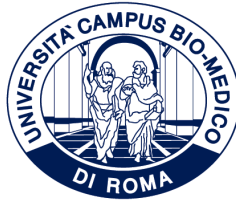


ID N. 32981



UNIVERSITÀ CAMPUS BIO-MEDICO DI ROMA

DEPARTMENT OF ENGINEERING

Italian National Ph.D. in Artificial Intelligence

Health and Life Sciences

XXXVII Cycle

**Resilient Multimodal Learning with
Incomplete Information in Biomedical
Applications**

Supervisors

Paolo Soda

Sara Ramella

Valerio Guarrasi

Candidate

Camillo Maria Caruso

July, 2025

To my Mom

Abstract

The deployment of Artificial Intelligence (AI) in clinical settings is increasingly viewed as a transformative step toward precision medicine and improved patient outcomes. However, the clinical reality of healthcare data is far from ideal. One of the most pervasive challenges remains the issue of missing and incomplete data, a systemic characteristic rather than a sporadic error. Conventional predictive models are often inadequately equipped to handle such irregularities, relying heavily on data preprocessing and imputation strategies that introduce biases and limit generalizability. This dissertation proposes a paradigm shift in the development of machine learning models, advocating for resilient learning systems, architectures inherently designed to operate under partial observability while maintaining high performance and stability. Rather than viewing missingness as a problem to be corrected, the proposed approach treats it as a fundamental property of clinical data that must be embraced and strategically incorporated into the learning process. A key focus is the application of these principles to multimodal learning, where information is drawn from heterogeneous sources, including clinical records, laboratory tests, and radiological images. In such contexts, the absence of entire modalities can significantly hinder model performance if not explicitly addressed. This thesis explores how to design multimodal systems resilient to missing and incomplete data.

To this end, the dissertation first presents NAIM (Not Another Imputation Method), a transformer-based model tailored for tabular data characterized by structural missingness. NAIM leverages a combination of feature-specific embeddings, a masked self-attention mechanism, and a novel regularization strategy that randomly masks inputs during training. This framework enables the model to learn directly from incomplete feature sets, avoiding the biases introduced by imputation. Indeed, unlike traditional methods that depend on complex imputation pipelines, NAIM bypasses the reconstruction step entirely, training directly on what is observed. The model is benchmarked against a wide spectrum of state-of-the-art machine learning and deep learning models across five public classification datasets, as well as on a real-world clinical application involving the prediction of overall survival in cancer patients. Across both experimental settings, NAIM demonstrates strong and consistent per-

formance, particularly under high levels of missingness. These results highlight the models capacity to generalize to realistic clinical scenarios and confirm its practical relevance for resilient AI systems in healthcare.

Building on these insights, the thesis then explores the transition to multimodal modeling by investigating whether combining radiological images (CT scans) and structured clinical data can enhance predictive performance in a real-world setting. A late fusion ensemble approach is employed, where each modality is processed independently and combined at the decision level. This setup reflects a common and practical design in clinical AI systems and ensures a certain degree of modularity. However, by treating modalities independently until the final stage, this approach cannot fully leverage the interdependencies and shared patterns across data types. This observation, supported by the findings of our systematic review on multimodal fusion strategies, underscores the need for more integrated fusion mechanisms capable of modeling cross-modal interactions throughout the learning process.

To overcome this limitation, the dissertation introduces MARIA (Multimodal Attention Resilient to Incomplete data), a novel transformer-based architecture designed for resilient multimodal data fusion. MARIA integrates modality-specific encoders with a shared attention-based fusion module, using a generalized masking mechanism to handle missing modalities. This intermediate fusion design allows the model to dynamically learn interactions between available inputs while maintaining resilience to missing components. Empirical evaluations across a range of simulated missing data scenarios confirm MARIA's ability to deliver accurate and stable predictions, even under severe data fragmentation.

Altogether, this work outlines a unified methodological framework for designing AI systems that are resilient to data incompleteness in both unimodal and multimodal settings. Rather than relying on preprocessing fixes, resilience is embedded directly into model architectures and training objectives. This design philosophy aligns more closely with the operational constraints of modern clinical environments, where decisions often need to be made based on partial information and data cannot always be recollected or completed. By embracing the imperfection and variability inherent in healthcare data, this dissertation lays the groundwork for the development of practical, adaptable, and resilient AI models, capable of supporting clinicians even in the most uncertain and data-sparse scenarios.

Contents

1	Toward Resilient AI in Clinical Practice	10
2	NAIM: Resilient Tabular Modeling under Structural Missingness	14
2.1	Introduction	14
2.2	State-of-the-art	16
2.2.1	SOTA techniques for missing values	16
2.2.2	Transformer for tabular data	17
2.2.3	Current limitations	18
2.3	Methods	19
2.3.1	Problem definition	19
2.3.2	NAIM architecture	19
2.3.3	Regularization technique	23
2.4	Experimental Configuration	25
2.4.1	Data and Evaluation	25
2.4.2	Competitors	27
2.4.3	Evaluation Metrics and Statistical Analysis	30
2.5	Results and Discussions	30
2.6	Conclusion	36
3	Survival Prediction from Incomplete Clinical Data	40
3.1	Introduction	40
3.2	Background	42
3.3	Materials	44
3.4	Methods	46
3.4.1	Model	46
3.4.2	Training and Testing	48
3.4.3	Evaluation Metric	49
3.4.4	Experimental setup	50

3.5	Results and Discussions	51
3.6	Conclusion	55
4	Fusing Imaging and Clinical Data for NSCLC Survival Prediction	57
4.1	Introduction	57
4.2	Materials	61
4.2.1	Imaging	62
4.2.2	Clinical Features	62
4.3	Methods	62
4.3.1	Training	64
4.3.2	Optimisation	66
4.3.3	Preprocessing	67
4.4	Results and Discussion	68
4.5	Conclusions	72
5	MARIA: a Multimodal Transformer Resilient to Missing Modalities	74
5.1	Introduction	74
5.2	State of the Art	76
5.2.1	Early Fusion	77
5.2.2	Late Fusion	78
5.2.3	Intermediate Fusion	79
5.2.4	Handling Incomplete Data	80
5.3	Methods	81
5.3.1	Model	81
5.3.2	Regularization Technique for Missing Data	84
5.4	Experimental Configuration	85
5.4.1	Data	85
5.4.2	Competitors	87
5.4.3	Preprocessing	89
5.4.4	Missingness Evaluation	89
5.4.5	Evaluation Metrics	90
5.4.6	Fusion Analysis	91
5.5	Results and Discussions	92
5.5.1	MARIA vs. ML	92
5.5.2	MARIA vs. DL	95
5.5.3	MARIA vs. NAIM	97
5.6	Conclusions	99

6	Conclusions and Future Perspectives on Resilient Clinical AI	101
6.1	Limitations and Future Works	103
	Appendices	120
A	Masking Example	120
B	Training of the Models	124
C	NAIM Complete Results	126
D	NAIM Statistical analysis	132
E	NAIM Computational analysis	134
F	MARIA Complete Results	136
G	MARIA Statistical Analysis	138

List of Figures

2.1	NAIM’s architecture.	20
2.2	Proposed <i>Feature Embedding</i> process for tabular data.	21
2.3	The proposed masked self-attention mechanism.	24
2.4	Proposed regularization strategy.	25
2.5	Comparison between NAIM and the competitor models.	32
2.6	Comparison between NAIM and the competitors implementing an intrinsic strategy.	33
2.7	Comparison between NAIM and the competitor imputers.	34
2.8	Analysis of the robustness of the models to increasing levels of missing data.	35
2.9	Ablation analysis evaluating the impact of the regularization technique and the masked self-attention mechanism.	37
3.1	Schematic representation of the proposed model.	47
3.2	Prediction errors made for different patient groups.	53
3.3	SHAP summary plots of features’ contributions in the 3 models implemented with the 3 time units.	54
3.4	Ablation study of the two terms of the loss function proposed in [49].	55
4.1	Schematic view of the pipeline.	64
5.1	Overview of multimodal fusion strategies in DL.	78
5.2	MARIA architecture.	82
5.3	MARIA vs. ML in the “missing modalities” scenario.	93
5.4	MARIA vs. ML in the “all missing” scenario.	94
5.5	MARIA vs. DL in the “missing modalities” scenario.	95
5.6	MARIA vs. DL in the “all missing” scenario.	96
5.7	MARIA vs. NAIM in the “missing modalities” scenario.	97
5.8	MARIA vs. NAIM in the “all missing” scenario.	98

List of Tables

2.1	Datasets’ details and references.	26
2.2	Combinations of competitors’ models and imputers used in the experiments.	28
3.1	Patients’ characteristics.	45
3.2	Performance of tested models.	52
4.1	Summary of the background on the multimodal learning to predict the overall survival in NSCLC.	60
4.2	Patients’ characteristics.	63
4.3	Performance of all the tested models.	70
5.1	Datasets’ details.	86
5.2	Combinations of models, missing techniques and fusion strategies used as competitors.	88
C.1	Average performance of the experiments across the 5 different datasets.	126
C.2	Average performance of the experiments on the ADULT dataset.	127
C.3	Average performance of the experiments on the BankMarketing dataset.	128
C.4	Average performance of the experiments on the OnlineShoppers dataset.	129
C.5	Average performance of the experiments on the SeismicBumps dataset.	130
C.6	Average performance of the experiments on the Spambase dataset.	131
D.1	Percentages of wins and losses of NAIM correct predictions compared with the competitors using the Wilcoxon signed-rank test.	133
E.1	Computational complexity analysis of the evaluated deep learning models in terms of FLOPs and number of trainable parameters across datasets.	134
F.1	Average AUC and MCC performance of the experiments across the respective tasks in the “missing modalities” setting.	136

F.2	Average AUC and MCC performance of the experiments across the respective tasks in the “all missing” setting.	137
G.1	Percentages of wins and losses of MARIA correct predictions compared with the competitors in the “missing modalities” scenario using the Wilcoxon signed-rank test.	139
G.2	Percentages of wins and losses of MARIA correct predictions compared with the competitors in the “all missing” scenario using the Wilcoxon signed-rank test.	140

Chapter 1

Toward Resilient AI in Clinical Practice

Artificial intelligence (AI) is rapidly redefining the landscape of medical research and clinical practice, offering transformative opportunities to enhance the quality, speed, and precision of healthcare delivery. In recent years, the convergence of increasing computational power, an increasing volume of healthcare data, and breakthroughs in deep learning has fueled the development of predictive models capable of extracting complex, non-linear patterns from multimodal biomedical information. These models have demonstrated remarkable performance in a wide array of applications, including disease diagnosis, prognosis, patient stratification, and personalized treatment planning. In high-stakes domains such as oncology, neurology, and cardiology, AI-driven tools have shown the potential to augment clinical decision-making, reduce diagnostic delays, and optimize therapeutic interventions.

This technological progress has been accompanied by a growing interest in deploying AI not only as a research tool but as a core component of modern clinical workflows. Algorithms capable of processing heterogeneous data sources, such as laboratory tests, electronic health records, radiological images, and genomic sequences, are now envisioned as pivotal enablers of precision medicine. The ability to synthesize multimodal data promises a deeper and more holistic understanding of patient health, allowing clinicians to anticipate risks, tailor interventions, and allocate resources more effectively. However, despite this promising scenario, the real-world implementation of AI systems in clinical settings remains filled with substantial challenges. Among the most pervasive and underappreciated obstacles is the incompleteness of clinical data. Unlike curated research datasets, data generated in routine clinical practice are often messy, irregular, and partially observed. Missing values in structured variables, irregularly sampled time-series, incomplete follow-up records, unavail-

able imaging modalities, and inconsistent documentation practices are widespread. These issues stem from a multitude of factors, including diverse institutional protocols, patient heterogeneity, non-standardized data collection, and resource limitations. Additionally, data availability is often hindered by privacy regulations, variations in technology infrastructure, and clinical priorities that do not align with optimal data acquisition. As such, incompleteness is not merely an incidental data quality problem, but rather a fundamental characteristic of clinical information systems.

The presence of missing data poses significant threats to the performance, reliability, and fairness of machine learning models. When not properly addressed, it can lead to biased estimates, reduced statistical power, and impaired generalizability. Traditional methods for handling missingness, such as mean imputation or model-based reconstruction, aim to “complete” the dataset before training. However, these techniques often overlook the underlying structure of the missingness and may introduce artifacts or spurious correlations. Furthermore, they rarely capture the uncertainty associated with imputed values, which can propagate into downstream predictions and decision-making processes. These limitations become particularly pronounced in modern multimodal settings, where different data modalities, such as clinical records, imaging, genomics, and time-series, may be partially or entirely unavailable for subsets of the population.

In light of these challenges, this thesis proposes a paradigm shift: to design AI models that are intrinsically robust to missing data, embedding resilience directly into the model architecture rather than relying on preprocessing or imputation. The goal is to treat missingness not as an error to be corrected, but as a property to be modeled and leveraged. To achieve this, the work presented herein explores a progressive methodological pathway, beginning with the development of models for structured tabular data and advancing toward the design of robust multimodal architectures that can integrate heterogeneous sources of information even in the presence of incomplete observations.

This thesis seeks to bridge the gap between theoretical advances in machine learning and their practical application in clinical contexts by developing systems capable of adapting to the often messy and incomplete nature of real-world healthcare data. Rather than focusing solely on maximizing accuracy under ideal conditions, the models proposed here are designed to remain reliable and effective when faced with the variability, irregularity, and missingness that characterize clinical datasets. Ensuring that these models can function accurately and consistently in such challenging settings is essential for their safe and meaningful use in high-stakes medical decision-making.

A central theme of this thesis is the use of masked self-attention mechanisms and modality-aware fusion strategies. These techniques enable models to condition their predictions on

the available data while ignoring the missing components, thus avoiding the pitfalls of naive imputation. Additionally, several training procedures are introduced to simulate missingness and encourage the model to learn representations that generalize well to incomplete scenarios. These include regularization methods that artificially mask inputs during training, as well as learning objectives that reward the model for performance on partially observed data. Such strategies ensure that the model is not only evaluated but also trained under conditions that mirror real-world incompleteness.

By jointly considering architecture, training strategy, and evaluation framework, this thesis aims to provide a comprehensive and principled approach to robust learning under partial observability. This effort culminates in a series of models and experimental studies that demonstrate how AI systems can be made more resilient to one of the most persistent barriers in medical data science: missingness.

The structure of the thesis reflects this methodological progression and is composed of a collection of research articles, each tackling a key aspect of the robust modeling pipeline:

Chapter 2 introduces NAIM (Not Another Imputation Method), a transformer-based model designed to handle incomplete tabular data. NAIM incorporates masked self-attention and a data corruption regularization scheme to promote robustness during both training and inference. The model is benchmarked against state-of-the-art approaches on multiple datasets, demonstrating competitive performance under varying levels and patterns of missingness.

Chapter 3 shifts focus to a real-world clinical application: the prediction of overall survival (OS) in cancer patients using routinely collected clinical data. This chapter investigates the challenges of irregular temporal sampling, partial histories, and the interpretability of model predictions. It provides insights into how models can be made clinically useful while dealing with highly imperfect and noisy data.

Chapter 4 marks the transition into multimodal modeling, combining clinical features and radiological images (CT scans) to improve survival prediction in non-small cell lung cancer (NSCLC). This study evaluates several fusion strategies and illustrates the potential benefits of multimodal systems.

Chapter 5 introduces MARIA (Multimodal Attention Resilient to Incomplete data), a novel multimodal transformer that generalizes the masked attention concept to multiple input types. By integrating intermediate fusion and modality-specific masking, MARIA enables predictive modeling even when certain modalities are completely absent. Extensive experiments under various missing-modality scenarios confirm the model's ability to maintain strong performance and stability, establishing it as a benchmark for robust multimodal learning in healthcare.

Finally, the concluding chapter synthesizes the key insights gathered across all contributions, discussing their implications for the development of clinically useful AI systems. It also outlines open challenges and future directions, emphasizing the need for continued progress toward resilient, multimodal, and real-world-ready machine learning models in healthcare.

Chapter 2

NAIM: Resilient Tabular Modeling under Structural Missingness

2.1 Introduction

Tabular data, characterized by data structured in tables with rows and columns, i.e., samples and features, respectively, poses unique challenges in training and testing artificial intelligence (AI) models because they differ from other data structures, such as text and speech that have a sequential nature or the images, which are characterized by spatial coherence. These challenges arise because, unlike in other domains, tabular data often consists of a *heterogeneous mix of categorical and numerical features*. Additionally, each feature can have a *different distribution and scale*, necessitating careful preprocessing steps such as normalization or embedding. Furthermore, tabular datasets frequently contain *missing values*, requiring models that can robustly handle such inconsistencies. To date, machine learning (ML) ruled over deep learning (DL) in many fields involving tabular data, offering methods capable of meeting each of these challenges [1, 2].

More specifically, the missing data issue poses a significant challenge in the realm of tabular data and common reasons for this event are: *i*) human error, where inaccuracies during data entry or collection lead to gaps in information; *ii*) non-response, common in surveys, where participants might skip questions for reasons like privacy or lack of interest; *iii*) data corruption, caused by technical failures or errors; *iv*) attrition, particularly in longitudinal studies, caused by participants drop out; *v*) systematic loss, where data is missing under specific conditions. Each of these factors could play a significant role in the challenge associated with handling missing values in tabular data: this absence of information may affect either training, testing, or both sets at the same time, whereas most of the state-

of-the-art approaches need a complete dataset to function properly. This brings the need for approaches that ignore the missing values and use all the available data to perform the task at hand. The state-of-the-art faced the problem using two main approaches [1, 2]: on the one hand, some methods fill in the missing entries, usually during the preprocessing phase before feeding the data to the model, by imputing the values using a predefined rule or an algorithm [1]; on the other hand, other approaches intrinsically handle missing data during inference time.

Recently, we have witnessed the successes of the transformer architecture in different domains. Since the transformer’s core mechanism of self-attention has been adapted to computer vision [3] and speech recognition [4], some studies have started adapting this architecture to tabular data [5, 6, 7]; although, none has proposed specific solutions to handle missing values. This is why in this manuscript we present “*Not Another Imputation Method*” (NAIM), a novel transformer-based model for tabular data, specifically designed to face the missing values problem. The main contributions of this work are:

- The development of a transformer model that integrates feature-specific embeddings and a novel masked self-attention mechanism able to learn only from the available information without any need for imputation of missing values.
- The proposal of a novel regularization technique, which randomly masks each sample at every epoch. This method, presenting a different sample version at each epoch, enables the model to generalize from incomplete data, ensuring more resilient and accurate predictive performance.
- A wide experimental assessment, showing that NAIM achieves better performance than state-of-the-art models in 5 publicly available classification tasks. State-of-the-art approaches consist of 6 ML and 5 DL models, paired with 3 different imputation techniques, when necessary. Results indicate a noteworthy advancement of transformer and not transformer-based DL approaches over ML in addressing the domain of tabular data, highlighting its efficacy and potential in this area.

The manuscript is organized as follows: Section 2.2 presents the state-of-the-art of data imputation techniques and those models able to handle missing values (Section 2.2.1) and the main adaptations of transformers to the tabular data domain (Section 2.2.2). Section 2.3 introduces the proposed model, whereas Sections 2.4 and 2.5 discuss the experimental configuration and the results, respectively. Section 2.6 provides concluding remarks.

2.2 State-of-the-art

Despite the advancements in handling missing values in tabular data, the state-of-the-art shows the following limitations. Current methods primarily focus on imputing missing entries or adapting models to work with incomplete data. However, these strategies often fail to fully leverage the inherent patterns in the data or require extensive preprocessing, potentially leading to information loss or biased predictions. Moreover, while the transformer architecture has shown promising results in various domains, its application to tabular data with missing values has not been realized, indicating a gap in effectively addressing this challenge within the deep learning domain.

This section first presents the most established strategies to address missing values in tabular data (Section 2.2.1) and then we explore existing transformer-based methods for tabular data (Section 2.2.2). These methodologies are then employed as benchmarks in the experiments detailed in Sections 2.4 and 2.5.

2.2.1 SOTA techniques for missing values

The simplest way to handle missing values is the *Complete Case Analysis* [8] that completely excludes samples and/or features with missing values; despite its simplicity, its use can cause a considerable loss of information for the model, which cannot be tolerated to make AI models resilient to be used in practice. Consequently, several approaches have been developed to retain as much information as possible.

ML approaches

State-of-the-art techniques for managing missing values mainly focus either on imputing missing values [1, 8, 9] or on applying a specific strategy for missing entries [1, 9]. In [1], where the authors benchmark different approaches to handle missing values, the imputation techniques are distinguished in *Constant Imputation* and *Conditional Imputation*. The first, despite its simplicity, can effectively support models by replacing missing values with a central tendency measure, e.g., mean imputation; the latter comprises techniques recognized for leveraging patterns within the data to predict missing values accurately, e.g., Multiple Imputation by Chained Equations (MICE) [10] and K-Nearest Neighbors imputation (KNN) [11]. Furthermore, the authors in [1] also introduce the *Missing Incorporated in Attributes* (MIA) strategy, specifically advantageous for tree-based models, which incorporate missingness directly into the model, thus preserving all available data. By presenting some of the main possible approaches to handle missing values, this analysis underlines the multitude of avail-

able options and the need to identify the most suitable one for the specific task at hand. Furthermore, the prevalence of imputation approaches over models that can handle missing values, which are tree-based only, highlights the need for approaches that can ignore missing values rather than impute them. Indeed, identifying the most suitable imputation method for a specific downstream task in advance is a considerable challenge. Ideally, it would be better to rely solely on the data at hand, avoiding the introduction of potentially misleading information. This approach, unlike imputation strategies that rely on the training set, would ensure a robust analysis, ignoring missing values entirely and making the most of the available data.

DL approaches

DL approaches have introduced new perspectives and possibilities compared to traditional ML methods, particularly in handling complex and structured data scenarios. A representative example of this advancement in the missing data handling problem is GRAPE [12], a graph neural network (GNN)-based architecture specifically designed for tabular data analysis able to simultaneously perform feature imputation and label prediction. GRAPE, remapping the dataset onto a bipartite graph composed of sample nodes and feature nodes, intrinsically handles missing values by simply omitting edges between nodes where data is unavailable. This setting lets GRAPE formulate feature imputation as an edge-level prediction problem and classification as a node-level prediction task.

Nevertheless, in recent years, transformer models, originally developed for natural language processing (NLP) [13], have become predominant in various application domains, including computer vision [3] and speech recognition [4]. Although some extensions have been proposed to analyze tabular data, many transformer-based approaches still struggle to effectively manage missing data.

2.2.2 Transformer for tabular data

Currently, there are 3 main transformer-based approaches, each developing a different aspect of the transformer architecture to adapt it to tabular data: TabNet [5], TabTransformer [6], and FTTransformer [7].

TabNet [5] leverages the self-attention mechanism to selectively focus on specific features for decision-making, which is crucial in interpreting high-dimensional data. Learning a sparse masking matrix and mimicking the boosting technique by adding a layer at each step, TabNet can perform a dynamic feature selection. Its design facilitates interpretability, allowing it to offer insights into feature importance and decision pathways, a critical requirement in fields

like finance and healthcare where understanding model decisions is as important as their accuracy.

TabTransformer [6] addresses the embedding of categorical features within tabular data. It applies a transformer-based self-attention mechanism to create contextual embeddings for categorical features, similar to how transformers process words in NLP. By doing so, it captures complex inter-feature relationships and dependencies, enhancing the model’s predictive performance. Given $x \in \mathbb{R}^{n \times m}$ as a mixture of categorical x^{cat} and numerical x^{num} features, each categorical feature x_i^{cat} , where i denotes the i -th feature, is encoded twice: on the one hand, a feature-specific lookup table E_i^{pos} encodes the value with respect only to the possible categorical values of that feature; on the second hand all categorical values are encoded using a shared lookup table E that embeds all the categorical features’ values. The resulting embedding vector for the i -th categorical feature is denoted as e_i^{cat} :

$$e_i^{\text{cat}} = [E_i^{\text{pos}}(x_i^{\text{cat}}), E(x_i^{\text{cat}})]. \quad (2.1)$$

Instead, the numerical features x^{num} simply undergo a normalization step.

FTTransformer [7] further explores the potential of transformer models in extracting patterns and interactions within tabular data, which are often overlooked by traditional machine learning methods. Specifically, this model addresses the embedding of numerical features applying two different approaches for the numerical and categorical features. The embedding e_i^{num} for a given numerical feature x_i^{num} is computed as follows:

$$e_i^{\text{num}} = b_i + x_i^{\text{num}} \cdot W_i^{\text{num}}, \quad x_i^{\text{num}} \in \mathbb{R}, \quad e_i^{\text{num}} \in \mathbb{R}^{d_e} \quad (2.2)$$

where the d_e -dimensional embedding vector is computed as the sum of the i -th feature bias, $b_i \in \mathbb{R}^{d_e}$, and the multiplication of the feature value with the feature-specific vector $W_i^{\text{num}} \in \mathbb{R}^{d_e}$. Instead, the embedding e_i^{cat} for a given categorical feature x_i^{cat} is implemented as the lookup table E_i^{cat} , which associates a d_e -dimensional trainable vector to the k_i entries of the specific categorical feature:

$$e_i^{\text{cat}} = b_i + E_i^{\text{cat}}(x_i^{\text{cat}}), \quad x_i^{\text{cat}} \in \mathbb{N}_{k_i}, \quad e_i^{\text{cat}} \in \mathbb{R}^{d_e}. \quad (2.3)$$

2.2.3 Current limitations

As illustrated in the previous sections, current state-of-the-art methods still have several limitations. Among traditional ML approaches, none of the methods truly ignore missing values; the MIA strategy, which does not impute missing data, attempts to incorporate missingness

into the learning process rather than bypassing it completely. Regarding DL, GRAPE is the only architecture that can explicitly ignore missing features, but its dependence on GNN structures limits its integration and scalability. In contrast, transformer-based models, not having dedicated embeddings and a masking mechanism, cannot handle missing data. This shortcoming highlights the need for transformer-based solutions that can robustly ignore missing features while maintaining architectural flexibility and compatibility with broader modeling frameworks.

2.3 Methods

This section first presents a formal definition of the problem; then it examines the proposed NAIM architecture (Section 2.3.2), presenting the innovation we made to the embedding of the features and the novel masked self-attention mechanism to dynamically handle missing values; finally, in section 2.3.3, we introduce a novel regularization strategy that aims to enable the model to learn how to handle missing data even if they are not present in the training set, and to extract the optimal representation for each feature by preventing any co-adaptation between them.

2.3.1 Problem definition

Let $X \in \mathbb{R}^{n \times m}$ represent a data matrix with n samples and m features, where the element x_i^{num} corresponds to a numerical feature of the i -th sample, while x_i^{cat} represents a categorical feature of the same sample. In scenarios involving missing data, some of these feature values are unavailable. This is captured, for each sample separately, by a binary mask matrix $M \in \{-\infty, 0\}^{m \times m}$, where the entry $M_{ij} = 0$ indicates that the i -th feature is observed and can attend to the j -th feature, while $M_{ij} = -\infty$ prevents attention from the i -th feature to the j -th feature when the latter is missing. Typically, datasets are also associated with labels for a specific downstream task, which we denote as $Y \in \mathbb{R}^n$. Our focus is on predicting labels for test instances.

2.3.2 NAIM architecture

The transformer architecture [13], originally developed for text translation, includes two main blocks: the first, known as the *encoder*, extracts a hidden representation from the data, whilst the second, named the *decoder*, reconstructs the hidden representation in the context of the specific domain. For instance, in the NLP domain, the decoder generates the next token or performs complex interactions given an input sequence. Nevertheless, the

decoder may not be required in scenarios where decision-making is the primary objective. Given that most of the tasks related to tabular data are classification or regression, the transformer-based models mostly make use of the *encoder* part only with a fully-connected (*FC*) module for the final prediction [5, 6, 7]. However, to perform other tasks, e.g., features imputation or self-supervised pretraining, the *decoder* part could also be used to reconstruct the feature representation in another domain [5]. Indeed, when developing NAIM to get a model that handles missing features and learns from the available information without any need for imputation to perform a classification task, we opted for the encoder-only architecture followed by a *FC* module (Figure 2.1), which receives as input the normalized and concatenated version of the *encoder* output (grey block).

Our approach exploits the *Feature Embedding* (yellow block) and *Masked Multi-Head Attention* (red block) steps reported in the figure. More specifically, our idea stems from observing that the use of the padding index in the lookup table paired with the masked self-attention mechanism, could be an interesting and unexplored way to handle missing entries for both the categorical and numerical types of features in tabular data.

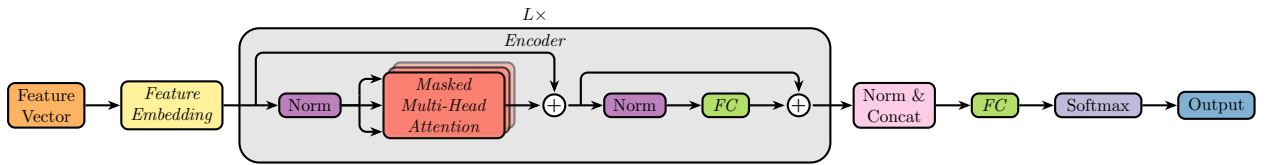


Figure 2.1: The architecture of NAIM, composed of the *Feature Embedding*, the *Encoder* equipped with the *Masked Multi-Head Attention* mechanism, and the final classification head [13].

Focusing on the *Feature Embedding*, which is a fundamental step to find a richer representation of the input data useful in the downstream task, we observed that the use of the padding index is reported in literature [6, 7], but none are able to handle missing values for both categorical and numerical features. Indeed, with the embedding lookup table E_i paired with the masked self-attention, we can completely exclude the missing values from the gradient calculations defining the $k_i + 1^{th}$ entry, named *padding index* or $\langle \text{pad} \rangle$, which is assigned to a not-trainable vector of zeros. Using the notation introduced in section 2.2.2, the embedding of categorical features e_i^{cat} is given by:

$$e_i^{cat} = b_i + E_i^{cat}(x_i^{cat}), \quad x_i^{cat} \in \mathbb{N}_{k_i}, \quad e_i^{cat} \in \mathbb{R}^{d_e} \quad (2.4)$$

whereas we defined an embedding lookup table E_i^{num} for each numerical feature with 2 possible entries, named *present* and *missing*. Using the padding index of the lookup table, these entries are associated with a trainable and a not-trainable d_e -dimensional vector, re-

spectively. Then, we scaled the numerical embedding vector e_i^{num} by multiplying it by the feature value x_i^{num} when present:

$$e_i^{num} = b_i + x_i^{num} \cdot E_i^{num}(x_i^{num}), \quad x_i^{num} \in \mathbb{R}, \quad e_i^{num} \in \mathbb{R}^{d_e}. \quad (2.5)$$

This setup, further described with an example reported in Figure 2.2, ensures that both types of features can be handled when presenting missing values without impacting the learning process. The figure provides an example of a feature vector with 2 categorical (x^{cat}) and 2 numerical (x^{num}) features. For each of these 2 types of features, there are one missing and one non-missing value, which are then embedded using the respective look-up tables in the *Feature Embedding* block. In particular, we can notice that both the missing features are encoded using the entry relative to the missing value of the look-up tables, represented with the **?** symbol. After obtaining the embedding vectors of the features, these are concatenated in the embedded representation e , which will be fed to the encoder later in the model architecture.

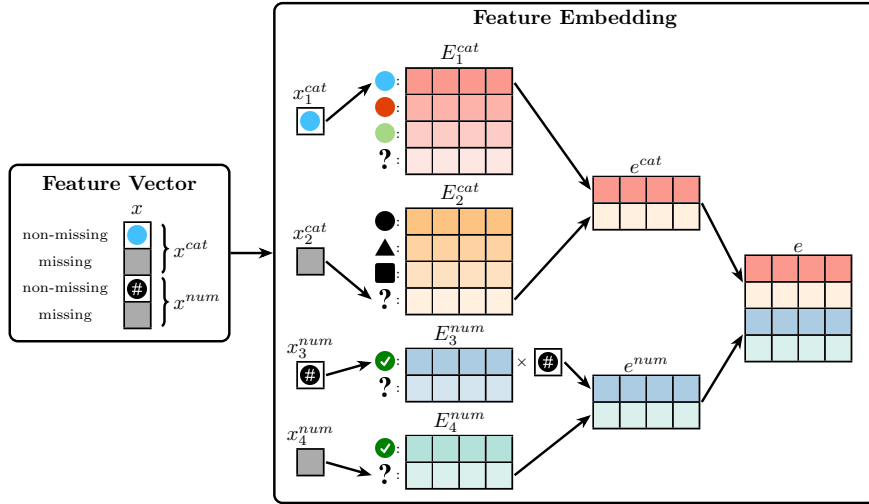


Figure 2.2: The proposed *Feature Embedding* process for tabular data. In the example, the feature vector x has 4 features: 2 categorical (x^{cat}) and 2 numerical (x^{num}). The colors (●, ●, ●) and the shapes (●, ▲, ■) are examples of possible values for the first categorical feature x_1^{cat} and for the second one x_2^{cat} , respectively, whilst ✓ stands for a non-missing numerical feature and # for its value. Finally, ? indicates the padding index related to missing features for both types of features. In the *Feature Embedding* block, we can see how the embedding e of the feature vector x is composed of the concatenation of embedded representations of the categorical and numerical features, denoted as e^{cat} and e^{num} , respectively. These representations are composed by the concatenation of the vectors associated with each feature value, selected using the feature-specific lookup tables E_i^{cat} and E_i^{num} .

Focusing on the masked self-attention mechanism to adapt it effectively for tabular data

and to completely mask out the contribution of the missing features, we first argued that it is necessary to modify it. Indeed, masked self-attention, a variant of the standard attention mechanism, plays a key role: it allows the model to focus only on certain parts of the input sequence while ignoring others. Traditionally, after the *Feature Embedding* step, the transformer calculates the query, key, and value matrices, denoted as Q , K , and V respectively, through linear transformations of e , i.e., mapping the embedding into a smaller space d_h based on the chosen number of heads h :

$$\begin{cases} Q = e \cdot W^Q, & d_h = d_e/h \\ K = e \cdot W^K, & W^Q, W^K, W^V \in \mathbb{R}^{d_e \times d_h} \\ V = e \cdot W^V, & Q, K, V \in \mathbb{R}^{n \times d_h} \end{cases} \quad (2.6)$$

where W^Q , W^K and W^V are weights matrices learned during training. Then, in the standard masked self-attention,

$$Attention(Q, K, V) = \underbrace{\text{softmax} \left(\frac{QK^T}{\sqrt{d_h}} + M \right)}_{A \in \mathbb{R}^{n \times n}} V \quad (2.7)$$

the attention matrix A uses the masking matrix $M \in \mathbb{R}^{n \times n}$ that sums $-\infty$ to the weights that should be ignored, zeroing their influence after the *softmax* operation. In the NLP domain, this technique is mainly employed to achieve 2 different purposes: on the one hand, the use of a mask M such as

$$M_{ij} = \begin{cases} 0 & \text{if } i \leq j \\ -\infty & \text{if } i > j \end{cases}, \quad M = \begin{bmatrix} 0 & -\infty & \dots & -\infty & -\infty \\ 0 & 0 & \dots & -\infty & -\infty \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & -\infty \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix} \quad (2.8)$$

allows the model to ignore future positions in the sequence, avoiding any leakage of information when training the model with the same sentence at different stages of completion; on the other hand, it could be used a mask M that cancels out the contributions of the last

tokens

$$M_{ij} = \begin{cases} 0 & \text{if } x_j \neq \langle \text{pad} \rangle \\ -\infty & \text{if } x_j = \langle \text{pad} \rangle \end{cases}, \quad M = \begin{bmatrix} 0 & 0 & \dots & 0 & -\infty \\ 0 & 0 & \dots & 0 & -\infty \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & -\infty \\ 0 & 0 & \dots & 0 & -\infty \end{bmatrix} \quad (2.9)$$

allowing the model to analyze sequences shorter than the dimension n used for training.

Considering the structure of these 2 possible masks, it is straightforward that only the latter could be adapted to mask missing features, since the former would mask large portions of the samples, even related to features that are non-missing. On this ground, the use of the Eq. 2.7 paired with the mask M in Eq. 2.9 used to mask the columns related to the missing values, might seem a viable solution, but, as shown in Figure 2.3 with colors and reported in Appendix A with an example, it does not cancel out the undesired value. Indeed, in Figure 2.3 we observe both the application of traditional masked attention and our new proposal: starting from left to right we have the matrices Q , K and V generated from the example feature vector shown in Figure 2.2, where the second and fourth features were missing. Immediately afterward, we report the attention matrix resulting from the calculation of the QK^T product, which is useful for understanding the next steps and especially how the contributions of the different features mix up together. In the next step, we show that the traditional masked attention mechanism, by masking the second and fourth columns, eliminates some of the contributions of the missing features, but leaves others on the second and fourth rows. For this reason, we here propose a new masked self-attention designed to completely mask out the contributions of the missing values using twice the same mask M from equation 2.9:

$$Attention(Q, K, V) = ReLU \left(softmax \left(\frac{QK^T}{\sqrt{d_h}} + M \right) + M^T \right) V \quad (2.10)$$

In this way, the attention matrix’s rows and columns related to missing values are now assigned zero attention, ignoring their contribution.

2.3.3 Regularization technique

Given the capacity of NAIM to effectively handle missing features, we introduce a novel regularization technique designed to enhance model robustness by simulating missing data scenarios. This approach, inspired by the Cutout technique [14] randomly masks elements

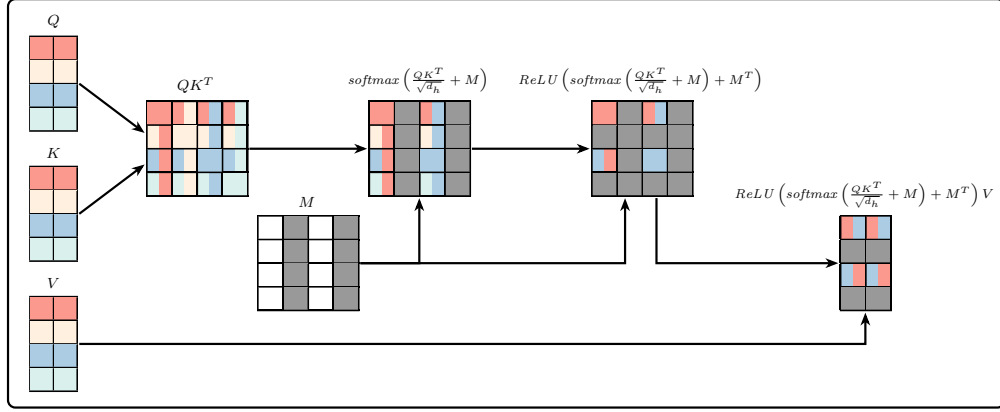


Figure 2.3: The proposed masked self-attention mechanism, designed to effectively ignore the impact of missing entries within the attention matrix. In the example, the QK^T matrix, obtained by the multiplication of the Q and K representations, is reported as an example of how the contributions of the different features, identified by different colors, mix up together. Next, the classic masked self-attention mechanism is applied, and some of the contributions of the missing features (indicated with \square and \square) remain. Then, the proposed attention mechanism ensures that the influence of these missing values is completely masked out, before the multiplication by the representation V of the sample.

within feature vectors (Figure 2.4), taking into account real-world data’s inherent variability and incompleteness. By integrating this regularization technique, we aim to improve the model’s ability to generalize from incomplete patterns, thus enhancing predictive performance on datasets with varying degrees of missing information.

