# Masked Autoencoder Joint Learning for Robust Spitzoid Tumor Classification

Ilán Carretero [1], Roshni Mahtani [1], Silvia Perez-Deben [2], José Francisco González-Muñoz [2],
Carlos Monteagudo [2], Valery Naranjo [1,3], Rocío del Amor [1,3]

[1] HUMAN-tech, Universitat Politècnica de València (UPV), Valencia, Spain

[2] INCLIVA, Universitat de València (UV), Valencia, Spain

[3] Artikode Intelligence S.L., Valencia, Spain

✉{ilcarjuc, rmahvas, vnaranjo, madeam2}@upv.es

## Abstract

*Accurate diagnosis of spitzoid tumors (ST) is critical to ensure a favorable prognosis and to avoid both under- and over-treatment. Epigenetic data, particularly DNA methylation, provide a valuable source of information for this task. However, prior studies assume complete data, an unrealistic setting as methylation profiles frequently contain missing entries due to limited coverage and experimental artifacts. Our work challenges these favorable scenarios and introduces ReMAC, an extension of ReMasker designed to tackle classification tasks on high-dimensional data under complete and incomplete regimes. Evaluation on real clinical data demonstrates that ReMAC achieves strong and robust performance compared to competing classification methods in the stratification of ST. Code is available: https://github.com/roshni-mahtani/ReMAC.*

## 1. Introduction

Spitzoid tumors (ST) are melanocytic neoplasms defined by large spindle and epithelioid cell morphology, coupled with unpredictable clinical behavior [1]. ST are typically classified into three categories: the benign form, Spitz Nevus (SN); an intermediate entity with uncertain malignant potential, Atypical Spitz Tumor (AST) or Spitz melanocytoma; and the malignant form, Spitz Melanoma (SM) [1]. Accurate diagnosis is essential, as misclassification may result in severe clinical consequences and inappropriate treatment [2].

Epigenetic profiling, particularly through Reduced Representation Bisulfite Sequencing (RRBS), has become a widely adopted next-generation sequencing technique that enriches CpG-dense regions to generate genome-wide methylation maps at single-base resolution [3]. Building on this foundation, several studies have leveraged DNA methylation (DNAm) to advance in the stratification of spitzoid tumors [4, 5]. However, most approaches have generally overlooked the issue of missing values, a recurrent challenge in methylation data stemming from limited sequencing depth or technical variability [6]. This limitation highlights the need for computational frameworks that remain accurate and robust when operating on incomplete epigenetic data.

Several state-of-the-art imputers have been developed to address missing values, including discriminative approaches such as MICE [7] and MIRACLE [8], as well as generative models such as GAIN [9] and HI-VAE [10]. More recently, ReMasker [11] introduced a self-supervised masked autoencoding (MAE) framework for tabular data imputation, demonstrating strong performance and the ability to learn missingness-invariant representations. Nevertheless, this model has not been tested on high-dimensional data such as epigenetics, nor has it been extended to produce latent representations tailored for classification tasks.

