

# Generate, Then Retrieve: Addressing Missing Modalities in Multimodal Learning via Generative AI and MoE

Sukwon Yun<sup>1\*</sup>, Jiayi Xin<sup>2\*</sup>, Inyoung Choi<sup>2</sup>, Jie Peng<sup>1</sup>, Ying Ding<sup>3</sup>, Qi Long<sup>2</sup>, Tianlong Chen<sup>1</sup>

<sup>1</sup>University of North Carolina at Chapel Hill <sup>2</sup>University of Pennsylvania <sup>3</sup>University of Texas at Austin  
{swyun, tianlong}@cs.unc.edu, {jiayixin, inyoungc}@seas.upenn.edu, ying.ding@ischool.utexas.edu, qlong@upenn.edu

## Abstract

In multimodal machine learning, effectively addressing the missing modality scenario is crucial for improving performance in downstream tasks such as in medical contexts where data may be incomplete. Although some attempts have been made to effectively retrieve embeddings for missing modalities, two main bottlenecks remain: the (1) consideration of both intra- and inter-modal context, and the (2) cost of embedding selection, where embeddings often lack modality-specific knowledge. In response, we propose MoE-Retriever, a novel framework inspired by the design principles of Sparse Mixture of Experts (SMoE). First, MoE-Retriever define a supporting group for intra-modal inputs, i.e., samples that commonly lack the target modality. This group is formed by selecting samples with complementary modality combinations for the target modality. It is then integrated with inter-modal inputs—i.e., inputs from different modalities of a sample—thereby establishing both intra- and inter-modal contexts. These inputs are processed by Multi-Head Attention, generating context-aware embeddings that serve as inputs to the SMoE Router, which automatically selects the most relevant experts, i.e., the embedding candidates to be retrieved. Comprehensive experiments on both medical and general multimodal datasets demonstrate the robustness and generalizability of MoE-Retriever, marking a significant step forward in embedding retrieval methods for incomplete multimodal data. The source code of MoE-Retriever is available here: <https://github.com/UNITES-Lab/moe-retriever>

## Introduction

In the era of generative AI and multimodal learning, effectively addressing the *missing modality scenario* has become a pivotal challenge for enhancing downstream task performance (Baltrušaitis, Ahuja, and Morency 2018; Guo, Wang, and Wang 2019; Wu, Wang, and Chen 2024). In practical cases such as clinical and biological settings, modalities such as imaging, genetic, and clinical data often contain missing entries due to varying acquisition times, costs, or patient-specific factors (Ma et al. 2021; Zhang et al. 2022a,b; Wang et al. 2023). To address this, prior approaches primarily focus on two strategies: imputing missing features directly within

the input feature space or employing learnable embedding to represent missing features in the latent space. The former often involves some rule-based prior, such as using the population mean to perform imputation. This method does not scale with data, as the imputation method remain fixed when the underlying distribution changes. In contrast, recent research has increasingly turned toward the latter — leveraging learnable embedding to provide more adaptive and context-aware representations for missing modalities (Zhang et al. 2022b,a; Wu et al. 2024; Han et al. 2024). However, despite their promise, these learnable embedding-based methods still face several critical limitations.

**Intra- & Inter-Modal Context.** As illustrated in Figure 1 (a), current methods inadequately address both intra-modal and inter-modal contexts when supplementing missing modalities, often focusing on one or the other. In intra-modal scenarios, the goal is to retrieve embeddings for the missing (target) modality by identifying similar samples (Malitesta et al. 2024). However, existing works often choose unimodal approaches that primarily address intra-modal context, failing to personalize the sample’s heterogeneous context. Conversely, in inter-modal scenarios, it is assumed that modality-invariant and modality-specific information exists across input modalities, suggesting that missing modalities can be imputed from the sample’s specific observed modalities (Zhang et al. 2022b; Wang et al. 2023). However, these works do not carefully consider intra-sample information while proceeding with multi-modal fusion. As a result, focusing solely on either intra-modal or inter-modal context leads to incomplete supplementation and limits the model’s ability to effectively leverage the rich multimodal information available in real-world datasets. This highlights the need for a more holistic approach that integrates both perspectives for more accurate and robust imputation of missing modalities.

**Embedding Selection.** Figure 1 (b) illustrates the current state of embedding retrieval. Current methods either treat the learnable or retrieved embeddings as a single embedding (Wang et al. 2023; Han et al. 2024) or use diverse embeddings but require activating all candidates every time a retrieval is performed, using operations like summation, averaging, or attention mechanisms. These methods can incur a high computational cost as the number of samples or modalities grows, and they lack the ability to adapt to diverse observed modality combinations, treating all potential scenarios equally re-

\*These authors contributed equally.

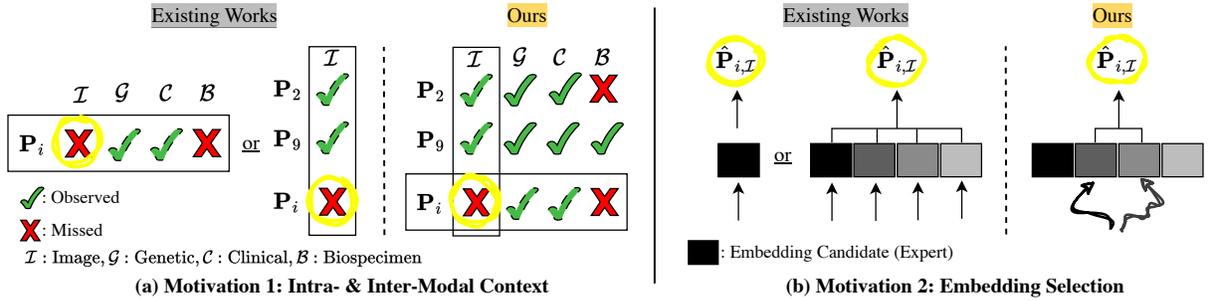


Figure 1: Motivation of this work. (a) **Motivation 1: Intra- & Inter-Modal Context**: Existing works typically consider either the intra-modal context (between samples with the same missing modality, such as  $P_{2,\mathcal{I}}$ ,  $P_{9,\mathcal{I}}$ ) or the inter-modal context (within a sample’s observed modalities, such as  $P_{i,\mathcal{G}}$ ,  $P_{i,\mathcal{C}}$ ). In contrast, our work incorporates both contexts simultaneously, generating context-aware embeddings to enhance missing modality retrieval. (b) **Motivation 2: Embedding Selection**: When retrieving the most relevant embedding ( $\hat{P}_{i,\mathcal{I}}$ ), existing approaches either use a single static embedding or combine multiple embeddings with simple methods (e.g., sum, average, attention), which makes it difficult to obtain specialized knowledge and requires activation of each embedding candidate every time. In contrast, our work leverages the design principles of SMoE, using a router to automatically select the most relevant experts through top-k selection in a sparse and efficient manner.

regardless of the specific context. This uniformity in handling observed modalities limits the capacity for more nuanced and context-specific supplementation. For instance, specific knowledge may be required when certain input modality combinations are present, which is crucial for improving downstream task performance.

**Our Approach.** To address these challenges, we propose MoE-Retriever, a novel framework for embedding retrieval given a incomplete multimodal data. The main idea of MoE-Retriever is to borrow the design principle from the Sparse Mixture of Experts (SMoE), which activates most relevant experts (i.e., embedding candidates) given a specific intra- and inter-modal context within in a router in a sparse manner. To achieve this, we first generate context-aware embeddings, which serve as inputs to SMoE. More specifically, we define a supporting group based on the given modality combination, aiming to reconstruct the target (missing) modality by sampling intra-modal examples. Next, by incorporating inter-modal samples and applying Multi-Head Attention, we generate context-aware embeddings. Finally, the SMoE router retrieves the missing embedding by selecting the most relevant experts, which include both shared and modality-specific experts. Extensive experiments on two medical datasets (ADNI, MIMIC) and two general machine learning datasets (ENRICO, CMU-MOSI) validate the efficacy and generalizability of MoE-Retriever, demonstrating its robust performance across various multimodal settings.

- We highlight that current intra- or inter- modal or single or multiple-but-lacking specialized knowledge brings the bottleneck into incomplete multimodal embedding retrieval.
- We propose MoE-Retriever, borrowing the design principle of Sparse Mixture of Experts design, which inputs the both intra-modal inter-sample and inter-modal intra-sample contexts and retrieve most relevant embedding from modality-specific and shared experts.
- Our comprehensive experimental evaluations on the medicinal dataset and machine learning datasets, showcase the effectiveness and portability of MoE-Retriever.

