# 411 A Appendix / supplemental material

## 412 A.1 Computational Results

Subgroup	Single Task Scores	Single-task Confidence Interval	Multi Task Scores	Multi-task Confidence Interval
Hospital Operations				
Hospital Aquired Infection		0.78 (0.74, 0.82)		0.88 (0.87, 0.89)
Length of Stay		0.56 (0.56, 0.57)	<ul> <li>•</li> </ul>	0.58 (0.57, 0.59)
Mortality in 48hr	Here	0.81 (0.77, 0.84)	(H)	0.85 (0.83, 0.87)
Patient Phenotyping				
Patient Phenotyping	•	0.41 (0.41, 0.42)		0.67 (0.67, 0.68)
Thoracic Testing				
Thoracic Testing		0.77 (0.76, 0.78)		0.79 (0.78, 0.8)
Blood Disorder				
Anemia		0.79 (0.78, 0.79)		0.81 (0.81, 0.81)
Cardiology		0.13 (0.10, 0.13)		0.01 (0.01, 0.01)
Heart Eailura		0.80 (0.80, 0.0)		0.01 (0.0.0.01)
Inchamic Heart Disease		0.03 (0.03, 0.3)		0.91 (0.9, 0.91)
Critical Care		0.79 (0.78, 0.8)		0.03 (0.02, 0.04)
Critical Care		0.00 (0.00, 0.0)		0.0 (0.0, 0.04)
Sepsis		0.89 (0.88, 0.9)		0.9 (0.9, 0.91)
Stroke	-	0.83 (0.82, 0.84)	-	0.87 (0.86, 0.87)
Dermatology				
Psoriasis		0.69 (0.6, 0.77)		0.93 (0.9, 0.97)
Endocrinology				
Diabetes		0.88 (0.87, 0.89)		0.89 (0.88, 0.9)
Obesity		0.87 (0.86, 0.88)		0.91 (0.9, 0.92)
Osteoporosis	101	0.76 (0.74, 0.78)		0.87 (0.86, 0.88)
Gastroenterology and Hepatology				
Chronic Liver Disease		0.9 (0.89, 0.91)		0.93 (0.92, 0.94)
Inflammatory Bowel Disease		0.67 (0.57, 0.77)	101	0.84 (0.82, 0.87)
Infectious Diseases				
Antimicrobial Resistance		0.7 (0.69, 0.72)	-	0.86 (0.85, 0.88)
Bacterial Intestinal Infections	• • • • • • • • • • • • • • • • • • •	0.78 (0.77, 0.78)	100 C 100 C	0.79 (0.79, 0.8)
Hepatitis BC	-	0.89 (0.87, 0.9)		0.91 (0.9. 0.92)
HIV		0.95 (0.94, 0.97)		0.97 (0.97, 0.98)
Tuberculosis		0.75 (0.7, 0.8)		0.95 (0.94, 0.96)
Internal Medicine		0.10 (0.1, 0.0)		0.00 (0.04, 0.00)
Eibromvalaia		0.91 (0.9. 0.92)		0.96 (0.94, 0.97)
		0.91 (0.91 0.92)		0.96 (0.94, 0.97)
Hypertension		0.01 (0.01, 0.02)		0.65 (0.65, 0.65)
Nephrology		0.00 (0.07, 0.00)		0.0 (0.0, 0.04)
Chronic Kidney Disease		0.88 (0.87, 0.88)		0.9 (0.9, 0.91)
Neurology				
Alzheimer Disease		0.97 (0.96, 0.99)		1.0 (1.0, 1.0)
Epilepsy		0.8 (0.76, 0.83)		0.86 (0.83, 0.88)
Multiple Sclerosis		0.62 (0.57, 0.66)		0.87 (0.79, 0.95)
Parkinson Disease		0.84 (0.75, 0.93)		0.99 (0.98, 1.0)
Oncology				
Breast Cancer		0.89 (0.85, 0.92)		0.9 (0.88, 0.92)
Leukemia	He-1	0.81 (0.78, 0.84)	10-1	0.87 (0.85, 0.9)
Lung Cancer	Hel	0.79 (0.77, 0.81)		0.89 (0.88, 0.9)
Lymphoma	Here .	0.78 (0.75, 0.81)		0.9 (0.88, 0.92)
Melanoma		0.9 (0.88, 0.92)		0.98 (0.97, 0.99)
Prostate Cancer		0.89 (0.87, 0.9)		0.9 (0.9, 0.91)
Ophthalmology				
Glaucoma	Herei	0.68 (0.65, 0.72)		0.85 (0.83, 0.87)
Macular Degeneration		0.85 (0.83. 0.88)		0.94 (0.93. 0.95)
Psychiatry and Psychology				
Bipolar Disorder	144	0.73 (0.7.0.75)		0.89 (0.87, 0.91)
Major Depressive Disorder	-	0.88 (0.86, 0.80)	-	0.9 (0.88 0.01)
Schizophrenia		0.79 (0.65, 0.93)		0.97 (0.94 0.90)
Bulmonology		0.10 (0.00, 0.00)		0.07 (0.04, 0.00)
Asthene		0.70 (0.77, 0.70)		0.04 (0.02, 0.04)
Asuma		0.78 (0.77, 0.79)		0.84 (0.83, 0.84)
Chronic Obstructive Pulmonary Disease		0.87 (0.85, 0.89)		0.89 (0.88, 0.9)
Pneumonia	-	0.85 (0.84, 0.86)	-	0.88 (0.87, 0.89)
Rheumatology				
Lupus	101	0.78 (0.75, 0.81)		0.9 (0.89, 0.92)
Urology				
Urinary Incontinence		0.66 (0.61, 0.71)		0.86 (0.85, 0.88)
	0.5 0.6 0.7 0.8 0.9 1		0.5 0.6 0.7 0.8 0.9 1	