Given the feature vector $x \in \mathbb{R}^n$, where n represents the dimensionality of the vector, let $v \leq n$ be the number of non-missing elements. The vector may be fully populated with $v = n$ or contain missing values with $v < n$. The process is governed by a binary decision variable $B \sim \text{Bernoulli}(0.5)$, which determines whether the masking will be applied to the instance x . If a sample is selected for masking, a random count of c elements to mask is chosen uniformly from the set $\{1, 2, \dots, v - 1\}$, ensuring that at least one element remains unmasked. Finally, c non-missing elements within x are randomly chosen and set their values to *missing*, resulting in the augmented vector.

The masking operation simulates scenarios of incomplete data, a common challenge in practical applications. This methodological approach is systematically applied during the training phase, ensuring that each instance x underwent the described random masking process with a probability of 50%: this introduces variability in the input data, encouraging the model to learn more robust features that are not overly dependent on the presence of specific elements within the feature vector.

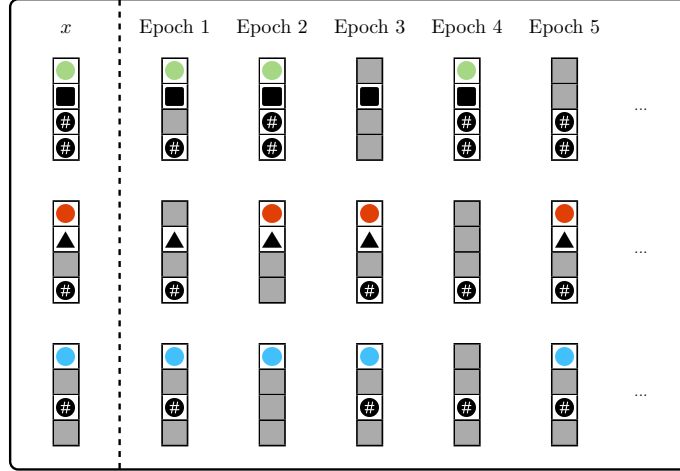


Figure 2.4: The proposed regularization strategy performed at every epoch before feeding the sample to the model. The colors (●, ●, ●) and the shapes (▲, ■) are examples of possible values for the categorical features, whilst # stands for a numerical feature. In the example are reported 3 feature vectors and their masked versions created in 5 different epochs. It should be noted that, when some features are originally missing, only the non-missing entries can be masked.

2.4 Experimental Configuration

In this section, we first report the data used in the experiments, the preprocessing applied and the metric used in the evaluation (Section 2.4.1), then we list the combinations of models and imputers used as competitors (Section 2.4.2). We provide some implementation details about the parameters used during training in Appendix B.

2.4.1 Data and Evaluation

We evaluate NAIM and competitor models on 5 publicly available datasets (reported in Table 2.1) from the UCI repository [15]. We selected these datasets as they are used to benchmark models for tabular data in [6]:

- **Adult** [16]: it is a dataset focusing on income prediction, which tries to predict people that make over $50k$ a year, from those that earn less than that amount, based on personal attributes such as education, occupation, and hours per week worked;
- **BankMarketing** [17]: its task is to predict whether the client will subscribe to a term deposit or not, and it comprises data from direct marketing campaigns of a Portuguese banking institution;

Dataset	# of samples	# of features	# of categorical features	Class distribution
Adult [16]	48842	14	8	0 : 37155 1 : 11687
BankMarketing [17]	41188	20	10	0 : 36548 1 : 4640
OnlineShoppers [18]	12330	17	7	0 : 10422 1 : 1908
SeismicBumps [19]	2584	18	4	0 : 2414 1 : 170
Spambase [20]	4601	57	0	0 : 2788 1 : 1813

Table 2.1: Datasets’ details and references. Datasets’ information consists of the number of samples, the number of features, how many of these features are categorical, and the classes’ distribution.

- **OnlineShoppers** [18]: it involves predicting whether a visitor will make a purchase based on session information, like page views and time spent on the site;
- **SeismicBumps** [19]: it aims at predicting seismic hazards in coal mines, i.e., whether any seismic bump with an energy higher than $10^4 J$ was registered in the next 8 hours, based on seismic activity and energy measurements;
- **Spambase** [20]: it involves identifying whether an email is spam or not based on word frequencies and other email attributes.

We note that this selection of datasets shares the same type of task, i.e., classification, and that all of them are complete, except for Adult, which has only 1% of missing data. The datasets exemplify well-curated data; however, in practical real-world scenarios, this level of completeness may not always be achievable.

For each dataset, we normalize the numerical features in the range $[0, 1]$ and apply one-hot encoding to the categorical features, before feeding them to the models. In the experiments featuring NAIM, HistGradientBoost, TabNet, TabTransformer, FTTransformer and XGBoost, we do not apply the one-hot encoding, since their implementations can handle categorical features. The preprocessing step is calibrated on the training data and then applied to validation and testing sets.

Our experiments focus on artificially generated *Missing Completely At Random* (MCAR) values, the most common evaluation regime used in missing data papers [12]. Indeed, MCAR’s lack of systematic bias means that the remaining data can be considered a fair representation of the whole. We aim to test our models under various missing data scenarios by introducing missing values at different percentages (denoted as p) across both training and testing sets. More specifically, we artificially generated missing percentages in the training and test sets separately, equal to 0%, 5%, 10%, 25%, 50% and 75% of the data in the set, for a total of 36 possible variations with repetitions. This allows us to methodically delineate varying scenarios of data completeness and missingness: *i*) a scenario with a well-curated

data collection for both training and testing sets, i.e., 0% of missing data in training and in testing; *ii*) a well-curated collection for the training set followed by a testing set with missing values, e.g., 0% of missing values in training and 25% in testing; *iii*), a scenario characterized by faulty data collection practices leading to missing values in the training phase, yet a completely curated test set without missing values, e.g., more common cases with 5%, 10%, or 25% of missing data in training and 0% in testing, or even more extreme cases where 50% or 75% of data are missing in the training set, but in the test set they are all not missing; *iv*) scenarios, which could represent domains where complete data collection is not possible, characterized by missing values in both training and testing sets, e.g., situations in which low percentages such as 5%, 10% or 25% of the data are missing in both sets, or even more extreme circumstances in which 50% or 75% of the values are missing. Given a targeted missing data percentage p and the total number of samples j in the considered set, we calculated the total number of values to be masked ($j \cdot n \cdot p$), considering any pre-existing missing values ($j \cdot n \cdot p - \sum_j (n - v_j)$). Following this, a random masking matrix of dimension $j \times n$ is generated to match the set structure. We ensured that at least one value in any fully masked row or column was replaced, avoiding complete data loss in any specific dimension. This approach leads to samples and features exhibiting varied percentages of missing values, all adhering to the MCAR paradigm.

2.4.2 Competitors

We conduct an extensive comparison of our methodology against approaches that incorporate missing data imputation as preprocessing followed by model training, and against those models able to intrinsically handle missing values. In this respect, our analysis exploits 11 unique model competitors, each integrated with the 3 main imputation techniques, i.e., the mean constant imputation, the KNN and the MICE with their default parameters [1]. This is summarized in Table 2.2, where the first column distinguishes between ML and DL approaches, the second lists the base learners, columns 3 – 5 report the imputation techniques, and finally last column shows the use of an intrinsic strategy. In this table, each of the 35 competitors is therefore marked by an “×”.

We can categorize the competing models based on the subsequent analyses we intend to conduct. Our objective is to benchmark NAIM against leading methodologies in the field of missing tabular data, specifically machine learning models combined with imputation techniques. Listed alphabetically, these models include:

- **AdaBoost** [21]: it is employed as a cascade of classifiers, which is particularly effective for its ability to adaptively focus on hard-to-classify instances, enhancing overall model

Type	Model	Imputers			Intrinsic
		Mean Constant	KNN	MICE	
Machine Learning	AdaBoost [21]	×	×	×	
	Decision Tree [22]	×	×	×	×
	HistGradientBoost [1]	×	×	×	×
	Random Forest [23]	×	×	×	×
	SVM [24]	×	×	×	
	XGBoost [25]	×	×	×	×
Deep Learning	MLP [26]	×	×	×	
	GRAPE [12]				×
	TabNet [5]	×	×	×	
	TabTransformer [6]	×	×	×	
	FTTransformer [7]	×	×	×	

Table 2.2: Combinations of models and imputers used as competitors in the experiments. Each competitor is represented by an “×”, given by the combination of a learner and an imputation technique, further to intrinsic strategies.

performance. This characteristic makes it a robust choice for handling diverse datasets, including those with imputed values.

- **Decision Tree** [22]: it is a versatile machine learning algorithm particularly useful for analyzing tabular datasets. It excels in its interpretability, as it visually represents decision-making processes and variable importance, enabling straightforward insights into complex data relationships and patterns.
- **HistGradientBoost** [1]: a variation of a gradient-boosted tree model, recognized for its efficiency in handling large datasets and its ability to improve upon the limitations of traditional gradient boosting by optimizing for speed and memory usage. Its inclusion allows for the assessment of advanced boosting techniques in missing data scenarios.
- **Random Forest** [23]: an ensemble of decision trees, it is noted for its robustness against overfitting. Leveraging a multitude of decision trees enhances the model’s accuracy and reliability, making it an exemplary model for evaluating the strengths of ensemble methods in the tabular data domain.
- **Support Vector Machine (SVM)** [24]: it is included for its versatility in dealing with non-linear data separations, offering, as a kernel machine equipped with the RBF kernel, a contrast to tree-based models in handling imputed datasets. The SVMs capacity to project data into higher dimensions where classes are more easily separable makes it a valuable model for comparison.

- **XGBoost** [25]: a sophisticated variation of AdaBoost that employs a gradient descent procedure to minimize the loss when adding weak learners, it stands out for its exceptional performance with tabular data. Its advanced handling of regularization and scalability positions XGBoost as the leading model for comparison in studies involving tabular datasets [2].

This selection of models, each with unique approaches for handling tabular data, provides a comprehensive background for evaluating the performance of our proposed methodology. Furthermore, considering that Decision Tree, HistGradientBoost, Random Forest and XGBoost offer implementations incorporating the MIA strategy, we assessed their effectiveness in managing missing values.

Lastly, we aim to showcase the performance of our method in comparison to other DL models specifically designed for tabular data. Our comparative analysis extends to 5 advanced architectures, each selected for its novel approach to handling tabular data, particularly when combined with imputation techniques for managing missing values. These models, discussed in sections 2.2.1 and 2.2.2, bring unique perspectives and methodologies to the challenges of tabular data analysis:

- **Multilayer Perceptron (MLP)** [26]: a foundational DL model, it stands out for its simplicity and versatility. It consists of multiple layers of neurons, each fully connected to those in the next layer, enabling the model to capture complex nonlinear relationships between features. Despite its straightforward structure, MLP’s performance hinges on the quality of input data, making it crucial to pair with effective imputation techniques to address missing values.
- **GRAPE** [12]: designed explicitly to tackle missing data within tabular datasets, this model employs a graph-based approach, representing each data instance as a node and modeling feature interactions as edges. Leveraging GNNs, it propagates information effectively across the nodes, allowing the model to learn directly from incomplete data without preliminary imputation. This structure enables GRAPE to capture complex relationships inherent to tabular data, enhancing robustness and predictive performance even in the presence of extensive missingness.
- **TabNet** [5]: leveraging the self-attention mechanism, it dynamically selects which features to focus on for making decisions. Its ability to perform dynamic feature selection mimics the boosting technique, enhancing its interpretability. This model is particularly suited to interpreting high-dimensional data, offering insights into feature importance and decision pathways.

- **TabTransformer** [6]: this model innovates by embedding categorical features within tabular data, applying transformer-based self-attention mechanisms to create contextual embeddings. It captures complex inter-feature relationships and dependencies, significantly improving the predictive performance on tabular datasets.
- **FTTransformer** [7]: it further explores the potential of transformers in extracting patterns and interactions within tabular data, employing distinct embedding strategies for numerical and categorical features.

These DL models, chosen for their innovative handling of tabular data with respect to classical ML approaches, are integrated with 3 main imputation techniques since none of them can handle missing values for both categorical and numerical features.

2.4.3 Evaluation Metrics and Statistical Analysis

For evaluation, each dataset is divided into 5 stratified cross-validation splits, to maintain the original class distribution, and for each fold, 20% of the training samples are selected for validation. We evaluate each experiment averaging out the values of Area Under the ROC curve (AUC) computed in the different cross-validation folds. AUC is a valuable metric to evaluate classification tasks, since it represents the degree to which a model can correctly classify positive and negative instances across all possible thresholds, making it a comprehensive measure of model performance, even in case of imbalanced data [27].

To ensure the statistical validity of our results, we investigated the statistical differences between the predictions made by NAIM and those by the competitors using the Wilcoxon signed-rank test. To this end, we assessed the rate of experiments for each dataset individually, among the 36 variations with repetitions of the 6 percentages of missing values artificially generated, in which NAIM achieved statistically superior and inferior results compared to the competitors, setting $p = 0.05$.

2.5 Results and Discussions

As outlined in the previous section, we compare NAIM with 35 leading competitors in ML and DL for tabular data, and we test the performance with 36 different combinations of percentages of missing values both in the training and testing sets, across the various datasets, resulting in a total of 6480 experiments, 1296 per dataset. Therefore, in analyzing the performance of the different approaches, we considered each level of missing data within the training set separately, delineating performance metrics as the percentages of missing

data in the testing set increased. In Table C.1 we reported for each combination of model and strategy to handle missing values (reported in different rows), the average performance, in terms of AUC and standard error, across the 5 datasets under consideration, obtained at the different percentages of missing values (reported in the different columns). In particular, the first rows denote the percentage of missing values used in training and the specific percentage of missing data in the test set. Similarly, in Appendix C, we also reported the detailed tables per dataset (Tables C.2 - C.6). As we can see in Table C.1, NAIM achieves the best performance, highlighted in bold, in most scenarios (23 out of 36). To be more specific, we conducted the statistical analysis described in section 2.4.3 and reported in D, to compare the predictions made by NAIM in comparison with those of its competitors. This analysis shows that on average the proposed model achieves better performance in 58.7% of the cases, while it only loses in 1.6% of the cases.

To further analyze the results we divide the competitors into groups and plot the performance constructing 6 separate charts, one for each level of missing data within the training set, which report the performance metrics as the percentage of missing data in the test set increases. We use this representation to compare NAIM with different groups of competitors, namely ML and DL models paired with the imputation techniques (Figure 2.5), model implementing some intrinsic strategy (Figure 2.6) and the various imputation techniques regardless of the model used (Figure 2.7). To present our findings, these charts report the average performance values.

As a first analysis, we compare NAIM with ML and DL models coupled with the 3 different imputers under consideration. In Figure 2.5, in panels A and B for the ML and DL models respectively, we reported the average performance across all 5 datasets and the 3 imputation techniques under consideration. These charts show the decreasing trend in model’s performance as the percentage of missing values in the testing set increases, a finding that aligns with our expectations. Interestingly, NAIM stands out for its consistently superior average performance across all levels of missing data, maintaining its high ranking even in the optimal scenario where no data is missing. This observation not only underscores the distinct advantage of NAIM over ML and DL models but also highlights the unexplored potential of DL methods in improving both the model’s performance when handling missing data. Furthermore, the performance difference is particularly evident when the training set has 0% missing data (first chart on the left in the first row for each panel A and B). This highlights a significant challenge that current methodologies face: the need to learn how to handle missing data during the training phase, to then correctly infer on testing data with missing values. This aspect emphasizes a critical limitation of conventional approaches and the need for innovative strategies that can effectively address data incompleteness. In this

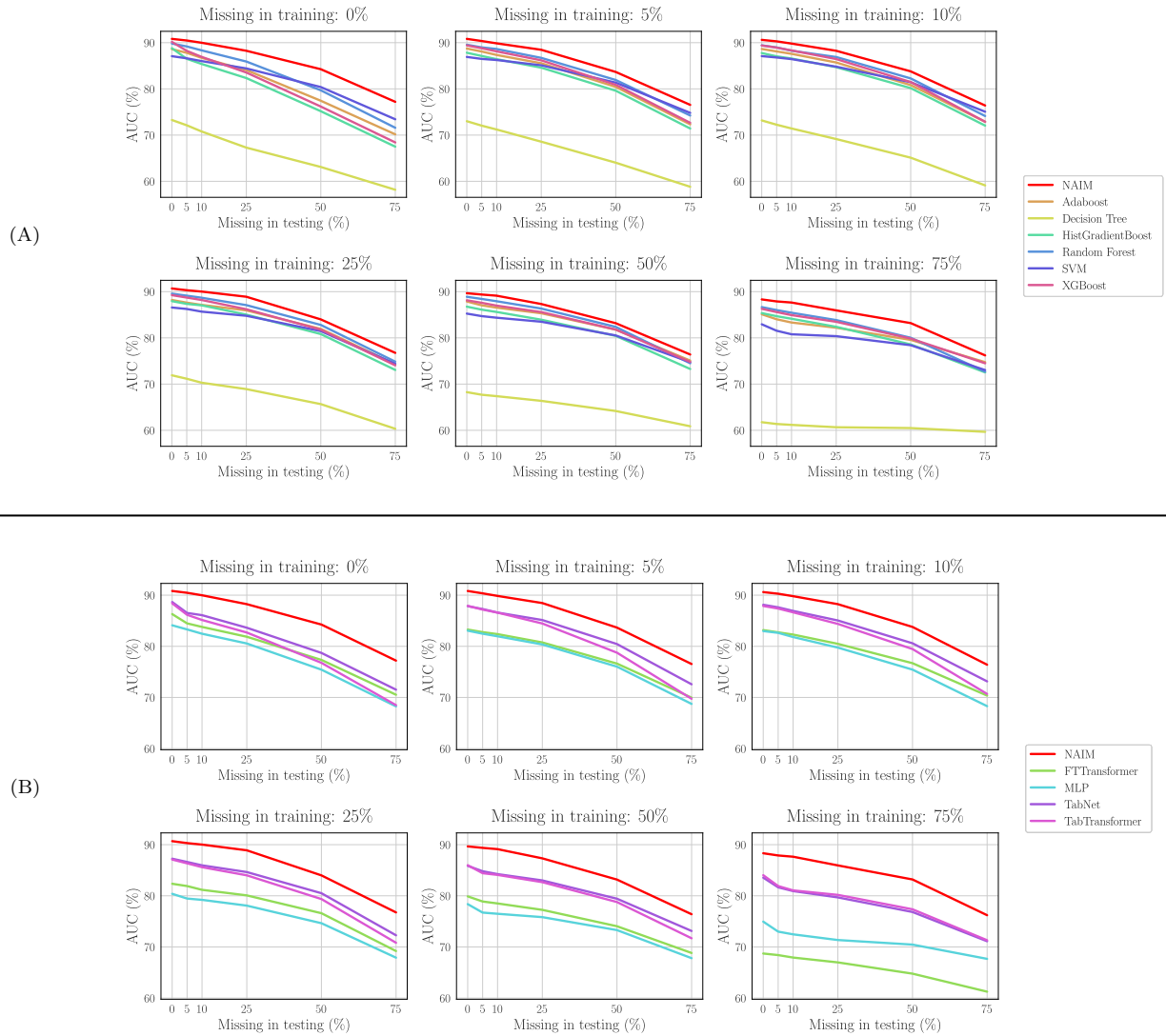


Figure 2.5: Comparison between the test performance of NAIM and the competitor models averaged over the 5 dataset and the 3 imputation strategies: panel A reports the comparison between NAIM and the ML models, whereas panel B compares NAIM with the DL models.

respect, we deem that the better performance of NAIM is due to our regularization technique that enables it to learn how to handle missing values even if all the training data are present.

Then, as a second analysis, we compare NAIM against the models capable of intrinsically handling missing values. In Figure 2.6, which reports the performance of these approaches, the average performances are computed across the 5 datasets. Turning our focus on the top left chart, which plots the performance in the 0% missing values in training scenario, we notice that the difference in performance as the percentage of missing values in the test set increases is higher compared to the other 5 charts where the training is performed on

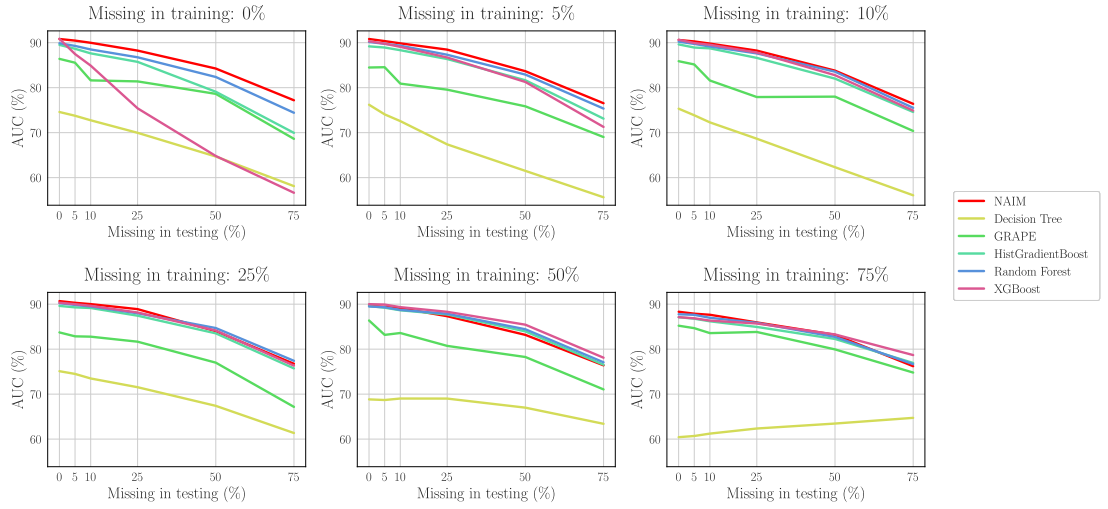


Figure 2.6: Comparison between the test performance of NAIM and the competitors implementing an intrinsic strategy averaged over the 5 dataset.

data containing missing values. This performance drop confirms that neither the imputation techniques nor the intrinsic strategies are effective if no missing values have been seen during training, whereas the stable performances of NAIM validate our novel regularization technique, which allows it to adapt to the presence of missing data in all training scenarios. These charts also show the competitive performance obtained by XGBoost and Random Forest, each delivering comparable performance across a spectrum of scenarios characterized by missing data in the training set, including the optimal scenario of 0% missing in both sets. It is worth noting that in the scenario with 50% missing data in training (second row, second chart from the left), they are even able to outperform NAIM regardless of the percentage of missing values in the test set. Moreover, it is noteworthy that HistGradientBoost achieves performance similar to those of XGBoost and Random Forest, although slightly inferior. On the opposite, Decision Tree and GRAPE, despite being theoretically capable of handling missing data, consistently fail to achieve comparable performance to other models in both missing data scenarios during training and testing.

Subsequently, as a third analysis, to further assess NAIM’s ability to handle missing data, we compare it against imputation strategies. In contrast to the previous analyses, here the average performance for each of the imputation methods, reported in Figure 2.7, are computed by averaging the results across the 5 datasets and the 10 models available. The figure illustrates how, in scenarios involving complete training data (first row, first chart from the left), models trained with data imputed using the KNN imputer exhibit slightly

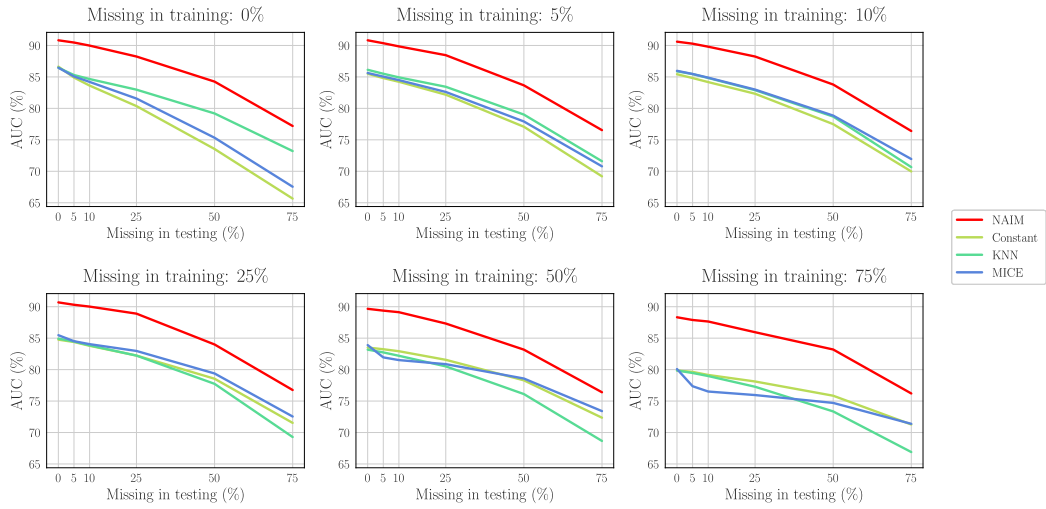


Figure 2.7: Comparison between the test performance of NAIM and the competitor imputers averaged over the 5 dataset and the 10 models available.

superior performance compared to other imputation methods. However, it is important to note that, overall, there are no significant and clear differences among the various imputation techniques, thus requiring extensive experimentation to find the optimal configuration for the data and task at hand.

Upon observing the superior performance of NAIM across all missing data scenarios, we examine its performance robustness to varying levels of missing data (Figure 2.8). To this end, we measure the performance drop presented in experiments with missing data either in the training or in the test set compared to the optimal scenario of 0% of missing data in training and test sets. Subsequently, we average these values to present them in a single plot. To quantify the robustness of the models when dealing with increasing missing values in testing, we evaluate the performance differences across scenarios with complete training data and various missing data percentages in the test sets. In other words, to compute the y coordinates of Figure 2.8, we measure the average performance drop in the experiments reported in the first chart of the first row of panels A and B of Figure 2.5 and in Figure 2.6. Conversely, to assess the robustness of the models when dealing with increasing missing values in training, we examine the performance differences in scenarios with various missing data percentages in training and complete test data. Therefore, we measure the differences among the first performance reported in each of the 6 charts shown in panels A and B of Figure 2.5 and in Figure 2.6 and report their average as x coordinates in Figure 2.8. To facilitate interpretation and effectively rank models based on their robustness, we also plot

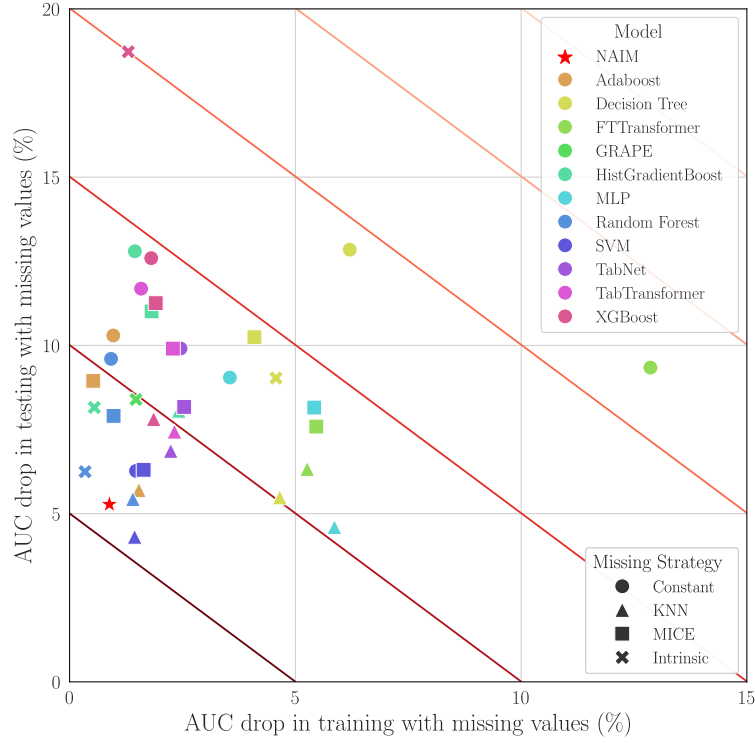


Figure 2.8: Analysis of the robustness of the models to increasing levels of missing data. Each experiment (combination of model and missing strategy) is positioned based on the average drop in performance obtained in scenarios with increasing percentages of missing values in the training and testing sets, separately.

equi-drop lines connecting data points that exhibit equivalent overall performance reductions from the no-missing data state (0% of missing values in both sets). For instance, an experiment exhibiting an average 10% drop in the presence of missing data in the test set and 0% drop in the training set, would be ranked equivalently to a model with the opposite behavior of 0% drop in the presence of missing data in the training set and 10% drop in the set set. Straightforwardly, a model completely robust to missing data in both training and testing sets would be placed on the point (0, 0). We note that NAIM’s robustness to missing data is indicative of minimal performance degradation with an increasing number of missing values, since it presents a 0.88% drop when the rate of missing values increases in the training set and a 5.27% drop when the rate of missing values increases in the test set. Notably, despite this analysis positions NAIM just behind the SVM model paired with the KNN imputer, when considering the comprehensive performance evaluations presented in panel A of Figure 2.5 and in Table C.1, NAIM significantly surpasses the competitor. Furthermore, looking

at Table D.1, we can see that NAIM performance is statistically larger in 67.8% of the cases, whereas it is never inferior. This means that SVM is a more robust model but with lower performance compared to NAIM, confirming that our approach presents a superior capacity to maintain high-performance levels despite the presence of missing data. Moreover, from Figure 2.8 we can see that the third model in terms of robustness in the presence of missing values is the Random Forest that exploits the MIA strategy, which also performs statistically worse than NAIM as depicted in Table D.1. These observations confirm the robustness of the proposed model and regularization technique in handling missing values.

As a further analysis to evaluate the contributions of the regularization technique and the masking mechanism, we conducted ablation studies in which the model was trained without regularization, both in scenarios with missing values (NAIM w/o reg) and when combined with imputation techniques (NAIM w/o reg + Imputer). Figure 2.9 shows the average performance of these experiments conducted across the five datasets, while a comprehensive presentation of the results is provided in Tables C.1C.6 in Appendix C. Notably, the contribution of the regularization technique is evident throughout the 6 plots. Indeed, in the first plot (top left), NAIM w/o reg exhibits a performance trend comparable to that of XGBoost in Figure 2.6, suggesting that exposure to missing data during training is essential for the model to learn effectively how to handle them. Conversely, in the final plot (bottom right), where the missing rate is highest, the benefit of regularization becomes less pronounced. When evaluating NAIM combined with the three imputation strategies, it is evident that the masking self-attention mechanism consistently outperforms the imputation-based baselines across all epochs. These findings indicate that each component of our approach, particularly the missing data regularization, contributes significantly to the overall performance of NAIM. Finally, to complement the performance analysis, we also conducted a computational assessment of the DL models considered in this study. Specifically, we report in E the number of floating-point operations (FLOPs) and trainable parameters for each model, highlighting the balance achieved by NAIM between representational power and computational efficiency.

2.6 Conclusion

In this work, we introduced NAIM, an innovative transformer-based architecture for modeling tabular data in supervised learning environments, specifically designed to handle missing values. This architecture simplifies the analysis process, essentially avoiding the need for traditional imputation strategies. Our empirical analysis demonstrates NAIM’s superiority over existing state-of-the-art solutions in handling missing data within tabular datasets. Our experiments spanned across NAIM and a combination of 11 ML and DL models, each inte-

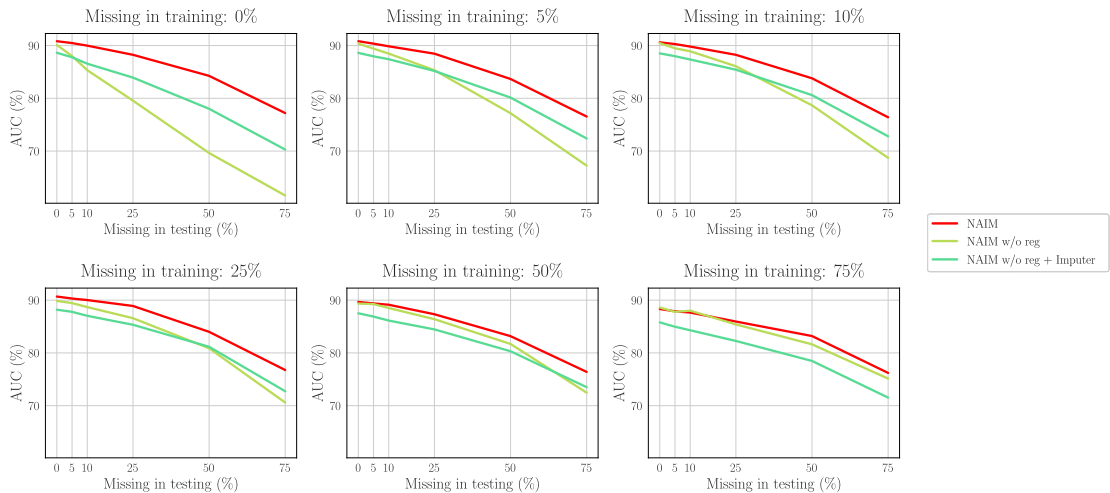


Figure 2.9: Ablation analysis evaluating the impact of the regularization technique (NAIM w/o reg) and the masked self-attention mechanism (NAIM w/o reg + Imputer). The plots report test performance averaged over the 5 datasets. For NAIM w/o reg + Imputer, results are additionally averaged across the 3 different imputers.

grated with 3 distinct imputation methods, and also against those 5 models that intrinsically handle missing values. We performed the analyses on 5 publicly available datasets, which differ in the number of samples and features.

NAIM leverages specialized embedding mechanisms for both categorical and numerical features, coupled with an innovative self-attention mechanism, to maximize the information of the available data at hand. Additionally, we introduced a novel regularization strategy designed to address a significant challenge identified in our research: the need for state-of-the-art approaches to learn how to handle missing values in the training process. This strategy, which involves random masking of each sample in every epoch, equips NAIM with the robust capability to handle missing data under all circumstances, including scenarios where there are no missing values in the training set.

While we have attained promising results, this work has some limitations that are now discussed introducing the related future directions:

- **Extension to MAR and MNAR mechanisms:** In this work we focused on the MCAR setting, which is the most commonly adopted evaluation regime in the literature and allowed us to systematically explore a wide range of scenarios. Nevertheless, real-world applications often involve more complex missingness mechanisms, such as Missing At Random (MAR) or Missing Not At Random (MNAR), where the probability of missingness depends on observed or unobserved variables. Future work will therefore extend

the evaluation of NAIM to MAR and MNAR conditions, in order to further validate its robustness and generalizability in more realistic settings.

- **Multimodal learning with missing modalities:** We plan to extend NAIM beyond the tabular domain by developing a multimodal variant capable of handling heterogeneous data types and the complete absence of entire modalities (e.g., missing imaging or tabular information). This extension would enable a unified framework that leverages available data across different domains while degrading as little as possible in performance when some sources are absent.
- **Efficient attention mechanisms:** Although transformers are powerful, they are known for their high computational cost. Future investigations should explore alternative attention formulations that can reduce memory and computational requirements. This would make NAIM more scalable, particularly for large datasets with high-dimensional features.
- **Model interpretability and transparency:** To foster trust and adoption in critical domains like healthcare and finance, future studies should deepen the interpretability of NAIM. This includes leveraging attention weights to understand feature importance, analyzing learned embeddings to uncover latent data structure [28], and developing attribution methods tailored to transformers with missing inputs. These efforts would help demystify the internal reasoning process of the model, especially in the presence of incomplete data.
- **Self-supervised and semi-supervised learning:** Since labeled tabular data is often scarce, particularly in specialized domains, future works should investigate self-supervised pre-training strategies for NAIM. Techniques such as contrastive learning, masked feature modeling, or denoising objectives could help the model learn robust representations from raw data, improving downstream performance when only limited supervision is available.
- **Temporal modeling for incomplete time-series:** Another natural extension of our framework involves adapting NAIM to handle longitudinal or sequential data. Time-series datasets, especially in domains like sensor monitoring or electronic health records, frequently exhibit irregular sampling and missingness. Future work should aim to develop a time-aware version of NAIM that integrates temporal encodings.
- **Domain-specific evaluation and adaptation:** While our evaluation included diverse public datasets, future work should assess NAIM in domain-specific real-world appli-

cations. Particular attention should be given to healthcare datasets, where missingness is often informative and non-random. This would not only validate the models utility in complex scenarios but also drive the development of domain-adaptive mechanisms that incorporate domain knowledge into the learning process.

These potential directions not only underscore the versatility and expansiveness of our approach but also highlight the fertile ground for future advancements in the domain of tabular data analysis. In conclusion, when faced with tabular datasets exhibiting missing data either in the training or testing phase, a promising and competitive choice among the tested methods is NAIM, as it enables the extraction of useful information from the available data without any need for imputation.

In this chapter, we presented NAIM, a transformer-based model specifically designed to handle missingness in tabular data by integrating a masked attention mechanism and a missingness-aware regularization strategy. The proposed approach demonstrated strong performance and robustness across various benchmark datasets, outperforming existing methods in most scenarios. However, the experiments were conducted on synthetic or structured datasets and focused solely on tabular data. These limitations hinder the clinical transferability of the model, as real-world medical datasets often include complex temporal dynamics, heterogeneous patient populations, and partially observed clinical histories. To address this gap, the next chapter focuses on applying machine learning models to real clinical data for survival prediction in cancer patients, emphasizing the interpretability and stability of predictions in the presence of naturally occurring missing information.

External evidence. In a subsequent study [29], NAIM was evaluated in the Pembreal 5Y registry, a global real-world cohort of 1050 patients with PD-L1 \geq 50% treated with first-line pembrolizumab across 61 institutions in 14 countries. NAIM handled heterogeneous missingness without imputation and supported time-dependent risk modelling; reported performance included a c-index around 0.63 (risk of death) and an $AUC \simeq 0.61$ for 5-year survival, accompanied by SHAP-based explanations highlighting evolving prognostic factors. This multicentre use case substantiates NAIM’s clinical robustness while also showing headroom for further optimisation and calibration.

Chapter 3

Survival Prediction from Incomplete Clinical Data

Building upon the findings of the previous chapter, we now shift our attention from generic benchmark datasets to real-world clinical scenarios. In particular, we explore the prediction of overall survival in oncology patients using routinely collected clinical data, where missingness naturally arises due to irregular follow-up schedules and incomplete records. Unlike the controlled settings of Chapter 2, this application highlights the practical challenges of modeling clinical features in the presence of partial observations. Furthermore, the emphasis moves beyond pure performance to include interpretability, a crucial aspect in clinical decision-making. This chapter thus represents a key step towards translational applications of robust modeling strategies in healthcare settings.

