# Harnessing Heterogeneous Healthcare Data: An Attention Neural Network Approach

**Alan Devkota[1], Xuqing Wu[1], Xin Fu[1], Amany Farag[2], and Omar Ahmed[3], and Renjie Hu[1*]**

[1]Cullen College of Engineering, University of Houston, Houston, TX

[2]College of Nursing, The University of Iowa, Iowa City, IA

[3]Division of Rhinology, Sinus, Sleep & Skull Base Surgery, Houston Methodist Hospital, Houston, TX    *: rhu7@central.uh.edu

## Introduction

**Heterogeneous Healthcare Data:** In the evolving landscape of healthcare, the surge of diverse data types - from Electronic Health Records (EHR) and wearable sensor data to genomic data and medical images - presents both opportunities and challenges. Each of these data type has its distinct format and complexity.

**Comprehensive Decision-Making:** While physicians routinely integrate insights from multiple sources, many current machine learning (ML) models remain narrowly focused on a single data modality, which limits the ML models' ability to make decisions from a holistic perspective.

**Missing Modality:** An added layer of complexity arises from the prevalent issue of data completeness. It's a frequent occurrence to encounter missing modalities in healthcare data, introducing inconsistencies across patient datasets. For instance, different data modalities may be collected for different patients. Such disparities pose significant challenges for both ML and statistical analysis.
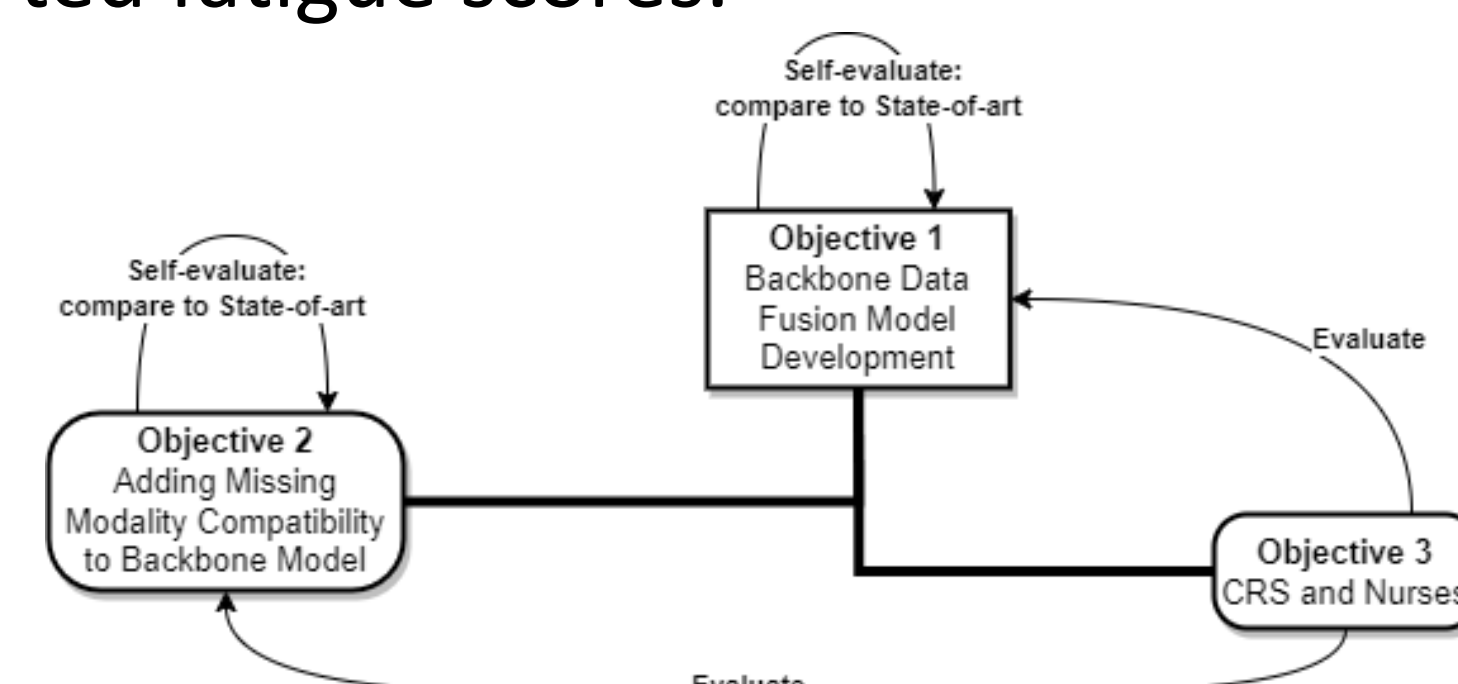
To truly harness the vast potential of heterogeneous data, there's a pressing need to develop advanced ML models that can seamlessly integrate multiple data sources, ensuring comprehensive decision-making and addressing the challenges posed by missing modalities.
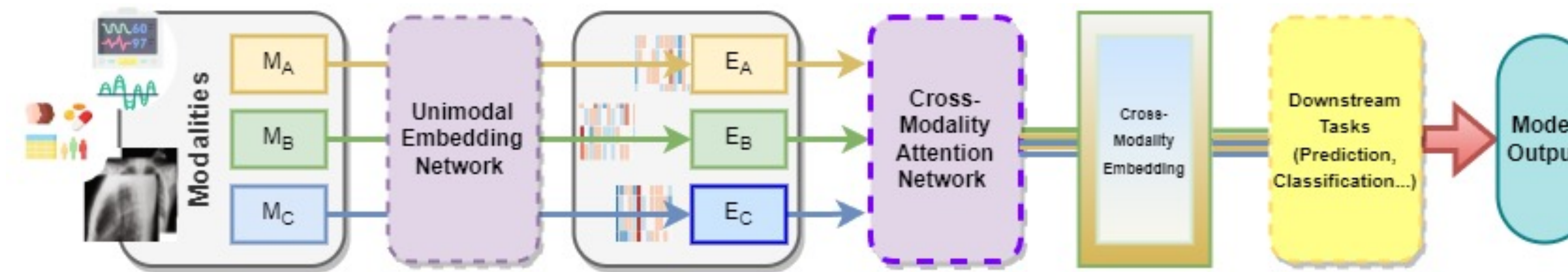
## Objectives