**Supplemental Figure A1** Comparison of the Performance of Single-task and Multi-task Models Across Important Healthcare Tasks.

### 413 A.2 Architecture Details

Modality-Specific Network	Binary Classification	1	Multiclass Classificatio	ı	Regression		Cluster Encoder	
Dense Layer	Dense Layer (4845)	ŧ	Dense Layer (4845)	t	Dense Layer (4845)	<b>†</b>	Dense Layer (4845)	1 B
(image = 2084) (language = 2304)	Activation Layer (ReLU)		Activation Layer (ReLU)		Activation Layer (ReLU)		Activation Layer (ReLU)	-Attent
Automation Lawrer (Ded 1.0	Dense Layer (1024)		Dense Layer (1024)		Dense Layer (1024)		Dense Layer (256)	¥ S
Dense Lever (1024)	Activation Layer (ReLU)	2	Activation Layer (ReLU)		Activation Layer (ReLU)		Activation Layer (ReLU)	Po
Antonine Layer (2024)	Dense Layer (2048)	e-Atte	Dense Layer (2048)		Dense Layer (2048)		Dense Layer (128)	st-Atte
Dennas Laws (2010)	Activation Layer (ReLU)	ntion	Activation Layer (ReLU)	Pre-At	Activation Layer (ReLU)	Pre-AL	Activation Layer (ReLU)	ntion
Dense Layer (2048)	Dense Layer (1024)		Dense Layer (1024)	tentior	Dense Layer (1024)	tentior	Dense Layer (32)	ł
Dense Lever (1024)	Activation Layer (ReLU)		Activation Layer (ReLU)		Activation Layer (ReLU)			
Antonina Laws (2024)	Dense Layer (256)	t	Dense Layer (512)		Dense Layer (512)			
Desce Leves (200)	Activation Layer (ReLU)		Activation Layer (ReLU)		Activation Layer (ReLU)		Cluster Decoder	
Dense Layer (250)	Dense Layer (64)		Dense Layer (256)	ł	Dense Layer (256)	+ I	Dense Layer (32)	
	Activation Layer (ReLU)	Post-	Activation Layer (ReLU)	Î	Activation Layer (ReLU)	1	Activation Layer (ReLU)	
	Dense Layer (32)	Attentic	Dense Layer (64)		Dense Layer (128)	Pos	Dense Layer (128)	
	Activation Layer (ReLU)	×	Activation Layer (ReLU)	Post-	Activation Layer (ReLU)	st-Atter	Activation Layer (ReLU)	
	Dense Layer (1)		Dense Layer (32)	Attenti	Dense Layer (32)	ntion	Dense Layer (256)	
	Activation Layer (Sigmoid)	ţ	Activation Layer (ReLU)	0n	Activation Layer (ReLU)		Activation Layer (Sigmoid)	
			Dense Layer (8)		Dense Layer (1)	ŧ	Dense Layer (4845)	
			Activation Layer (Log Softmax)	+				

Supplemental Figure A2-1. Modality-specific and task-specific network architectures.



Supplemental Figure A2-2. Overall Pipeline of the M3H Architecture.

