Article

https://doi.org/10.1038/s41467-025-56273-3

STICI: Split-Transformer with integrated convolutions for genotype imputation

Received: 23 April 2024

Accepted: 8 January 2025

Published online: 31 January 2025

Check for updates

Mohammad Erfan Mowlaei¹, Chong Li[®]¹, Oveis Jamialahmadi[®]², Raquel Dias³, Junjie Chen⁴, Benyamin Jamialahmadi⁵, Timothy Richard Rebbeck ^{® 6.7}, Vincenzo Carnevale ^{® 8,9}, Sudhir Kumar ^{® 1,8,10} & Xinghua Shi ^{® 1,8} ⊠

Despite advances in sequencing technologies, genome-scale datasets often contain missing bases and genomic segments, hindering downstream analyses. Genotype imputation addresses this issue and has been a cornerstone pre-processing step in genetic and genomic studies. Although various methods have been widely adopted for genotype imputation, it remains challenging to impute certain genomic regions and large structural variants. Here, we present a transformer-based framework, named STICI, for accurate genotype imputation. STICI models automatically learn genome-wide patterns of linkage disequilibrium, evidenced by much higher imputation accuracy in regions with highly linked variants. Our imputation results on the human 1000 Genomes Project and non-human genomes show that STICI can achieve high imputation accuracy comparable to the state-of-the-art genotype imputation methods, with the additional capability to impute multi-allelic variants and various types of genetic variants. STICI can be trained for any collection of genomes automatically using self-supervision. Moreover, STICI shows excellent performance without needing any special presuppositions about the underlying patterns in collections of non-human genomes, pointing to adaptability and applications of STICI to impute missing genotypes in any species.

Genetic and genomic studies, such as linkage analysis, genome-wide association study (GWAS), and polygenic risk score (PRS) estimation, enable us to dissect the genetic architecture of complex traits and diseases¹. In recent years, whole genome sequencing (WGS) platforms and genotyping techniques have become increasingly cost-effective, resulting in the accumulation of large collections of genotypes and deeper insights into the genetic architecture of individual traits and diseases. The resolution of genotyping has been improving due to advances in genotyping and sequencing technologies, accompanied by algorithmic improvement in analyzing such data. However, genotype data still contains missing values and untyped loci². Such missing genotypes may decrease the statistical power in disease association studies and causal variant discovery³⁻⁵. Causes of missing genotypes include the difficulty in sequencing rare alleles⁶⁻⁸, failure of experimental assays, genotype calling errors, and differences in densities and properties of

¹Computer & Information Sciences, College of Science and Technology, Temple University, Philadelphia, PA, USA. ²Department of Molecular and Clinical Medicine, Institute of Medicine, Sahlgrenska Academy, Wallenberg Laboratory, University of Gothenburg, Gothenburg, Sweden. ³Department of Microbiology and Cell Science, University of Florida, Gainesville, FL, USA. ⁴School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen, Guangdong, China. ⁵David R. Cheriton School of Computer Science, University of Waterloo, ON, Canada. ⁶Division of Population Sciences, Dana-Farber Cancer Institute, Boston, MA, USA. ⁷Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, MA, USA. ⁸Institute for Genomics and Evolutionary Medicine, Temple University, Philadelphia, PA, USA. ⁹Institute for Computational Molecular Science, Temple University, Philadelphia, PA, USA. ¹⁰Department of Biology, Temple University, Philadelphia, PA, USA. ¹⁰Computer Science, USA.



a) Sporadic missingness

Fig. 1 | Two categories of genotypes missingness. a Sporadic missingness: It arises due to genotype calling errors and assay failures. Prediction of sporadic missingness is typically done during the pre-phasing step of imputation pipelines.
b Systematic missingness: Differences in sequencing resolution are common



b) Systematic missingness

causes of systematic missingness because a subset of genomic positions are assayed. The inference of missing variants in untyped regions is a major focus of imputation pipelines.

genotyping platforms³. As such, genotype missingness, as depicted in Fig. 1, can be classified into two distinct categories: sporadic missingness, where for each site/segment, some values could be absent; and systematic missingness, in which some genomic loci or segments are not genotyped. Challenges in handling missing data are further compounded when considering different types of genetic variation in addition to Single Nucleotide Variants (SNVs). Compared with SNVs, Structural Variations (SVs) pose a greater challenge in genotype calling and imputation due to their increased complexity, limitations of current sequencing technologies, extensive allelic diversity, and their variable frequencies among populations^{9,10}. Moreover, SVs can have a more significant impact on genetic diseases than SNVs, so their accurate imputation can lead to enhancements in disease association studies¹¹.

Consequently, there is a common need for reliable imputation of genotypes using computational methods. Imputation is the process of inferring missing values in the data based on the information already present in the dataset, such as the density and distribution of variants of various complexities within and among sequences in the dataset. Imputation of missing data in genomics needs specialized methods because genomic information is inherently different from data in many other domains, such as vision or natural language processing. For example, genotype data have some unique characteristics, including high data dimensionality when the number of variants is usually larger than the sample size, linear and non-linear correlations among the variants¹², and shared segments of sequences due to common ancestries. Furthermore, given the multifaceted nature of genotype imputation, it is well-recognized that no single method can serve as a universal solution. Consequently, it is a common practice in the field to utilize multiple imputation tools for a particular study.

Widely used imputation methods often require a reference panel of genomes to impute missing genotypes in other genomes, which assumes that the missing information comes from the same ancestry patterns as those in the reference panel¹³. These methods utilize Hidden Markov Models (HMMs), graphical models, and haplotype-cluster algorithms to impute missing values¹⁴. For example, Minimac4¹⁵, the most recent version of MACH¹⁶, uses an HMM. For each individual, Minimac4 updates the phase iteratively in both directions based on haplotypes in the reference panel and neighboring loci in the individual. It splits sequences into overlapping chunks in order to reduce memory consumption and make the model scalable. Similarly, SHAPEIT5¹⁷, Impute5¹⁸, and Beagle¹⁹ also employ HMMs to perform imputation. The Haplotype-clustering algorithm is utilized in the fastPHASE²⁰ in order to cluster haplotypes in an SNV-wise manner and impute missing values per locus. GLIMPSE2²¹ uses an HMM to genotype low-coverage whole-genome sequencing (WGS) data.

Deep Learning (DL) methods have lately been applied for genotype imputation. Sparse Convolutional Denoising autoencoder (SCDA) is used in ref. 14 to impute missing data in the Human Leukocyte Antigen (HLA) region on chromosome 6 and yeast²² genotypes. In ref. 3 an improvement in SCDA training is proposed to improve the model's performance. Similarly²³, used autoencoders on identified linkage disequilibrium (LD) blocks as well as focal loss to improve the performance. RNN-IMP²⁴ utilizes recurrent neural networks (RNNs) and augments the samples using recombination and mutation in order to impute systematic missingness in genotype data. In a follow-up study²⁵, this model was updated with a complementary denoising autoencoder to perform pre-phasing in addition to systematic imputation. GRUD²⁶ utilizes RNN-IMP in an adversarial training schema which substantially reduces the model training time. DEEP*HLA²⁷ uses a convolutional neural network to perform imputation on pre-phased genotypes at the gene level. Inspired by DEEP*HLA, HLARIMNT²⁸ uses transformers for the same task.

While the overall performance of existing imputation methods for genomic data is generally good for common variants and alleles with high-frequency alleles, it is still challenging to impute rare variants, especially in scenarios where reference panels are limited. Also, many current solutions cannot directly handle multi-allelic variants. Moreover, existing methods have not been systematically benchmarked for imputing SVs, particularly complex variants with multiple alleles, to the best of our knowledge. Also, the training of DL models^{3,14,23,27} is generally slow, and they usually need many more samples to perform as well as classical methods based on HMMs, such as Minimac415 and Beagle5.4¹⁹. Although RNN-IMP and GRUD^{24,26} addressed this performance disparity, these methods require retraining when the variant sets in the target are different from those in the training set. Although GRUD has the benefit of faster model training, it is not designed to handle sporadic missingness (Fig. 1a). Finally, the rest of the existing DL methods rely on convolutional neural networks (CNNs), which excel at exploiting local patterns but do not exhibit a robust mechanism to effectively capture pairwise correlations among local and distant markers simultaneously, such as the presence of LD blocks in genotypes.