In this work, we propose ReMAC, a framework that extends the **Re**Masker model by incorporating a **M**ean **A**ttribution **C**lassification (MAC) branch. As a result, ReMAC emerges as a state-of-the-art method capable of producing robust discriminative representations on DNA methylation data to classify spitzoid tumors.
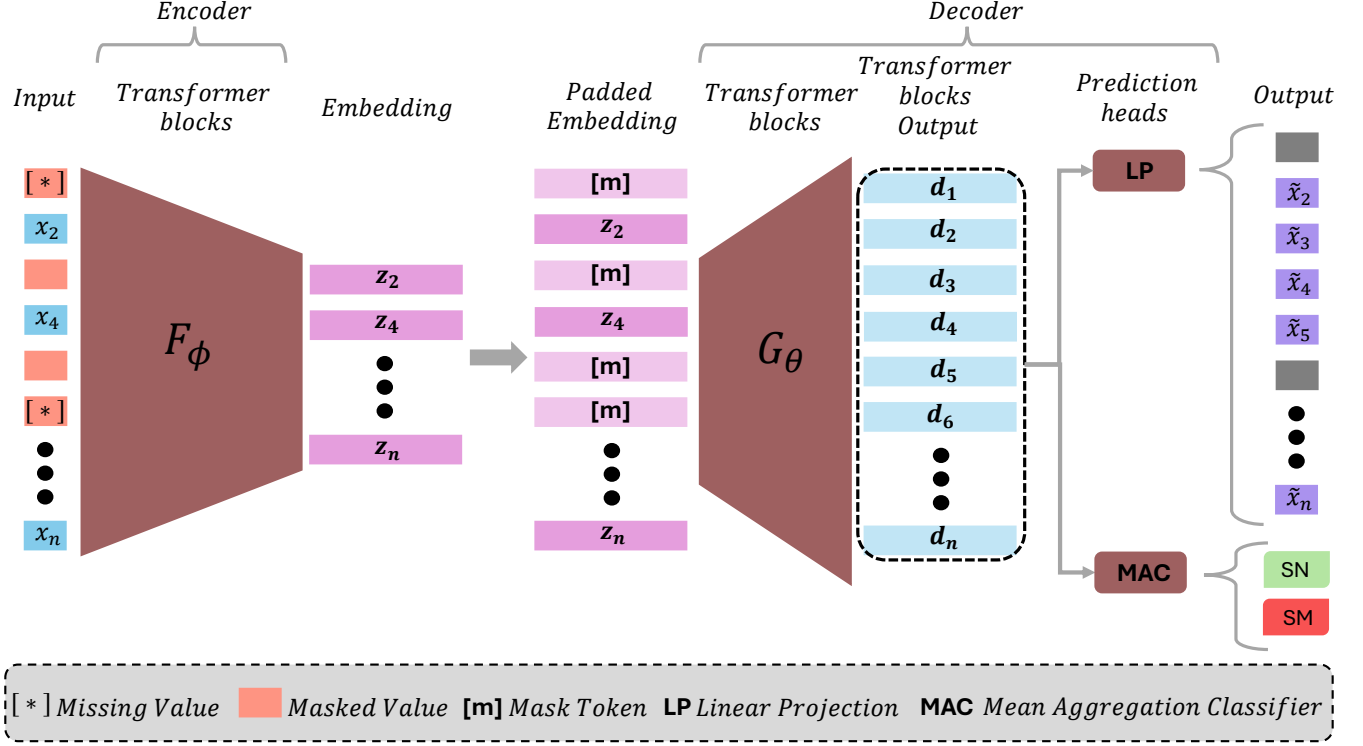
## 2. Methodology

An overview of the proposed method is illustrated in Fig. 1. The problem formulation and the different components implemented are described in the following subsections.

### 2.1. Problem formulation

Let $\mathcal{S} = \{(x^{(j)}, y^{(j)})\}_{j=1}^M$ denote the set of tumor samples $(M)$, where each sample $x^{(j)} = (x_1^{(j)}, x_2^{(j)}, \ldots, x_n^{(j)}) \in \mathbb{R}^n$ is the epigenetic profile represented by $n$ features, and $y^{(j)} \in \{0, 1\}$ is the corresponding label. The objective is to learn a mapping $h_\psi : \mathbb{R}^n \to \{0, 1\}$, where label 0 corresponds to SN and label 1 corresponds to SM.

### 2.2. Masked Autoencoding for tabular data

ReMasker [11] is a masked autoencoding framework for tabular data. An encoder $F_\phi$, implemented as a stack of Transformer blocks, maps the input into a set of latent embeddings $Z = \{z_i\}_{i=1}^n$. To model missing patterns, positions selected for masking are added in $Z$ by a learned mask token $\mathbf{m}$. This design enables the encoder to capture contextual feature dependencies through multi-head self-attention.