## Method

**Motivation behind bringing SMoE design.** In the context of incomplete multimodal data, only the observed features in the raw feature space can pass through the modality-specific encoder. This raises a critical question: *how can we effectively handle samples with missing modalities to provide robust embeddings for the missing features?* Ensuring that the embedding space, followed by the fusion and prediction layers, remains trainable through continuous gradient flow is essential. It is important to note that different samples exhibit varying combinations of observed modalities, which necessitates a personalized approach capable of handling each sample’s unique environment, such as its specific modality combination.

To address this challenge, we introduce the design principles of SMoE (Shazeer et al. 2017). Given a pool of diverse experts (i.e., trainable feed-forward networks), the SMoE architecture enables the automatic and sparse activation of different experts, each specializing in certain knowledge, based on the input scenario. This dynamic routing mechanism effectively mitigates the limitations of static, one-size-fits-all designs, where learnable embeddings are constrained to a single expert or a fixed combination of embeddings without a router. In such static setups, embeddings for missing modalities are often selected at random, leading to suboptimal performance for downstream tasks.

**Notation.** Formally, SMoE consists of multiple experts, denoted as  $\mathcal{E}_1, \dots, \mathcal{E}_{|\mathcal{E}|}$ , where  $|\mathcal{E}|$  represents the total number of experts, and a router,  $\mathcal{R}$ , which governs the routing mechanism, sparsely selecting the top- $k$  experts. For a given embedding or token  $\mathbf{x}$ , the router  $\mathcal{R}$  activates the top- $k$  experts based on the highest scores derived from a softmax function applied to the outputs of a learnable gating function,  $g(\cdot)$ , typically modeled as a one or two-layer MLP. The router’s output,  $\mathcal{R}(\mathbf{x})_i$ , indicates the selection of the  $i$ -th expert. This process is formally described as follows:

$$\mathbf{y} = \sum_{i=1}^{|\mathcal{E}|} \mathcal{R}(\mathbf{x})_i \cdot \mathcal{E}_i(\mathbf{x}),$$

$$\mathcal{R}(\mathbf{x}) = \text{Top-K}(\text{softmax}(g(\mathbf{x})), k), \quad (1)$$

$$\text{TopK}(\mathbf{v}, k) = \begin{cases} \mathbf{v}, & \text{if } \mathbf{v} \text{ is in the top } k, \\ 0, & \text{otherwise.} \end{cases}$$

**MoE-Retriever.** The overall framework of MoE-Retriever, along with the detailed procedure, is illustrated in Figure 2. In essence, the key idea behind MoE-Retriever is to retrieve the most relevant embedding for the missing modality by leveraging two contexts: (1) **Intra-Modal Context**, which samples similar examples from a well-defined supporting group based on the observed modality combination (Sec ), and (2) **Inter-Modal Context**, which considers the sample-specific heterogeneous combination of observed modalities (Sec ). The next step is (3) **Context-Aware Routing**, where the expert pool is designed modality-specifically, using both contexts to effectively supplement the target (i.e., missing) modality. Finally, the selected experts and their linear combination with the inputs are integrated into a single embedding, which is regarded as the final retrieved embedding (Sec ).

**Intra-Modal Context** We begin with the *intra-modal context* (column-wise context in Figure 2), where intra-modal refers to the homogeneous modality that matches the target modality we aim to supplement. The rationale for incorporating this context is that, by forming a pool of similar samples, we can capture patterns directly observed across patients, without requiring any additional preprocessing. The observed pattern can be represented as a *modality combination*, which reflects similar trends or patterns, i.e., knowledge observed across the samples. Empirically, samples (e.g., patients) with similar observed modality combinations have shown exhibit analogous characteristics. For instance, patients who lack the image modality but possess both genetic and clinical modalities may be more likely to display correlations with certain domain-specific traits, such as early-stage diagnosis or slower progression rates, often associated with genetic risk factors like the APOE  $\epsilon 4$  allele (Dubois et al. 2023; Jack Jr et al. 2018; Lambert et al. 2013).

To effectively sample from an intra-modal sample pool, we first need to generate a modality combination-specific pool, which we denote as the *supporting group*. The core idea behind the supporting group is that, given an observed modality combination and a target (missing) modality, the corresponding group must include the observed modalities as well as the target modality to support the patient’s intra-modal pool. For example, if a sample contains the modalities ‘GC’ and we aim to impute the modality ‘T’ (as illustrated in Figure 2), the supporting group should include samples with ‘GC’ as well as the missing modality ‘T’. Consequently, the supporting group would comprise samples with modality combinations such as ‘TGC’ or ‘TGC $\mathcal{B}$ ’.

Formally, as an example from Figure 2, let the set of modalities be  $\mathcal{M} = \{\mathcal{T}, \mathcal{G}, \mathcal{C}, \mathcal{B}\}$ . With a specific modality combination  $mc \in \mathcal{MC} =$

$\{\mathcal{T}, (\mathcal{T}, \mathcal{G}), (\mathcal{T}, \mathcal{G}, \mathcal{C}), \dots, \mathcal{G}, (\mathcal{G}, \mathcal{C}), \dots, (\mathcal{T}, \mathcal{G}, \mathcal{C}, \mathcal{B})\}$ , where the total number of combinations in  $\mathcal{MC}$  is  $|\mathcal{MC}| = \sum_{m=1}^{|\mathcal{M}|-1} \binom{|\mathcal{M}|}{m} = 2^{|\mathcal{M}|} - 1$ , the supporting group  $G$  consists of the samples that satisfy the following constraints:

$$G(j | \mathcal{T}, mc) = \{j \in \{1, \dots, N\} \mid mc_j \in \mathcal{X}(\mathcal{S} | \mathcal{T}, mc)\}$$

where  $\mathcal{X}(\mathcal{S} | \mathcal{T}, mc) = \{S \subseteq \mathcal{M} \mid (mc \subseteq S) \wedge (\mathcal{T} \in S)\}$ ,

$$\forall \mathcal{T} \in \mathcal{M}, \forall mc \in \mathcal{MC}. \quad (2)$$

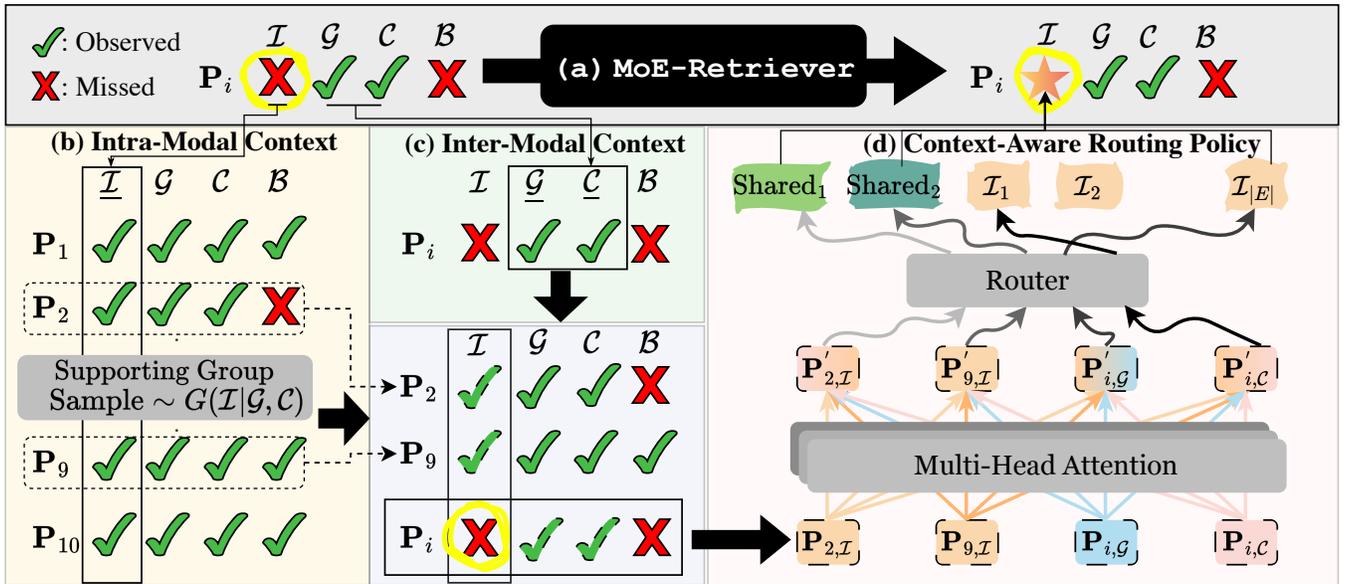
where  $G(j | \mathcal{T}, mc)$  denotes the set of sample indices among total sample size  $N$ , derived from the set of possible modality combinations  $\mathcal{X}(\mathcal{S} | \mathcal{T}, mc)$  for a given target modality  $\mathcal{T}$  and modality combination  $mc$ . In this context, the satisfying  $\mathcal{S}$  denotes any arbitrary set of modality combinations that satisfies the constraint of including both  $mc$  (i.e.,  $(mc \subseteq S)$ ) and (i.e.,  $\wedge$ ) the target modality  $\mathcal{T}$  as subsets (i.e.,  $(\mathcal{T} \in S)$ ). Given the supporting group  $G$ , we sample<sup>1</sup> intra-modal examples that assist in the final retrieval by SMoE by referring to similar examples within the homogeneous modality.

