

M³Care: Learning with Missing Modalities in Multimodal Healthcare Data

Chaohe Zhang*
choc@pku.edu.cn
Key Lab of High Confidence Software
Technologies, Ministry of Education
School of Computer Science,
Peking University
Beijing, China

Xu Chu*
Department of Computer Science and
Technology, Tsinghua University
Key Lab of High Confidence Software
Technologies, Ministry of Education
Beijing, China

Liantao Ma
Yinghao Zhu
Key Lab of High Confidence Software
Technologies, Ministry of Education
National Engineering Research
Center of Software Engineering,
Peking University
Beijing, China

Yasha Wang†
wangyasha@pku.edu.cn
Key Lab of High Confidence Software
Technologies, Ministry of Education
National Engineering Research
Center of Software Engineering,
Peking University
Beijing, China

Jiangtao Wang
Center for Intelligent Healthcare,
Coventry University
Coventry, UK

Junfeng Zhao
Key Lab of High Confidence Software
Technologies, Ministry of Education
School of Computer Science,
Peking University
Beijing, China

ABSTRACT

Multimodal electronic health record (EHR) data are widely used in clinical applications. Conventional methods usually assume that each sample (patient) is associated with the unified observed modalities, and all modalities are available for each sample. However, missing modality caused by various clinical and social reasons is a common issue in real-world clinical scenarios. Existing methods mostly rely on solving a generative model that learns a mapping from the latent space to the original input space, which is an unstable ill-posed inverse problem. To relieve the underdetermined system, we propose a model solving a direct problem, dubbed learning with Missing Modalities in Multimodal healthcare data (M³Care). M³Care is an end-to-end model compensating the missing information of the patients with missing modalities to perform clinical analysis. Instead of generating raw missing data, M³Care imputes the task-related information of the missing modalities in the latent space by the auxiliary information from each patient's similar neighbors, measured by a task-guided modality-adaptive similarity metric, and thence conducts the clinical tasks. The task-guided modality-adaptive similarity metric utilizes the uncensored modalities of the patient and the other patients who also have the

same uncensored modalities to find similar patients. Experiments on real-world datasets show that M³Care outperforms the state-of-the-art baselines. Moreover, the findings discovered by M³Care are consistent with experts and medical knowledge, demonstrating the capability and the potential of providing useful insights and explanations.¹

CCS CONCEPTS

• **Applied computing** → **Health informatics**; • **Information systems** → **Data mining**.

KEYWORDS

healthcare informatics, multimodal data, electronic health record

ACM Reference Format:

Chaohe Zhang, Xu Chu, Liantao Ma, Yinghao Zhu, Yasha Wang, Jiangtao Wang, and Junfeng Zhao. 2022. M³Care: Learning with Missing Modalities in Multimodal Healthcare Data. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22)*, August 14–18, 2022, Washington, DC, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3534678.3539388>

1 INTRODUCTION

Multimodal data can provide complementary information from various modalities that reveal the fundamental characteristics of real-world subjects [6, 8, 21, 52, 53]. Thus, many clinical applications, such as disease diagnosis and mortality prediction [21, 39, 44, 53, 59], require multimodal electronic health record (EHR) data to achieve good diagnostic or prognostic results. Conventional approaches usually assume that each sample is associated with the unified uncensored modalities, and all modalities are available for each sample [40, 58]. However, missing modality is a common issue

*Both authors contributed equally to this research.

†Corresponding Author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '22, August 14–18, 2022, Washington, DC, USA

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9385-0/22/08...\$15.00
<https://doi.org/10.1145/3534678.3539388>

¹published as a conference paper in ACM SIGKDD 2022 (modified a few mistakes)

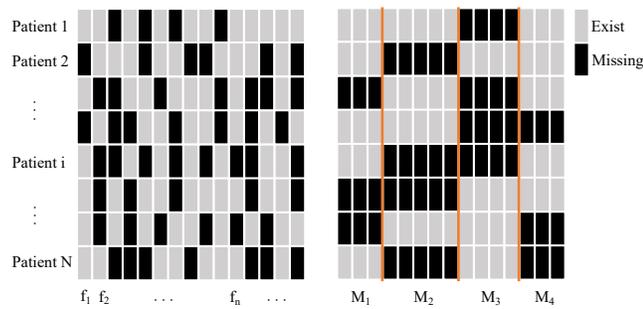


Figure 1: Left: missing features; right: missing modalities. Each row refers to a patient. In the left figure, each column refers to a feature (f_n). In the right figure, each group of columns refers to a modality (M_n), meaning that a modality contains many features. These features have high correlations, since they belong to the same modality (e.g., medical images, medical notes, etc.). The boxes in gray indicate the features exist, and others in black represent missing ones.

in real-world clinical scenarios [23]. For example, different types of examinations are usually conducted for different patients [58]. Also, patients may lack some specific modalities due to patient dropout [45], sensor damage, data corruption [9], safety considerations [47, 60] and high cost [16]. Formally, we define modality-missing in EHR data as, for a sample, data for at least one modality is missing. The absence of a modality means that all features in this modality are missing. Moreover, the modality-missing patterns (i.e., combinations of available modalities) make it more complex for the data with more modalities [58].

Thus, some pioneering research works are proposed to handle the modality-missing issue. Some researchers drop the incomplete samples [44, 51] and achieve some improvements. However, this approach cannot be applied in areas where data is scarce and contains rigid requirements, such as healthcare. Also, it will escalate the small-sample-size issue and over-fitting [9, 45].

The complementary way to dropping methods is the imputation-based method. As shown in the left part of Figure 1, some methods assume that the entries of the data matrix are missing at random (or some more specific assumption on the matrix space, e.g., incoherence, confer [14] for a survey), and the missingness can be imputed via modeling correlations between the columns (features) [35]. Whereas, as illustrated in the right part of Figure 1, missing modalities manifest themselves by column-wise consecutive missingness, where the most correlated information inside the same modalities is missing entirely. On the other hand, the features inside the same modality are naturally more correlated than thereof in different modalities, exhibiting a coherent behavior in the matrix space. Thus the traditional imputation methods do not work well [51]. Additionally, block-wise missingness [54], such as image [22] or geosensory data [54], often assumes a sample realization is from the matrix space, such that the row vectors in a matrix are not permutable. Different from that, the rows in missing modalities refer to the samples, which are permutable. This results in the prior spatial (or spatio-temporal) correlations required to

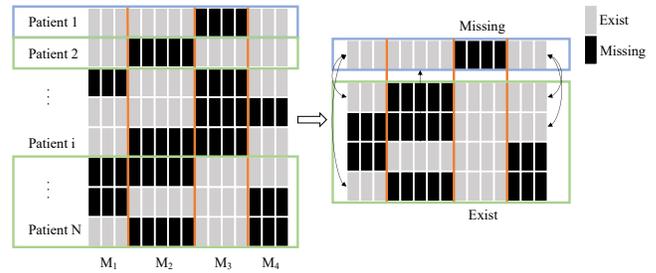


Figure 2: Intuition: For a patient with missing modalities, we utilize the other uncensored modalities of the patient and the other patients who also have the same uncensored modalities to find similar patients and estimate the missing information.

complete the block-wise missingness are not present in the missing modalities. In a word, conceptually, there is a gap between the missing modalities (column-wise consecutive missingness) and the existing imputation-based methods for random or block-wise missingness.

In methodology, some deep generative methods [8, 33, 43, 45] are proposed for missing modalities. Essentially, such methods are usually based on the manifold assumption: the probability mass of real-world objects is supported on low-dimensional manifolds [1]. In terms of EHR data imputation, the manifold assumption can be interpreted as the low-rankness and stability of the feature covariance matrix. In other words, there exist a set of low-dimensional basis vectors and a deterministic mapping \mathcal{T} subject to: (a) the basis vectors span the pre-image of the mapping \mathcal{T} , (b) and the observed EHR feature vectors live on the image of the mapping \mathcal{T} . Existing EHR generative-based completion methods tackle the problem by solving a generative model which learns a mapping from a latent space (spanned by basis vectors) to the original input space [8, 33, 43, 45]. Solving such mapping is essentially an ill-posed inverse problem (low- to high-dimensional underdetermined system [13]), whose solution is often non-unique and unstable [26]. On the other hand, such complex auxiliary models may introduce extra noise, which has negative impacts [9, 15, 51]. To this end, the problem of completing missing modalities requires a different way.