**Figure 1:** *Method Overview. In this article, we introduce ReMAC, a framework that extends the ReMasker approach by incorporating an aggregation and classification head to learn missingness-invariant discriminative representations for the stratification of spitzoid tumors into Spitz nevus and Spitz melanoma.*

A decoder $G_\theta$, also composed of Transformer blocks, processes the augmented set of embeddings $Z$, and a subsequent linear projection (LP) maps the decoder embeddings $d_i$ to the feature space, producing reconstructions $\tilde{x}_i = \mathrm{LP}(d_i)$. The training objective minimizes the mean squared error (MSE) restricted to the union of observed indices $\Omega_o$ and masked-but-originally-observed indices $\Omega_m$:

$$\mathcal{L}_{\mathrm{REC}} = \frac{1}{|\Omega_o \cup \Omega_m|} \sum_{i \in \Omega_o \cup \Omega_m} \left(x_i - \tilde{x}_i\right)^2.$$

This formulation forces the model to reconstruct true feature values from contextual information, yielding robust latent representations invariant to missingness patterns.

### 2.3. Extending ReMasker for classification

ReMasker learns latent representations invariant to missingness patterns; however, its original design is limited to the imputation task. To enable classification, we introduce a Mean Aggregation Classification (MAC) head. Let $D = \{d_i\}_{i=1}^n$ denote the sequence of decoder embeddings, i.e., the token representations produced by the decoder $G_\theta$. These embeddings are aggregated via mean pooling, $\bar{d} = \frac{1}{n} \sum_{i=1}^n d_i$, and a dense layer maps the global representation to class probabilities, $\tilde{y} = \mathrm{softmax}(W\bar{d} + b)$, with $y \in \{0, 1\}$. The classification loss is defined as the binary cross-entropy $\mathcal{L}_{\mathrm{CLF}} = \mathrm{BCE}(y, \tilde{y})$, and the overall training objective combines reconstruction and classification:

$$\mathcal{L}_{\mathrm{ReMAC}} = \mathcal{L}_{\mathrm{REC}} + \mathcal{L}_{\mathrm{CLF}}.$$

This joint objective ensures that the learned representations remain robust to missingness while being explicitly shaped for discriminative classification.

## 3. Experimental setting

### 3.1. Dataset

For this study, a total of 21 formalin–fixed paraffin–embedded (FFPE) tumor samples from patients with spitzoid lesions were analyzed. These included 12 cases classified as Spitz Nevus (SN) and 9 cases classified as Spitz Melanoma (SM). Tumor specimens were collected at the time of surgery and reported to the Department of Anatomic Pathology of the Hospital Clínico Universitario, Valencia (Spain) between 1990 and 2018. The protocol was approved by the Ethical and Scientific Committees of the Hospital Clínico Universitario, and written informed consent was obtained from all patients.

### 3.2. Bioinformatic preprocessing

Sequencing quality was assessed with *FastQC*, and adapters were removed using *Trim Galore!*. Reads were aligned to the human reference genome (hg19) with *Bismark*, followed by methylation calling. Additional quality control metrics were summarized with *MultiQC*, and coverage normalization was carried out using the *methylKit* R package. The resulting methylation values ranged from 0 to 1 with no missing values.

### 3.3. Implementation and validation protocol

All experiments were conducted under a 4-fold stratified cross-validation regime. The proposed method, as well as competing state-of-the-art approaches, were implemented in Python 3.10 using standard libraries for data processing and modeling (e.g., *NumPy* 1.22.2, *pandas* 1.5.3, *scikit-learn* 1.2.0). To ensure reproducibility of both code and results, pseudo-random seeds were fixed and experiments were executed within Docker, specifically using the *PyTorch* 23.10 container image.

## 4. Results

### 4.1. Results on the complete dataset

The results of different classification models on the complete dataset, together with ReMAC, are reported in Table 1. The proposed method is compared with conventional machine learning approaches, including K-Nearest Neighbors (KNN), Logistic Regression (LR), Support Vector Machine (SVM), Random Forest (RF), and eXtreme Gradient Boosting (XGB), as well as with deep learning models such as a Multilayer Perceptron (MLP), an Autoencoder with a classification head (AE+MLP), and TabNet [12]. ReMAC achieves strong performance in the classification of spitzoid tumors under the complete data regime, outperforming the competing models in two out of the four reported metrics.