3.1 Introduction

In recent years, Artificial Intelligence (AI) has rapidly become a part of our daily lives, including healthcare [30]. There have been notable advances in applying quantitative methods in clinical practice, which have paved the way for precision medicine. Cancer research is one of the most promising areas where AI can be applied [31]. According to the World Health Organization, cancer is responsible for nearly 10 million yearly deaths globally, with lung cancer accounting for 18% of these [32] and non-small cell lung cancer (NSCLC) is the most frequent type, being approximately 82% of all cases [33].

In cancer research, a key factor is the overall survival (OS) of the patient, which refers to the time from the initial cancer diagnosis to the time of death. Identifying subgroups of patients with a higher or lower chance of survival is critical in developing effective strategies

to improve OS rates. For example, in the case of NSCLC, the 5-year survival rate is only 26%, and this rate drops to a mere 7% when cancer returns locally or spreads to distant organs [33]. Therefore, although some AI-based methods have been proposed in healthcare [34], accurately predicting OS remains a major challenge.

In medicine, it can be quite challenging to obtain a dataset, i.e., tabular dataset, with no missing features: how to handle this situation effectively is a key question in AI research. In fact, traditional methods require complete data samples, so standard practices involve either excluding samples with missing values or utilizing imputation strategies. However, it is important to note that these methods can compromise the findings, introducing bias and reducing statistical power. Therefore, research should go towards the development of models which can cope with missing features without any imputation.

To address these issues, we introduce a novel approach that tackles the problem of handling missing features in tabular data in the study of OS analysis for patients affected by NSCLC. Indeed, tabular data are the most simple and diffused data type in AI applications [35]. Our method stems from the transformer architecture [13] and leverages the idea of the mask inside the self-attention module to learn from incomplete input data. By masking missing features, it learns only from available data, avoiding any type of imputation of missing features, with the goal of improving the performance of survival prediction of patients affected by NSCLC.

To validate the proposed model, we compared its predictive performance with feature imputation and OS prediction state-of-the-art approaches. Indeed, even though some methods are capable of handling missing features in the test phase, to the best of our knowledge, none has yet addressed the issue of how to deal with them during training without the need for any imputation method in the context of OS prediction [34].

The main contributions of this work are the following:

- We propose a specialized transformer-based decision support system on tabular data, for predicting lung cancer OS using clinical data.
- Our approach efficiently handles missing data, eliminating the need for imputation strategies.
- Our approach uses ad-hoc designed loss functions allowing the use of both censored and uncensored patients, without the need to exclude any patient.
- We assessed the performance of the proposed model at various time granularities, confirming its robustness across different time windows. The results obtained on the

OS prediction in NSCLC patients outperform OS state-of-the-art models, regardless of the imputation method used.

The manuscript is organized as follows: section 3.2 presents the state-of-the-art of survival analysis and data imputation techniques; section 3.3 reports the details about the data employed in the analyses; section 3.4 introduces the proposed model and explains the metrics employed; section 3.5 discusses the experimental results; section 3.6 provides concluding remarks.

3.2 Background

Survival analysis, also known as time-to-event analysis, plays a crucial role in various fields, especially medicine [36]. It aims to understand the relationship between covariates, such as patient features, and the distribution of survival times. In the lung cancer setting, this type of analysis helps to identify risk factors that affect survival and to compare risks among different subjects, with the aim of tailoring the right therapy for the right patient".

One important aspect to consider in this particular field is censored data, specifically right-censoring. This occurs when a patient withdraws from the study, is lost to follow-up, or is still alive without experiencing the event of interest at the last follow-up. In such cases, a special analysis is required since it is not known when and if the event occurred, and thus these samples cannot be considered together with uncensored data, i.e., patients for which it is known when they have experienced the event. However, most of the studies that predict OS in NSCLC have approached the survival problem as a classification task, dividing patients into high and low risks based on a survival time threshold [37, 38, 39, 40, 41, 42]. This approach does not fully exploit the information of the censored patients, since it excludes from the analysis those patients with a survival time shorter than the threshold. In contrast, other studies predict the risks faced by patients, which can be used to evaluate the correct ordering of the patients through the *C-index* metric [43, 44, 45, 46, 47, 48]. The downside of these approaches is that they do not make full use of the available temporal information, not taking into account possible changes in risk over time. On the contrary, another possible approach is to use several output nodes for each time interval, which represent the risk the patient experienced at the specific time interval. Indeed, in [49] the authors present ad-hoc designed loss functions to exploit information from both uncensored and censored patients, enabling learning from the latter that the event of interest did not occur up until the last follow-up. This was achieved by defining the survival problem as identifying the first time an underlying stochastic process hits a specific boundary, also known as the first hitting time.

Although this approach still does not enable predictions to be made about censored samples, it does permit learning from them.

To tackle the survival task, a few models in the literature have been proposed. The most standard used approach is the Kaplan-Meier estimator, which can take censored samples into account, but does not incorporate patient covariates, making it useful to estimate the survival rate at the population level but not at the patient level. Instead, the most commonly used method is the Cox proportional hazard (CPH), which incorporates the features of the patient, but assumes that the hazard rate, i.e., the probability of experiencing the event within a short time interval, is constant and that the log of the hazard rate is a linear function of the covariates. These two assumptions are known as the proportional hazard assumption, and a few methods have been proposed to address the limitations it introduces. The Survival Tree (ST) and the Random Survival Forest (RSF) are extensions of the CART and Random Forest algorithms that deal with censored survival data. ST accomplished this by a recursive partitioning procedure based on maximizing the dissimilarity in the survival distributions of patients between different regions of the covariate space, thus taking into account not only the homogeneity of predictor variables but also the survival probability and time of the events. RSF estimates the cumulative hazard function for each case's terminal nodes by leveraging out-of-bag data and then averages the cumulative hazard functions of all trees in the forest. The cumulative hazard function is a measure that describes the cumulative probability of experiencing a particular event, such as death or failure, up to a specified time point. Instead, DeepHit (DH) [49] offers a new approach that predicts the probability of the first hitting time of an event using a deep neural network, employing a loss function that exploits survival times and relative risks to predict the cumulative incidence function, a measure useful in understanding the risk of developing a particular health outcome, as it takes into account both the risk of developing the outcome and the time period over which the risk is evaluated.

When working with medical data, an important aspect to consider is the presence of incomplete records. This phenomenon may be due to different reasons and can hinder the employment of AI models. Although some strategies have been proposed to address this issue, only a few methods are commonly used in a healthcare setting [34]: complete case analysis, overall mean imputation, k-nearest neighbors (kNN) imputation, multiple imputations by chained equations (MICE) and MissForest. The first two methods for dealing with missing data, complete case analysis, and overall mean imputation, involve either discarding incomplete samples or imputing the missing values with the mean value of the corresponding feature computed from the available data, which are straightforward but have some drawbacks. Specifically, the former may discard too many samples, which can prevent the learning

of deep learning techniques, whereas the latter has a high likelihood of introducing bias into the model’s final outcome. As a result, kNN, MICE and MissForest have been proposed as alternative solutions. The first method identifies similar samples to the one being analyzed by calculating the distance based on other non-missing features in order to impute a more appropriate value. The second one imputes missing values by modeling each feature with missing values as a function of other features in a round-robin fashion. The third one uses a Random Forest model to predict missing data, which starts by imputing missing values using a simple strategy, e.g., mean and mode imputation, and then iteratively updates the imputations using the model predictions, refining these estimates until the changes between iterations are minimal.

In this context, we present our approach, which combines the strengths of DH and transformer architecture [13]. We use the transformer architecture to exploit the power of the mask in the self-attention module, which allows us to avoid any imputation of missing values in the survival analysis task.

3.3 Materials

To validate our clinical decision support system, presented in the next section, we used clinical data from the CLARO dataset [50]. This consists of 297 patients affected by NSCLC, who underwent concurrent chemoradiation for locally advanced NSCLC and systemic treatment for metastatic disease. The OS of the entire population, which included 184 censored and 113 uncensored patients, has a mean of 20.74 ± 42.45 months (95% CI). The population was enrolled in the study under two separate Ethical Committee approvals, including a retrospective phase that was approved on October 30th, 2012, and registered on ClinicalTrials.gov on July 12th, 2018 with the identifier NCT03583723. The prospective phase was approved with the identifier 16/19 OSS. The Institutional Review Board approved this analysis, and all patients provided written informed consent.

The clinical data that we have collected contains 8 clinical descriptors, as outlined in Table 3.1, which reports the distribution and the amount of missing values for each feature. As mentioned above, our features are composed of personal information, i.e., *Age* and *Sex* (assigned at birth), and details about tumor histopathology, i.e., clinical target volume (*CTV*), *Overall Stage*, tumor (*T*), nodule (*N*), metastasis (*M*) stages and *Histology*. To determine the stage of the tumor, we had two radiation oncologists with extensive experience independently review CT scans and assign staging scores for the tumor: *Overall Stage*, *T*, *N*, *M*. If there was any disagreement between the two experts, they would review the patient’s CT images together until they reached a consensus. It is worth noting that some patients in

our study did not undergo a histopathological examination, this is why the unknown" class is included as one of the categories for the *Histology* feature.

To prepare the data for analysis, we applied one-hot encoding for categorical features and z-score normalization for continuous features¹. For our approach, we generated empty vectors for the one-hot encoding of the missing features.

Feature	Missing Data	Categories	Distribution
<i>Age</i> *	0 (0.0%)	< 68 <i>years</i>	138 (46.46%)
		≥ 68 <i>years</i>	159 (53.54%)
<i>Sex</i>	0 (0.0%)	F	97 (32.66%)
		M	200 (67.34%)
<i>CTV</i> *	112 (37.71%)	< 146.51 <i>cm</i> ³	115 (38.72%)
		≥ 146.51 <i>cm</i> ³	70 (23.57%)
<i>Overall Stage</i>	0 (0.0%)	II	8 (2.69%)
		III	188 (63.31%)
		IV	96 (32.32%)
		Recurrence	5 (1.68%)
<i>T</i>	116 (39.06%)	1	11 (3.70%)
		2	40 (13.46%)
		3	74 (24.92%)
		4	56 (18.86%)
<i>N</i>	104 (35.02%)	0	17 (5.72%)
		1	39 (13.13%)
		2	112 (37.71%)
		3	18 (6.06%)
		Recurrence	7 (2.36%)
<i>M</i>	94 (31.65%)	0	201 (67.68%)
		1	2 (0.67%)
<i>Histology</i>	3 (1.01%)	Adenocarcinoma	153 (51.52%)
		Squamous	73 (24.58%)
		Other	20 (6.73%)
		Unknown	48 (16.16%)

Table 3.1: Patients' characteristics. For each feature is reported the number and percentage of missing values, along with the distribution (number and percentage) of the possible categorical values. Note that the variables marked with *, are continuous, but, for the sake of clarity, we have presented their distribution using their mean values as thresholds. It is important to note that the model used the original continuous values of these variables.

¹This operation is performed on the test set based on the parameters computed using the training set.

3.4 Methods

In this section, we first describe the architecture of the proposed approach with an example of its main blocks, next we illustrate the loss function used during training and the evaluation metric employed, and finally we report the experimental setup used to perform the experiments.

3.4.1 Model

Given that AI methods for OS prediction require complete data [34], and that in healthcare it can be quite challenging to obtain a complete data set without missing values, we need to effectively handle the data without the necessity of imputing or removing any information. Taking inspiration from the cutting-edge transformer architecture [13], we present here a novel model that takes into account only the available features. Our approach involves adapting the transformer’s encoder architecture to tabular data, via a novel positional encoding for tabular features, and utilizing padding to mask any missing features within the attention module, enabling the model to ignore them effectively.

The schematic representation of the proposed model is shown in Figure 3.1: left panel, denoted by the letter A, represents the overall architecture from the input (the orange block shown at the bottom of the figure) to the output (the blue block at the top of the figure), whilst right panel, marked as B, offers an example of four blocks of panel A, which can be identified by the name as well as by the color.

The input to the model, shown in the first block represented in orange in Figure 3.1, is a feature vector \mathbf{x} of dimension d composed of the preprocessed patient information. The second block, represented in yellow, is the positional encoding, which is used for identifying the feature itself without explicitly encoding the feature order. To achieve this, we explored the use of a one-hot encoding vector representing the position of each feature. In this respect panel B of Figure 3.1 shows an example in the case of the CLARO dataset: here the initial feature vector is turned into a $d \times (d + 1)$ matrix, where the columns from the first to the second-last represent the positional encoding, whilst the last column reports the feature vector.

The third block is the transformer encoder (light gray block), which consists of a stack of M identical layers. Each layer has two sub-layers: the first is a multi-head self-attention mechanism (red block), and the second is a position-wise fully connected feed-forward network (green block). As in the original architecture [13], we employ a layer normalization (purple block) and a residual connection (denoted with \oplus) around each of the two sub-layers.

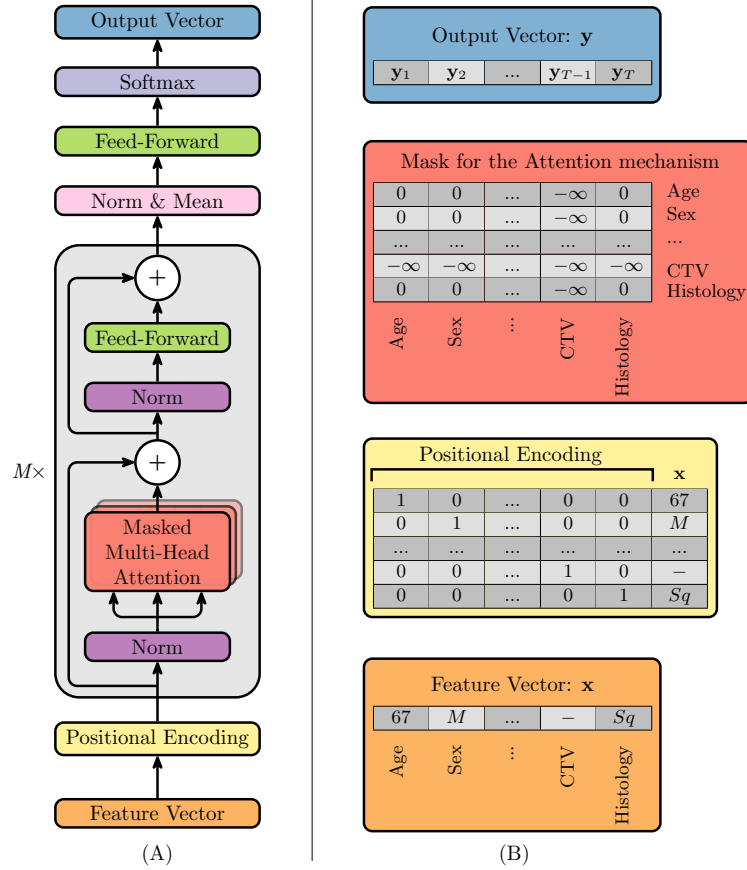


Figure 3.1: Schematic representation of the proposed model: (A) Architecture of the proposed approach and (B) example of positional encoding, mask and output vector, where the $-$ symbol represents a missing feature. Note that for simplicity of representation, just a few features are reported and none of the preprocessing procedures are applied.

The normalization step precedes each of these blocks to reduce the risk of the vanishing or exploding gradients problem, which can occur during the training process.

To handle the presence of missing features, we drew inspiration from the padding mask technique used in natural language processing [13]. This technique extends the capabilities of the attention mechanism, preventing leftward information flow, which exploits the use of the mask within the scaled dot-product attention mechanism to mask out, setting to $-\infty$, any values that would result in illegal connections. In this regard, the red block of panel B of Figure 3.1 shows an example of the mask employed in the attention mechanism. Using such a mask, we were able to effectively ignore any missing features, without requiring any imputation strategy.

Next, we applied a normalization and averaging step (pink block) in order to get a latent representation vector, respective to the non-missing features, to be fed to the final classifier, composed of a feed-forward module (green block) and softmax function (violet block). The

feed-forward module that succeeds the encoder maps the encoder’s embedding to the output vector’s dimension T . Hence, each element \mathbf{y}_t of the output vector \mathbf{y} (blue block) represents the probability that the event occurs at the time point t .

3.4.2 Training and Testing

When training the model, the goal is to correctly predict the probability that the event of interest, formally denoted as $k = 1$ and corresponding to the patient’s death, occurs at time t , given the patient’s feature vector \mathbf{x} and under the constraint that t is smaller or equal to s , the true time when the event $k = 1$ occurs. Straightforwardly, we identify with $k = 0$ the event of censoring for patients who did not experience the event. The cumulative incidence function $F(t|\mathbf{x})$ is therefore defined as:

$$F(t|\mathbf{x}) = P(s \leq t, k = 1|\mathbf{x}) = \sum_{s=0}^t P(s, k = 1|\mathbf{x})$$

However, since the true cumulative incidence function $F(t|\mathbf{x})$ is not known, it can be estimated with $\hat{F}(s|\mathbf{x})$, which is the cumulative sum of the outputs of the model, defined as:

$$\hat{F}(s|\mathbf{x}) = \sum_{t=0}^s \mathbf{y}_t.$$

Therefore we can now compute this measure and thus compare the risks faced by different patients at a specific time.

We train the proposed model using the loss function L presented in [49], specifically designed to handle censored patients. It is composed of two terms, L_1 and L_2 , so that

$$L = L_1 + L_2$$

where the first term captures the death and censoring times of patients, evaluating respectively the output \mathbf{y} and the estimated cumulative incidence function \hat{F} at patient-specific event times, whereas the second term measures the correct ordering of the patients based on their relative risk. Indeed, L_1 is the log-likelihood of the distribution of the first hitting time of an event, defined as:

$$\begin{aligned} L_1 = & - \sum_{i=1}^N [\mathbb{1}(k^{(i)} = 1) \cdot \log(\mathbf{y}_{s^{(i)}}^{(i)}) + \\ & + \mathbb{1}(k^{(i)} = 0) \cdot \log(1 - \hat{F}(s^{(i)}|\mathbf{x}^{(i)}))] \end{aligned}$$

where $\mathbb{1}(\cdot)$ is an indicator function, and the summation is performed on all N samples, where the apex (i) indicates that the information is related to the specific patient i . Through two terms, it exploits the information of both uncensored and censored patients: the first term captures the first hitting time to maximize the patient's risk at the specific time of occurrence of the event; the second one maximizes the survival function, defined as $1 - F(t)$ [51], evaluated at the patient's last follow-up.

L_2 employs a ranking loss function based on the concept of concordance. Essentially, it indicates that a patient who has passed away at a certain time should be considered at a greater risk than a patient who is still alive at that time. It is defined as:

$$L_2 = \sum_{i=1}^N \sum_{j \neq i}^N A_{i,j} \cdot \exp \left(- \frac{\hat{F}(s^{(i)}|\mathbf{x}^{(i)}) - \hat{F}(s^{(i)}|\mathbf{x}^{(j)})}{0.1} \right)$$

where $A_{i,j} = \mathbb{1}(k^{(i)} = 1, s^{(i)} < s^{(j)})$. Using these terms the loss function penalizes the incorrect ordering of pairs. Indeed, on the one hand, $A_{i,j}$ is an indicator function that identifies pairs of patients whose comparison is meaningful, known as acceptable pairs. These pairs consist of patients in whom the first patient experienced the event at a specific time point, while the second patient has a longer survival time. On the other hand, the exponential function compares the risks faced by the two patients, i and j , at the time the first one experienced the event. In particular, it assumes small values in the case of a correct ordering, i.e. if the first patient has a higher risk than the second, whilst it assumes large values in the case of a wrong ordering.

3.4.3 Evaluation Metric

Once a survival analysis model is trained, we need to evaluate its outputs taking time into account, as does the time-dependent concordance index (*Ct-index*) does. It is an evolution of the commonly used *C-index*; this latter is based on the assumption that patients who lived longer should be assigned lower risks than those who lived a shorter period. However, the *C-index* does not account for any potential changes in risk over time, whereas the *Ct-index* is calculated by comparing pairs of patients, where one patient has experienced the event at a specific time while the other patient has neither experienced the event nor been censored

by that time. Formally, *Ct-index* is defined as:

$$\begin{aligned} Ct\text{-index} &= P\left(\hat{F}(s^{(i)}|\mathbf{x}^{(i)}) > \hat{F}(s^{(i)}|\mathbf{x}^{(j)})|_{s^{(i)} < s^{(j)}}\right) \\ &\approx \frac{\sum_{i=1}^N \sum_{j \neq i}^N A_{i,j} \cdot \mathbb{1}\left(\hat{F}(s^{(i)}|\mathbf{x}^{(i)}) > \hat{F}(s^{(i)}|\mathbf{x}^{(j)})\right)}{\sum_{i=1}^N \sum_{j \neq i}^N A_{i,j}} \end{aligned}$$

This index is not based on a fixed time, unlike the *C-index*, which considers only the ordering of subjects at a fixed time point. Instead, the *Ct-index* takes into account the timing of events for each subject over time. In fact, using $A_{i,j}$ and the indicator function that verifies whether, at the time the first patient experienced the event, he/she presents a higher risk than the second one, this metric computes the fraction of correctly ordered patients out of the total number of acceptable pairs.

3.4.4 Experimental setup

Our architecture adopts 12 consecutive encoder layers (M), with each layer utilizing 17 attention heads and a feed-forward module composed of one hidden layer of 3072 neurons. We do not further investigate any other hyperparameters configuration, since their tuning is out of the scope of this manuscript. Nevertheless, the "No Free Lunch" Theorem for optimization states that there is no universal set of hyperparameters that will optimize the performance of a model across all possible datasets [52].

To have a fair comparison between the models employed in the analysis, we applied 5-fold stratified cross-validation, maintaining the distribution of censored and uncensored patients among the different folds. Thus, we divided the data into test set (20%) and train set (80%), part of which was used as validation set (20%).

We compared our approach with each pair of state-of-the-art imputation strategies and models for OS analysis presented in section 3.2. Regarding the imputation methods, we employed overall mean imputation strategy, kNN imputer, MICE and MissForest with their default parameters [53, 54]. More specifically, we opted for 5 neighbors and the euclidean distance for the kNN imputer, for a maximum number of iterations equal to 10 and a mean initial strategy for the MICE strategy, and for a maximum number of iterations equal to 10, 100 estimators and the squared error as a criterion for the growth of the MissForest. We opted not to use the complete case analysis since in our application the number of patients would be almost halved, passing from 297 to 158. Focusing on the models, we tested the CPH, the ST, the RSF and the DH. For training the CPH, ST and RSF, we used their default parameters [55]. More specifically, we opted for the breslow method to handle tied

event times and 100 iterations in the training of CPH, for the best splitter, a minimum of 6 samples to split a node and 3 samples to define leaf nodes in the growth of the ST, and 100 estimators, a minimum of 6 samples to split a node and 3 samples to define leaf nodes in the growth of the RSF. In training both our model and DH we opted for the Adam optimizer, a batch size of 32, an initial learning rate of 10^{-4} and a Xavier initialization. The training was set to a maximum of 1500 epochs, in conjunction with an early stopping criterion and a learning rate scheduler both based on the validation loss with patience of 200 and 100 epochs, respectively.

To conduct the comparisons, we used different units of time, i.e., one month, one year, and two years, covering a period of six years. Note that we defined this time limit in order to include at least 95% of the patients' survival times without any modification, whereas we considered all the patients with a longer survival time as censored. Thus each element of the output vector indicates the level of risk that the patient faces within the corresponding time interval. It is worth emphasizing that the output vector's size N , changes across the different experiments based on the specific unit of time employed: 72 elements for the 1-month, 6 for the 1-year and 3 for the 2-year time unit.

3.5 Results and Discussions

As reported in section 3.4.4, we compared the performance of our approach with the state-of-the-art OS models, trained on the CLARO dataset [50], using different imputation strategies. The results, in terms of *Ct-index*, are presented in Table 3.2, where the first two columns report the combination of model and imputation method employed, whereas the rest of the columns represent the units of time applied. The mean value and standard error reported are calculated on the different folds and the best mean performance for each unit of time is marked in bold. As we can see, our approach always outperforms the benchmarks.

These results, independently from the imputation strategy applied, confirm the considerations made in section 3.2. It appears that CPH performs the worst, likely due to the limitations of the proportional hazard assumption in this context. Conversely, all other methods, which eliminate this constraint, perform better than CPH. Notably, DH and our approach outperform RSF, showcasing the potential of deep learning once again. We deem that the improvement in the performance of our approach compared with DH is due to its ability to handle a high degree of missing data without the need for imputation, which has the potential to bias the final prediction. Indeed, the missing data imputation itself poses some challenges in the selection of the most appropriate approach to the task at hand; furthermore, most of the existing imputation methods struggle to handle high levels of missing

Model	Imputation	1-month	1-year	2-year
CPH	Mean	61.10 ± 3.05*	60.92 ± 3.59*	47.64 ± 8.72
	kNN	60.72 ± 3.90	59.72 ± 5.46	48.51 ± 9.14
	MICE	60.80 ± 3.16*	60.01 ± 4.29*	47.16 ± 8.74
	MissForest	61.34 ± 3.05*	58.95 ± 4.11*	44.99 ± 9.44
ST	Mean	23.46 ± 4.24*	39.27 ± 4.02*	43.82 ± 2.78*
	kNN	27.05 ± 3.14*	22.51 ± 3.20*	35.28 ± 4.29*
	MICE	22.03 ± 4.17*	45.51 ± 4.02*	37.29 ± 7.47*
	MissForest	22.31 ± 3.48*	28.45 ± 8.12*	37.15 ± 10.88*
	–	26.49 ± 2.05*	37.61 ± 7.32*	42.71 ± 6.19*
RSF	Mean	65.88 ± 2.37	72.16 ± 3.02	57.59 ± 7.08
	kNN	61.92 ± 3.27*	68.24 ± 1.34*	54.66 ± 8.10
	MICE	67.10 ± 3.31	71.53 ± 4.06	57.45 ± 8.21
	MissForest	62.86 ± 2.66	67.69 ± 5.99	58.79 ± 11.16
DH	Mean	69.72 ± 3.67	75.19 ± 5.38	78.23 ± 4.11
	kNN	68.24 ± 3.94	71.28 ± 6.54	77.94 ± 3.66
	MICE	71.04 ± 2.92	75.26 ± 5.81	78.39 ± 3.54
	MissForest	68.39 ± 3.41	72.64 ± 5.42	74.64 ± 5.34
Ours	–	71.97 ± 2.39	77.58 ± 1.82	80.72 ± 9.11

Table 3.2: Performance of tested models in terms of *Ct-index* (mean ± standard error). The asterisks (*) denote experiments with performances that are statistically different (p -value < 0.05) from our proposed model.

data, leading to a drop in performance. On the contrary, our approach smoothly copes with this situation by learning only from the available features.

When we examine the *Ct-index* of the different models at various units of time, we observe that those of DH and our model improve as unit time increases. This can be attributed to the reduction of complexity for the task of correctly ordering pairs, as the granularity of the problem reduces, suggesting that our method effectively interprets the available features to estimate the OS time.

Focusing on the differences between the various imputation strategies, Table 3.2 highlights how difficult is to determine the most appropriate method, since the outcome depends both on the data and the models employed in the analysis. Therefore, our approach not only proves to be the best in terms of achieving the highest performance, but it also enables us to eliminate one of the variables from the problem by simply disregarding the missing features instead of searching for the most appropriate imputation strategy for the task at hand.

To further validate our results, we performed a paired t-test to evaluate the statistical

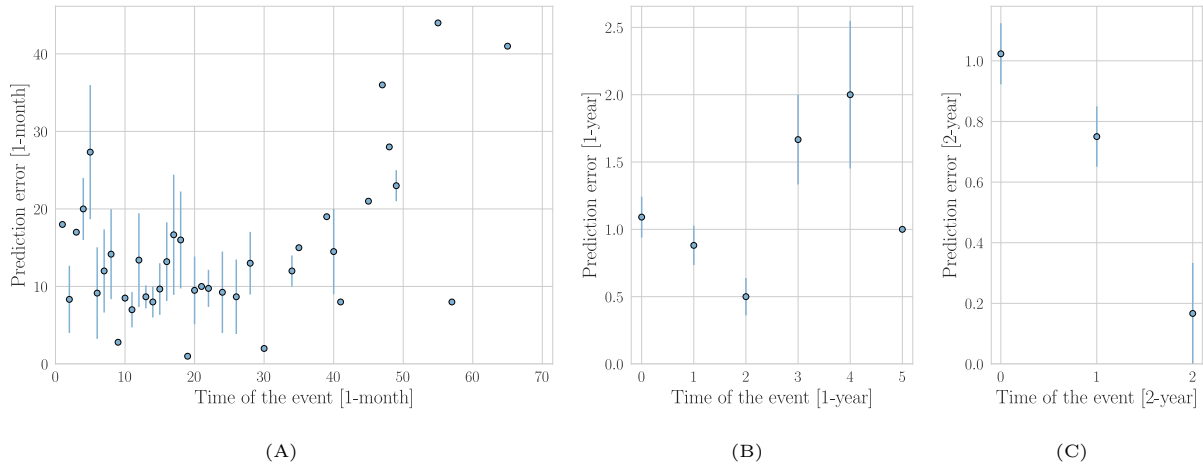


Figure 3.2: Prediction errors made for different patient groups by the proposed method, taking into account the actual times of occurrence of the event. The graph reports the mean and standard error of the prediction of uncensored patients, as for censored patients no valid prediction can be made since the event did not occur.

difference between the performance obtained in different folds by our model compared with all competitors. In Table 3.2, denoted by the asterisks (*), we report the results with a $p\text{-value} < 0.05$. These results show that our approach is significantly better than ST for all combinations with imputers and with all time granularities. Most notable, among the various combinations, is the case without imputers, showing that our model can handle missing values better. With regards to CPH, our model is statistically better than the configurations in combination with overall mean imputation, MICE and MissForest. Moreover, when the time granularity is increased to 2 years, the performance is no longer statistically distinguishable from those of our model, again indicating the reduction in complexity of the problem as granularity increases. In contrast, looking at the results obtained by RSF and DH, only the combination of RSF with kNN imputer for the time granularities of 1 month and 1 year are statistically worse than those of our model. Despite this, the results obtained from our model are still worthwhile since they are on average superior or comparable to the others, but having simplified the overall analysis, since no imputation of missing data is required.

Moreover, we performed an analysis to better understand on which patients the proposed model makes most errors. In Figure 3.2, in panels A, B, and C for the time granularities of 1 month, 1 year, and 2 years, respectively, we reported the mean errors and relative standard errors for uncensored patients grouped by unit of time in which the event occurred. From panel A we can observe that most uncensored patients, who present the event in the first few years after the diagnosis, have the lowest error, which is below 30 months, whereas, with regard to those patients who have survived longer, the errors are on average larger, even

surpassing 40 months. When the granularity of time increases, on the other hand, the trend reverses, leading to a minimum error of 4.8 months (0.2 of a 2-year period) in the case of patients who died in the third two-year period, as reported in panel C. This trend inversion, which probably also affects performance explaining its increase as time granularity increases, is perhaps due to the difference in the numbers of patients in the different groups of patients.

Additionally, we performed an explainability analysis to better understand which features most influence the model’s decisions. We used the SHAP method [56], an explainability technique based on Shapley values, to estimate the contribution of each feature to the final output. Figure 3.3 shows the SHAP summary plots for the three models implemented with the 3 time units, i.e., 1 month, 1 year, and 2 years. In these plots, the order of the features is based on their importance in distinguishing the uncensored class. The higher the position of a feature in the plot, the more important it is. To represent the global feature importance, we averaged the absolute SHAP values obtained for each patient and time represented in the output vector. Moreover, to get a general overview of the importance of the categorical features, we averaged the contributions of each category. As we can see, in all time granularity models, the most relevant features for the task at hand refer to disease-related values such as the volume of the *CTV* and the *T*, *N* and *M* stages, which are measures of the severity of the tumour [57].

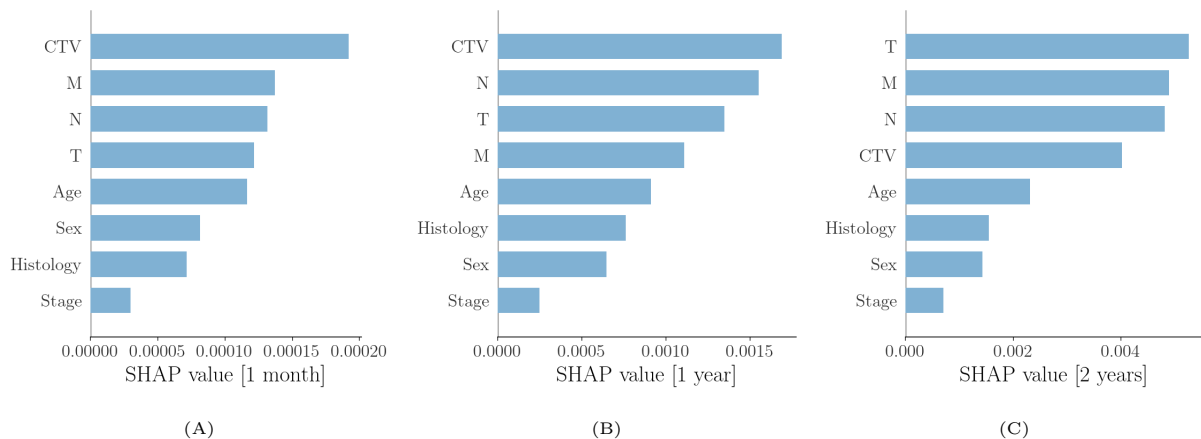


Figure 3.3: SHAP summary plots of features’ contributions in the 3 models implemented with the 3 time units: Panel A) 1 month; Panel B) 1 year; Panel C) 2 years. The plots show the global feature importance by averaging the absolute SHAP values obtained for each patient and time represented in the output vector.

Furthermore, we conducted an ablation study to gain a better understanding of the individual contributions of its two terms, L_1 and L_2 . We examined whether both terms played a role in the model’s final performance (Figure 3.4.A) and in achieving convergence

(Figure 3.4.B). The figures depict the *Ct-index* and the number of epochs, respectively, which clearly demonstrate that utilizing both terms together positively impacted the model’s performance and the convergence time.

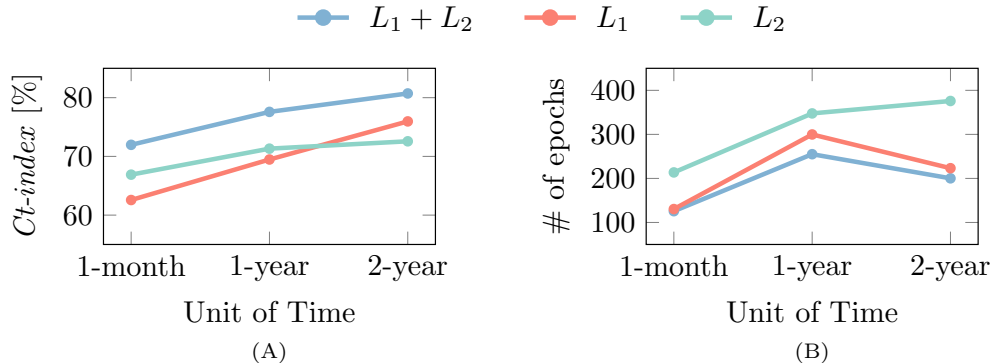


Figure 3.4: Ablation study of the two terms of the loss function proposed in [49]: (A) Average performance (*Ct-index*) and (B) mean number of epochs to achieve convergence.

Finally, considering that the model is designed for a medical application, we analyzed its computational burden and inference speed. We evaluated the former in terms of the model’s weight count, approximately 2.7 million, which implies significantly longer training times compared to ML competitors, despite not requiring imputation. However, since training would occur prior to deployment, the critical factor for practical application is the inference time, approximately 5×10^{-3} seconds per sample, making it feasible for clinical practice.

3.6 Conclusion

In this manuscript, we proposed a novel approach to address the problem of missing data in the context of survival analysis of NSCLC. This task is usually tackled, in the literature, either by discarding the incomplete data samples or employing some imputation strategy, which could bias the final prediction. Conversely, our method is able to ignore those values, taking into account the available ones only. All the experiments are performed on the in-house dataset, CLARO [50], consisting of the clinical information of 297 patients suffering from NSCLC.

From a clinical perspective, the development of reliable and accurate prognosis prediction tools, to be used prior to treatment initiation, is a critical and unmet need in clinical practice. The availability of such tools would allow clinicians to tailor treatment strategies to the anticipated response, intensifying or descaling therapy as necessary based on the patient’s prognosis.

Experimental results show that our model outperforms conventional OS methods at various time units with different imputation techniques, bringing an increase in performance when missing features are present.

The results described so far and the limitations of this work suggest future directions worthy of investigation. These include conducting additional experiments of our method on other OS datasets, as we are aware that validating a model on one dataset only, even if it is representative of a large population since it comes from a large Italian metropolis, is somewhat limiting, but obtaining medical datasets, preferably multicentric, is not always a straightforward process, particularly the raw version of the data in which imputation or exclusion of missing samples and features is not performed. Therefore, in future works, we plan to extend the range of datasets, artificially generating different missing value scenarios on which to test our approach. Furthermore, this work considers clinical features only, thus having a reduced view of the patient’s status, which could be expanded by taking into account other sources of information and types of data. Hence, another future work could expand our approach to include various types of data beyond tabular data (e.g., imaging data). For a more seamless integration in the clinical setting, it could also be interesting to conduct further development to attempt to reduce the training time, by reducing the number of trainable weights, so as to accommodate eventual updates over time as new patient data are acquired. Moreover, it could be interesting to explore the generalization ability of the proposed approach beyond the time-to-event analysis, thus directing future works in testing it in other task domains such as classification and regression. In the meanwhile, we opted to make publicly available our GitHub repository² in a way to make it possible to test our approach on other datasets.

In this chapter, we investigated the prediction of OS in oncology patients using real-world clinical data. We demonstrated how DL models can provide meaningful and interpretable prognostic insights, even in the presence of naturally occurring missing values. The analysis highlighted the challenges posed by incomplete follow-up, and heterogeneity in clinical records. While these models are informative and clinically relevant, their applicability is constrained by their monomodality, only tabular data were utilized. However, modern clinical practice increasingly relies on imaging data alongside structured information. To advance towards a more holistic modeling approach, the next chapter explores the integration of radiological and clinical data through multimodal pipelines for survival prediction.

²<https://github.com/cosbidev/OSTransformer>

Chapter 4

Fusing Imaging and Clinical Data for NSCLC Survival Prediction

The growing availability of both imaging and structured clinical data offers new opportunities to develop predictive models that reflect the full complexity of medical decision-making. In this chapter, we take a first step toward multimodal modeling by combining CT scans and clinical variables to predict overall survival in patients with non-small cell lung cancer (NSCLC). We investigate whether integrating heterogeneous modalities yields tangible benefits over single-modality models and explore fusion strategies to effectively combine imaging and tabular features. This study represents a preliminary but essential attempt to bridge the gap between unimodal robustness and multimodal expressiveness in real-world clinical tasks.