In fact, for missing modalities, solving the generative model is not necessary. By assuming the low-rankness and the stability of the covariance matrix of EHR features, the locally similar row vectors (patients with similar features from some uncensored modalities) in a sub-matrix imply globally similar row vectors (those patients should have similar missing features) in the data matrix. Moreover, if a local row vector X falls in a convex hull spanned by a set of local row vectors, then the global row vector corresponding to X is likely to fall into the convex hull spanned by the particular set of global row vectors. Thus, instead of solving the inverse problem of the low- to high-dimensional mapping, we can solve a less underdetermined problem: comparing similarities of local row vectors and impute the missing entries in a row by referring to the uncensored entries of locally similar rows. The similarity comparison can be conducted in the original input space with a sophisticated metric on data manifolds, or in a learned latent space with a more straightforward

metric. More importantly, modeling the similarity relationship in the low-dimensional latent space is a direct (or forward, namely) problem of solving a mapping from a high-dimensional space to a low-dimensional space, which is less complicated than the inverse problem of solving a generative model.

On the other hand, this intuition is also in agreement with the real-world clinical practice, i.e., how doctors use the relationships between patients to assist the clinical analysis². If two patients are similar in one clinical modality, they are more likely to be similar in another one [18, 19, 25, 32]. Thus, as shown in Figure 2, for a patient with missing modalities, we can utilize the other complete modalities of the patient (local row vector) and other patients who also have these modalities (local uncensored row vectors), to find similar patients. Although seeming straightforward, applying this intuition to clinical tasks will face the following challenges:

Challenge 1. Which space can be used to perform imputation? There are at least two options: the original input space or a learned latent space. Existing methods usually estimate missing data in the original input space [8, 33, 43]. However, the probability distribution in the original input space contains task-relevant and task-irrelevant information. Imputation in the original input space treats task-irrelevant information equally, weakening the task-specific information conveyed in EHR data. This results in an indiscriminate loss of task-relevant and task-irrelevant information, leading to inferior performance.

Challenge 2. What metric(s) should be used to model the similarity relation? T-LSTM [7] uses autoencoders to generate patient representations, based on which the similarity is obtained. SMIL [40] uses multivariate Gaussian to assign similarity weights. However, they did not associate the connotation of similarity with clinical tasks. In different clinical tasks (e.g., mortality prediction and disease diagnosis), the patient characteristics that need attention are different, so two patients considered similar in one clinical task may be considered not so similar in another [57]. More importantly, the similarity metric might vary for different modalities. The manifolds in different modalities in the original input spaces are naturally equipped with different metrics. A learned deterministic mapping from the original input spaces to the learned latent space is not guaranteed to result in a unified metric for different mapped modalities in the latent space.

Challenge 3. How to infer the local-to-global similarities of patients? Metrics computing the similarities between local row vectors is not sufficient to describe the similarities of global row vectors. Thence modeling the intra- and inter-correlations of features from various modalities is challenging but indispensable to aggregate the information from different modalities.

By jointly considering the above issues, we propose a model learning with Missing Modalities in Multimodal healthcare data (M³Care), an end-to-end inductive learning model to compensate for the missing modalities and perform clinical tasks. In summary, our main contributions are summarized as follows:

- We propose M³Care to compensate for the modality-missing patient in the latent space and perform clinical tasks with EHR data.

Since the latent representations are highly compressed and task-supervised, this results in less loss of task-relevant information and is thus more beneficial for subsequent tasks in an end-to-end learning schema (Response to Challenge 1).

- Methodologically, (a) M³Care uses task-guided deep kernels in the latent space of each modality as the metric to compute patient similarities (Response to Challenge 2). (b) M³Care captures intra-correlations within each modality and inter-correlations between modalities by a self-attentive multi-modal interaction module so that the local metrics are aggregated to calculate the similarities of global row vectors. (Response to Challenge 3).
- Extensive experiments show that M³Care outperforms all state-of-the-art models under multiple levels of incompleteness in different evaluation metrics. Besides, the findings discovered by M³Care are in accord with experts and medical knowledge, which shows it can provide useful insights and explanations.

2 RELATED WORK

Multimodal learning for healthcare. With the advancement of medical technology, comprehensive healthcare is burgeoning to meet the demands of patients. This has allowed for multiple medical modalities (e.g., medical image, clinical notes, etc.) to be analyzed to offer patients with feedback, as well as physicians with insights on clinical applications [17, 21, 37, 44, 53].

To this end, multimodal learning for healthcare has attracted the interest of researchers. For example, RAIM [53] is proposed for jointly analyzing continuous monitoring data (e.g., ECG, heart rate) and discrete clinical events (e.g., intervention, lab test) to predict patient decompensation. Gao et al. [17] utilize a multimodal inference model to jointly encode trial criteria text and patient EHR tabular data for patient-trial inference. Huang et al. [24] develop and compare different multimodal fusion architectures to classify Pulmonary Embolism (PE) cases. Ma et al. [38] and Hoang et al. [21] develop distillation frameworks to leverage the multimodal EHR data to enhance the prognosis. Although the methods above work well, one common drawback is that they can only handle samples with complete modalities. Limitations exist while modeling multimodal interactions with the presence of missing modalities.

Methods for missing modalities. Currently, there have been research interests in handling missing modalities, which are mainly divided into two types: deleting incomplete samples or imputing missing modalities. For the first type, FitRec [44] performs workout profile forecasting based on multimodal user data, which discards samples with missing modalities. Wang et al. [51] propose a knowledge distillation framework on samples with complete modalities, while distilling the supplementary information from the incomplete ones. However, such methods can not handle the samples in need with missing modalities and have limitations to be applied in rigid demand domains like healthcare. Besides, such methods dramatically reduce training data and result in over-fitting of deep learning models [9, 45], especially when there are many modalities and many different missing combination patterns³.

The second type is generating the missing modalities at first [8, 40, 45]. However, the incompleteness of modalities leads to column-wise consecutive missingness of features, which makes traditional

²We also substantiate this intuition by mining the real-world clinical datasets, please refer to the Intuition discovery experiments in Appendix.

³e.g., five modalities can result in $2^5 - 1 = 31$ missing patterns.

methods like matrix completion can not be used [51]. Some advanced generative methods such as autoencoders [33, 43] and generative adversarial networks (GAN) [8, 45] have been proposed. These solutions, however, may introduce unwanted extra noise [9]. Especially when the size of samples with complete modalities is small, yet the number of modalities is large, the modalities imputed by such methods may have a negative effect [15, 51]. Moreover, while facing data with many modalities or missing patterns, the number of generators required is also large, which is difficult to train. Chen and Zhang [9] propose a method to enable multimodal fusion of incomplete data and get good performance. However, the method is transductive and needs pre-training, indicating difficulty in applying it when a new sample comes. Therefore, in this paper, we propose a new inductive framework to tackle the above limitations in an end-to-end schema.

3 PROBLEM FORMULATION

In this section, we define the input data and the modeling problem in this paper. Besides, the necessary notations used in the paper are listed in Table 1 for ease of understanding.

Definition 1 (Patient multimodal EHR data). In multimodal EHR data, each patient can be represented as a collection of observations from multiple modalities (data sources), e.g., medical images, clinical notes, lab tests, etc. Suppose M is the number of modalities, N is the number of patients (samples), and let n be the subscript referring to a specific patient, the patient multimodal EHR dataset can be denoted as: $\mathbb{X} = \{X_n\}_{n=1}^N = \{(x_n^1, x_n^2, \dots, x_n^M)\}_{n=1}^N$.

Definition 2 (Patient data with missing modalities). For a specific patient, as mentioned in Section 1, various clinical and social reasons cause the absence of some modalities. Thus, the observed data of a patient are represented as: $X_n = \{x_n^1, x_n^2, \dots, x_n^{M'}\}$, where $0 < M' < M$. It should be noted that, we used the most relaxed setting, i.e., the modality missingness is irregular across patients. In all the training and test (validation) sets, each modality is potentially missing, but at least one modality is present for each patient.

Problem 1 (Disease diagnosis). Given a patient's multimodal EHR data X_n with some missing modalities, we formulate the disease diagnosis task as a *binary* or *multi-label* classification problem, which is diagnosing the disease $y \in \{0, 1\}^{|C|}$ of the patient, where $|C|$ is the number of the unique number of disease categories.

Table 1 shows the notations used in the paper.