**Inter-Modal Context** Beyond intra-modal context, we now consider another critical dimension: inter-modal context (illustrated row-wise in Figure 2). This approach allows us to incorporate personalized context specific to a given sample that would be missed by only considering intra-modal context. As a real-world example, this perspective is particularly meaningful in multimodal medical scenarios such as Alzheimer’s diagnosis. When genetic (G) and clinical (C) data are available but imaging (I) is missing (case of Figure 2), it may suggest the patient is in the early stages of the disease, where less invasive and more accessible modalities are prioritized. Imaging, typically more expensive, may be reserved for later stages when symptoms progress (Dubois et al. 2023; Li et al. 2022). Additionally, genetic and clinical data alone can provide valuable early insights, guiding initial interventions before resorting to costly imaging techniques (Kim 2023).

Formally, to consider inter-modal context, we directly focus on the observed modalities, i.e.,  $mc$  (e.g.,  $(\mathcal{G}, \mathcal{C})$ ) for a sample index,  $i$ . By doing so, we integrate these sample-specific heterogeneous modality combinations, which will serve as input for the inter-modal examples in the final retrieval by SMoE, referring to the personalized context within the heterogeneous modalities.

**Context-Aware Routing Policy** Now, given two contexts, i.e., intra-modal and inter-modal, we proceed with context-aware routing via the SMoE design. The goal of this routing is to retrieve the most relevant expert given an input combination that includes both homogeneous and heterogeneous modality information. For each embedding (i.e., token) input to the router, the router is trained to select the most relevant

<sup>1</sup>For the number of samples, we used a count that matches the observed modalities of the samples (i.e.,  $|mc|$ ) to ensure a balanced impact of both. They may vary and can be treated as a hyperparameter for flexibility. However, empirical observations indicate that varying the number of intra-modal samples has only a marginal effect on model performance.



$\checkmark$ : Observed  
 $\times$ : Missed  
 $\mathcal{I}$ : Image,  $\mathcal{G}$ : Genetic,  $\mathcal{C}$ : Clinical,  $\mathcal{B}$ : Biospecimen |  $\text{orange}$ : Intra-Modal Inputs  $\text{blue}$   $\text{pink}$ : Inter-Modal Inputs  $\text{grey}$ : Experts

Figure 2: Overall illustration of MoE-Retriever. (a) **The role of MoE-Retriever.** Given a sample ( $P_i$ ) with a missing modality,  $\mathcal{I}$  (Image), our goal is to retrieve the most relevant embedding ( $P_{i,\mathcal{I}}$ ) by considering two contextual factors. First, we focus on (b) **Intra-Modal Context**, which seeks to find embeddings within the same modality as the missing one ( $\mathcal{I}$ ) to reflect similar contextual knowledge. To achieve this, we define a supporting group ( $G(\mathcal{I}|\mathcal{G}, \mathcal{C})$ ), where the target modality ( $\mathcal{I}$ ) and the sample’s observed modalities ( $\mathcal{G}, \mathcal{C}$ ) form a sufficient context for grouping. After sampling from this group, we incorporate the sample’s specific (c) **Inter-Modal Context**, leveraging the observed modalities. We then proceed to (d) **Context-Aware Routing Policy**, which first applies multi-head attention and adopts the SMoE framework. Here, the router (top-1 selection in this example) selects the most relevant experts given two intra- and inter-modal inputs. After integrating all the embeddings, the final embedding is regarded as the retrieved embedding for the sample  $i$ ’s missing modality  $\mathcal{I}$ , denoted as  $P_{i,\mathcal{I}}$ . For retrieving an embedding for another missing modality,  $\mathcal{B}$ , the supporting group would be updated to  $G(\mathcal{B}|\mathcal{G}, \mathcal{C})$ , and the intra-modal embeddings would consist of  $P_{i,\mathcal{B}}$ , with the expert selection adapted accordingly to  $\{\mathcal{B}_1, \dots, \mathcal{B}_E\}$ .

expert that can benefit the downstream task. The selected experts are expected to specialize in handling the specific input modalities.

The context-aware router design is detailed as follows:

$$\hat{P}_{i,\mathcal{T}} = \sum_{e=1}^{|\mathcal{E}|} \mathcal{R}(\mathbf{x})_e \cdot \mathcal{E}_e^T(\mathbf{x}) \text{ where } \mathbf{x} \in \{P'_{i_{\text{intra}},\mathcal{T}} \cup P'_{i,mc}\},$$

$$\forall i_{\text{intra}} \in G(\mathcal{T} | mc), \forall \mathcal{T} \in \mathcal{M}, \forall mc \in \mathcal{MC} \quad (3)$$

where  $\hat{P}_{i,\mathcal{T}}$  is the predicted retrieved embedding for sample  $i$ ’s missing modality  $\mathcal{T}$ .  $\mathcal{R}(\cdot)$  denotes the router responsible for top-k expert selection, as defined in Equation 1, given an input embedding or token. Here, the input of SMoE, context-aware embedding,  $\mathbf{x}$  includes (i.e.,  $\cup$ ) both intra-modal examples ( $P'_{i_{\text{intra}},\mathcal{T}}$ ) and inter-modal examples ( $P'_{i,mc}$ ).  $P' = \text{MHA}(P)$ , where  $P$  represents the embedding after passing through the modality-specific encoder from raw feature space. This denotes the embedding or token after undergoing Multi-Head Attention (MHA), i.e., Cross-Attention, enabling interaction between tokens. Thus, tokens are endowed with not only self-modality knowledge but also inter-modal harmonization before being passed to the SMoE router.

For the expert design,  $\mathcal{E}_e^T(\mathbf{x})$  represents the modality-

specific expert, where each expert corresponds to a distinct FFN layer, is distinct and newly introduced in MoE-Retriever to enhance context-awareness, particularly in handling missing modality scenarios. Notably, the retrieval target differs for each modality combination in various samples, leading us to allocate specific expert indices for each target modality. For instance, if there are 32 experts and four modalities, each modality will have its own pool of 8 experts. Additionally, to enhance flexibility and generalizability, we include shared experts (denoted as ‘Shared’ in Figure 2), expecting that common knowledge can be leveraged across different modalities. The number of shared experts is controlled by the hyperparameter  $b$ , and we elaborate on this design in Appendix D.

After retrieving the most relevant embedding for each missing modality, we proceed to the subsequent fusion layer<sup>2</sup>, followed by the prediction head for the downstream task. Since gradients flow continuously from the input features to the output predictions, this enables end-to-end training. For the overall algorithm, please refer to Appendix A.

<sup>2</sup>The fusion layer can be based on diverse architectures, such as Transformers or even an SMoE layer. To ensure generalizability, we choose a vanilla Transformer encoder as our fusion layer and explore alternative backbones in the Experiments section.

Table 1: Performance comparison in ADNI and MIMIC Datasets. Image ( $\mathcal{I}$ ), Genetic ( $\mathcal{G}$ ), Clinical ( $\mathcal{C}$ ), and Biospecimen ( $\mathcal{B}$ ) modalities are used for ADNI dataset. For ADNI dataset, we use the image modality as a central reference, and sequentially added genetic, clinical, and finally all four modalities. Lab ( $\mathcal{L}$ ), Notes ( $\mathcal{N}$ ), and Code ( $\mathcal{C}$ ) modalities are used in MIMIC dataset.