| MODEL | ACC | SEN | SPE | AUC |
|---|---|---|---|---|
| KNN | 0.73 ± 0.15 | 0.63 ± 0.48 | 0.83 ± 0.19 | 0.77 ± 0.13 |
| LR | 0.73 ± 0.22 | 0.63 ± 0.48 | 0.83 ± 0.19 | 0.82 ± 0.14 |
| SVM | 0.68 ± 0.15 | 0.63 ± 0.48 | 0.75 ± 0.17 | **0.86 ± 0.10** |
| RF | 0.68 ± 0.15 | 0.63 ± 0.48 | 0.75 ± 0.17 | 0.85 ± 0.11 |
| XGB | 0.68 ± 0.15 | 0.50 ± 0.41 | 0.83 ± 0.19 | 0.58 ± 0.10 |
| MLP | 0.81 ± 0.16 | 0.79 ± 0.25 | 0.83 ± 0.19 | **0.86 ± 0.10** |
| AE+MLP | 0.81 ± 0.02 | **0.92 ± 0.17** | 0.75 ± 0.17 | **0.86 ± 0.10** |
| TabNet | 0.72 ± 0.22 | 0.71 ± 0.30 | 0.75 ± 0.43 | 0.73 ± 0.18 |
| ReMAC (*Ours*) | **0.86 ± 0.10** | 0.67 ± 0.24 | **1.00 ± 0.00** | **0.86 ± 0.10** |

**Table 1:** *Comparison of machine learning and deep learning models on the complete dataset (no missing values). Reported metrics include accuracy (ACC), sensitivity (SEN), specificity (SPE), and area under the curve (AUC). The best value for each metric is highlighted in bold, and the proposed method is shown with a gray background.*
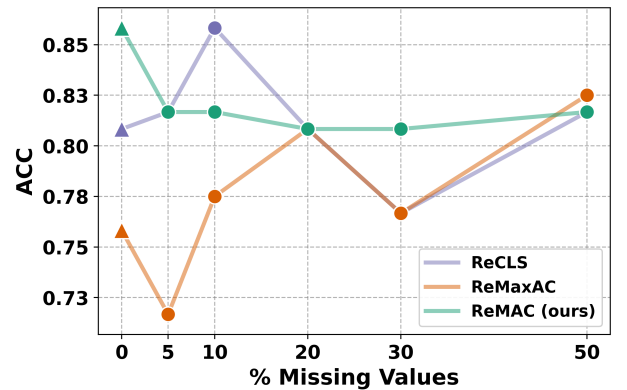
### 4.2. Results on incomplete datasets

Table 2 reports the performance of classification models robust to missing data, together with ReMAC, under pseudo-randomly generated missing-value (%MV) regimes. We considered widely used models for classification with incomplete data, namely Histogram-based Gradient Boosting (HistGB), XGB, and CatBoost. Across most metrics, ReMAC surpasses its competitors. Notably, by learning missingness-invariant representations, our method maintains strong performance regardless of the percentage of missing values.

| %MV | MODEL | ACC | SEN | SPE | AUC |
|---|---|---|---|---|---|
| 0% | HistGB | 0.68 ± 0.15 | 0.50 ± 0.41 | 0.83 ± 0.19 | 0.63 ± 0.16 |
| | XGB | 0.68 ± 0.15 | 0.50 ± 0.41 | 0.83 ± 0.19 | 0.58 ± 0.10 |
| | CatBoost | 0.72 ± 0.10 | **0.71 ± 0.34** | 0.75 ± 0.17 | 0.76 ± 0.14 |
| | ReMAC (*Ours*) | **0.86 ± 0.10** | 0.67 ± 0.24 | **1.00 ± 0.00** | **0.86 ± 0.10** |
| 10% | HistGB | 0.68 ± 0.28 | 0.63 ± 0.48 | 0.75 ± 0.17 | 0.71 ± 0.28 |
| | XGB | 0.68 ± 0.28 | 0.63 ± 0.48 | 0.75 ± 0.17 | 0.69 ± 0.29 |
| | CatBoost | 0.77 ± 0.18 | **0.71 ± 0.34** | 0.83 ± 0.19 | 0.81 ± 0.18 |
| | ReMAC (*Ours*) | **0.82 ± 0.14** | 0.67 ± 0.24 | **0.92 ± 0.17** | **0.86 ± 0.10** |
| 20% | HistGB | 0.68 ± 0.28 | 0.63 ± 0.48 | 0.75 ± 0.32 | 0.69 ± 0.28 |
| | XGB | 0.58 ± 0.17 | 0.50 ± 0.41 | 0.67 ± 0.27 | 0.60 ± 0.18 |
| | CatBoost | 0.72 ± 0.10 | **0.71 ± 0.34** | 0.75 ± 0.17 | 0.83 ± 0.19 |
| | ReMAC (*Ours*) | **0.81 ± 0.16** | 0.67 ± 0.24 | **0.92 ± 0.17** | **0.86 ± 0.10** |
| 30% | HistGB | 0.73 ± 0.22 | **0.71 ± 0.34** | 0.75 ± 0.17 | 0.67 ± 0.24 |
| | XGB | 0.68 ± 0.15 | **0.71 ± 0.34** | 0.67 ± 0.00 | 0.67 ± 0.19 |
| | CatBoost | 0.72 ± 0.10 | 0.58 ± 0.29 | 0.83 ± 0.19 | **0.89 ± 0.08** |
| | ReMAC (*Ours*) | **0.81 ± 0.16** | 0.54 ± 0.42 | **1.00 ± 0.00** | 0.86 ± 0.10 |