4 METHODOLOGY

4.1 Overview

Figure 3 and Figure 4 show the architecture of M³Care. It consists of two main sub-models. The first one is used to compensate for the missing information in the latent space (Corresponding to Figure 3). The other utilizes the processed representations to perform the clinical tasks (Corresponding to Figure 4). Specifically, M³Care includes the following detailed components:

- The *Unimodal Representation Extraction* module maps the original input features of a patient in each modality to the latent space by encoders with various backbones. The backbones are different due to different input modalities (Left part in Figure 3).

Table 1: Notations for M³Care

Notation	Definition
$y \in \{0, 1\}^{ C }$	Ground truth of the classification target
$\hat{y} \in \{0, 1\}^{ C }$	Classification result
\mathbb{X}	The multimodal EHR dataset
$X_n \in \mathbb{X}$	The n -th patient in the dataset
x_n^m	Modality m 's raw data of the n -th patient
$h_n^m \in \mathbb{R}^{N_h}$	Learned representation of modality m of patient n
$H^m \in \mathbb{R}^{B \times N_h}$	Learned representations of modality m for the patient batch
$\hat{H}^m \in \mathbb{R}^{B \times N_h}$	Modality m 's auxiliary information representation aggregated from similar patients
$\bar{H}^{\text{seq}} \in \mathbb{R}^{N_{\text{seq}} \times N_h}$	Representations of a sequential modality with positional encoding added
$\delta \in (0, 1)$	Parameter to control learnable kernel
$\Pi^m \in \mathbb{R}^{B \times B}$	The patient similarity matrix for modality m
$\text{mask}^m \in \mathbb{R}^{B \times B}$	Matrix of booleans that determines each element of the associated value is valid or not (i.e., similarity for missing modality is invalid)
$\Lambda \in \mathbb{R}$	A learnable threshold to filter out dissimilar pairs
$\bar{\Pi} \in \mathbb{R}^{B \times B}$	The comprehensive similarity matrix across all modalities
z^l	The output representations of the l -th layer in the model
\hat{z}^l	The middle representations inside the l -th layer in the model
$\alpha_m \in \mathbb{R}^{B \times 1}$	The importance of self-information H^m
$\beta_m \in \mathbb{R}^{B \times 1}$	The importance of similar patients information \hat{H}^m

- The *Similar Patients Discovery and Information Aggregation* module computes patient similarities of each modality with learned task-guided deep kernels. The similarities induce patient graphs. With graph information propagation, the information from similar patients is aggregated (Middle and right parts in Figure 3)
- The *Adaptive Modality Imputation* module imputes the missing modality in the latent space with the aggregated information, and fuses the existing modality and the auxiliary information to enhance the representation learning (Right part in Figure 3).
- The *Multimodal Interaction Capture* module takes the intra- and inter-modality dynamics into consideration to perform the final clinical tasks (Figure 4).

4.2 Unimodal Representation Extraction

For a specific patient n , it is hard to model the interactions among the raw data, since his/her data X_n is high-dimensional and inconsistent with respect to different data structures in different modalities [9]. Therefore the unimodal representation extraction models are in need to extract the task-relevant feature latent representation in the latent space of each modality. Here, suppose $f_m(\cdot; \Theta_m)$ be the modality m 's unimodal representation extraction model with learnable parameter Θ_m . For modality m 's raw input, the corresponding latent representation can be obtained via:

$$h_n^m = f_m(x_n^m; \Theta_m), \quad (1)$$

where $h_n^m \in \mathbb{R}^{N_h}$ and N_h is the dimension of the representation for modality m . We use the lowercase letter x_n^m to denote a modality for a single patient n .

Three types of $f_m(\cdot; \Theta_m)$ are taken into consideration in this paper: 1) ResNet [20] for embedding image modalities; 2) Transformer Encoder [50] for embedding sequential modalities, such as time-series data like patient laboratory test, medication, and free texts like clinical notes; and 3) multi-layer perceptron (MLP) for embedding vector-based modalities, such as demographic information. As shown in the left part in Figure 3, the black boxes denote the missing information.

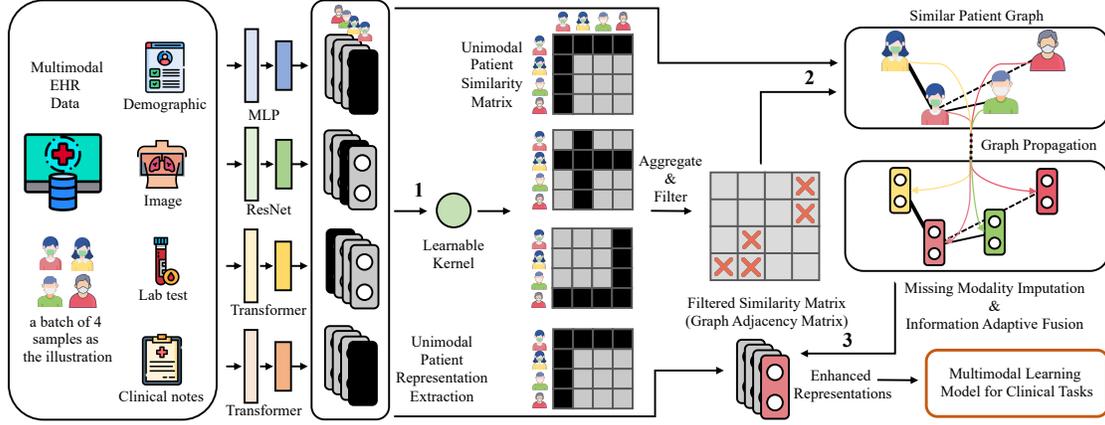


Figure 3: Framework of M³Care, the black boxes denote missing.

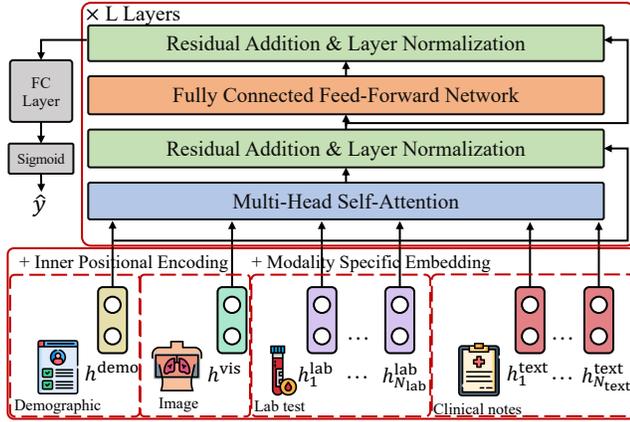


Figure 4: Multimodal Learning Model for Clinical Tasks. Continued from the bottom right corner in Figure 3.

4.3 Similar Patients Discovery and Information Aggregation

The above section describes how to deal with a single sample, and now we focus on a group of samples. Given a batch of patients, the representations of modality m for them are denoted as:

$$H^m = [h_1^m, h_2^m, \dots, h_B^m]^T \in \mathbb{R}^{B \times N_h},$$

where B is the batch size. For sequential modalities, we select the representation of the last timestamp as the one.

For a specific patient with missing modalities, as shown in Figure 1, the representations from the missing modalities cannot be obtained, which results in a lack of modality information. As discussed in Section 1, we can compare similarities of local patient data and impute the missing modalities of the patient by referring to the uncensored modalities of locally similar patients. There is thus a problem with similar patient discovery: What metric(s) should be used to model the similarity relation? In practical applications, one often adopted strategy is to try different kinds of similarity measures on the learned representations in each modality space, such

as Cosine, Euclidean distance, and so on, and then select the best similarity measure [28]. However, this approach is time-consuming, and even if one tries different similarity metrics, it is found that those traditional similarity metrics often fail to consider the local environment of data points, and may learn incomplete and inaccurate relationships. In this case, it is unlikely that the traditional similarity function will be adequate to capture the local manifold structure precisely [27]. Moreover, complex relationships such as higher-order statistics are failed to capture in that way [27].

To this end, we extend this idea to kernel spaces and select the RBF kernel. Given two samples h_i^m and h_j^m , the similarity calculated from the RBF kernel is defined as:

$$k(h_i^m, h_j^m) = \exp\left(-\frac{\|h_i^m - h_j^m\|_2^2}{2\sigma^2}\right), \quad (2)$$

where σ is the bandwidth to control the extent to which similarity of small distances is emphasized over large distances. Following [12], we set σ as a fraction of the mean distance between examples. Expanding the exponential via Taylor series, we can see that the RBF kernel implies an infinite dimension mapping, capturing the higher-order statistics.

