spatial information of features in data that have higher-order structure (Kolda and Bader 2009).

Feature extraction is one of the foremost steps for classifying data, and tensor decomposition has emerged as an efficient approach to extract a small number of features from high-dimensional data. When extracting features from images of raw genomic data, the curse of dimensionality emerges as a problem for which traditional dimensionality reduction approaches (e.g. principal component analysis) are unideal solutions as they do not retain the spatial structure of the images. Also, many classical machine learning algorithms, such as support vector machines (SVMs), take only feature vectors as input for image data (feature matrices) must be converted into first-order tensors (vectors), which not only compromises the spatial structure of the input data but is also prone to classification errors (Liu 2021).

In this article, we introduce a set of methods termed T-REx (Tensor decomposition-based Robust feature Extraction and classification) that utilize tensor decomposition for automatic feature extraction and classification of genomic image data with an aim toward distinguishing sweep footprints from neutrality. We decompose genomic data obtained from images of haplotypes using CANDECOMP/PARAFAC (CP) decomposition (Carroll and Chang 1970; Harshman 1970), which is a popular model for tensor decomposition. After decomposition, the tensor is expressed as an outer product of three factors, each of which are vectors, resulting in retention of spatial structure. We feed these extracted features as input to classical linear and nonlinear classifiers to predict whether genomic regions represented as images show properties consistent with positive natural selection or neutrality. We also performed an empirical analysis using variant calls from whole-genomes of a central European (CEU) population curated from the 1000 Genomes Project (The 1000 Genomes Project Consortium 2015), in which we found novel candidate sweep genes (e.g. MIR6874, ZNF815P, OCM, and SNHG17) as well as recapitulated prior findings from the literature (e.g. LCT, MCM6, SLC45A2, and EMC7). Finally, we implemented T-REx as open-source software, which is available at https://github.com/RuhAm/T-REx.

Results

The objective of *T-REx* is to automatically extract features from high-dimensional genomic data using tensor decomposition (Kolda and Bader 2009), and to use these features to build a model to detect patterns of adaptation in genomes. To explore the efficacy of *T-REx* for detecting sweeps, we considered a diverse array of factors that can ultimately influence method power, accuracy, and robustness. We first evaluated how machine learning architecture affected accuracy and power, exploring both linear and nonlinear modeling frameworks (Hastie et al. 2009). We then considered how the confounding effects of background selection, nonequilibrium demographic history,

sweep completeness, mutation and recombination rate variability, allele polarization and mispolarization, and missing genomic segments alter relative classification ability. We also directly compared *T-REx* with a leading sweep classifier, *ImaGene* (Torada et al. 2019), which also uses images of haplotype alignments as input. Finally, based on these simulation results, we apply the best strategy to whole-genome sequences from central European human individuals (The 1000 Genomes Project Consortium 2015), and compare our findings with previously reported results from the literature.

Feature Extraction and Model Training

To generate training and testing data for *T-REx*, we created two datasets of varying degrees of constraint associated with them. These datasets are simulated under a constant population size demographic history of 10,000 diploid individuals (Takahata 1993; Excoffier et al. 2013) with the coalescent simulator discoal (Kern and Discoal 2016) using a uniform per-site per-generation mutation rate of 1.25×10^{-8} (Scally and Durbin 2012) and per-site pergeneration recombination rate of 10^{-8} (Payseur and Nachman 2000) drawn from an exponential distribution and truncated at three times the mean (Schrider and Kern 2016). The length of the sequences was set to 1.1 megabases (Mb), and we sampled 200 haplotypes from each simulation under this setting.

In addition to these parameters, to simulate selective sweeps we introduced a beneficial mutation at the center of the simulated sequences and set the per-generation selection coefficient $s \in [0.005, 0.5]$, which was sampled uniformly at random on a logarithmic scale. We set the initial frequency of the beneficial allele at the time of selection to be $f \in [0.001, 0.1]$, which was also sampled uniformly at random on a logarithmic scale. This range for f allowed us to explore both hard and soft sweeps (Hermisson and Pennings 2017). The beneficial mutation became fixed t generations prior to sampling, and we created two datasets based on the distribution of t that are of varying difficulty to discriminate sweeps from neutrality. In the first dataset (denoted by **constant** 1), we set t = 0, and in the second more challenging dataset (denoted by **constant** 2), we draw $t \in [0, 1200]$ uniformly at random, thereby permitting greater overlap between sweep and neutral classes. Using this protocol, we independently generated 10,000 training and 1,000 test observations per class for each dataset. We developed an approach for processing haplotype alignments that may make the structure of input images easier to discern by CP decomposition. Full details of this alignment processing strategy are provided in the Methods.

For each dataset (<code>constant_1</code> or <code>constant_2</code>), using the <code>rTensor</code> package (Li et al. 2018), we performed a rank R CP tensor decomposition across a set of 20,000 training observations (10,000 per class) to obtain a low-dimensional representation of the observations in R-dimensional space. Using Equation (3) in the <code>Methods</code>

section, we projected the 2,000 (1,000 per class) test observations of processed image alignments onto the R-dimensional subspace learned from the training set. The CP tensor decomposition subsection of the Methods provides a detailed overview of CP tensor decomposition, including learning the low-dimensional representation of the training set and projection of the test observations onto this subspace. Identifying an appropriate rank or number of components (R) is a key task for performing CP decomposition, yet an exact algorithm does not exist for finding the optimum R that gives the best approximation to the original tensor (Kolda and Bader 2009). Because the performances of our classifiers vary greatly across different ranks, we evaluated different values of rank $R \in$ {50, 100, 150, 200, 250, 300} until we identified a rank that yielded excellent power and accuracy whereas remaining computationally efficient.

After extracting the factor matrices **A**, **B**, and **C** upon performing CP tensor decomposition, we fed the extracted features from the factor matrix **A** (details are provided in the *Methods*) into both classical linear (elastic net [EN] logistic regression) and nonlinear (SVM with a radial basis kernel and random forest [RF]) models. We refer to these EN, SVM, and RF algorithms integrated within *T-REx* as *T-REx*(EN), *T-REx*(SVM), and *T-REx*(RF), respectively (details on training each classifier in *Methods* section). The pipeline outlining the overall procedure, from feature extraction via CP tensor decomposition to classification of genomic regions as neutral or sweep, is illustrated in Fig. 1.

Power and Accuracy for Detecting Sweeps

We first evaluate the performance of T-REx under the constant 1 and constant 2 datasets (details are provided above in the Feature extraction and model training subsection of the Results) across different CP decomposition ranks $R \in \{50, 100, 150, 200, 250, 300\}$. We selected the model resulting from the best-performing rank for each of the methods based on the smallest crossvalidation loss across the ranks (supplementary Figs. S1 and S2, Supplementary Material online). We find that across different ranks, T-REx(EN) has the lowest error among the three methods and T-REx(RF) showed lower error than T-REx(SVM). For the constant 1 dataset, T-REx(EN) achieves an accuracy of 93.15% and maintains relatively balanced classification rates across neutral and sweep settings, with a slight, yet conservative skew toward prediction of neutrality (Fig. 2). T-REx(SVM) and T-REx(RF) have lower accuracies (87.15 and 89.70%, respectively), with T-REx(SVM) reaching 98.2% accuracy on neutral settings (Fig. 2). For the more challenging constant 2 dataset, T-REx(EN) attains accuracy of 91.55% with high classification accuracies for both sweep and neutral scenarios, and with minimal misclassification of neutral regions as sweeps. Upon a closer look at the classification rates, we find that T-REx(SVM) has a high accuracy of 97.0% on neutral settings, but suffers from greater sweep misclassification than T-REx(EN) (Fig. 2). The high power displayed by the receiver operating characteristic (ROC) curves echos the high accuracy evidenced by the confusion matrices, showing that *T-REx*(EN) has high true positive rates for low false positive rates (Fig. 2).

By comparing our methods to the CNN-based classifier *ImaGene* (Torada et al. 2019), we find that *T-REx*(EN) surpasses *ImaGene* in terms of power, accuracy, and classification balance on both datasets (Fig. 2). However, *ImaGene* outperforms *T-REx*(SVM) in terms of power, accuracy, and classification balance whereas *T-REx*(RF) has slightly more balanced classification rates than *ImaGene*. Though all methods of *T-REx* and *ImaGene* mistakes sweeps for neutrality more often than *T-REx*(EN) (Fig. 2), which drives the lower accuracy and power of *ImaGene* relative to *T-REx*(EN).

Robustness to Background Selection

We have shown that some of our T-REx models can accurately distinguish selective sweeps from neutrality through patterns of lost genomic diversity. However, a pervasive force acting in genomes is negative selection, which imposes long-term selective constraint on functional genomic regions such as genes (Loewe 2008), and reduces genetic diversity at selected sites, much that like of positive selection, by removing deleterious alleles from a population. Moreover, similar to positive selection, alleles at nearby neutral loci are also purged from the population in a manner akin to hitchhiking for sweeps through a process termed background selection (Charlesworth et al. 1995; Comeron 2014; Charlesworth and Jensen 2021). The effects of background selection on variation across the genome in diverse sets of lineages have been reported (e.g. McVicker et al. 2009; Comeron 2014), studies have shown that distributions of allelic diversity under background selection may resemble those of sweeps (Charlesworth et al. 1993, 1995, 1997; Keinan and Reich 2010; Seger et al. 2010; Nicolaisen and Desai 2013; Huber et al. 2016), and some methods can mistake background selection for selective sweeps (DeGiorgio et al. 2016; Huber et al. 2016). However, more recent evidence suggests that sweeps and background selection leave distinct footprints of genetic variation (Schrider 2020) and that background selection is unlikely to be a problem when using haplotype data (Fagny et al. 2014; Schrider 2020; Lauterbur et al. 2022). Though T-REx employs images of haplotype variation as input and is therefore unlikely to be negatively swayed by background selection, it is nevertheless critical that we demonstrate that T-REx is robust to this common force affecting genomes.

To evaluate whether *T-REx* is misled by genetic variation deriving from background selection, we have simulated 1,000 test replicates with background selection under a constant-size demographic history using the forward-time simulator SLiM (Haller and Messer 2019). Similarly to our training data, for each replicate we allowed the recombination rate to be drawn from an exponential distribution

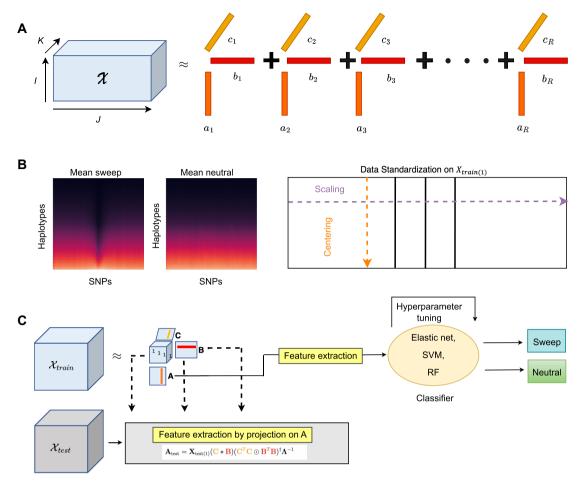
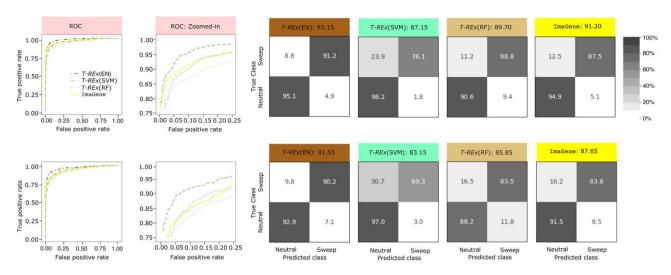


FIG. 1. A) CANDECOMP/PARAFAC (CP) decomposition of three-way training tensor $\mathcal{X} \in \mathbb{R}^{|\mathcal{X}| \times K}$ reduces the tensor into R rank-one composition. nents where $\mathbf{a}_r \in \mathbb{R}^l$, $\mathbf{b}_r \in \mathbb{R}^l$, and $\mathbf{c}_r \in \mathbb{R}^K$ for $r = 1, 2, \ldots, R$. B) Heatmaps illustrate mean images for sweep and neutral class simulations with haplotypes along rows and SNPs along columns, with mean taken across I/2 training observations for each class ($I = N_{\text{train}}$ is the total number of training observations across classes). Each cell of the image is a minor allele frequency value ranging from zero (darker colors) to one (brighter colors) representing the mean number of copies of the minor allele for the haplotype on row $j \in \{1, 2, ..., J\}$ at SNP in column $k \in \{1, 2, ..., K\}$, where the average is taken across overlapping windows during image processing (see Methods). Rows are sorted from top to bottom of the image with increasing L₂-norm taken across the K columns. Therefore, haplotypes toward the top of the image have on average a greater number of SNPs with the major allele than haplotypes toward the bottom. This sorting demonstrates that near the center of the K columns (where selection occurs in sweep simulations), there is a greater number of haplotypes with the major allele (darker colors) at many SNPs. The right figure in panel B) illustrates the standardization process, where the mode-1 unfolded (matricized) data is centered and scaled along the columns and rows, respectively. C) Feature extraction from the training data and the testing data is based on factor matrix A from the CP decomposition. For training data, the matrix **A** is obtained from CP decomposition on the training tensor $\mathcal{X}_{\text{train}}$, whereas the corresponding \mathbf{A}_{test} factor matrix for the test dataset is obtained by projecting the test observations onto this factor learned from the training dataset. This projection is accomplished using the displayed equation, where $\mathbf{X}_{\text{test}(1)}$ is the mode-1 unfolding (matricization) of the tensor $\mathcal{X}_{\text{test}}$, the superscript T denotes transpose, the symbol * denotes the Khatri–Rao product, the \odot symbol denotes the Hadamard (element-wise) product, the superscript † denotes the Moore–Penrose pseudoinverse, and where Λ^{-1} represents the inverse of the diagonal matrix $\Lambda \in \mathbb{R}^{R \times R}$ of scaling terms $\lambda_1, \lambda_2, ..., \lambda_R$. The extracted features are fed to a classifier, which outputs the class predictions.

with mean r and truncated at 3r, but rather than using only $r=10^{-8}$ per site per generation as in our training replicates, we instead considered $r=10^{-8}$, 10^{-9} , or 10^{-10} per site per generation so that we can evaluate background selection in low recombining regions as well as the recombination rates used to train T-REx. These low recombination rates are important to consider, as background selection in such regions can create allele frequencies distributions at long physical distances that might mimic those of sweeps and potentially mislead sweep detectors (DeGiorgio et al. 2016; Huber et al. 2016). To ensure proper simulation burn-in, we allowed a constant-size

population of $N_e = 10^4$ diploid individuals to evolve for $12N_e$ generations, where $10N_e$ generations were devoted to burn-in and a sample of 200 haplotypes were drawn from each replicate after $12N_e$ generations, with mutation rate identical to that of the $constant_1$ and $constant_2$ datasets. Each simulation evolved sequences of length 1.1 Mb and introduced deleterious mutations that are distributed at the center of the sequence within a 55 kb structure that mimics the architecture of a protein-coding gene. The protein-coding gene consists of 50 exons each having length of 100 bases, 49 introns each having a length of 1,000 bases, and 5 ifinmath and



3ifinmath untranslated regions (UTRs), respectively, having lengths 200 and 800 bases, where these lengths approximate the mean values derived from humans (Mignone et al. 2002; Sakharkar et al. 2004). Following Cheng et al. (2017), we set the percentage of deleterious mutations arising within elements of this gene to 75, 10, and 50% for exons, introns, and UTRs, respectively. We drew selection coefficients for deleterious mutations from a gamma distribution with mean of -0.0294 and shape parameter of 0.184 following Schrider and Kern (2017) who based their protocol on the empirical estimates from the African human model of Boyko et al. (2008).