Compared with CNNs, the emerging attention mechanism in Transformers offers an effective solution to capturing local and distal interactions²⁹ in genomes. The attention mechanism in DL mimics visual attention to focus on specific parts of pictures^{30,31} by calculating importance scores among genomic loci. Therefore, attention can capture global interactions amongst markers. Transformers utilize multi-head attention to capture intricate and multi-level interactions among the variants. AlphaFold2³² and ESM-2³³ are successful examples of transformers in biological sequence analysis.

In this article, we present a novel genotype imputation model, STICI, based on the attention mechanisms in a transformer framework and convolutional blocks. Our model utilizes attention to capture patterns among SNVs and SVs in the genome collections analyzed. We found STICI to achieve high imputation accuracy at a modest memory consumption cost, achieved by dividing the data into chunks (following¹⁵) that enable efficient application of STICI to long sequences. Furthermore, STICI needs to be trained only once, unlike other DL models. After this training, the imputation time in STICI is comparable to or faster than using classical methods (Supplementary Tables 10 and 11). In brief, our study makes the following key contributions.

- We propose a DL transformer framework termed STICI, designed to specifically address the genotype imputation problem.
- STICI imputation does not need a standard reference panel, which makes it more generally applicable to any data set.
- STICI excels at SV imputation, where the variants harbor a higher degree of complexity while achieving performance comparable to competing models for SNV imputation.
- We analyze the effect of different masking rates on building better imputation models and explain the reasons for STICI's performance.

Results

Overview of the study

In this section, first, we present the results of our empirical study to find the optimal masking percentage (masking rate, MaskR) for STICI training in order to investigate the need for building imputation models for specific missing rate (MissR) in the target dataset. After that, we present results for sporadic missingness imputation on the yeast dataset, SVs in human chromosome 22, and extensive SVs dataset (see Data in Methods for the dataset details). This is followed by systematic missingness imputation for selected regions in human chromosome 22, simulated human chromosome 19. rat, and chicken datasets. We benchmarked STICI's performance against classical pre-phasing methods (SHAPEIT5 and Eagle2³⁴), classical imputation methods (Impute5, Beagle4, Beagle5.4, Minimac3, and Minimac4.1.4), deep learning models (SCDA and AE), and two variants of STICI, one that uses no embedding (STICI-NE) and another that uses an additional MaCH-Rsq loss term (STICI-Rsq). The details of STICI architecture and aforementioned methods are discussed in STICI architecture and Baseline models, respectively.

Optimal masking rate analysis

For this analysis, we used the HLA region on chromosome 6 from the human 1000 Genomes Project and performed a 3-fold cross-validation on the data. The aim was to examine the relationship of MaskR for the training set with varying MissR in target sets. The results are presented in Fig. 2 in which the results on the left/right column belong to the validation and test set, respectively.

Figure 2a, b shows that the performances of STICI models trained using MaskR of 0.5 and 0.7 were high for imputations in which MissRs were up to 0.5. Therefore, a single STICI model, trained with a MaskR of 0.5, could be used for a variety of research datasets as long as the MissR is less than 0.5. However, STICI models trained with MaskR > 0.5 are needed for reliable imputations when the target datasets have more than 50 percent missing variants. Therefore, we recommend STICI models trained with a MaskR of 0.5 for imputing sporadically missing variants and a higher MaskR (preferably matching the MissR in the target dataset indicated by our empirical studies data) for other datasets.

Figure 2a, b also shows that when the target MissR is sufficiently low, the performance gap of the imputation models is not discernible. The performance gap becomes evident with a MissR of 0.2 or higher. The underlying cause of this observation is that when the MissR is significantly low, a sufficient number of variants in LD with the target variant are readily available, making predictions less challenging for all the models. Conversely, a large MissRs means that the amount of information from LD blocks diminishes, presenting a greater challenge to the imputation model.

Figure 2c, d shows that, generally, STICI models trained with lower MaskR will produce poor performance for imputing missing SNVs located in regions with high LD. For instance, variants in regions with LD = 0.01 have the lowest accuracy for all the masking percentages. In addition, these results indicate that it is easier to predict missing data in high LD regions compared to low LD regions, which aligns well with biological expectations that low LD regions do not benefit from additional information (LD) available for better imputation of high LD regions. These trends suggest that the use of a low MaskR prevents the model from learning LD patterns, resulting in a worse performance. In other words, the model training needs to effectively disturb the LD blocks (and other latent patterns among variants) to capture direct and indirect correlations and haplotypes. Consequently, MaskR of 0.5 and higher provides robust results across a large range of target MissR values.

The relative performance of STICI for sporadic missingness

For each dataset in this experiment, we performed a 3-fold crossvalidation where missing values in the test were introduced using fixed random seeds to ensure reproducibility of results across experiments and methods, i.e., every method imputed the same test set. To ensure that all the methods have access to the same set of information for imputation, in each split of the cross-validation, the training set is used to train deep learning models including STICI, while the same training set is also used as a reference panel for Minimac4, Beagle5.4, SHA-PEIT5, and Eagle2.4.1. We verified that no two samples are identical in a training set and the corresponding test set for a given split. However, there are potentially shared haplotypes among the training and the test sets, specially in strong LD blocks, which provide information for making correct imputations. Shared haplotypes may introduce potential information leakage in the imputation process for secure and privacy-preserving genotype imputation, which will be a future direction to explore regarding secure and privacy-preserving genotype imputation.

The missing values were distributed randomly according to one of three strategies: uniformly, based on Minor Allele Frequency (MAF), or based on LD. These methods were chosen to ensure that missing values are representative of the data distribution in different biological aspects. In all of the experiments, missing positions in the test sets were the same for all the methods. Further details on these procedures can be found in the methods section.

The overall results for the yeast and chromosome 22 datasets are presented in Supplementary Table 7. The numerical values in this table indicate the average of the metric values on the test sets in a 3-fold cross-validation. We used maximum LD bins and/or MAF bins (Fig. 3b, c) to distribute missing positions in the datasets extracted from the human 1000 Genomes Project. If the bins had too few positions (e.g., at a 0.01 MissR on chromosome 22 datasets), we excluded this MissR for the experiments related to these datasets. We used a consistent approach to introduce missing values in chromosomes 6, 10, 16, and 20 based on LD distributions and a single test MissR of 0.2. In this experiment, we focused on comparing Minimac4.1.4, Beagle5.4, SHA-PEIT5, and STICI, because they were identified as the top performers from classical and DL methods in prior experiments. We also added Eagle2³⁴ to the competing methods in this experiment. We employed 3-fold cross-validation for both methods, training and imputing each chromosome separately. R² was calculated for each variant, and the results were averaged over each fold, chromosome, and SV type. Figure 4 presents the experimental results for the extensive structural



Fig. 2 | **Average accuracy over 3-fold cross-validation for validation and test sets in the HLA dataset using different masking rate (MaskR) values during STICI training.** Bars indicate a 95% confidence interval per experiment. **a**, **b** A breakdown of average accuracy for various missing rate (MissR) values of validation/test set when the model is trained using different MaskR values. The patterns show that a model trained using a higher MaskR is more robust across different target MissRs. **c**, **d** Average accuracy for validation/test sets over 3 folds and different MissR values calculated for various LD bins. The trend suggests that a higher

variation datasets where the top plot shows the improvement that STICI provides compared to the best of other methods for each SV type, and is calculated as follows:

Improvement(%) =
$$\frac{R_{\text{STICI}}^2 - R_{\text{Best}}^2}{R_{\text{Best}}^2} * 100$$
 (1)



MaskR increases the performance across LD bins, which could be attributed to the impact of MaskR on STICI to learn LD patterns comprehensively. When MaskR is low, STICI imputations do not benefit from the LD patterns present and thus, STICI does not learn the majority of pairwise correlations (LD) among the variants. Consequently, STICI is not able to infer the missing value using all possible information in the respective LD block of the target variant. Source data are provided as a Source Data file.