4.1 Introduction

Lung cancer is the second most common type of tumour worldwide, accounting for approximately 11.4% of all cases [32], and it is the first in terms of number of deaths. Non-small-cell lung cancer (NSCLC) is the most frequent, with approximately 82% of all cases [33]. The most common treatment options, selected according to patients' characteristics, include radiotherapy, chemotherapy, surgical resection, and immunotherapy but also targeted therapy [33, 58].

Overall survival (OS), a measure of the time elapsed from the date of diagnosis until the patient's death, allows the identification of subgroups of patients with a better or worse prognosis. Nevertheless, the 5-year survival rate for NSCLC is 26%, and it drops further to 7% when local recurrence or distant metastases occur [33]; in this respect, strategies to

improve OS are urgently needed.

Over the last few years, there has been a growing interest in the development and application of Artificial Intelligence (AI) methods to oncology to help personalised medicine make further progress by facilitating the identification of the correct treatment for each patient. This has fostered the emergence of radiomics, which represents the bridge between medical imaging and personalised medicine since it computes, in a non-invasive manner, quantitative characteristics from medical images, such as CT, MRI, X-ray, and PET, representing tumour phenotype [59, 60, 61, 62]. In addition to radiomics, researchers have attempted to extract prognostic information from other modalities, e.g., genome sequencing, whole-slide images (WSI), etc. [63, 64, 65]. For example, genomics data from a tumour allow the identification of cancer driver genes, whilst a WSI from a biopsy provides insight into the morphology and microenvironment of the tumour.

Several learning methods exist to perform these prognostic tasks, which can be roughly divided into model-based and data-based approaches. The former assume a model to describe the data trend, whilst the latter, exploiting the current large availability of digital repositories and using increasingly high-performance AI algorithms, learn directly from the data. In lung cancer predictive applications, such learning methods usually exploit one modality only [66, 67, 68, 69, 70, 71], but the availability of multimodal data, which provide complementary information about the phenomenon under investigation, has led to the development of multimodal learning techniques able to cope with different information and to perform significantly better than unimodal models [72, 73, 74, 75, 76]. From an AI perspective, early, joint, and late fusion are the three main fusion techniques to merge different modalities' information. In the first technique, the features of each modality are merged according to a rule into a feature vector to be given to the learner; in the second, the different modalities are merged at hidden and embedded levels, whilst in the last technique, the predictions made using the individual modalities are aggregated according to an aggregation rule.

In NSCLC, several studies have searched for a set of quantitative biomarkers, also referred to as a signature, to predict the overall survival. Among them, Table 4.1 summarises those using multimodal approaches [77, 78, 79, 80], which are also now shortly overviewed.

In [77], the authors used the NSCLC dataset available on The Cancer Imaging Archive (TCIA) [81] to present an early fusion-like approach which fuses PET and CT images, using a technique based on 3D discrete wavelet transform to combine spatial and frequency features, and then it extracts radiomic features (first-order, textural, and moment invariant features). After performing feature selection via univariate Cox analysis, the authors applied the Kaplan–Meier method. The proposed approach obtained a concordance index (C-index) of 0.708, measured with 1000-time bootstraps, which is higher than the results they achieved

from unimodal and traditional early fusion approaches (concatenation and averaging of the feature vectors separately extracted for each modality).

In [78], the authors used another NSCLC dataset also available on TCIA [82], and they performed an early fusion of deep features extracted from CT images and clinical data. The former were extracted using a 3D-ResNet34, whilst the latter using a Multilayer Perceptron (MLP). The concatenation of these features fed an MLP. In 5-fold cross-validation with a patient-level split, the authors tested different configurations by varying the structure of the ResNet, the depth of the final MLP, and the ratio between the number of the two types of deep features, achieving a C-index equal to 0.658 as best result.

In [79], the authors developed a hierarchical multicriterion fusion strategy to combine the predictions made by various classifiers working with different modalities. Even this study is based on the same data available on TCIA [82] used by [78], and it only takes into account 316 patients in whom the gross tumour volume was delineated. This permitted to extract clinical features and radiomic features (textural and non-textural) for each patient that, after a feature selection step separately performed for the two modalities were fed into the system. The modular architecture allows each modality to be analysed separately with a set of classifiers (Support Vector Machine, k-Nearest Neighbours, Decision Tree, Random Forest, and Extreme Gradient Boosting). By means of a sequence of aggregation rules that weight the contribution of each classifier to the output probability of each modality and then combine the probabilities of each modality, the system produces the final prediction. The experiments, run in 5-fold cross-validation, return an Area Under the ROC Curve (AUC) equal to 0.81.

In [80], the authors used the data in the National Cancer Institute's Genomic Data Commons database [83] to develop a multimodal deep learning method for long-term pan-cancer survival prediction, called MultiSurv, which works with six different modalities, namely clinical data, gene expression, microRNA expression, DNA methylation, gene copy number variation data, and WSI. In this modular architecture, each input data modality is handled by a dedicated submodel. For the clinical and omics submodels, they used an MLP, whilst for the imaging submodel a ResNeXt-50. The data fusion layer aggregates the multimodal feature representations by taking the element-wise maxima across the set of representation vectors, allowing any missing modalities to be handled as well. The fusion vector is the input to an MLP, which returns as output a vector of probabilities, one for each time interval of a set of predefined follow-up time intervals. This system was trained in an end-to-end fashion, applying an holdout cross-validation stratified by cancer type. The authors evaluated the model with different numbers and combinations of the six modalities, and the best performance was obtained with bimodal inputs combining clinical data with gene expression

Author	Modalities	Study Population	Number of Patients	Data Representation	Fusion Modality	Learning Model	Performance
[77]	CT, PET	NSCLC I-IV stages	182	Radiomic features extracted from an image obtained by merging PET and CT scans through a technique based on 3D discrete wavelet transform	Early	Kaplan–Meier method	C-index: 0.708
[78]	CT, clinical data	NSCLC I-III stages	422	Concatenation of deep features extracted by a 3D-ResNet34 and an MLP for CT images and clinical data, respectively	Early	MLP	C-index: 0.658
[79]	CT, clinical data	NSCLC I-III stages	316	Clinical data and radiomic features	Late	Modular architecture with SVM, DT, KNN, RF, and XGBoost as base classifiers	AUC: 0.81
[80]	Clinical data, gene expression, microRNA expression, DNA methylation, gene copy number variation data, and WSI	33 different cancer types	11.081	Element-wise maxima across the set of representation vectors of single-modality submodels	Joint	Modular architecture, with dedicated input data modality submodels, a data fusion layer, and a final survival prediction MLP submodel	Time-dependent C-index: best 0.822 lung squamous cell carcinoma 0.554
<i>Putting our work in the background</i>	CT, clinical data	NSCLC II-IV stages	191	Clinical data and CT slices	Optimisation-driven late	Multimodal ensemble of learners trained on different modalities and selected by a multiobjective optimisation algorithm	ACC: 0.75

Table 4.1: Summary of the background on the multimodal learning to predict the overall survival in NSCLC. For the sake of completeness, the last section puts our contribution in the context of the literature.

(time-dependent C-index: 0.822).

Although the works in the literature achieved promising results, they are few in number, despite the importance of predicting the overall survival in NSCLC cancer that, in turn, may open the chance to develop personalised therapeutic approaches. Furthermore, two out four of such contributions explored early fusion, one investigated late fusion, and the other joint fusion. In particular, the one using late fusion computes handcrafted features from CT images that feed well-established classifiers. Nevertheless, in the last decade, deep learning

has shown its potential in several fields, medical imaging included [84, 85, 86], to automatically learn discriminative features directly from images, without being limited to using predefined features or other descriptors whose definition come from researchers' experience. In particular, Convolutional Neural Networks (CNNs) are a well-established set of network architectures exploiting convolutional layers (and their variations) to learn a compact hierarchical representation of the input that well fits the specific task to solve. In this respect, and as an evolution of the state-of-the-art shown in Table 4.1, in this work, we present a method to algorithmically optimise the way to set up a multimodal ensemble of deep networks, which are then combined by a late fusion approach. Such an ensemble uses image and clinical data to tackle the challenge to predict the overall survival in a cohort of 191 patients affected by NSCLC cancer. Exploiting the classifications of different unimodal models, we propose an optimised multimodal late fusion approach, whose performance is shown in Section 4.4. In particular, our method addresses a key and open question in multimodal deep learning [73, 87], i.e., which should be the deep networks for each modality to be combined in the ensemble among the many available.

The manuscript is organised as follows: the next section describes the materials, and Section 4.3 introduces the methods. Section 4.4 presents and discusses the experimental results; finally, Section 4.5 provides concluding remarks.

4.2 Materials

Our clinical decision support system uses image and clinical data available within the CLARO dataset, which includes 191 NSCLC patients treated with concurrent chemoradiation for locally advanced NSCLC (86% of cases) and systemic treatment in the metastatic setting (14%). During treatment, all patients underwent weekly chest Computed Tomography (CT) scans, without intravenous contrast, to assess acute toxicity and tumour shrinkage, which were reviewed by two radiation oncologists independently. For all CTs, each physician was able to judge whether reduction was: (a) present and clinically significant, (b) present and clinically non-significant, or (c) absent. In the case of physician agreement for the (a) category, a contrast-enhanced CT was performed to better visualise node reduction, a new target volume was delineated, and a new treatment plan performed. Patients were treated without any time break.

The population was enrolled under two different approvals (the retrospective and prospective phases) of the Ethical Committee. The former was approved on 30 October 2012 and registered at ClinicalTrials.gov on 12 July 2018 with Identifier NCT03583723, whilst the latter was approved on 16 April 2019 with Identifier 16/19 OSS, and it was closed on April

2022. The Institutional Review Board approved this review. Written informed consent was obtained in all patients. The authors confirm that all ongoing and related trials for this intervention are registered.

The median OS for the entire population was 15.64 months, with a mean of 23.85 ± 77.22 (95% CI). The patients were then clinically followed until they were divided into two classes based on the median OS of all the patients: 95 dead and 96 alive.

4.2.1 Imaging

The characteristics investigated were extracted from CT scans collected at the time of patient diagnosis, on which expert radiation oncologists delineated the Clinical Target Volume (CTV).

For each patient, the CT images were acquired before the treatment using a Siemens Somatom Emotion, with 140 Kv, 80 mAs, and 3 mm for slice thickness. The scans were preprocessed applying a lung filter (kernel B70) and a mediastinum filter (kernel B31).

4.2.2 Clinical Features

Clinical data contained different information, which are listed in Table 4.2 together with the number of missing values and the distribution for each tabular feature among the different discrete values. To define the stage of the tumour, two experienced radiation oncologists (ROs) independently reviewed CT scans and assigned the staging scores of the tumour (T, N, and tumour stage); in case of disagreement, they reviewed the CT images together until consensus was reached. In addition to staging, age, and sex, Table 4.2 shows that we also collected features describing the histology of the tumour and the initial CTV, so that the clinical data account for seven descriptors in total.

In the imputation of missing values, the median value and the mode of the training set data were assigned for the numerical and categorical features, respectively. Furthermore, it should be noted that not all patients underwent a histopathological examination. Nevertheless, since on the one side it was not possible to impute the histology of the tumour and, on the other side, this feature could be informative, we add a virtual category named *unknown*.

4.3 Methods

To predict the prognosis in terms of binary classification task over the OS, we exploited both the images and clinical data described before, which were processed by a multimodal DL pipeline that, in the training phase, finds the optimal combination of models of different

Feature	Missing Data	Categories	Distribution
<i>Age</i> *	26 (13.62%)	<71 years	82 (42.93%)
		≥ 71 years	83 (43.46%)
<i>CTV</i> *	37 (19.37%)	<114.88 cm ³	77 (40.31%)
		≥ 114.88 cm ³	77 (40.31%)
<i>Sex</i>	0 (0.00%)	Male	133 (69.63%)
		Female	58 (30.37%)
<i>Histology</i>	0 (0.00%)	Adenocarcinoma	95 (49.74%)
		Squamous	59 (30.89%)
		Other	11 (5.76%)
		Unknown	26 (13.61%)
<i>Stage</i>	0 (0.00%)	II	4 (2.09%)
		III	160 (83.77%)
		IV	27 (14.14%)
<i>T stage</i>	36 (18.85%)	T0	1 (0.52%)
		T1	9 (4.71%)
		T2	32 (16.75%)
		T3	65 (34.03%)
		T4	48 (25.13%)
<i>N stage</i>	26 (13.61%)	N0	15 (7.85%)
		N1	33 (17.28%)
		N2	93 (48.69%)
		recurrence N2	6 (3.14%)
		N3	18 (9.42%)

Table 4.2: Patients’ characteristics. As marked by *, note that, although *age* and *CTV* are continuous variables, for the sake of synthesis we report here their distribution considering their median values as thresholds, whilst the model used the continuous values. The division into stages is further defined by letters (a, b, and c), which are not reported for the sake of brevity, but the model uses the actual stages.

modalities via multiobjective optimisation. The idea stems from observing that today many deep neural networks are available, both in terms of architectures as well as of pretrained weights. This allows researchers to train or fine-tune them to search for the most suitable for the task at hand. Furthermore, it is well-known that, in many cases, ensembles of classifiers combined in late fusion provide better performance than unimodal models [88], but, at the same time, the learners in the ensemble have to complement each other, i.e., they have to make wrong decisions on different samples. Therefore, the abundance of available models asks for methods to support researchers in determining which is the best multimodal ensemble, a challenge that we address using an algorithmic and multimodal approach, schematically represented in Figure 4.1. It works with m different modalities and M different models, so that M_m is the number of models available for the m th modality. Furthermore, we denote with E an ensemble built using one or more models per modality, whose outputs are

combined by majority voting. Figure 4.1 shows that our method essentially consists of three main steps:

- Training all the available models for every single modality using the training sets defined by the bootstrap validation approach;
- Finding the multimodal set of unimodal models solving a multiobjective optimisation problem working with evaluation and diversity scores, which are computed on the validation sets defined by the same bootstrap approach;
- Computing the performance on the test sets defined by bootstrap, which are then averaged out (block “Average performance evaluation”).

These steps are now detailed in the next subsections.

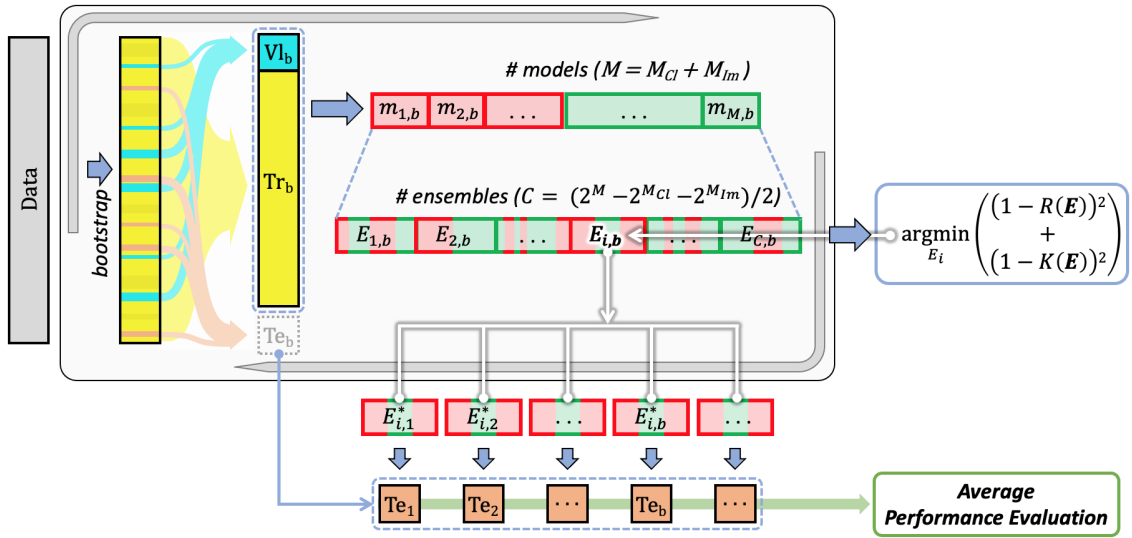


Figure 4.1: Schematic view of the pipeline. Symbols: Tr : training set, Vl : validation set, Te : test set, m : model, M : number of models (\cdot_{Cl} : for clinical data and \cdot_{Im} : for data from images), \cdot_b : a generic bootstrap fold, E : (models’) ensemble, \cdot_i : a generic ensemble, C : number of ensembles, R : function of recall, K : function of diversity.

4.3.1 Training

To obtain the optimal ensemble E^* of models, the first step is to independently train and evaluate the different M unimodal models on the respective m modalities. In our scenario, we had $M = M_{Cl} + M_{Im}$, where M_{Cl} and M_{Im} denote the number of models for the clinical data and the imaging modality, respectively.

With respect to the clinical data, we worked with $M_{Cl} = 7$ different ML and DL models, which are acknowledged in the literature as those that best work with this modality [89]. In alphabetical order they are:

- AdaBoost as a cascade of classifiers;
- Decision Tree (DT) as tree model;
- Multilayer perceptron (MLP) as neural architecture with one hidden layer with 13 neurons and 1 neuron in the output layer, which use the ReLU and Sigmoid activation functions, respectively;
- Random forest (RF) as an ensemble of trees;
- Support Vector Machine (SVM) as a kernel machine;
- TABNET [5] as a neural architecture;
- XGBoost a variation of the AdaBoost that uses a gradient descent procedure to minimise the loss when adding weak learners.

Let us now turn the attention to see image modality. We worked with $M_{Im} = 30$ different CNNs from 8 architecture families, which have proved to have promising results in many biomedical applications [90]. They are:

- AlexNet [91];
- VGG [92]: VGG11, VGG11-BN, VGG13, VGG13-BN, VGG16, VGG16-BN, VGG19, VGG19-BN, where the suffix BN means that batch normalization is used;
- ResNet [93]: ResNet18, ResNet34, ResNet50, ResNet101, ResNet152, ResNeXt50, ResNeXt101, Wide-ResNet50-2, Wide-ResNet101-2;
- DenseNet [94]: DenseNet121, DenseNet169, DenseNet161, DenseNet201;
- GoogLeNet [95];
- ShuffleNet [96]: ShuffleNet-v2-x0-5, ShuffleNet-v2-x1-0, ShuffleNet-v2-x1-5, ShuffleNet-v2-x2-0;
- MobileNetV2 [97];
- MNasNet [98]: MNasNet0-5, MNasNet1-0.

All the CNNs were pretrained on the ImageNet dataset [99]. The architectures, layer organisation, and complexity of such models gave us the opportunity to investigate how different models perform on the task at hand.

4.3.2 Optimisation

To answer the question of which architectures should be used to construct the best multi-modal ensemble, we solved a multiobjective optimisation problem that works with two scores capturing different views of the ensemble performance. Indeed, given an ensemble E , on one side we measured its recall (R) using, straightforwardly, the labels computed by applying the aforementioned majority voting scheme. R is defined as

$$R = \frac{TP}{P} \quad (4.1)$$

where TP is the number of true positive classifications and P is the number of positive instances, and it measures the sensitivity of the model, a desirable property in our application ensuring that no positive patients get excluded before treatment. On the other side, the optimisation algorithm also works with the kappa diversity (K), a pairwise score measuring to what extent two models provide the same errors. It is defined as

$$K = 1 - \frac{2(N_{11}N_{00} - N_{01}N_{10})}{(N_{11} + N_{10})(N_{01} + N_{00}) + (N_{11} + N_{01})(N_{10} + N_{00})} \quad (4.2)$$

where N_{11} and N_{00} are the number of instances classified correctly and incorrectly by each of the two models under consideration, respectively, and N_{10} and N_{01} are the number of instances classified correctly by the first model and incorrectly by the second and vice versa, respectively. The overall ensemble diversity is given by

$$\frac{2}{|E|(|E| - 1)} \sum_{i=1}^{|E|-1} \sum_{j=i+1}^{|E|} k \quad (4.3)$$

where $|E|$ is the number of models in E . Given these premises, let us notice that both R and K range in $[0, 1]$, and the higher the values, the more accurate and diverse the models. Hence, our algorithm solves the following multiobjective problem to determine the best ensemble E^* :

$$E^* = \arg \min_E [(1 - R(E))^2 + (1 - K(E))^2] \quad (4.4)$$

s.t.

$$\begin{cases} |E^*| > 1 \\ |E^*|_{mod2} = 1 \\ |E^*|_{Cl} \geq 1 \\ |E^*|_{Im} \geq 1 \end{cases} \quad (4.5)$$

where $R(E)$, $K(E)$ represent the average values of R and K , respectively, of an ensemble E computed across all the validation sets given by bootstrap, a choice that avoids any bias. Looking at the constraints, $|E^*|$ denotes the number of models in E^* , whilst $|E^*|_{Cl}$ and $|E^*|_{Im}$ stand for the number of models in E^* working with clinical and imaging data, respectively; finally, mod is the modulo operation. The first two conditions imply that the number of models in E^* is odd to prevent ties in the majority voting, whilst the third and fourth conditions ensure that at least one model for each modality is present in E^* . Note also that finding E^* is equivalent to finding the Pareto optimum of this optimisation problem, as we showed in [100, 101, 102]; nevertheless, here, we are extending our previous unimodal approach [100] to multimodal learning, as guaranteed by the last two conditions in Equation (4.5).

Hence, this optimisation algorithm performs an exhaustive search for the ensemble E^* that, among the $C = \frac{2^M - 2^{M_{Cl}} - 2^{M_{Im}}}{2}$ combinations of learners, returns the best classification performance and reduces the incidence and effect of coincident errors among its members, thus considering possible relationships between models and modalities. Furthermore, the simple minimisation of only one of the objective functions (R or K) is not the best approach, since some models may degrade the performance of the ensemble, and they may have redundant classifications between each other, not exploiting the trade-off between performance and diversity [103].

Finally, in the test phase, each input instance is given to all the learners in E^* , whose outputs are combined by majority voting to obtain the final prediction.

4.3.3 Preprocessing

Before feeding the data to the models, a preprocessing phase was executed for both modalities.

With reference to the clinical data, we applied one hot encoding to categorical features, so that the original 7 features were mapped to 27 descriptors, which in practice were used as input to all the classifiers mentioned before for the clinical data. Furthermore, numerical features were normalised in $[0, 1]$. No data augmentation was applied to the clinical data. For all clinical models listed in section 4.3.1, the default parameters of the libraries were

used.

With reference to the imaging modality, we used a U-Net to automatically align the images by detecting the region of interest of the scans by including the bounding cuboid segmenting the lungs. The U-Net architecture has proved to obtain good performance in many biomedical applications [104]. We trained this network on the TCIA publicly available dataset [82], which comprises 422 patients, and on a subset of our dataset, 125 patients whose lungs had already been delineated, with the goal of segmenting the lung pixels of each 2D slice. From this segmentation, we extracted the minimum bounding cuboid of the segmented volume, preventing any deformation once re-scaled. As input, the U-Net received 224 \times 224 images, and it was trained with an Adam optimiser and with a Dice loss function. The batch size was set to 32, and the number of epochs was equal to 50, but an early stop criterion was triggered at 13th epoch. We assessed the performance of this network in holdout cross-validation, obtaining a Dice score and an intersection over union equal to 98.5 and 97.0, respectively, which we considered satisfactory for our task.

Let us now focus on the image classification stage. All the CNNs work with 2D images, i.e., at CT slice level, and they need an input of size 224x224. To this end, each slice of the segmented lungs was resized to 224x224 and normalised with a min-max scaler, bringing the pixel values between 0 and 1. Random data augmentation was applied to prevent overfitting of the CNNs: horizontal or vertical shift ($-22 \leq \text{pixels} \leq 22$), random zoom ($0.9 \leq \text{factor} \leq 1.1$), vertical flip, random rotation ($-15^\circ \leq \text{angle} \leq 15^\circ$), and elastic transform ($20 \leq \alpha \leq 40$, $\sigma = 7$). The cross-entropy loss was used and was regulated by an Adam optimiser with an initial learning rate of 0.001, which is scheduled to reduce by an order of magnitude every time the minimum validation loss does not change for 10 consecutive epochs. For all the nets, a maximum of 300 epochs was fixed, with an early stopping of 25 epochs following the validation loss.

Given the fact that the CNNs work at slice level and the clinical data at patient level, to uniform the classifications, we aggregated the predictions of the slices of each patient via a majority-voting rule, thus obtaining a final outcome for each modality at patient level.

All the training processes were executed using an NVIDIA TESLA V100 GPU with 16 GB of memory, using PyTorch and Scikit-learn as the main coding library.

4.4 Results and Discussion

All the experiments were performed in bootstrap, performing five random extractions of the samples, where in each fold the proportions between the training, validation, and testing sets are 80%-10%-10%, respectively. Straightforwardly, in the imaging modality, all the slices

coming from the same patient were always in the same set.

Table 4.3 shows the results: each row corresponds to a classifier in the case of unimodal learners reported in the uppermost section; it corresponds to a multimodal ensemble in the middle section, and it corresponds to a competitor in the bottom most section. The columns report the performance measured in terms of accuracy, F-score, and recall to have a complete view of how the different models perform on the test sets. With reference to the unimodal learners, the values in Table 4.3 show that the best classifier working with clinical data is the AdaBoost, whilst in the case of image data, the best CNN is the VGG11-BN. Both achieve the largest accuracy and F-score among the pool of unimodal models, whilst the latter is also the best in terms of recall.

The ensemble returned by our algorithm, denoted by E^* in the table, achieves larger performance in terms of accuracy, F-score, and recall with respect to the unimodal classifiers. Whilst this could be expected in the case of the recall, as it is built maximising a function including this metric, it is interesting to note that this happens also in the case of the accuracy and F-score.

It is worth noting that the Pareto optimum E^* is composed of three models (two from the imaging modality and one from the clinical modality): ResNet34, VGG11-BN, and TAB-NET, which belong to different families, suggesting that each model interprets its modality in a different way to address the classification task. We notice that E^* has better performance (for all metrics) than the two best unimodal models. This finding implies that it is useful to fuse different modalities, each carrying useful and distinct information for the prognosis task whilst, at the same time, it is important to consider the diversity also, since it offers complementary points of view to the ensemble.

To assess the optimisation function, we also investigated which are the performances of the ensembles maximising only R or K , denoted as E^R and E^K , respectively (middle section of Table 4.3). The former consists of DT, RF, and AlexNet, whereas the latter comprises DT, RF, and DenseNet161. Moreover, in this case, the class predictions on the test set revealed that the outputs provided by E^* are better than those returned by E^R and E^K . This finding supports the importance of satisfying the proposed multiobjective optimisation condition. This agrees with the literature that in other fields, and in the case of the majority voting rule, reports that a necessary and sufficient condition for an ensemble to be more accurate than any of its models is if the models are accurate and diverse [105].

To further prove the efficacy of the proposed approach, the middle section of Table 4.3 presents the performance of the following other experiments:

- \bar{E} : it denotes the average performance for all the possible ensembles;

Classifier	Modality	Accuracy	F-Score	Recall
AdaBoost	Clinical	65.00 ± 5.00	67.35 ± 6.53	74.00 ± 16.73
DT	Clinical	60.00 ± 3.54	59.42 ± 9.15	62.00 ± 20.49
MLP	Clinical	61.00 ± 5.48	54.37 ± 23.57	60.00 ± 38.08
RF	Clinical	60.00 ± 6.12	60.72 ± 9.74	64.00 ± 16.73
SVM	Clinical	59.00 ± 2.24	55.46 ± 10.29	54.00 ± 18.17
TABNET	Clinical	63.00 ± 10.37	64.68 ± 11.69	70.00 ± 22.36
XGBoost	Clinical	54.00 ± 8.22	49.67 ± 16.74	50.00 ± 24.49
AlexNet	Imaging	50.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
DenseNet121	Imaging	62.00 ± 19.24	59.97 ± 27.79	66.00 ± 35.07
DenseNet161	Imaging	69.00 ± 6.52	68.28 ± 8.88	70.00 ± 20.00
DenseNet169	Imaging	71.00 ± 17.82	72.28 ± 17.44	76.00 ± 20.74
DenseNet201	Imaging	63.00 ± 16.05	65.95 ± 16.37	74.00 ± 23.02
GoogLeNet	Imaging	60.00 ± 6.12	50.04 ± 19.69	48.00 ± 31.14
MNasNet0-5	Imaging	51.00 ± 13.42	45.65 ± 19.37	44.00 ± 23.02
MNasNet1-0	Imaging	62.00 ± 7.58	65.11 ± 9.94	74.00 ± 20.74
MobileNetV2	Imaging	67.00 ± 17.18	68.61 ± 17.17	74.00 ± 23.02
ResNet101	Imaging	51.00 ± 5.48	49.97 ± 20.44	60.00 ± 38.08
ResNet152	Imaging	71.00 ± 7.42	63.65 ± 19.16	60.00 ± 30.82
ResNet18	Imaging	64.00 ± 18.84	58.74 ± 29.30	60.00 ± 33.91
ResNet34	Imaging	70.00 ± 11.73	71.71 ± 10.51	78.00 ± 22.80
ResNet50	Imaging	69.00 ± 11.40	69.45 ± 17.58	78.00 ± 27.75
ResNeXt101	Imaging	69.00 ± 7.42	68.95 ± 8.46	70.00 ± 15.81
ResNeXt50	Imaging	63.00 ± 10.37	64.35 ± 19.82	78.00 ± 33.47
ShuffleNet-v2-x0-5	Imaging	74.00 ± 10.25	74.66 ± 11.07	78.00 ± 16.43
ShuffleNet-v2-x1-0	Imaging	67.00 ± 17.18	67.14 ± 20.74	72.00 ± 26.83
ShuffleNet-v2-x1-5	Imaging	74.00 ± 13.87	72.3 ± 19.53	74.00 ± 27.02
ShuffleNet-v2-x2-0	Imaging	73.00 ± 9.08	71.23 ± 11.99	70.00 ± 20.00
VGG11	Imaging	50.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
VGG11-BN	Imaging	74.01 ± 16.36	75.03 ± 16.37	78.00 ± 19.24
VGG13	Imaging	50.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
VGG13-BN	Imaging	64.00 ± 8.22	61.58 ± 25.24	72.00 ± 35.64
VGG16	Imaging	50.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
VGG16-BN	Imaging	71.00 ± 13.42	72.19 ± 10.95	74.00 ± 13.42
VGG19	Imaging	50.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
VGG19-BN	Imaging	59.00 ± 15.17	51.68 ± 32.01	58.00 ± 38.99
Wide-ResNet101-2	Imaging	68.00 ± 10.95	69.84 ± 9.55	76.00 ± 20.74
Wide-ResNet50-2	Imaging	64.00 ± 13.87	66.02 ± 12.41	70.00 ± 18.71
E^*	Multimodal	75.00 ± 16.20	77.70 ± 13.83	84.00 ± 15.17
E^R	Multimodal	60.00 ± 6.12	58.15 ± 9.40	58.00 ± 17.89
E^K	Multimodal	61.00 ± 5.48	62.02 ± 9.58	66.00 ± 16.73
\bar{E}	Multimodal	66.58 ± 11.30	61.44 ± 15.13	62.35 ± 22.00
E_{post}^3	Multimodal	72.00 ± 12.04	75.41 ± 10.68	83.00 ± 15.17
E_{post}^{2+*}	Multimodal	70.94 ± 10.90	71.79 ± 10.21	74.91 ± 13.86
E^ℓ	Multimodal	61.00 ± 2.24	61.09 ± 8.11	64.00 ± 18.17
DeepMMSA [78]	Multimodal	59.00 ± 6.52	58.07 ± 12.32	52.00 ± 32.71
MCF [79]	Multimodal	62.00 ± 2.74	61.04 ± 10.53	64.00 ± 23.02

Table 4.3: Performance of all the tested models with the best for each modality reported in bold. Each column shows the mean value of a performance metric followed by the standard deviation. E^* is our optimum ensemble; E^R and E^K are the ensembles which maximise R and K , respectively; \bar{E} is the average performance for all the possible ensembles; E_{post}^3 is the ensemble consisting of the unimodal models with the largest recall, $\overline{E_{post}^{2+*}}$ is the ensemble with the two unimodal classifiers with the largest recall per modality whilst varying the remaining experts included in the ensemble; E^ℓ is the ensemble obtained relaxing the multimodality constraints.

- E_{post}^3 : it denotes the performance of the ensemble consisting of the unimodal models with the largest recall, i.e., AdaBoost, ResNet34, and VGG11-BN. In this case, we

adopt the subscript *post* to specify that such three models were *a posteriori* selected, i.e., they provide the largest performance on the test set, and not on the validation set;

- $\overline{E_{post}^{2+*}}$: it denotes the average performance attained by all the possible ensembles, including the two unimodal classifiers with the largest *a posteriori* recall, i.e., Adaboost and VGG11-BN, whilst varying the remaining experts included in the ensemble;
- E^{\dagger} : it denotes the performance of the ensemble obtained relaxing the multimodality constraints, and it is composed of AdaBoost, DT and RF.

It is worth noting that the performances returned by such four experiments are always lower than the performance of E^* and even lower than several unimodal learners. This confirms, again, that maximising recall and diversity together is a useful driver to guide the ensemble set-up, i.e., to select which are the unimodal learners to be included. Furthermore, the fact that such ensembles in some cases provide lower performance than some unimodal learners confirms that handcrafted ensemble definitions can lead to sub-optimal results.

The last section of Table 4.3 presents a direct comparison of our approach with two state-of-the-art studies [78, 79], which are denoted as DeepMMSA and MCF. As described in section 4.1, they work with clinical and imaging modalities so that we can apply them to our data, computing the same scores we used for the other architectures under consideration. Note also that we do not experimentally evaluate [77, 80] because, on the one side, [77] works with CT and PET images and it is not designed to handle clinical data, whereas on the other side, [80] does not work with CT images. The results show that such competitors perform worse than our method; we deem that this happens because such papers manually define the composition of the multimodal architectures, whilst our solution relies on an optimisation process.

As a further issue in our discussion, let us recall that three types of fusion exist in multimodal learning: early, joint, and late fusion. The latter is the one we used in this work, whilst the other two are other possible ways to proceed, which we considered as possible competitors for our method. To this end, we set up an early fusion learner using the best model per modality, i.e., AdaBoost and VGG11-BN, as already mentioned. Furthermore, we used the VGG11-BN as a feature extractor from the CT images, which we then concatenated with the clinical features and feed to the AdaBoost. We tested both slice-level and patient-level early fusion. The former consists of repeating the clinical features for each individual patient slice, whilst in the latter we averaged the CNN output of each slice to obtain a single feature vector per patient. These approaches got an accuracy equal to $62.92 \pm 8.31\%$ and $70.00 \pm 9.35\%$, and an F-score equal to $62.69 \pm 6.41\%$ and $69.83 \pm 12.25\%$, respectively, for the

slice- and patient-level fusion, which are lower than the proposed approach. We did not perform any joint fusion since the best unimodal model (AdaBoost) has larger performance than a fully connected network (MLP). This, in turn, makes it not possible to apply joint fusion between the adaptive boosted ensembles and the VGG11-BN, although this is an issue worthy of investigation in a future work.

4.5 Conclusions

In this manuscript, we proposed a multimodal method for survival analysis of NSCLC. NSCLC has been already studied in a few other works employing multimodal learning but, differently from the literature, we propose an algorithm able to identify the optimal set of classifiers to be added to the multimodal ensemble in a late fusion approach. Our study is based on two modalities, clinical and CT imaging data, of a cohort of 191 patients suffering from locally advanced non-small-cell lung cancer.

From a clinical point of view, the possibility of having prognosis prediction tools in addition to clinical data, and especially before starting treatment, represents an unmet need of particular interest. If this data are available at the start of therapy, the treatment itself could be modified, adapting it to the expected response, thus intensifying or descaling therapy in patients with poor or good prognosis, respectively.

Indeed, we presented an optimised late fusion ensemble search method that finds the optimal combination of multimodal models considering both a metric of performance and a diversity score. Experimental results show that our method outperforms conventional unimodal models, bringing significant increase in performance in the multimodal ensemble. Among the different combinations of classification algorithms, the proposed approach achieves an accuracy of 75.00%, an F-score of 77.70%, and a recall of 84.00%, achieved using a ResNet34 and a VGG11-BN for the imaging modality and a TABNET for the clinical modality. A limitation of our approach is the need to train all models before the optimal set can be selected, which certainly represents a high computational cost.

The results described so far suggest four future directions worthy of investigation:

- Retrieving data at 1-, 2-, and 3-year time points as well as the progression free survival, which would add useful information;
- Provide more complementary information by adding other modalities to improve performance, such as WSI, genome sequencing, etc.;
- Perform different multimodality fusion approaches, such as joint fusion to obtain a

end-to-end trainable system able to exploit the inherent correlations between multiple modalities;

- Search for an approach that a priori selects the models to be included in the ensemble, without the need to train them all individually;
- Switch from a classification to a regression task, which will allow predicting the actual survival time, also integrating the “Input doubling method” [106] as a preprocessing tool to augment the training set size.

In this chapter, we explored the feasibility and benefits of combining imaging and clinical data through multimodal learning strategies for survival prediction in NSCLC patients. The experimental results confirmed that multimodal models can outperform unimodal baselines, especially when the complementary information between modalities is effectively exploited. However, the proposed architectures assumed the availability of all modalities at both training and inference time, a condition that is rarely satisfied in real-world clinical scenarios. Missing data, including entire missing modalities, remains a major obstacle to the deployment of multimodal AI systems in healthcare. To address this challenge, the next chapter introduces MARIA, a novel transformer-based architecture specifically designed to enable robust multimodal fusion even in the presence of partially missing information.

Chapter 5

MARIA: a Multimodal Transformer Resilient to Missing Modalities

The deployment of multimodal models in clinical practice requires architectures that are not only accurate but also resilient to missing modalities. As highlighted in our systematic review of multimodal learning in healthcare [107], current strategies often struggle with real-world constraints such as partial observability and data heterogeneity. Inspired by the limitations observed in both literature and previous experiments, this chapter presents MARIA (Multimodal Attention Resilient to Incomplete data), a novel architecture specifically designed to address incomplete multimodal data scenarios. MARIA leverages intermediate fusion and a masked self-attention mechanism to selectively integrate information from available modalities while ignoring missing components during training and inference. We evaluate its performance across diverse settings, including “missing modality” and “all missing” scenarios, and benchmark it against state-of-the-art strategies for both data imputation and modality fusion. This work establishes a robust foundation for multimodal learning under the imperfect conditions typical of real-world clinical applications.

5.1 Introduction

In recent years, multimodal learning has emerged as a powerful approach for leveraging diverse data sources to achieve a comprehensive understanding of complex systems [108]. This is particularly relevant in domains such as healthcare, where integrating multiple data modalities, such as clinical assessments, imaging, laboratory tests, and patient histories, can significantly enhance diagnostic accuracy and treatment outcomes [109]. The human experience itself exemplifies multimodality, as it relies on diverse sensory inputs to form a

unified perception of the environment. Similarly, deep learning (DL) models have been developed to synthesize and analyze disparate data sources, thereby enhancing their predictive capabilities and enabling more informed decision-making in multifaceted domains such as healthcare [110, 111].