Furthermore, as mentioned in Challenge 2, since the data are in multiple modalities, M³Care is required to calculate similarities in each modality space and associate the connotation of similarity with clinical tasks. Thus, a task-guided modality-semantic-adaptive similarity metric is needed. We extend the standard RBF kernel to deep kernel [36] to build an adaptive kernel with a learnable network to fit the representations in a modality data-driven way. Specifically, the kernel is denoted as:

$$k_{\omega_m}(h_i^m, h_j^m) = [(1 - \delta_m)k(\phi_{\omega_m}(h_i^m), \phi_{\omega_m}(h_j^m)) + \delta_m]q(h_i^m, h_j^m), \quad (3)$$

where ϕ_{ω_m} is a network with parameters ω_m for modality m . k and q are different RBF kernels with different σ . The $\delta_m \in (0, 1)$ is a learnable *safeguard* to preventing the learned kernel from going extremely far-away from the right direction.

Now, back to the batch of patient representations, the pairwise similarities with respect to each modality are calculated as (i.e., the

No.1 arrow in Figure 3):

$$\Pi^m = k_{\omega_m}(H^m, H^m), \quad (4)$$

where $\Pi^m \in \mathbb{R}^{B \times B}$ is the patient similarity matrix for modality m . Meanwhile, to ensure the stability of the similarity measure and prevent collapse, we restrict the norm of the difference between the deep representation and the original representation as the optimization objective:

$$\mathcal{L}_{\text{stab}} = \sum_{m=1}^M \left| \|\phi_{\omega_m}(H^m)\|_F - \|H^m\|_F \right|, \quad (5)$$

where the outer $|\cdot|$ means the absolute value, and the inner $\|\cdot\|_F$ is the Frobenius norm.

Because of the characteristic of the kernel method, Π^m is a totally positive symmetric matrix, and each cell $\Pi_{i,j}^m$ in Π^m ranges from 0 to 1, denoting the similarity between the i -th patient and j -th patient for modality m . However, in this way, all the patients are considered similar since the positive similarity. There should be dissimilar ones to be filtered out. A straightforward way is setting a threshold, and the similarities below the threshold are considered dissimilar. Nevertheless, in the early training phase, we notice that when the model is not convergent and the representations are not fully learned, all similarities are unstable. The threshold may filter out some similar patients and lead to inferior performance. Moreover, the determination of the value of the threshold is not trivial. Thus, we utilize a more flexible learnable threshold here. By comprehensively considering similarity from each modality, the filtered similarity matrix can be obtained as:

$$\tilde{\Pi} = \frac{\sum_1^M \Pi^m \cdot \text{mask}^m}{\sum_1^M \text{mask}^m + \epsilon} \quad (6)$$

$$\tilde{\Pi}_{i,j} = \begin{cases} \tilde{\Pi}_{i,j} & \text{if } \tilde{\Pi}_{i,j} > \Lambda \\ 0 & \text{if } \tilde{\Pi}_{i,j} \leq \Lambda \end{cases} \quad (7)$$

where Λ is the learnable threshold to filter out dissimilar pairs and ϵ is used to prevent unstable division by zero. $\text{mask}^m \in \mathbb{R}^{B \times B}$ is the mask matrix of booleans that determines whether each element of the associated value is valid or not (i.e., similarity for missing modality is invalid). For example, in modality m , if both the data of the i -th sample and j -th sample exist, $\text{mask}_{i,j}^m = 1$, and otherwise $\text{mask}_{i,j}^m = 0$, which masks the invalid similarity cell.

Our aim is to impute the modality-missing sample by incorporating auxiliary information from the similar patients. Thus, to aggregate the information from the similar ones, we formulate the batch of patients' representations as a graph in each modality, with the similarity matrix $\tilde{\Pi}$ as the graph adjacency matrix (i.e., the No.2 arrow in Figure 3). Then the graph convolutional layers (GCN) [31] are applied to enhance the representation learning by leveraging the structural information:

$$\hat{H}^m = [\hat{h}_1^m, \hat{h}_2^m, \dots, \hat{h}_B^m]^\top = \text{GCN}(H^m, \tilde{\Pi}) \\ = \text{ReLU}(\tilde{\Pi} \text{ReLU}(\tilde{\Pi} H^m W^0) W^1), \quad (8)$$

where \hat{H}^m is the aggregated auxiliary information from similar patients in the space of modality m . W^0 and $W^1 \in \mathbb{R}^{N_h \times N_h}$ are the projection matrices. We ignore the bias term here and after.

4.4 Adaptive Modality Imputation

Now we obtain two different representations for the batch of patients in each modality: H^m and \hat{H}^m . The former focuses on the patients themselves for modality m , while the latter refers to the information aggregated from similar patients. For a specific patient i , if the data of modality m are missing, we can directly impute the representation with the auxiliary information aggregated from similar patients. While for a patient whose modality m is complete, such auxiliary information can also be fused into the original representation, making the representation smoother to reduce noise, thus enhancing the representation learning. Here, we use an attention fusion to adaptively extract the proper amount of information from them (i.e., H^m and \hat{H}^m). Specifically, two weights $\alpha^m, \beta^m \in \mathbb{R}^{B \times 1}$ are introduced to determine the importance of the above two representations, which are obtained by fully connected layers:

$$\alpha^m = \text{Sigmoid}(H^m W_o), \beta^m = \text{Sigmoid}(\hat{H}^m W_s), \quad (9)$$

where $W_o, W_s \in \mathbb{R}^{N_h \times 1}$ are the weight matrices. α and β indicate the importances of self-information and information of similar patients. We add a constraint $\alpha + \beta = 1$ by calculating $\alpha = \frac{\alpha}{\alpha + \beta}$, $\beta = 1 - \alpha$. The final imputed and enhanced representations can be obtained as (i.e., the No.3 arrow in Figure 3):

$$h_i^m = \begin{cases} \hat{h}_i^m & \text{if modality } m \text{ of sample } i \text{ is missing} \\ \alpha_i^m \cdot h_i^m + \beta_i^m \cdot \hat{h}_i^m & \text{otherwise} \end{cases} \quad (10)$$

where h_i^m and \hat{h}_i^m are the i -th sample of H^m and \hat{H}^m , respectively.

4.5 Multimodal Interaction Capture

Back to a specific patient, so far, the representations of the missing modalities have been imputed through the above sections. These representations are used to perform the clinical tasks. Thus, we need to consider complex correlations among multimodal EHR, including intra-correlations within each modality and inter-correlations between modalities. Inspired by Akbari et al. [2], Kim et al. [30], a context-aware multimodal interaction capture is built. Specifically, for sequential modalities, the internal positional encoding is added: $\bar{H}^{\text{seq}} = [h_1^{\text{seq}}, h_2^{\text{seq}}, \dots, h_{N_{\text{seq}}}^{\text{seq}}] + \text{PE}^{\text{seq}}$, where PE^{seq} is the positional encoding for sequential modality seq and N_{seq} is the length of the sequence. Next, the representations are added with the corresponding modality type embeddings and concatenated to form the input:

$$z^0 = [\bar{H}^1 + \text{TE}^1; \bar{H}^2 + \text{TE}^2; \dots; \bar{H}^M + \text{TE}^M], \quad (11)$$

where TE^m is the corresponding type embedding to identify each modality. And the multimodal interactions are captured through:

$$\bar{z}^l = \text{LayerNorm}(z^{l-1} + \text{MHSA}(z^{l-1})), \\ z^l = \text{LayerNorm}(\bar{z}^l + \text{FFN}(\bar{z}^l)), \quad (12)$$

where $l = 1, \dots, L$ refers to the number of such stacked layers. MHSA refers to the Multi-Head Self-Attention [50], FFN refers to a feed-forward network and LayerNorm is the layer normalization [4]. The predictor is built via:

$$\hat{y}_i = \text{Sigmoid}(z_0^L W_{\text{final}}), \quad (13)$$

where $W_{final} \in \mathbb{R}^{N_h \times 1}$ is the weight matrix. The Sigmoid function is used to turn the output into the probability. In this case, the cross-entropy loss is used as the prediction loss function:

$$\mathcal{L}_{pre} = -\frac{1}{B} \sum_{i=1}^B (y_i^\top \log(\hat{y}_i) + (1 - y_i)^\top \log(1 - \hat{y}_i)), \quad (14)$$

where B is the batch size. $\hat{y}_i \in [0, 1]^{|C|}$ is the predicted probability, and $y_i \in \{0, 1\}^{|C|}$ is the ground truth. The overall loss function is:

$$\mathcal{L} = \mathcal{L}_{pre} + \lambda \mathcal{L}_{stab}, \quad (15)$$

where λ is the hyperparameter to control the loss. For ease of understanding, we summarize M³Care in Algorithm 1 in Appendix.