Yeast dataset. Missing positions in samples were selected randomly, as the LD analysis showed that the maximum LD for all the SNVs was high in the [0.8, 1.0] range. As mentioned, Minimac4.1.4, SHAPEIT5, and Beagle5.4 cannot be used to impute variants of the yeast dataset due to the lack of a reference panel. However, STICI could be applied and outperformed other methods, achieving a minimum average imputation accuracy of 99.86%. Overall, all the applicable models



Fig. 3 | MAF and LD distributions of benchmark datasets from the 1000 Genomes Project. MAF and maximum LD distributions are presented using kernel density estimation plots for SNVs and SVs in (a). HLA region on chromosome 6, (b) deletions in chromosome 22, (c) SVs in chromosome 22, (e) SVs in chromosome 6, (f) SVs in chromosome 10, (g) SVs in chromosome 16, and (h) SVs in chromosome 20. Overall, SVs exhibit a low LD value, posing a significant challenge to imputation

methods. Plot (**d**) LD among different SV types in chromosome 22 shows that structural events are commonly correlated with deletions. Furthermore, deletion, copy number variation, and duplication events appear in different ranges of LD, while the rest of the events are limited to LD \leq 0.1. Lastly, the majority of correlated SVs to deletions are of the same event, making deletions a good separate dataset for our experiment. Source data are provided as a Source Data file.

performed well on the yeast dataset, which we attribute to the presence of high LD among SNVs in this dataset.

Deletions in chromosome 22. For this dataset, we introduced missing positions proportional to the maximum LD/MAF distribution Fig. 3b. Overall, STICI emerged as the best or the second-best model for imputation across all the metrics. STICI was more accurate than others for LD/MAF missingness distribution schema. Furthermore, SCDA + demonstrates a substantial performance advantage over AE in terms of IQS and R² in the majority of the cases. Supplementary Table 8 shows the accuracy trends for different maximum LD values for this dataset when missing values are distributed proportional to variant density in maximum LD bins. The results in this table indicate the performance of the competing approaches in the presence of shared haplotypes

among the reference panel (training set) and the test set, as the occurrence of similar haplotypes across the samples is expected to be higher in stronger LD blocks. Minimac4.1.4 and Beagle5.4 were less accurate for SNVs with lower maximum LD compared to AE, SCDA +, STICI-NE, and STICI. Since HMMs and graphical models rely on conditional probabilities, we suggest that they would perform relatively weak due to a low correlation between the events (states).

All SVs in chromosome 22. Similar to the previous dataset, missing positions were distributed among SVs based on maximum LD/MAF (Fig. 3c). Despite having a reference panel, Minimac4.1.4 and SHAPEIT5 cannot be directly used for some missing variants for this dataset because they can only handle bi-allelic events. Furthermore, IQS is not well-defined for multi-allelic events.



Fig. 4 | **Comparison sporadic imputation results of competing methods across SV types.** Average R² of ground-truth genotypes in the test sets and respective predictions over 3-fold cross-validations on chromosomes 6, 10, 16, and 20. The experiments are performed on each chromosome separately, and the results are averaged over chromosomes and folds. Vertical lines indicate standard deviations.

The improvement plot shows the R² score difference between STICI and the best of other methods, normalized by the best R² scores for each SV type. We only report biallelic imputation results for SHAPEIT5 because we faced issues with imputing normalized multi-allelic variants using this software. Source data are provided as a Source Data file.

Supplementary Table 7 shows that STICI outperforms all other methods on average accuracy and F1-score. STICI performance in terms of R^2 is much better than the competing methods at high MissRs. R^2 considers the correlation among genotypes encoded as categorical values. As such, depending on the difference in encoded values for the predicted and the ground truth genotypes, the penalty can be severe. For example, if O|O, O|1, and 1|1 are encoded as O, 1, and 2 in genotypes and the ground truth for a given genotype is O|O, the model is punished moderately(severely) for predicting O|1(1|1). In addition, SCDA + outperforms AE in most comparisons, indicating the effectiveness of our proposed training procedure.

Extensive structural variation datasets. In this experiment, we focus on R² between the predicted and ground truth genotypes as R² was the most discriminating metric for comparing the performance in imputing SVs. For estimating R², predictions are converted into categorical values, e.g., 0|0, 0|1, and 1|1 are encoded as 0, 1, and 2. Any discrepancy between the model's prediction and the ground truth leads to a substantial penalty on the correlation, enabling us to see differences more clearly. STICI consistently outperformed Beagle5.4 and Minimac4 across various SV types, often by a noticeable margin. The underlying cause of this observation is the lack of high LD in this dataset (Fig. 3) and fundamental differences between HMMs and the Transformer model. In HMMs, information propagation between two distant variants occurs sequentially through intermediate sites. However, this mechanism falters when the LD block is sparse, leading to reduced performance. In contrast, STICI employs a direct variant-to-variant attention mechanism within each chunk without needing to model an intermediate site, which effectively mitigates the limitations posed by a weak LD. Furthermore, the multi-head attention mechanism equips STICI to discern higher-order and complex patterns among variants, which appear to be crucial for better imputations in the absence of strong LD patterns. These capabilities highlight STICI's superiority in managing SV imputation challenges where traditional HMM-based approaches may be suboptimal. This is particularly the case for duplications (DUP) and insertions (INS) where STICI is able to attain a very high R² value. This observation matches our expectations since these two types of SVs are relatively challenging in genotype calling as well³⁵.

The relative performance of STICI for systematic missingness

In order to evaluate STICI against the competing methods for systematic missingness imputation, we curated four datasets that are missing ~ 90% of the variants in the test set. The first dataset contains the Infinium Omni 2.5 BeadChip microarray dataset on human chromosome 22 (12,725 variants) as the test set and WGS genotypes from 1000 genomes project of the same region (99,314 variants) as the reference panel. We used the same individuals as ref. 24 (100 samples from various populations) for the test set and the rest (2404) for the reference panel. The second dataset was generated using stdpopsim³⁶ using msprime simulation engine³⁷. There were 45,000 samples with 30,720 variants on human chromosome 19 in the reference panel and 5000 samples with 3044 variants for the same region in the test set. The third dataset contains 5147 samples on a selected region on rat^{38,39} chromosome 20 with 61,440 variants as the reference panel and 1000 samples with 6140 variants scattered thorough the reference panel variants. The fourth dataset consists of 2258 reference samples and 55,255 variants on Sasso chicken⁴⁰ chromosome 20 and 100 test samples with 5488 variants selected among the reference variants. More details about these datasets are provided in Data.

We used accuracy, Impute info score (INFO score)⁴¹, MaCH-Rsq¹⁵, and Pearson correlation coefficient R² as evaluation metrics. The reason we included additional metrics for these experiments was that we would like to utilize multiple metrics so that the evaluation of model performance is more comprehensive and less biased. For example, accuracy is not considered a good metric for highly imbalanced data. In this case, accuracy for the variants with rare alleles is misleading because a method that always predicts the majority allele can retain a high accuracy. INFO score indicates the certainty of a model for alternative allele dosage prediction. In an imputation pipeline, the INFO score is used to discard unreliable predictions. The MaCH-Rsq (Equation (6)) metric is designed based on the idea that poorly imputed genotypes will shrink towards their expectations based on population allele frequencies at a given site. In imputation pipelines, MaCH-Rsq is used as a quality control measure in an imputation pipeline and variants with a MaCH-Rsq of less than 0.3 or 0.4 are dropped^{42,42,43}. In addition, we calculated the Pearson correlation coefficient R² (Supplementary Equation 6) of the alternative allele dosages. The results for this metric are available in the Supplementary Fig. 2.