Despite the promise of multimodal learning, integrating multiple data sources presents unique challenges, with data incompleteness being one of the most significant [112]. Missing data is a common feature of real-world datasets, arising from issues such as sensor failures, patient non-compliance, technical limitations during data collection, or privacy restrictions [8]. Whether the missing information relates to features within a modality or the complete absence of a modality, such gaps can severely degrade the performance of machine learning models unless effectively addressed [113]. Thus, the development of multimodal learning models resilient to incomplete data is critical to ensuring reliability and robustness, especially in critical fields like healthcare [114].

To address these challenges, multimodal fusion strategies such as early fusion, late fusion, and intermediate fusion have been extensively studied [108, 111]. Early fusion, which combines features at the raw data level into a unified representation, is straightforward but highly susceptible to the effects of missing data, as it requires the availability of all feature vectors. Late fusion, which merges outputs from independently trained models, offers flexibility when modalities are missing but often fails to capture the intricate interactions across modalities. Intermediate fusion strikes a balance by integrating modality-specific features after initial processing. This forms a shared representation that enhances the ability to capture cross-modal dependencies, ultimately improving performance [107]. Therefore, especially in healthcare, it is essential to develop methods that leverage the potential of intermediate fusion while maintaining robustness to missing data [115].

The MARIA (Multimodal Attention Resilient to Incomplete data) model introduced in this work is designed to address the challenges of incomplete multimodal data. It builds upon our previous work, NAIM [116], by extending its masked attention mechanism to a multimodal scenario, where each modality is processed through a dedicated NAIM module and subsequently integrated via a shared transformer encoder. By employing an intermediate fusion strategy, MARIA combines modality-specific NAIM-based encoders with a shared attention-based encoder, enabling the model to effectively handle both missing features within each modality and entirely missing modalities without relying on data imputation. Unlike traditional methods that rely on data imputation to fill in missing entries, MARIA focuses exclusively on the available features, utilizing a modified masked self-attention mechanism to process observed information only, while completely masking out the contribution of missing information. This approach enhances both robustness and accuracy, while reducing

biases typically introduced by imputation techniques.

To summarize, the main contributions of this work are as follows:

- We propose MARIA, a novel multimodal transformer architecture specifically designed to handle incomplete tabular data without relying on imputation techniques.
- MARIA employs an intermediate fusion strategy with masked self-attention mechanisms at both the modality-specific and shared encoder levels, allowing the model to dynamically ignore missing features and modalities.
- We introduce regularization strategies tailored for missing data, including stochastic masking at both the feature and modality levels during training, to enhance generalization under incomplete conditions.
- We conduct an extensive benchmark comparison with 32 competing configurations across 8 diagnostic and prognostic tasks using two real-world clinical datasets, demonstrating MARIA’s robustness and superior performance under various missingness scenarios.

The manuscript is organized as follows: Section 5.2 reviews related work on multimodal learning and data handling methods; Section 5.3 introduces the MARIA model and its architecture; Section 5.4 explains the experimental setup and evaluation methodology; Section 5.5 presents the obtained results, comparing MARIA with other models under various missing data conditions; Section 5.6 summarizes the key findings and suggests directions for future research.

5.2 State of the Art

Multimodal learning combines information from diverse data sources to achieve a more comprehensive understanding of complex systems. This mirrors the inherently multimodal nature of human perception: we rely on multiple senses, e.g., sight, sound, and touch, to develop a complete understanding of our environment. Similarly, DL models must be designed to integrate diverse data sources to better comprehend intricate systems. This is particularly relevant in healthcare, where clinicians utilize multimodal data, including patient histories, imaging, laboratory results, and physical examinations, to make informed decisions. By effectively integrating such diverse information, multimodal learning models can enhance decision-making and predictive accuracy, leading to improved diagnostic outcomes and more effective treatment plans [115, 117, 118, 119, 120].

However, one of the primary challenges in multimodal learning is handling missing data, which frequently arises due to factors such as sensor failures, survey non-responses, or technical issues during data collection [121]. Effectively managing missing data, whether it involves incomplete features within a modality or entirely absent modalities, is critical to ensuring the reliability and robustness of multimodal models.

Multimodal fusion techniques play a vital role in successfully integrating diverse data sources. These techniques are typically categorized into three main strategies: early fusion, late fusion, and intermediate fusion (Figure 5.1). Each approach has distinct characteristics, making it suitable for different scenarios [115, 122].

Early fusion integrates features at the raw data level, combining them into a unified representation before any significant processing. Late fusion merges outputs from independently trained models at the decision level, offering flexibility when dealing with missing modalities. Intermediate fusion takes a balanced approach, integrating modality-specific features after initial processing to create a shared representation. Each of these fusion strategies has specific advantages and limitations in terms of performance, computational complexity, and their ability to manage missing data. This is particularly significant in healthcare, where data quality and completeness often vary. The subsequent sections provide a detailed analysis of these fusion techniques, focusing on their applications and limitations in healthcare scenarios.

5.2.1 Early Fusion

Early fusion, as illustrated in Figure 5.1.a, involves integrating multiple data modalities at the feature level. In this approach, the raw features X_i from each modality are concatenated (\oplus in the figure) to form a single feature vector, which is then fed into the learning model. This method facilitates the early combination of information from all available modalities, making it particularly advantageous when different data sources are highly complementary. Early fusion is conceptually straightforward and often enables the learning model to effectively exploit cross-modal correlations [115, 122].

However, early fusion presents several challenges, especially when dealing with incomplete data. Because this approach relies on the availability of all feature vectors, missing data from even a single modality can significantly degrade model performance. Imputation is a common strategy for addressing these gaps, but it introduces potential risks such as bias and information loss [123]. Additionally, early fusion typically requires extensive preprocessing to harmonize features from different modalities, which often vary in scale and distribution.

In healthcare, early fusion can be particularly beneficial when modalities are guaran-

teed to be complete or when missing data is minimal and can be addressed through robust preprocessing techniques. However, given the variability and incompleteness commonly encountered in real-world medical datasets, early fusion may struggle to perform effectively without sophisticated data-handling strategies [115, 122].

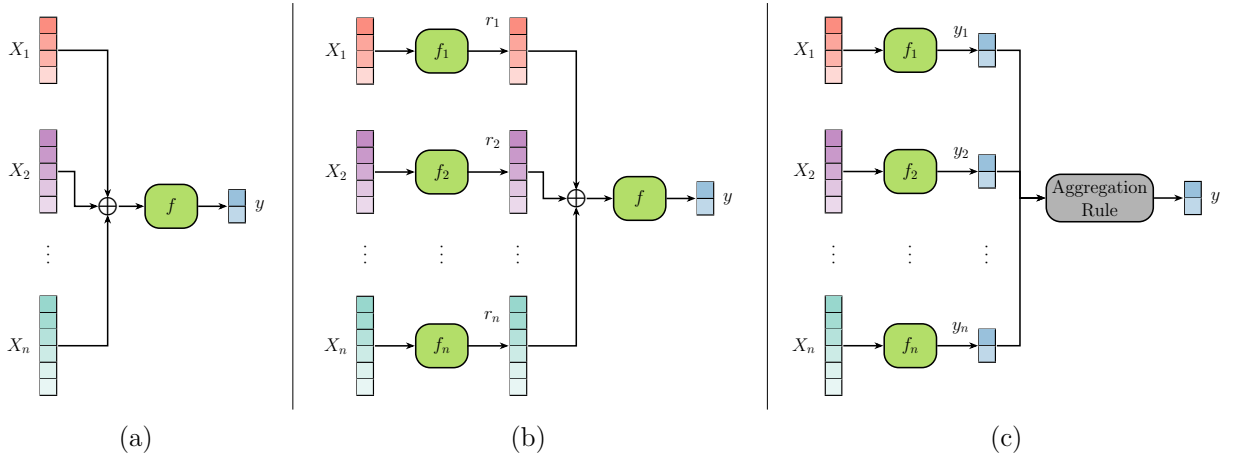


Figure 5.1: Overview of multimodal fusion strategies in DL: (a) Early fusion; (b) Intermediate fusion; (c) Late fusion. X_1, \dots, X_n represent the modalities; $f_{(i)}$ is a generic function representing a module or an entire model; r_1, \dots, r_n stand for the latent representations of the modalities; $y_{(i)}$ indicates an output of a model; \oplus represent a fusion, e.g. concatenation or average, of the inputs.

5.2.2 Late Fusion

Late fusion, in contrast to early fusion, integrates modalities at the decision level, as illustrated in Figure 5.1.c. In this approach, separate models f_i are trained for each modality X_i , and their predictions y_i are subsequently aggregated using a predefined rule to generate the final output y . This method is particularly advantageous when the modalities differ significantly in data type or exhibit varying levels of reliability. By training separate models, late fusion allows each modality to be optimally utilized before combining their outputs.

One of the primary advantages of late fusion is its flexibility in handling missing modalities. Since each model operates independently, missing data from one modality does not prevent predictions from being made using the available modalities. However, late fusion has a notable limitation: it fails to fully exploit cross-modal interactions. Because modalities are only integrated at the decision level using a static, predefined rule, rather than a dynamically learned process, potentially rich correlations between features from different data sources may remain untapped. This drawback makes late fusion less suitable for tasks that require deep integration of modality-specific features, particularly in scenarios demanding

high levels of diagnostic precision.

5.2.3 Intermediate Fusion

Intermediate fusion, also known as joint or hybrid fusion, strikes a balance between early and late fusion. This approach combines modality-specific features at an intermediate stage of the learning process (Figure 5.1.b), typically after each modality has undergone initial, independent processing. In this setup, modality-specific modules f_i generate latent representations r_i for each modality. These representations are then fused using a defined technique (denoted by \oplus in Figure 5.1.b) to form a shared representation r_{sh} . Finally, this shared representation is processed by a final module f to produce the desired output y . This approach facilitates a richer integration of features, retaining modality-specific information while capturing inter-modal relationships during the feature extraction phase.

One of the major advantages of intermediate fusion is its ability to handle incomplete data more flexibly and effectively. Various techniques are available for fusing latent representations from different modalities, and we recommend readers refer to the review [107] for an in-depth exploration of these methods. However, intermediate fusion comes with certain challenges, including increased computational complexity and training difficulty. The model must learn both unimodal and multimodal representations simultaneously, requiring significant computational and data resources. These demands can pose a barrier in resource-constrained environments, such as many healthcare settings.

Despite these challenges, the dynamic nature of intermediate fusion offers significant advantages. By allowing the model to learn how to fuse information from different sources dynamically, it enhances robustness and adaptability. This integration of complementary information from multiple modalities enables the model to leverage the strengths of each data source while mitigating their individual weaknesses. Such adaptability is particularly valuable in real-world healthcare scenarios, where the quality and availability of data often vary across modalities. Intermediate fusion's ability to handle partially missing or noisy modalities can result in more reliable predictions. Moreover, the shared representation created through intermediate fusion fosters a deeper understanding of correlations between different data types, which is critical for complex tasks such as medical diagnosis and prognosis [115, 122].

Despite its resource demands, intermediate fusion represents a promising direction for the development of DL models that are both comprehensive and resilient. This makes it a powerful approach for enhancing decision-making in healthcare environments [107, 115, 122].

5.2.4 Handling Incomplete Data

Real-world multimodal data are often imperfect due to missing features or modalities. Therefore, there is a pressing need for multimodal models robust in the presence of incomplete data. Missing data, whether involving individual features or entire modalities, is a common challenge across various fields and is often caused by issues, such as human error, survey non-responses, data corruption, or systematic loss. Traditional approaches to address missing data typically rely on imputation techniques, which attempt to fill these gaps but can introduce biases or fail to capture underlying complexities. For instance, we employed the k-Nearest Neighbors (kNN) imputer, which has demonstrated effectiveness in handling missing values in tabular data [123, 124, 116]. Additionally, we tested the Missing in Attributes (MIA) strategy, used by tree-based models to dynamically manage missing features without requiring imputation.

Healthcare settings are particularly vulnerable to the problem of incomplete data, as patients may follow different treatment plans, discontinue care for reasons such as transferring facilities, voluntarily ceasing treatment, or even passing away. Moreover, privacy concerns further exacerbate data incompleteness in these settings [124]. Many methods simply exclude patients with missing values, which can significantly reduce data availability and compromise the robustness of analyses. Other approaches often involve imputing missing information using data from available modalities for the same subject or from other patients with similar characteristics.

Several advanced methodologies have been proposed to address the issue of missing data in multimodal contexts:

The Contrastive Masked-Attention Model integrates a Generative Adversarial Network (GAN)-based augmentation mechanism to synthesize data for missing modalities and employs contrastive learning to enhance cross-modal representations. Masked attention ensures that only interactions between observed modalities are captured, thereby minimizing the introduction of extraneous noise [125].

The Cascaded Multi-Modal Mixing Transformers implement a cascaded cross-attention architecture to effectively integrate multiple available modalities, enabling robust performance even when some modalities are missing. This approach offers flexibility and adaptability in fusing modality-specific information [126].

The Missing Modalities in Multimodal healthcare framework employs task-guided, modality-adaptive similarity metrics to identify similar patients and impute missing data. By leveraging auxiliary information from comparable patients, this method preserves the underlying relationships in multimodal healthcare data [127].

Shared-Specific Feature Modeling disentangles shared features from modality-specific ones, enabling efficient handling of missing data during both training and inference. By learning shared features across all available modalities, this approach ensures the retention of essential representations while maintaining model performance [128].

The Severely Missing Modality model uses a Bayesian meta-learning framework to approximate latent representations for missing modalities. This method is designed to handle incomplete data during both training and testing, offering robust generalization capabilities even when data availability is severely limited [121].

These methodologies highlight diverse strategies for compensating for missing information, including identifying shared latent representations, generating synthetic data, and leveraging auxiliary patient information. By reducing dependence on complete multimodal datasets, these approaches improve the practicality of multimodal models in real-world clinical and resource-constrained settings. However, both traditional and DL-based approaches share a common limitation: they rely on artificially filling data gaps, which can introduce bias and compromise task accuracy.

To address this limitation, we propose a model that exclusively utilizes the available features and modalities, avoiding the generation of synthetic data. By focusing solely on effectively leveraging observed information, our approach aims to enhance robustness and reliability, even in scenarios with severely missing data.

5.3 Methods

In this work, we propose MARIA (Multimodal Attention Resilient to Incomplete data), a multimodal model specifically designed to address the challenges of incomplete features and modalities in multimodal healthcare data. The model effectively integrates data modalities that may be incomplete or entirely absent, offering a robust solution for predictive analysis without relying on traditional data imputation techniques or synthetic data generation.

This section first provides an overview of the MARIA model design. We then focus on the architecture of the modality-specific encoders, the strategies employed for handling missing data, and the regularization techniques implemented during training to enhance generalizability under incomplete input conditions.

5.3.1 Model

MARIA is specifically designed to be resilient to incomplete data and modalities without relying on imputation. It employs intermediate fusion, using modality-specific encoders and

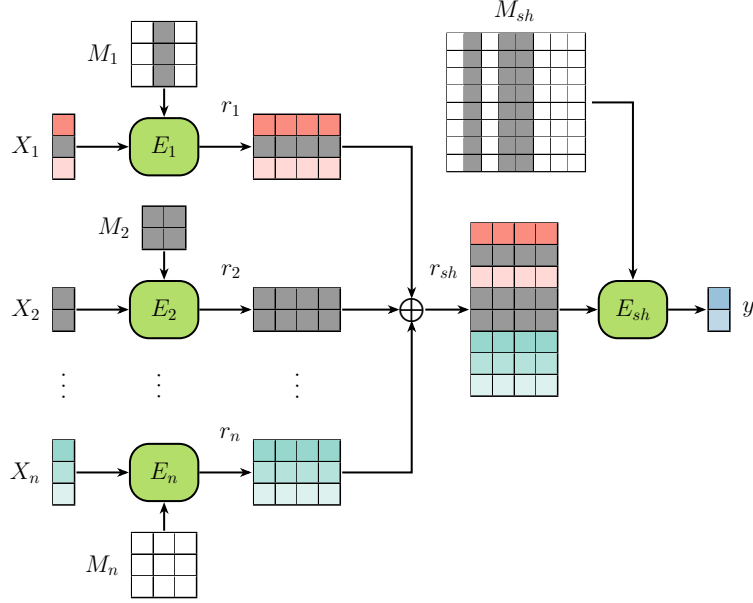


Figure 5.2: MARIA architecture. Each modality-specific encoder E_i takes the modality X_i as input to generate the latent representation r_i . Then, the shared encoder E_{sh} elaborates the concatenation of the latent representations r_{sh} to get the final output y . In the figure, a gray square represents a missing feature or its respective element in the masking matrix.

a shared encoder with a masked self-attention mechanism to combine latent representations while effectively managing missing data. The architecture incorporates multiple modality-specific encoders for each data modality. In this work, we focus on multimodal problems where the modalities represent tabular data describing various aspects of a patients condition. Thus, MARIA utilizes separate NAIM [116] modules as modality-specific encoders (Figure 5.2). These encoders integrate a modified masked multi-head attention mechanism that selectively focuses on available features within each modality while completely ignoring missing ones. This mechanism ensures robustness by excluding absent features from attention computations.

Each tabular modality X_i , where $i = 1, \dots, n$, is encoded into embeddings using look-up tables [116], which represent missing entries with a specific non-trainable embedding. The modality embeddings are then processed by their corresponding encoders E_i , which compute query, key, and value matrices, denoted as Q_i , K_i , and V_i , using linear transformations:

$$\begin{cases} Q_i = X_i \cdot W_i^Q \\ K_i = X_i \cdot W_i^K \\ V_i = X_i \cdot W_i^V \end{cases} \quad \begin{cases} d^h = d^e/h \\ W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d^e \times d^h} \\ Q_i, K_i, V_i \in \mathbb{R}^{|X_i| \times d^h} \end{cases} \quad (5.1)$$

where W_i^Q , W_i^K and W_i^V are learnable weight matrices. These transformations reduce di-

dimensionality to d^h , determined by the token dimensions d^e and the number of heads h in the model. Next, a modified masked self-attention mechanism is applied:

$$MSA(Q_i, K_i, V_i) = ReLU \left(softmax \left(\frac{Q_i K_i^T}{\sqrt{d^h}} + M_i \right) + M_i^T \right) V_i \quad (5.2)$$

where the masking matrix M_i ensures that missing features do not influence the latent representation r_i . The elements of M_i are defined as follows:

$$M_i^{kj} = \begin{cases} -\infty & \text{if } X_i^j \text{ is missing} \\ 0 & \text{otherwise} \end{cases}, \quad M_i = \begin{bmatrix} 0 & -\infty & \dots & 0 & 0 \\ 0 & -\infty & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & -\infty & \dots & 0 & 0 \\ 0 & -\infty & \dots & 0 & 0 \end{bmatrix} \quad (5.3)$$

This operation effectively zeroes out weights associated with missing features after applying *softmax* and *ReLU*. Each modality-specific encoder E_i generates a latent representation r_i , of dimensions $|X_i| \times d^e$, where $|X_i|$ represents the number of tokens in the modality, e.g., the number of features in the i -th modality. These latent representations are then concatenated to form a joint representation r_{sh} , composed only of available information, with null vectors representing missing features. This multimodal representation is passed to the shared encoder E_{sh} , which computes its own query, key, and value matrices, denoted as Q_{sh} , K_{sh} , and V_{sh} , as follows:

$$\begin{cases} Q_{sh} = r_{sh} \cdot W_{sh}^Q & W_{sh}^Q, W_{sh}^K, W_{sh}^V \in \mathbb{R}^{d^e \times d^h} \\ K_{sh} = r_{sh} \cdot W_{sh}^K & Q_{sh}, K_{sh}, V_{sh} \in \mathbb{R}^{\sum_i |X_i| \times d^h} \\ V_{sh} = r_{sh} \cdot W_{sh}^V \end{cases} \quad (5.4)$$

where W_{sh}^Q , W_{sh}^K and W_{sh}^V are weights matrices learned during training. As with the modality-specific encoders, dimensionality is reduced to d^h through linear transformations.

A similar modified masked self-attention mechanism is then applied:

$$MSA(Q_{sh}, K_{sh}, V_{sh}) = ReLU \left(softmax \left(\frac{Q_{sh} K_{sh}^T}{\sqrt{d^h}} + M_{sh} \right) + M_{sh}^T \right) V_{sh} \quad (5.5)$$

where the masking matrix M_{sh} ensures that missing modalities do not impact the final shared representation. This matrix operates in the same manner as M_i , zeroing out weights associated with missing modalities. The elements of the masking matrix are defined as

follows:

$$M_{sh}^{kj} = \begin{cases} -\infty & \text{if } r_{sh}^j \text{ is missing} \\ 0 & \text{otherwise} \end{cases}, \quad M_{sh} = \begin{bmatrix} 0 & -\infty & \dots & 0 & 0 \\ 0 & -\infty & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & -\infty & \dots & 0 & 0 \\ 0 & -\infty & \dots & 0 & 0 \end{bmatrix} \quad (5.6)$$

M_{sh} sums the $-\infty$ values to weights that need to be ignored, effectively zeroing them after applying the *softmax* and *ReLU* operations.

The masking matrices M_i and M_{sh} play a crucial role in MARIA's ability to handle incomplete data. These matrices are incorporated into the attention mechanism to explicitly prevent the model from attending to missing values. Specifically, the matrix M_i is applied within each modality-specific encoder to suppress the influence of missing features. Similarly, the matrix M_{sh} is applied in the shared encoder to control attention across modalities. This dual-masking mechanism allows MARIA to dynamically adapt its attention based on the availability of data, without resorting to imputation or synthetic representations.

Finally, the joint representation is passed through a fully connected layer, which predicts the output y . The training process minimizes prediction error, updating the weights of both the modality-specific encoders (E_i) and the shared encoder (E_{sh}) via end-to-end backpropagation.

This architecture qualifies MARIA as an intermediate fusion model, performing fusion at the latent representation level. By dynamically optimizing all encoders, MARIA balances the contributions of different modalities and adapts to maximize their utility during training and inference. The use of a masked multi-head attention mechanism ensures the model focuses adaptively on informative parts of the input, completely ignoring missing data. This approach allows each modality to contribute based on its completeness, resulting in accurate and reliable multimodal representations.

5.3.2 Regularization Technique for Missing Data

To enhance the model's generalizability under incomplete input conditions, we employ regularization strategies during training that improve its resilience to varying degrees of data incompleteness. These strategies ensure that even when some modalities or features are unavailable, the model can still generate accurate and meaningful outputs [116, 126]. During training, we simulate a relaxed missing data setting where each modality or feature is treated as potentially missing, while maintaining at least one available data point per patient. This

approach allows the model to learn how to handle different levels of missingness effectively, making it particularly well-suited to the variability typically found in clinical datasets. By encouraging the model to extract meaningful representations from each available modality, these masking strategies promote robustness against incomplete information during both training and inference.

Modality Dropout During training, the model uses a stochastic masking procedure to simulate incomplete data scenarios. Given a sample $X = \{X_1, \dots, X_n\}$, where n represents the number of modalities, let $v_m \leq n$ denote the number of non-missing modalities in the sample (where $v_m = n$ for fully populated data or $v_m < n$ for partially missing data). A binary decision variable determines whether masking will be applied to the sample X . If the sample is selected for masking, a random count c_m of modalities to mask is chosen uniformly from the set $\{1, 2, \dots, v_m - 1\}$, ensuring that at least one modality remains unmasked. Finally, c_m non-missing modalities are randomly selected, and their values are set to *missing*, producing the augmented sample.

Feature Dropout Similarly, when a tabular modality X_i is set as present, a similar stochastic masking procedure is applied at the feature level. A binary decision variable determines whether masking will be applied to the features of the modality X_i . If masking is applied, a random count c_i of features to mask is chosen uniformly from the set $\{1, 2, \dots, v_i - 1\}$, where v_i is the number of non-missing features of the modality X_i . This ensures that at least one feature remains unmasked. Finally, c_i non-missing features within X_i are randomly chosen and set to *missing*, resulting in the augmented modality.

5.4 Experimental Configuration

In this section, we first describe the datasets used in the experiments and the preprocessing applied (Section 5.4.1). We then outline the combinations of models and imputers employed as competitors (Section 5.4.2). Finally, we present the metrics used for evaluation (Section 5.4.5).

5.4.1 Data

We evaluated MARIA and its competitor models on two publicly available datasets across eight diagnostic and prognostic tasks (details in Table 5.1). These tabular datasets represent real-world scenarios where patient data is often incomplete, highlighting the need for methods that are resilient to missing information.

Dataset	Task	# of samples	Class distribution		Modalities Info (% of missing features and modalities)
ADNI [129]	Diagnosis Binary	953	CN: 542	AD: 411	Assessment: 37 (f: 35% - m: 0%) Biospecimen: 47 (f: 57% - m: 5%) Image Analysis: 14 (f: 38% - m: 1%) Subject Characteristics: 17 (f: 37% - m: 0%)
	Diagnosis Multiclass	2066	CN: 542 LMCI: 690	EMCI: 423 AD: 411	Assessment: 37 (f: 32% - m: 0%) Biospecimen: 47 (f: 56% - m: 7%) Image Analysis: 14 (f: 36% - m: 1%) Subject Characteristics: 17 (f: 38% - m: 0%)
	Prognosis 12 months	1340	CN: 427 MCI: 797	Dementia: 116	Assessment: 126 (f: 33% - m: 0%) Biospecimen: 51 (f: 47% - m: 0%) Image Analysis: 18 (f: 27% - m: 0%) Subject Characteristics: 21 (f: 29% - m: 0%)
	Prognosis 24 months	1159	CN: 428 MCI: 535	Dementia: 196	Assessment: 126 (f: 34% - m: 0%) Biospecimen: 51 (f: 46% - m: 0%) Image Analysis: 18 (f: 27% - m: 0%) Subject Characteristics: 21 (f: 28% - m: 0%)
	Prognosis 36 months	856	CN: 239 MCI: 420	Dementia: 197	Assessment: 126 (f: 27% - m: 0%) Biospecimen: 51 (f: 37% - m: 0%) Image Analysis: 18 (f: 27% - m: 0%) Subject Characteristics: 21 (f: 23% - m: 0%)
	Prognosis 48 months	693	CN: 269 MCI: 280	Dementia: 144	Assessment: 126 (f: 36% - m: 0%) Biospecimen: 51 (f: 47% - m: 0%) Image Analysis: 18 (f: 28% - m: 0%) Subject Characteristics: 21 (f: 30% - m: 0%)
AIforCOVID [130]	Mild/Severe	1585	Mild: 839	Severe: 746	Blood Analysis: 14 (f: 31% - m: 1%) Clinical History: 13 (f: 24% - m: 9%) Personal Info: 2 (f: 0% - m: 0%) Admission State: 5 (f: 8% - m: 0%)
	Death	1585	Censored: 1336	Uncensored: 249	

Table 5.1: Datasets’ details consisting of: dataset name and reference, the task name, the number of samples, the classes’ distribution, and the different tabular modalities, with the corresponding number of features and the respective rates of original missing features (f) and modalities (m).

The first dataset was obtained from the Alzheimers Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu) [129]. The ADNI was launched in 2003 as a public-private partnership, led by Principal Investigator Michael W. Weiner, MD. The primary goal of ADNI has been to test whether serial magnetic resonance imaging, positron emission tomography, other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimers disease (AD), with respect to the Cognitively Normal (CN) group. For our study, we used tabular data derived from four baseline modalities: Assessment (cognitive and neuropsychological scores), Biospecimen (CSF, ApoE genotyping, and lab data), Image analysis (MRI and PET neuroimaging biomarkers), and Subject Characteristics (family history, demographics). We analyzed data from ADNI 1, GO, 2, and 3 phases. Diagnostic tasks included binary classification (CN vs. AD) and ternary classification (CN vs. AD vs. MCI), reflecting real clinical differentiation scenarios. Additionally, prognostic tasks aimed to predict whether treatment intervention might be necessary at specific future points (12, 24, 36, and 48 months post-recruitment). These tasks classified patients into CN, MCI, and

Dementia categories. Table 5.1 provides details on patient distributions for both baseline and follow-up classifications.

The second dataset, AIforCOVID [130], contains clinical data from six Italian hospitals, collected during the first wave of the COVID-19 pandemic (March-June 2020). Data was recorded at the time of hospitalization of symptomatic patients and subsequently anonymized and reviewed. All patients tested positive for SARS-CoV-2 via RT-PCR, with 5% confirmed only after a second test. Each patient was classified as either mild (discharged or hospitalized without ventilatory support) or severe (requiring non-invasive ventilation, ICU care, or deceased). Additionally, we evaluated the proposed approach on a death prediction task, classifying patients as either censored (alive) or uncensored (deceased).

MARIA’s performance across diverse medical tasks provides valuable insights into its resilience and adaptability in real-world healthcare scenarios. By leveraging the ADNI and AIforCOVID datasets, we demonstrate the model’s ability to handle challenges such as missing modalities and heterogeneous data distributions.

These experiments emphasize the importance of multimodal fusion techniques that are not only robust to missing data but also capable of learning from complex, interrelated medical datasets. Such advancements provide the foundation for more resilient and flexible DL models in healthcare, ultimately supporting clinicians in making better-informed decisions under real-world constraints.

5.4.2 Competitors

We conduct an extensive comparison of our methodology against early, late, and intermediate fusion approaches that use missing data imputation as a preprocessing step before model training. Additionally, we benchmark against tree-based models that manage missing values using the Missing In Attributes (MIA) strategy. We choose not to include generative approaches, as these are primarily developed for imaging modalities rather than tabular data, and they introduce additional complexity and randomness, which can hinder reproducibility. Instead, we focus on interpretable, efficient approaches widely adopted in clinical settings, where model simplicity and reliability are crucial. Our analysis includes a total of 10 distinct competitor models, each combined with the k-Nearest Neighbors (kNN) imputation technique, as this method outperformed other imputation strategies in prior studies [123, 116, 1]. For all models, we used default hyperparameters. Indeed, as shown in [52], in line with the “No Free Lunch” theorem, the authors empirically observe that, in many cases, tuning hyperparameters does not lead to significantly better performance compared to the default values suggested in the literature.

Table 5.2 provides an overview of the competitors. The first column categorizes models as ML or DL approaches, the second specifies the base learners, the subsequent columns indicate the techniques used for handling missing data, and the final columns describe the fusion strategies employed. This results in 32 competitor configurations, each marked by an “×” in the respective columns.

Type	Model	Imputer		Multimodal Strategies		
		With	Without	Early	Intermediate	Late
Machine Learning	AdaBoost	×		×		×
	Decision Tree	×	×	×		×
	HistGradientBoost	×	×	×		×
	Random Forest	×	×	×		×
	SVM	×		×		×
	XGBoost	×	×	×		×
Deep Learning	MLP	×		×	×	×
	TabNet	×		×	×	×
	TabTransformer	×		×	×	×
	FTTransformer	×		×	×	×

Table 5.2: Combinations of models, missing techniques and fusion strategies used as competitors, represented by an “×”, in the experiments.

To thoroughly evaluate the performance of our proposed methodology, we designed experiments comparing it against a diverse set of baseline models. These experiments involve both ML models, which may rely on imputation or employ the MIA strategy, and DL models paired with imputation techniques.

The ML models include AdaBoost, Decision Trees, HistGradientBoost, Random Forests, Support Vector Machines (SVM), and XGBoost. AdaBoost [21] is a cascading ensemble model that prioritizes hard-to-classify instances, offering robustness across diverse datasets. Decision Trees [22] are highly interpretable models that visually represent decision-making processes, providing insights into complex data relationships. HistGradientBoost [1] offers an efficient variation of gradient boosting, optimized for handling large datasets with improved speed and memory usage. Random Forests [23] is an ensemble of decision trees known for robustness against overfitting and enhanced reliability. SVM [24], equipped with an RBF kernel, is included for its versatility in handling non-linear data separations, providing a contrast to tree-based models. XGBoost [25], an advanced boosting model, employs a gradient descent procedure to minimize loss and is highly effective for tabular datasets.

Additionally, we evaluate DL models paired with imputation methods. These approaches include Multilayer Perceptron (MLP), TabNet, TabTransformer, and FTTransformer. MLP [26] is a foundational DL model that captures complex relationships between features, offering a

baseline for comparison. TabNet [5] leverages self-attention to dynamically select features, improving interpretability and decision-making. TabTransformer [6] uses transformer-based self-attention mechanisms to embed categorical features and capture complex inter-feature relationships. FTTransformer [7] further explores transformers’ potential, using distinct embedding strategies for numerical and categorical features.

The selected DL models were evaluated using both early and late fusion approaches. In addition, we developed intermediate fusion variations of these models, where the respective architectures were employed for both modality-specific encoders and shared encoder settings. These intermediate fusion configurations assess the models’ ability to concurrently handle multiple input types, leveraging the shared encoder to effectively combine information from various modalities.

These experiments were designed to comprehensively assess the strengths and limitations of each competitor model across various settings. Our comparisons aim to benchmark the performance of our intermediate fusion methodology against both traditional ML approaches and advanced DL competitors. By exploring a wide range of techniques, we highlight the effectiveness of our approach in handling incomplete and heterogeneous multimodal healthcare data.

5.4.3 Preprocessing

For each dataset, we normalize the numerical features using a Min-Max scaler and apply one-hot encoding to the categorical features before feeding them into the models. However, for models such as MARIA, NAIM, HistGradientBoost, TabNet, TabTransformer, FTTransformer, and XGBoost, one-hot encoding is not applied, as their implementations can directly handle categorical features. The preprocessing steps are calibrated using the training data and then applied to the validation and testing sets.

5.4.4 Missingness Evaluation

Our experiments center on generating *Missing Completely At Random* (MCAR) values artificially as it represents the broadest class of missing data type without the introduction of any bias. Our goal is to test our model under diverse missing data conditions by introducing missing values and modalities at various rates, denoted as p , across both the training and testing sets. Specifically, we generate separate missing rates for the training and test sets, set to 0%, 5%, 10%, 30%, 50%, and 75%. No additional missing values were introduced if the generated rate was lower than the dataset’s pre-existing level of missingness, resulting in a variable number of experiments per dataset.

Moreover, we performed two types of experiments to evaluate our model under different missingness scenarios:

- **Missing Modalities:** where entire modalities for each patient within the dataset are masked, simulating scenarios where certain data sources were absent.
- **All Missing:** where a certain percentage of individual elements across the entire dataset are masked, thereby affecting all modalities simultaneously and eventually obtaining both missing features and modalities.

This approach allows us to explore a wide range of data completeness scenarios, starting from the dataset’s original missing rate (Ω) to extreme cases where up to 75% of features or modalities are missing. Given a target missing rate p , the total number of samples N in the dataset, and the number of elements per sample (either the number of features $|X_i|$ or the number of modalities n), the total number of values to be masked is calculated as $N \cdot |X_i| \cdot p$ or $N \cdot n \cdot p$, respectively. This calculation also takes into account any pre-existing missing entries m_j for each sample j , representing either the number of missing features or missing modalities. Thus, the adjusted number of values to be masked was computed as $N \cdot |X_i| \cdot p - \sum_j m_j$ or $N \cdot n \cdot p - \sum_j m_j$. We then generated a random masking matrix of dimensions $N \times |X_i|$ or $N \times n$ to reflect the structure of the dataset, ensuring that at least one value in any fully masked row or column was replaced to avoid complete data loss in a specific dimension. As a result, the samples and features exhibited varying degrees of missingness, all conforming to the MCAR framework.

5.4.5 Evaluation Metrics

To evaluate the models, each dataset is divided into five stratified cross-validation splits, ensuring the original class distribution is preserved. Within each fold, 20% of the training samples is reserved for validation. The performance of each experiment is assessed by averaging the values of the Area Under the Receiver Operating Characteristic Curve (AUC) or the Matthews Correlation Coefficient (MCC) computed across the cross-validation folds.

The AUC is a widely used metric for evaluating classification tasks, as it measures the model’s ability to distinguish between positive and negative instances across all possible thresholds. This makes it a comprehensive indicator of performance, especially for tasks with balanced or slightly imbalanced class distributions. We employed the AUC for tasks such as the ADNI diagnostic and AIforCOVID prognostic evaluations, where the class distributions were relatively balanced. The AUC effectively captures the trade-off between true positive and false positive rates, providing a robust measure of classifier performance.

For highly imbalanced datasets, the AUC may not reliably reflect model performance, as it can be disproportionately influenced by the majority class. In such cases, we use the MCC, which accounts for all four categories of the confusion matrix (true positives, false positives, true negatives, and false negatives). The MCC provides a more balanced view, making it particularly suitable for scenarios with significantly unbalanced positive and negative classes. We employ the MCC for tasks such as the ADNI prognostic evaluations and the AIforCOVID death prediction task, where class imbalance is pronounced.

5.4.6 Fusion Analysis

As an additional analysis of our proposed approach, we also compared MARIA with NAIM [116], our unimodal baseline model, which serves as the foundation for our intermediate fusion methodology. This comparison illustrates the progression from unimodal analysis to the more sophisticated fusion mechanism that supports our proposed methodology. By conducting this comparative analysis, we provide deeper insights into how our approach builds on and improves upon the unimodal baseline, emphasizing the advantages and advancements made through the intermediate fusion technique.

Specifically, we evaluated our approach against the NAIM model under both early and late fusion configurations, which are representative of different strategies for leveraging multimodal information. In the early fusion scenario, we concatenated the features from all modalities before inputting them into the model, effectively treating all available information as a unified input space. This allowed us to explore the interactions between different modalities at an early stage of the modeling process. In the late fusion scenario, we train separate models for each modality, subsequently combining their predictions by averaging the corresponding decision profiles. This approach enables each modality to be independently modeled, allowing the strengths of each individual modality to contribute to the final decision in an aggregated manner.

The results of this extended comparison help demonstrate the value of our proposed fusion strategies and their ability to extract meaningful insights from multimodal data. By examining the performance differences between early and late fusion configurations, we better understand the unique strengths of each strategy and how our intermediate fusion approach effectively balances them. This balance allows MARIA to extract meaningful insights from multimodal data, achieving optimal performance by leveraging both modality-specific information and cross-modal interactions.

5.5 Results and Discussions

As described in the previous section, we compare MARIA with 32 leading competitor models across both ML and DL approaches for tabular data. We evaluated performance under two distinct experimental configurations: missing modalities and all missing. Each configuration involved 36 combinations of missing value percentages in the training and testing sets across 8 tasks, resulting in a total of 18432 experiments (2304 per task). However, due to the pre-existing missing rates in the datasets, the final results were derived from 12192 experiments across all tasks.

To analyze the results, we grouped the competitors into categories for clearer comparisons and visualized the average performance achieved during five-fold cross-validation. The results are presented in eight separate plots (Figures 5.3, 5.4, 5.5, 5.6), each corresponding to a specific task. Each task-specific plot includes charts, one for each level of missing data in the training set, separately showing performance metrics (y-axis) as the missing rate in the testing set increased (x-axis). Note that the number of charts and points along the x-axis may vary due to the initial missing rates (Ω). These charts allow for a detailed comparison of MARIA against different groups of competitors, specifically ML (Figures 5.3 and 5.4) and DL models (Figures 5.5 and 5.6).