Dataset	Modality	Metric	mmFormer	ShaSpec	M3Care	MUSE	FuseMoE	MoE-Retriever
ADNI	$\mathcal{I}+\mathcal{G}$	Acc.	50.42±4.98	54.81±4.47	48.69±4.03	43.90±2.59	52.19±4.25	<b>61.09±2.12</b>
		F1	46.66±2.40	54.43±4.11	40.29±6.49	26.83±2.68	48.22±6.28	<b>62.10±1.12</b>
	$\mathcal{I}+\mathcal{G}+\mathcal{C}$	Acc.	51.73±1.40	58.36±1.65	48.97±2.45	45.04±2.65	60.97±1.32	<b>63.12±1.19</b>
		F1	49.97±1.89	52.69±4.99	43.55±6.24	37.21±2.61	52.21±3.87	<b>62.17±2.90</b>
	$\mathcal{I}+\mathcal{G}+\mathcal{C}+\mathcal{B}$	Acc.	55.46±1.05	59.94±2.25	54.68±0.70	52.24±2.61	59.52±1.00	<b>64.52±2.55</b>
		F1	46.94±0.31	59.94±1.88	46.09±2.29	43.07±2.01	54.63±1.16	<b>63.80±2.96</b>
MIMIC	$\mathcal{L}+\mathcal{N}$	Acc.	77.37±0.00	77.37±0.15	76.14±0.46	<b>77.40±1.12</b>	60.50±3.82	76.82±3.02
		F1	43.62±0.00	55.19±1.52	45.26±0.44	51.53±1.90	51.58±1.32	<b>58.06±2.19</b>
	$\mathcal{L}+\mathcal{C}$	Acc.	77.37±0.00	77.37±0.13	76.76±0.59	<b>77.40±1.12</b>	63.31±3.21	77.20±0.47
		F1	43.62±0.00	57.32±0.52	43.92±0.52	51.53±1.90	51.24±0.60	<b>57.73±0.64</b>
	$\mathcal{N}+\mathcal{C}$	Acc.	77.37±0.00	77.40±0.03	77.26±0.35	77.32±1.13	64.77±0.36	<b>77.45±0.14</b>
		F1	43.62±0.00	54.59±0.65	45.31±1.22	51.53±1.90	48.11±1.05	<b>56.65±1.23</b>
$\mathcal{L}+\mathcal{N}+\mathcal{C}$	Acc.	77.37±0.00	<b>77.40±0.09</b>	76.04±0.70	77.40±1.12	63.90±1.72	76.59±0.07	
	F1	43.62±0.00	55.79±0.94	45.43±1.17	51.25±1.87	55.72±1.03	<b>59.74±0.81</b>	

## Experiments

**Datasets.** We evaluate MoE-Retriever on four multimodal datasets across medical and general domains. For medical data, we use the ADNI dataset (Weiner et al. 2010, 2017), which integrates imaging (MRI), genetics, clinical metrics, and biospecimens to classify Alzheimer’s Disease stages (Dementia, Cognitively Normal, or Mild Cognitive Impairment) from 2,380 samples, and the MIMIC-IV dataset (Johnson et al. 2023), consisting of structured (labs, vitals) and unstructured data (clinical notes, ICD-9 codes) from 9,003 critical care patients for one-year mortality prediction. For general multimodal datasets, we use CMU-MOSI (Zadeh et al. 2016), a video sentiment analysis dataset with 2,199 annotated clips, and ENRICO (Leiva, Hota, and Oulasvirta 2020), a collection of 1,460 app screens classified into 20 design categories. Detailed preprocessing steps for each dataset and implementation details are provided in Appendix B.

**Baselines.** We compare MoE-Retriever against various state-of-the-art baselines from three categories. (1) *feature modeling methods*: mmFormer (Zhang et al. 2022b) and ShaSpec (Wang et al. 2023). (2) *graph-based approaches*: MUSE (Wu et al. 2024) and M3Care (Zhang et al. 2022a). (3) *MoE-based method*: FuseMoE (Han et al. 2024). For details on the modality-specific encoder settings, please refer to Appendix B. For a more comprehensive discussion of related works, including these baselines is provided at Appendix E.

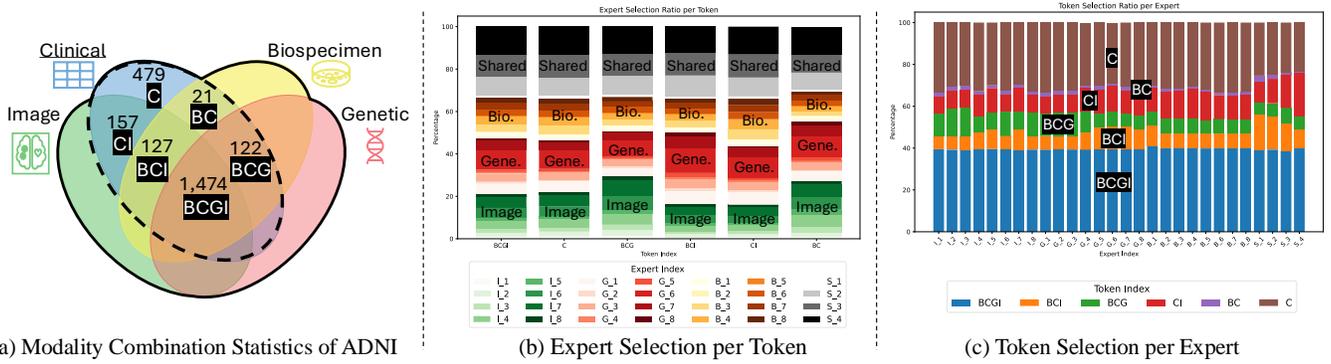
**Primary Results.** Table 1 presents several insights: **1)** On the ADNI dataset, among all modality combinations, MoE-Retriever outperforms all baselines by a notable margin. **2)** Notably, as the number of available modalities increases (e.g.,  $\mathcal{I} + \mathcal{G} + \mathcal{C} + \mathcal{B}$ ), the potential of MoE-Retriever grows, providing a large margin of improvement (7.64% gain compared to the best-performing model, ShaSpec, and 8.40% gain compared to the state-of-the-art model, FuseMoE). This shows that with more modalities, there is greater room for improvement, which can be attributed to the fact that a larger number of intra- and inter-

modal samples facilitate the retrieval process. **3)** The two graph-based methods, M3Care (Zhang et al. 2022a) and MUSE (Wu et al. 2024), perform the worst on the ADNI dataset. This suggests that while graph-based approaches capture intra-modal relationships between samples, they struggle due to the lack of handling inter-modal interactions, highlighting the importance of these interactions. **4)** FuseMoE (Han et al. 2024), a mixture-of-experts (MoE)-based method, achieves the best performance on the ADNI dataset but significantly underperforms on the MIMIC dataset<sup>3</sup>. This can be attributed to FuseMoE’s reliance on a single random embedding to impute missing modalities. **5)** On the MIMIC dataset, all baseline models suffer from the label imbalance problem, resulting in either Acc or F1 scores being biased. However, MoE-Retriever appears to be a well-balanced model, where the F1 score, being more significant than Acc in imbalanced cases, consistently outperforms all baselines. All in all, MoE-Retriever achieves notable performance gains on both datasets, thanks to its ability to model intra- and inter-modal contexts and its context-aware routing policy via the SMOE design, showcasing that better-retrieved embeddings for missing modalities lead to downstream performance improvements. For the results on ENRICO and CMU-MOSI datasets, please refer to Appendix C.

### How MoE-Retriever Contributes?

**In-depth Analysis.** To gain a deeper understanding of how MoE-Retriever functions and contributes to embedding retrieval, we provide an in-depth analysis using the ADNI dataset in Figure 3. First, as shown in Figure 3 (a), we observe six unique modality combination regions. Interestingly, the clinical modality is present in all combinations, indicating that the input token will always include the clinical

<sup>3</sup>We attempted to use the authors’ code but observed unstable performance. Thus, we borrowed FuseMoE’s performance on these datasets from the recent Flex-MoE paper (Yun et al. 2024).