### 414 A.3 Step-by-Step Data Integration and Modeling Procedure

### Algorithm A3 End-to-End Data Integration and Modeling Pipeline

**Input**: Tabular data  $X^{tabular}$ , Time-series data  $X^{time-series}$ , Image data  $X^{vision}$ , Language data  $X^{language}$ , Feature extractor  $f_i$  of modality i, Outcome vector  $y_k$  for task  $k \in \mathcal{K}$ ,  $\hat{k} \in \hat{\mathcal{K}}$  indicates a set of task combinations.  $\mathcal{L}_{\hat{k}}$  as the aggregated loss function of each task combination  $\hat{k}$ .  $p \in \mathcal{P}$  is the set of hyperparameter combinations.  $\epsilon = 10^{-6}$  to avoid numerical precision error during computation, **Output**: Trained model and evaluation scores

### Step 0 – Data pre-processing and cleaning

• Impute missing values for all modalities, where here x is a generic data entry:

 $x = \begin{cases} 0 & \text{if } x \text{ is numerical or image data} \\ "" (empty string) & \text{if } x \text{ is text data} \end{cases}$ 

• Rescale image size:

$$X^{vision} \leftarrow \text{resize}(X^{vision}, 224 \times 224)$$

Step 1 – Embedding generation of each modality, an example of difference sources with the same modality is EKG notes vs. radiology notes:

$$E_j^i = f_i(X_j^i) \quad \forall i \in \{tabular, time-series, vision, language\}, \\ \forall j \in \{different \ sources \ in \ each \ modality\}$$

### Step 2 - Concatenate embeddings of all sources of the same modality into a single flattened vector:

 $E^{i} = \operatorname{vec}(E_{1}^{i}, E_{2}^{i}, \cdots, E_{n}^{i}) \quad \forall i \in \{tabular, time-series, vision, language\}$ 

Step 3 - Data Normalization

$$E^{i} = rac{E^{i} - \operatorname{mean}(E^{i})}{\operatorname{STD}(E^{i}) + \epsilon} \quad \forall i \in \{tabular, time-series, vision, language\}$$

Step 4 - Structure input data with outcomes for a task combination

$$E_{\hat{k}} = \operatorname{vec}(E^{tabular}, E^{time-series}, E^{vision}, E^{language}, y_1, y_2, \cdots, y_k)$$

### Step 5 - Model Training, Validation and Evaluation

For task combination  $\hat{k}$  in the set of all prediction tasks  $\hat{\mathcal{K}}$ :

• Split data into train and test datasets with fixed seed.

$$E_{\hat{k}}^{train}, E_{\hat{k}}^{test}, y_{\hat{k}}^{train}, y_{\hat{k}}^{test} \leftarrow train\_test\_split(E_{\hat{k}})$$

Perform 5-fold cross-validation with grid search to select the best parameter combination p<sup>\*</sup> ∈ P on the training data that has the best cumulative performance across all tasks inside the task combination k̂.

$$M3H_{\hat{k}}^{*} \leftarrow \underset{M3H_{p} \forall p \in \mathcal{P}}{\operatorname{argmin}} \mathcal{L}_{\hat{k}}(M3H_{p}(E_{\hat{k}}^{train}), y_{\hat{k}}^{train})$$

• Evaluate the optimal M3H model on the test set data:

test\_set\_score = 
$$M3H_{\hat{k}}^*(E_{\hat{k}}^{test}, y_{\hat{k}}^{test})$$

For each potential number of task combinations (i.e., single task = 1, 3-combined multitask = 3) and each task k, report the best model performance for each task.

### 415 A.4 Explainability of Input Space: by SHAP

We demonstrate below that by using SHAP values, we can effectively understand the magnitude 416 and directionality of each input clinical variable's contribution to the outcome prediction and thus 417 provide actionable insights for the physicians. Specifically, we analyze an M3H-framework- trained 418 multi-task model between diabetes and heart failure and study the effects of tabular features on 419 diabetes outcomes. We sampled 100 patients and studied their mean-standard deviation normalized 420 tabular features using two types of analysis: feature importance and feature interaction. We observe 421 that patients with lower age are less likely to have diabetes (blue dots for age have mostly negative 422 SHAP values). 423



**Supplemental Figure A4-1.** SHAP feature importance plot: each dot indicates a single sample among the 100 test set samples. Higher values of the feature are indicated in red, and lower values in blue. The most important feature is ranked at the top, followed by other features. A higher SHAP value (right-hand side of the axis) indicates a higher likelihood of a positive outcome (has diabetes), and a lower SHAP value indicates a negative outcome (does not have diabetes).