In addition to the aggregate performance visualizations, we conducted a statistical significance analysis using the Wilcoxon signed-rank test to assess whether the differences in prediction performance between MARIA and its competitors were statistically meaningful. This analysis was carried out across all experimental conditions, comparing the binary correctness vectors (i.e., correct vs. incorrect predictions) on the concatenated test sets of the five folds, i.e., the entire dataset. The number of experiments in which MARIA significantly outperformed or underperformed each competitor was computed and summarized. The results, presented in tables G.1 and G.2 in Appendix G, further support the robustness of MARIA, particularly in scenarios involving missing modalities, where MARIA outperforms all the competitors in at least 40% of the experiments, while underperforming in no more than 2%. In the more challenging all missing scenario, MARIA still achieves statistically significant wins in at least 30% of the experiments and loses in no more than 10%, confirming its ability to handle severe data incompleteness.

5.5.1 MARIA vs. ML

As an initial analysis, we compared MARIA with traditional ML approaches, as depicted in Figures 5.3 and 5.4. These figures present the average performance of the models under

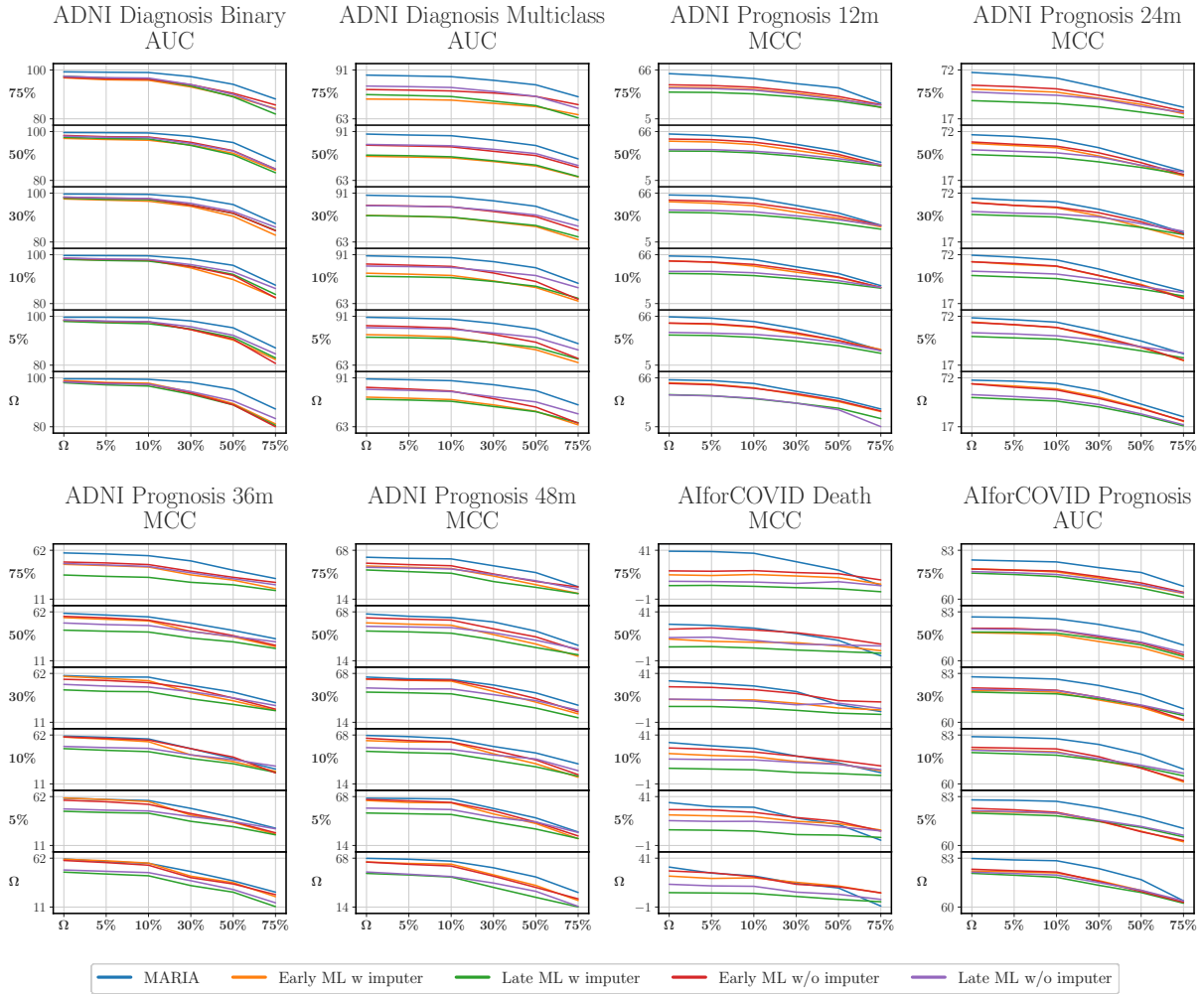


Figure 5.3: MARIA vs. ML in the “missing modalities” scenario. Each plot, one for each task, reports different charts, one for each of the missing rates in the training set, showing how the performance (y-axis) changes as the missing rate in the test set increases.

investigation, categorized by their use of either an imputer or the MIA strategy (without an imputer) and further grouped into early and late fusion approaches. Notably, MARIA consistently demonstrates superior average performance across all levels of missing data in both the training and testing phases, maintaining its advantage even under ideal conditions where no additional missing data is introduced.

This consistent superiority highlights not only MARIA’s distinct advantages over traditional ML models but also underscores the largely unexplored potential of DL methods in addressing incomplete data. Moreover, the performance gap between MARIA and its competitors widens as the missing data rate during training increases, whether in the “missing modalities” or “all missing” configurations. This trend suggests that MARIA is particu-

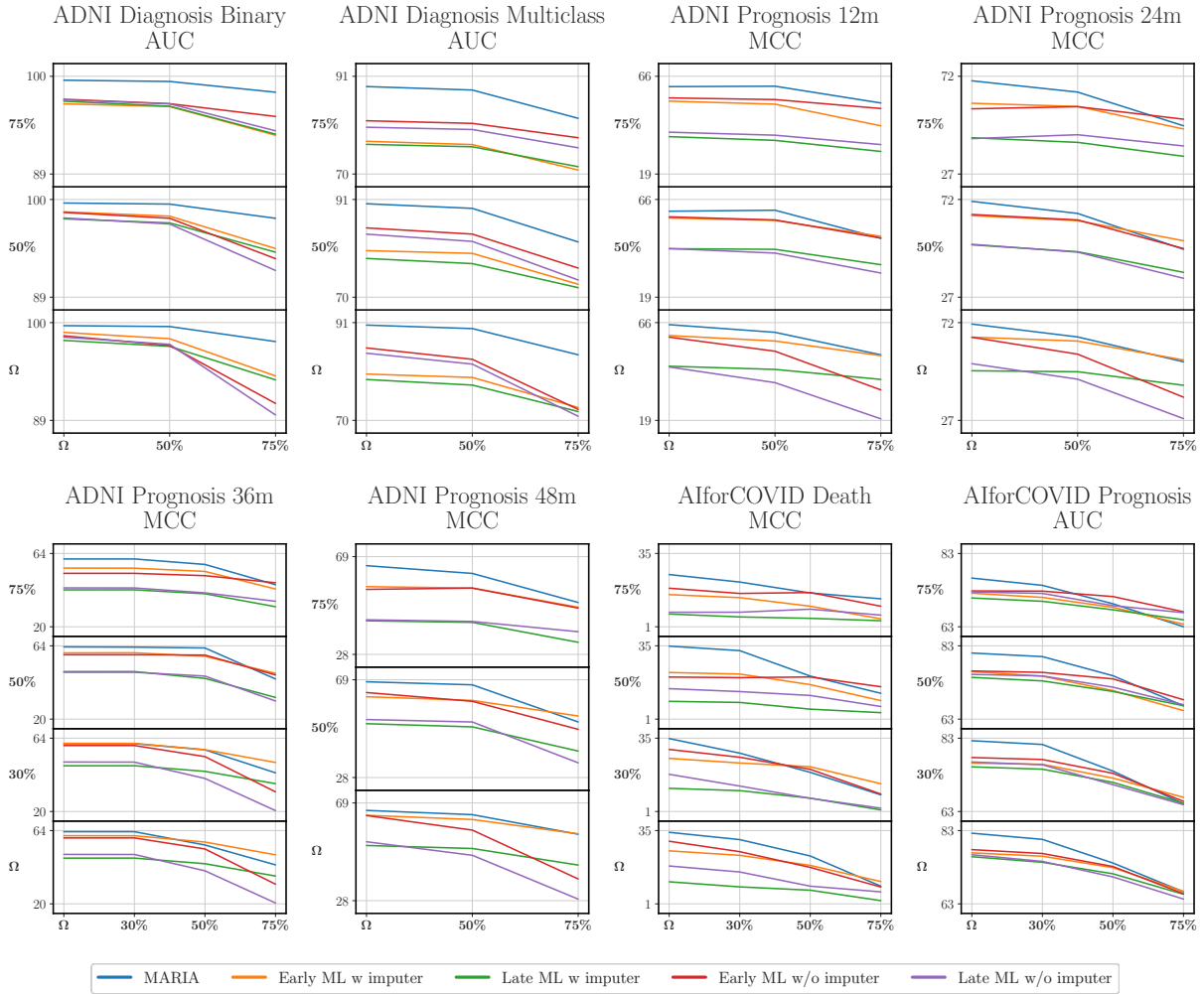


Figure 5.4: MARIA vs. ML in the “all missing” scenario. Each plot, one for each task, reports different charts, one for each of the missing rates in the training set, showing how the performance (y-axis) changes as the missing rate in the test set increases.

larly resilient to varying missing data scenarios. We attribute the improved performance of MARIA to its robust regularization techniques, which enable the model to effectively learn how to handle diverse missing rates during training.

Additionally, an analysis of Figure 5.4 reveals that early fusion approaches generally outperform late fusion ones. This finding underscores the limitations of late fusion in capturing intercorrelations between features from different modalities. Indeed, while late fusion can be compared to an ensemble of many experts, each specialized in a single modality and making independent predictions without communicating with each other, early fusion has a comprehensive view of the patient’s information, similar to how a physician integrates all available data about a patient, allowing it to achieve superior performance as a result.

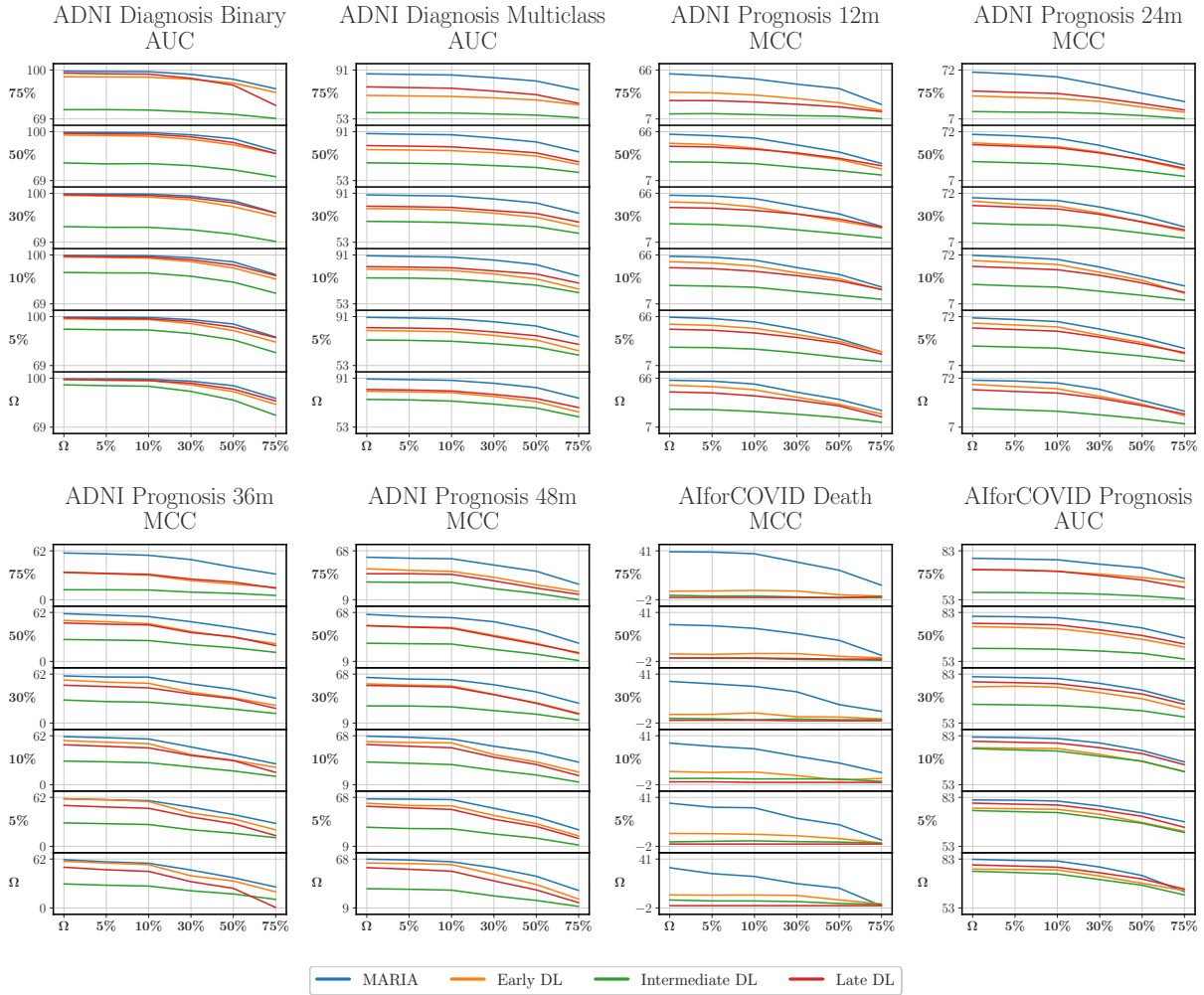


Figure 5.5: MARIA vs. DL in the “missing modalities” scenario. Each plot, one for each task, reports different charts, one for each of the missing rates in the training set, showing how the performance (y-axis) changes as the missing rate in the test set increases.

5.5.2 MARIA vs. DL

We also compared MARIA to leading competitors from the DL domain, specifically designed to analyze tabular data. In Figures 5.5 and 5.6, we reported the average performance of these models when paired with an imputer and grouped into early and late fusion approaches. Additionally, since no existing multimodal approaches are tailored for tabular data, we compared our method to intermediate fusion versions based on these models, as described in section 5.4.2. Once again, MARIA demonstrated superior performance across all levels of missing data in both experimental configurations. Similar to the observations made in comparisons with ML-based approaches, the performance gap between MARIA and its competitors widened as the missing rate in the training set increased. However, unlike

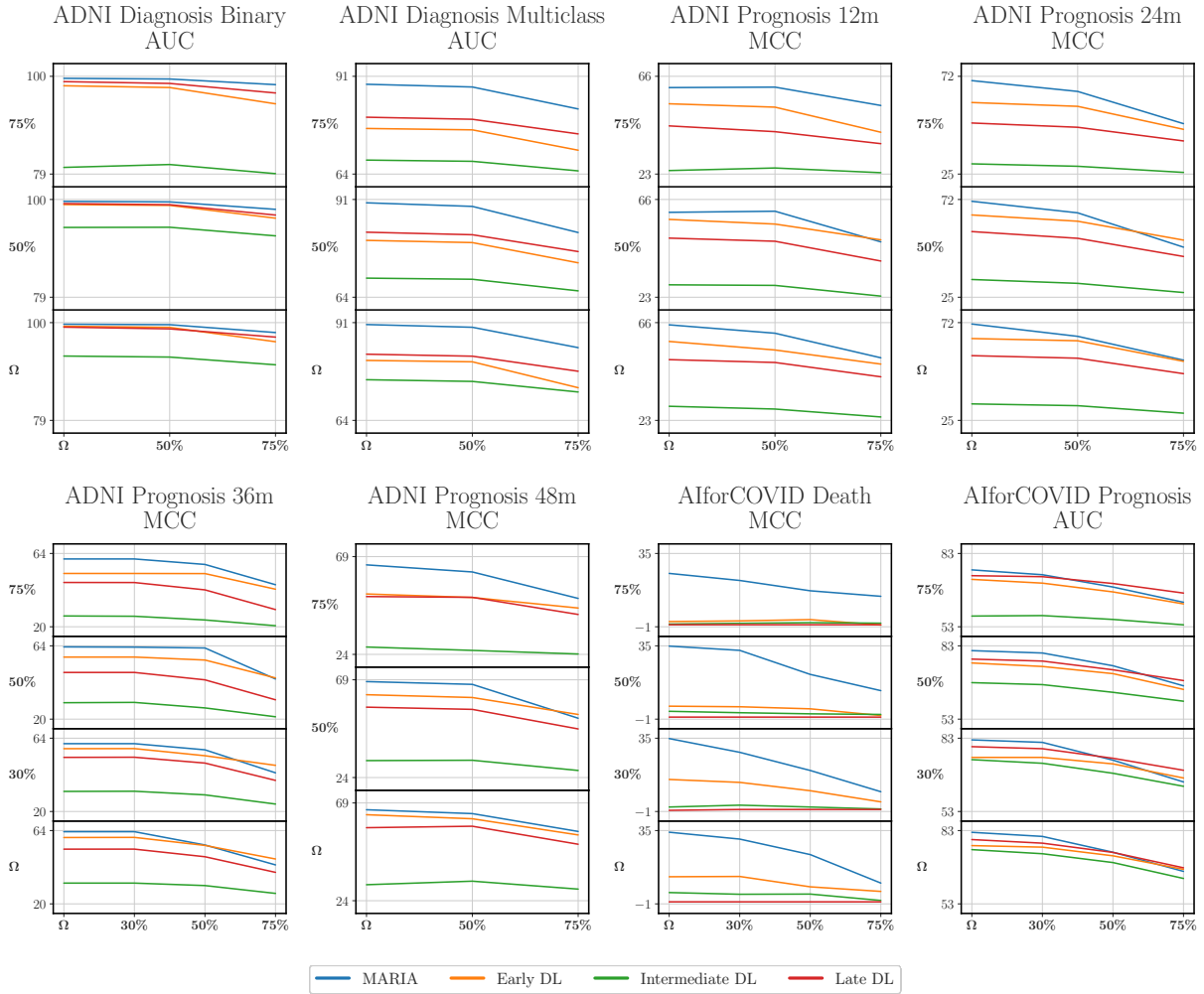


Figure 5.6: MARIA vs. DL in the “all missing” scenario. Each plot, one for each task, reports different charts, one for each of the missing rates in the training set, showing how the performance (y-axis) changes as the missing rate in the test set increases.

previous experiments, early and late fusion approaches exhibited comparable performance, indicating that both methodologies possess similar robustness to missing data. This outcome highlights the capacity of DL techniques to derive meaningful and informative representations from data, even in incomplete scenarios.

In contrast, the intermediate fusion approaches struggled to match the performance of other methods, consistently failing to perform as well as early fusion models. This suggests that intermediate fusion learning is unlikely to outperform its corresponding early and late fusion configurations when all modalities are tabular, particularly in scenarios with high missing rates.

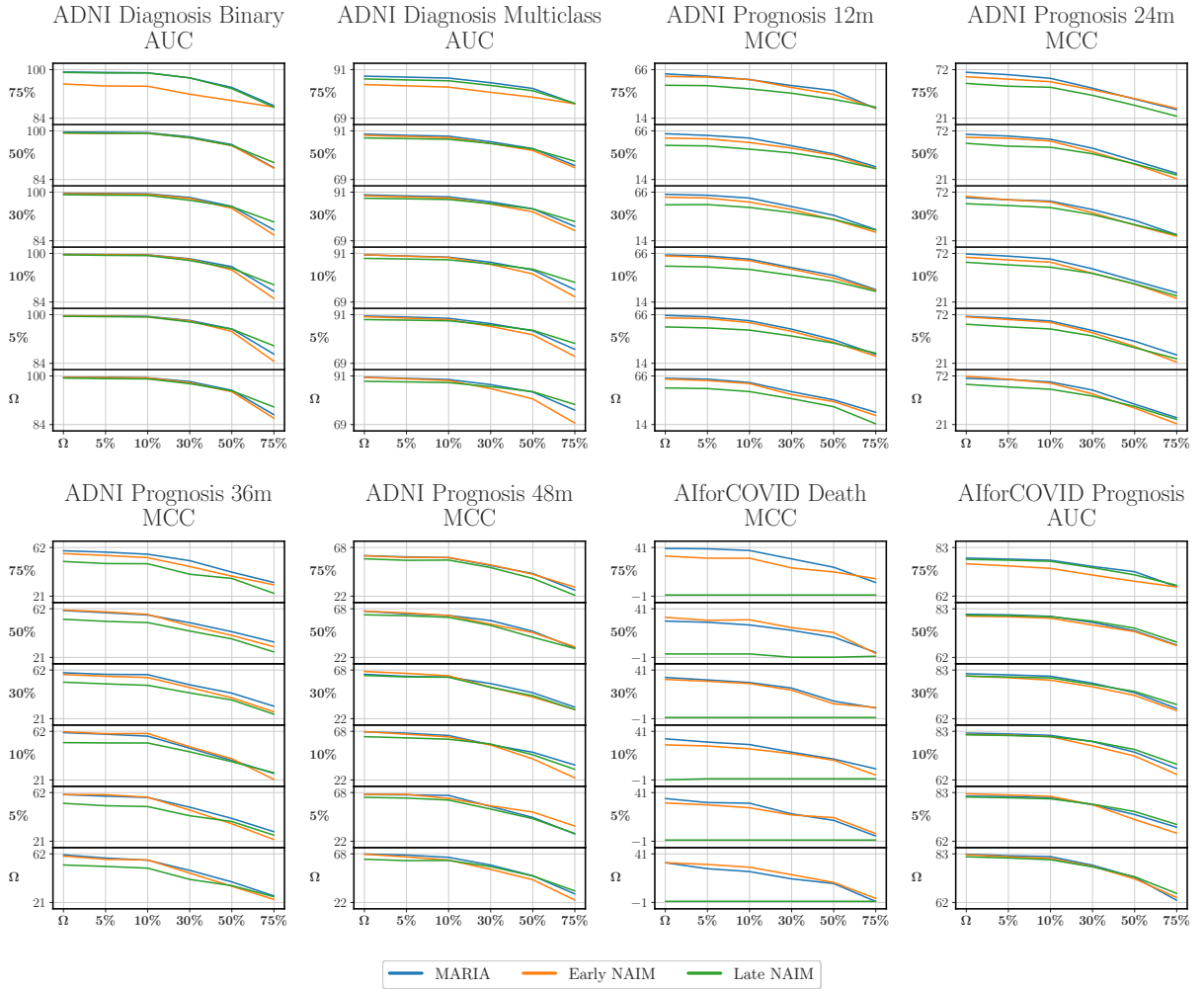


Figure 5.7: MARIA vs. NAIM in the “missing modalities” scenario. Each plot, one for each task, reports different charts, one for each of the missing rates in the training set, showing how the performance (y-axis) changes as the missing rate in the test set increases.

5.5.3 MARIA vs. NAIM

As a final analysis, we compared the proposed intermediate fusion approach with its respective early and late fusion counterparts, all based on the NAIM model [116]. As in the previous analyses, we used the same evaluation framework to compare the different approaches (Figures 5.7 and 5.8).

Interestingly, and as previously noted but now more pronounced, the intermediate fusion approach struggles to consistently outperform the early fusion approach, particularly in the context of tabular data. This suggests that early fusion may offer specific advantages when handling highly structured tabular datasets, where the model benefits from a unified representation of all input features from the beginning. In such contexts, the expected benefits of

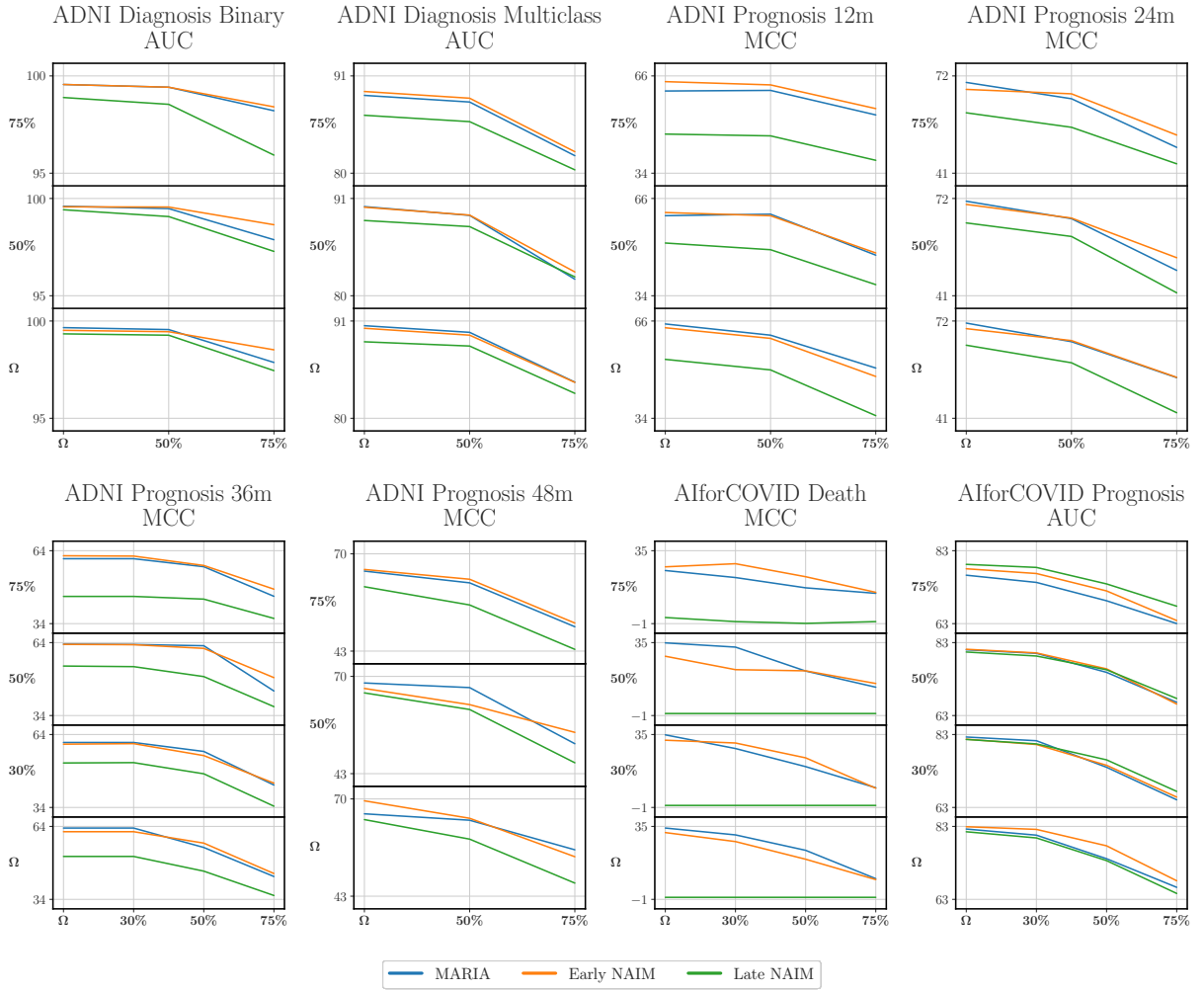


Figure 5.8: MARIA vs. NAIM in the “all missing” scenario. Each plot, one for each task, reports different charts, one for each of the missing rates in the training set, showing how the performance (y-axis) changes as the missing rate in the test set increases.

intermediate fusion appear less significant, potentially due to the inherent heterogeneity of tabular features, which may not require the additional representational flexibility provided by intermediate fusion.

By contrast, the late fusion approach, especially in the “all missing” configuration, rarely achieves performance comparable to the other two approaches. This limitation highlights the challenges faced by late fusion in capturing intercorrelations between features from different modalities, particularly when data completeness is severely compromised.

Overall, these results underscore the trade-offs between fusion strategies, suggesting that the choice between early, intermediate, and late fusion should be guided by the specific characteristics of the data modalities. For structured tabular data, early fusion seems to

provide an optimal balance between simplicity and performance. By contrast, intermediate and late fusion approaches may be more beneficial in scenarios involving heterogeneous or unstructured data sources.

In Appendix F, we present detailed tables for both types of experiments, showing the average performance in terms of AUC and MCC computed across the respective tasks under consideration. The tables include results at different percentages of missing values (reported in the columns). The first rows indicate the percentage of missing values used during training and testing, while the subsequent rows report the performance for each combination of fusion strategy, model, and missing data handling technique. To facilitate readability, we highlighted in bold the best performance for each combination of missing rates.

As presented in Table F.1, which summarizes the results for the “missing modalities” setting, the first table shows the average AUC scores for the ADNI diagnosis and AIforCOVID prognosis tasks. In this context, MARIA outperforms its competitors in 69.4% of cases (25 out of 36). Similarly, the second table reports the average MCC scores for the ADNI prognosis and AIforCOVID death tasks, where MARIA achieves superior performance in 61.1% of cases (22 out of 36).

In contrast, Table F.2, which details the results for the “all missing” setting, indicates that MARIA exhibits stronger performance in fewer instances: 37.5% (6 out of 16) for tasks evaluated using AUC (first table) and 12.5% (2 out of 16) for tasks assessed with MCC (second table). This difference is likely due to the “all missing” setting generally involving less severe information loss compared to the “missing modalities” setting. These findings suggest that the advantages of MARIA are particularly pronounced in scenarios characterized by significant modality loss.

These analyses demonstrate the robustness of the MARIA model, particularly under conditions with high missing rates, a common challenge in medical applications where patient data are often incomplete or inconsistently available. By effectively leveraging only the available information, the MARIA model enhances diagnostic accuracy and decision support in healthcare, ultimately leading to improved patient outcomes. Across a wide range of experimental configurations and degrees of missingness, the model consistently outperforms competing approaches, including both traditional ML and DL models.

5.6 Conclusions

In this study, we introduced MARIA, a novel transformer-based model designed to tackle the challenges of incomplete multimodal data, especially in healthcare. MARIA employs an intermediate fusion strategy, integrating data from multiple incomplete modalities through

a masked self-attention mechanism that selectively focuses on available information while ignoring missing parts. This approach not only avoids the disadvantages of synthetic imputation but also ensures a robust predictive performance even in the presence of severe data missingness. The results show that MARIA is effective across multiple diagnostic and prognostic tasks, consistently outperforming traditional ML and DL models by adapting to various missing data scenarios during both training and inference.

Despite these promising results, some limitations need to be addressed. One of the key challenges of MARIA lies in its computational complexity. Although the model eliminates the need for imputers, the use of masked attention mechanisms and the intermediate fusion strategy requires substantial computational resources, which may limit its scalability in low-resource environments or when applied to extremely large datasets. To mitigate this, future work will focus on enhancing MARIA’s scalability and efficiency, potentially by incorporating more lightweight attention mechanisms or by developing hybrid techniques that blend intermediate fusion with other fusion strategies in a computationally feasible manner.

Furthermore, while MARIA’s design is effective for tabular data, its generalization to other types of multimodal inputs, such as imaging or textual data, has not been fully explored [131, 132]. To address this, future research will aim to expand MARIA’s applicability to other forms of multimodal data, including medical images and clinical notes [118, 133, 87, 134, 135], which would ensure broader usability across diverse healthcare datasets.

Moreover, a limitation of the current work lies in the exclusive focus on MCAR patterns during training and evaluation. While this setting allows for unbiased benchmarking and controlled experimentation, it does not capture the complexity of Missing At Random (MAR) or Missing Not At Random (MNAR) scenarios, which are more reflective of real-world clinical datasets. In future work, we plan to extend the masking strategies used in training to emulate more realistic missingness mechanisms, by introducing feature- or outcome-dependent dropout schemes. This would allow us to better assess and enhance MARIA’s robustness under structured or biased missingness conditions.

Chapter 6

Conclusions and Future Perspectives on Resilient Clinical AI

The development of methods resilient to incomplete clinical data is a crucial step toward ensuring that AI can be safely, effectively, and meaningfully integrated into real-world healthcare systems. In modern clinical environments, where data is collected under varying protocols and with inconsistent frequency, the presence of missing or incomplete data is not a rare anomaly, it is a defining characteristic. These gaps often stem from differences in clinical workflows, limitations in data acquisition tools, patient variability, or institutional constraints. Consequently, designing AI models that can operate under these imperfect conditions is essential to building systems that can perform reliably in real-world clinical practice.

In this thesis, we have taken a structured and progressive approach to tackling the challenges posed by data incompleteness. We began with foundational research on structured tabular data and gradually advanced to more complex scenarios involving multimodal inputs. Each step in this progression has aimed to build a deeper understanding of how deep learning systems can be designed to handle the uncertainty and fragmentation that are inherent to healthcare datasets.

In Chapter 2, we introduced NAIM, a transformer-based model that incorporates masked self-attention to learn from tabular datasets with incomplete features. Unlike conventional approaches that rely heavily on data imputation techniques before modeling, NAIM directly attends to the available inputs, ignoring missing elements in a principled way. This method was benchmarked against multiple existing approaches and demonstrated consistent performance across varying missingness levels. The findings highlighted that integrating resilience directly into the models design, rather than treating it as an external correction, can lead to better adaptability and more trustworthy predictions, especially in real-world scenarios

where clean, complete datasets are rarely available.

Chapter 3 extended this idea to a more complex and realistic clinical task: predicting OS in NSCLC patients using routinely collected clinical data. We employed a variant of NAIM tailored for time-dependent outputs, incorporating custom loss functions capable of handling censored survival data across multiple prediction horizons. Although the presence of missing values posed clear difficulties, the model demonstrated strong predictive stability and reliability. These results show that transformer-based models, when properly adapted, can serve as powerful tools for survival analysis in routine clinical environments where complete data is rarely the norm.

The thesis then transitioned to the domain of multimodal learning in Chapter 4, where we investigated whether combining structured clinical data with medical imaging (specifically CT scans) could enhance predictive performance. This initial attempt employed a late fusion strategy, in which each modality was processed separately before combining the resulting representations at the decision level. While this design offers robustness to missing modalities, since decisions can still be made from the available components, it lacks the ability to capture deep inter-modality relationships. As such, it limits the models ability to fully leverage the complementary nature of different data types. These observations align with the findings of our prior systematic review on multimodal fusion techniques in clinical AI, which identified late fusion as a commonly adopted yet fundamentally limited approach [107]. In contrast, intermediate fusion approaches appear more promising, as they allow for richer interactions between modalities during the learning process, enabling the model to better capture complementary information and dependencies across data types.

To address these limitations, Chapter 5 introduced MARIA, a novel transformer-based architecture designed to perform masked integration across multiple modalities. MARIA employs an intermediate fusion strategy, which enables interaction between modalities, enhancing the models capacity to capture cross-modal synergies while also accounting for missing data. The architecture uses masks to attend selectively to available inputs, making it resilient to various forms of structural incompleteness. Our experiments, conducted across numerous simulated missing-modality scenarios, demonstrated that MARIA outperformed traditional fusion models in both performance and stability. This advancement underscores the importance of designing fusion architectures that are not only effective in ideal conditions but also robust enough to handle the inherent variability of clinical data.

From a methodological standpoint, the central argument of this thesis is that resilience to missing data must be treated as a core design principle in AI for healthcare. Rather than viewing incomplete data as a minor problem to be fixed before analysis, we propose that models should be built with an awareness of what is missing. This perspective leads

to architectures that are inherently more flexible and adaptive, capable of making the most out of whatever information is available without introducing additional bias or uncertainty through arbitrary imputations.

Clinically, the implications of this work are substantial. AI tools that fail when data is missing are unlikely to gain acceptance in day-to-day practice, where clinicians often make decisions based on incomplete records. The models presented in this thesis provide a framework for developing AI systems that are both technically sound and practically deployable. By embracing the reality of missing data rather than working around it, these systems demonstrate a level of robustness and adaptability that is essential for meaningful integration into clinical workflows. Moreover, their ability to operate under uncertainty makes them more aligned with the realities of medical decision-making, where perfect information is rarely available.

6.1 Limitations and Future Works

Despite the promising outcomes and contributions presented in this thesis, several important limitations remain that shape the boundaries of the work and highlight opportunities for future investigation. Addressing these limitations will be essential to broaden the clinical utility, adaptability, and impact of resilient AI models developed for healthcare applications.

Evaluation in a Truly Multimodal Context Currently, the evaluation of MARIA has been limited to tasks involving only tabular data coming from different sources. Although the architecture was explicitly designed to generalize across heterogeneous modalities, such as radiological images and structured clinical features, its capabilities in genuinely multimodal settings have yet to be tested. A logical and necessary extension of this work would be to validate MARIA using datasets that include both imaging and tabular components. Such testing would reflect more realistic clinical environments where multiple types of patient data must be processed and interpreted simultaneously, often under conditions of partial or missing modalities. Demonstrating MARIA’s ability to learn meaningful representations across different data types in these scenarios will be critical to assessing its broader clinical applicability.

Expanding Across Modalities and Clinical Domains While the models in this thesis are tested on structured tabular data and CT imaging, healthcare data often spans a far broader spectrum. Modalities such as clinical free-text notes, genomics, laboratory time series, biosensor streams, and electronic health records introduce new challenges and oppor-

tunities. Evaluating the proposed models, and especially MARIA, on tasks involving diverse data types will offer insights into their generalizability, flexibility, and capacity for transfer learning. Furthermore, applying these models in different clinical domains beyond oncology, such as cardiology, neurology, and intensive care, will be crucial for establishing their real-world relevance.

Incorporating Temporal and Longitudinal Information Much of the work in this thesis treats patient data as static inputs, without fully modeling how information evolves over time. In clinical reality, however, healthcare decision-making often relies on interpreting sequences of events, such as changes in clinical measurements, the timing of treatments, or disease progression over multiple encounters. The integration of temporal dynamics across various modalities, especially when faced with missing or irregularly sampled data, remains an underexplored but essential area of investigation. Extending architectures like MARIA to handle temporal data streams, such as visit sequences, diagnostic timelines, or evolving treatment patterns, could enable more accurate modeling of patient trajectories and support early detection, risk stratification, and personalized interventions based on longitudinal trends.

Improving Scalability and Efficiency A further point concerns computational efficiency. As previously discussed, transformer-based models inherently demand greater training resources than classical ML pipelines coupled with imputation. Nonetheless, this cost is partially offset by the elimination of extensive imputation and hyperparameter optimization procedures, which can themselves be computationally and methodologically demanding. Future research should therefore focus on designing lightweight attention mechanisms and hybrid fusion strategies that preserve predictive performance while reducing training and memory overhead. Such advances would help make NAIM and MARIA more practical for deployment in large-scale or resource-constrained clinical environments, where scalability and efficiency are as critical as accuracy.