(a) Modality Combination Statistics of ADNI

(b) Expert Selection per Token

(c) Token Selection per Expert

Figure 3: (a) Statistics of modality combinations observed in the ADNI dataset. We observe that although the ADNI dataset comprises four modalities, the modality combinations are not as diverse, showing only six unique regions. Notably, all modality combinations include the clinical modality. (b) Given an input token (i.e., modality combination), we track the expert selection ratio based on the modality combination. Alternatively, (c) from the expert’s perspective, we provide how each expert selects the input token and their relative ratio. The backbone illustration of (a) is adapted from (Yun et al. 2024).

(C) modality. This also suggests that the missing modality, i.e., the target modality, will often include  $\mathcal{I}$ ,  $\mathcal{G}$ , or  $\mathcal{B}$ , depending on its interaction with other modalities.

Next, after training MoE-Retriever, we track the activation ratio from both token and expert perspectives. In Figure 3 (b), we observe: **1)** MoE-Retriever successfully learns which modality should be selected and imputed. For example, when the token index is given as  $BCG$ , which lacks the  $\mathcal{I}$  modality, the majority of tokens select image-specific experts, ranging from  $I_1$  to  $I_8$ . **2)** This imputation tendency is also observed when the input token is  $BCI$  or  $BC$ , naturally incorporating the missing modality. This indicates that both the router and the experts are equipped with the knowledge of how to handle different input modality combinations. **3)** It is also notable that shared experts are frequently selected among activated experts, suggesting that these shared experts have learned and contain common knowledge that can interact with various modalities, aligning with the motivation behind designing shared experts as a buffer.

**4)** In Figure 3 (c), which shows the token selection ratio from the expert’s perspective, it is expected that  $BCGI$  is widely chosen by the experts, as this full modality combination is the majority in the ADNI dataset. This combination is frequently sampled through the supporting group, serving as a reference for missing cases. **5)** We also observe that experts select the necessary inputs, such as  $B_3, B_4, B_5$ , which most often select tokens like  $CI$ . **6)** In summary, by equipping the router and experts with the knowledge to select the most relevant embedding candidates, missing embeddings are effectively retrieved to interact with other modalities. This, in turn, boosts performance in downstream tasks by leveraging intra- and inter-modal context and a context-aware routing policy. For the ablation study and variants of each module, please refer to Appendix D.

## Computational Efficiency

In Figure 4, we compare the inference time for a single epoch, computational cost, and the number of parameters for each model across different modality configurations in the ADNI

dataset. The results show that MoE-Retriever outperforms in all three computational dimensions: **1)** Mean Time, **2)** GFLOPs, and **3)** Number of Parameters, thanks to the adoption of the SMOE design. Notably, as the modality combinations increase, the efficiency is maintained, highlighting the advantage of SMOE, which sparsely activates the relevant parameters. This represents a significant step forward in embedding retrieval design.

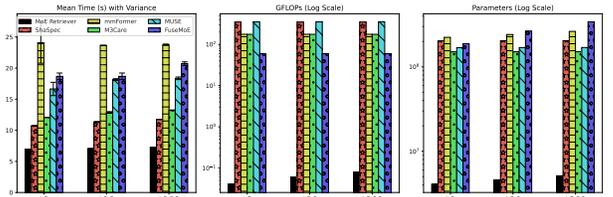


Figure 4: Comparison of computational efficiency of different methods. The left figure displays the averaged inference time for a single epoch of testing data, with error bar showing the variance. The middle plot illustrates the computational cost in GFLOPs (floating-point operations per second divided by  $10^9$ ), while the right figure shows the number of parameters on a logarithmic scale. The FLOPs and GFLOPs are computed using the fvc core package.

## Conclusion

In this work, we propose MoE-Retriever, a novel framework inspired by the SMOE design that uniquely integrates both intra-modal and inter-modal contexts. MoE-Retriever first generates context-aware embeddings from a modality combination-based supporting group for intra-modal and inter-modal contexts. Then, SMOE router selects the most relevant experts—i.e., embeddings tailored to specific missing modality scenarios. Our extensive experiments on both medical and general domain datasets demonstrate that MoE-Retriever not only enhances accuracy and robustness in missing modality scenarios but also exhibits scalability and computational efficiency.

## Acknowledgement

This research was, in part, funded by the National Institutes of Health (NIH) under other transactions 1OT2OD038045-01. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing official policies, either expressed or implied, of the NIH.

## References

- Ahmed, K.; Baig, M. H.; and Torresani, L. 2016. Network of experts for large-scale image categorization. In *European Conference on Computer Vision*, 516–532. Springer.
- Aoki, R.; Tung, F.; and Oliveira, G. L. 2021. Heterogeneous Multi-task Learning with Expert Diversity. *CoRR*, abs/2106.10595.
- Baltrušaitis, T.; Ahuja, C.; and Morency, L.-P. 2018. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2): 423–443.
- Chen, K.; Xu, L.; and Chi, H. 1999. Improved learning algorithms for mixture of experts in multiclass classification. *Neural networks*, 12(9): 1229–1252.
- Chen, Z.; Shen, Y.; Ding, M.; Chen, Z.; Zhao, H.; Learned-Miller, E. G.; and Gan, C. 2023. Mod-Squad: Designing Mixtures of Experts As Modular Multi-Task Learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 11828–11837.
- Doshi, J.; Erus, G.; Ou, Y.; Resnick, S. M.; Gur, R. C.; Gur, R. E.; Satterthwaite, T. D.; Furth, S.; Davatzikos, C.; Initiative, A. N.; et al. 2016. MUSE: MULTI-atlas region Segmentation utilizing Ensembles of registration algorithms and parameters, and locally optimal atlas selection. *Neuroimage*, 127: 186–195.
- Dubois, B.; von Arnim, C. A.; Burnie, N.; Bozeat, S.; and Cummings, J. 2023. Biomarkers in Alzheimer’s disease: role in early and differential diagnosis and recognition of atypical variants. *Alzheimer’s Research & Therapy*, 15(1): 175.
- Esmailzadeh, S.; Belivanis, D. I.; Pohl, K. M.; and Adeli, E. 2018. End-to-end Alzheimer’s disease diagnosis and biomarker identification. In *Machine Learning in Medical Imaging: 9th International Workshop, MLMI 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 9*, 337–345. Springer.
- Guo, W.; Wang, J.; and Wang, S. 2019. Deep multimodal representation learning: A survey. *Ieee Access*, 7: 63373–63394.
- Han, X.; Nguyen, H.; Harris, C.; Ho, N.; and Saria, S. 2024. Fusemoe: Mixture-of-experts transformers for fleximodal fusion. *arXiv preprint arXiv:2402.03226*.
- Hazimeh, H.; Zhao, Z.; Chowdhery, A.; Sathiamoorthy, M.; Chen, Y.; Mazumder, R.; Hong, L.; and Chi, E. H. 2021. DSelect-k: Differentiable Selection in the Mixture of Experts with Applications to Multi-Task Learning. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y. N.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, 29335–29347.
- Jack Jr, C. R.; Bennett, D. A.; Blennow, K.; Carrillo, M. C.; Dunn, B.; Haeberlein, S. B.; Holtzman, D. M.; Jagust, W.; Jessen, F.; Karlawish, J.; et al. 2018. NIA-AA research framework: toward a biological definition of Alzheimer’s disease. *Alzheimer’s & dementia*, 14(4): 535–562.
- Jacobs, R. A.; Jordan, M. I.; Nowlan, S. J.; and Hinton, G. E. 1991. Adaptive mixtures of local experts. *Neural computation*, 3(1): 79–87.
- Jiang, H.; Zhan, K.; Qu, J.; Wu, Y.; Fei, Z.; Zhang, X.; Chen, L.; Dou, Z.; Qiu, X.; Guo, Z.; et al. 2021. Towards More Effective and Economic Sparsely-Activated Model. *arXiv preprint arXiv:2110.07431*.
- Johnson, A. E.; Bulgarelli, L.; Shen, L.; Gayles, A.; Shammout, A.; Horng, S.; Pollard, T. J.; Hao, S.; Moody, B.; Gow, B.; et al. 2023. MIMIC-IV, a freely accessible electronic health record dataset. *Scientific data*, 10(1): 1.
- Jordan, M. I.; and Jacobs, R. A. 1994. Hierarchical mixtures of experts and the EM algorithm. *Neural computation*, 6(2): 181–214.
- Khader, F.; Kather, J. N.; Müller-Franzes, G.; Wang, T.; Han, T.; Tayebi Arasteh, S.; Hamesch, K.; Bressemer, K.; Haarbuerger, C.; Stegmaier, J.; et al. 2023. Medical transformer for multimodal survival prediction in intensive care: integration of imaging and non-imaging data. *Scientific Reports*, 13(1): 10666.
- Kim, S. Y. 2023. Personalized Explanations for Early Diagnosis of Alzheimer’s Disease Using Explainable Graph Neural Networks with Population Graphs. *Bioengineering*, 10(6): 701.
- Kudugunta, S.; Huang, Y.; Bapna, A.; Krikun, M.; Lepikhin, D.; Luong, M.; and Firat, O. 2021. Beyond Distillation: Task-level Mixture-of-Experts for Efficient Inference. In Moens, M.; Huang, X.; Specia, L.; and Yih, S. W., eds., *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, 3577–3599. Association for Computational Linguistics.
- Lambert, J.-C.; Ibrahim-Verbaas, C. A.; Harold, D.; Naj, A. C.; Sims, R.; Bellenguez, C.; Jun, G.; DeStefano, A. L.; Bis, J. C.; Beecham, G. W.; et al. 2013. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer’s disease. *Nature genetics*, 45(12): 1452–1458.
- Leiva, L. A.; Hota, A.; and Oulasvirta, A. 2020. Enrico: A dataset for topic modeling of mobile UI designs. In *22nd International Conference on Human-Computer Interaction with Mobile Devices and Services*, 1–4.
- Lepikhin, D.; Lee, H.; Xu, Y.; Chen, D.; Firat, O.; Huang, Y.; Krikun, M.; Shazeer, N.; and Chen, Z. 2021. GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