During these experiments we noticed that the original implementation of STICI is not as accurate as classical models for rare alleles. To alleviate this problem, we developed a variant of STICI, namely STICI-Rsq, that used MaCH-Rsq as a new loss term to be used to improve imputation (more details in Loss function). We trained STICI and STICI-Rsq, respectively, using a random MaskR (per sample seen in each training iteration) between 0.85 and 0.95. In other words, in each epoch, the model would see various MaskRs in the aforementioned range for different samples. We find this masking strategy more useful because in real human data, different segments had different MissRs, but the average was around 0.9 MissR.

The experimental results are presented in Fig. 5, where each row is dedicated to the results of one dataset (experiment), and columns show accuracy, INFO score, and MaCH-Rsq from left to right, respectively. While STICI-Rsq achieved high accuracy (> 0.95), high info score (> 0.93), and high MaCH-Rsq (> 0.99), some other methods produced almost perfect results for the simulated human data unlike that seen for the real data. We also investigate the performance of all the methods on non-human genotype data, including a data set from another mammal (rat) and a bird (Sasso chicken). For the mentioned non-human genomes, STICI-Rsq performed comparable to the other methods. We also recorded the time taken to impute for the competing methods. These timings can be found in Supplementary Tables 10 and 11.

Discussion

More accurate genotype imputation will improve the performance of downstream functional and biomedical genomic studies. Scientists frequently need to employ multiple tools, adapted based on the degree of missingness and types of variants missing, within individual pipelines to carry out genotype imputation. To address this problem, we have presented STICI, a masked DL framework, which appears to be one of the first uses of transformer architecture in imputing genomic data. While STICI is currently limited in a few ways, we believe that it represents a step toward developing a unified approach for successfully imputing missing values for a range of datasets, from small to large amounts of missingness, as well as SNVs and SVs. We explored STICI's performance for various masking rates (training) and missing rates (application). Our experiments revealed that a single STICI model, trained with a masking rate of 0.5, could be applied for imputing sporadic missingness of SNVs and SVs, while for systematic missingness which generally constitutes higher missing rates (above 80%), we can employ previously trained models using a masking rate similar to the missing rate in the test set. STICI's performance in imputing SNVs and SVs was comparable to many other methods and approaches for SNVs and SVs found in low and high LD regions. That is, STICI is capable of effectively capturing short and long-range correlations among different genetic variants for genotype imputation.

Currently, training an STICI model is still resource- and timedemanding as transformers are computationally intensive. Hence, the performance evaluation of STICI on large biobank datasets remains unexplored due to constraints of computing and data access. In the case of imputation of synthetic data STICI cannot impute synthetic data as well as the HMMs based on Li and Stephens model, although the performance gap is very small. More extensive synthetic datasets will be further explored in future work, along with the improvement of STICI training on larger and real datasets. To make STICI more applicable in real-world applications and scenarios, in the future, we plan to scale up STICI so that it can be trained on all of the available human reference panels stratified by populations or other features, just as used for classical settings of genotype imputation. This trained STICI model can then be readily used to impute missing genotypes on the samples in real-world datasets. Moreover, the transformer blocks utilized in STICI are known to require a large number of samples to achieve optimal performance. Therefore, we expect STICI's performance to greatly improve even for variants in low LD regions as the number of samples is increased due to the release of new cohorts and data sources. For example, we plan to explore the potential of improving STICI utilizing current and emerging genotypes for an increasing number of individuals in projects like UK Biobank⁴⁴, Trans-Omics for Precision Medicine (TOPMed)⁴⁵, and All of Us⁴⁶. In addition, we expect the STICI framework to spur the next generation of approaches to advance generalized and efficient imputing, which could serve as online imputation servers because of the low computational burden of imputation after intensive training.

Methods

In this section, we first introduce the datasets we used in this study and discuss their characteristics. This is followed by the architectural design of STICI and the procedure for model training.

Data

We used eight datasets from four sequencing projects^{10,22,38-40} and simulated a human genotype dataset using stdpopsim³⁶ with msprime simulation engine³⁷ to fine-tune and benchmark STICI against the baselines. The 1000 Genomes project datasets are pre-phased using SHAPEIT2. Thus, the phasing information of the test sets is propagated to the training sets, but the process is identical for all the methods so it will not bias the results in favor of any imputation method. Scikit-allel package⁴⁷ is employed to compute LD and MAF for the datasets. In the sporadic missingness experiments, we use 3-fold cross-validation to assess the performance of the methods. In a 3-fold cross-validation, each dataset is separated into three distinct partitions where there is no sample overlap. Each time, one of the partitions is used as the test set while the remaining two partitions are used for the training process. For the DL models, a validation set is selected from the training samples for early stopping. To ensure that the same training/validation/test set is used across different methods, we used fixed random seeds for splitting the data into folds and introducing missingness into the test sets. For systematic missingness experiments, we selected a fixed number of individuals as the test set. We used Python scripts to ensure that no parental strands (haploids) are shared among the training and test sets in both sporadic and systematic experiments.

For sporadic missingness imputation, we used three schemes to introduce sporadic missingness: random selection, MAF-distributed





from left to right respectively contain accuracy, INFO score, and MaCH-Rsq results. The lines show the average of the metrics, while the bars around each line indicate a 95% confidence interval. Source data are provided as a Source Data file.

selection, and LD-distributed selection. For the latter two, we computed MAF and LD on the whole data, as depicted in Fig. 3, and selected the missing SNVs/SVs for the test set evaluation proportional to those. That is, if 10% of the variants have an MAF in the range of [0.2, 0.3), we selected 10% of the missing values from these specific variants. This approach ensures that the missing data are imputed based on the distribution of MAF or LD in the data, providing a representative imputation strategy. Regardless of the scheme, we used fixed random

seeds per sample to decide the missing genotypes. Therefore, the missing genotypes across the test samples are not identical. However, the coordinates of missing values are identical across the competing methods.

For systematic missingness imputation datasets, excluding Omni 2.5 BeadChip microarray imputation, we decided on a missing rate and randomly selected and removed the variants in the following predefined MAF bins:

The characteristics of the datasets we used in our experiments are as follows:

HLA dataset. This dataset contains human leukocyte antigen genotypes, covering a 3 Mbp region at chromosome 6p21.31 and sitting at a major histocompatibility complex (MHC) region. HLA region regulates the immune system in humans⁴⁸. It is highly polymorphic and heterogeneous among individuals; i.e., it harbors various alleles, enabling the adaptive immune system to be fine-tuned⁴⁹. In this study, we used the genotypes of this region, obtained from phase 3 of the 1000 Genomes Project¹⁰, which contained 7161 unique genetic variants for 2504 individuals from five super-populations across the world: American (AMR), East Asian (EAS), European (EUR), South Asian (SAS), and African (AFR). The majority of SNVs in this dataset exhibit maximum LD values in the range of [0.9, 1.0]. We used this dataset for our masking study and fine-tuning the hyper-parameters of SCDA+, AE, and STICI.

Yeast dataset. The second dataset is the comprehensively assayed yeast dataset²², representing a simple genetic background and high correlation among genotypes. This dataset contains 4390 genotyped profiles for 28,220 genetic variants. The samples were obtained by sequencing crosses between two strains of yeast, namely an isolate from a vineyard (RM) and a popular laboratory strain (BY). In the original dataset, the data is encoded as -1/1 for BY/RM, which are mapped to 0/1 in our code, respectively, before one-hot encoding.