**Table 2:** *Performance of classification models capable of handling missing values, together with ReMAC, under different missing-value (%MV) regimes. Metrics include accuracy (ACC), sensitivity (SEN), specificity (SPE), and area under the curve (AUC). The best value for each metric is highlighted in bold, and the proposed method is shown with a gray background.*

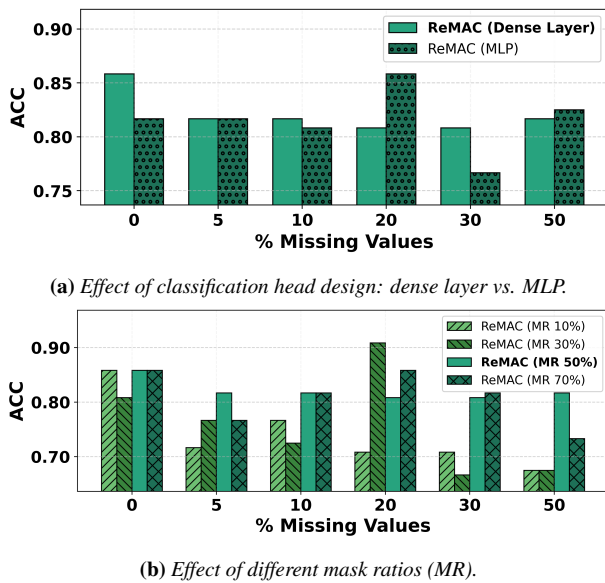### 4.3. Analysis of representation strategies

Different strategies for aggregating the latent representations of the decoder were explored for the classification task. Specifically, we considered the use of a learnable embedding appended to the sequence of Transformer decoder blocks (ReCLS), similar to the approach commonly adopted in Vision Transformers, as well as aggregation of decoder embeddings through max pooling (ReMaxAC) and mean pooling (ReMAC). The results are depicted in Fig. 2. It is worth noting that ReMAC not only achieves the best performance in the absence of missing values but also exhibits the highest stability across different missing-value regimes. Consequently, in this high-dimensional setting, mean pooling yields a more robust and discriminative latent representation space compared to other, seemingly more straightforward, alternatives such as ReCLS.



**Figure 2:** *Comparison of different representation strategies for decoder embeddings under varying missing-value (%MV) regimes. Accuracy (ACC) is reported for ReCLS (learnable token), ReMaxAC (max pooling), and ReMAC (mean pooling, ours).*

### 4.4. Impact of mask ratio and classification architecture

Figure 3 presents the results of the ablation studies evaluating the architectural design of the classification head and the mask ratio (MR) used to mask input features randomly. Regarding Fig. 3a, results indicate slightly better and more consistent performance when employing a single dense layer compared to a multilayer perceptron (MLP) with one hidden layer. This shows that a simple classification head is sufficient, providing stable generalization across multiple scenarios due to an already discriminative latent space. Concerning Fig. 3b, we observe that, across different missing-value regimes, ReMAC generally performs better with higher mask ratios (MR). These results are align with the intuition that larger mask ratios, by exposing the model to a greater number of masked variables, naturally encourage more consistent representation learning under missing values.