We applied the three *T-REx* classifiers to these simulated test datasets to ascertain what happens to the neutral classification rate in comparison with the neutral classification rate derived from the *constant_1* and *constant_2* test datasets having no background selection. Our results indicate that for models trained using the *constant_1* dataset, for test background selection replicates generated under mean recombination rate of 10⁻⁸ that matches our original simulation protocol, all *T-REx* classifiers show a negative proportional change that signifies increased neutral classification rate under background selection (supplementary Fig. S3, Supplementary Material online). For the same trained models, we find that decreasing the mean recombination rate under background selection to either 10⁻¹⁰ or 10⁻⁹ leads to a positive proportional change in

neutral classification rate, which indicates decreased neutral classification rate under background selection though this elevated misclassification is slight for T-REx(EN) and T-REx(SVM), with proportion of change in neutral classification rates ranging from 0.02 to 0.075 (supplementary Fig. S3, Supplementary Material online). For classifiers trained using the constant 2 dataset, we find that T-REx(EN) under background selection with a mean recombination rate of 10⁻⁸ shows negative proportional change in neutral classification rate similar to what we observed using constant 1 dataset for the same mean recombination rate (supplementary Fig. S3, Supplementary Material online). In concordance with the results for constant 1, we also find that when we reduce the mean recombination rate to 10^{-10} or 10⁻⁹, all methods exhibit decreased neutral classification rates under background selection, with T-REx(EN) and T-REx(SVM) having proportional changes slightly elevated compared with what was observed for the constant 1 dataset (supplementary Fig. S3, Supplementary Material online), which likely stems from the fact that the sweep and neutral distributions used to train the T-REx models overlap more for the constant 2 dataset than for constant 1. In comparison to T-REX(EN) and T-REx(SVM), T-REx(RF) exhibits substantially decreased neutral classification rates under background selection with low mean recombination rates, with proportion of change in neutral classification rate as

MBE

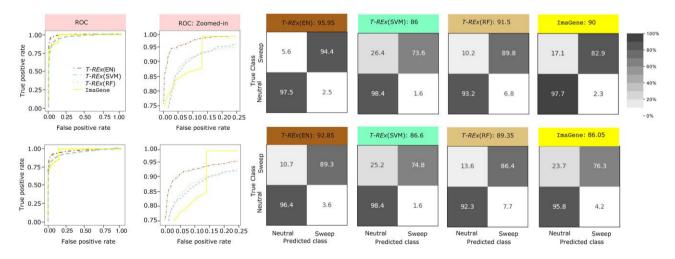
high as 0.20 and 0.27 for the <code>constant_1</code> and <code>constant_2</code> datasets, respectively (supplementary Fig. S3, Supplementary Material online). These results highlight that some of our *T-REx* classifiers are robust against background selection, even under settings of recombination rates that fall outside the domain for which *T-REx* models were trained. We revisit these results in light of neutral simulations under low mean recombination rates within the *Robustness to recombination rate variation* subsection.

Performance under Population Size Changes

The constant-size demographic history underlying the constant 1 and constant 2 datasets is an idealistic model and does not capture the fluctuations in population size often experienced by real populations (Beichman et al. 2018). In particular, demographic scenarios, such as strong and recent population bottlenecks, which lead to an overall loss of haplotypic diversity across the genome as well as an increase in the variance of coalescence times, have been shown to generate false signatures of sweeps as well as reduce the power of sweep detection (Jensen et al. 2005). Therefore, to investigate the performance of T-REx on a nonequilibrium setting with population size fluctuations and a strong, recent population bottleneck, we simulated data under a demographic history inferred (Terhorst et al. 2017) from the central European (CEU) human individuals of the 1000 Genomes Project dataset (The 1000 Genomes Project Consortium 2015).

The distributions that selection parameters were drawn from and the number of simulated replicates per class were identical to the constant-size setting (details regarding the constant-size setting are provided in the *Feature extraction and model training* subsection of the *Results*). Analogous to the two constant-size models, we generated a dataset (denoted by CEU_1) where we set t=0 as well as a second dataset (denoted by CEU_2) representing a more complicated setting where we draw $t \in [0, 1200]$. For each dataset, we consider an array of ranks $R \in \{50, 100, 150, 200, 250, 300\}$ and compared T-REX with the CNN-based sweep classifier ImaGene.

Similar to the evaluation of the two constant-size datasets, we chose the best model through cross-validation, and T-REx(EN) generally showed the lowest error, followed by T-REx(RF) and T-REx(SVM) across different ranks. Among all the methods considered, we find that T-REx(EN) generally has the highest accuracy and power on both the CEU 1 and CEU 2 datasets (Fig. 3). Additionally, T-REx(EN) showed the lowest error in general among the three models selected from their optimal ranks. On either dataset, T-REx(EN) generally exhibits an increase in accuracy with the increase in R, whereas the opposite tendency holds for T-REx(SVM) and T-REx(RF) in which their highest accuracies were attained with a small R value (supplementary Figs. S4 and S5, Supplementary Material online). This trend in the accuracy of T-REx(EN) with increasing R appears to be primarily driven by decreases in the rate of misclassifying sweeps as neutral, leading to



more balanced classification rates. However, *T-REx*(EN) also achieves higher accuracy on neutral settings with increasing *R*, which is desirable as it limits false discovery of sweeps. Finally, we find, as expected, that accuracies for all methods tend to be lower for the more complex *CEU 2* dataset compared with *CEU 1* (Fig. 3).

In general, *T-REx*(EN) and *T-REx*(RF) outperform *ImaGene* for both the *CEU_1* and *CEU_2* datasets in terms of power and accuracy, and *T-REx*(SVM) has similar (on *CEU_2*) or worse (on *CEU_1*) accuracy compared with *ImaGene* due to it incurring higher misclassification rates of sweeps (Fig. 3). Moreover, though *ImaGene* has a low misclassification rate for neutral regions, its overall accuracy suffers due to the high misclassification rate of sweeps as neutral, similar to *T-REx*(SVM). These imbalances in classification rates, however, are conservative as *ImaGene* and *T-REx*(SVM) are not prone to false discovery of sweeps. These results reiterate the strength of CP decomposition to extract features from images, even when prediction is made with a linear model (i.e. *T-REx*(EN)).

The high power of *T-REx*(EN) on the two datasets reflects its strong accuracy evidenced by the confusion matrices, with high true positive rates for low false positive rates (see ROC curves in Fig. 3). We note that *ImaGene* displays a spike in power at a false positive rate of about 15% for both datasets (Fig. 3), which is due to approximately 19% of the predicted sweep probabilities for *ImaGene* being exactly one. The excellent classification performance of *T-REx*(EN) on a complex bottleneck setting (the *CEU_2* dataset) is promising, and so we will apply it to whole-genome data from individuals derived from the same population to scan for sweeps as a proof of concept of our prediction framework (see *Application to human genome variation* subsection of the *Results*).

Feature Maps for Model Interpretability

In addition to its capacity to extract features for prediction problems. CP tensor decomposition provides a low-rank representation of the original tensor, thereby allowing a mechanism for visualizing the spatial components of the images we have collected within our training tensor through factor matrices. These feature maps provide a depiction of the image characteristics that are then fed to classification models. We generated feature maps for R =250 components under the CEU 2 dataset and these feature maps reveal part of the latent structure of the tensor, with the rows and columns of these feature maps representing haplotypes and loci, respectively. Close examination of each of the components (supplementary Figs. S6-S10, Supplementary Material online) reveals gradients in each of the individual feature matrices that represent the separation of features characterized by clusters of similar colors. Though some of the components show gradients in each of the individual feature matrices and clusters of similar colors where we might expect there to be signal in the haplotype alignments to discriminate between sweeps and neutrality, creating a lucid picture of the underlying features is difficult from the set of R =250 images. Moreover, these feature maps only convey information about what characteristics of images were used to separate out observations from the training set, and therefore are not guaranteed to be informative about what characteristics are important for prediction.

To address this issue, we created model-informed feature maps for both datasets through a linear combination of the *R* feature maps, weighing each map by its component's regression coefficient in the trained *T-REx*(EN) model (Fig. 4). Displaying the feature maps in this fashion enables visualization of the characteristics of haplotype

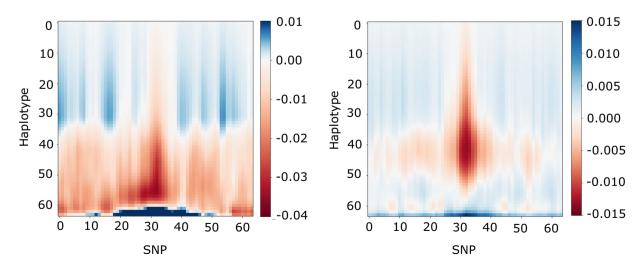


Fig. 4. Model-informed feature maps illustrating emphasis put on different genomic regions of interest by the T-REX(EN) classifier trained to differentiate sweeps from neutrality under a demographic history inferred from CEU humans (Terhorst et al. 2017). Model-informed feature maps were generated through a linear combination of the R feature maps (created using factor matrices R and R) from the training set, where feature map R, R, is weighted by the regression coefficient of component R from a trained logistic regression model with EN penalty. The number of components R0 was selected as in Fig. 3 for the R1-R2-R4. CEU 2 dataset (right panel).

alignments the trained T-REx(EN) models place most emphasis. The pronounced red region around the center of the SNPs alludes to the expected location of lost diversity in sweeps, which the models use to distinguish sweeps from neutrality (Fig. 4). A closer look at the heatmaps suggests that the models place negative weight on these features near the center of the alignment. In contrast, there is also a large dark blue region at the bottom of each heatmap, in which the models place positive emphasis to distinguish sweeps from neutrality. Differences between sweeps and neutrality in this region are expected to be due to the most recent, strongest, and hardest sweeps in our training sets (based on the procedure that we used to process haplotype alignments; see Methods). Another interesting observation we can discern from Fig. 4 is the white, light blue, and light red shading surrounding the dark red region, signifying that T-REx(EN) puts little emphasis on these areas. This lack of emphasis suggests that diversity in this region provides little extra information for discriminating between sweeps and neutrality in the T-REx(EN) model.

Ability to Detect Incomplete Sweeps

We have demonstrated the accuracy and power of *T-REx* under settings where the training and testing were performed on complete sweeps for which the beneficial mutation reached fixation at the time of sampling. However, most realistic scenarios encountered when analyzing empirical data would likely not involve sweeps for which the beneficial allele reached fixation (Burke 2012; Kelly et al. 2013; Ferrer-Admetlla et al. 2014; Vy and Kim 2015; Xue et al. 2021), which can result from a variety of reasons including diminished selective advantage (Pritchard et al. 2010). It is therefore critical that we assess the efficacy of *T-REx* in detecting incomplete sweeps under settings in which the models were trained with complete sweeps.

To evaluate whether our T-REx models trained on complete sweeps have sufficient capacity to detect incomplete sweeps, we simulated an additional 1,000 sweep test replicates with discoal (Kern and Discoal 2016) using idenprotocols for each of the constant 1, constant 2, CEU 1, and CEU 2 datasets with the difference that the beneficial mutation at the time of sampling has frequency 0.5, 0.6, 0.7, 0.8, or 0.9 rather than a frequency of one. As expected, for all four datasets and for each of the three T-REx models, we find that accuracies of detecting incomplete sweeps have an upward trend with increasing frequency of the beneficial mutation at sampling (supplementary Fig. S11, Supplementary Material online). We generally find that T-REx(EN) and T-REx(RF) have higher accuracies than T-REx(SVM) on all four datasets for each frequency of beneficial mutation considered at the time of sampling, with T-REx(EN) showing an edge over T-REx(RF) for all settings aside from when the beneficial allele frequency at the time of sampling is 0.5, which is difficult for all methods considered. Moreover, when the frequency of the beneficial mutation

is 0.9, *T-REX*(EN) shows the highest accuracy among the three *T-REx* models for all datasets, with values ranging from 90 to 97%.

We also considered the power (true positive rate) of T-REx models to detect incomplete sweeps at a 5% false positive rate. Our results show that T-REx(EN) exhibits a similar upward trend in terms of power for all datasets, reaching power in the range from 0.90 to 0.97 across datasets for sweeps to a frequency of 0.9 (supplementary Fig. S12, Supplementary Material online). T-REx(SVM) also demonstrates an upward trend in terms of power for all datasets, reaching values as high as 0.95 (supplementary Fig. \$12C, Supplementary Material online) when the beneficial allele frequency at sampling is 0.9. In contrast to the findings regarding accuracy (supplementary Fig. S11, Supplementary Material online), T-REx(RF) lags in power when compared with T-REx(EN) and T-REx(SVM) with the values reaching only as high as 0.80 (supplementary Fig. S12D, Supplementary Material online). Overall, T-REX models, especially T-REx(EN), hold excellent power to detect incomplete sweeps at moderately high frequencies, even though they were trained to only detect complete sweeps. Training T-REx models with incomplete sweep replicates would likely further improve their accuracies and powers under such settings.

Performance under Mutation Rate Variation

Mutation rate varies across chromosomes and taxa, leading to a variable landscape of genetic diversity within genomes across the tree of life (Bromham 2011; Bromham et al. 2015; Harpak et al. 2016; Bergeron et al. 2023; Danovi 2023). Reductions in polymorphic sites within genomes due to low mutation rates can result in lower haplotype variation, which may mimic signatures of adaptive processes such as selective sweeps. On the other hand, elevated mutation rates can erode footprints of past selective sweeps, making such events more difficult to detect. Thus, it is important to evaluate whether sweep detectors are adversely affected by mutation rate variation.

To evaluate the performance of *T-REx* under mutation rate variation, we simulated 1000 neutral and 1000 sweep test replicates with discoal (Kern and Discoal 2016) using identical protocols for each of the *constant_1*, *constant_2*, *CEU_1*, and *CEU_2* datasets with the difference that mutation rate for a given replicate was drawn uniformly at random within the interval $[\mu/2, 2\mu]$, where $\mu = 1.25 \times 10^{-8}$ per site per generation, instead of a fixed value of $\mu = 1.25 \times 10^{-8}$ per site per generation used to train the *T-REx* classifiers. When applied to these test data, all *T-REx* methods have excellent power and accuracy to detect sweeps, with *T-REx*(EN) outperforming *T-REx*(SVM), and *T-REx*(RF) in terms of accuracy and power (supplementary Fig. S13, Supplementary Material online) as observed from prior experiments.

In terms of correctly classifying neutrally evolving regions, *T-REx*(SVM) has higher accuracy (98.1–99.1%) than *T-REx*(EN) (93.3–98.1%) and *T-REx*(RF) (84.8–95.3%), but

T-REx(EN) has highest overall accuracy on this setting compared with the other two approaches (supplementary Fig. \$13, Supplementary Material online). This increased overall accuracy for T-REx(EN) is driven by its low misclassification rate of sweeps (6.3-14.6%) compared with that of T-REx(SVM) (24.1-33.4%) and T-REx(RF) (12.2-21%) (supplementary Fig. S13, Supplementary Material online). Consistent with its overall high accuracy under mutation rate variation, T-REx(EN) also has substantially higher power than T-REx(SVM) and T-REx(RF) at low false positive rates, exhibiting a quicker ascent to the upper left-hand corner of the ROC curve (supplementary Fig. \$13, Supplementary Material online). Thus, we find that our T-REx models, particularly T-REx(EN), showcase both high power to detect sweeps and robustness to false detection of sweeps under mutation rate variation.

Robustness to Recombination Rate Variation

Recombination rate varies within and between the genomes of different species, and this variable recombination landscape influences patterns of haplotype diversity across genomes (Smukowski and Noor 2011; Cutter and Payseur 2013; Singhal et al. 2015). For instance, genomic regions with low recombination rates may be associated with low haplotype diversity, and thus the observed haplotypic variation may masquerade as sweep signature. In contrast, high recombination rate regions can more quickly eliminate the footprint of lost haplotype diversity, which is characteristic of past selective sweeps, similar to the effects of high mutation. Furthermore, a variety of organisms harbor regions with extreme levels of recombination, which may lead to localized coldspots and hotspots of recombination (Hey 2004; Myers et al. 2005; Galetto et al. 2006; Baudat et al. 2010; Singhal et al. 2015; Booker et al. 2020; Lauterbur et al. 2023). Given the challenges associated with detection of sweeps under recombination rate variation, it is important that we evaluate the relative robustness of our T-REx models to scenarios involving recombination rate variation and to settings with recombination hotspots and coldspots.

To explore the performance of T-REx under recombination rate variation for low recombination regions, we simulated 1,000 neutral test replicates with discoal (Kern and Discoal 2016) using identical protocols for each of the constant 1, constant 2, CEU 1, and CEU 2 datasets with the difference that recombination rate (r) for a given replicate was drawn from an exponential distribution with mean of 10⁻¹⁰ or 10⁻⁹ and truncated at three times the mean. Moreover, to simulate recombination coldspots and hotspots, we also simulated 1,000 neutral test replicates with the coalescent-based simulator msHOT (Hellenthal and Stephens 2007) using identical protocols for each of the constant 1, constant 2, CEU 1, and CEU 2 datasets with the exception that recombination rate r for a replicate was drawn from an exponential distribution with mean of 10⁻⁸ per site per generation and truncated at three times

the mean but with a central 100 kb region of the sequence evolving as r/10 (coldspot) or 10r (hotspot). For each setting, we compared the proportional change in neutral classification rates of T-REx models under recombination rate variation with respect to those under the usual protocols for each of the $constant_1$, $constant_2$, CEU_1 , and CEU_2 neutral test datasets.