Chromosome 22 datasets. We used structural variation data from the 1000 Genomes Project in two settings. In the first, we only selected deletions (DEL), excluding ALU/SVA/LINE1 deletions, among all SVs. This resulted in 573 positions harboring bi-allelic events in the dataset. In the second, a total of 848 SVs including, but not limited to deletions, insertions, duplications, inversions (INV), and copy number variations (CNV) in chromosome 22 are selected. As shown in Fig. 3b, c the majority of SVs in chromosome 22 exhibit a low LD, rendering these datasets challenging for imputation compared to SNVs. According to Fig. 3d, deletions cover a wide range of LD among them and other SVs, making them a good target for a separate bi-allelic dataset.

Extensive structural variation datasets. In the concluding experiment, we undertook a thorough investigation of SV imputation using the human 1000 Genomes Project, selecting 4187, 3126, 2062, and 1569 SVs located in chromosomes 6, 10, 16, and 20, respectively. This selection strategy was informed by the aim to encompass chromosomes of different lengths, providing a representative cross-section of the genome. This diverse chromosome selection allows for a broader understanding of the genomic distribution and characteristics of SVs, facilitating a more nuanced analysis of their presence and impact across different regions of the human genome. These SVs include deletions, duplications, insertions, inversions, and copy number variations. Among these SVs, 469 of them are multi-allelic (CNVs). For each model, we train on and impute each chromosomes, and SV type. The distributions of SVs in these chromosomes in terms of MAF and LD

are presented in Fig. 3e-h, indicating low LD and diverse MAF in general for the mentioned SV datasets.

Systematic missingness imputation. For the first dataset (Fig. 5a–c), we used the SNVs (MAF > 0.01) in the chromosome 22 dataset from the human 1000 genomes project, used PLINK2⁵⁰ and bcftools⁵¹ to preprocess the data and converted multi-allelic events, and selected the first 99,314 variants in this chromosome. Following the instructions for data collection in ref. 24 and a script provided by the authors, we created a microarray dataset using Infinium Omni 2.5 BeadChip manifest for chromosome 22 containing 12,725 variants in the same region. We followed²⁴ for selecting the exact individuals for the microarray data (test set), and the rest for the reference panel.

The second dataset (Fig. 5d–f) was generated using stdpopsim³⁶ and msprime simulation engine³⁷ using a reference panel and demographic model (four population out-of-Africa history) integrated into stdpopsim package⁵²⁻⁵⁶. This dataset contained simulated CEU population samples with a minimum MAF of 0.01 on chromosome 19. We selected 45,000 samples as the reference panel and 5000 samples with unique parental strands (haploids) for the test set and shortlisted the first 30,720 variants for the reference panel. Out of 30,720 variants present in the reference panel, 90% of them in each MAF bin (described at Data) were removed, leaving 3044 variants in the test samples. It is worth mentioning that there were shared parental strands (haploids) in the reference panel but the parental strands in the test set were all unique within the test set and among the test set and the reference panel.

The rat dataset^{38,39} (Fig. 5g–i) contains 5147 outbred samples from more than 10 projects on a selected region at rat chromosome 20 with 61,440 variants as the reference panel and 1000 samples with 6140 variants scattered thorough the reference panel variants. To preserve parental strand uniqueness, we used variants with a minimum MAF of 0.01.

The Sasso chicken dataset⁴⁰ (Fig. 5j–l) was already pre-processed by the curators, and non-biallelic SNVs and SNVs with MAF lower than 0.02 were removed. This dataset constitutes of 2258 pre-processed samples and 55,255 variants on chicken chromosome 20 and 100 test samples with 5488 variants selected among the reference variants. The test variants were obtained by randomly removing 90% of the reference panel variants in each MAF bin described earlier.

We used python scripts, PLINK2⁵⁰ and bcftools⁵¹ to pre-process the data and we used SHAPEIT5¹⁷ to impute the sporadic missing data (pre-phasing) for the rat and chicken datasets.

In previous studies, the training data is masked using different rates to match the test set, e.g., refs. 3,14. In our experiments, we observed improved performance of the model with 50% dynamic and random masking of the variants in the training data. So we trained the DL model once and reused it multiple times for sporadic missingness (MissR < 0.5). For higher MissRs, we can train multiple models for different MissR bins (e.g., 0.8 < MissR < 0.9 and 0.9 < MissR < 0.99) and use the proper saved model based on the missing rate in the target data. Notably, this masking is similar to the masking performed in modern large language models. The benefit of such a technique in genomic data imputation is the notable reduction in the inference (imputation) times when compared to the fastest traditional methods. Consequently, a DL model trained in this manner becomes particularly advantageous for deployment on imputation servers, where retraining needs to be avoided for quick and efficient processing.

Another improvement we achieved was by representing phased diploids into haploids, followed by one-hot encoding. That is, instead of feeding (one-hot encoded) phased diploids to the models, we fed them haploids. This idea is proposed in ref. 27, but there is no discussion about the merits of this procedure. We surmised that predicting haploids would be easier because mutations in paternal and maternal haploids are independent of each other. In the output,



Fig. 6 | **The architecture of STICI. a** Overall pipeline of the proposed framework: the data is separated into paternal and maternal haplotypes in the case of diplotypes, and it remains the same for haplotypes. While the figure shows phased genotypes, STICI can handle unphased data as well (though the performance degrades). Next, the data is one-hot encoded and fed into our Cat-Embedding layer, followed by splitting the data vertically into k chunks. The chunks overlap in order to capture information for the SNVs residing around the chunks' edges. Each branch passes through a unique set of attention, convolution, and fully connected layers. In the self-attention block, the flanking variants that come from the neighboring chunks are removed after applying multi-head attention. Finally, the results of all branches are assembled to generate the final sequence. **b** The workflow of proposed Categorical Embedding: we consider a unique vector space for each

diploid genotypes were reconstructed by combining corresponding haploids together.

STICI architecture

Split-Transformer Impute is an extended transformer model²⁹ specifically tailored for genotype imputation. STICI models do not require any additional information provided by a reference panel, except for the genotypes and their relative positions. This makes STICI adaptable to any genotype data and allows it to be applied to a wider range of datasets with less effort and fewer preparations. Moreover, although here we focus on sporadic missingness, once STICI is trained on a dataset, it can predict both sporadic missingness and systematic missingness in genotype data as long as the target variants are a subset of the training variants. An overview of STICI is presented in Fig. 6. We implemented STICI and the rest of the DL models using the Tensorflow framework⁵⁷ in Python. In order to train the models, we used tensor processing units (TPU) provided by the Google Colaboratory platform and GPU resources in Temple University's HPC servers. A learning rate scheduler and early stopping are employed in order to reduce the loss and training duration.

Cat-Embedding. One important part of STICI is categorical embedding (Fig. 6b), termed Cat-Embedding, which enables it to learn embedding representation per allele in each position. For the imputation task, we consider missing values as another allele that is equivalent to special tokens in natural language processing. The

unique categorical value in each SNV/feature. To save computational resources, instead of pre-allocating these vectors, we use the addition of positional embedding and categorical value embeddings in order to generate unique embedding vectors for each categorical value in each SNV/feature. We consider a missing (or masked) value as another categorical value (allele) in our model. Here, 2 (high-lighted in red) represents the missing value. **c** Convolution blocks: two parallel convolutional branches with varying kernel sizes are used in our convolution blocks. These multi-scale convolutional blocks allow STICI to capture information at multiple spatial scales in the input data, similar to the pattern-matching idea used in classical computer vision methods using convolution. Given the variable sizes of LD blocks, multi-scale convolutions.

corresponding vector for each allele is added to the respective positional variant embedding vector to generate the final embedding. The idea is similar to a natural language processing embedding layer that accepts word indices, except that Cat-Embedding accepts one-hot encoded data.