5 EXPERIMENT

We evaluate M³Care on the following datasets: Ocular Disease Intelligent Recognition (ODIR) Dataset and Ophthalmic Vitrectomy (OV) Dataset. The model code is provided in ⁴.

5.1 Data Description and Task Formulation

We use the following datasets and tasks to evaluate our model.

- **Ocular Disease Intelligent Recognition (ODIR) Dataset** comes from an ophthalmic database, which is meant to represent real-life set of patients collected from hospitals [34]. 3,500 patients are extracted to construct this dataset to diagnose ocular diseases. This dataset contains the following modalities (incomplete modalities exist): demographic information, clinical text for both eyes, and fundus images for both eyes. The detailed statistics are presented in Table 7 in Appendix.

Task. The ocular diseases diagnosis task on this dataset is defined as a multi-label classification task. Following existing works [5, 55], we assess the performance using micro-averaged of the area under the receiver operating characteristic curve (i.e., micro-AUC), macro-AUC, and the average test loss value. We divide the dataset into the training set, validation set, and test set with a proportion of 0.8 : 0.1 : 0.1, and report the performance with the standard deviation of bootstrapping for 1,000 times.

- **Ophthalmic Vitrectomy (OV) Dataset** comes from an ophthalmic hospital⁵. In clinical practice, after vitrectomy, the intraocular pressure (IOP) may increase abnormally. This symptom cannot be predicted by doctors. We collect 832 patients to predict whether the IOP will increase abnormally. This dataset contains six modalities (incomplete modalities exist): demographic, clinical notes, medications, admission records, discharge records, and surgical consumables. The detailed statistics are presented in Table 7.

Task. The task on the dataset are binary classification tasks. Following existing works [21, 57], we assess performance using the area under the precision-recall curve (AUPRC), the area under the ROC curve (AUROC), and accuracy (ACC). AUPRC is the most informative and primary evaluation metric, especially while dealing with skewed real-world data [11, 14, 57]. Due to the size of the dataset, we employ 10-fold cross-validation to assure the consistency of the performance and report the average performance with standard deviations.

⁴<https://github.com/choczhang/M3Care>

⁵This study was approved by the Research Ethical Committee.

5.2 Experimental Setup and Baselines

To conduct the experiment, we use the Adam optimization algorithm in Pytorch 1.5.1. More details are in the Appendix. We include these state-of-the-art models as our baseline models:

- **MFN** [56] captures view-specific and cross-view interactions, and summarizes them with a multi-view gated memory module.
- **MuT** [49] utilizes directional pairwise cross-modal transformers to attend to interactions between multimodal data.
- **ViLT** [30] commissions the transformer module to extract and process all the multimodal features simultaneously.
- **CM-AEs** [43]: The cross-modal autoencoders, which generate missing modalities first and make predictions.
- **SMIL** [40] approximates the missing modality using a weighted sum of manually defined modality priors learned from the dataset.
- **HGMF** [9] fuses incomplete multimodal data within a heterogeneous graph structure, and we modify it to an inductive version. The following ablation studies are also conducted:
- **M³Care₁₋** does not use the task-guided deep kernels of each modality. It directly calculates similarity via cosine similarity.
- **M³Care₂₋** does not consist of the Information Aggregation and the Adaptive Modality Imputation module. It directly computes the mean similarity from each modality and approximates the missing-modality representations via the similar patients.

It should be noted that some of the above models' embedding networks of raw data are a little bit weak. Thus, to perform a fair comparison, we upgrade their embedding layers to the same ones as ours (e.g., Transformer Encoder [50] and ResNet18 [20]) and we do not include any pre-trained parameters.

5.3 Experimental Results

As shown in Table 2 and 3, we can see that M³Care can outperform all the baselines in terms of different evaluation metrics⁶.

Table 2: Results on the ODIR Dataset

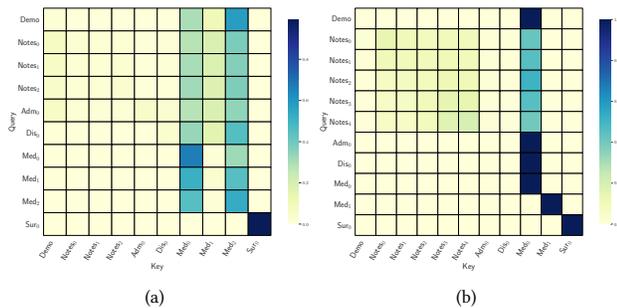
ODIR Dataset (Multi-label Classification)			
Methods	micro-AUC ↑	macro-AUC ↑	test loss ↓
MFN [56]	0.7877 (0.030)	0.7766 (0.029)	0.1772 (0.020)
MuT [49]	0.7944 (0.028)	0.8032 (0.026)	0.2339 (0.019)
ViLT [30]	0.7966 (0.031)	0.7624 (0.029)	0.1731 (0.016)
CM-AEs [43]	0.8028 (0.030)	0.7672 (0.027)	0.1878 (0.031)
SMIL [40]	0.8092 (0.032)	0.7978 (0.025)	0.2278 (0.032)
HGMF [9]	0.8080 (0.030)	0.8103 (0.031)	0.1810 (0.022)
M ³ Care ₁₋	0.8130 (0.031)	0.8059 (0.032)	0.1781 (0.020)
M ³ Care ₂₋	0.8030 (0.031)	0.8138 (0.029)	0.1631 (0.018)
M³Care	0.8490** (0.025)	0.8245** (0.026)	0.1543** (0.018)

Specifically, on the ODIR Dataset, the number in () denotes the standard deviation of bootstrapping for 1,000 times. The results show that, compared with the best baseline method, M³Care achieves relative improvements of 4.9% in micro-AUC. On the OV Dataset, the number in () denotes the standard deviation of 10-fold cross-validation. We can see that, compared with the best baseline method, M³Care achieves relative improvements of 6.1% in AUPRC and 6.0% in AUROC. Among these baseline methods, some ones

⁶** : $p < 0.01$, * : $p < 0.05$

Table 3: Results on the OV Dataset

OV Dataset (Binary Classification)			
Methods	AUPRC \uparrow	AUROC \uparrow	ACC \uparrow
MFN [56]	0.6456 (0.038)	0.6789 (0.032)	0.6627 (0.032)
MuLT [49]	0.6814 (0.047)	0.6891 (0.043)	0.6988 (0.031)
VILT [30]	0.6987 (0.051)	0.7245 (0.048)	0.6627 (0.033)
CM-AEs [43]	0.6891 (0.031)	0.6927 (0.040)	0.6747 (0.029)
SML [40]	0.7109 (0.045)	0.7041 (0.033)	0.6867 (0.032)
HGMF [9]	0.7037 (0.050)	0.7544 (0.027)	0.7100 (0.032)
M ³ Care ₁₋	0.6849 (0.054)	0.7472 (0.052)	0.7080 (0.064)
M ³ Care ₂₋	0.7110 (0.044)	0.7562 (0.057)	0.7234 (0.072)
M³Care	0.7549** (0.065)	0.7998** (0.049)	0.7438** (0.058)

**Figure 5: Attention weight visualization for two patients with abnormally increased intraocular pressure in the test set of the OV dataset. Best viewed in color.**

like CM-AEs [43], SML [40], HGMF [9] use various mechanisms to handle the missing modalities, and thus they achieve relative higher performance. However, the performance boost demonstrates the effectiveness of M³Care. Besides, it is worth mentioning that the OV Dataset only contains the multimodal EHR data of 832 patients yet has six modalities, demonstrating that M³Care performs well on a small dataset while the number of modalities is large, which is suitable for the real-world scenario, where data has a large number of modalities or missing patterns.

The superior performance of M³Care than the M³Care₁₋ (i.e., calculating similarity via cosine similarity) verifies the efficacy of the task-guided deep kernels. Moreover, M³Care outperforms M³Care₂₋, which demonstrates the superiority of the Information Aggregation and the Adaptive Modality Imputation module.

