

# Inferring Missing Genotypes from Partially Observed SNP Data Using Deep Learning-Based Architectures

W. Zawadzka<sup>1\*</sup>, J. Szyda<sup>1,2</sup>, and M. Frąszczak<sup>1</sup>

<sup>1</sup> *Biostatistics Group, Department of Genetics, Wrocław University of Environmental and Life Sciences, 51-631 Wrocław, Poland;* <sup>2</sup> *Wrocław University of Science and Technology, Department of Biomedical Engineering, Wrocław, Poland;* \*Presenting and corresponding author: [130732@student.upwr.edu.pl](mailto:130732@student.upwr.edu.pl)

## Introduction

The goal of genomic imputation is to predict missing genotypes from partially observed single nucleotide polymorphism (SNP) data. Missing genotypes commonly arise when individuals genotyped using oligonucleotide arrays are imputed to whole genome sequence level. Accurate imputation is crucial for downstream analyses, particularly for the detection and interpretation of rare variants. Traditional imputation tools, such as Beagle, infer missing genotypes using Hidden Markov Models that exploit local haplotype structure (Browning and Browning, 2007; Browning and Browning, 2016). These methods are widely used due to their efficiency and computational speed. Our study explores the use of deep learning (DL) for SNP imputation, investigating how different model architectures affect imputation accuracy taking the unbalanced genotypes into account. The motivation for using DL-based architectures is their flexibility in modeling complex patterns, which may allow improved reconstruction of rare genotypes by more effectively exploring the parameter space. We evaluate four AutoEncoder based approaches: a naïve AutoEncoder using dense neural networks, a dense AutoEncoder incorporating input embeddings, a convolutional neural network based AutoEncoder (CAE), and a transformer based AutoEncoder utilising the BERT model. Imputation accuracy obtained using Beagle serves as a baseline for comparison.

## Materials & Methods

**Materials.** The material consisted of a subset of bulls from the 1000 Bull Genomes Project (Hayes et al., 2019). A fragment of chromosome 28 containing the first 500 SNPs was used. In total, 928 bulls were included, with 742 bulls assigned to the training set and 186 bulls to the testing set. The training set was further randomly split into training and validation subsets using an 80:20 ratio. Minor allele frequencies (MAFs) for the test data were calculated using PLINK2 (Chang et al., 2026), and SNPs were considered rare if their MAF was below 0.05.

To mimic different real life complexity scenarios, multiple test datasets with missing genotypes were generated from the testing set. For each scenario, a random subset of SNP columns was masked for all individuals, with missingness ranging from 10% to 80%. DL models were trained using a missingness level of 60%, which was selected as a compromise between task difficulty and model stability. One advantage of such models is that they can be trained directly on the actual missingness pattern; fixed missingness levels here are used only for comparison purposes. Under the 60% missingness scenario applied to the test data, the corresponding training set shows a strongly imbalanced genotype distribution: genotype 0 accounts for 75.57% of observations, genotype 1 for 19.54%, and genotype 2 for 4.89%.

**Autoencoder models.** Autoencoder based approaches have been chosen for this comparison due to their ability to learn latent representations that capture dependencies between SNPs and to use this information to accurately reconstruct missing genotypes. An autoencoder consists of an encoder, which maps the input genotypes (0/1/2) into a lower dimensional

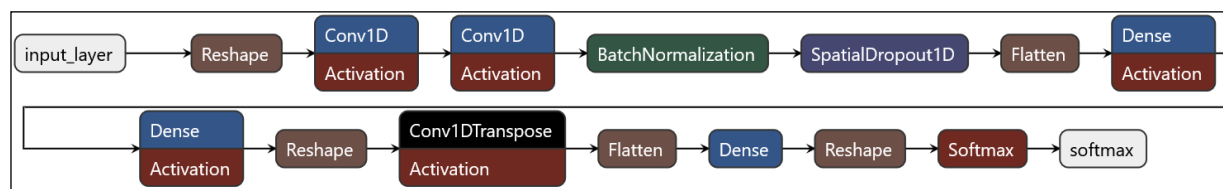
latent space, and a decoder, which reconstructs the original genotype matrix from this representation. ReLU activations enable learning of complex patterns, and training minimizes weighted (using inverse class frequency) categorical cross-entropy at masked positions to focus on imputing missing genotypes.

The dense feedforward autoencoder (Figure 1) encodes 500 SNP genotypes through two fully connected layers down to a 64-unit bottleneck, then decodes back to 500 SNPs, outputting three probabilities per SNP via SoftMax. In the DenseEmb variant, genotypes are first mapped to learnable 3D embeddings, processed with a 1D convolution and global average pooling, then passed through an encoder expanding to 1024 units and compressing to a 64-unit bottleneck.