We observe that for recombination rate variation with mean rates at one or two orders of magnitude below what T-REx models were trained under, T-REx(EN) and T-REx(SVM) exhibited up to an approximately 10% increase in misclassification of such regions relative to the setting on which the models were trained (supplementary Fig. \$14, Supplementary Material online). However, T-REx(RF) performs comparatively poorly on this setting, with an increase in misclassification error by up to 30% in some cases (supplementary Fig. S14, Supplementary Material online). Notably, for these same mean recombination rates, the methods performed similarly under background selection (supplementary Fig. S3, Supplementary Material online), which highlights that the altered neutral detection rate within regions of low mean recombination rate is likely driven by recombination reducing the diversity of haplotypes rather than a significant influence of background selection in such regions.

For the case of recombination hotspots and coldspots in neutrally evolving regions, the neutral classification rates for all T-REx models are close to those found under test neutral replicates without hotspots or coldspots (supplementary Fig. S15, Supplementary Material online, small magnitude proportional changes). Moreover, all T-REx models actually have improved neutral classification rates under recombination hotspots (supplementary Fig. S15, Supplementary Material online, negative proportional changes). Thus, we find that when our T-REx models are applied to recombination rates of orders of magnitude different from what they were trained, their performance is dependent on the size of the region of altered recombination rate, with smaller regions (e.g. 100 kb) leading to coldspots and hotspots having minimal impact, whereas large regions (e.g. over a megabase) leading to slight changes in robustness of T-REx(EN) and T-REx(SVM) and generally comparatively poorer performance for T-REx(RF).

Effect of Ancestral Allele Polarization and Mispolarization

Knowledge of the ancestral and derived allele at SNPs can often be helpful in detecting natural selection, as these states encode more information than simply assigning alleles as major or minor (Vitti et al. 2015; Bitarello et al. 2023). To perform this polarization of alleles as ancestral and derived in practice, one or more outgroup species is used to establish the likely ancestral and derived (mutant) allelic states. However, incorrect assignment of the ancestral and derived alleles can lead to false signatures of natural selection, and such allele mispolarization becomes more common as more distant outgroups are used to decide on the allelic states (Hernandez et al. 2007). It is

therefore useful to evaluate whether coding alleles as ancestral and derived provides *T-REx* with significant performance gains, whereas also exploring the robustness of such polarization when it is misspecified.

To this end, under the CEU_2 dataset, we performed an experiment in which we coded ancestral alleles as zero and derived alleles as one in place of our original coding of major and minor alleles, and performed all other haplotype alignment processing as in our original experiments. We trained the T-REx(EN) classifier on these new alignments (10,000 observations per class) estimating the two EN hyperparameters and the number of components of tensor decomposition (R=250) through cross-validation. We then applied this trained model to test data (1,000 observations per class) when the haplotype alignments were processed in an identical manner, and when 5% of SNPs were chosen uniformly at random to be mispolarized (i.e. ancestral and derived alleles swapped).

We find out that T-REx(EN) has excellent power and accuracy in identifying sweep signatures when ancestral and derived alleles are used in place of our original coding of major and minor alleles (supplementary Fig. \$16A, Supplementary Material online). In particular, T-REx(EN) is able to correctly classify 97.6 and 88.6% of neutral and sweep observations, respectively (supplementary Fig. \$16A, Supplementary Material online). These results reflect a slight improvement in classification of neutrality although reduced accuracy on sweeps, relative to our original coding as major and minor allele (Fig. 3). To understand how allele mispolarization influences classification accuracy when employing derived and ancestral allele information, we compared the proportional change in neutral and sweep classification rates of T-REx(EN) under allele mispolarization with respect to those under correct polarization. We find that T-REx(EN) exhibits an approximate 44% reduction in correct classification of neutrally evolving regions relative to the setting on which the models were trained (supplementary Fig. S16B, Supplementary Material online). In contrast, we observe a roughly 10% increase in correct classification of sweeps (supplementary Fig. S16B, Supplementary Material online). These results point to allelic mispolarization leading to a skew in more often predicting sweeps, regardless of the true class label, relative to correct polarization. We therefore warrant caution when using ancestral and derived alleles over major and minor alleles within T-REx, as allelic mispolarization may have deleterious effects on model performance.

Robustness to Missing Data

Many genomic regions contain segments with missing SNPs, which may arise due to artifacts in the data, mapping and alignment problems, and sequencing errors. An issue that missing genomic segments poses to methods for detecting sweeps is the problem of false discovery, in which a method erroneously detects a neutrally evolving region as a sweep (Mallick et al. 2009; Mughal and DeGiorgio 2019). These false signals result from the loss of SNPs in missing

segments decreasing haplotypic diversity (see schematic in Fig. 5), which has been shown to mislead some machine learning classifiers to call such neutral regions with confidence as sweeps if such data issues are not accounted for during model training (e.g. Kern and Schrider 2018; Mughal and DeGiorgio 2019). Thus, it is important to demonstrate that *T-REx* not only has high accuracy and power to detect sweeps on idealistic data, but is robust also to common technical artifacts posed by the presence of missing genomic segments.

The haplotype images used for training and testing sets so far have assumed no missing data, and so we seek to examine the effectiveness of our methods when test data have missing segments that may ultimately reduce observed haplotypic variation. To this end, we followed the protocol in Mughal and DeGiorgio (2019) by removing 30% of the SNPs from each test replicate to evaluate the impact of missing data on method accuracy, power, and robustness. The removal of 30% of the SNPs is accomplished in 10 non-intersecting chunks, each accounting for roughly three percent of the total SNPs in the replicate, and with starting position of each chunk chosen uniformly at random. In cases of overlap with previously drawn missing chunks, a new starting location for the current chunk is redrawn.

Using T-REx models trained with nonmissing data and assuming the rank R of CP decomposition that gave each method (T-REx(EN), T-REx(SVM), and T-REx(RF)) their smallest cross-validation loss, we find that on both the CEU 1 and CEU 2 datasets T-REx(EN) continues to show greater power and accuracy compared with competing approaches (center and bottom rows in Fig. 5). Specifically, for both the CEU 1 and CEU 2 datasets, T-REx(EN) outperforms ImaGene with a margin of around 6% in terms of accuracy (center and bottom rows in Fig. 5). Moreover, under both datasets, ImaGene is more prone to false discovery of sweeps than T-REX(EN), as it displays a skew toward falsely classifying neutrally evolving regions as sweeps. In the case of the CEU 1 dataset, T-REx(RF) marginally outperforms ImaGene in terms of accuracy. However, the accuracy of T-REx(SVM) suffers, as 26.4% of sweeps are misclassified (center row in Fig. 5). In contrast, on the CEU 2 dataset, *ImaGene* outperforms both *T-REx*(SVM) and *T-REx*(RF) in terms of accuracy, but falsely classifies 22.5% of the neutral observations as sweeps. This result illustrates that when presented with data containing missing genomic segments, the CNN-based classifier ImaGene may mistake the reduced haplotypic diversity as a sweep footprint. We expand upon this issue in the Discussion section, and detail procedures that can be taken to alleviate the issue of missing segments (e.g. Kern and Schrider 2018).

To further evaluate whether *T-REx* is robust to false discovery of sweeps in neutral regions with missing data, we compute the proportion of false signals, based on the distribution of sweep probabilities of neutral replicates with missing segments, as a function of false positive rate, based on the distribution of sweep probabilities of neutral replicates without missing segments. For this purpose, we

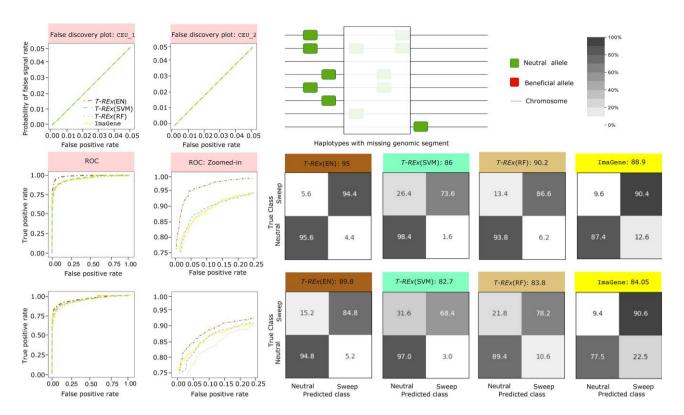


Fig. 5. Powers, accuracies, and robustness to detect sweeps when faced with missing data for the linear *T-REx*(EN) and nonlinear *T-REx*(SVM) and *T-REx*(RF) classifiers in comparison with the CNN-based classifier *ImaGene* under a demographic history inferred from the CEU human population (Terhorst et al. 2017) history using two datasets (*CEU_1* and *CEU_2*) of varying difficulty. For training and testing purposes, the number of observations used for each class was 10,000 and 1,000, respectively where 30% of the total SNPs from each test observation were removed using protocol in Mughal and DeGiorgio (2019). (Top row) Performance of *T-REx* in comparison with *ImaGene* under missing data to ascertain whether the classifiers are robust against false discovery of sweeps, that is, erroneously detecting neutrally evolving regions as sweeps. First and second panel shows probability of false discovery of sweeps when classifying neutral genomic regions containing missing data on the *CEU_1* and *CEU_2* datasets, respectively. Third panel shows how missing genomic segment can masquerade as sweep due to apparent lack of haplotype diversity. (Middle and bottom rows) Powers to detect sweeps of all four methods are compared using receiver operating characteristic (ROC) curves (first column) and ROC curves zoomed in to the upper left-hand corner with false positive rate less than 0.25 and true positive rate greater than 0.75 (second column). Classification accuracy and rates of all four methods are depicted using confusion matrices in columns three through six for *T-REx*(EN), *T-REx*(SVM), *T-REx*(RF), and *ImaGene*, respectively. For both the *CEU_1* and *CEU_2* dataset, *R* = 250, 50, and 50 were chosen for *T-REx*(EN), *T-REx*(SVM), and *T-REx*(RF), respectively, as these ranks yielded the small validation loss on nonmissing data.

generated an additional 1,000 neutral replicates each having 30% missing SNPs so that these two distributions were generated from independent neutral replicates. Sweep classifiers that are robust to neutral missing segments will have the curve relating the proportion of false signals (on the *y*-axis) as a function of the false positive rate (on the *x*-axis) fall on or below the y = x line. Our results show that for both variations of the simulated CEU dataset, curves for all tested methods fall on the y = x line, considering relevant false positive rates between 0 and 5% (top row in Fig. 5). We therefore conclude that all methods considered here are robust to false discovery of sweeps due to missing data when conditioning on reasonable false positive rates.

Application to Human Genome Variation

In addition to evaluating the performance of *T-REx* under simulated scenarios, we also embarked on an empirical application to whole-genome variant calls from a European

human population as a proof of concept (details regarding processing of the empirical data are provided in the Application to empirical data subsection of the Methods). Using the identical protocol as in our assessment of model performance, we trained T-REx(EN) on 10,000 simulated replicates per class with parameters identical to those that generated the CEU_2 dataset, with the exception of sampling 198 haplotypes per simulation to match the 99 diploid individuals sampled for the CEU population of the 1,000 Genomes Project dataset (The 1000 Genomes Project Consortium 2015). We opted to apply T-REx(EN) for our empirical analysis, as it emerged as the best-performing model among the three T-REx methods evaluated across a range of simulated settings.

To uncover candidate genes that show evidence of sweep signatures, we evaluated whether each gene harbored a high predicted sweep probability and a sweep probability peak, observed by computing a moving average computed as the mean of sweep probabilities at 11 contiguous genomic windows. This 11-window mean

Table 1. Autosomal regions showing high predicted sweep probability in the CEU population as predicted by T-REx(EN)

| Chromosome | Start | Stop | Genes |
|------------|-------------|-------------|---|
| 1 | 115,397,483 | 116,311,335 | SYCP1,CASQ2 |
| 1 | 37,940,044 | 38,422,646 | SF3A3, MIR4255 |
| 2 | 136,545,419 | 136,634,013 | LCT, MCM6 |
| 5 | 33,936,490 | 33,984,798 | SLC45A2 |
| 6 | 29,640,259 | 30,594,169 | HLA-F, HLA-F-AS1, IFITM4P, HCG4, HLA-V, HLA-G, HLA-H, HCG4B, HLA-A, HCG9 |
| 7 | 5,751,470 | 6,369,041 | MIR6874, ZNF815P, OCM, CCZ1, RSPH10B |
| 7 | 27,132,611 | 27,287,449 | HOXA1, HOXA2, HOXA3, HOXA9, HOXA10, HOXA-AS2, HOXA-AS3 |
| 10 | 15,253,641 | 15,761,921 | FAM171A1, ITGA8 |
| 15 | 76,507,693 | 77,474,268 | ETFA, TISL2, TYRO3P, SCAPER, RCN2, MIR3713, TSPAN3 |
| 15 | 34,376,217 | 34,649,936 | EMC7, PGBD4, KANTBL1, EMC4, SLC12A6, NUTM1 |
| 15 | 38,988,798 | 41,248,710 | LINC02694, C15orf54, RMDN3, GCHFR, DNAJC17, C15orf62, ZFYVE19, PPP1R14D, SPIT1-AS1, SPIT1, VPS18, LOC105370943, DLL4, CHAC1 |
| 17 | 29,861,900 | 29,902,540 | MIR4724, MIR193A, MIR4725, MIR365B |
| 17 | 41,453,295 | 41,864,988 | LINC00910, ARL4D, MIR2117HG, DHX8, MEOX1, SOST, DUSP3, CFAP97D1 |
| 20 | 37,230,451 | 37,401,163 | ARHGAP40, SLC32A1, ACTR5 |
| 20 | 50,700,549 | 51,266,965 | ZFP64, LINC01524 |
| 20 | 37,049,234 | 37,358,015 | SNHG17, SNORA71B, SNORA71C, SNORA71D, SNORA71E, SNORA60, RALGAPB, ADIG, ARHGAP40, SLC32A1 |
| 22 | 40,139,048 | 40,439,538 | ENTHD1, GRAP2, FAM83F |

approach provides a smoothed representation of the probabilities and helps us observe the underlying trend of probability as a function of genomic position. We identified 17 regions from eight autosomes displaying pronounced peaks in predicted sweep probability, which we list together with associated genes in Table 1 and depict within supplementary Figs. S17 and S18, Supplementary Material online. In particular, we found candidate genes that have been supported by previous studies (e.g. LCT, MCM6, SLC45A2, and EMC7; Bersaglieri et al. 2004; Oleksyk et al. 2010; López et al. 2014; Racimo 2016) as well as novel candidates (e.g. MIR6874, ZNF815P, OCM, and SNHG17).

Sweep Candidates Supported by the Literature

On chromosome 2, we find a peak surrounding the region containing the genes LCT and MCM6 (Fig. 6A). In particular, we see a clear peak that reaches an 11-window mean sweep probability close to one near LCT and MCM6 and decays in value with distance from these genes. This trend of reduction in sweep probability with distance from a putative adaptive locus is consistent with the footprint of a selective sweep, and is due to the action of recombination breaking down linkage disequilibrium and shaping haplotypic diversity across the chromosome (Slatkin 2008). LCT encodes the enzyme lactase that aids in lactose digestion in humans, and is a strong selection candidate, especially across European populations as the ability to digest lactose persists into adulthood within individuals of European ancestry (Scrimshaw and Murray 1988). This lactose tolerance is an outcome of positive selection owing to the advent of farming that resulted in an infusion of milk as part of regular consumption within particular cultures in the last 1,000 years (Sabeti et al. 2006). Moreover, near LCT, we also detect the gene MCM6 with high confidence, which has been hypothesized to have undergone positive selection by previous studies (e.g. Shatin 1968; Bersaglieri

et al. 2004; Nielsen et al. 2005; Harris and Meyer 2006; Sabeti et al. 2007; Tishkoff et al. 2007; Ingram et al. 2009; Itan et al. 2009; Schlebusch et al. 2012; Fan et al. 2016; Cheng et al. 2017). MCM6 contains two introns, one of which harbors an enhancer that acts as a regulatory mechanism for *LCT* and therefore may contribute to lactase persistence and have been positively selected in the past (Anguita-Ruiz et al. 2020).