5.4 Further Analysis

We conduct several further experiments. Due to the limitation of pages, some of the experiments are in the Appendix.

5.4.1 Clinical implications. To intuitively show the implication of M³Care, we visualize the attention weights of the prediction process. Due to the limitation of pages, we report two cases in the test set here. As shown in Figure 5, the intraocular pressure (IOP) of the two patients increase abnormally and M³Care successfully predicts the outcome. The rows and columns show the Query and Key multimodal records, which are the abbreviations of each modality, i.e.,

demographic information, clinical notes, medications, admission records, discharge records, and surgical consumables, respectively.

We notice that M³Care gives strong focus on Med₀ and Med₂ (i.e., the first and third medications) of patient *a*, and Med₀ of patient *b*. In all these medications, the patients received the two drugs: *Tropicamide Phenylephrine Eye Drops (Mydrin-P)* and *Prednisolone Acetate Ophthalmic Suspension (Pred Forte)*. These drugs are used for ophthalmologic examinations, prior to ocular surgery [42] or treat eye swelling caused by allergy, infection, injury, or other conditions [48]. Our model discovers that these two drugs may have a strong relationship with the abnormally increased intraocular pressure. This is highly consistent with medical literature [3, 29, 41, 46] and clinician experience, which confirm that the two drugs can lead to adverse reactions like elevation of IOP and should be used with caution for specific patients in clinical practice.

6 CONCLUSIONS

In this paper, we propose M³Care, an end-to-end model to compensate for the missing information of the patients with missing modalities and perform clinical prediction as well as analysis. For a patient with missing modalities, M³Care finds similar patients with a task-guided modality-adaptive similarity metric. Instead of generating raw missing data, M³Care imputes the hidden representations of the missing modalities in the latent space by the auxiliary information from these similar ones, and conducts the clinical tasks. Experiments show that M³Care outperforms all baseline models. Besides, the findings are in accord with experts and medical knowledge, which shows it can provide useful insights.

ACKNOWLEDGMENTS

This work is supported by the National Natural Science Foundation of China (No.62172011). L. Ma is supported by the China Postdoctoral Science Foundation (2021TQ0011). J. Wang is supported by EPSRC New Investigator Award under Grant No.EP/V043544/1.

REFERENCES

- [1] Arvind Agarwal, Samuel Gerber, and Hal Daume. 2010. Learning multiple tasks using manifold regularization. *Advances in neural information processing systems* 23 (2010).
- [2] Hassan Akbari, Linagzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. 2021. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text. *arXiv preprint arXiv:2104.11178* (2021).
- [3] Eray Atalay, Nevbahar Tamçelik, Ceyhan Arici, Ahmet Özkök, and Metin Dastan. 2015. The change in intraocular pressure after pupillary dilation in eyes with pseudoexfoliation glaucoma, primary open angle glaucoma, and eyes of normal subjects. *International ophthalmology* 35, 2 (2015), 215–219.
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).
- [5] Tian Bai and Slobodan Vucetic. 2019. Improving medical code prediction from clinical text via incorporating online knowledge sources. In *The World Wide Web Conference*. 72–82.
- [6] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2018. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence* 41, 2 (2018), 423–443.
- [7] Inci M Baytas, Cao Xiao, Xi Zhang, Fei Wang, Anil K Jain, and Jiayu Zhou. 2017. Patient subtyping via time-aware LSTM networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 65–74.
- [8] Lei Cai, Zhengyang Wang, Hongyang Gao, Dinggang Shen, and Shuiwang Ji. 2018. Deep adversarial learning for multi-modality missing data completion. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1158–1166.