**Figure 1.** Architecture of deep autoencoder models.

The convolutional autoencoder (Figure 2) encodes 500 SNPs via two 1D convolutional layers with ReLU, batch normalization, and spatial dropout, then flattens to a 16-unit dense bottleneck that captures global dependencies. The decoder mirrors this with dense and transposed convolution layers to reconstruct the SNPs. Convolutional filters use a window size of 21 and stride 3 to capture local SNP dependencies. All AE models were implemented in TensorFlow 2.20.0.



**Figure 2.** Architecture of convolutional autoencoder model.

**BERT model.** This transformer-based model was trained using the masked language modeling task (MLM), where the objective is to predict missing or masked tokens in a sequence based on their surrounding context. Genotypes were treated as words in the vocabulary, forming sequences of length 500. Implemented with BertForMaskedLM in PyTorch 2.10, training used focal loss on masked positions with inverse class-frequency weights. The architecture included 2 transformer layers, 1 attention head, a hidden size of 256, and a feedforward size of 1024.

**Mask overrepresentation.** Given the high global imbalance of missing genotypes, a controlled mask upsampling strategy was applied when creating masked inputs. For each training sample, SNPs were initially masked according to the same positions missing in the test set. To mitigate the extreme imbalance between genotype classes, masked SNPs with genotypes 1 and 2 were preferentially included. A multiplier controlled the relative frequency of each genotype class in the masked positions, and its optimal value 1-3 was determined via hyperparameter optimization. This approach does not constitute traditional oversampling, as no new or repeated examples are added to the mask; rather, it selectively increases the representation of less frequent genotypes in those positions.

**Optimization.** Model parameters were tuned with Optuna (Akiba et al., 2019) over 50 trials, maximizing macro F1 on the validation set. Macro F1 was chosen to account for the imbalanced genotype classes. Since the 0 genotype is overwhelmingly common, overall accuracy can be misleadingly high, whereas macro F1 ensures that performance across all three genotype classes is evaluated. Training was performed on AMD Ryzen 7 4800HS.

## Results & Discussion

Performance was evaluated in three scenarios: (i) 60% missingness, matching the training condition (in-distribution); (ii) the same evaluation restricted to rare genotypes; and (iii) application to other missingness levels to test generalization. However, in realistic settings, the missingness pattern used during training is typically reflective of the target data to be imputed. The additional missingness settings are included to evaluate robustness.

Overall, in the 60% missingness scenario, the transformer-based BERT achieved the highest F1 macro score of 0.868 (Table 1). Despite having the fewest parameters, it required the longest training time of approximately 3 hours. Faster alternatives, all training in under a minute on this dataset, included dense AEs, with the version using an embedding layer before the AE performing slightly better (F1 micro 0.853 vs 0.847). The convolutional autoencoder (trained for approx. 8 minutes) performed worse than the other three models. Its large number of parameters did not translate into better performance, making it less efficient and less accurate. Beagle requires no training, with near real-time inference.

**Table 1.** All missing positions imputation performance comparison of trained models and Beagle on the dataset with 60% missingness.

Model	Model performance on test set					Number of parameters
	F1 micro	F1 macro	F1 genotype 0	F1 genotype 1	F1 genotype 2	
<b>BERT</b>	<b>0.917</b>	<b>0.868</b>	<b>0.953</b>	<b>0.811</b>	<b>0.840</b>	<b>1779976</b>
DenseEmb	0.906	0.853	0.946	0.777	0.835	3166748
Dense	0.909	0.847	0.950	0.782	0.810	3133948
CAE	0.906	0.838	0.949	0.757	0.809	73035868
Beagle	0.897	0.807	0.948	0.676	0.797	-

While focusing specifically on the rare genotypes inside the missing positions, we can observe the change of the ranking (see Table 2). BERT has remained the highest scoring model (F1-macro=0.563) but the Beagle is right after it (F1-macro=0.556). The Dense autoencoder with embedding and the CAE failed to correctly identify any genotype 2 examples. This table also highlights the limitations of relying on overall accuracy (or F1-micro), which exceed 0.96 for all models and therefore obscure the true challenge of imputing rare genotypes.