- Li, L.; Yu, X.; Sheng, C.; Jiang, X.; Zhang, Q.; Han, Y.; and Jiang, J. 2022. A review of brain imaging biomarker genomics in Alzheimer's disease: implementation and perspectives. *Translational Neurodegeneration*, 11(1): 42.
- Liu, M.; Li, S.; Yuan, H.; Ong, M. E. H.; Ning, Y.; Xie, F.; Saffari, S. E.; Shang, Y.; Volovici, V.; Chakraborty, B.; et al. 2023. Handling missing values in healthcare data: A systematic review of deep learning-based imputation techniques. *Artificial intelligence in medicine*, 142: 102587.
- Lou, Y.; Xue, F.; Zheng, Z.; and You, Y. 2021. Cross-token Modeling with Conditional Computation. *arXiv preprint arXiv:2109.02008*.
- Ma, J.; Zhao, Z.; Yi, X.; Chen, J.; Hong, L.; and Chi, E. H. 2018. Modeling Task Relationships in Multi-task Learning with Multi-gate Mixture-of-Experts. In Guo, Y.; and Farooq, F., eds., *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, 1930–1939. ACM.
- Ma, M.; Ren, J.; Zhao, L.; Tulyakov, S.; Wu, C.; and Peng, X. 2021. Smil: Multimodal learning with severely missing modality. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 2302–2310.
- Malitesta, D.; Rossi, E.; Pomo, C.; Di Noia, T.; and Malliaros, F. D. 2024. Do We Really Need to Drop Items with Missing Modalities in Multimodal Recommendation? *arXiv preprint arXiv:2408.11767*.
- Mustafa, B.; Riquelme, C.; Puigcerver, J.; Jenatton, R.; and Houlsby, N. 2022. Multimodal Contrastive Learning with LIMoE: the Language-Image Mixture of Experts. *CoRR*, abs/2206.02770.
- Ou, Y.; Sotiras, A.; Paragios, N.; and Davatzikos, C. 2011. DRAMMS: Deformable registration via attribute matching and mutual-saliency weighting. *Medical image analysis*, 15(4): 622–639.
- Pan, Y.; Liu, M.; Xia, Y.; and Shen, D. 2021. Disease-image-specific learning for diagnosis-oriented neuroimage synthesis with incomplete multi-modality data. *IEEE transactions on pattern analysis and machine intelligence*, 44(10): 6839–6853.
- Riquelme, C.; Puigcerver, J.; Mustafa, B.; Neumann, M.; Jenatton, R.; Pinto, A. S.; Keysers, D.; and Houlsby, N. 2021. Scaling Vision with Sparse Mixture of Experts. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y. N.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, 8583–8595.
- Shazeer, N.; Mirhoseini, A.; Maziarz, K.; Davis, A.; Le, Q.; Hinton, G.; and Dean, J. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*.
- Steyaert, S.; Pizurica, M.; Nagaraj, D.; Khandelwal, P.; Hernandez-Boussard, T.; Gentles, A. J.; and Gevaert, O. 2023. Multimodal data fusion for cancer biomarker discovery with deep learning. *Nature machine intelligence*, 5(4): 351–362.
- Wang, H.; Chen, Y.; Ma, C.; Avery, J.; Hull, L.; and Carneiro, G. 2023. Multi-modal learning with missing modality via shared-specific feature modelling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15878–15887.
- Wang, X.; Yu, F.; Dunlap, L.; Ma, Y.-A.; Wang, R.; Mirhoseini, A.; Darrell, T.; and Gonzalez, J. E. 2020. Deep mixture of experts via shallow embedding. In *Uncertainty in artificial intelligence*, 552–562. PMLR.
- Weiner, M. W.; Aisen, P. S.; Jack Jr, C. R.; Jagust, W. J.; Trojanowski, J. Q.; Shaw, L.; Saykin, A. J.; Morris, J. C.; Cairns, N.; Beckett, L. A.; et al. 2010. The Alzheimer's disease neuroimaging initiative: progress report and future plans. *Alzheimer's & Dementia*, 6(3): 202–211.
- Weiner, M. W.; Veitch, D. P.; Aisen, P. S.; Beckett, L. A.; Cairns, N. J.; Green, R. C.; Harvey, D.; Jack Jr, C. R.; Jagust, W.; Morris, J. C.; et al. 2017. The Alzheimer's Disease Neuroimaging Initiative 3: Continued innovation for clinical trial improvement. *Alzheimer's & Dementia*, 13(5): 561–571.
- Wu, R.; Wang, H.; and Chen, H.-T. 2024. A Comprehensive Survey on Deep Multimodal Learning with Missing Modality. *arXiv preprint arXiv:2409.07825*.
- Wu, Z.; Dadu, A.; Tustison, N.; Avants, B.; Nalls, M.; Sun, J.; and Faghri, F. 2024. Multimodal patient representation learning with missing modalities and labels. In *The Twelfth International Conference on Learning Representations*.
- Yao, W.; Yin, K.; Cheung, W. K.; Liu, J.; and Qin, J. 2024. DrFuse: Learning Disentangled Representation for Clinical Multi-Modal Fusion with Missing Modality and Modal Inconsistency. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 16416–16424.
- Yuksel, S. E.; Wilson, J. N.; and Gader, P. D. 2012. Twenty Years of Mixture of Experts. *IEEE Transactions on Neural Networks and Learning Systems*, 23(8): 1177–1193.
- Yun, S.; Choi, I.; Peng, J.; Wu, Y.; Bao, J.; Zhang, Q.; Xin, J.; Long, Q.; and Chen, T. 2024. Flex-MoE: Modeling Arbitrary Modality Combination via the Flexible Mixture-of-Experts. *arXiv preprint arXiv:2410.08245*.
- Zadeh, A.; Zellers, R.; Pincus, E.; and Morency, L.-P. 2016. Mosi: multimodal corpus of sentiment intensity and subjectivity analysis in online opinion videos. *arXiv preprint arXiv:1606.06259*.
- Zhang, C.; Chu, X.; Ma, L.; Zhu, Y.; Wang, Y.; Wang, J.; and Zhao, J. 2022a. M3care: Learning with missing modalities in multimodal healthcare data. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2418–2428.
- Zhang, Y.; He, N.; Yang, J.; Li, Y.; Wei, D.; Huang, Y.; Zhang, Y.; He, Z.; and Zheng, Y. 2022b. mmformer: Multimodal medical transformer for incomplete multimodal learning of brain tumor segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 107–117. Springer.
- Zhang, Y.; Peng, C.; Wang, Q.; Song, D.; Li, K.; and Zhou, S. K. 2024. Unified multi-modal image synthesis for missing modality imputation. *IEEE Transactions on Medical Imaging*.