The region surrounding LCT and MCM6 represents a positive control, as we expect most sweep detection methods to uncover this region with high confidence. We next went on to probe for other well-studied candidates of natural selection, and found evidence for sweeps in the major histocompatibility complex (MHC) region on chromosome 6 (supplementary Fig. S17E, Supplementary Material online). Specifically, T-REx identified high sweep support for the genes HLA-H, HCG4B, HLA-A, and HCG9, which had 11-window mean sweep probabilities close to one. Other candidate genes with moderate support in the region include HLA-F, HLA-F-AS1, IFITM4P, HCG4, HLA-V, and HLA-G, with 11-window mean sweep probabilities ranging from 0.65 to 0.81. Many genes located in the MHC region code for proteins that aid in pathogen immune defense through peptide binding (Mladkova and Kiryluk 2017). Loci in such genes tend to be highly polymorphic, and have long been hypothesized as evolving under balancing selection, likely due to the evolution of the host in the face of pathogens and parasites (Lederberg 1999; Bernatchez and Landry 2003). The high structural variation coupled with extreme polymorphism in this region makes variant calling difficult (Stipoljev et al. 2020), and potentially poor genotype calls may have contributed toward the ambiguity in detecting sweeps in this region. Though often having different genomic footprints to positive selection, balancing selection is a clear deviation from neutrality and T-REx was able to identify the lack of neutrality at the MHC region. The classification of this region

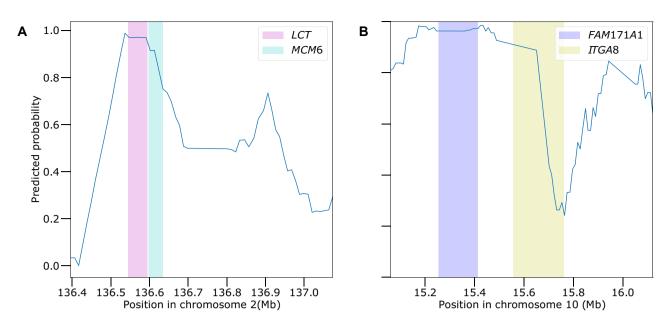


Fig. 6. Detection of two example genomic regions containing sweep signatures within the CEU population of the 1,000 Genomes Project dataset. *T-REx*(EN) predicted sweep probabilities as a function of chromosomal position surrounding the *LCT* and *MCM6* regions on chromosome 2 panel A) and the *FAM171A1* region on chromosome 10 panel B). The probabilities are calculated as 11-window moving averages, computed with five windows before and five windows after a given central window. The genomic intervals containing each gene are shaded using colors in accordance with the order of their appearance in the labels.

as positive selection by *T-REx* may be partially due to its extreme levels of linkage disequilibrium (Stipoljev et al. 2020), consistent with expectations of sweeps. However, our results are also consistent with prior studies, which have found evidence for sweep-like signals at the MHC region in humans (e.g. Campbell et al. 2019).

The gene SLC45A2 (supplementary Fig. Supplementary Material online) on chromosome 5 has moderate sweep support with 11-window mean sweep probabilities around 0.75. This gene encodes a protein that plays a crucial role in melanin synthesis that affects skin pigmentation in humans (López et al. 2014). The frequencies of alleles in this gene that are associated with pigmentation in Europeans demonstrate a latitudinal cline across Europe, resulting in lighter skin pigmentation in northern Europe (Norton et al. 2007). Patterns of variation mimicking footprints of positive selection near SLC45A2 in European humans are supported by numerous studies (e.g. Hider et al. 2013; Laayouni et al. 2014; López et al. 2014; Wilde et al. 2014; Goodwin and de Guzman Strong 2017). A number of the candidate genes identified by T-REx have also been uncovered using ancient DNA studies, which employ additional temporal information on allele frequency trajectories. These candidate genes include LCT (Souilmi et al. 2022), SLC45A2 (Mathieson et al. 2015), and MCM6 (Skoglund and Mathieson 2018).

Further investigation into the regions with high sweep support revealed *EMC7* (supplementary Fig. S18B, Supplementary Material online), which codes for a protein that is an important part of the endoplasmic reticulum membrane and acts as a molecular tether enabling the transport of viruses between different cellular

compartments (Bagchi et al. 2020). *T-REx* detects *EMC7* with an 11-window mean sweep probability of 0.86, which has prior support for positive selection (Racimo 2016). Moreover, with 11-window mean sweep probabilities ranging from 0.95 to 0.99, *T-REx* captured the genomic region containing the protein-coding gene *SF3A3* (supplementary Fig. S17B, Supplementary Material online) on chromosome 1. García-Cárdenas et al. (2022) demonstrated a possible connection between *SF3A3* and breast cancer and a network of cancer-driving genes. Though potentially associated with the harmful disorder of cancer in contemporary environments, Racimo et al. (2014) also suggested that *SF3A3* may have been subjected to past positive selection.

Novel Sweep Candidates

In addition to these previously identified sweep candidates, we uncovered a number of novel candidates. On chromosome 1, we found SYCP1 (supplementary Fig. S17A, Supplementary Material online) as a possible sweep candidate with 11-window mean sweep probabilities reaching 0.88. This protein-coding gene is part of the synaptonemal complex, which is a protein structure that forms between homologous chromosomes (Seo et al. 2016). Hosoya and Miyagawa (2021) highlight that some of the proteins coded by SYCP1 are abnormally expressed in 13 different cancer tissues, including breast and stomach cancer, and acute myelogenous lukemia. Moreover, mutations in SYCP1 have been associated with male infertility (Nabi et al. 2022).

On chromosome 7, we found candidate genes belonging to the HOXA-family that exhibit high sweep support with 11-window mean sweep probabilities ranging from

0.80 to 0.95 (supplementary Fig. S17H, Supplementary Material online). HOXA genes are part of the homeobox cluster, which encode proteins that play an important part in early development of humans by performing embryo segmentation (Shah and Sukumar 2010), and it has been suggested that HOXA-family genes are involved in the inception and development of human cancers (Ge et al. 2021). Specifically, HOXA9 (supplementary Fig. S17H, Supplementary Material online) is responsible for the pathogenesis of acute myelogenous leukemia, which is a cancer of the bloods and bones (Chen et al. 2019).

Additionally, SHNG17 on chromosome 20 (supplementary Fig. S18H, Supplementary Material online) has high sweep support with 11-window mean sweep probabilities reaching 0.95. SHNG17 is known to be an important factor behind gastric cancer in humans, as it is upregulated in gastric cancer tissues (Zhang et al. 2019). Furthermore, on chromosome 10, we identified a strong peak with high sweep support at the protein-coding gene FAM171A1 (supplementary Fig. S17I, Supplementary Material online), which is also associated with breast cancer survival and plays an important role in immune system regulation (Parada et al. 2017). Among our highlighted novel candidates, as well as those that are previously identified (SF3A3), there is an intriguing connection between these sweep candidates and cancer proliferation and suppression. This pattern of selective sweeps at genes related to cancer was also found by other studies that developed machine learning approaches for detecting sweep (e.g. Lou et al. 2014; Schrider and Kern 2017, 2018; Mughal et al. 2020; Arnab et al. 2023). Detection of cancer-related genes by T-REx as well as methods from previous studies, provides an interesting pattern that many past positively selected genes may drive current carcinogenesis in humans.

Discussion

In this article, we have introduced a tensor decompositionbased feature extraction and classification method termed T-REx that is able to differentiate sweeps from neutrality with a high degree of power and accuracy. Specifically, we found that our linear model (T-REx(EN)) demonstrated overall superior performance to the nonlinear models (T-REx(SVM) and T-REx(RF)) across an array of different settings, including demographic history, positive selection regime, and technical artifacts due to missing genomic segments (Figs. 2, 3, and 5). Moreover, in addition to its high power and accuracy to detect sweeps, this modeling framework facilitated easy interpretation of the fitted model by providing feature maps for visualization, which convey the particular location in the haplotype alignments that the models place emphasis when discriminating sweeps from neutrality (Fig. 4).

From our experiments, an unexpected observation was that the linear *T-REx*(EN) model had higher power and accuracy than the nonlinear *T-REx*(SVM) and *T-REx*(RF) models (Figs. 2 and 3). It is possible that the linear model performs better here because it yields a better decision

boundary between the neutral and sweep classes. However, it is more likely that other factors have played a more critical role in leading T-REx(EN) to have the best performance. First, the R components resulting from the CP tensor decomposition are not required to be independent, and may, in fact, be highly correlated (Kolda and Bader 2009). The EN regularization employed by T-REX(EN) has both L_1 - and L_2 -norm penalties, which are both meant to handle correlated features (Hastie et al. 2009). In particular, the L_2 -norm penalty reduces the effective number of features in the model, but encourages a dense model by ensuring that all features remain included in the fitted model (Hastie et al. 2009). In contrast, the L_1 -norm penalty encourages a sparse model by emphasizing fewer features and selecting out those that are redundant or irrelevant for prediction (Hastie et al. 2009). Therefore, the L_1 -norm penalty employed by T-REx(EN) method is particularly useful in reducing the overall dimension of the input data by removing irrelevant and redundant features. This hypothesis is supported by the fact that T-REx(EN)tends to have nondecreasing power and accuracy with increasing R (supplementary Figs. S4 and S5, Supplementary Material online). Second, though T-REx(SVM) also has an L_2 -norm penalty (Hastie et al. 2009), this penalty does not encourage sparsity in the set of input features like the L_1 -norm penalty. Moreover, we employ the radial basis kernel within the T-REx(SVM) classifier, which requires a distance be taken between observations, and distances in high-dimensional space may not behave well due to the curse of dimensionality (Verleysen and François 2005). This hypothesis related to the curse of dimensionality is supported by power and accuracy of T-REx(SVM) tending to diminish with increasing R, and hence has decreasing performance with increasing numbers of input features (supplementary Figs. S4 and S5, Supplementary Material online).

To put forth a better perspective on the utility of the haplotype alignment processing method T-REx uses, we experimented with another protocol for processing haplotype alignments, which is similar to that of Torada et al. (2019). As in Torada et al. (2019), we sorted the haplotypes along the entire 1.1 Mb genomic region, which is in contrast to the alignment processing method employed by T-REx, where haplotypes were sorted in a sliding window. This key difference between these two protocols may be an important factor behind the decreased false discovery of sweeps by T-REx(EN) (compare supplementary Figs. S19 to S1 and supplementary Figs. S20 to S2, Supplementary Material online). Overall, our experiments under the constant-size demographic history across different ranks (compare supplementary Figs. S19 to S1 and supplementary Figs. S20 to S2, Supplementary Material online) show that our unique alignment processing method has a distinct advantage in terms of downstream classification accuracy and power over another contemporary approach for processing haplotype alignments. If ImaGene adopted this local alignment processing approach, then it would have potentially resulted in performance that is more close to that exhibited by *T-REx*. Another factor that has likely impacted the performance of *ImaGene* in our study is that it is CNN-based, and CNNs typically require large training sets to achieve optimal performance (Luo et al. 2018). In the original *ImaGene* article, Torada et al. (2019) employed 50,000 observations per class for training. In contrast, we used 10,000 observations per class for comparison purposes with *T-REx*, which may have influenced the results shown by *ImaGene*. Moreover, a key distinction between *ImaGene* and *T-REx* is that *ImaGene* uses larger resized 128 × 128-dimensional images as input, which have the potential for reduced robustness to noise compared with *T-REx*, as more noise is averaged out with its smaller 64 × 64-dimensional input images.

When analyzing modern genomic data, it is common to encounter regions with missing segments due to artifacts or sequence alignment problems, making it critical that machine learning tools remain robust to the challenge such missing data poses. In our tests with missing segments, we found that T-REx(EN) was fairly robust, but ImaGene was deleteriously affected by an increase in the misclassification rate of neutral regions—though for reasonable false positive rates, ImaGene was also robust (Fig. 5). An avenue to alleviate this problem is to train classifiers with missing random segments (Kern and Schrider 2018), which allows classifiers to learn the underlying patterns associated with missing data. Randomly removing chunks from alignments in non-overlapping windows from the training data before training classifiers has been shown to offset the deleterious effects of such missing data (Mughal and DeGiorgio 2019; Mughal et al. 2020). Also, filling in missing values in test data through genotype imputation (e.g. Li et al. 2010; Moritz and Bartz-Beielstein 2017; Browning et al. 2021; Davies et al. 2021) may be another direction to combat the problem of missing data. Classifiers that are fed test data after imputing the missing values tend to be robust when faced with missing data in genomes and may achieve better prediction accuracy (Sarkar et al. 2021).

We have implemented T-REx as a binary classifier to differentiate sweeps from neutrality, but this modeling strategy can also be employed for broader classification problems in evolutionary genomics. For example, using multiclass extensions to the machine learning models discussed here, the T-REx framework could accommodate classifiers for jointly discriminating among other evolutionary processes, such as balancing selection and adaptive introgression, in addition to neutrality and sweeps from de novo mutations or standing variation. To illustrate, twodimensional representations of genomic data have been employed in multiclass models for robustly determining whether a genomic region is neutrally-evolving or has undergone a hard or soft sweep (Kern and Schrider 2018), as well as been shown to improve discrimination of adaptive introgression from sweeps and neutrality (Mughal et al. 2020). Moreover, Gower et al. (2021) employed images of sorted haplotype alignments as input

to a CNN with the aim to detect adaptive introgression—a setting that Mughal et al. (2020) still had trouble with based on two-dimensional images derived from hand-engineered population-genetic summary statistics. Indeed, Isildak et al. (2021) showed that CNNs applied to extract features from images of haplotype alignments outperformed feed-forward neural networks applied to hand-engineered population-genetic features in discriminating between recent balancing selection and incomplete sweeps, which are two evolutionary settings that can yield similar distributions of haplotype variation and are thus difficult to tease apart. These examples highlight the promise that automatic feature extraction from image representations of haplotypic variation has for probing genomes for diverse forms of natural selection.

Throughout this article, we have explored the problem of identifying natural selection as a classification task. However, the machine learning models employed by T-REx are flexible, and changing from a qualitative to a quantitative output would shift the problem from a classification to a regression problem. By using a regression framework, T-REx could predict underlying sweep parameters, such as selection strength, frequency of the selected allele when it became beneficial, and time at which a sweep completed (Mughal and DeGiorgio 2019). Moreover, as in Flagel et al. (2019), framing the prediction problem as regression would allow for estimation of key demographic quantities, such as the timing and magnitude of population size changes, as well as genetic parameters, such as recombination rate. Hence, tensor decomposition represents a complementary tool for tackling an array of inference problems within population genomics that CNNs have already been demonstrated to be highly effective.

Another interesting avenue that can be explored and could potentially increase the accuracy and robustness of T-REx is the incorporation of ancient DNA data. Because information on temporal trajectories of genetic variation can be exploited when using ancient DNA, such additional data could enhance the detection and characterization of adaptive footprints. Indeed, recent studies have incorporated genetic variation from ancient samples to detect adaptive loci (e.g. Mathieson and McVean 2013a, 2013b; Field et al. 2016; Dehasque et al. 2020; Mathieson 2020; Rees et al. 2020; Whitehouse and Schrider 2023), and these studies have highlighted that use of such temporal data can aid in better detection of adaptive events. Our T-REx framework is amenable to direct incorporation of ancient DNA data sampled across time by including images of haplotype variation consecutively ordered across sampled time points along a fourth dimension of the tensor structure prior to applying tensor decomposition to perform feature extraction. Because this tensor decomposition will naturally preserve the autocovariation in diversity not only expected spatially in the genome, but also temporally through the addition of the fourth dimension, we believe this is a viable avenue for future exploration to boost the performance and scope of T-REx models.

Important limitations of T-REx are the runtime and memory-usage associated with larger training sets (N) and higher ranks (R). In our experiments, we found that tensor decomposition took substantially greater time and memory even for modest increases in R. Downsampling each observation to a 64 × 64-dimensional matrix helped in reducing the complexity, and also likely aided in robustness of our models by averaging some of the noise in the input images. Moreover, we have been concerned with three-way tensors only, but if we were to consider increasing the number of dimensions, it would render the process computationally costlier than a three-way case, as the number of elements in the tensor would increase exponentially with each added dimension (Kruppa 2017). Also, the alternating least squares algorithm (see Methods section) for learning the factors matrices for CP tensor decomposition will need to find the factor matrices associated with each added dimension. For example, if we were to include ancient DNA data sampled over time as the fourth dimension in our existing pipeline, then it would be a four-way tensor where we would have an extra factor matrix D, which the ALS algorithm has to estimate through iteration and will incur greater runtime before reaching convergence.