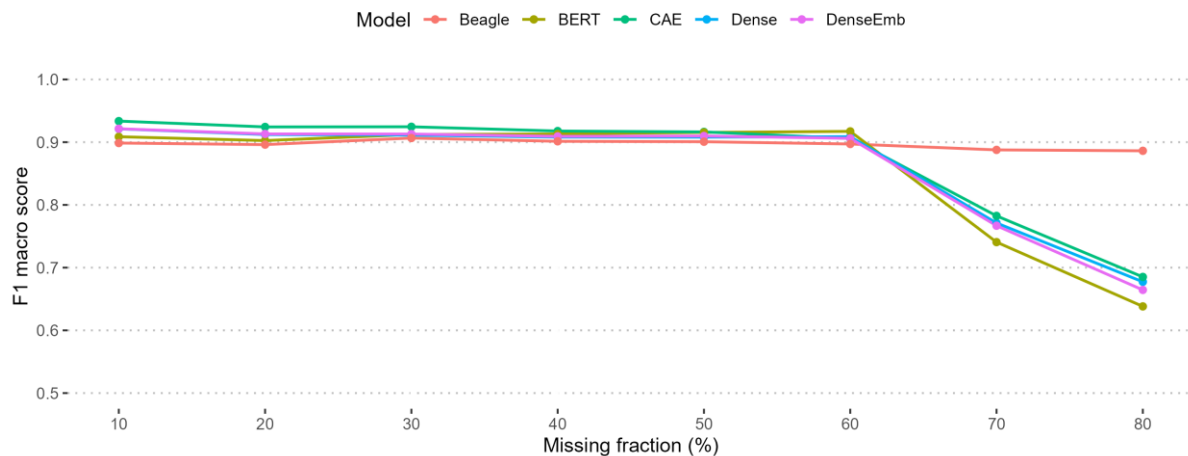
**Table 2.** Rare<sup>1</sup> genotype imputation performance comparison of trained models and Beagle on the dataset with 60% missingness.

Model	Model performance on test set for rare genotypes				
	F1 micro	F1 macro	F1 genotype 0	F1 genotype 1	F1 genotype 2
<b>BERT</b>	<b>0.954</b>	<b>0.563</b>	<b>0.976</b>	<b>0.541</b>	<b>0.173</b>
Beagle	0.971	0.556	0.986	0.533	0.148
Dense	0.966	0.490	0.983	0.421	0.066
DenseEmb	0.965	0.450	0.982	0.369	0
CAE	0.961	0.356	0.980	0.088	0

<sup>1</sup> Rare defined as MAF < 0.05.

**Comparison across missingness levels.** Models trained on 60% missingness performed similarly on lower missing fractions, but their performance dropped significantly on higher

fractions beyond the training range (Figure 3). Beagle remained consistent across the full range of missingness, showing the greatest adaptability.



**Figure 3.** F1-macro comparison across models and missing fractions.

## Conclusions

Deep learning models, especially transformer-based BERT, achieve high F1 scores overall and for rare variants. BERT balances high accuracy with a small parameter count but requires longer training. Dense autoencoders, with or without embeddings, are faster and competitive but struggle with rare (low MAF) genotypes. Beagle is fast and consistent across missingness levels, performing well on rare genotypes but with lower overall accuracy at 60% missingness. Accurate imputation of rare variants is critical for downstream analyses not only in humans for accurate prediction of Polygenic Risk Scores or treatment design, but also in livestock for mate selection or accurate estimation of breeding values. Genome-Wide Association Studies in all species also profit from accurate SNP genotype imputation that allows for processing very large phenotyped cohorts in which not all individuals are sequenced with the genome-wide resolution. However, high class imbalance between common and rare variant complicates imputation. Strategies include using imbalance-aware loss functions, controlled masking during training, and careful evaluation metrics, as standard metrics can be misleading. More experimentation is needed, particularly on larger datasets where longer linkage disequilibrium patterns could potentially improve rare genotype imputation. It would also be useful to train models on a range of missingness levels to assess true performance, since in real scenarios the missingness fraction of the test data is known.

## References

- Browning B.L., and Browning S.R. (2007) Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am J Hum Genet* 81:1084-1097. doi:10.1086/521987
- Browning B.L., and Browning S.R. (2016) Genotype imputation with millions of reference samples. *Am J Hum Genet* 98:116-126. doi:10.1016/j.ajhg.2015.11.020
- Chang C.C., Chow C.C., Tellier L.C.A.M., Vattikuti S., Purcell S.M., Lee J.J. (2015) Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience* 4:s13742-015-0047-8. doi:10.1186/s13742-015-0047-8
- Akiba T., Sano S., Yanase T., Ohta T., Koyama M. (2019) Optuna: a next-generation hyperparameter optimization framework. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2623–2631. doi:10.1145/3292500.3330701
- Hayes B.J., Daetwyler H.D. (2019) 1000 Bull Genomes Project to map simple and complex genetic traits in cattle: applications and outcomes. *Annu Rev Anim Biosci* 7:89-102. doi:10.1146/annurev-animal-020518-115024