- [9] Jiayi Chen and Aidong Zhang. 2020. Hgmf: heterogeneous graph-based fusion for multimodal data with incompleteness. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1295–1305.
- [10] Pei-Qing Chen, Xue-Mei Han, Ya-Nan Zhu, and Jia Xu. 2016. Comparison of the anti-inflammatory effects of fluorometholone 0.1% combined with levofloxacin 0.5% and tobramycin/dexamethasone eye drops after cataract surgery. *International Journal of Ophthalmology* 9, 11 (2016), 1619.
- [11] Edward Choi, Cao Xiao, Walter Stewart, and Jimeng Sun. 2018. Mime: Multilevel medical embedding of electronic health records for predictive healthcare. In *Advances in Neural Information Processing Systems*. 4547–4557.
- [12] Xu Chu, Yang Lin, Yasha Wang, Xiting Wang, Hailong Yu, Xin Gao, and Qi Tong. 2020. Distance Metric Learning with Joint Representation Diversification. (2020).
- [13] Biswa Nath Datta. 2010. *Numerical linear algebra and applications*. Vol. 116. Siam.
- [14] Jesse Davis and Mark Goadrich. 2006. The relationship between Precision-Recall and ROC curves. In *Proceedings of the 23rd international conference on Machine learning*. 233–240.
- [15] Craig K Enders. 2010. *Applied missing data analysis*. Guilford press.
- [16] FM Ford and J Ford. 2000. Non-attendance for Social Security medical examination: patients who cannot afford to get better? *Occupational medicine* 50, 7 (2000), 504–507.
- [17] Junyi Gao, Cao Xiao, Lucas M Glass, and Jimeng Sun. 2020. Compose: Cross-modal pseudo-siamese network for patient trial matching. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 803–812.
- [18] Jen J Gong and John V Guttag. 2018. Learning to Summarize Electronic Health Records Using Cross-Modality Correspondences. In *Machine Learning for Healthcare Conference*. PMLR, 551–570.
- [19] BARI 2D Study Group. 2009. A randomized trial of therapies for type 2 diabetes and coronary artery disease. *New England Journal of Medicine* 360, 24 (2009), 2503–2515.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [21] Trong Nghia Hoang, Shenda Hong, Cao Xiao, Bryan Low, and Jimeng Sun. 2021. AID: Active Distillation Machine to Leverage Pre-Trained Black-Box Models in Private Data Settings. In *Proceedings of the Web Conference 2021*. 3569–3581.
- [22] Xin Hong, Pengfei Xiong, Renhe Ji, and Haoqiang Fan. 2019. Deep fusion network for image completion. In *Proceedings of the 27th ACM international conference on multimedia*. 2033–2042.
- [23] Shih-Cheng Huang, Anuj Pareek, Saeed Seyyedi, Imon Banerjee, and Matthew P Lungren. 2020. Fusion of medical imaging and electronic health records using deep learning: a systematic review and implementation guidelines. *NPJ digital medicine* 3, 1 (2020), 1–9.
- [24] Shih-Cheng Huang, Anuj Pareek, Roham Zamanian, Imon Banerjee, and Matthew P Lungren. 2020. Multimodal fusion with deep neural networks for leveraging CT imaging and electronic health record: a case-study in pulmonary embolism detection. *Scientific reports* 10, 1 (2020), 1–9.
- [25] Yawen Huang. 2018. *Cross-Modality Feature Learning for Three-Dimensional Brain Image Synthesis*. Ph. D. Dissertation. University of Sheffield.
- [26] Sergei Igorevich Kabanikhin. 2008. Definitions and examples of inverse and ill-posed problems. (2008).
- [27] Zhao Kang, Chong Peng, and Qiang Cheng. 2017. Kernel-driven similarity learning. *Neurocomputing* 267 (2017), 210–219.
- [28] Zhao Kang, Chong Peng, Ming Yang, and Qiang Cheng. 2016. Top-n recommendation on graphs. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. 2101–2106.
- [29] Joon Mo Kim, Ki Ho Park, So Young Han, Kwan Soo Kim, Dong Myung Kim, Tae Woo Kim, and Joseph Caprioli. 2012. Changes in intraocular pressure after pharmacologic pupil dilation. *BMC ophthalmology* 12, 1 (2012), 1–5.
- [30] Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. ViLT: Vision-and-Language Transformer Without Convolution or Region Supervision. In *Proceedings of the 38th International Conference on Machine Learning*. 5583–5594.
- [31] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [32] Jacek Kwiecinski, Evangelos Tzolos, Timothy RG Cartledge, Alexander Fletcher, Mhairi K Doris, Rong Bing, Jason M Tarkin, Michael A Seidman, Gaurav S Gulsin, Nicholas L Cruden, et al. 2021. Native aortic valve disease progression and bioprosthetic valve degeneration in patients with transcatheter aortic valve implantation. *Circulation* 144, 17 (2021), 1396.
- [33] Linchao Li, Bowen Du, Yonggang Wang, Lingqiao Qin, and Huachun Tan. 2020. Estimation of missing values in heterogeneous traffic data: Application of multimodal deep learning model. *Knowledge-Based Systems* 194 (2020), 105592.
- [34] Ning Li, Tao Li, Chunyu Hu, Kai Wang, and Hong Kang. 2021. A Benchmark of Ocular Disease Intelligent Recognition: One Shot for Multi-disease Detection.
- [35] Wei-Chao Lin and Chih-Fong Tsai. 2020. Missing value imputation: a review and analysis of the literature (2006–2017). *Artificial Intelligence Review* 53, 2 (2020), 1487–1509.
- [36] Feng Liu, Wenkai Xu, Jie Lu, Guangquan Zhang, Arthur Gretton, and Dougal J Sutherland. 2020. Learning Deep Kernels for Non-Parametric Two-Sample Tests. *arXiv preprint arXiv:2002.09116* (2020).
- [37] Liantao Ma, Junyi Gao, Yasha Wang, Chaohe Zhang, Jiangtao Wang, Wenjie Ruan, Wen Tang, Xin Gao, and Xinyu Ma. 2020. AdaCare: Explainable Clinical Health Status Representation Learning via Scale-Adaptive Feature Extraction and Recalibration. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- [38] Liantao Ma, Xinyu Ma, Junyi Gao, Xianfeng Jiao, Zhihao Yu, Chaohe Zhang, Wenjie Ruan, Yasha Wang, Wen Tang, and Jiangtao Wang. 2021. Distilling Knowledge from Publicly Available Online EMR Data to Emerging Epidemic for Prognosis. In *Proceedings of the Web Conference 2021*. 3558–3568.
- [39] Liantao Ma, Chaohe Zhang, Yasha Wang, Wenjie Ruan, Jiangtao Wang, Wen Tang, Xinyu Ma, Xin Gao, and Junyi Gao. 2020. ConCare: Personalized Clinical Feature Embedding via Capturing the Healthcare Context. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- [40] Mengmeng Ma, Jian Ren, Long Zhao, Sergey Tulyakov, Cathy Wu, and Xi Peng. 2021. Smil: Multimodal learning with severely missing modality. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 2302–2310.
- [41] Cynthia Matossian, John Hovanesian, Jason Bacharach, Dario Paggiarino, and Keyur Patel. 2020. Impact of dexamethasone intraocular suspension 9% on intraocular pressure after routine cataract surgery: post hoc analysis. *Journal of Cataract and Refractive Surgery* 47, 1 (2020), 53.
- [42] ndrugs. 2022. Mydrin P Uses. <https://www.ndrugs.com/?s=mydrin%20p> [Online; January 2022].
- [43] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y Ng. 2011. Multimodal deep learning. In *ICML*.
- [44] Jianmo Ni, Larry Muhlstain, and Julian McAuley. 2019. Modeling heart rate and activity data for personalized fitness recommendation. In *The World Wide Web Conference*. 1343–1353.
- [45] Yongsheng Pan, Mingxia Liu, Yong Xia, and Dinggang Shen. 2021. Disease-image-specific Learning for Diagnosis-oriented Neuroimage Synthesis with Incomplete Multi-Modality Data. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [46] Sharmila Rajendrababu, Srilekha Pallamparthi, Anusha Arunachalam, Mohammed Sithiq Uduman, Senthilkumari Srinivasan, SR Krishnadas, and Vijayalakshmi A Senthilkumar. 2021. Incidence and risk factors for postoperative intraocular pressure response to topical prednisolone eye drops in patients undergoing phacoemulsification. *International Ophthalmology* 41, 12 (2021), 3999–4007.
- [47] Pedro Ramos-Cabrer, JPM Van Duynhoven, A Van der Toorn, and K Nicolay. 2004. MRI of hip prostheses using single-point methods: in vitro studies towards the artifact-free imaging of individuals with metal implants. *Magnetic resonance imaging* 22, 8 (2004), 1097–1103.
- [48] RxList. 2022. Pred Forte (prednisolone acetate ophthalmic suspension). <https://www.rxlist.com/pred-forte-side-effects-drug-center.htm> [Online; January 2022].
- [49] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2019. Multimodal Transformer for Unaligned Multimodal Language Sequences. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 6558–6569.
- [50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [51] Qi Wang, Liang Zhan, Paul Thompson, and Jiayu Zhou. 2020. Multimodal Learning with Incomplete Modalities by Knowledge Distillation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1828–1838.
- [52] Chang Xu, Dacheng Tao, and Chao Xu. 2013. A survey on multi-view learning. *arXiv preprint arXiv:1304.5634* (2013).
- [53] Yanbo Xu, Siddharth Biswal, Shriprasad R Deshpande, Kevin O Maher, and Jimeng Sun. 2018. Raim: Recurrent attentive and intensive model of multimodal patient monitoring data. In *Proceedings of the 24th ACM SIGKDD international conference on Knowledge Discovery & Data Mining*. 2565–2573.
- [54] Xiufen Yi, Yu Zheng, Junbo Zhang, and Tianrui Li. 2016. ST-MVL: filling missing values in geo-sensory time series data. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence*.
- [55] Wenhao Yu, Wei Peng, Yu Shu, Qingkai Zeng, and Meng Jiang. 2020. Experimental evidence extraction system in data science with hybrid table features and ensemble learning. In *Proceedings of The Web Conference 2020*. 951–961.
- [56] Amir Zadeh, Paul Pu Liang, Navonil Mazumder, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018. Memory fusion network for multi-view sequential learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [57] Chaohe Zhang, Xin Gao, Liantao Ma, Yasha Wang, Jiangtao Wang, and Wen Tang. 2021. GRASP: Generic Framework for Health Status Representation Learning Based on Incorporating Knowledge from Similar Patients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 715–723.
- [58] Changqing Zhang, Zongbo Han, Yajie Cui, Huazhu Fu, Joey Tianyi Zhou, and Qinghua Hu. 2019. CPM-nets: cross partial multi-view networks. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 559–569.

- [59] Xi Zhang, Jingyuan Chou, and Fei Wang. 2018. Integrative analysis of patient health records and neuroimages via memory-based graph convolutional network. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 767–776.
- [60] Tao Zhou, Huazhu Fu, Geng Chen, Jianbing Shen, and Ling Shao. 2020. Hi-net: hybrid-fusion network for multi-modal MR image synthesis. *IEEE transactions on medical imaging* 39, 9 (2020), 2772–2781.

A INTUITION DISCOVERY EXPERIMENT

In each dataset, for each modality, we first take the samples containing the complete modality data and divide them into training, validation, and test sets. Then we use a unimodal classifier (e.g., multilayer perceptron, transformer encoder) to classify the data for each modality. We take the best model from the validation set, apply it to the test set, and collect the representations in the latent space of each modality for each sample. Next, we compute pair-wise similarity matrix Π_m between samples in each modality, where m is the corresponding modality.

We want to justify the intuition in our datasets: if two patients are similar in one modality, they are more likely to be similar in another modality with regard to the clinical task. We call intuition the *cross-modal transfer of sample similarity*. Given modalities a and b , the intuition holds if their difference of the pair-wise similarity matrix $\|\Pi_a - \Pi_b\|_{\text{norm}}$ is small enough, where norm is a type of matrix norm. In this experiment, we try different similarity metrics such as normalized Euclidean distance, cosine similarity and RBF kernel. And we also try different norms as metrics of difference, such as Frobenius norm, 2-norm and mean value of all entries in $\|\Pi_a - \Pi_b\|$.

For comparison, we add noise to the representations in the latent space in one modality (we select modality b here), and calculate the difference $\|\Pi_a - \Pi'_b\|_{\text{norm}}$. If the intuition holds, the difference should be bigger than the above original $\|\Pi_a - \Pi_b\|_{\text{norm}}$. We perform this experiment 1,000 times and calculate the average difference to avoid chance.

Furthermore, we also shuffle the representations in the latent space in one modality (we select modality b here), and calculate the difference $\|\Pi_a - \Pi''_b\|_{\text{norm}}$. If the intuition holds, the difference should also be bigger than the first original $\|\Pi_a - \Pi_b\|_{\text{norm}}$. In the same way, we perform this experiment 1,000 times and calculate the average difference to avoid chance. The results are shown in Table 4 and 5.