Broadening the Types of Missingness Considered This thesis primarily addresses missing completely at random (MCAR) conditions, which are analytically tractable but often overly simplistic for real-world clinical data. In practice, clinical missingness frequently follows more complex mechanisms, such as missing at random (MAR) or missing not at random (MNAR), where the probability of missingness depends on observed or unobserved variables, respectively. These forms of missingness introduce both statistical and computational challenges that must be addressed for models to be resilient in operational healthcare

settings. Future work should consider how models can be explicitly trained under MAR and MNAR conditions, possibly by modifying loss functions or introducing auxiliary variables that encode the likelihood of missingness. A promising direction involves embedding these more realistic missingness mechanisms into the models regularization schemes, allowing the system to internalize and adapt to various missingness patterns during training. This could enable more detailed and accurate reasoning in the presence of biased or structured data gaps.

Leveraging Self-Supervised and Semi-Supervised Learning The current models are trained under fully supervised paradigms, relying on labeled datasets curated for specific tasks. However, in many healthcare contexts, annotated data is limited or prohibitively expensive to obtain. Large portions of available data remain unlabeled, especially for rare conditions or outcomes. Incorporating self-supervised learning approaches, such as masked prediction tasks or contrastive learning, could allow the models to pretrain on unlabeled data and adapt to new tasks with minimal supervision. Similarly, semi-supervised methods that leverage both labeled and unlabeled samples may improve data efficiency, accelerate model development, and enable better generalization in low-resource scenarios.

Toward Federated and Privacy-Preserving Learning A significant limitation of the current thesis lies in its implicit assumption of centralized data access. In practice, healthcare data is distributed across multiple institutions, each bound by legal, ethical, and infrastructural constraints on data sharing. Future research should explore how the architectural principles introduced here can be translated into federated learning settings. In such frameworks, models can be trained collaboratively across sites without exchanging raw patient data, thus preserving privacy. Moreover, federated models can be adapted to heterogeneous data environments, where each institution may have different patient populations and data collection practices. This adaptation is essential for building fair and inclusive clinical AI systems.

Enhancing Interpretability and Clinical Trust Interpretability remains a central challenge for deploying AI in medicine. Although the transformer models developed in this thesis incorporate attention mechanisms, which can provide insight into feature importance, the clinical validity of such interpretations is not well established. Indeed, attention maps may not align with human clinical reasoning, and their interpretive value remains context-dependent. Future efforts should investigate whether attention weights can be reliably used for explanation in healthcare tasks, or whether additional techniques, such as saliency maps, SHAP values, or counterfactual example generation, are necessary. Embedding these expla-

nations into clinical workflows could enhance transparency, support clinical validation, and ultimately increase trust and acceptance among healthcare professionals.

Beyond these specific directions, the general theme of this thesis supports the vision of resilient, multimodal AI as a transformative force in healthcare. By designing models that can intelligently integrate incomplete and heterogeneous information, we move closer to systems that reflect the complexity and imperfection of real-world medicine. Such systems have the potential to offer not only higher accuracy but also increased fairness and adaptability, qualities that are essential for clinical tools deployed in diverse and unpredictable environments. As research in this field advances, resilient multimodal learning could become a cornerstone of next-generation clinical decision support tools, capable of handling noisy data, adapting to new settings, and empowering healthcare professionals to make more informed, context-sensitive decisions.

The limitations identified in this chapter, rather than undermining the value of the work, serve to highlight where future efforts should be directed. They outline a clear direction for future research aimed at developing models that are better suited to the complexities of real-world medical data and more aligned with the practical challenges of clinical AI applications.

Ultimately, this thesis argues that embracing the imperfections of clinical data is not a compromise, but a necessary evolution in the development of AI for healthcare. By shifting the emphasis from data completeness to model-level resilience, we move toward systems that are grounded in the realities of clinical decision-making, where information is often partial, but the need to act remains urgent.

“In medicine, waiting for complete data is often not an option. Learning to act with what we have: that’s where AI must rise.”

Bibliography

- [1] Alexandre Perez-Lebel et al. “Benchmarking missing-values approaches for predictive models on health databases”. In: *GigaScience* 11 (2022), giac013.
- [2] Ravid Shwartz-Ziv and Amitai Armon. “Tabular data: Deep learning is not all you need”. In: *Information Fusion* 81 (2022), pp. 84–90.
- [3] Alexey Dosovitskiy et al. “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”. In: *arXiv preprint arXiv:2010.11929* (2020).
- [4] Alec Radford et al. “Robust Speech Recognition via Large-Scale Weak Supervision”. In: *International Conference on Machine Learning*. PMLR. 2023, pp. 28492–28518.
- [5] Sercan Ö Arik and Tomas Pfister. “TabNet: Attentive Interpretable Tabular Learning”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 35. 2021, pp. 6679–6687.
- [6] Xin Huang et al. “TabTransformer: Tabular Data Modeling Using Contextual Embeddings”. In: *arXiv preprint arXiv:2012.06678* (2020).
- [7] Yury Gorishniy et al. “Revisiting Deep Learning Models for Tabular Data”. In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 18932–18943.
- [8] Mingxuan Liu et al. “Handling missing values in healthcare data: A systematic review of deep learning-based imputation techniques”. In: *Artificial Intelligence in Medicine* (2023), p. 102587.
- [9] H Cevallos Valdiviezo and Stefan Van Aelst. “Tree-based prediction on incomplete data using imputation or surrogate decisions”. In: *Information Sciences* 311 (2015), pp. 163–181.
- [10] Stef Van Buuren and Karin Groothuis-Oudshoorn. “mice: Multivariate Imputation by Chained Equations in R”. In: *Journal of statistical software* 45 (2011), pp. 1–67.
- [11] Olga Troyanskaya et al. “Missing value estimation methods for DNA microarrays”. In: *Bioinformatics* 17.6 (2001), pp. 520–525.

-
- [12] Jiaxuan You et al. “Handling missing data with graph representation learning”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 19075–19087.
- [13] Ashish Vaswani et al. “Attention is All you Need”. In: *Advances in neural information processing systems* 30 (2017).
- [14] Terrance DeVries and Graham W Taylor. “Improved Regularization of Convolutional Neural Networks with Cutout”. In: *arXiv preprint arXiv:1708.04552* (2017).
- [15] D.; Dua and C. Graff. *UCI Machine Learning Repository*. 2017. URL: <http://archive.ics.uci.edu/ml>.
- [16] Barry Becker and Ronny Kohavi. *Adult*. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5XW20>. 1996. DOI: <https://doi.org/10.24432/C5XW20>.
- [17] S. Moro, P. Rita, and P. Cortez. *Bank Marketing*. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5K306>. 2012. DOI: <https://doi.org/10.24432/C5K306>.
- [18] C. Sakar and Yomi Kastro. *Online Shoppers Purchasing Intention Dataset*. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5F88Q>. 2018. DOI: <https://doi.org/10.24432/C5F88Q>.
- [19] Marek Sikora and Lukasz Wrobel. *seismic-bumps*. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5W902>. 2013. DOI: <https://doi.org/10.24432/C5W902>.
- [20] Mark Hopkins et al. *Spambase*. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C53G6X>. 1999. DOI: <https://doi.org/10.24432/C53G6X>.
- [21] Yoav Freund, Robert Schapire, and Naoki Abe. “A Short Introduction to Boosting”. In: *Journal-Japanese Society For Artificial Intelligence* 14.771-780 (1999), p. 1612.
- [22] Leo Breiman et al. “Classification and regression trees”. In: *Pacific Grove, Wadsworth* (1984).
- [23] Tin Kam Ho. “Random Decision Forests”. In: *Proceedings of 3rd international conference on document analysis and recognition*. Vol. 1. IEEE. 1995, pp. 278–282.
- [24] Corinna Cortes and Vladimir Vapnik. “Support-Vector Networks”. In: *Machine learning* 20.3 (1995), pp. 273–297.
- [25] Tianqi Chen and Carlos Guestrin. “XGBoost: A Scalable Tree Boosting System”. In: *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*. 2016, pp. 785–794.

- [26] Hassan Ramchoun et al. “Multilayer Perceptron: Architecture Optimization and Training”. In: *International Journal of Interactive Multimedia and Artificial Intelligence* 4.1 (Special Issue on Artificial Intelligence Underpinning 2016), pp. 26–30.
- [27] Philipp Thölke et al. “Class imbalance should not throw you off balance: Choosing the right classifiers and performance metrics for brain decoding with imbalanced data”. In: *NeuroImage* 277 (2023), p. 120253.
- [28] Guy Dar et al. “Analyzing Transformers in Embedding Space”. In: *Annual Meeting of the Association for Computational Linguistics*. 2023, pp. 16124–16170.
- [29] A Cortellini et al. “Transformer-based AI approach to unravel long-term, time-dependent prognostic complexity in patients with advanced NSCLC and PD-L1 \geq 50%: insights from the pembrolizumab 5-year global registry”. In: *J Immunother Cancer* 13 (2025). DOI: [10.1136/jitc-2025-012423](https://doi.org/10.1136/jitc-2025-012423).
- [30] Abdullah Alanazi. “Using machine learning for healthcare challenges and opportunities”. In: *Informatics in Medicine Unlocked* (2022), p. 100924.
- [31] Muhammad Javed Iqbal et al. “Clinical applications of artificial intelligence and machine learning in cancer diagnosis: looking into the future”. In: *Cancer cell international* 21.1 (2021), pp. 1–11.
- [32] World Health Organisation. *LUNG*. Available online: <https://gco.iarc.fr/today/data/factsheets/cancers/15-Lung-fact-sheet.pdf>. 2020.
- [33] Cancer.net. *Lung Cancer Non-Small Cell: Statistics*. Available online: <https://www.cancer.net/cancer-types/lung-cancer-non-small-cell/statistics>. 2022.
- [34] Georgios Kantidakis, Audinga-Dea Hazewinkel, Marta Fiocco, et al. “Neural Networks for Survival Prediction in Medicine Using Prognostic Factors: A Review and Critical Appraisal”. In: *Computational and Mathematical Methods in Medicine* 2022 (2022).
- [35] Daochen Zha et al. “Data-centric artificial intelligence: A survey”. In: *arXiv preprint arXiv:2303.10158* (2023).
- [36] Mogana Darshini Ganggayah et al. “Predicting factors for survival of breast cancer patients using machine learning techniques”. In: *BMC medical informatics and decision making* 19 (2019), pp. 1–17.
- [37] Sumeet Hindocha et al. “A comparison of machine learning methods for predicting recurrence and death after curative-intent radiotherapy for non-small cell lung cancer: Development and validation of multivariable clinical prediction models”. In: *EBioMedicine* 77 (2022), p. 103911.

- [38] Jason C Hsu et al. “Development and Validation of Novel Deep-Learning Models Using Multiple Data Types for Lung Cancer Survival”. In: *Cancers* 14.22 (2022), p. 5562.
- [39] Yang Yang et al. “Machine learning application in personalised lung cancer recurrence and survivability prediction”. In: *Computational and Structural Biotechnology Journal* 20 (2022), pp. 1811–1820.
- [40] Yash Dagli, Saumya Choksi, and Sudipta Roy. “Prediction of two year survival among patients of non-small cell lung cancer”. In: *Computer Aided Intervention and Diagnostics in Clinical and Medical Images*. Springer. 2019, pp. 169–177.
- [41] Yougen Wu et al. “Using machine learning for mortality prediction and risk stratification in atezolizumab-treated cancer patients: Integrative analysis of eight clinical trials”. In: *Cancer Medicine* 12.3 (2023), pp. 3744–3757.
- [42] Magdalena Oguszka et al. “Evaluate Cutpoints: Adaptable continuous data distribution system for determining survival in Kaplan-Meier estimator”. In: *Computer Methods and Programs in Biomedicine* 177 (2019), pp. 133–139.
- [43] Cary Oberije et al. “A validated prediction model for overall survival from stage III non-small cell lung cancer: toward survival prediction for individual patients”. In: *International Journal of Radiation Oncology* Biology* Physics* 92.4 (2015), pp. 935–944.
- [44] Hugo Loureiro et al. “Artificial intelligence for prognostic scores in oncology: a benchmarking study”. In: *Frontiers in Artificial Intelligence* 4 (2021), p. 625573.
- [45] Bora Lee et al. “DeepBTS: prediction of recurrence-free survival of non-small cell lung cancer using a time-binned deep neural network”. In: *Scientific reports* 10.1 (2020), pp. 1–10.
- [46] Qianyu Yuan et al. “Performance of a machine learning algorithm using electronic health record data to identify and estimate survival in a longitudinal cohort of patients with lung cancer”. In: *JAMA Network Open* 4.7 (2021), e2114723–e2114723.
- [47] Gaetano Manzo et al. “Breast cancer survival analysis agents for clinical decision support”. In: *Computer Methods and Programs in Biomedicine* 231 (2023), p. 107373.
- [48] Yiling Wang et al. “Deep learning based time-to-event analysis with PET, CT and joint PET/CT for head and neck cancer prognosis”. In: *Computer Methods and Programs in Biomedicine* 222 (2022), p. 106948.

- [49] Changhee Lee et al. “DeepHit: A Deep Learning Approach to Survival Analysis With Competing Risks”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 32.1 (Apr. 2018). DOI: [10.1609/aaai.v32i1.11842](https://doi.org/10.1609/aaai.v32i1.11842). URL: <https://ojs.aaai.org/index.php/AAAI/article/view/11842>.
- [50] *CLARO - CoLLaborative multi-sources Radiopathomics approach for personalized Oncology in non-small cell lung cancer*. <http://www.cosbi-lab.it/claro/>. Accessed: 2023-03-20.
- [51] Peter C Austin, Douglas S Lee, and Jason P Fine. “Introduction to the analysis of survival data in the presence of competing risks”. In: *Circulation* 133.6 (2016), pp. 601–609.
- [52] A. Arcuri and G. Fraser. “Parameter tuning or default values? An empirical investigation in search-based software engineering”. In: *Empirical Software Engineering* 18.3 (2013), pp. 594–623.
- [53] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [54] Daniel J Stekhoven and Peter Bühlmann. “MissForestnon-parametric missing value imputation for mixed-type data”. In: *Bioinformatics* 28.1 (2012), pp. 112–118.
- [55] Sebastian Pölsterl. “scikit-survival: A Library for Time-to-Event Analysis Built on Top of scikit-learn”. In: *Journal of Machine Learning Research* 21.212 (2020), pp. 1–6. URL: <http://jmlr.org/papers/v21/20-729.html>.
- [56] Scott M Lundberg and Su-In Lee. “A Unified Approach to Interpreting Model Predictions”. In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon et al. Curran Associates, Inc., 2017, pp. 4765–4774.
- [57] Myra Van Laar et al. “Prognostic factors for overall survival of stage III non-small cell lung cancer patients on computed tomography: A systematic review and meta-analysis”. In: *Radiotherapy and Oncology* 151 (2020), pp. 152–175.
- [58] National Cancer Institute. *Non-Small Cell Lung Cancer Treatment*. Available online: <https://www.cancer.gov/types/lung/hp/non-small-cell-lung-treatment-pdq>. (accessed on 18 May 2022). 2022.
- [59] Philippe Lambin et al. “Radiomics: Extracting more information from medical images using advanced feature analysis”. en. In: *European Journal of Cancer* 48.4 (Mar. 2012), pp. 441–446. ISSN: 09598049. DOI: [10.1016/j.ejca.2011.11.036](https://doi.org/10.1016/j.ejca.2011.11.036). URL: <https://linkinghub.elsevier.com/retrieve/pii/S0959804911009993> (visited on 05/18/2022).

- [60] Madeleine Scrivener et al. “Radiomics applied to lung cancer: a review”. en. In: *Translational Cancer Research* 5.4 (Aug. 2016), pp. 398–409. ISSN: 2218676X, 22196803. DOI: [10.21037/tcr.2016.06.18](https://doi.org/10.21037/tcr.2016.06.18). URL: <http://tcr.amegroups.com/article/view/8536/7742> (visited on 05/24/2022).
- [61] Sara Ramella et al. “A radiomic approach for adaptive radiotherapy in non-small cell lung cancer patients”. en. In: *PLOS ONE* 13.11 (Nov. 2018). Ed. by Aamir Ahmad, e0207455. ISSN: 1932-6203. DOI: [10.1371/journal.pone.0207455](https://doi.org/10.1371/journal.pone.0207455). URL: <https://dx.plos.org/10.1371/journal.pone.0207455> (visited on 04/03/2022).
- [62] Rosa Sicilia et al. “Exploratory radiomics for predicting adaptive radiotherapy in non-small cell lung cancer”. In: *2018 IEEE 31st International Symposium on Computer-Based Medical Systems (CBMS)*. IEEE. 2018, pp. 250–255.
- [63] Irantzu Anzar et al. “NeoMutate: An ensemble machine learning framework for the prediction of somatic mutations in cancer”. In: *BMC Medical Genomics* 12.1 (2019). DOI: [10.1186/s12920-019-0508-5](https://doi.org/10.1186/s12920-019-0508-5). URL: <https://link.springer.com/article/10.1186/s12920-019-0508-5>.
- [64] Rajarsi Gupta et al. “The Emergence of Pathomics”. en. In: *Current Pathobiology Reports* 7.3 (Sept. 2019), pp. 73–84. ISSN: 2167-485X. DOI: [10.1007/s40139-019-00200-x](https://doi.org/10.1007/s40139-019-00200-x). URL: <http://link.springer.com/10.1007/s40139-019-00200-x> (visited on 05/24/2022).
- [65] Charles Z. Liu et al. “Exploring Deep Pathomics in Lung Cancer”. en. In: *2021 IEEE 34th International Symposium on Computer-Based Medical Systems (CBMS)*. Aveiro, Portugal: IEEE, June 2021, pp. 407–412. ISBN: 978-1-66544-121-6. DOI: [10.1109/CBMS52027.2021.00092](https://doi.org/10.1109/CBMS52027.2021.00092). URL: <https://ieeexplore.ieee.org/document/9474667/> (visited on 04/04/2022).
- [66] Hugo J. W. L. Aerts et al. “Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach”. en. In: *Nature Communications* 5.1 (Sept. 2014), p. 4006. ISSN: 2041-1723. DOI: [10.1038/ncomms5006](https://doi.org/10.1038/ncomms5006). URL: <http://www.nature.com/articles/ncomms5006> (visited on 04/04/2022).
- [67] Elizabeth Huynh et al. “CT-based radiomic analysis of stereotactic body radiation therapy patients with lung cancer”. In: *Radiotherapy and Oncology* 120.2 (2016), pp. 258–266. ISSN: 0167-8140. DOI: <https://doi.org/10.1016/j.radonc.2016.05.024>. URL: <https://www.sciencedirect.com/science/article/pii/S0167814016311380>.

- [68] Xenia Fave et al. “Delta-radiomics features for the prediction of patient outcomes in nonsmall cell lung cancer”. en. In: *Scientific Reports* 7.1 (Dec. 2017), p. 588. ISSN: 2045-2322. DOI: [10.1038/s41598-017-00665-z](https://doi.org/10.1038/s41598-017-00665-z). URL: <http://www.nature.com/articles/s41598-017-00665-z> (visited on 05/20/2022).
- [69] Hongming Li et al. “Unsupervised machine learning of radiomic features for predicting treatment response and overall survival of early stage non-small cell lung cancer patients treated with stereotactic body radiation therapy”. In: *Radiotherapy and Oncology* 129.2 (2018). Special Issue: Radiotherapy in Asia - Part 2, pp. 218–226. ISSN: 0167-8140. DOI: <https://doi.org/10.1016/j.radonc.2018.06.025>. URL: <https://www.sciencedirect.com/science/article/pii/S0167814018333401>.
- [70] Chintan Parmar et al. “Machine Learning methods for Quantitative Radiomic Biomarkers”. In: *Scientific reports* 5 (Aug. 2015), p. 13087. DOI: [10.1038/srep13087](https://doi.org/10.1038/srep13087).
- [71] Natascha Claudia DAMico et al. “Radiomics-Based Prediction of Overall Survival in Lung Cancer Using Different Volumes-Of-Interest”. en. In: *Applied Sciences* 10.18 (Sept. 2020), p. 6425. ISSN: 2076-3417. DOI: [10.3390/app10186425](https://doi.org/10.3390/app10186425). URL: <https://www.mdpi.com/2076-3417/10/18/6425> (visited on 04/04/2022).
- [72] Lars Heiliger et al. “Beyond Medical Imaging: A Review of Multimodal Deep Learning in Radiology”. en. In: (), p. 13.
- [73] Sören Richard Stahlschmidt, Benjamin Ulfenborg, and Jane Synnergren. “Multimodal deep learning for biomedical data fusion: a review”. en. In: *Briefings in Bioinformatics* 23.2 (Mar. 2022), bbab569. ISSN: 1467-5463, 1477-4054. DOI: [10.1093/bib/bbab569](https://doi.org/10.1093/bib/bbab569). URL: <https://academic.oup.com/bib/article/doi/10.1093/bib/bbab569/6516346> (visited on 04/04/2022).
- [74] Yiyang Zhang et al. “A Novel Multimodal Radiomics Model for Preoperative Prediction of Lymphovascular Invasion in Rectal Cancer”. en. In: *Frontiers in Oncology* 10 (Apr. 2020), p. 457. ISSN: 2234-943X. DOI: [10.3389/fonc.2020.00457](https://doi.org/10.3389/fonc.2020.00457). URL: <https://www.frontiersin.org/article/10.3389/fonc.2020.00457/full> (visited on 04/04/2022).
- [75] Matteo Tortora et al. “RadioPathomics: Multimodal Learning in Non-Small Cell Lung Cancer for Adaptive Radiotherapy”. In: *arXiv preprint arXiv:2204.12423* (2022).
- [76] Stefano Cipollari et al. “Convolutional neural networks for automated classification of prostate multiparametric magnetic resonance imaging based on image quality”. In: *Journal of Magnetic Resonance Imaging* 55.2 (2022), pp. 480–490.

- [77] Mehdi Amini et al. “Multi-Level PET and CT Fusion Radiomics-based Survival Analysis of NSCLC Patients”. en. In: *2020 IEEE Nuclear Science Symposium and Medical Imaging Conference (NSS/MIC)*. Boston, MA, USA: IEEE, Oct. 2020, pp. 1–4. ISBN: 978-1-72817-693-2. DOI: [10.1109/NSS/MIC42677.2020.9507759](https://doi.org/10.1109/NSS/MIC42677.2020.9507759). URL: <https://ieeexplore.ieee.org/document/9507759/> (visited on 04/07/2022).
- [78] Yujiao Wu et al. “DeepMMSA: A Novel Multimodal Deep Learning Method for Non-small Cell Lung Cancer Survival Analysis”. en. In: *arXiv:2106.06744 [cs, eess]* (June 2021). arXiv: 2106.06744. URL: <http://arxiv.org/abs/2106.06744> (visited on 04/04/2022).
- [79] Qiang He et al. “Feasibility study of a multi-criteria decision-making based hierarchical model for multi-modality feature and multi-classifier fusion: Applications in medical prognosis prediction”. en. In: *Information Fusion* 55 (Mar. 2020), pp. 207–219. ISSN: 15662535. DOI: [10.1016/j.inffus.2019.09.001](https://doi.org/10.1016/j.inffus.2019.09.001). URL: <https://linkinghub.elsevier.com/retrieve/pii/S1566253518309059> (visited on 04/04/2022).
- [80] Luís A. Vale-Silva and Karl Rohr. “Long-term cancer survival prediction using multimodal deep learning”. en. In: *Scientific Reports* 11.1 (Dec. 2021), p. 13505. ISSN: 2045-2322. DOI: [10.1038/s41598-021-92799-4](https://doi.org/10.1038/s41598-021-92799-4). URL: <http://www.nature.com/articles/s41598-021-92799-4> (visited on 04/04/2022).
- [81] Shaimaa Bakr et al. “Data for NSCLC Radiogenomics Collection”. In: The Cancer Imaging Archive, 2017. DOI: [10.7937/K9/TCIA.2017.7hs46erv](https://doi.org/10.7937/K9/TCIA.2017.7hs46erv).
- [82] H. J. W. L. Aerts et al. “Data From NSCLC-Radiomics”. In: The Cancer Imaging Archive, 2019. DOI: [10.7937/K9/TCIA.2015.PF0M9REI](https://doi.org/10.7937/K9/TCIA.2015.PF0M9REI).
- [83] Robert L. Grossman et al. “Toward a Shared Vision for Cancer Genomic Data”. In: *New England Journal of Medicine* 375.12 (2016). PMID: 27653561, pp. 1109–1112. DOI: [10.1056/NEJMp1607591](https://doi.org/10.1056/NEJMp1607591). URL: <https://doi.org/10.1056/NEJMp1607591>.
- [84] Ravi Aggarwal et al. “Diagnostic accuracy of deep learning in medical imaging: A systematic review and meta-analysis”. In: *NPJ digital medicine* 4.1 (2021), pp. 1–23.
- [85] Mehmet A Gulum, Christopher M Trombley, and Mehmed Kantardzic. “A review of explainable deep learning cancer detection models in medical imaging”. In: *Applied Sciences* 11.10 (2021), p. 4573.
- [86] Matteo Tortora et al. “Deep Reinforcement Learning for Fractionated Radiotherapy in Non-Small Cell Lung Carcinoma”. In: *Artificial Intelligence in Medicine* 119 (2021), p. 102137.

-
- [87] Valerio Guarrasi and Paolo Soda. “Multi-objective optimization determines when, which and how to fuse deep networks: An application to predict COVID-19 outcomes”. In: *Computers in Biology and Medicine* 154 (2023), p. 106625.
- [88] Víctor Aceña et al. “Minimally overfitted learners: A general framework for ensemble learning”. In: *Knowledge-Based Systems* 254 (2022), p. 109669. ISSN: 0950-7051. DOI: <https://doi.org/10.1016/j.knosys.2022.109669>. URL: <https://www.sciencedirect.com/science/article/pii/S0950705122008450>.
- [89] Vadim Borisov et al. “Deep neural networks and tabular data: A survey”. In: *arXiv preprint arXiv:2110.01889* (2021).
- [90] Jiuxiang Gu et al. “Recent advances in convolutional neural networks”. In: *Pattern Recognition* 77 (2018), pp. 354–377.
- [91] Alex Krizhevsky. “One weird trick for parallelizing convolutional neural networks”. In: *arXiv preprint arXiv:1404.5997* (2014).
- [92] Karen Simonyan et al. “Very deep convolutional networks for large-scale image recognition”. In: *arXiv preprint arXiv:1409.1556* (2014).
- [93] Kaiming He et al. “Deep residual learning for image recognition”. In: *IEEE Conf. on computer vision and pattern recognition*. 2016, pp. 770–778.
- [94] Gao Huang et al. “Densely connected convolutional networks”. In: *IEEE Conf. on computer vision and pattern recognition*. 2017, pp. 4700–4708.
- [95] Christian Szegedy et al. “Going deeper with convolutions”. In: *IEEE Conf. on computer vision and pattern recognition*. 2015, pp. 1–9.
- [96] Ningning Ma et al. “ShuffleNet V2: Practical guidelines for efficient CNN architecture design”. In: *European Conference on Computer Vision (ECCV)*. 2018, pp. 116–131.
- [97] Mark Sandler et al. “MobileNetV2: Inverted residuals and linear bottlenecks”. In: *IEEE Conf. on computer vision and pattern recognition*. 2018, pp. 4510–4520.
- [98] Mingxing Tan et al. “MnasNet: Platform-aware neural architecture search for mobile”. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019, pp. 2820–2828.
- [99] Jia Deng et al. “Imagenet: A large-scale hierarchical image database”. In: *2009 IEEE conference on computer vision and pattern recognition*. Ieee. 2009, pp. 248–255.
- [100] Valerio Guarrasi et al. “Pareto optimization of deep networks for COVID-19 diagnosis from chest X-rays”. In: *Pattern Recognition* 121 (2022), p. 108242.

-
- [101] Valerio Guarrasi et al. “A Multi-Expert System to Detect COVID-19 Cases in X-ray Images”. In: *2021 IEEE 34th International Symposium on Computer-Based Medical Systems (CBMS)*. IEEE. 2021, pp. 395–400.
- [102] Valerio Guarrasi and Paolo Soda. “Optimized Fusion of CNNs to Diagnose Pulmonary Diseases on Chest X-Rays”. In: *International Conference on Image Analysis and Processing*. Springer. 2022, pp. 197–209.
- [103] Thomas G Dietterich. “An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization”. In: *Machine learning* 40.2 (2000), pp. 139–157.
- [104] Olaf Ronneberger et al. “U-Net: Convolutional networks for biomedical image segmentation”. In: *International Conf. on Medical image computing and computer-assisted intervention*. Springer. 2015, pp. 234–241.
- [105] Gavin Brown et al. “Diversity creation methods: a survey and categorisation”. In: *Information fusion* 6.1 (2005), pp. 5–20.
- [106] Ivan Izonin et al. “The Additive Input-Doubling Method Based on the SVR with Nonlinear Kernels: Small Data Approach”. In: *Symmetry* 13.4 (2021). ISSN: 2073-8994. DOI: [10.3390/sym13040612](https://doi.org/10.3390/sym13040612). URL: <https://www.mdpi.com/2073-8994/13/4/612>.
- [107] Valerio Guarrasi et al. “A systematic review of intermediate fusion in multimodal deep learning for biomedical applications”. In: *Image and Vision Computing* 158 (2025), p. 105509. ISSN: 0262-8856. DOI: <https://doi.org/10.1016/j.imavis.2025.105509>.
- [108] Tadas Baltruaitis, Chaitanya Ahuja, and Louis-Philippe Morency. “Multimodal machine learning: A survey and taxonomy”. In: *IEEE transactions on pattern analysis and machine intelligence* 41.2 (2018), pp. 423–443.
- [109] Michael Moor et al. “Foundation models for generalist medical artificial intelligence”. In: *Nature* 616.7956 (2023), pp. 259–265.
- [110] Jiquan Ngiam et al. “Multimodal deep learning.” In: *ICML*. Vol. 11. 2011, pp. 689–696.
- [111] Fei Zhao, Chengcui Zhang, and Baocheng Geng. “Deep multimodal data fusion”. In: *ACM computing surveys* 56.9 (2024), pp. 1–36.
- [112] Renjie Wu et al. “Deep multimodal learning with missing modality: A survey”. In: *arXiv preprint arXiv:2409.07825* (2024).

-
- [113] Md Kaykobad Reza, Ashley Prater-Bennette, and M Salman Asif. “Robust multi-modal learning with missing modalities via parameter-efficient adaptation”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2024).
- [114] Zihui Zhang et al. “Synergistic Prompting for Robust Visual Recognition with Missing Modalities”. In: *arXiv preprint arXiv:2507.07802* (2025).
- [115] Sören Richard Stahlschmidt, Benjamin Ulfenborg, and Jane Synnergren. “Multimodal deep learning for biomedical data fusion: a review”. In: *Briefings in Bioinformatics* 23.2 (2022), bbab569.
- [116] Camillo Maria Caruso, Paolo Soda, and Valerio Guarrasi. “Not Another Imputation Method: A Transformer-based Model for Missing Values in Tabular Datasets”. In: *arXiv preprint arXiv:2407.11540* (2024).
- [117] Arianna Francesconi et al. “Class balancing diversity multimodal ensemble for Alzheimers disease diagnosis and early detection”. In: *Computerized Medical Imaging and Graphics* 123 (2025), p. 102529.
- [118] Filippo Ruffini et al. “Benchmarking Foundation Models and Parameter-Efficient Fine-Tuning for Prognosis Prediction in Medical Imaging”. In: *arXiv preprint arXiv:2506.18434* (2025).
- [119] Giulia Di Teodoro et al. “A graph neural network-based model with Out-of-Distribution Robustness for enhancing Antiretroviral Therapy Outcome Prediction for HIV-1”. In: *Computerized Medical Imaging and Graphics* (2025), p. 102484.
- [120] Valerio Guarrasi et al. “Multimodal explainability via latent shift applied to COVID-19 stratification”. In: *Pattern Recognition* (2024), p. 110825.
- [121] Mengmeng Ma et al. “Smil: Multimodal learning with severely missing modality”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 35. 2021, pp. 2302–2310.
- [122] Fatemeh Behrad and Mohammad Saniee Abadeh. “An overview of deep learning methods for multimodal medical data mining”. In: *Expert Systems with Applications* 200 (2022), p. 117006.
- [123] Marziyeh Afkanpour, Elham Hosseinzadeh, and Hamed Tabesh. “Identify the most appropriate imputation method for handling missing values in clinical structured datasets: a systematic review”. In: *BMC Medical Research Methodology* 24.1 (2024), p. 188.
- [124] Javier E Flores et al. “Missing data in multi-omics integration: Recent advances through artificial intelligence”. In: *Frontiers in Artificial Intelligence* 6 (2023), p. 1098308.

-
- [125] Shuwei Qian and Chongjun Wang. “COM: Contrastive Masked-attention model for incomplete multimodal learning”. In: *Neural Networks* 162 (2023), pp. 443–455.
- [126] Linfeng Liu et al. “Cascaded multi-modal mixing transformers for alzheimers disease classification with incomplete data”. In: *NeuroImage* 277 (2023), p. 120267.
- [127] Chaohe Zhang et al. “M3care: Learning with missing modalities in multimodal health-care data”. In: *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2022, pp. 2418–2428.
- [128] Hu Wang et al. “Multi-modal learning with missing modality via shared-specific feature modelling”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 15878–15887.
- [129] Clifford R Jack Jr et al. “The Alzheimer’s disease neuroimaging initiative (ADNI): MRI methods”. In: *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine* 27.4 (2008), pp. 685–691.
- [130] Paolo Soda et al. “AIforCOVID: Predicting the clinical outcomes in patients with COVID-19 applying AI to chest-X-rays. An Italian multicentre study”. In: *Medical Image Analysis* 74 (2021), p. 102216. ISSN: 1361-8415. DOI: <https://doi.org/10.1016/j.media.2021.102216>.
- [131] Filippo Ruffini et al. “Multi-dataset multi-task learning for COVID-19 prognosis”. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer. 2024, pp. 251–261.
- [132] Alice Natalina Caragliano et al. “Doctor-in-the-Loop: An explainable, multi-view deep learning framework for predicting pathological response in non-small cell lung cancer”. In: *Image and Vision Computing* (2025), p. 105630.
- [133] Camillo Maria Caruso et al. “A multimodal ensemble driven by multiobjective optimisation to predict overall survival in non-small-cell lung cancer”. In: *Journal of Imaging* 8.11 (2022), p. 298.
- [134] Tao Wang et al. “Multi-view imputation and cross-attention network based on incomplete longitudinal and multimodal data for conversion prediction of mild cognitive impairment”. In: *Expert Systems with Applications* 231 (2023), p. 120761.
- [135] Wei Xiong et al. “Disentanglement and codebook learning-induced feature match network to diagnose neurodegenerative diseases on incomplete multimodal data”. In: *Pattern Recognition* 165 (2025), p. 111597.

Appendices

Appendix A

Masking Example

In this section, we aim to illustrate in more depth how the proposed masked self-attention mechanism is able to mask and ignore the contribution of missing values via an example.