We have focused on CP tensor decomposition (Hitchcock 1927; Harshman 1970). However, other algorithms for decomposing tensors exist, each with their own advantages and disadvantages relative to CP decomposition. Examples are multilinear principal component analysis (MPCA) (Lu et al. 2008), Tucker decomposition (Tucker 1966), higher-order singular value decomposition (HOSVD) (Lathauwer et al. 2000), and tensor train (TT) decomposition (Oseledets 2011), which are widely used alternative approaches for performing tensor decomposition (e.g. Sidiropoulos et al. 2017; Yuwang et al. 2019). Methods such as CP decomposition, MPCA, HOSVD, and TT are closely related to Tucker decomposition (Zare et al. 2018; Yuwang et al. 2019) in their working procedures, which is based on finding the linear combination of outer products of vectors. Among these different techniques, Tucker decomposition (Tucker 1966) is the most similar in operation to CP decomposition, as it also hinges on the idea of using alternating least squares to estimate a core tensor and factor matrices, though the core tensor produced by Tucker decomposition is not necessarily diagonal like the one CP decomposition outputs (Yuwang et al. 2019) and the ranks of the factor matrices are not constrained to be identical. Despite their similarities, CP decomposition is able to produce unique solutions unlike Tucker decomposition, where factor matrices change as the core tensor is changed (Kim et al. 2014; Zare et al. 2018). Also, the rank-one factors generated by Tucker decomposition are orthonormal, which is not the case for CP tensor decomposition (Kim et al. 2014).

The *T-REx* methodology introduced here represents complementary approach to CNNs for automatic feature extraction of haplotype alignment images. This framework is flexible, as it permits learned features to be used in both

linear and advanced nonlinear models, and can be extended into multiclass and quantitative prediction problems within evolutionary genomics. Moreover, we demonstrated that T-REx has an edge over a current leading CNN-based architecture in terms of power and accuracy, partially due to its unique alignment processing strategy for easier feature detection. Moreover, T-REx identified previously hypothesized and novel candidate sweeps in our empirical application, highlighting its efficacy in practice. Despite the promising performance metrics of T-REx, computation time of T-REx increases with increasingly higher ranks and sample sizes. However, excellent power and accuracy were achieved for modest numbers of features and training set sizes, and so we do not see this as a major hurdle for T-REx. Given the rapidly changing landscape of computational approaches for learning about and uncovering evolutionary mechanisms, T-REx provides a bridge between modern methodologies for feature extraction and well-established classical machine learning prediction techniques.

Methods

CP Tensor Decomposition

Consider a tensor $\mathcal{X} \in \mathbb{R}^{|x| \times K}$ of order three, where the first dimension will collect I observations of two-dimensional images each with $J \times K$ pixel values. The idea behind CP tensor decomposition is to express such a tensor as a sum of R tensors, where each of these tensors is expressed as the outer product of three rank one tensors. That is, we wish to estimate \mathcal{X} as

$$\widehat{\mathcal{X}} = \sum_{r=1}^{R} \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r,$$

where the symbol ° denotes the outer product and where $\mathbf{a}_r \in \mathbb{R}^I$, $\mathbf{b}_r \in \mathbb{R}^J$, and $\mathbf{c}_r \in \mathbb{R}^K$ such that $\mathcal{X} \approx \widehat{\mathcal{X}}$. For our setting, I will represent the number of training observations, J a proxy for the number of haplotypes, and K a proxy for the number of loci (see *Alignment processing* subsection of the *Methods* for details). Because we are working with tensors of order three, which is a higher-order tensor, we have column, row, and tube *Fibers*, which are, respectively, termed mode-1, mode-2, and mode-3 of the tensor.

Preprocessing Tensors

Prior to application of CP decomposition, we need to preprocess the input tensors through centering and scaling operations. Because the data are represented as a three-way tensor, preprocessing is different from conventional methods (Kolda and Bader 2009). Let value x_{ijk} denote elements i, j, and k, respectively, for the first, second, and third dimensions of the tensor $\mathcal{X} \in \mathbb{R}^{|\mathcal{X}| \times K}$. This tensor is centered as

$$x_{ijk}^{\text{centered}} = x_{ijk} - \overline{x}_{jk} \tag{1}$$

where

$$\overline{x}_{jk} = \frac{1}{l} \sum_{i=1}^{l} x_{ijk}$$

is the sample mean across the I training observations. Here the index i is related to the first mode, so it runs from 1 to I. Similarly index j runs from 1 to J and index k runs from 1 to K. This kind of centering is called single centering across the first mode (Bro 1997), and causes the mean of each pixel of an image to be zero across the training samples. We could have centered on multiple modes simultaneously, which is called double or triple centering depending on the number of modes on which to simultaneously center. However, centering one mode at a time is appropriate for CP decomposition, as any other kind of centering would destroy the multilinear properties of the data (Bro 1997)

In addition to centering, scaling should be performed on only one mode at a time, and we have chosen to scale in the first mode for our application (Kolda and Bader 2009). Scaling is performed as

$$x_{ijk}^{\text{scaled}} = \frac{x_{ijk}}{s_i} \tag{2}$$

where

$$s_i = \sqrt{\sum_{j=1}^{J} \sum_{k=1}^{K} x_{ijk}^2}.$$

This kind of scaling ensures that the overall intensity of values across pixels in an image are identical for each training sample. The order of scaling and centering is not arbitrary, as the operations are not commutative (Kolda and Bader 2009). Centering across a particular mode after scaling disturbs scaling across all modes. On the other hand, scaling across a particular mode after centering destroys centering across that mode. For these reasons, the order of centering and scaling is important. Centering is performed after scaling so that the scaled mode variance is not exactly one, but any large differences across the mode are mostly equalized (Kolda and Bader 2009). Centering is then performed, which ensures that the mode to be centered has a mean of zero.

Computing the CP Decomposition

After performing tensor decomposition on the training tensor $\mathcal{X} \in \mathbb{R}^{|\mathbf{X}|\mathbf{X}|K}$ to obtain a rank R CP decomposition, we obtain the three-factor matrices

$$\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_R] \in \mathbb{R}^{I \times R}$$

$$\mathbf{B} = [\mathbf{b}_1 \ \mathbf{b}_2 \ \cdots \ \mathbf{b}_R] \in \mathbb{R}^{I \times R}$$

$$\mathbf{C} = [\mathbf{c}_1 \ \mathbf{c}_2 \ \cdots \ \mathbf{c}_R] \in \mathbb{R}^{K \times R}$$

which yield an approximation of the tensor through the outer product

$$\widehat{\mathcal{X}} = \sum_{r=1}^{R} \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r,$$

where λ_r , $r=1, 2, \ldots, R$, scales the rth tensor to have unit norm. From the factor matrices \mathbf{B} and \mathbf{C} , we can depict the features extracted by component r of the CP model from the training data with the expression $\mathbf{b}_r \circ \mathbf{c}_r \in \mathbb{R}^{J \times R}$ (Papastergiou et al. 2018).

The key algorithm behind computing the CP decomposition is alternating least squares (Carroll and Chang 1970), which is a minimization algorithm. For a tensor of order three, given a rank R to approximate the training tensor (\mathcal{X}), alternating least-squares fixes two of the factor matrices, whereas solving for the remaining factor matrix that minimizes the sum of the squared differences in the elements of the estimated tensor ($\widehat{\mathcal{X}}$) and the training tensor. For example, if factor matrices \mathbf{B} and \mathbf{C} are fixed, then we seek to find \mathbf{A} that has this minimal sum of squared errors.

Denote the best factor matrices **A**, **B**, and **C** at iteration $t \in \{0, 1, 2, ...\}$ of the alternating least-squares algorithm by $\mathbf{A}^{(t)}$, $\mathbf{B}^{(t)}$, and $\mathbf{C}^{(t)}$, respectively. Given these factor matrices, let the current estimate of the training tensor be

$$\widehat{\mathcal{X}}^{(t)} = \sum_{r=1}^{R} \mathbf{a}_r^{(t)} \circ \mathbf{b}_r^{(t)} \circ \mathbf{c}_r^{(t)}.$$

Define the element-wise squared difference between two order-three tensors $\mathcal{X} \in \mathbb{R}^{|\mathcal{X}| \times K}$ and $\mathcal{Y} \in \mathbb{R}^{|\mathcal{X}| \times K}$ as

$$D^{2}(\mathcal{X}, \mathcal{Y}) = \sum_{i=1}^{J} \sum_{j=1}^{K} \sum_{k=1}^{K} (\mathcal{X}_{ijk} - \mathcal{Y}_{ijk})^{2}.$$

Alternating least squares on this tensor of order three is given by the following three steps:

1) Step 1: fix $\mathbf{A}^{(t)}$ and $\mathbf{B}^{(t)}$ and solve for $\mathbf{C}^{(t+1)}$

$$\mathbf{C}^{(t+1)} = \underset{\mathbf{C} = [\mathbf{c}_1 \ \mathbf{c}_2 \ \cdots \ \mathbf{c}_R]}{\text{arg min}} D^2 \left(\mathcal{X}, \ \sum_{r=1}^R \mathbf{a}_r^{(t)} \circ \mathbf{b}_r^{(t)} \circ \mathbf{c}_r \right)$$

2) Step 2: fix $\mathbf{A}^{(t)}$ and $\mathbf{C}^{(t)}$ and solve for $\mathbf{B}^{(t+1)}$

$$\mathbf{B}^{(t+1)} = \underset{\mathbf{B} = [\mathbf{b}_1 \ \mathbf{b}_2 \ \cdots \ \mathbf{b}_R]}{\text{arg min}} D^2 \left(\mathcal{X}, \ \sum_{r=1}^R \mathbf{a}_r^{(t)} \circ \mathbf{b}_r \circ \mathbf{c}_r^{(t)} \right)$$

3) Step 3: fix $\mathbf{B}^{(t)}$ and $\mathbf{C}^{(t)}$ and solve for $\mathbf{A}^{(t+1)}$

$$\mathbf{A}^{(t+1)} = \underset{\mathbf{A}=[\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_R]}{\text{arg min}} D^2 \left(\mathcal{X}, \ \sum_{r=1}^R \mathbf{a}_r \circ \mathbf{b}_r^{(t)} \circ \mathbf{c}_r^{(t)} \right)$$

Steps 1 to 3 are repeated until convergence, and final estimated factor matrices are identical to those from the final iteration—i.e. $\mathbf{A} = \mathbf{A}^{(t+1)}$, $\mathbf{B} = \mathbf{B}^{(t+1)}$, and $\mathbf{C} = \mathbf{C}^{(t+1)}$. At each step, we incorporate the values for λ_r , $r = 1, 2, \ldots, R$ into the estimated tensor.

Projecting Test Observations onto Identified Factor Matrices Given a new test tensor $\mathcal{X}_{\text{test}} \in \mathbb{R}^{I_{\text{test}}J,K}$ of I_{test} test observations, we can project the test observations onto the learned factor A so that it falls within the subspace learned by decomposing the training tensor \mathcal{X} . However, before doing so, we must ensure that the test dataset lies in the same input space as the training set. Thus, we preprocess the test dataset by applying Equations (1) and (2) for centering and scaling. It is important to note that Equation (1) refers to centering with respect to the training set (i.e. subtracting \bar{x}_{ik}), and so the test set must be centered with the mean training pixel value \bar{x}_{ik} and not a similar quantity for the test set. Thus, the centering values for the training set must be retained so that the test set is centered with identical values. Assuming $\mathcal{X}_{\mathsf{test}}$ has now been properly preprocessed, we can project the test data onto the learned features representing each input image by Kolda and Bader (2009)

$$\mathbf{A}_{\text{test}} = \mathbf{X}_{\text{test}(1)}(\mathbf{C} * \mathbf{B})(\mathbf{C}^{\mathsf{T}} \mathbf{C} \odot \mathbf{B}^{\mathsf{T}} \mathbf{B})^{\dagger} \Lambda^{-1}, \tag{3}$$

where $\mathbf{X}_{\text{test}(1)}$ is the mode-1 unfolding (matricization) of the tensor $\mathcal{X}_{\text{test}}$, the superscript T denotes transpose, the symbol * denotes the Khatri–Rao product, the \odot symbol denotes the Hadamard (element-wise) product, the superscript \dagger denotes the Moore–Penrose pseudoinverse, and $\mathbf{\Lambda}^{-1}$ represents the inverse of the diagonal matrix $\mathbf{\Lambda} \in \mathbb{R}^{R \times R}$ of scaling terms $\lambda_1, \lambda_2, ..., \lambda_R$.

Alignment Processing

We used a novel approach for processing the haplotype alignments in a way that helps the classifiers detect the footprint of a selective sweep. For each simulated 1.1 Mb region, we locally sorted haplotypes in windows of 100 single-nucleotide polymorphisms (SNPs), moving the window along the region with a stride of 10 SNPs, where values at SNPs were averaged for all windows that overlapped them. This method of alignment processing can help classifiers identify signals of lost haplotypic diversity if sweeps are weak or old, whereas also retaining power for strong and recent sweeps. To reduce the complexity of the tensor decomposition and noise in the sorted haplotype alignments, we downsampled the alignment images to 64 × 64-dimensional matrices using the scikitimage library (Pedregosa et al. 2011), where Gaussian smoothing was employed to preserve the spatial relationships of pixels within the images and to avoid aliasing artifacts. We highlight the advantage of our alignment processing approach by pitting the results obtained after employing our unique alignment processing strategy against those of an alignment processing approach that is similar to that used by *ImaGene* (Torada et al. 2019) (compare supplementary Figs. S19 to S1 and supplementary Figs. S20 to S2, Supplementary Material online).

T-REx Model Training and Hyperparameter Tuning

We have implemented three classical linear and nonlinear machine learning models with different R packages into our T-REx framework. For performing tensor decomposition, we used the R package rTensor (Li et al. 2018). Additionally, we employed the R packages glmnet (Friedman et al. 2010), liquidsym (Steinwart and Thomann 2017), and ranger (Wright and Ziegler 2017) for implementing T-REx(EN), T-REx(SVM), and T-REx(RF), respectively. During the training of each classifier, we have 10⁴ observations in each class, with each observation consisting of sorted haplotype alignments (details provided in the Alignment processing subsection of the Methods). We then applied a rank R tensor decomposition (see CP tensor decomposition subsection of the Methods for details) to obtain a set of R derived features for each observation in each class to be used as input for our T-REx classifiers.

Before the testing phase commences, we tuned hyperparameters, which control certain components of the model training process, of each model by selecting optimal hyperparameters through the cross-validation procedure. Hyperparameter tuning is a way of selecting suitable hyperparameter values from a range of possible values. Specifically, we performed 10-fold cross-validation such that on each of the 10-folds we selected 10% of the samples (1,000 observations per class) from the dataset to be reserved for model validation and the remaining 90% of the samples (9,000 observations per class) to be employed for model training. This procedure allowed us to evaluate how well the model would perform on unseen data (from the validation set) for a given set of hyperparameter values. For each of the classifiers, we chose the model structure that yielded the smallest cross-validation error after performing hyperparameter tuning.

For hyperparameter tuning of T-REx(EN), we explored a grid of values $\alpha \in \{0, 0.1, ..., 1.0\}$, where α denotes the proportion of the model for which the parameters are penalized with an L_2 -norm penalty, whereas $1 - \alpha$ is the proportion penalized with an L_1 -norm penalty (Hastie et al. 2009). In addition to α , we tuned another hyperparameter $\lambda \geq 0$, which modulates the complexity of the fitted model by controlling the influence of the L_1 - and L_2 -norm penalties during model training. By tuning both λ and α , we are controlling the complexity of the fitted model, whereas simultaneously performing feature selection by inclusion of the L₁-norm penalty (Hastie et al. 2009). We find the optimal λ and α combination that gives the minimum 10-fold cross-validation error. For implementing T-REx(SVM), we used the radial basis kernel for nonlinear modeling, which has a hyperparameter γ that is inversely proportional to the variance (width) of the radial basis kernel, which has a shape similar to a Gaussian function (Hastie et al. 2009). To implement T-REx(RF), we chose a large number (5,000) of random trees to use in the RF ensemble, as test error stabilizes with enough trees in the forest (Hastie et al. 2009), and used the default number of 10 random splits within ranger for growing each decision tree within the RF.

Finally, regardless of machine learning method, another important hyperparameter is the rank R of the tensor decomposition. For each value of $R \in \{50, 100, 150, 200, 250, 300\}$, we computed the 10-fold cross-validation error for T-REx(EN) and T-REx(SVM) and the out-of-bag error for T-REx(RF) (Hastie et al. 2009). We chose the (R, λ, α) triple that resulted in the smallest 10-fold cross-validation error for T-REx(EN), the (R, γ) pair that results in the smallest 10-fold cross-validation for T-REx(SVM), and the value of R that results in the smallest out-of-bag error for T-REx(RF). After selecting the set of optimal hyperparameters of each method, the three T-REx models were each trained on the full dataset of 10^4 training observations per class conditional on their optimal hyperparameters, and these models were deployed on further testing data.