**Data Fusion Model Development:** Develop an attention neural network-based method for fusing heterogeneous healthcare data, emphasizing cross-modality attention transformer blocks for optimal modality integration.

**Adaptation to Missing Modalities:** Integrate Prompt Learning to enhance the model's performance with datasets that have missing modalities, preserving the core structure of the model and optimizing computational resources.
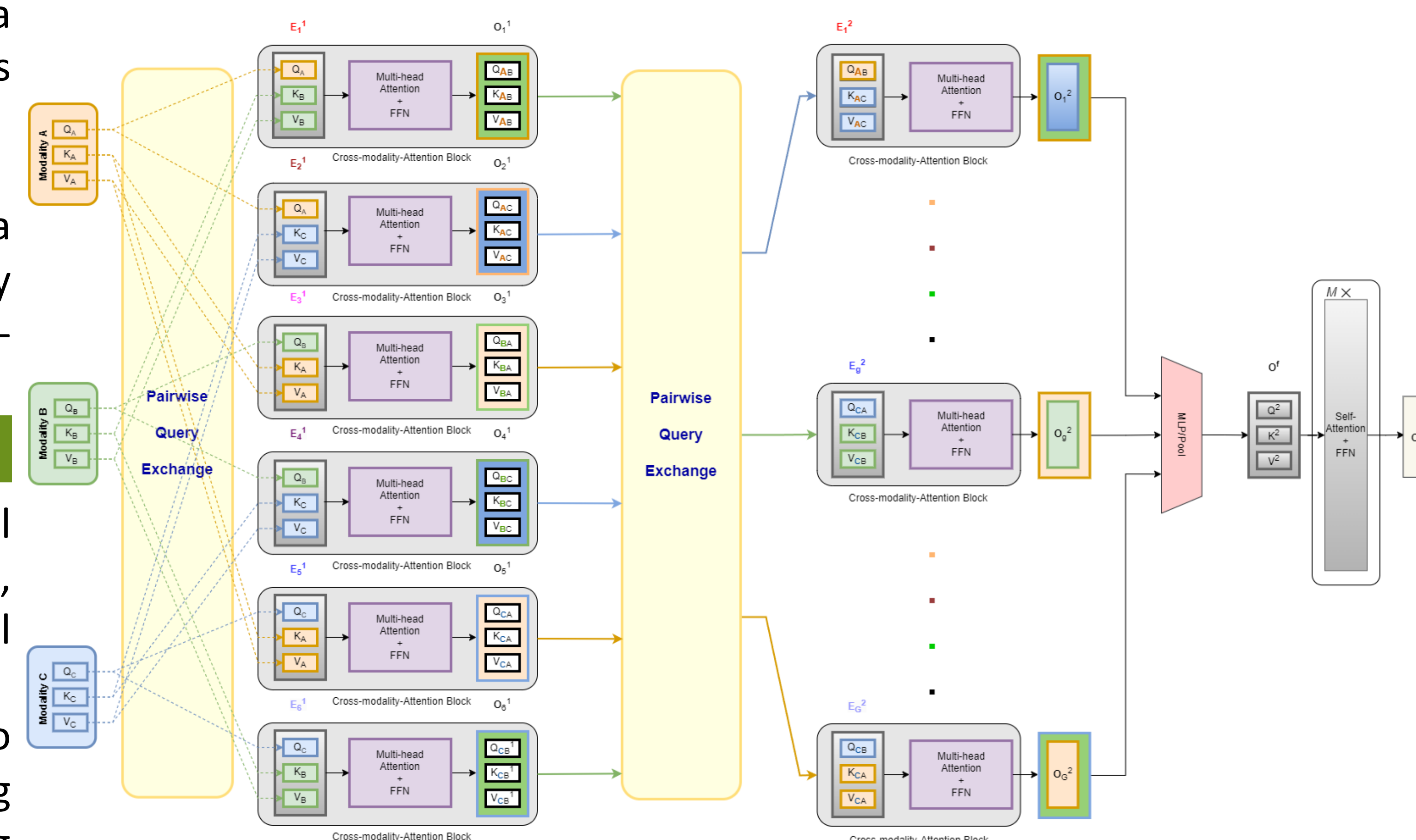