Splitting. While the multi-headed attention in a transformer offers significant advantages, a major drawback is quadratic memory cost for computations, which becomes important in genomic analysis since the number of variants in a sample is normally in the thousands. In genotypes, the majority of interactions are local⁵⁸. Therefore, it is of great importance to limit the scope of attention to save computational resources. To do so, we split the variants into chunks (vertical partitioning). The chunk size and overlap size are employed in a comparable manner in Minimac4.1.4 and analogous software applications. In order to prevent loss of imputation accuracy at chunk borders, we include flanking variants from neighboring chunks and discard them after applying self-attention to get the original variants in the chunk. Though the average LD block size in the dataset can be used to decide the size of overlap, we do not use LD blocks directly to decide the chunk size in the current version.

Each chunk passes through a dedicated branch inside the model, leading to increased imputation quality. Ideally, having a vast number of samples allows training a single model with attention across the whole genome. However, when the number of samples is not enough, the model is left with untrained parameters, resulting in poor performance. Hence, chunking regulates the number of parameters. In a vanilla transformer, the cost of computing global attention is quadratic with respect to the number of SNVs (m^2); however, the amount is lowered to (m/w) × (w + o)² = mw in STICI, considering that the overlaps of chunks are negligible compared to the chunk size. For instance, for $m = 10^4$ and a chunk size of 10^3 , STICI uses 10 times less memory for attention computations compared to a vanilla transformer.

Attention. The attention blocks are implemented similarly to those of other transformers, such as self-attention blocks in Vision Transformer (ViT)⁵⁹. There is a difference between the first and second attention blocks in the branches. The first block is a self-attention block, meaning that the query, key, and value of the attention layer are the same. The output of multi-head attention in Tensorflow has the same dimensions as the query. By excluding the neighboring variants of a chunk from the query and only including them in the key and value, we involve them in the attention mechanism and, at the same time, shrink the output of a chunk to the target size (chunk size without counting flanking/overlap variants) after applying multi-headed attention. In the second block, the query is the output of the previous layer, while the key and value are the outputs of the first self-attention block. This skip connection considerably affects the overall performance of the model.

Convolutional block. Convolutional blocks, as illustrated in Fig. 6c, are also crucial components of STICI. Through empirical studies, we found that using two parallel convolutional branches with varying kernel sizes, similar to the Inception module⁶⁰, is the best trade-off between accuracy gain and increase in a number of model parameters, compared to using a single branch or more than two branches. Furthermore, a Depth-wise convolutional layer at the end of the block helps STICI extract local information without mixing channel information and substantially improves imputation accuracy.

Output formation. Finally, the outputs of all branches are concatenated to form the output, that is, either maternal or paternal haplotype in the case of 1000 Genomes Project datasets or the genotypes in the case of yeast. For the former, by assembling maternal and paternal haplotypes, we obtain imputed genotypes, and the latter needs no further post-processing. Since genetic variations in parents are independent, directly encoding and imputing the genotypes in diploid life forms results in lower imputation accuracy compared to imputing their haplotypes. Hence we undergo extra steps to separate diplotypes into haplotypes in pre-processing and combining respective predicted haplotypes into diplotypes in post-processing for the human, chicken, and rat datasets.

Loss function. For the loss function, we used a combination of Kullback-Leibler divergence (D_{KL}) and categorical cross-entropy (CCE), similar to the loss function of variational autoencoder⁶¹, as follows:

$$\operatorname{Loss}(\mathbf{y}, \hat{\mathbf{y}}) = (\theta) \operatorname{CCE}(\mathbf{y}, \hat{\mathbf{y}}) + (1 - \theta) \operatorname{D}(\mathbf{y} \| \hat{\mathbf{y}})$$
(2)

where θ is the weight parameter. The first term, representing categorical cross-entropy, and the second term, representing Kullback-Leibler divergence loss, are calculated as follows:

$$CCE(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{C} \mathbf{y}_{ij} \log\left(\mathbf{p}(\mathbf{y}_{ij})\right)$$
(3)

$$\mathbf{D}_{\mathrm{KL}}(\mathbf{y}\|\hat{\mathbf{y}}) = \sum_{i=1}^{N} \mathbf{p}(\mathbf{y}_i) \frac{\mathbf{p}(\mathbf{y}_i)}{\mathbf{p}(\hat{\mathbf{y}}_i)} \tag{4}$$

We set θ to 0.5, meaning that STICI minimizes Equations (3), (4) equally. CCE captures reconstruction error between the input and the

output, while D_{KL} measures asymmetric distance, with **y** as the base, between their probability distributions. In our experiments, omitting any of these losses resulted in reduced model performance. Theoretically, KL-divergence and cross-entropy are related and using both might not seem to contribute to the performance of the model. However, adding KL-divergence to the loss term helps the model to retain the probability/dosage distribution of alleles per variant. In other words, while cross-entropy focuses on predicting the correct genotype, KL-divergence acts as a regularization factor and penalizes the model whenever the shape of the predicted probability distribution (allele probability/dosage) shows divergence from the ground truth. Moreover, the mathematical relation of D_{KL} and CCE can be summarized as follows:

$$D_{KL}(\mathbf{y}\|\hat{\mathbf{y}}) = CCE(\mathbf{y}, \hat{\mathbf{y}}) - CCE(\mathbf{y})$$
(5)

where CCE(**y**) is the entropy of the ground truth. According to Equation (5) minimizing $D_{KL}(\mathbf{y} \parallel \hat{\mathbf{y}})$ is equivalent to minimizing CCE($\mathbf{y}, \hat{\mathbf{y}}$) under the condition that the entropy of the ground truth remains constant. However, in deep learning models, data is typically processed in mini-batches. This means that the entropy of each mini-batch may not accurately represent the entropy of the entire ground truth. As a result, Equation (5) does not hold for the DL models in general.

We also used MaCH-Rsq⁶² as an additional loss term for STICI-Rsq. MaCH-Rsq metric is calculated for each variant/site in each sample as follows:

$$MaCH - Rsq \ loss_i = \frac{\frac{1}{n} \sum_{j=2}^{n_a} (\mathbf{D}_i^j - \hat{p})^2}{\hat{p}(1 - \hat{p})}$$
(6)

where \hat{p} is the alternate allele frequency for the ground truth at variant i in the current batch, N_a is the number of alleles, \mathbf{D}_i^j is the imputed allele probability at the i^{th} haplotype, and $\sum_{j=2}^{n_a} \mathbf{D}_i^j$ is the sum of predicted probabilities of alternative alleles for the i^{th} site. In the original MaCH-Rsq formula \hat{p} is calculated using the imputed genotypes. Notably, adding this loss causes the model to put more weight on the loss generated by rare alleles. This loss is sensitive to batch size and with an increase in batch size, this loss term chips away at other loss terms. We found out that using a batch size of 4 for training STICI-Rsq presents us with the best trade-off among accuracy and other metrics.

Baseline models

In order to benchmark our model, we compare STICI to state-of-the-art imputation models capable of imputing sporadic missingness: SCDA¹⁴, AE³, ImputeS¹⁸, SHAPEITS¹⁷, Eagle2³⁴ Beagle5.4¹⁹, and Minimac4.1.4¹⁵. In ref. 14, experimental results indicate that SCDA outperforms shallow ML models for genotype imputation. Hence, we do not include shallow ML models in our benchmarking analyses. In addition, in order to assess the contribution of Cat-Embedding, we replaced it with a convolution layer in STICI, named the resulting model STICI-NE, fine-tuned it, and applied it to the benchmark datasets. Lastly, we trained SCDA, in addition to STICI, using our proposed pre-processing and training procedure, and compared it to AE. Since AE and original SCDA are the same and only differ in the pre-processing step (which results in AE outperforming SCDA), we believe that this comparison can demonstrate the effective-ness of our proposed pre-processing and training procedure.