Training and Evaluating ImaGene

To fully evaluate the performance *T-REx*, we compared it with the CNN-based sweep classifier ImaGene (details are provided in the *Results*). Although both *T-REx* and ImaGene use haplotype alignments in the form of images, there are differences in the procedure used to process the images and perform model training.

For training, ImaGene employs a "simulation-on-thefly" approach of using newly generated data at each training epoch (iteration of gradient descent). This simulation-on-the-fly approach prevents ImaGene from overfitting. For consistency and fairness in comparison between T-REx and ImaGene, we deviated from this default setting of ImaGene so that it is pitted against T-REx on identical simulation data. Specifically, we used the same 10⁴ training observations per class when training ImaGene as we employed for training T-REx for each simulation setting (details regarding the simulation protocol are provided in the Results). To prevent overfitting, we employed early stopping (Goodfellow et al. 2016), by setting the number of epochs to train ImaGene as the point at which the validation loss starts to rise, which suggests overfitting, where the validation loss was computed across 1,000 observations per class that were held out for validation. supplementary Figure S21, Supplementary Material online displays the validation and training loss curves over 200 training epochs, showing that the validation curve begins to increase at approximately 25 epochs. We therefore retrained the ImaGene model on the full dataset of 10⁴ observations per class for 25 epochs.

Application to Empirical Data

With the aim of detecting novel candidate genes that may be subject to positive natural selection and previously hypothesized candidates of positive natural selection, we used empirical data of the CEU human population from the 1000 Genomes Project dataset (The 1000 Genomes Project Consortium 2015). We first filtered variant calls to include biallelic SNPs. Second, we removed SNPs with minor allele count less than three, as Mughal et al. (2020) demonstrated the frequencies of singleton and doubleton SNPs in the CEU population from the 1000 Genomes Project dataset differed from those predicted by the inferred demographic model (Terhorst et al. 2017) that we used to train our classifiers. Moreover, because regions of the genome that are harder to map and align may lead to technical artifacts affecting observed genomic variation (Derrien et al. 2012), we removed sites that could have problematic mapping or alignability to circumvent such potential artifacts. Specifically, we used the CRG score to measure mappability and alignability of a genomic region and removed sites falling within 100 kb windows for which the mean CRG100 score within the window was less than 0.9 (Mughal et al. 2020). We then applied our unique alignment processing approach to further process the data before supplying it to T-REx.

Supplementary Material

Supplementary material is available at Molecular Biology and Evolution online.

Acknowledgments

This work was supported by National Institutes of Health grant R35GM128590 and by National Science Foundation grants DEB-1949268, BCS-2001063, and DBI-2130666. Computations for this research were performed using the services provided by Research Computing at the Florida Atlantic University.

Conflict of interest: None declared.

References

The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* 2015:**526**:68–74. https://doi.org/10.1038/nature15393.

Anguita-Ruiz A, Aguilera CM, Gil Á. Genetics of lactose intolerance: an updated review and online interactive world maps of phenotype and genotype frequencies. *Nutrients* 2020:**12**:2689. https://doi.org/10.3390/nu12092689.

Arnab SP, Amin MR, DeGiorgio M. Uncovering footprints of natural selection through time-frequency analysis of genomic summary statistics. *Mol Biol Evol*. 2023:**40**(7):msad157. https://doi.org/10. 1093/molbey/msad157.

Bagchi P, Torres M, Tsai B, Qi L. Selective EMC subunits act as molecular tethers of intracellular organelles exploited during viral entry. *Nat Commun*. 2020:**11**(1):1127. https://doi.org/10.1038/s41467-020-14967-w.

Baudat F, Buard J, Grey C, Fledel-Alon A, Ober C, Przeworski M, Coop G, de Massy B. PRDM9 is a major determinant of meiotic recombination hotspots in humans and mice. *Science* 2010:**327**: 836–840. https://doi.org/10.1126/science.1183439.

Beall CM, Cavalleri GL, Deng L, Elston RC, Gao Y, Knight J, Li C, Chuan Li J, Liang Y, McCormack M, et al. Natural selection on *EPAS1* (*HIF2a*) associated with low hemoglobin concentration in Tibetan highlanders. *Proc Natl Acad Sci U S A*. 2010:**107**:11459–11464. https://doi.org/10.1073/pnas.1002443107.



- Bedford T, Cobey S, Pascual M. Strength and tempo of selection revealed in viral gene genealogies. *BMC Evol Biol.* 2011:11:220. https://doi.org/10.1186/1471-2148-11-220.
- Beichman AC, Huerta-Sanchez E, Lohmueller KE. Using genomic data to infer historic population dynamics of nonmodel organisms. *Annu Rev Ecol Evol Syst.* 2018:**49**(1):433–456. https://doi.org/10.1146/ecolsys.2018.49.issue-1.
- Bellman R. Dynamic programming. Science 1966:**153**(3731):34–37. https://doi.org/10.1126/science.153.3731.34.
- Bergeron LA, Besenbacher S, Zheng J, Li P, Bertelsen MF, Quintard B, Hoffman JI, Li Z, Leger JSt., Shao C, et al. Evolution of the germline mutation rate across vertebrates. *Nature* 2023:**615**:285–291. https://doi.org/10.1038/s41586-023-05752-y.
- Bernatchez L, Landry C. MHC studies in nonmodel vertebrates: what have we learned about natural selection in 15 years? *J Evol Biol.* 2003:**16**:363–377. https://doi.org/10.1046/j.1420-9101.2003. 00531.x.
- Bersaglieri T, Sabeti PC, Patterson N, Vanderploeg T, Schaffner SF, Drake JA, Rhodes DE, Reich M, Hirschhorn JN. Genetic signatures of strong recent positive selection at the lactase gene. Am J Hum Genet. 2004;**74**:1111–1120. https://doi.org/10.1086/421051.
- Bigham A, Bauchet M, Pinto D, Mao X, Akey JM, Mei R, Scherer S, Julian CG, Wilson MJ, Herráez DL, Brutsaert T, Parra EJ, Moore LG, Schriver MD. Identifying signatures of natural selection in Tibetan and Andean populations using dense genome scan data. PLoS Genet. 2010:6:e1001116. https://doi.org/10.1371/journal.pgen.1001116.
- Bitarello BD, Brandt DYC, Meyer D, André AM. Inferring balancing selection from genome-scale data. *Genome Biol Evol*. 2023:**15**: evad032. https://doi.org/10.1093/gbe/evad032.
- Booker TR, Yeaman S, Whitlock MC. Variation in recombination rate affects detection of outliers in genome scans under neutrality. *Mol Ecol.* 2020:**29**:4274–4279. https://doi.org/10.1111/mec.v29.22.
- Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE, Adams MD, Schmidt S, Sninsky JJ, Sunyaev SR, et al. Assessing the evolutionary impact of amino acid mutations in the human genome. *PLoS Genet*. 2008:**30**:e1000083. https://doi.org/10.1371/journal.pgen.1000083.
- Bro R. PARAFAC. Tutorial and applications. *Chemometr Intell Lab Syst.* 1997:**38**:149–171. https://doi.org/10.1016/S0169-7439(97) 00032-4
- Bromham L. The genome as a life-history character: why rate of molecular evolution varies between mammal species. *Philos Trans R Soc Lond B Biol Sci.* 2011:**366**:2503–2513. https://doi.org/10.1098/rstb.2011.0014.
- Bromham L, Hua X, Lanfear R, Cowman PF. Exploring the relationships between mutation rates, life history, genome size, environment, and species richness in flowering plants. *Am Nat.* 2015:**185**: 507–524. https://doi.org/10.1086/680052.
- Browning BL, Tian X, Zhou Y, Browning SR. Fast two-stage phasing of large-scale sequence data. *Am J Hum Genet.* 2021:**108**(10): 1880–1890. https://doi.org/10.1016/j.ajhg.2021.08.005.
- Burke MK. 2012. How does adaptation sweep through the genome? Insights from long-term selection experiments. Proc R Soc. 279: 5029–5038.
- Campbell MC, Ashong B, Teng S, Harvey J, Cross CN. Multiple selective sweeps of ancient polymorphisms in and around LTα located in the MHC class III region on chromosome 6. BMC Evol Biol. 2019:19(1):218. https://doi.org/10.1186/s12862-019-1516-y.
- Carroll JD, Chang JJ. Analysis of individual differences in multidimensional scaling via an n-way generalization of "Eckart-Young" decomposition. *Psychometrika* 1970:**35**:283–319. https://doi.org/10.1007/BF02310791.
- Chan J, Perrone V, Spence JP, Jenkins PA, Mathieson S, Song YS. A likelihood-free inference framework for population genetic data using exchangeable neural networks. Adv Neural Inf Process Syst. 2018:**31**:8594.
- Charlesworth D, Charlesworth B, Morgan MT. The pattern of neutral molecular variation under the background selection model.

- Genetics 1995:**141**:1619–1632. https://doi.org/10.1093/genetics/141.4.1619.
- Charlesworth B, Jensen JD. Effects of selection at linked sites on patterns of genetic variability. *Annu Rev Ecol Evol Syst.* 2021:**52**: 177–197. https://doi.org/10.1146/ecolsys.2021.52.issue-1.
- Charlesworth B, Morgan MT, Charlesworth D. The effect of deleterious mutations on neutral molecular variation. *Genetics* 1993:**134**: 1289–1303. https://doi.org/10.1093/genetics/134.4.1289.
- Charlesworth B, Nordborg M, Charlesworth D. The effects of local selection, balanced polymorphism and background selection on equilibrium patterns of genetic diversity in subdivided populations. *Genet Res.* 1997:**70**:155–174. https://doi.org/10.1017/S0016672397002954.
- Chen H, Patterson N, Reich D. Population differentiation as a test for selective sweeps. *Genome Res.* 2010:**20**:393–402. https://doi.org/10.1101/gr.100545.109.
- Chen SL, Qin ZY, Hu F, Wang Y, Dai YJ, Liang Y. The role of the HOXA gene family in acute myeloid leukemia. *Genes* 2019:**10**(8):621. https://doi.org/10.3390/genes10080621.
- Cheng X, Xu C, DeGiorgio M. Fast and robust detection of ancestral selective sweeps. *Mol Ecol*. 2017:**26**:6871–6891. https://doi.org/10.1111/mec.2017.26.issue-24.
- Comeron JM. Background selection as baseline for nucleotide variation across the drosophila genome. *PLoS Genet*. 2014:**10**: e1004434. https://doi.org/10.1371/journal.pgen.1004434.
- Cutter AD, Payseur BA. Genomic signatures of selection at linked sites: unifying the disparity among species. *Nat Rev Genet*. 2013:**14**:262–274. https://doi.org/10.1038/nrg3425.
- Danovi S. 2023. Mutation rates across species. *Nat Genet.* **54**:285–291. Davies RW, Kucka M, Su D, Shi S, Flanagan M, Cunniff CM, Chan Y, Myers S. Rapid genotype imputation from sequence with reference panels. *Nat Genet.* 2021:**53**:1104–1111. https://doi.org/10.1038/s41588-021-00877-0.
- DeGiorgio M, Huber CD, Hubisz MJ, Hellmann I, Nielsen R. SweepFinder2: increased sensitivity, robustness, and flexibility. Bioinformatics 2016:32:1895–1897. https://doi.org/10.1093/bioinformatics/btw051.
- DeGiorgio M, Szpiech ZA. A spatially aware likelihood test to detect sweeps from haplotype distributions. *PLoS Genet.* 2022:**18**: e1010134. https://doi.org/10.1371/journal.pgen.1010134.
- Dehasque M, Avila-Árcos MC, Díez-del-Molino D, Fumagalli M, Guschanski K, Lorenzen ED, Malaspinas AS, Marques-Bonet T, Martin MD, Murray GGR, et al. Inference of natural selection from ancient DNA. *Evol Lett.* 2020:**4**:94–108. https://doi.org/10.1002/evl3.165.
- Derrien T, Estellé J, Marco Sola S, Knowles DG, Raineri E, Guigó R, Ribeca P. Fast computation and applications of genome mappability. *PLoS ONE*. 2012:**7**(1):e30377. https://doi.org/10.1371/journal.pone.0030377.
- Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. Robust demographic inference from genomic and SNP data. *PLoS Genet*. 2013:9(10):1–17. https://doi.org/10.1371/journal.pgen.1003905.
- Fagny M, Patin E, Enard D, Barreiro LB, Quintana-Murci L, Laval G. Exploring the occurrence of classic selective sweeps in humans using whole-genome sequencing data sets. Mol Biol Evol. 2014:31(7):1850-1868. https://doi.org/10.1093/molbev/msu118.
- Fan S, Hansen ME, Lo Y, Tishkoff SA. Going global by adapting local: a review of recent human adaptation. *Science* 2016:**354**:54–59. https://doi.org/10.1126/science.aaf5098.
- Feder AF, Pennings PS, Petrov DA. The clarifying role of time series data in the population genetics of HIV. *PLoS Genet.* 2021:17: e1009050. https://doi.org/10.1371/journal.pgen.1009050.
- Feder AF, Rhee S-Y, Holmes SP, Shafer RW, Petrov DA, Pennings PS. More effective drugs lead to harder selective sweeps in the evolution of drug resistance in HIV-1. *eLife* 2016:**5**:e10670. https://doi.org/10.7554/eLife.10670.
- Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R. On detecting incomplete soft or hard selective sweeps using haplotype structure. *Mol Biol Evol.* 2014:**31**:1275–1291. https://doi.org/10.1093/molbev/msu077.

- Field Y, Boyle EA, Telis N, Gao Z, Gaulton KJ, Golan D, Yengo L, Rocheleau G, Froguel P, McCarthy MI, et al. Detection of human adaptation during the past 2000 years. Science 2016:354: 760–764. https://doi.org/10.1126/science.aag0776.
- Flagel L, Brandvain Y, Schrider DR. The unreasonable effectiveness of convolutional neural networks in population genetic inference. *Mol Biol Evol*. 2019:**36**:220–238. https://doi.org/10.1093/molbev/msy224.
- Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *J Stat Softw.* 2010 **33**(5):1–22.
- Galetto R, Giacomoni V, Véron M, Negroni M. Dissection of a circumscribed recombination hot spot in HIV-1 after a single infectious cycle. *J Biol Chem.* 2006:**281**:2711–2720. https://doi.org/10.1074/jbc.M505457200.
- García-Cárdenas JM, Armendáriz-Castillo I, Pérez-Villa A, Indacochea A, Jácome-Alvarado A, López-Corté A, Guerrero S. Integrated in silico analyses identify PUF60 and SF3A3 as new spliceosome-related breast cancer RNA-binding proteins. *Biology* 2022:11: 481. https://doi.org/10.3390/biology11040481.
- Garud NR, Messer PW, Buzbas EO, Petrov DA. Recent selective sweeps in North American *Drosophila melanogaster* show signatures of soft sweeps. *PLoS Genet*. 2015:11:e1005004. https://doi.org/10.1371/journal.pgen.1005004.
- Ge F, Tie W, Zhang J, Zhu Y, Fan Y. Expression of the HOXA gene family and its relationship to prognosis and immune infiltrates in cervical cancer. *J Clin Lab Anal*. 2021:**35**:e24015. https://doi.org/10.1002/jcla.24015.
- Gillespie JH. Population genetics: a concise guide. 2nd ed. Baltimore (MD): The Johns Hopkins University Press; 2004.
- Goodfellow I, Bengio Y, Courville A. 2016, editors. *Deep learning*. Cambridge (MA): MIT Press.
- Goodwin ZA, de Guzman Strong D. Recent positive selection in genes of the mammalian epidermal differentiation complex locus. *Front Genet.* 2017:**7**:227. https://doi.org/10.3389/fgene.2016.00227.
- Gower G, Iáñez Picazo PI, Fumagalli M, Racimo F. Detecting adaptive introgression in human evolution using convolutional neural networks. eLife 2021:10:e64669. https://doi.org/10.7554/eLife.64669.
- Graham AM, McCracken KG. Convergent evolution on the hypoxia-inducible factor (HIF) pathway genes EGLN1 and EPAS1 in high-altitude ducks. Heredity. 2019:122:819–832. https://doi.org/10. 1038/s41437-018-0173-z.
- Gu J, Wang Z, Kuen J, Ma L, Shahroudy A, Shuai B, Liu T, Wang X, Wang G, Cai J, et al. Recent advances in convolutional neural networks. *Pattern Recogn*. 2018:**77**:354–377. https://doi.org/10.1016/ j.patcog.2017.10.013.
- Haller BC, Messer PW. SLiM 3: forward genetic simulations beyond the Wright–Fisher model. *Mol Biol Evol.* 2019:**36**:632–637. https://doi.org/10.1093/molbev/msy228.
- Harpak A, Bhaskar A, Pritchard JK. Mutation rate variation is a primary determinant of the distribution of allele frequencies in humans. PLoS Genet. 2016:12:e1006489. https://doi.org/10.1371/journal.pgen.1006489.
- Harris AM, DeGiorgio M. A likelihood approach for uncovering selective sweep signatures from haplotype data. *Mol Biol Evol*. 2020a:37: 3023–3046. https://doi.org/10.1093/molbev/msaa115.
- Harris AM, DeGiorgio M. Identifying and classifying shared selective sweeps from multilocus data. *Genetics* 2020b:**215**:143–171. https://doi.org/10.1534/genetics.120.303137.
- Harris AM, Garud NR, DeGiorgio M. Detection and classification of hard and soft sweeps from unphased genotypes by multilocus genotype identity. *Genetics* 2018:**210**:1429–1452. https://doi.org/10.1534/genetics.118.301502.
- Harris EE, Meyer D. The molecular signature of selection underlying human adaptations. *Am J Phys Anthropol*. 2006:**43**:89–130. https://doi.org/10.1002/ajpa.20518.
- Harshman RA. Foundations of the PARAFAC procedure: models and conditions for an explanatory multimodal factor analysis. UCLA Working Papers in Phonetics. Vol. 16. 1970. p. 1–84.