Let the feature vector x be composed of 4 different features, out of which the third feature is missing:

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \quad (\text{A.1})$$

Let us now compute the embedding matrix e of the features, using equations 2.4 and 2.5 for the categorical and numerical features respectively, using a dimension of the embeddings $d_e = 4$:

$$e = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & a_3 & a_4 \\ b_1 & b_2 & b_3 & b_4 \\ c_1 & c_2 & c_3 & c_4 \\ d_1 & d_2 & d_3 & d_4 \end{bmatrix} \quad (\text{A.2})$$

As x_3 was missing, we can remark that embedding c denotes the missing feature and the

mask M can be defined:

$$M = \begin{bmatrix} 0 & 0 & -\infty & 0 \\ 0 & 0 & -\infty & 0 \\ 0 & 0 & -\infty & 0 \\ 0 & 0 & -\infty & 0 \end{bmatrix} \quad (\text{A.3})$$

Then, following the steps in the transformer pipeline, we compute Q , K and V using Eq. 2.6 and a number of heads $h = 2$:

$$Q = \begin{bmatrix} a^Q \\ b^Q \\ c^Q \\ d^Q \end{bmatrix} = \begin{bmatrix} a_1^Q & a_2^Q \\ b_1^Q & b_2^Q \\ c_1^Q & c_2^Q \\ d_1^Q & d_2^Q \end{bmatrix}, \quad K = \begin{bmatrix} a^K \\ b^K \\ c^K \\ d^K \end{bmatrix} = \begin{bmatrix} a_1^K & a_2^K \\ b_1^K & b_2^K \\ c_1^K & c_2^K \\ d_1^K & d_2^K \end{bmatrix}, \quad V = \begin{bmatrix} a^V \\ b^V \\ c^V \\ d^V \end{bmatrix} = \begin{bmatrix} a_1^V & a_2^V \\ b_1^V & b_2^V \\ c_1^V & c_2^V \\ d_1^V & d_2^V \end{bmatrix} \quad (\text{A.4})$$

We now compute the QK^T product to better understand how the contributions from various features (a , b , c and d) are distributed across the attention matrix.

$$QK^T = \begin{bmatrix} a^Q a^K & a^Q b^K & a^Q c^K & a^Q d^K \\ b^Q a^K & b^Q b^K & b^Q c^K & b^Q d^K \\ c^Q a^K & c^Q b^K & c^Q c^K & c^Q d^K \\ d^Q a^K & d^Q b^K & d^Q c^K & d^Q d^K \end{bmatrix} \quad (\text{A.5})$$

Now, using the QK^T product, the vector V and the mask M , we can calculate the output of an attention head for the classical masked self-attention mechanism applied to

tabular data:

$$\begin{aligned}
 & \text{softmax}(QK^T + M)V = \\
 & = \text{softmax} \left(\begin{bmatrix} a^Q a^K & a^Q b^K & a^Q c^K & a^Q d^K \\ b^Q a^K & b^Q b^K & b^Q c^K & b^Q d^K \\ c^Q a^K & c^Q b^K & c^Q c^K & c^Q d^K \\ d^Q a^K & d^Q b^K & d^Q c^K & d^Q d^K \end{bmatrix} + \begin{bmatrix} 0 & 0 & -\infty & 0 \\ 0 & 0 & -\infty & 0 \\ 0 & 0 & -\infty & 0 \\ 0 & 0 & -\infty & 0 \end{bmatrix} \right) \begin{bmatrix} a_1^V & a_2^V \\ b_1^V & b_2^V \\ c_1^V & c_2^V \\ d_1^V & d_2^V \end{bmatrix} = \\
 & = \begin{bmatrix} \text{softmax}(a^Q a^K & a^Q b^K & -\infty & a^Q d^K) \\ \text{softmax}(b^Q a^K & b^Q b^K & -\infty & b^Q d^K) \\ \text{softmax}(c^Q a^K & c^Q b^K & -\infty & c^Q d^K) \\ \text{softmax}(d^Q a^K & d^Q b^K & -\infty & d^Q d^K) \end{bmatrix} \begin{bmatrix} a_1^V & a_2^V \\ b_1^V & b_2^V \\ c_1^V & c_2^V \\ d_1^V & d_2^V \end{bmatrix} = \\
 & = \begin{bmatrix} (a^Q a^K)_S & (a^Q b^K)_S & 0 & (a^Q d^K)_S \\ (b^Q a^K)_S & (b^Q b^K)_S & 0 & (b^Q d^K)_S \\ (c^Q a^K)_S & (c^Q b^K)_S & 0 & (c^Q d^K)_S \\ (d^Q a^K)_S & (d^Q b^K)_S & 0 & (d^Q d^K)_S \end{bmatrix} \begin{bmatrix} a_1^V & a_2^V \\ b_1^V & b_2^V \\ c_1^V & c_2^V \\ d_1^V & d_2^V \end{bmatrix} = \\
 & = \begin{bmatrix} (a^Q a^K)_S \cdot a_1^V + (a^Q b^K)_S \cdot b_1^V + (a^Q d^K)_S \cdot d_1^V + 0 & (a^Q a^K)_S \cdot a_2^V + (a^Q b^K)_S \cdot b_2^V + (a^Q d^K)_S \cdot d_2^V + 0 \\ (b^Q a^K)_S \cdot a_1^V + (b^Q b^K)_S \cdot b_1^V + (b^Q d^K)_S \cdot d_1^V + 0 & (b^Q a^K)_S \cdot a_2^V + (b^Q b^K)_S \cdot b_2^V + (b^Q d^K)_S \cdot d_2^V + 0 \\ (c^Q a^K)_S \cdot a_1^V + (c^Q b^K)_S \cdot b_1^V + (c^Q d^K)_S \cdot d_1^V + 0 & (c^Q a^K)_S \cdot a_2^V + (c^Q b^K)_S \cdot b_2^V + (c^Q d^K)_S \cdot d_2^V + 0 \\ (d^Q a^K)_S \cdot a_1^V + (d^Q b^K)_S \cdot b_1^V + (d^Q d^K)_S \cdot d_1^V + 0 & (d^Q a^K)_S \cdot a_2^V + (d^Q b^K)_S \cdot b_2^V + (d^Q d^K)_S \cdot d_2^V + 0 \end{bmatrix}
 \end{aligned} \tag{A.6}$$

As can be seen, some contributions of the missing feature are still present on the third line of the output matrix. Using as a starting point, for the sake of simplicity, the attention matrix computed on the third step of the Eq. A.6, let us compute the output matrix obtained

using the presented masked self-attention mechanism:

$$\begin{aligned}
 & ReLU(\text{softmax}(QK^T + M) + M^T)V = \\
 & = ReLU \left(\begin{bmatrix} (a^Q a^K)_S & (a^Q b^K)_S & 0 & (a^Q d^K)_S \\ (b^Q a^K)_S & (b^Q b^K)_S & 0 & (b^Q d^K)_S \\ (c^Q a^K)_S & (c^Q b^K)_S & 0 & (c^Q d^K)_S \\ (d^Q a^K)_S & (d^Q b^K)_S & 0 & (d^Q d^K)_S \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -\infty & -\infty & -\infty & -\infty \\ 0 & 0 & 0 & 0 \end{bmatrix} \right) \begin{bmatrix} a_1^V & a_2^V \\ b_1^V & b_2^V \\ c_1^V & c_2^V \\ d_1^V & d_2^V \end{bmatrix} = \\
 & = ReLU \left(\begin{bmatrix} (a^Q a^K)_S & (a^Q b^K)_S & 0 & (a^Q d^K)_S \\ (b^Q a^K)_S & (b^Q b^K)_S & 0 & (b^Q d^K)_S \\ -\infty & -\infty & -\infty & -\infty \\ (d^Q a^K)_S & (d^Q b^K)_S & 0 & (d^Q d^K)_S \end{bmatrix} \begin{bmatrix} a_1^V & a_2^V \\ b_1^V & b_2^V \\ c_1^V & c_2^V \\ d_1^V & d_2^V \end{bmatrix} \right) = \\
 & = \begin{bmatrix} (a^Q a^K)_S & (a^Q b^K)_S & 0 & (a^Q d^K)_S \\ (b^Q a^K)_S & (b^Q b^K)_S & 0 & (b^Q d^K)_S \\ 0 & 0 & 0 & 0 \\ (d^Q a^K)_S & (d^Q b^K)_S & 0 & (d^Q d^K)_S \end{bmatrix} \begin{bmatrix} a_1^V & a_2^V \\ b_1^V & b_2^V \\ c_1^V & c_2^V \\ d_1^V & d_2^V \end{bmatrix} = \\
 & = \begin{bmatrix} (a^Q a^K)_S \cdot a_1^V + (a^Q b^K)_S \cdot b_1^V + (a^Q d^K)_S \cdot d_1^V + 0 & (a^Q a^K)_S \cdot a_2^V + (a^Q b^K)_S \cdot b_2^V + (a^Q d^K)_S \cdot d_2^V + 0 \\ (b^Q a^K)_S \cdot a_1^V + (b^Q b^K)_S \cdot b_1^V + (b^Q d^K)_S \cdot d_1^V + 0 & (b^Q a^K)_S \cdot a_2^V + (b^Q b^K)_S \cdot b_2^V + (b^Q d^K)_S \cdot d_2^V + 0 \\ 0 & 0 \\ (d^Q a^K)_S \cdot a_1^V + (d^Q b^K)_S \cdot b_1^V + (d^Q d^K)_S \cdot d_1^V + 0 & (d^Q a^K)_S \cdot a_2^V + (d^Q b^K)_S \cdot b_2^V + (d^Q d^K)_S \cdot d_2^V + 0 \end{bmatrix}
 \end{aligned} \tag{A.7}$$

By comparing the final outputs, we can note that our method allows the contributions of the missing feature to be completely masked, making the most of the available data.

Appendix B

Training of the Models

For NAIM, we specified its architecture with an embedding dimension $d_e = 6$, encoder layers $L = 6$, and attention heads $h = 3$ [13]. We configured the feed-forward layer size to 1000 neurons, and we eliminated the bias in the embedding of numerical and categorical features, to represent missing values as zero vectors.

During the training of DL models, since all the datasets work on a classification task, we opted to use the cross-entropy as loss function. To enhance the training process, we employed a dynamic learning rate schedule that reduced the rate by a factor of 10 whenever the loss stagnates for 25 epochs. The training incorporated the Adam optimization algorithm with a consistent batch size of 32 across all models and datasets. For all the DL models we used a Glorot uniform initialization of the weights, which were trained all the models for 1500 epochs, using an early stopping rule with patience of 50 epochs. We set the initial 50 epochs as a warm-up period for parameter tuning for all models, except GRAPE for which 500 warm-up epochs are assigned due to its lack of batch sampling. Furthermore, we incorporated L1 and L2 regularization techniques to prevent overfitting.

All models and techniques under consideration are evaluated using their publicly available implementations. For XGBoost¹, GRAPE², TabNet³, TabTransformer³, and FTTransformer³, these implementations can be found on GitHub, whereas for the remaining models and imputation techniques, we used their implementations available in the Sklearn library⁴. For all the competitors, we set their default parameters or those provided in the public implementations by authors because we are not interested in fine-tuning the models. Therefore, for the same reason, even in the NAIM implementation we did not perform parameter

¹<https://github.com/dmlc/xgboost/tree/master>

²<https://github.com/maxiaoba/GRAPE/tree/master>

³<https://github.com/dreamquark-ai/tabnet>

³<https://github.com/lucidrains/tab-transformer-pytorch>

⁴<https://scikit-learn.org/stable>

tuning, but we selected them according to the dimensionality of the data used. Indeed, in [52], confirming the "No Free Lunch" theorem, the authors empirically observe that in many cases the use of tuned parameters cannot significantly outperform the default values of a classifier suggested in the literature.

Appendix D

NAIM Statistical analysis

Table D.1 presents the percentages of experiments in which NAIM is significantly better than the competitors. Additionally, the averages across rows and columns are provided at the right and bottom margins, respectively, to identify an average trend. The analysis of the row and column averages reveals that the win rate of NAIM significantly exceeds the loss rate, with a considerable minimum advantage of 7.8% over the competitors. This observation, reinforced by the data in the bottom right-hand corner, which indicates that NAIM, on average, wins against its competitors in 58.7% of the cases, while it only loses in 1.6% of the cases, underlines NAIM's ability to provide competitive performance on tabular data without the need for any imputation strategy.

Model	Strategy	Datasets										Mean	
		ADULT		BankMarketing		OnlineShoppers		SeismicBumps		Spambase		% Win	% Loss
		% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss
Adaboost	Constant	19.4	11.1	13.9	0.0	69.4	0.0	0.0	0.0	38.9	0.0	28.3	2.2
	KNN	36.1	2.8	19.4	0.0	100.0	0.0	0.0	0.0	75.0	0.0	46.1	0.6
	MICE	55.6	0.0	8.3	2.8	97.2	0.0	13.9	0.0	86.1	0.0	52.2	0.6
Decision Tree	Constant	100.0	0.0	100.0	0.0	100.0	0.0	97.2	0.0	100.0	0.0	99.4	0.0
	KNN	100.0	0.0	100.0	0.0	100.0	0.0	100.0	0.0	100.0	0.0	100.0	0.0
	MICE	100.0	0.0	100.0	0.0	100.0	0.0	97.2	0.0	100.0	0.0	99.4	0.0
	Intrinsic	100.0	0.0	100.0	0.0	97.2	0.0	100.0	0.0	100.0	0.0	99.4	0.0
FTTransformer	Constant	100.0	0.0	19.4	0.0	52.8	0.0	0.0	0.0	100.0	0.0	54.4	0.0
	KNN	100.0	0.0	38.9	0.0	94.4	0.0	0.0	0.0	100.0	0.0	66.7	0.0
	MICE	100.0	0.0	19.4	0.0	80.6	0.0	0.0	0.0	100.0	0.0	60.0	0.0
GRAPE	Intrinsic	100.0	0.0	88.9	0.0	94.4	0.0	0.0	0.0	100.0	0.0	76.7	0.0
HistGradientBoost	Constant	100.0	0.0	2.8	19.4	88.9	0.0	5.6	0.0	19.4	0.0	43.3	3.9
	KNN	100.0	0.0	8.3	13.9	100.0	0.0	0.0	0.0	36.1	2.8	48.9	3.3
	MICE	100.0	0.0	8.3	8.3	100.0	0.0	0.0	0.0	38.9	0.0	49.4	1.7
	Intrinsic	8.3	80.6	0.0	13.9	19.4	0.0	2.8	0.0	5.6	2.8	7.2	19.4
MLP	Constant	100.0	0.0	100.0	0.0	100.0	0.0	36.1	0.0	88.9	0.0	85.0	0.0
	KNN	100.0	0.0	100.0	0.0	100.0	0.0	27.8	0.0	83.3	0.0	82.2	0.0
	MICE	100.0	0.0	100.0	0.0	100.0	0.0	36.1	0.0	97.2	0.0	86.7	0.0
Random Forest	Constant	100.0	0.0	16.7	0.0	50.0	0.0	2.8	0.0	44.4	0.0	42.8	0.0
	KNN	100.0	0.0	16.7	0.0	88.9	0.0	0.0	0.0	44.4	2.8	50.0	0.6
	MICE	100.0	0.0	11.1	0.0	75.0	0.0	0.0	0.0	52.8	0.0	47.8	0.0
	Intrinsic	66.7	0.0	19.4	0.0	13.9	0.0	2.8	0.0	2.8	2.8	21.1	0.6
SVM	Constant	100.0	0.0	72.2	0.0	88.9	0.0	0.0	0.0	77.8	0.0	67.8	0.0
	KNN	100.0	0.0	52.8	0.0	100.0	0.0	0.0	0.0	86.1	0.0	67.8	0.0
	MICE	100.0	0.0	63.9	0.0	88.9	0.0	0.0	0.0	72.2	0.0	65.0	0.0
TabNet	Constant	100.0	0.0	13.9	5.6	72.2	0.0	0.0	0.0	100.0	0.0	57.2	1.1
	KNN	100.0	0.0	25.0	5.6	100.0	0.0	0.0	0.0	97.2	0.0	64.4	1.1
	MICE	100.0	0.0	38.9	0.0	77.8	0.0	0.0	0.0	100.0	0.0	63.3	0.0
TabTransformer	Constant	100.0	0.0	13.9	0.0	30.6	0.0	0.0	0.0	94.4	0.0	47.8	0.0
	KNN	100.0	0.0	13.9	0.0	80.6	0.0	0.0	0.0	83.3	0.0	55.6	0.0
	MICE	100.0	0.0	38.9	0.0	80.6	0.0	0.0	0.0	88.9	0.0	61.7	0.0
XGBoost	Constant	100.0	0.0	0.0	11.1	22.2	0.0	0.0	0.0	41.7	0.0	32.8	2.2
	KNN	100.0	0.0	16.7	5.6	88.9	0.0	0.0	0.0	38.9	0.0	48.9	1.1
	MICE	100.0	0.0	22.2	8.3	83.3	0.0	0.0	0.0	63.9	0.0	53.9	1.7
	Intrinsic	16.7	63.9	13.9	8.3	25.0	0.0	22.2	0.0	33.3	0.0	22.2	14.4
Mean		88.7	4.5	39.4	2.9	78.9	0.0	15.6	0.0	71.2	0.3	58.7	1.6

Table D.1: Percentages of wins and losses of NAIM correct predictions compared with the competitors using the Wilcoxon signed-rank test.

Appendix E

NAIM Computational analysis

To provide a quantitative comparison of the computational complexity of the evaluated DL approaches, we computed the number of floating-point operations (FLOPs) and the total number of trainable parameters for each model. It is important to note that ML models were excluded from this analysis, as their training and inference procedures are inherently less computationally demanding and not directly comparable with DL architectures. Similarly, the computational contribution of the imputation methods was not considered, since imputers mainly affect the preprocessing and training stages, whereas their cost during inference is negligible and varies significantly across implementations. The results are summarized in Table E.1, which reports the average number of operations and parameters across the five datasets. The reported FLOPs correspond to the total number of operations required for a single forward pass, while the number of parameters reflects the count of all trainable weights within the model.

Model	Adult		BankMarketing		OnlineShoppers		SeismicBumps		Spambase	
	FLOPs	# Params	FLOPs	# Params	FLOPs	# Params	FLOPs	# Params	FLOPs	# Params
NAIM	1.03×10^6	7.99×10^4	1.48×10^6	7.98×10^4	1.26×10^6	7.98×10^4	1.33×10^6	7.95×10^4	4.29×10^6	8.04×10^4
FTTransformer	2.96×10^6	1.78×10^5	4.33×10^6	1.77×10^5	3.63×10^6	1.77×10^5	3.86×10^6	1.76×10^5	1.53×10^7	1.78×10^5
GRAPE	3.78×10^4	3.73×10^3	5.73×10^4	3.93×10^3	4.73×10^4	3.83×10^3	5.06×10^4	3.86×10^3	2.28×10^5	5.15×10^3
MLP	2.12×10^4	1.06×10^4	1.28×10^4	6.40×10^3	1.52×10^4	7.60×10^3	4.80×10^3	2.40×10^3	1.18×10^4	5.90×10^3
TabNet	1.00×10^4	6.67×10^3	1.16×10^4	7.02×10^3	1.08×10^4	6.84×10^3	1.10×10^4	6.85×10^3	2.61×10^4	9.41×10^3
TabTransformer	1.53×10^6	2.20×10^5	1.96×10^6	2.44×10^5	1.33×10^6	2.11×10^5	7.32×10^5	1.87×10^5	2.30×10^3	1.77×10^5

Table E.1: Computational complexity analysis of the evaluated deep learning models in terms of FLOPs and number of trainable parameters across datasets.

As expected, transformer-based models (NAIM, FTTransformer, and TabTransformer) exhibit higher computational demands compared to simpler architectures such as MLP and TabNet, reflecting the additional cost introduced by the self-attention mechanism and feature embedding layers. Nevertheless, NAIM maintains a balanced trade-off between expressivity and efficiency, achieving competitive FLOPs and parameter counts despite its specialized

mechanisms for handling missing data. Overall, this analysis confirms that NAIMs ability to robustly manage incomplete tabular data does not come at a prohibitive computational cost. In fact, it remains within the same order of magnitude as comparable transformer-based approaches while significantly outperforming them in predictive performance, as discussed in the main text.

APPENDIX F. MARIA COMPLETE RESULTS

AUC		Train missing percentage: Ω				Train missing percentage: 30%				Train missing percentage: 50%				Train missing percentage: 75%			
Approach	Model	Test missing percentage:				Test missing percentage:				Test missing percentage:				Test missing percentage:			
		Ω	30%	50%	75%	Ω	30%	50%	75%	Ω	30%	50%	75%	Ω	30%	50%	75%
Early ML with imputer	AdaBoost	83.89	77.53	82.18	78.04	76.65	76.10	73.36	68.80	84.44	78.21	81.80	78.00	81.66	72.82	80.07	76.00
	DecisionTree	76.40	65.40	74.30	67.49	65.10	65.63	60.00	58.39	74.89	62.44	73.13	67.49	68.93	58.40	68.47	65.15
	HistGradientBoost	88.05	78.35	86.40	80.06	79.34	78.84	75.44	69.79	87.90	77.65	85.73	79.81	85.73	74.23	83.66	78.90
	RandomForest	88.75	80.64	86.95	81.84	80.97	80.64	77.47	70.50	88.13	78.84	86.35	81.36	85.85	75.24	84.33	78.91
	SVM	86.43	75.79	84.62	79.75	76.49	75.93	72.47	65.13	86.66	74.99	84.51	79.54	85.57	74.34	83.10	78.37
Early ML without imputer	XGBoost	88.00	78.30	86.12	79.19	78.65	78.14	73.92	68.52	87.15	76.14	85.15	78.99	84.42	71.01	82.86	78.53
	DecisionTree	78.33	64.65	75.27	65.65	66.07	66.74	64.49	55.37	76.77	62.93	75.63	67.68	71.45	59.06	70.62	69.75
	HistGradientBoost	89.99	80.59	87.27	78.82	81.67	80.54	75.29	67.26	89.93	80.70	88.07	82.05	89.10	78.54	87.56	82.62
	RandomForest	90.78	82.27	89.22	83.63	82.96	82.35	78.85	72.10	90.17	81.20	89.04	84.85	88.09	78.87	87.12	83.94
	XGBoost	90.07	79.44	86.55	77.30	80.16	79.14	75.04	68.34	89.22	78.17	87.83	82.61	86.89	74.42	86.76	84.16
Early DL with imputer	FTTransformer	87.10	79.70	85.50	80.29	79.19	78.76	74.88	69.07	86.34	77.20	84.49	79.37	83.52	71.44	80.79	76.96
	MLP	85.08	75.44	83.14	77.86	72.99	73.90	71.12	65.51	83.94	72.07	82.60	77.55	82.80	72.24	80.88	76.93
	TabNet	85.11	75.86	84.11	78.72	76.03	75.67	74.13	68.04	84.93	75.74	82.95	77.78	80.52	68.61	79.22	73.15
	TABTransformer	84.87	73.84	83.03	77.86	72.35	72.23	69.82	64.22	84.24	73.18	82.53	78.17	82.43	71.14	80.54	75.68
	XGBoost	76.82	69.54	74.18	69.98	72.40	70.00	65.30	60.30	68.47	64.18	66.80	64.03	50.00	50.00	50.00	50.00
Intermediate DL with imputer	MLP	86.61	78.94	84.96	81.16	80.63	79.60	75.42	69.67	82.58	67.02	81.44	78.70	75.61	50.00	75.61	73.68
	TabNet	75.58	69.67	73.57	70.89	69.88	68.47	64.16	59.89	73.29	63.37	72.14	70.10	67.17	62.60	67.72	63.89
	TABTransformer	85.32	75.91	83.66	79.54	73.95	72.88	69.73	63.43	83.90	73.93	82.22	78.21	81.46	67.70	79.43	76.11
	MARIA	90.80	80.60	87.81	82.75	82.31	81.29	74.02	65.10	90.24	80.06	87.80	82.13	88.21	74.30	85.58	81.07
	AdaBoost	83.94	78.32	81.73	76.69	77.61	77.39	73.46	67.79	83.50	76.31	81.90	78.41	81.86	73.34	80.38	77.72
Late ML with imputer	DecisionTree	77.36	66.37	75.27	70.16	67.78	67.27	64.59	57.76	76.02	65.82	74.34	70.11	72.44	60.34	71.09	67.16
	HistGradientBoost	83.29	73.04	80.87	75.24	73.52	73.26	69.43	64.36	82.30	71.83	80.34	75.68	80.74	69.54	78.87	75.05
	RandomForest	84.92	71.79	82.80	78.52	73.02	71.86	68.06	63.66	84.12	70.63	82.52	77.68	81.98	65.25	80.92	78.38
	SVM	88.13	78.98	86.20	81.69	80.33	79.16	76.01	70.70	87.82	78.98	85.96	82.00	86.78	78.21	85.44	81.93
	XGBoost	87.65	77.55	85.37	79.86	78.79	78.32	74.07	67.20	87.32	77.01	85.18	80.59	85.06	72.99	83.55	79.69
Late ML without imputer	DecisionTree	83.18	70.14	79.75	69.13	71.52	70.37	65.09	59.79	80.85	68.42	78.47	70.73	77.18	65.68	75.43	73.72
	HistGradientBoost	84.66	73.00	80.95	73.02	74.27	73.65	66.82	62.13	84.25	73.65	81.71	74.83	82.48	71.05	79.93	75.12
	RandomForest	88.70	77.01	86.52	80.82	79.67	78.78	74.39	69.19	88.79	79.40	87.22	82.78	86.87	76.65	85.64	83.03
	XGBoost	89.05	78.21	86.13	76.76	80.63	80.26	75.00	68.33	88.34	77.70	87.20	81.87	86.68	75.02	85.97	83.27
	FTTransformer	87.78	79.09	85.87	81.40	80.98	80.66	76.09	71.17	86.64	76.37	84.79	80.89	85.20	73.03	83.96	80.11
Late DL with imputer	MLP	87.13	79.12	85.09	81.02	81.01	80.29	76.56	71.65	86.76	78.90	84.89	81.23	85.31	77.10	83.98	80.23
	TabNet	87.12	78.08	85.03	81.09	78.46	77.40	73.67	69.37	86.43	76.42	84.90	80.81	83.60	72.95	82.20	79.62
	TABTransformer	85.46	75.05	83.28	79.44	77.69	77.03	72.83	67.32	85.03	75.37	83.20	79.77	82.54	70.90	80.95	77.77

MCC		Train missing percentage: Ω				Train missing percentage: 30%				Train missing percentage: 50%				Train missing percentage: 75%			
Fusion Strategy	Model	Test missing percentage:				Test missing percentage:				Test missing percentage:				Test missing percentage:			
		Ω	30%	50%	75%	Ω	30%	50%	75%	Ω	30%	50%	75%	Ω	30%	50%	75%
Early ML with imputer	AdaBoost	54.79	42.40	51.87	46.94	44.12	44.11	40.29	30.76	53.53	45.20	50.99	42.47	51.08	42.22	46.91	37.64
	DecisionTree	47.80	38.64	47.23	40.77	38.40	36.56	36.00	32.32	43.48	32.36	43.23	37.08	33.69	24.19	33.26	28.50
	HistGradientBoost	58.61	47.79	53.01	44.53	47.86	45.95	40.78	32.98	57.44	46.68	54.51	44.51	54.05	41.73	50.22	40.28
	RandomForest	54.79	38.09	51.59	44.54	39.81	39.01	38.88	32.20	53.59	38.16	50.69	42.98	44.33	25.86	44.49	37.26
	SVM	57.10	41.61	53.19	44.17	43.28	41.26	39.00	29.67	54.38	37.35	51.54	42.02	52.83	34.05	50.11	37.68
Early ML without imputer	XGBoost	57.46	44.66	53.51	45.20	46.22	46.28	41.62	32.04	56.68	45.40	52.90	45.54	53.47	40.92	51.37	40.75
	DecisionTree	50.50	41.29	43.22	28.99	41.48	39.92	34.16	20.07	45.78	37.01	44.76	37.41	35.60	26.52	36.33	33.04
	HistGradientBoost	57.42	43.88	49.58	29.49	47.14	45.64	41.46	22.76	57.88	44.28	55.38	44.82	56.04	44.07	54.26	45.10
	RandomForest	56.06	38.42	50.28	33.00	42.12	40.37	31.33	20.58	52.59	34.93	51.20	41.21	46.08	25.57	45.29	42.13
	XGBoost	58.14	46.00	47.47	27.78	48.00	45.58	40.06	18.42	56.21	41.52	54.83	42.64	53.35	40.86	53.05	45.43
Early DL with imputer	FTTransformer	54.17	38.05	50.58	43.36	37.30	36.66	31.88	28.88	52.11	32.33	50.00	42.52	47.66	30.32	46.48	38.45
	MLP	51.21	34.08	48.72	42.34	36.02	35.18	32.03	26.29	48.59	28.66	47.22	40.75	44.42	25.81	44.60	36.96
	TabNet	46.36	30.30	43.46	37.48	31.89	30.90	28.81	21.53	45.93	29.30	43.93	35.34	36.05	23.08	36.13	27.17
	TABTransformer	54.49	42.26	50.37	42.80	39.70	38.38	32.66	26.18	50.57	34.41	48.52	40.81	46.82	28.53	44.49	37.93
	FTTransformer	1.51	3.78	0.95	0.54	1.06	1.06	0.34	0.36	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Intermediate DL with imputer	MLP	49.13	30.24	49.06	43.20	30.27	30.27	29.74	25.54	49.06	29.98	47.34	40.56	44.52	26.83	43.10	37.52
	TabNet	3.96	5.72	3.85	1.51	0.90	2.22	1.12	1.13	4.68	3.10	3.34	1.10	1.18	2.51	1.65	3.08
	TABTransformer	49.81	32.64	48.40	42.15	34.27	34.99	31.17	22.42	47.44	31.33	45.46	38.93	41.22	24.66	39.51	35.37
	MARIA	59.95	47.06	53.88	42.60	47.77	44.40	38.09	25.97	59.57	48.00	55.39	40.96	56.42	41.22	52.40	42.27
	AdaBoost	46.42	30.37	45.47	41.85	31.14	30.32	28.16	25.78	45.78	32.32	44.54	38.01	37.39	29.56	37.58	32.85
Late ML with imputer	DecisionTree	30.16	27.84	27.18	21.15	28.41	27.06	22.40	16.46	27.99	25.08	25.78	20.66	18.89	17.24	19.61	16.24
	HistGradientBoost	35.13	22.53	32.12	26.85	25.70	26.32	22.32	13.38	35.17	27.19	30.50	22.35	30.55	19.08	26.95	22.00
	RandomForest	43.04	32.86	41.00	33.13	33.82	33.07	29.92	19.00	41.99	30.20	39.18	30.88	34.38	22.08	30.81	24.97
	SVM	46.09	27.90	44.58	35.87	27.80	27.90	26.28	17.45	45.15	27.70	43.31	33.60	42.64	26.28	39.93	31.09
	XGBoost	44.84	27.55	42.74	37.98	30.74	29.78	24.76	23.25	46.09	29.40	43.86	33.62	42.39	28.90	40.64	33.60
Late ML without imputer	DecisionTree	38.39	35.36	31.03	21.85	34.26	33.29	25.21	16.56	33.85	32.32	31.72	21.83	23.22	20.32	23.73	23.37
	HistGradientBoost	38.64	27.10	29.57	16.51	30.44	24.6										

Appendix G

MARIA Statistical Analysis

To assess the statistical significance of the differences between MARIA and competing methods, we performed a Wilcoxon signed-rank test based on binary correctness vectors (correct vs. incorrect predictions) on the concatenated test sets of the five folds, i.e., the entire dataset. Each test was conducted independently for every experimental setting, defined by a specific combination of missingness levels in training and test sets. For each competitor, we computed the rate of experiments in which MARIA performed significantly better ($p < 0.05$) or significantly worse ($p < 0.05$). The results for the “missing modalities” and “all missing” settings are summarized in Tables G.1 and G.2, respectively.

Across both scenarios, MARIA consistently outperforms models relying on late fusion strategies, regardless of whether they rely on ML or DL backbones, and whether they use imputation or not. This is particularly noteworthy given that late fusion architectures are, by design, naturally equipped to handle missing modalities. However, our results show that such models often fail to fully exploit the complementary information available across modalities, leading to inferior performance compared to MARIA.

When compared against other intermediate fusion methods, MARIA achieves statistically significant wins in more than 60% of the experiments, and loses at most in 5% of cases, in the most challenging all missing setting. These results underscore the models robustness at all levels of incompleteness, effectively handling both missing features and entirely missing modalities.

Early fusion approaches present a more complex scenario. MARIA consistently outperforms both early fusion models that rely on imputation, whether based on ML or DL backbones, although the margin is narrower compared to what is observed against late and intermediate fusion strategies. Notably, early fusion methods that do not employ imputation emerge as some of the most competitive baselines: in the challenging all missing setting, MARIA still achieves statistically significant wins in over 30% of the experiments and loses

in fewer than 10%, yet the gap is considerably smaller than with other competitor groups. This result highlights the robustness of deterministic, imputation-free early fusion strategies under certain conditions, while at the same time reinforcing MARIA’s superior flexibility in handling both partially and fully missing data, thanks to its attention-based architecture, which dynamically adjusts to the available information at both the feature and modality levels.

Overall, these results confirm that MARIA maintains a clear statistical advantage across all fusion strategies and backbone types, with the most notable gains observed against late fusion and imputation-dependent methods. Its ability to simultaneously manage missing features and entire modalities without relying on external imputers is a key factor behind its consistent performance across a wide range of experimental settings.

Approach	Model	ADNI Diagnosis				ADNI Prognosis				AlforCovid				Mean					
		Binary		Multiclass		12m		24m		36m		48m		Death		Prognosis		% Win	% Loss
		% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss
Early ML with imputer	AdaBoost	28.00	0.00	100.00	0.00	100.00	0.00	92.00	0.00	64.00	0.00	96.00	0.00	76.00	0.00	80.00	4.00	79.50	0.50
	DecisionTree	92.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	92.00	0.00	100.00	0.00	100.00	0.00	96.00	0.00	97.50	0.00
	HistGradientBoost	24.00	0.00	100.00	0.00	24.00	0.00	32.00	0.00	32.00	0.00	48.00	0.00	4.00	0.00	84.00	4.00	43.50	0.50
	RandomForest	12.00	0.00	100.00	0.00	8.00	12.00	20.00	0.00	4.00	12.00	20.00	4.00	12.00	0.00	28.00	4.00	25.50	4.00
	SVM	4.00	0.00	100.00	0.00	24.00	4.00	32.00	0.00	28.00	8.00	36.00	0.00	0.00	0.00	84.00	4.00	38.50	1.50
	XGBoost	12.00	0.00	100.00	0.00	0.00	8.00	16.00	0.00	8.00	0.00	44.00	0.00	4.00	0.00	64.00	4.00	31.00	1.50
	Mean	28.67	0.00	100.00	0.00	42.67	4.00	48.67	0.00	38.00	3.33	57.33	0.67	32.67	0.00	72.67	2.67	52.58	1.33
Early ML without imputer	DecisionTree	68.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	96.00	4.00	95.50	0.50
	HistGradientBoost	32.00	0.00	84.00	0.00	28.00	0.00	44.00	0.00	32.00	0.00	44.00	0.00	20.00	0.00	68.00	0.00	44.00	0.00
	RandomForest	4.00	0.00	76.00	0.00	8.00	12.00	8.00	0.00	0.00	4.00	8.00	4.00	16.00	4.00	0.00	4.00	15.00	3.50
	XGBoost	12.00	0.00	20.00	0.00	0.00	8.00	20.00	0.00	4.00	4.00	24.00	0.00	0.00	0.00	0.00	4.00	10.00	2.00
	Mean	29.00	0.00	70.00	0.00	34.00	5.00	43.00	0.00	34.00	2.00	44.00	1.00	34.00	1.00	41.00	3.00	41.12	1.50
Early DL with imputer	FTTransformer	24.00	4.00	100.00	0.00	40.00	0.00	28.00	0.00	32.00	0.00	48.00	0.00	0.00	0.00	40.00	4.00	39.00	1.00
	MLP	28.00	0.00	100.00	0.00	80.00	0.00	92.00	0.00	72.00	0.00	68.00	0.00	0.00	0.00	76.00	0.00	64.50	0.00
	TabNet	84.00	0.00	100.00	0.00	4.00	0.00	84.00	0.00	76.00	0.00	100.00	0.00	0.00	0.00	96.00	0.00	68.00	0.00
	TABTransformer	24.00	0.00	100.00	0.00	48.00	0.00	92.00	0.00	72.00	0.00	60.00	0.00	4.00	0.00	76.00	4.00	59.50	0.50
	Mean	40.00	1.00	100.00	0.00	43.00	0.00	74.00	0.00	63.00	0.00	69.00	0.00	1.00	0.00	72.00	2.00	57.75	0.38
Intermediate DL with imputer	FTTransformer	96.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	0.00	0.00	96.00	0.00	86.50	0.00
	MLP	0.00	4.00	100.00	0.00	20.00	0.00	52.00	0.00	28.00	0.00	64.00	0.00	0.00	0.00	84.00	0.00	43.50	0.50
	TabNet	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	96.00	0.00	100.00	0.00	0.00	0.00	96.00	0.00	86.50	0.00
	TABTransformer	16.00	0.00	100.00	0.00	48.00	0.00	92.00	0.00	76.00	0.00	96.00	0.00	0.00	0.00	76.00	0.00	63.00	0.00
	Mean	53.00	1.00	100.00	0.00	67.00	0.00	86.00	0.00	75.00	0.00	90.00	0.00	0.00	0.00	88.00	0.00	69.88	0.12
Late ML with imputer	AdaBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	DecisionTree	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	HistGradientBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	RandomForest	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	SVM	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	XGBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	Mean	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
Late ML without imputer	DecisionTree	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	HistGradientBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	RandomForest	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	XGBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	Mean	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
Late DL with imputer	FTTransformer	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	MLP	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	TabNet	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	TABTransformer	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	Mean	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00

Table G.1: Percentages of wins and losses of MARIA correct predictions compared with the competitors in the “missing modalities” scenario using the Wilcoxon signed-rank test.

APPENDIX G. MARIA STATISTICAL ANALYSIS

Approach	Model	ADNI Diagnosis				ADNI Prognosis								AlforCovid				Mean	
		Binary		Multiclass		12m		24m		36m		48m		Death		Prognosis		% Win	% Loss
		% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss	% Win	% Loss
Early ML with imputer	AdaBoost	33.33	0.00	100.00	0.00	100.00	0.00	77.78	0.00	37.50	0.00	44.44	0.00	31.25	0.00	18.75	12.50	55.38	1.56
	DecisionTree	88.89	0.00	100.00	0.00	100.00	0.00	88.89	0.00	87.50	0.00	88.89	0.00	100.00	0.00	81.25	0.00	91.93	0.00
	HistGradientBoost	22.22	0.00	100.00	0.00	11.11	11.11	22.22	11.11	0.00	12.50	22.22	0.00	0.00	25.00	12.50	18.75	23.78	9.81
	RandomForest	11.11	0.00	100.00	0.00	22.22	11.11	11.11	11.11	12.50	31.25	0.00	11.11	0.00	25.00	0.00	37.50	19.62	15.88
	SVM	11.11	0.00	100.00	0.00	22.22	0.00	0.00	11.11	0.00	12.50	0.00	11.11	0.00	25.00	43.75	0.00	22.13	7.46
	XGBoost	33.33	0.00	100.00	0.00	0.00	11.11	22.22	11.11	0.00	18.75	11.11	11.11	0.00	18.75	31.25	12.50	24.74	10.42
	Mean	33.33	0.00	100.00	0.00	42.59	5.55	37.04	7.41	22.92	12.50	27.78	5.55	21.88	15.62	31.25	13.54	39.60	7.52
Early ML without imputer	DecisionTree	77.78	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	87.50	0.00	68.75	0.00	91.75	0.00
	HistGradientBoost	11.11	0.00	33.33	0.00	22.22	0.00	11.11	11.11	12.50	6.25	22.22	0.00	0.00	18.75	0.00	43.75	14.06	9.98
	RandomForest	0.00	0.00	33.33	0.00	11.11	22.22	11.11	11.11	0.00	12.50	11.11	11.11	0.00	25.00	0.00	62.50	8.33	18.05
	XGBoost	11.11	0.00	22.22	11.11	22.22	22.22	33.33	22.22	25.00	6.25	22.22	0.00	0.00	12.50	0.00	18.75	17.01	11.63
	Mean	25.00	0.00	47.22	2.78	38.89	11.11	38.89	11.11	34.38	6.25	38.89	2.78	21.88	14.06	17.19	31.25	32.79	9.92
Early DL with imputer	FTTransformer	11.11	0.00	100.00	0.00	33.33	0.00	22.22	11.11	12.50	25.00	22.22	11.11	0.00	25.00	0.00	12.50	25.17	10.59
	MLP	44.44	0.00	100.00	0.00	66.67	0.00	55.56	0.00	50.00	0.00	44.44	0.00	0.00	25.00	37.50	12.50	49.83	4.69
	TabNet	88.89	0.00	100.00	0.00	33.33	0.00	66.67	0.00	56.25	0.00	55.56	0.00	0.00	18.75	56.25	12.50	57.12	3.91
	TABTransformer	22.22	0.00	100.00	0.00	66.67	11.11	55.56	0.00	56.25	0.00	44.44	0.00	0.00	25.00	43.75	6.25	48.61	5.29
	Mean	41.66	0.00	100.00	0.00	50.00	2.78	50.00	2.78	43.75	6.25	41.66	2.78	0.00	23.44	34.38	10.94	45.18	6.12
Intermediate DL with imputer	FTTransformer	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	0.00	25.00	93.75	0.00	86.72	3.12
	MLP	0.00	0.00	100.00	0.00	22.22	0.00	33.33	11.11	31.25	12.50	44.44	0.00	0.00	25.00	43.75	12.50	34.37	7.64
	TabNet	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	0.00	25.00	87.50	0.00	85.94	3.12
	TABTransformer	11.11	0.00	100.00	0.00	77.78	0.00	66.67	0.00	56.25	0.00	44.44	0.00	0.00	25.00	56.25	6.25	51.56	3.91
	Mean	52.78	0.00	100.00	0.00	75.00	0.00	75.00	2.78	71.88	3.12	72.22	0.00	0.00	25.00	70.31	4.69	64.65	4.45
Late ML with imputer	AdaBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	DecisionTree	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	HistGradientBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	RandomForest	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	SVM	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	XGBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	Mean	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
Late ML without imputer	DecisionTree	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	HistGradientBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	RandomForest	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	XGBoost	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	Mean	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
Late DL with imputer	FTTransformer	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	MLP	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	TabNet	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	TABTransformer	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00
	Mean	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00	100.00	0.00

Table G.2: Percentages of wins and losses of MARIA correct predictions compared with the competitors in the “all missing” scenario using the Wilcoxon signed-rank test.