Zhang, Z.; Lin, Y.; Liu, Z.; Li, P.; Sun, M.; and Zhou, J. 2021. Moefication: Conditional computation of transformer models for efficient inference. *arXiv preprint arXiv:2110.01786*.

Zhou, T.; Ruan, S.; and Hu, H. 2023. A literature survey of MR-based brain tumor segmentation with missing modalities. *Computerized Medical Imaging and Graphics*, 104: 102167.

Zuo, S.; Liu, X.; Jiao, J.; Kim, Y. J.; Hassan, H.; Zhang, R.; Gao, J.; and Zhao, T. 2022. Taming Sparsely Activated Transformer with Stochastic Experts. In *International Conference on Learning Representations*.

## A Overall Algorithm

To summarize, the overall algorithm of MoE-`Retriever` is detailed in Algorithm 1.

---

Algorithm 1: The overall procedure of MoE-`Retriever`.

---

```

1: Input: Samples,  $i \leq N$ , Supporting Group,  $G(\mathcal{T} \mid mc)$ ,
   Modality Set,  $\mathcal{M}$ , Modality Combination Set,  $mc$ 
2: Output: Retrieved Embedding for Missing Modality,  $\mathcal{T}$ 
3: for  $i = 1, \dots, N$  do
4:   if  $|mc_i| < |\mathcal{M}|$  :
5:     for  $t \in \mathcal{T}_i$  do
6:        $\mathbf{x} = []$ 
7:       /* Intra-Modal Context */
8:       Samples  $\sim G(t \mid mc_i)$ 
9:       for  $j \in \text{Samples}$  do
10:         $\mathbf{x}.\text{append}(\mathbf{P}_{j,\mathcal{T}})$ 
11:       end for
12:       /* Inter-Modal Context */
13:       for  $mc \in mc_i$  do
14:         $\mathbf{x}.\text{append}(\mathbf{P}_{i,mc})$ 
15:       end for
16:       /* Context-Aware Routing Policy */
17:        $\mathbf{x} \leftarrow \text{MHA}(\mathbf{x})$ 
18:        $\hat{\mathbf{P}}_{i,\mathcal{T}} \leftarrow \text{SMoE}(\mathbf{x}, \mathcal{R}, \mathcal{E}^{\mathcal{T}}, \text{top-k})$ 
19:     end for
20: end for

```

---

## B Details on Datasets

We followed the same preprocessing procedure of the ADNI dataset and MIMIC dataset, as described in Flex-MoE (Yun et al. 2024).

### Detailed Data Preprocessing in ADNI

**Image Modality** To preprocess the image data, we first applied a correction for magnetic field intensity inhomogeneity to ensure consistency and reliability across MRI images. Next, we used the MUSE (Multiatlas Region Segmentation Utilizing Ensembles of Registration Algorithms and Parameters) method to segment gray matter tissue, the primary focus of this study (Doshi et al. 2016). This technique involves utilizing multiple atlases to extract the most accurate region-of-interest values from the segmented gray matter. Afterward, voxel-wise volumetric maps of tissue regions were created by spatially aligning skull-stripped images to a template in the Montreal Neurological Institute (MNI) space, using a registration method (Ou et al. 2011).

**Genetic Modality** We obtained SNP (single nucleotide polymorphisms) data from the ADNI 1, GO/2, and 3 studies, and pre-processed it as follows. First, SNP data from these studies were aligned to a unified reference build using Liftover <https://liftover.broadinstitute.org/>, converting all data to NCBI build 37 (UCSC hg19). Next, we aligned strands based on the 1000 Genome Project phase 3, using McCarthy Group Tools <https://www.well.ox.ac.uk/~wrayner/tools/>. Linkage disequilibrium (LD) pruning was then applied with parameters (50, 5, 0.1) to remove highly correlated SNPs, reducing the total SNPs from 565,989 to 144,746. Imputation was performed on this pruned set using the Michigan Imputation

Server <https://imputationserver.sph.umich.edu/index.html#!>, and the resulting SNP data was recoded as  $\{0, 1, 2\}$ .

**Biospecimen Modality** Biospecimen data was extracted from several ADNI-provided csv files. CSF A $\beta$ 1-42 and A $\beta$ 1-40 data were taken from ISOPROSTANE\_09May2024.csv, Total Tau and Phosphorylated Tau from UPENNBIOMK\_ROCHE\_ELECSYS\_09May2024.csv, Plasma Neurofilament Light Chain data from batemanlab\_20221118\_09May2024.csv, and ApoE genotype data from APOERES\_09May2024.csv. Numerical data was scaled using a MinMax scaler to a range of -1 to 1, while categorical data was one-hot encoded. For missing values, we imputed the mean for numerical fields and the mode for categorical fields.

**Clinical Modality** Clinical data was extracted from ADNI csv files, including MEDHIST\_09May2024.csv, NEUROEXM\_09May2024.csv, PTDEMOG\_09May2024.csv, RECCMEDS\_09May2024.csv, and VITALS\_09May2024.csv. During preprocessing, we excluded the columns 'PTCOGBEG,' 'PTADDX,' and 'PTADBEG,' which contain direct Alzheimer’s Disease diagnosis information. Numerical data was scaled using a MinMax scaler (-1 to 1), while categorical data was one-hot encoded. Missing values were imputed by using the mean for numerical columns and the mode for categorical columns.

### Detailed Data Preprocessing in MIMIC

**Lab, Notes, Codes Modalities.** For the MIMIC dataset, we use the Medical Information Mart for Intensive Care IV (MIMIC-IV) database, which contains de-identified health data for patients who were admitted to either the emergency department or stayed in critical care units of the Beth Israel Deaconess Medical Center in Boston, Massachusetts24. MIMIC-IV excludes patients under 18 years of age. We take a subset of the MIMIC-IV data, where each patient has at least more than 1 visit in the dataset as this subset corresponds to patients who likely have more serious health conditions. For each datapoint, we extract ICD-9 codes, clinical text, and labs and vital values. Using this data, we perform binary classification on one-year mortality, which foresees whether or not this patient will pass away in a year. We drop visits that occur at the same time as the patient’s death.

**Missingness in MIMIC dataset. Code Modality:** This combines diagnosis and procedure data. There are 4 records with missing diagnoses and 1777 with missing procedures. **Note Modality:** Derived from the “text” column of the original CSV file, there are 108 records with missing notes. **Lab Modality:** This presents a more complex scenario, as it includes 2172 different measurements. If we consider all 2172 measurements as potentially missing, then technically, there is no missing data since essential measurements, like heart rate, are consistently collected for each patient. However, if we evaluate the proportion of missing values in the (9003, 2172) matrix, we find that 94.216% of the entries are NaN.

**Implementations.** To ensure a fair comparison with other baselines, we utilized the optimal hyperparameter settings provided in the original papers. For dataset split, we choose 70% for training, 15% as validation set, and the remaining 15% for testing. Both the ADNI and MIMIC datasets contain

missing data. For the CMU-MOSI and ENRICO datasets, we applied random dropping with probability of 0.3 for each modality independently to simulate missing modality scenarios. Given the incomplete nature of the datasets, if a baseline implementation could impute or interact with other modalities, we leveraged those methods. Otherwise, we used zero-padding to support batch-wise training. All experiments were conducted on NVIDIA A100 GPUs. Each experiment was run three times with different seeds to ensure reproducibility, and the results were averaged.

### Modality-specific Encoder Settings

**ADNI Dataset.** For image modality, we used a customized 3D-CNN (Esmaeilzadeh et al. 2018) with hidden dimension 256 as encoder. For genomics, clinical, and biospecimen modalities, we used MLP with hidden dimension 256 as encoder. **MIMIC Dataset.** For all lab, note, and code modalities, we used LSTM with hidden dimension 256 as encoder. **ENRICO Dataset.** For both screenshot image and wireframe image modality, we used VGG11 from torchvision library with hidden dimension size 16 as encoder. **CMU-MOSI Dataset.** For both vision, audio, and text modality, we used Gated Recurrent Unit with hidden dimension 256 as encoder.