For SCDA and AE, hyper-parameter tuning information on the yeast dataset is present in the original papers. For SCDA, AE, and STICI, we conducted a grid search for optimal hyper-parameters on the HLA dataset using validation sets in a 3-fold cross-validation. We assessed the impact of these hyper-parameters on the performance of the models within the HLA dataset and applied these findings to select suitable hyper-parameters for the yeast dataset in the case of STICI and for the SV dataset across all four mentioned methods. The upper limit for the hyper-parameters was the resource limit of Google

Colaboratory using Nvidia Titan IV GPU with 16 GB of RAM size for AE, and roughly the same limitation for TPU RAM size. Classical imputation tools, such as Minimac, do not require fine-tuning for the experiments we run.

Experimental settings

The input to all DL models is one-hot encoded. While STICI can handle diploids, we found that the best performance was achieved when the inputs of the DL models were haplotypes, an analysis inspired by ref. 27. Therefore, for the HLA dataset and chromosome 22 datasets, we separated each diplotype into maternal and paternal haplotypes, fed them into the model, and reconstituted the resulting predictions for SCDA¹⁴ and STICI. We continue using diplotypes as inputs for AE³ since it is an improved version of SCDA in which the training process was modified, and we wanted to keep it intact. By doing so, we also compare the improvement in AE to our implementation of SCDA, called SCDA +, in which we use proposed pre-processing in conjunction with the changes to the training process as a contribution. The yeast dataset contains haplotypes, so there is no need for the aforementioned extra steps.

In this study, to evaluate the imputation power of the models, multiple evaluation metrics are used including imputation accuracy, imputation quality score (IQS)63, weighted F1-score, Pearson correlation coefficient between imputed and real genotypes in terms of R²⁶⁴, INFO score⁴¹, and MaCH-Rsq¹⁵. Accuracy and weighted *F*1-score are calculated only for positions with missing genotypes and for these metrics, heterozygous genotypes are encoded differently; i.e., 0|1 and 1|0 are encoded to two different categorical values. IQS adjusts the chance of concordance between predicted and ground truth genotypes and is defined for bi-allelic events only (thus not applicable to multi-allelic SVs on chromosome 22). R² is the squared Pearson correlation coefficient between the imputed genotypes and the true genotypes at a specific locus. INFO score is primarily used for quality control and indicates the quality of imputation. The definition of these metrics is provided in the Evaluation Metrics section of the Supplement. Lastly, MaCH-Rsg evaluates the guality of alternative allele imputation. For our experiments, we used a Python implementation of the INFO score provided in the GitHub repository of ref. 24.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

All data used in this study are publicly available. The yeast dataset can be found as the Supplementary Data 5 at https://www.nature.com/ articles/ncomms9712, the rest of datasets for sporadic missingness imputation are extracted from the 1000 Genomes Project phase 3 dataset available at http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/ 20130502/. Instructions on how to prepare the data for human chromosome 22 systematic missingness can be found in https://github. com/kanamekojima/rnnimp. The Rat dataset can be accessed at https://library.ucsd.edu/dc/using accession code *bb15123938*, and the chicken genotype data can be found in https://datashare.ed.ac.uk/ handle/10283/8761. The scripts used to simulate the human chromosome 19 dataset is stored in our GitHub repository. Source data are provided with this paper in Zenodo⁶⁵.

Code availability

The source code of STICI is publicly available on GitHub (https://github.com/shilab/STICI) and Zenodo⁶⁶.

References

 Lewis, C. M. & Vassos, E. Polygenic risk scores: from research tools to clinical instruments. *Genome Med.* 12, 1–11 (2020).

- 2. Torkamaneh, D., Belzile, F. Accurate imputation of untyped variants from deep sequencing data. *Methods Mol. Biol.* 271–281 https://doi. org/10.1007/978-1-0716-1103-6_13 (2021).
- Song, M. et al. An autoencoder-based deep learning method for genotype imputation. *Front. Artif. Intell.* 5, https://doi.org/10.3389/ frai.2022.1028978 (2022).
- 4. Das, S., Abecasis, G. R. & Browning, B. L. Genotype imputation from large reference panels. *Annu. Rev. Genomics Hum. Genet.* **19**, 73–96 (2018).
- Graffelman, J., Nelson, S., Gogarten, S. & Weir, B. Exact inference for hardy-weinberg proportions with missing genotypes: Single and multiple imputation. G3 5, 2365–2373 (2015).
- Wigginton, J. E., Cutler, D. J. & Abecasis, G. R. A note on exact tests of hardy-weinberg equilibrium. *Am. J. Human Genet.* 76, 887–893 (2005).
- Pei, Y.-F., Li, J., Zhang, L., Papasian, C. J. & Deng, H.-W. Analyses and comparison of accuracy of different genotype imputation methods. *PloS ONE* 3, 3551 (2008).
- Auer, P. L., Wang, G., Project, N. E. S. & Leal, S. M. Testing for rare variant associations in the presence of missing data. *Genet. Epidemiol.* 37, 529–538 (2013).
- 9. Chiang, C. et al. The impact of structural variation on human gene expression. *Nat. Genet.* **49**, 692–699 (2017).
- 10. Consortium, G. P. et al. A global reference for human genetic variation. *Nature* **526**, 68 (2015).
- 11. Liu, Z. et al. Towards accurate and reliable resolution of structural variants for clinical diagnosis. *Genome Biol.* **23**, 68 (2022).
- Bartlett, J. W., Seaman, S. R., White, I. R., Carpenter, J. R. & Initiative*, A.D.N. Multiple imputation of covariates by fully conditional specification: accommodating the substantive model. *Stat. Methods Med. Res.* 24, 462–487 (2015).
- 13. Song, M. et al. A review of integrative imputation for multi-omics datasets. *Front. Genet.* **11**, 570255 (2020).
- Chen, J. & Shi, X. Sparse convolutional denoising autoencoders for genotype imputation. *Genes* 10, 652 (2019).
- Das, S. et al. Next-generation genotype imputation service and methods. Nat. Genet. 48, 1284–1287 (2016).
- Li, Y., Willer, C. J., Ding, J., Scheet, P. & Abecasis, G. R. Mach: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet. Epidemiol.* 34, 816–834 (2010).
- Hofmeister, R. J., Ribeiro, D. M., Rubinacci, S. & Delaneau, O. Accurate rare variant phasing of whole-genome and whole-exome sequencing data in the uk biobank. *Nat. Genet.* 55, 1243–1249 (2023).
- Rubinacci, S., Delaneau, O. & Marchini, J. Genotype imputation using the positional burrows wheeler transform. *PLoS Genet.* 16, 1009049 (2020).
- Browning, B. L., Zhou, Y. & Browning, S. R. A one-penny imputed genome from next-generation reference panels. *Am. J. Hum. Genet.* 103, 338–348 (2018).
- Scheet, P. & Stephens, M. A fast and flexible statistical model for large-scale population genotype data: applications to inferring missing genotypes and haplotypic phase. *Am. J. Hum. Genet.* 78, 629–644 (2006).
- Rubinacci, S., Hofmeister, R. J., Mota, B. & Delaneau, O. Imputation of low-coverage sequencing data from 150,119 uk biobank genomes. *Nat. Genet.* 55, 1088–1090 (2023).
- Bloom, J. S. et al. Genetic interactions contribute less than additive effects to quantitative trait variation in yeast. *Nat. Commun.* 6, 1–6 (2015).
- 23. Dias, R. et al. Rapid, reference-free human genotype imputation with denoising autoencoders. *Elife* **11**, 75600 (2022).
- 24. Kojima, K. et al. A genotype imputation method for de-identified haplotype reference information by using recurrent neural network. *PLoS Comput. Biol.* **16**, 1008207 (2020).