As shown in Table 4, the original difference of the pair-wise similarity matrices in the two modalities is smaller than both the Noise and Shuffle ones with regard to different similarity metrics and different norms. This justifies that if two patients are similar in one modality, they are more likely to be similar in another modality in different view. To this end, we come up with our intuition, i.e.,

Table 4: Intuition observation results for two modalities: admission records and clinical notes, on the OV Dataset

Metric	Norm	Original	Noise	Shuffle
Normalized	Frobenius	183.35	403.23	207.27
Euclidean	2	153.97	385.21	173.73
Distance	Mean	0.2334	0.4871	0.2724
Cosine similarity	Frobenius	402.96	462.59	452.72
	2	336.31	365.02	376.69
	Mean	0.5081	0.5821	0.5908
RBF kernel	Frobenius	176.35	261.05	199.60
	2	147.51	196.10	166.46
	Mean	0.2241	0.3024	0.2623

the *cross-modal transfer of sample similarity*. In another pair, as shown in Table 5, the same conclusion can be drew.

B ALGORITHM

Algorithm 1 shows the algorithm of M³Care.

Algorithm 1: Algorithm of M³Care

Input:
Multimodal EHR dataset \mathbb{X}

Output:
Prediction for the patient \hat{y}

Training:

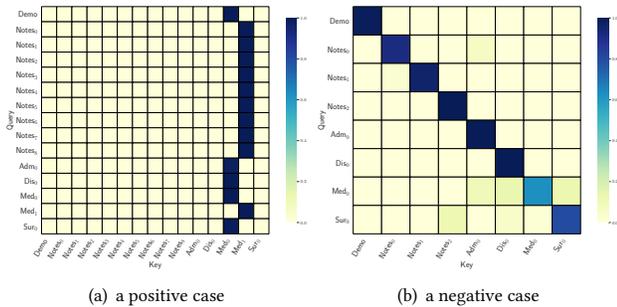
- 1 Initialize weights;
- 2 **while** training is not convergence **do**
- 3 **for** each batch of patient **do**
- 4 Extract h_n^m of each modality m via Eq. 1;
- 5 Form the batch-wise representation matrices H^m ;
- 6 Compute patient similarity matrix Π^m in each modality space via Eq. 4;
- 7 Compute comprehensive similarity $\tilde{\Pi}$ via Eq. 6, 7;
- 8 Form the similar patient graph with H^m as nodes and $\tilde{\Pi}$ as adjacency matrix;
- 9 Compute the aggregated information \hat{H}^m via Eq. 8 in the space of each modality m ;
- 10 **for** each patient in the batch **do**
- 11 **if** missing modalities exist **then**
- 12 | Impute the missing modalities via Eq. 10
- 13 **end**
- 14 **else**
- 15 | Enhance the representations via Eq. 9-10;
- 16 **end**
- 17 Model the multimodal dynamics via Eq. 11-12;
- 18 Make prediction via Eq. 13;
- 19 **end**
- 20 Update the parameters by optimizing Eq. 15;
- 21 **end**
- 22 **end**

Table 5: Intuition observation results for two modalities: medications and surgical consumables information, on the OV Dataset

Metric	Norm	Original	Noise	Shuffle
Normalized	Frobenius	152.96	159.02	370.39
Euclidean	2	131.37	145.77	360.48
Distance	Mean	0.2021	0.2289	0.5588
Cosine similarity	Frobenius	384.57	406.01	404.89
	2	326.61	353.72	367.75
	Mean	0.4984	0.5328	0.5671
RBF kernel	Frobenius	168.85	173.56	203.48
	2	144.43	157.96	188.48
	Mean	0.2081	0.2260	0.2825

Table 6: Micro-AUC on ODIR Dataset under different additional synthesizing multimodal missing rates.

Methods	30%	40%	50%	60%
MFN [56]	.7625 (.03)	.7334 (.03)	.7260 (.03)	.7166 (.03)
MuIT [49]	.7768 (.02)	.7637 (.03)	.7480 (.03)	.7348 (.03)
ViLT [30]	.7601 (.03)	.7583 (.03)	.7492 (.03)	.7355 (.02)
CM-AEs [43]	.7846 (.03)	.7707 (.00)	.7648 (.03)	.7477 (.03)
SMIL [40]	.7702 (.02)	.7595 (.02)	.7485 (.03)	.7396 (.03)
HGMF [9]	.7831 (.03)	.7711 (.03)	.7585 (.03)	.7427 (.02)
M³Care	.8119** (.02)	.7927** (.03)	.7795** (.03)	.7715** (.02)

**Figure 6: Attention weight visualization for two patients in the test set of OV dataset.**

C FURTHER ANALYSIS

C.1 Multiple levels multimodal incompleteness

To consider more realistic various settings and verify the generalizability of M³Care, experiments on multimodal data with multiple levels of multimodal incompleteness are conducted. We evaluate the influences of missing modalities by attaching additional synthesizing multimodal incompleteness rates on ODIR Dataset from 30% to 60% with an intermittent 10%. The experiments are repeated test with bootstrapping for 1,000 times, and the results (micro-AUC) are in Table 6.

We can see that all the models' micro-AUCs decrease as the missing rate increases, and M³Care still outperforms all baselines. When the missing rate is the biggest of all settings (60%), M³Care also demonstrated significantly better performance than the best baselines CM-AEs [43] and HGMF [9]. Specifically, M³Care achieves a micro-AUC of 0.7715, while the baseline models CM-AEs [43] and HGMF [9] achieve 0.7477 and 0.7427, showing 3.2% and 3.8% relative improvement, respectively.

C.2 Clinical implications

To intuitively show the implication of M³Care, similar to Section 5.4.1, we further visualize the attention weights of the prediction process for another two cases. As shown in Figure 6, the first case has a positive label, and the second one is negative. For the first one, M³Care gives a strong focus on Med₀ and Med₁. In both these medications, the patients received not only the drug mentioned above

in Section 5.4.1: *Tropicamide Phenylephrine Eye Drops (Mydrin-P)*, which have been proved by medical literature [3, 29] and clinician experience that they can lead to adverse reactions like elevation of IOP. The patients also received Tobramycin Dexamethasone Eye Drops (Tobradex), which is highly consistent with medical literature [10] and clinician. The drugs can lead to adverse reactions and have the tendency to increase intraocular pressure [10]. For the second one, a healthy patient with a negative label, M³Care gives relatively even attention to each data of each modality of the patient, which indicates that M³Care does not discover significant signs of elevation of IOP and finally makes a right prediction.

D DETAILS OF EXPERIMENTAL SETTINGS

D.1 Statistics of the Datasets

- **Ocular Disease Intelligent Recognition (ODIR) Dataset** contains the following modalities (incomplete modalities exist): demographic information, clinical text for both eyes, and fundus images for both eyes. The detailed statistics are presented in Table 7.
- **Ophthalmic Vitrectomy (OV) Dataset** contains six modalities (incomplete modalities exist): demographic, clinical notes, medications, admission records, discharge records, and surgical consumables. The detailed statistics are presented in Table 7.

Table 7: Statistics of the Datasets

Dataset	Statistic	Value
Ocular Disease Intelligent Recognition (ODIR) Dataset	# patients	3,500
	# modalities	3
	% missing per modality	[0%, 48.34%, 0%]
	% positive labels	[0.061, 0.060, .0046]
Ophthalmic Vitrectomy (OV) Dataset	% female	.461
	# patients	832
	# modalities	6
	% missing per modality ⁷	[0%, 0.60%, 15.38%, 10.33%, 10.33%, 4.08%]
	% positive labels	41.7%
	% female	.456

D.2 Model Implementation

The experiment environment is a machine equipped with CPU: Intel Xeon E5-2630, 256GB RAM, and GPU: Nvidia RTX8000. The code is implemented based on Pytorch 1.5.1. The hyper-parameter setting of the proposed M³Care is as follows: We set the embedding dimension and hidden dimension as 128/256 for Ocular Disease Intelligent Recognition (ODIR) / Ophthalmic Vitrectomy (OV) dataset, respectively. Since the clinical notes modality in OV dataset is too long per patient, and the number of samples is small (832 patients), we set the batch size as 32 while conducting experiments on OV dataset. For ODIR dataset, we set the batch size as 512.

⁷The order of the missing rates are corresponding to the above data description.