**Validation with Nurse Sleep Data and CRS Data:** Evaluate the model's performance 1) using the "All of Us" dataset and Houston Methodist's dataset focused on CRS patients, considering various data modalities like diagnosis, medication, procedures, demographics, and CT scans, 2) using the University of Iowa's Nursing College's data to analyze the correlation between nurse sleep patterns, fatigue, and shift strategies, utilizing survey data, wearable sensor data, and self-reported fatigue scores.
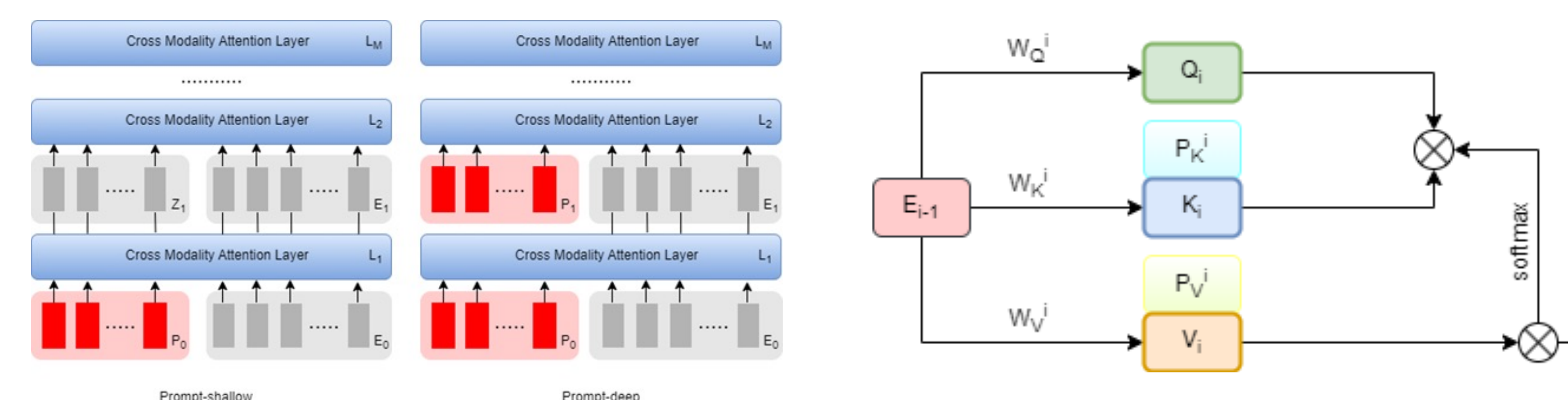