**(a)** *Effect of classification head design: dense layer vs. MLP.*



**(b)** *Effect of different mask ratios (MR).*

**Figure 3:** *Ablation Studies on ReMAC: impact of the classification head complexity and the mask ratio (MR) under varying missing-value regimes. In both subfigures, the configuration adopted in the main experiments is highlighted in bold in the legend.*

## 5. Conclusion

The diagnosis of spitzoid tumors is essential for ensuring timely and appropriate treatment. Epigenetic data, such as DNA methylation, have emerged as a valuable source of information for ST stratification. However, the frequent presence of missing values in data collection highlights the need for effective methodologies under such regimes. In this work, we introduced ReMAC, an extension of ReMasker tailored for high-dimensional classification tasks. Our results demonstrate that the proposed method achieves strong performance both on complete data and under varying degrees of incompleteness. The main limitations of this study include the small sample size, which is inherent to the problem under investigation, and the need to further assess the explainability of the method. These findings open a promising avenue for extending ReMAC to larger datasets and exploring its potential in providing interpretable predictions.

## References

[1] R. L. Barnhill, "The spitzoid lesion: rethinking spitz tumors, atypical variants, 'spitzoid melanoma' and risk assessment," *Modern pathology*, vol. 19, pp. S21–S33, 2006.

[2] D. C. Orchard, J. P. Dowling, and J. W. Kelly, "Spitz naevi misdiagnosed histologically as melanoma: prevalence and clinical profile," *Australasian journal of dermatology*, vol. 38, no. 1, pp. 12–14, 1997.

[3] D. Beck, M. Ben Maamar, and M. K. Skinner, "Genome-wide cpg density and dna methylation analysis method (medip, rrbs, and wgbs) comparisons," *Epigenetics*, vol. 17, no. 5, pp. 518–530, 2022.

[4] R. Del Amor, A. Colomer, C. Monteagudo, M. J. Garzón, J. L. García-Giménez, and V. Naranjo, "A deep embedded framework for spitzoid neoplasm classification using dna methylation data," in *2021 29th European Signal Processing Conference (EUSIPCO)*, pp. 1271–1275, IEEE, 2021.

[5] J. F. González-Muñoz, B. Sánchez-Sendra, and C. Monteagudo, "Diagnostic algorithm to subclassify atypical spitzoid tumors in low and high risk according to their methylation status," *International Journal of Molecular Sciences*, vol. 25, no. 1, p. 318, 2023.

[6] D. Seiler Vellame, I. Castanho, A. Dahir, J. Mill, and E. Hannon, "Characterizing the properties of bisulfite sequencing data: maximizing power and sensitivity to identify between-group differences in dna methylation," *BMC genomics*, vol. 22, no. 1, p. 446, 2021.

[7] S. Van Buuren and K. Groothuis-Oudshoorn, "mice: Multivariate imputation by chained equations in r," *Journal of statistical software*, vol. 45, pp. 1–67, 2011.

[8] T. Kyono, Y. Zhang, A. Bellot, and M. van der Schaar, "Miracle: Causally-aware imputation via learning missing data mechanisms," *Advances in Neural Information Processing Systems*, vol. 34, pp. 23806–23817, 2021.

[9] J. Yoon, J. Jordon, and M. Schaar, "Gain: Missing data imputation using generative adversarial nets," in *International conference on machine learning*, pp. 5689–5698, PMLR, 2018.

[10] A. Nazabal, P. M. Olmos, Z. Ghahramani, and I. Valera, "Handling incomplete heterogeneous data using vaes," *Pattern Recognition*, vol. 107, p. 107501, 2020.

[11] T. Du, L. Melis, and T. Wang, "Remasker: Imputing tabular data with masked autoencoding," *arXiv preprint arXiv:2309.13793*, 2023.

[12] S. Ö. Arik and T. Pfister, "Tabnet: Attentive interpretable tabular learning," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, pp. 6679–6687, 2021.