- Kojima, K., Tadaka, S., Okamura, Y., Kinoshita, K. Two-stage strategy using denoising autoencoders for robust reference-free genotype imputation with missing input genotypes. J. Hum. Genet. 69, 511–518 (2024)
- 26. Chi Duong, V. et al. A rapid and reference-free imputation method for low-cost genotyping platforms. *Sci.c Rep.* **13**, 23083 (2023).
- Naito, T. et al. A deep learning method for hla imputation and transethnic mhc fine-mapping of type 1 diabetes. *Nat. Commun.* 12, 1–14 (2021).
- Tanaka, K., Kato, K., Nonaka, N. & Seita, J. Efficient HLA imputation from sequential SNPs data by transformer. J. Hum. Genet. 69, 533–540 (2024).
- Vaswani, A. et al. Attention is all you need. In Advances in Neural Information Processing Systems 30, (2017).
- Desimone, R. & Duncan, J. Neural mechanisms of selective visual attention. Annu. Rev. Neurosci. 18, 193–222 (1995).
- Cho, K., Courville, A. & Bengio, Y. Describing multimedia content using attention-based encoder-decoder networks. *IEEE Trans. Multimedia* 17, 1875–1886 (2015).
- 32. Jumper, J. et al. Highly accurate protein structure prediction with alphafold. *Nature* **596**, 583–589 (2021).
- Lin, Z. et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. Science 379, 1123–1130 (2023).
- 34. Loh, P.-R. et al. Reference-based phasing using the haplotype reference consortium panel. *Nat. Genet.* **48**, 1443–1448 (2016).
- 35. Sudmant, P. H. et al. An integrated map of structural variation in 2,504 human genomes. *Nature* **526**, 75–81 (2015).
- 36. Adrion, J. R. et al. A community-maintained standard library of population genetic models. *Elife* **9**, 54967 (2020).
- Kelleher, J., Etheridge, A. M. & McVean, G. Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS Comput. Biol.* 12, 1004842 (2016).
- Gunturkun, M. H. et al. Genome-wide association study on three behaviors tested in an open field in heterogeneous stock rats identifies multiple loci implicated in psychiatric disorders. *Front. Psychiatry* **13**, 790566 (2022).
- Gileta, A. F. et al. Adapting genotyping-by-sequencing and variant calling for heterogeneous stock rats. G3 10, 2195–2205 (2020).
- Morris, K.M. et al. Genotype data of Sasso chicken. University of Edinburgh. Centre For Tropical Livestock Genetics and Health https://doi.org/10.7488/ds/7718 (2024).
- Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation method for the next generation of genomewide association studies. *PLoS Genet.* 5, 1000529 (2009).
- Bolormaa, S., Chamberlain, A., Werf, J., Daetwyler, H. & MacLeod, I. Evaluating the accuracy of imputed whole genome sequence in sheep. In Proceedings of the World Congress on Genetics Applied to Livestock Production (2018).
- Pistis, G. et al. Rare variant genotype imputation with thousands of study-specific whole-genome sequences: implications for costeffective study designs. *Eur. J. Hum. Genet.* 23, 975–983 (2015).
- 44. Sudlow, C. et al. Uk biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* **12**, 1001779 (2015).
- 45. Taliun, D. et al. Sequencing of 53,831 diverse genomes from the nhlbi topmed program. *Nature* **590**, 290–299 (2021).
- 46. Us Research Program Investigators, A. The "all of us" research program. *N. Engl. J. Med.* **381**, 668–676 (2019).
- Miles, A. et al. Cggh/scikit-allel: V1.3.3. https://doi.org/10.5281/ zenodo.4759368 (2021).
- Hillert, J. Human leukocyte antigen studies in multiple sclerosis. Ann. Neurol. 36, 15–17 (1994).
- Terasaki, P. I. & Cai, J. Human leukocyte antigen antibodies and chronic rejection: from association to causation. *Transplantation* 86, 377-383 (2008).

- 50. Chang, C. C. et al. Second-generation plink: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 13742–015 (2015).
- 51. Li, H. A statistical framework for snp calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
- Takahata, N. Allelic genealogy and human evolution. *Mol. Biol. Evol.* 10, 2–22 (1993).
- Tremblay, M. & Vézina, H. New estimates of intergenerational time intervals for the calculation of age and origins of mutations. *Am. J. Hum. Genet.* 66, 651–658 (2000).
- 54. Spence, J. P. & Song, Y. S. Inference and analysis of populationspecific fine-scale recombination maps across 26 diverse human populations. *Sci. Adv.* **5**, 9206 (2019).
- 55. Consortium, I. H. et al. A second generation human haplotype map of over 3.1 million snps. *Nature* **449**, 851 (2007).
- Jouganous, J., Long, W., Ragsdale, A. P. & Gravel, S. Inferring the joint demographic history of multiple populations: beyond the diffusion approximation. *Genetics* **206**, 1549–1567 (2017).
- Abadi, M. et al. TensorFlow: a system for large-scale machine learning. In Proceedings of the 12th USENIX conference on Operating Systems Design and Implementation (OSDI'16). 265–283 (USE-NIX Association, USA, 2016).
- Weir, B. Linkage disequilibrium and association mapping. Annu. Rev. Genom. Hum. Genet. 9, 129–142 (2008).
- 59. Dosovitskiy, A. et al. An image is worth 16 × 16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*. https://openreview.net/forum?id=YicbFdNTTy (2021).
- Szegedy, C. et al. Going deeper with convolutions. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 1–9 (2015).
- 61. Kingma, D. P. & Welling, M. Auto-encoding variational bayes. Preprint at https://doi.org/10.48550/arXiv.1312.6114 (2013).
- 62. Das, S. Minimac3 Info File. https://genome.sph.umich.edu/wiki/ Minimac3_Info_File. (2024).
- 63. Lin, P. et al. A new statistic to evaluate imputation reliability. *PloS ONE* **5**, 9697 (2010).
- 64. Deng, T. et al. Comparison of genotype imputation for snp array and low-coverage whole-genome sequencing data. *Front. Genet.* **12**, 704118 (2022).
- 65. Mowlaei, M. E. et al. Split-transformer with integrated convolutions for genotype imputation (STICI) SOURCE DATA FILES. https://zenodo.org/records/14649993 (2024).
- 66. Mowlaei, M. E. et al. shilab/STICI: STICI-V1.1.2. https://doi.org/10. 5281/zenodo.14451462 (2024).

Acknowledgements

This work is partially supported by the US National Science Foundation (DBI 1750632) and the National Institutes of Health (R35-GM139540-04, U24HG007497, R01GM093290). This research includes calculations carried out on HPC resources supported in part by the National Science Foundation through major research instrumentation grant number 1625061 and by the US Army Research Laboratory under contract number W911NF-16-2-0189. We appreciate the suggestions provided by Dr. Francisco McGee, John Allard, Rohan Alibutud, and Vahid Mahzoon that helped us improve the model performance and design the experiments. In addition, we would like to thank Dr. Kaname Kojima for helping us obtain the data for the Missing variant experiment and supplementary code for calculating some of the evaluation metrics and Emily Thyrum for proofreading the manuscript.

Author contributions

M.E.M. developed the method with the help of J.C., B.J., V.C., X.S., and M.E.M. implemented the code. C.L. and J.C. prepared the datasets. M.E.M. and R.D. performed the experiments. M.E.M., S.K., C.L., O.J., and

X.S. designed and/or conducted data analysis. M.E.M., C.L., S.K., T.R.R., V.C., and X.S. wrote the manuscript. All the authors read and approved the submitted manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at https://doi.org/10.1038/s41467-025-56273-3.

Correspondence and requests for materials should be addressed to Xinghua Shi.

Peer review information *Nature Communications* thanks Masaru Koido and the other anonymous reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

Reprints and permissions information is available at http://www.nature.com/reprints

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http:// creativecommons.org/licenses/by-nc-nd/4.0/.

© The Author(s) 2025