- Hastie T, Tibshirani R, Friedman J. The elements of statistical learning: data mining, inference, and prediction. 2nd ed. New York (NY): Springer; 2009.
- Hellenthal G, Stephens M. msHOT: modifying Hudson's ms simulator to incorporate crossover and gene conversion hotspots. *Bioinformatics* 2007:**23**:520–521. https://doi.org/10.1093/bioinformatics/btl622.
- Hermisson J, Pennings PS. Soft sweeps: molecular population genetics of adaptation from standing genetic variation. *Genetics* 2005: **169**:2335–2352. https://doi.org/10.1534/genetics.104.036947.
- Hermisson J, Pennings PS. Soft sweeps and beyond: understanding the patterns and probabilities of selection footprints under rapid adaptation. *Methods Ecol Evol*. 2017:**8**:700–716. https://doi.org/10.1111/mee3.2017.8.issue-6.
- Hernandez RD, Williamson SH, Bustamante CD. Context dependence, ancestral misidentification, and spurious signatures of natural selection. *Mol Biol Evol.* 2007:24:1792–1800. https://doi.org/10.1093/molbev/msm108.
- Hey J. What's so hot about recombination hotspots? *PLoS Biol.* 2004;**2**:e190. https://doi.org/10.1371/journal.pbio.0020190.
- Hider JL, Gittelman RM, Shah T, Edwards M, Rosenbloom A, Akey JM, Parra EJ. Exploring signatures of positive selection in pigmentation candidate genes in populations of east asian ancestry. *Evol Biol.* 2013:**13**:150.
- Hitchcock FL. The expression of a tensor or a polyadic as a sum of products. *J Math Phys.* 1927:**6**(1):164–189. https://doi.org/10.1002/sapm192761164.
- Hosoya N, Miyagawa K. Synaptonemal complex proteins modulate the level of genome integrity in cancers. *Cancer Sci.* 2021:**112**(3):989–996. https://doi.org/10.1111/cas.v112.3.
- Huber CD, DeGiorgio M, Hellmann I, Nielsen R. Detecting recent selective sweeps while controlling for mutation rate and background selection. *Mol Ecol.* 2016:**25**:142–156. https://doi.org/10.1111/mec.2016.25.issue-1.
- Huerta-Sánchez E, DeGiorgio M, Pagani L, Tarekegn A, Ekong R, Antao T, Cardona A, Montgomery HE, Cavalleri GL, Robbins PA, et al. Genetic signatures reveal high-altitude adaptation in a set of Ethiopian populations. *Mol Biol Evol*. 2013:**30**: 1877–1888. https://doi.org/10.1093/molbev/mst089.
- Huerta-Sánchez E, Jin X, Asan, Bianba Z, Peter BM, Vinckenbosch N, Liang Y, Yi X, He M, Somel M, et al. Altitude adaptation in Tibetans caused by introgression of Denisovan-like DNA. *Nature* 2014:512:194–197. https://doi.org/10.1038/nature13408.
- Ingram CJ, Mulcare CA, Itan Y, Thomas MG, Swallow DM. Lactose digestion and the evolutionary genetics of lactase persistence. *Hum Genet*. 2009:**124**:579–591. https://doi.org/10.1007/s00439-008-0593-6.
- Isildak U, Stella A, Fumagalli M. Distinguishing between recent balancing selection and incomplete sweep using deep neural networks. Mol Ecol Resour. 2021:21:2706–2718. https://doi.org/10.1111/men.v21.8.
- Itan Y, Powell A, Beaumont MA, Burger J, Thomas MG. The origins of lactase persistence in Europe. *PLoS Comput Biol.* 2009:5: e1000491. https://doi.org/10.1371/journal.pcbi.1000491.
- Jensen JD, Kim Y, DuMont VB, Aquadro CF, Bustamante CD. Distinguishing between selective sweeps and demography using DNA polymorphism data. *Genetics* 2005:**170**(3):1401-1410. https://doi.org/10.1534/genetics.104.038224.
- Kang L, AK Sharp GHe, Wang X, Brown AM, Michalak P, Weger-Lucarelli J. A selective sweep in the Spike gene has driven SARS-CoV-2 human adaptation. Cell. 2021:184:4392–4400.e4. https://doi.org/10.1016/j.cell.2021.07.007.
- Keinan A, Reich D. Human population differentiation is strongly correlated with local recombination rate. PLoS Genet. 2010:6(3): e1000886. https://doi.org/10.1371/journal.pgen.1000886.
- Kelly JK, Koseva B, Mojica JP. The genomic signal of partial sweeps in Mimulus guttatus. Genome Biol Evol. 2013:5:1457–1469. https:// doi.org/10.1093/gbe/evt100.
- Kern AD, Schrider DR. Discoal: flexible coalescent simulations with selection. *Bioinformatics* 2016:**32**(24):3839–3841. https://doi. org/10.1093/bioinformatics/btw556.



- Kern AD, Schrider DR. diploS/HIC: an updated approach to classifying selective sweeps. G3 (Bethesda). 2018:8:1959–1970. https://doi.org/10.1534/g3.118.200262.
- Kim B, Haotian L, Ngai W. 2014. A constructive algorithm for decomposing a tensor into a finite sum of orthonormal rank-1 terms. SIAM J Matrix Anal Appl. **36**(3):1315–1337.
- Kim K, Kim Y. Population genetic processes affecting the mode of selective sweeps and effective population size in influenza virus H₃N₂. BMC Evol Biol. 2016:**16**:156. https://doi.org/10.1186/s12862-016-0727-8.
- Kim Y, Nielsen R. Linkage disequilibrium as a signature of selective sweeps. *Genetics* 2004:**167**:1513–1524. https://doi.org/10.1534/genetics.103.025387.
- Kim Y, Stephan W. Detecting a local signature of genetic hitchhiking along a recombining chromosome. *Genetics* 2002:**160**:765–777. https://doi.org/10.1093/genetics/160.2.765.
- Kolda TG, Bader BW. Tensor decompositions and applications. SIAM Rev. 2009:**51**:455–500. https://doi.org/10.1137/07070111X.
- Kruppa K. Comparison of tensor decomposition methods for simulation of multilinear time-invariant systems with the MTI toolbox. IFAC-PapersOnLine. 2017:50(1):5610–5615. https://doi.org/10.1016/j.ifacol.2017.08.1107.
- Laayouni H, Oosting M, Luisi P, Ioana M, Alonso S, Ricaño-Ponce I, Trynka G, Zhernakova A, Plantinga TS, Cheng SC, et al. Convergent evolution in European and Rroma populations reveals pressure exerted by plague on toll-like receptors. *Proc Natl Acad Sci U S A*. 2014:**111**(7):2668–2673. https://doi.org/10.1073/pnas.1317723111.
- Lathauwer L, De Moor B, Vandewalle J. 2000. Multilinear singular value tensor decompositions. SIAM J Matrix Anal Apl. 24: 1253–1278.
- Lauterbur ME, Munch K, Enard D. Versatile detection of diverse selective sweeps with flex-sweep. bioRxiv, 2022.
- Lauterbur ME, Munch K, Enard D. Versatile detection of diverse selective sweeps with flex-sweep. *Mol Biol Evol*. 2023:**40**:msad139. https://doi.org/10.1093/molbev/msad139.
- LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature* 2015:**521**: 436–444. https://doi.org/10.1038/nature14539.
- LeCun Y, Bottou L, Bengio Y, Hafner P. Gradient-based learning applied to document recognition. *Proc IEEE*. 1998:**86**:2278–2324. https://doi.org/10.1109/5.726791.
- Lederberg J. J. B. S. Haldane (1949) on infectious disease and evolution. *Genetics* 1999:**153**:1–3. https://doi.org/10.1093/genetics/153.1.1.
- Lee KM, Coop G. Distinguishing among modes of convergent adaptation using population genomic data. *Genetics* 2017:**207**: 1591–1619. https://doi.org/10.1534/genetics.117.300417.
- Li JL, Bien J, Wells MT. rTensor: an R package for multidimensional array (tensor) unfolding, multiplication, and decomposition. J Stat Softw. 2018:87(10):1–31. https://doi.org/10.18637/jss.v087.i10.
- Li Y, Willer CJ, Ding J, Scheet P, Abecasis GR. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet Epidemiol*. 2010:**34**:816–834. https://doi.org/10.1002/gepi.20533.
- Lin K, Li H, Schlötterer C, Futschik A. Distinguishing positive selection from neutral evolution: boosting the performance of summary statistics. *Genetics* 2011:**187**:229–244. https://doi.org/10.1534/genetics.110.122614.
- Lindo J, Haas R, Hofman C, Apata M, Moraga M, Verdugo RA, Watson JT, Viviano Llvae C, Witonsky D, Beall C, et al. The genetic prehistory of the Andean highlands 7000 years BP through European contact. Sci Adv. 2018:4:eaau4921. https://doi.org/10.1126/sciadv.aau4921.
- Liu Y. 2021. Tensors for data processing: theory, methods, and applications. Cambridge (MA): Elsevier Science.
- Liu X, Zhang Y, Li Y, Pan J, Wang D, Chen W, Zheng Z, He X, Zhao Q, Pu Y, et al. EPAS1 gain-of-function mutation conributes to high-altitude adaptation in Tibetan horses. *Mol Biol Evol*. 2019:**36**: 2591–2603. https://doi.org/10.1093/molbev/msz158.

- Loewe L. Negative selection. Nat Educ. 2008:1(1):59.
- López S, García Ó, Yurrebaso I, Flores C, Acosta-Herrera M, Chen H, Gardeazabal J, Careaga JM, Boyano MD, Sánchez AAR, et al. The interplay between natural selection and susceptibility to melanoma on allele 374f of SLC45A2 gene in a south European population. *PLoS ONE*. 2014:9(8):1–12.
- Lou DI, McBee RM, Le UQ, Stone AC, Wilkerson GK, Demogines AM, Sawyer SL. Rapid evolution of BRCA1 and BRCA2 in humans and other primates. *BMC Evol Biol.* 2014:14:3136–3144. https://doi.org/10.1186/1471-2148-14-155.
- Lu H, Plataniotis KN, Venetsanopoulos AN. MPCA: multilinear principal component analysis of tensor objects. *IEEE Trans Neural Netw.* 2008:**19**(1):18–39. https://doi.org/10.1109/TNN.2007.901277.
- Luo C, Li X, Wang L, He J, Li D, Zhou J. How Does the Data set Affect
 CNN-based Image Classification Performance? 2018 5th
 International Conference on Systems and Informatics (ICSAI),
 Nanjing, China; 2018. p. 361–366.
- Mallick S, Gnerre S, Muller P, Reich D. The difficulty of avoiding false positives in genome scans for natural selection. *Genome Res.* 2009:**19**(5):922–933. https://doi.org/10.1101/gr.086512.108.
- Mathieson I. Limited evidence for selection at the fads locus in native American populations. *Mol Biol Evol*. 2020:**37**:2029–2033. https://doi.org/10.1093/molbev/msaa064.
- Mathieson I, Lazaridis I, Rohland N, Mallick S, Patterson N, Roodenberg SA, Harney E, Stewardson K, Fernandes D, Novak M, et al. Genome-wide patterns of selection in 230 ancient eurasians. *Nature* 2015:**528**:499–503. https://doi.org/10.1038/nature16152.
- Mathieson I, McVean G. Estimating selection coefficients in spatially structured populations from time series data of allele frequencies. *Genetics* 2013a:193:973–984. https://doi.org/10.1534/genetics.112.147611.
- Mathieson I, McVean G. Robust identification of local adaptation from allele frequencies. *Genetics* 2013b:**195**:205–220. https://doi.org/10.1534/genetics.113.152462.
- Maynard Smith J, Haigh J. The hitch-hiking effect of a favourable gene. *Genet Res.* 1974:**23**:23–35. https://doi.org/10.1017/S0016672300014634.
- McVicker G, Gordon D, Davis C, Green P. Widespread genomic signatures of natural selection in hominid evolution. *PLoS Genet*. 2009:**5**(5):e1000471. https://doi.org/10.1371/journal.pgen.1000471.
- Mignone F, Gissi C, Liuni S, Pesole G. Untranslated regions of mRNAs. *Genome Biol.* 2002:**3**:1–10. https://doi.org/10.1186/gb-2002-3-3-reviews0004.
- Mladkova N, Kiryluk K. Genetic complexities of the HLA region and idiopathic membranous nephropathy. *J Am Soc Nephrol.* 2017: **28**(5):1331–1334. https://doi.org/10.1681/ASN.2017030283.
- Moritz S, Bartz-Beielstein T. imputeTS: time series missing value imputation in R. R J. 2017:**9**(1):207–218. https://doi.org/10.32614/RJ-2017-009.
- Mughal MR, DeGiorgio M. Localizing and classifying selective sweeps with trend filtered regression. *Mol Biol Evol*. 2019:**36**:252–270. https://doi.org/10.1093/molbev/msy205.
- Mughal MR, Koch H, Huang J, Chiaromonte F, DeGiorgio M. Learning the properties of adaptive regions with functional data analysis. *PLoS Genet.* 2020:**16**:e1008896. https://doi.org/10.1371/journal.pgen.1008896.
- Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. A fine-scale map of recombination rates and hotspots across the human genome. Science 2005;310:321–324. https://doi.org/10.1126/science.1117196.
- Nabi S, Askari M, Rezaei-Gazik M, Salehi N, Almadani N, Tahamtani Y, Totonchi M. A rare frameshift mutation in SYCP1 is associated with human male infertility. *Mol Hum Reprod.* 2022:**28**:gaac009. https://doi.org/10.1093/molehr/gaac009.
- Nicolaisen LE, Desai MM. Distortions in genealogies due to purifying selection and recombination. *Genetics* 2013:**195**:221–230. https://doi.org/10.1534/genetics.113.152983.
- Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, Bustamante C. Genomic scans for selective sweeps using SNP data. *Genome Res.* 2005:**15**:1566–1575. https://doi.org/10.1101/gr.4252305.