## C More Results

**Results on ENRICO and CMU-MOSI Datasets.** Table 2 shows the performance across generalized domains: design motifs for the ENRICO dataset and sentiment analysis for the CMU-MOSI dataset. We observe that **1) MoE-Retriever** outperforms current multimodal baselines, demonstrating its generalizability across diverse multimodal domains. Specifically, in the CMU-MOSI dataset, we observe **2)** that as the number of modalities increases, the performance of existing baselines improves, but the increase does not surpass that of **MoE-Retriever**, highlighting its effectiveness as a strong benchmark model for various domains and modality combinations.

## D Ablation Study

To verify the effectiveness of **MoE-Retriever**, we conducted an extensive ablation study using the ADNI dataset in the  $\mathcal{I} + \mathcal{G} + \mathcal{C} + \mathcal{B}$  scenario in Table 3. Key observations include: **1)** Regarding the core module design in **MoE-Retriever**, involving inter-modal context is crucial as it personalizes the specific observed modality context of each sample. **2)** When designing shared experts ( $E_{sh}$ ), it is important to strike a balance in the number of shared experts. Having too many can deteriorate the acquisition of specialized knowledge required by modality-specific experts. **3)** For modality-specific experts, selecting too few or too many experts can lead to suboptimal results, emphasizing the need for a balanced number, such as eight. **4)** For the router design, utilizing a single router to handle both intra- and inter-modal contexts proved to be sufficient. The more examples it encounters during training, the more knowledge it is able to accumulate. **5)** In the subsequent fusion layer, we experimented with both a vanilla transformer design and

Dataset	Modality	mmFormer	ShaSpec	M3Care	MUSE	FuseMoE	MoE-Retriever
ENRICO	$\mathcal{S}+\mathcal{W}$	36.19 $\pm$ 0.98	21.03 $\pm$ 0.32	19.06 $\pm$ 5.17	36.01 $\pm$ 2.81	36.99 $\pm$ 6.83	<b>38.24<math>\pm</math>1.16</b>
CMU-MOSI	$\mathcal{V}+\mathcal{A}$	42.23 $\pm$ 0.00	50.91 $\pm$ 1.63	42.23 $\pm$ 0.00	44.64 $\pm$ 1.94	47.46 $\pm$ 2.36	<b>53.12<math>\pm</math>2.26</b>
	$\mathcal{V}+\mathcal{T}$	62.20 $\pm$ 0.90	60.01 $\pm$ 1.44	42.12 $\pm$ 0.14	52.54 $\pm$ 1.92	63.77 $\pm$ 1.62	<b>65.74<math>\pm</math>0.55</b>
	$\mathcal{A}+\mathcal{T}$	65.65 $\pm$ 0.63	65.09 $\pm$ 1.02	47.05 $\pm$ 6.83	50.82 $\pm$ 1.91	61.33 $\pm$ 0.93	<b>66.13<math>\pm</math>0.69</b>
	$\mathcal{V}+\mathcal{A}+\mathcal{T}$	62.75 $\pm$ 1.12	64.02 $\pm$ 0.65	42.23 $\pm$ 0.00	50.66 $\pm$ 1.93	60.67 $\pm$ 0.22	<b>65.21<math>\pm</math>2.72</b>

Table 2: Performance comparison in ENRICO and CMU-MOSI Datasets. Screenshot ( $\mathcal{S}$ ), and Wireframe ( $\mathcal{W}$ ) modalities are used for ENRICO dataset. Vision ( $\mathcal{V}$ ), Audio ( $\mathcal{A}$ ), and Text ( $\mathcal{T}$ ) modalities are used in CMU-MOSI dataset. We report Accuracy (Acc.) for both datasets.

a version with the SMOE layer attached. However, no significant performance gain was observed, suggesting that the utilization of SMOE in embedding retrieval was sufficient.

Table 3: Ablation Study.

Model Variants	Acc.	F1
MoE-Retriever	<b>64.52<math>\pm</math>2.55</b> ( $ E_{\mathcal{T}} =8,  E_{Sh.} =4,  \mathcal{R} =1$ )	<b>63.80<math>\pm</math>2.96</b>
w/o Intra-Modal Context	61.26 $\pm$ 2.33	61.80 $\pm$ 1.67
w/o Inter-Modal Context	60.97 $\pm$ 1.50	61.60 $\pm$ 0.78
w/o Context-Aware Routing	62.34 $\pm$ 1.25	63.11 $\pm$ 2.11
$ E_{\mathcal{T}} =8,  E_{Sh.} =1,  \mathcal{R} =1$	60.60 $\pm$ 1.32	59.70 $\pm$ 1.26
$ E_{\mathcal{T}} =8,  E_{Sh.} =2,  \mathcal{R} =1$	63.77 $\pm$ 1.35	62.92 $\pm$ 0.28
$ E_{\mathcal{T}} =8,  E_{Sh.} =8,  \mathcal{R} =1$	62.98 $\pm$ 0.79	62.75 $\pm$ 1.41
$ E_{\mathcal{T}} =4,  E_{Sh.} =4,  \mathcal{R} =1$	63.14 $\pm$ 2.47	60.88 $\pm$ 2.21
$ E_{\mathcal{T}} =16,  E_{Sh.} =4,  \mathcal{R} =1$	60.14 $\pm$ 2.97	59.91 $\pm$ 1.22
$ E_{\mathcal{T}} =8,  E_{Sh.} =4,  \mathcal{R} =2$	61.14 $\pm$ 1.85	61.04 $\pm$ 1.12
$ E_{\mathcal{T}} =8,  E_{Sh.} =4,  \mathcal{R} =4$	60.54 $\pm$ 2.52	60.23 $\pm$ 2.71
2 x Transformer Layer	63.34 $\pm$ 0.97	62.79 $\pm$ 1.31
Sparse MoE Fusion Layer	62.84 $\pm$ 2.85	63.11 $\pm$ 2.25

## E Related Work

**Multimodal Learning with Missing Modality.** Multimodal learning has garnered increasing attention in the machine learning community, particularly in the medical domain, where clinical data is inherently multimodal (Khader et al. 2023; Steyaert et al. 2023). However, in real-world clinical practice, missing modalities are a common challenge (Zhou, Ruan, and Hu 2023; Liu et al. 2023). To address this issue, one straightforward approach is to leverage generative models to impute the missing modalities (Pan et al. 2021; Zhang et al. 2024). Nonetheless, generative modeling of another distribution is a ill-posed problem (Zhang et al. 2022a). In contrast, non-generative approaches have emerged, utilizing techniques such as graph-based modeling (Wu et al. 2024), and modality fusion (Zhang et al. 2022b; Wang et al. 2023; Yao et al. 2024). While these methods can harness both inter-patient and intra-patient information, they face challenges related to scalability and struggle to handle fleximodal scenarios (Han et al. 2024), where any combination of modalities

may be present. To improve scalability, FuseMoE (Han et al. 2024) introduced a sparse Mixture-of-Experts (MoE) model aims to be robust to any combination of missing modality scenario. However, despite its scalability advantages, FuseMoE do not explicitly account both the inter-patient and intra-patient relationships simultaneously, limiting its ability to fully utilize the multimodal context of clinical data.

**Sparse Mixture-of-Experts (SMoE).** SMoE (Shazeer et al. 2017) builds on the traditional Mixture-of-Experts (MoE) model (Jacobs et al. 1991; Jordan and Jacobs 1994; Chen, Xu, and Chi 1999; Yuksel, Wilson, and Gader 2012) by introducing sparsity, which enhances both computational efficiency and model performance. By selectively activating only the most relevant experts for a specific task, SMoE minimizes overhead and improves scalability, making it particularly useful for complex, high-dimensional datasets across various applications. It has been widely applied in both vision (Riquelme et al. 2021; Lou et al. 2021; Ahmed, Baig, and Torresani 2016; Wang et al. 2020) and language processing (Lepikhin et al. 2021; Zhang et al. 2021; Zuo et al. 2022; Jiang et al. 2021). Its capacity to dynamically allocate different network parts to specific tasks (Ma et al. 2018; Chen et al. 2023) or data modalities (Kudugunta et al. 2021) has been explored for various applications (Mustafa et al. 2022). Research shows its effectiveness in areas like classification tasks for digital number recognition (Hazimeh et al. 2021) and medical signal processing (Aoki, Tung, and Oliveira 2021). However, the current use of SMoE is often biased toward its role as a backbone design, typically integrated into Transformer architectures to improve embedding representations in fusion or prediction layers. Its potential for more effective use, such as serving as a retriever or supplementing missing embeddings to bridge the feature space and encoder space, remains underexplored.