## Attention-Based NN Fusion Model



Unimodal Embedding Network (UMEN) learns embeddings for each data modality, ensuring the intrinsic characteristics of each are optimally captured. Cross-Modality Attention Network (CMAN) attends to different modalities' embedding and learns a cross-modalities information representation of the input. Post-fusion Downstream Network (DSN) serves as the final task-specific output layer after MCAN that contains the final loss function that enables the backpropagation process, ensuring cohesive training across all three modules.



## Prompt Learning for Missing Modalities



In this research, we present a data-driven approach that capitalizes on prompt learning to effectively address the challenge of missing modalities in healthcare data analysis. Prompting can be implemented by prepending a series of learnable tokens to the input of the transformer as shown in above left figure.

Prompt-based tuning methods can be classified into two categories: prompt-shallow and prompt-deep. During the training, we first train the transformer without prompts. Then only the prompts are updated through backpropagation during fine-tuning and the transformer backbone is kept frozen. Attention prompt prepends prompts to the "key" vector and "query" vector of each transformer layer as shown in above right figure.

## Nurse Sleep Data and CRS Data

| Nurse Fatigue and Sleep Dataset | |
|---|---|
| **Data Modalities** | **Data** |
| Text Data | Individual interviews (n=42) |
| Longitudinal Data | Ecological Momentary Assessment (Fatigue Score) (n = 675 (approximately 29,000 Fatigue Scores)) |
| Tabular Data | Surveys (n=1,136) |
| Time Series | Wearable Actigraphy (Fatigue Science) (n=117) |

| CRS Patient Dataset (n = 1380) | |
|---|---|
| **Data Modalities** | **Data** |
| Text Data | Chief Complaint and Symptoms, History of Present Illness, Pathology Reports, Patient Communications (Email/Telephone) |
| Longitudinal EHR | Physical Exam, Procedures, Medications, Immunizations, Appointment Schedules |
| Image Data | PET, CT, MRI, X-ray, ultrasound |
| Tabular Data | Demographic data, Social Determinants of Health, Insurance Status, REDCap Database* |

*The University of Iowa's Nurse Fatigue and Sleep dataset.* The dataset encompasses three distinct modalities. Survey Data: This encompasses 147 distinct measurements, capturing demographics, work environment dynamics, sleep quality metrics, fatigue levels, and instances of medication errors. Wearable Actigraphy Data: Leveraging Fatigue Science® technology, this data provides minute-by-minute monitoring of alertness and sleep quality over a 17-day period. Fatigue Score Reports: These are longitudinal scores, ranging from 0 to 10, self-reported by nurses through text messages during their shifts and rest periods. Depending on shift lengths, these scores are typically reported at one to two-hour intervals, resulting in over 29,000 individual fatigue report records.

*CRS Patient Dataset.* The Houston Methodist CRS Patient Dataset is a comprehensive multimodal dataset. The dataset is primarily sourced from the REDCap Database and EHR. It includes a wide range of variables, from basic demographic information and symptom questionnaires to clinical histories and social determinants of health.