- Norton HL, Kittles RA, Parra E, McKeigue P, Mao X, Cheng K, Canfield VA, Bradley DG, McEvoy B, Shriver MD. Genetic evidence for the convergent evolution of light skin in Europeans and east Asians. *Mol Biol Evol.* 2007:24:710-722. https://doi.org/10.1093/molbev/msl203.
- Oleksyk TK, Smith MW, O'Brien SJ. Genome-wide scans for footprints of natural selection. *Philos Trans R Soc Lond B, Biol Sci.* 2010:**365**:185–205. https://doi.org/10.1098/rstb.2009.0219.
- Oseledets IV. Tensor-train decomposition. SIAM J Sci Comput. 2011:33(5):2295-2317. https://doi.org/10.1137/090752286.
- Papastergiou T, Zacharaki El, Megalooikonomou V. Tensor Decomposition for Multiple-Instance Classification of High-Order Medical Data. Computer Engineering and Informatics Department, University of Patras, Rio, Achaia 26504, Greece. 2018.
- Parada H, Sun X, Fleming JM, Williams-DeVane CR, Kirk EL, Olsson LT, Perou CM, Olshan AF, Troester MA. Race-associated biological differences among luminal A and basal-like breast cancers in the Carolina Breast Cancer Study. *Breast Cancer Res.* 2017:19: 131. https://doi.org/10.1186/s13058-017-0914-6.
- Payseur BA, Nachman MW. Microsatellite variation and recombination rate in the human genome. *Genetics* 2000:**156**(3): 1285–1298. https://doi.org/10.1093/genetics/156.3.1285.
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, et al. Scikit-learn: machine learning in Python. *J Mach Learn Res.* 2011:**12**: 2825–2830.
- Peng Y, Yang Z, Zhang H, Cui C, Qi X, Luo X, Tao X, Wu T, Ouzhuluobu, Basang, et al. Genetic variations in Tibetan populations and high-altitude adaptation at the Himalayas. *Mol Biol Evol*. 2011;28:1075–1081. https://doi.org/10.1093/molbev/msq290.
- Pennings PS, Hermisson J. Soft sweeps II: molecular population genetics of adaptation from recurrent mutation or migration. *Mol Biol Evol.* 2006a:**23**:1076–1084. https://doi.org/10.1093/molbev/msj117.
- Pennings PS, Hermisson J. Soft sweeps III: the signature of positive selection from recurrent mutation. *PLoS Genet.* 2006b:**2**:e186. https://doi.org/10.1371/journal.pgen.0020186.
- Pritchard JK, Pickrell JK, Coop G. The genetics of human adaptation: hard sweeps, soft sweeps, and polygenic adaptation. *Curr Biol.* 2010:**20**:R208–R215. https://doi.org/10.1016/j.cub.2009.11.055.
- Przeworski M. The signature of positive selection at randomly chosen loci. *Genetics* 2002:**160**:1179–1189. https://doi.org/10.1093/genetics/160.3.1179.
- Racimo F. Testing for ancient selection using cross-population allele frequency differentiation. *Genetics* 2016:202:733–750. https://doi.org/10.1534/genetics.115.178095.
- Racimo F, Kuhlwilm M, Slatkin M. A test for ancient selective sweeps and an application to candidate sites in modern humans. *Mol Biol Evol.* 2014:**31**(12):3344–3358. https://doi.org/10.1093/molbev/msu255.
- Rambaut A, Pybus OG, Nelson MI, Viboud C, Taubenberger JK, Holmes EC. The genomic and epidemiological dynamics of human influenza A virus. *Nature* 2008:**453**:615–619. https://doi.org/10.1038/nature06945.
- Rees JS, Castellano S, André AM. The genomics of human local adaptation. *Trends Genet.* 2020:**36**:415–428. https://doi.org/10.1016/j.tig.2020.03.006.
- Sabeti PC, Reich DE, Higgins JM, Levine HZP, Richter DJ, Schaffner SF, Gabriel SB, Platko JV, Patterson NJ, McDonald GJ, et al. Detecting recent positive selection in the human genome from haplotype structure. *Nature* 2002:**419**:832–837. https://doi.org/10.1038/nature01140.
- Sabeti PC, Schaffner SF, Fry B, Lohmueller J, Varilly P, Shamovsky O, Palma A, Mikkelsen TS, Altshuler D, Lander ES. Positive natural selection in the human lineage. *Science* 2006:**312**:1614–1620. https://doi.org/10.1126/science.1124309.
- Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH, McCarroll SA, Gaudet R, et al. Genome-wide detection and characterization of positive selection in human

- populations. *Nature* 2007:**449**:913–918. https://doi.org/10.1038/nature06250.
- Sakharkar MK, Chow VTK, Kangueane P. Distributions of exons and introns in the human genome. *In Silico Biol.* 2004:**4**:387–393.
- Salem N, Hussein S. Data dimensional reduction and principal components analysis. *Procedia Comput Sci.* 2019:**163**:292–299. https://doi.org/10.1016/j.procs.2019.12.111.
- Sarkar E, Chielle E, Gürsoy C, Mazonka O, Gerstein M, Maniatakos M. Fast and scalable private genotype imputation using machine learning and partially homomorphic encryption. *IEEE Access*. 2021:9: 93097–93110. https://doi.org/10.1109/ACCESS.2021.3093005.
- Scally A, Durbin R. Revising the human mutation rate: implications for understanding human evolution. *Nat Rev Genet*. 2012:**13**: 745–753. https://doi.org/10.1038/nrg3295.
- Schlebusch CM, Sjödin P, Skoglund P, Jakobsson M. Stronger signal of recent selection for lactase persistence in Maasai than in Europeans. *Eur J Hum Genet*. 2012:**21**(5):550–553. https://doi.org/10.1038/ejhg.2012.199.
- Schrider DR. Background selection does not mimic the patterns of genetic diversity produced by selective sweeps. *Genetics* 2020: **216**(2):499–519. https://doi.org/10.1534/genetics.120.303469.
- Schrider DR, Kern AD. S/HIC: robust identification of soft and hard sweeps using machine learning. *PLoS Genet.* 2016:**12**:e1005928. https://doi.org/10.1371/journal.pgen.1005928.
- Schrider DR, Kern AD. Soft sweeps are the dominant mode of adaptation in the human genome. *Mol Biol Evol*. 2017:**34**(8): 1863–1877. https://doi.org/10.1093/molbev/msx154.
- Schrider DR, Kern AD. Supervised machine learning for population genetics: a new paradigm. *Trends Genet*. 2018:**34**:301–312. https://doi.org/10.1016/j.tig.2017.12.005.
- Scrimshaw NS, Murray EB. The acceptability of milk and milk products in populations with a high prevalence of lactose intolerance. *Am J Clin Nutr.* 1988:**48**(4):1079–1159. https://doi.org/10.1093/ajcn/48.4.1142.
- Seger J, Smith WA, Perry JJ, Hunn J, Kaliszewska ZA, Sala LL, Pozzi L, Rowntree VJ, Adler FR. Gene genealogies strongly distorted by weakly interfering mutations in constant environments. *Genetics* 2010:**184**:529–545. https://doi.org/10.1534/genetics.109.103556.
- Ségurel L, Bon C. On the evolution of lactase persistence in humans. Ann Rev Genomics Hum Genet. 2017:**18**:297–319. https://doi.org/10.1146/annurev-genom-091416-035340.
- Seo EK, Choi JY, Jeong JH, Kim YG, Park HH. Crystal structure of C-terminal coiled-coil domain of SYCP1 reveals non-canonical anti-parallel dimeric structure of transverse filament at the synaptonemal complex. PLoS ONE. 2016:11(8):e0161379. https://doi. org/10.1371/journal.pone.0161379.
- Setter D, Mousset S, Cheng X, Nielsen R, DeGiorgio M, Hermisson J. VolcanoFinder: genomic scans for adaptive introgression. PLoS Genet. 2020:16:e1008867. https://doi.org/10.1371/journal.pgen. 1008867.
- Shah N, Sukumar S. The Hox genes and their roles in oncogenesis. *Nat Rev Cancer*. 2010:**10**(5):361–371. https://doi.org/10.1038/nrc2826.
- Shatin R. Evolution and lactase deficiency. *Gastroenterology* 1968:**54**: 992–993. https://doi.org/10.1016/S0016-5085(68)80176-3.
- Sheehan S, Song YS. Deep learning for population genetic inference. *PLoS Comput Biol.* 2016:**12**:e1004845. https://doi.org/10.1371/journal.pcbi.1004845.
- Sidiropoulos ND, De Lathauwer L, Fu X, Huang K, Papalexakis EE, Faloutsos C. Tensor decomposition for signal processing and machine learning. *IEEE Trans Signal Process*. 2017:65(13):3551–3582. https://doi.org/10.1109/TSP.2017.2690524.
- Simonson TS, Yang Y, Huff CD, Yun H, Qin G, Witherspoon DJ, Bai Z, Lorenzo FR, Xing J, Jorde LB, et al. Genetic evidence for high-altitude adaptation in Tibet. *Science* 2010:**329**:72–75. https://doi.org/10.1126/science.1189406.
- Singhal S, Leffler EM, Sannareddy K, Turner I, Venn O, Hooper DM, Strand AI, Li Q, Raney B, Balakrishnan CN, et al. Stable recombination hotspots in birds. *Science* 2015:**350**:928–932. https://doi.org/10.1126/science.aad0843.

MBF

- Skoglund P, Mathieson I. Ancient genomics of modern humans: the first decade. Annu Rev Genomics Hum Genet. 2018:19:381–404. https://doi.org/10.1146/genom.2018.19.issue-1.
- Slatkin M. Linkage disequilibrium—understanding the evolutionary past and mapping the medical future. *Nat Rev Genet*. 2008:9(6):477–485. https://doi.org/10.1038/nrg2361.
- Smukowski CS, Noor MA. Recombination rate variation in closely related species. *Heredity*. 2011:**107**:496–508. https://doi.org/10.1038/hdy.2011.44.
- Souilmi Y, Tobler R, Johar A, Williams M, Grey ST, Schmidt J, Teixeira JC, Rohrlach A, Tuke J, Johnson O, et al. Admixture has obscured signals of historical hard sweeps in humans. *Nat Ecol Evol*. 2022:**6**: 2003–2015. https://doi.org/10.1038/s41559-022-01914-9.
- Steinwart I, Thomann P. liquidSVM: a fast and versatile SVM package. arXiv 1702.06899, 2017, e-prints.
- Stipoljev S, Bužan E, Rolečková B, Iacolina L, Šprem N. MHC genotyping by SSCP and amplicon-based NGS approach in chamois. *Animals* (*Basel*). 2020:**10**(9):1694. https://doi.org/10.3390/ani10091694.
- Sugden LA, Atkinson EG, Fischer AP, Rong S, Henn BM, Ramachandran S. Localization of adaptive variants in human genomes using averaged one-dependence estimation. *Nat Commun.* 2018:9:703. https://doi.org/10.1038/s41467-018-03100-7.
- Sun X, Liu Y, An L. Ensemble dimensionality reduction and feature gene extraction for single-cell RNA-seq data. *Nat Commun.* 2020:11:5853. https://doi.org/10.1038/s41467-020-19465-7.
- Szpiech ZA, Novak TE, Bailey NP, Stevison LS. Application of a novel haplotype-based scan for local adaptation to study high-altitude adaptation in rhesus macaques. *Evol Lett.* 2021:**5**:408–421. https://doi.org/10.1002/evl3.232.
- Takahata N. Allelic genealogy and human evolution. *Mol Biol Evol.* 1993:**10**(1):2–22. https://doi.org/10.1093/oxfordjournals.molbev. a039995.
- Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, Gagliano Taliun SA, Corvelo A, Gogarten SM, Kang HM, et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed program. *Nature* 2021:**590**:290–299. https://doi.org/ 10.1038/s41586-021-03205-y.
- Terhorst J, Kamm JA, Song YS. Robust and scalable inference of population history from hundreds of unphased whole-genomes. *Nat Genet*. 2017:**49**:303–309. https://doi.org/10.1038/ng.3748.
- Tishkoff SA, Reed FA, Ranciaro A, Voight BF, Babbitt CC, Silverman JS, Powell K, Mortensen HM, Hirbo JB, Osman M, et al. Convergent adaptation of human lactase persistence in Africa and Europe. *Nat Genet*. 2007:**39**:31–40. https://doi.org/10.1038/ng1946.
- Torada L, Lorenzon L, Beddis A, Isildak U, Pattini L, Mathieson S, Fumagalli M. Imagene: a convolutional neural network to quantify natural selection from genomic data. *BMC Bioinformatics*. 2019:**20**:337. https://doi.org/10.1186/s12859-019-2927-x.
- Torres R, Szpiech ZA, Hernandez RD. Human demographic history has amplified the effects of background selection across the genome. *PLoS Genet.* 2018:**14**(6):e1007387. https://doi.org/10.1371/journal.pgen.1007387.
- Tucker LR. Some mathematical notes on three-mode factor analysis. *Psychometrika* 1966:**31**(3):279-311. https://doi.org/10.1007/BF02289464.
- Verleysen M, François D. The curse of dimensionality in data mining and time series prediction. In: Proceedings of the 8th International Conference on Artificial Neural Networks: Computational Intelligence and Bioinspired Systems, IWANN'05. Berlin, Heidelberg: Springer-Verlag; 2005. p. 758–770.

- Vitti JJ, Grossman SR, Sabeti PC. Detecting natural selection in genomic data. *Annu Rev Genet*. 2015:**47**:97–120. https://doi.org/10. 1146/annurev-genet-111212-133526.
- Voight BF, Kudaravalli S, Wen X, Pritchard JK. A map of recent positive selection in the human genome. *PLoS Biol.* 2006:**4**:e72. https://doi.org/10.1371/journal.pbio.0040072.
- Vy HM, Kim YA. A composite-likelihood method for detecting incomplete selective sweep from population genomic data. *Genetics* 2015:**200**:633–649. https://doi.org/10.1534/genetics.115.175380.
- Wang B, Zhang Y-B, Zhang F, Lin H, Wang X, Wan N, Ye Z, Weng H, Zhang L, Li X, et al. On the origin of Tibetans and their genetic basis in adapting high-altitude environments. *PLoS ONE*. 2011:**6**:e17002. https://doi.org/10.1371/journal.pone.0017002.
- Wei C, Wang H, Liu G, Zhao F, Kijas JW, Ma Y, Lu J, Zhang L, Cao J, Wu M, et al. Genome-wide analysis reveals adaptation to high altitudes in Tibetan sheep. Sci Rep. 2016:**6**:26770. https://doi.org/10.1038/srep26770.
- Whitehouse LS, Schrider DR. Timesweeper: accurately identifying selective sweeps using population genomic time series. *Genetics* 2023:**224**:iyad084. https://doi.org/10.1093/genetics/iyad084.
- Wilde S, Timpson A, Kirsanow K, Kaiser E, Kayser M, Unterländer M, Hollfelder N, Potekhina ID, Schier W, Thomas MG, et al. Direct evidence for positive selection of skin, hair, and eye pigmentation in Europeans during the last 5,000 years. *Proc Natl Acad Sci U S A*. 2014:**111**(13):4832–4837. https://doi.org/10.1073/pnas.1316513111.
- Wright MN, Ziegler A. ranger: A fast implementation of random forests for high dimensional data in C++ and R. *J Stat Softw*. 2017:**77**(1):1–17. https://doi.org/10.18637/jss.v077.i01.
- Xu S, Li S, Yang Y, Tan J, Lou H, Jin W, Yang L, Pan X, Wang J, Shen Y, et al. A genome-wide search for signals of high-altitude adaptations in Tibetans. *Mol Biol Evol*. 2011:**28**:1003–1011. https://doi.org/10.1093/molbev/msq277.
- Xue AT, Schrider DR, Kern AD, Ag1000g Consortium. Discovery of ongoing selective sweeps within anopheles mosquito populations using deep learning. Mol Biol Evol. 2021;38:1168-1183. https://doi.org/10.1093/molbev/msaa259.
- Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE, Xu X, Jiang H, Vinckenbosch N, Korneliussen TS, et al. Sequencing of 50 human exomes reveals adaptation to high altitude. *Science* 2010:**329**: 75–78. https://doi.org/10.1126/science.1190371.
- Yuwang J, Qiang W, Xuan L, Jie L. A survey on tensor techniques and applications in machine learning. IEEE Access. 2019:7: 162950–162990. https://doi.org/10.1109/ACCESS.2019.2949814.
- Zare A, Ozdemir A, Iwen MA, Aviyente S. Extension of PCA to higher order data structures: an introduction to tensors, tensor decompositions, and tensor PCA. *Proc IEEE*. 2018:**106**(8):1341–1358. https://doi.org/10.1109/JPROC.2018.2848209.
- Zhang W, Fan Z, Han E, Hou R, Zhang L, Galverni M, Huang J, Liu H, Silva P, Li P, et al. Hypoxia adaptations in the grey wolf (*Canis lupus chanco*) from Qinghai-Tibet Plateau. *PLoS Genet*. 2014:**10**: e1004466. https://doi.org/10.1371/journal.pgen.1004466.
- Zhang X, Witt KE, Mañuelos MM, Ko A, Yuan K, Xu S, Nielsen R, Huerta-Sanchez E. The history and evolution of the Denisovan-EPAS1 haplotype in Tibetans. *Proc Natl Acad Sci U S A*. 2021:**118**: e2020803118. https://doi.org/10.1073/pnas.2020803118.
- Zhang G, Xu Y, Wang S, Gong Z, Zou C, Zhang H, Ma G, Zhang W, Jiang P. LncRNA SNHG17 promotes gastric cancer progression by epigenetically silencing of p15 and p57. *J Cell Physiol.* 2019:**234**: 5163–5174. https://doi.org/10.1002/jcp.v234.4.