**Supplemental Figure A4-2.** The SHAP interaction plot demonstrates the nonlinear interactions between features on the outcome prediction captured by the M3H model. Age impacts the risk of diabetes differently depending on the patient's gender.

### 424 A.5 Characteristics of HAIM-MIMIC-MM

### 425 A.5.1 Limitations:

HAIM-MIMIC-IV was developed from the MIMIC-IV database, and several inherent biases and
limitations should be addressed. The cohort is collected from a single-care hospital in Boston and
focuses on intensive-care unit patients. This could potentially restrict the demographics and clinical
conditions of the patients to this specific geographical location and hospital setting. We also note that
MIMIC-IV has recording errors, missing values, and other inconsistencies that are universal to all
medical datasets and could pose a challenge for model development.

### 432 A.5.2 Embedding Dimensionality and Corresponding Clinical Variables:

The embeddings used as input data for M3H come from the multimodal database HAIM-MIMIC-433 MM, where the dimensionality of the features is explained and summarized in the paper's original 434 supplemental tables 1 and 2, which are included below for reference. The size of time-series 435 embedding is computed as the number of raw features multiplied by 11 unique features extracted: 436 maximum, minimum, mean, variance, average piece-wise change over time, average absolute piece-437 wise change over time, maximum absolute piece-wise change over time, sum of absolute piece-wise 438 change over time, change from end-beginning magnitude, number of peaks, and slope of the original 439 time series sequence. There are three categories: chart event (9  $\times$  11 = 99 features), lab event (22  $\times$ 440 11 = 242 features), and procedure event ( $10 \times 11 = 110$  features). The size of note embedding comes 441 from the output shape of the pre-trained model ClinicalBERT, which is 768. Similarly, the size of 442 vision embeddings comes from the output shape of the pre-trained model Densenet121-res224-chex, 443 which is 1024 (the dimension of the second to last layer of the model), and 18 (the output/last layer 444 dimension). 445

### 446 A.5.3 Missing Data:

We also include here a table of the missing value distribution of the HAIM-MIMIC-MM dataset reported in the original paper (originally Supplemental Table 3) and how it was handled in that integration procedure.

#	Chart events	Laboratory events	Procedure events
1	Heart rate	Glucose	Foley Catheter
2	Non-invasive systolic blood pressure	Potassium	PICC Line
3	Non-invasive blood diastolic pressure	Sodium	Intubation
4	Non-invasive nominal blood pressure	Chloride	Peritoneal dialysis
5	Respiratory rate	Creatinine	Bronchoscopy
6	$O_2$ saturation by pulse oximetry	Urea nitrogen	EEG
7	Verbal GCS response	Bicarbonate	Dialysis CRRT
8	Eye opening GCS response	Anion gap	Dialysis catheter
9	Motor GCS response	Hemoglobin	Chest tube removed
10		Hematocrit	Hemodialysis
11		Magnesium	
12		Platelet count	
13		Phosphate	
14		White Blood Cells	
15		Total calcium	
16		MCH	
17		Red Blood Cells	
18		MCHC	
19		MCV	
20		RDW	
21		Platelet count	
22		Neutrophils	
23		Vancomycin	

Supplemental Table A5-1. Patient signals in MIMIC-IV-MM by type of event used as time-series
 for embedding extraction. Nine time-dependent signals were derived from procedures, twenty-three
 were derived from laboratories, and eight were derived from information included in the patient chart.
 CRRT=Continuous renal replacement therapy, EEG=Electroencephalogram, GCS=Glasgow Coma
 Scale, MCH=Mean corpuscular hemoglobin, MCHC=Mean corpuscular hemoglobin concentration,
 PICC=Peripherally inserted central catheter, RDW=Red blood cell distribution width.

456

#	Data Modalities		# Data Sources
1	Tabular	1	Demographics $(E_{de})$
2	Time-series	2	Chart events $(E_{ce})$
		3	Laboratory events $(E_{le})$
		4	Procedure events $(E_{pe})$
3	Text	5	Radiological notes $(E_{radn})$
		6	Electrocardiogram notes $(E_{ecgn})$
		7	Echocardiogram notes $(E_{econ})$
4	Images	8	Visual probabilities $(E_{yp})$
	•	9	Visual dense-layer feature $(E_{vd})$
		10	Aggregated visual probabilities $(E_{ymp})$
		11	Aggregated visual dense-layer features $(E_{\rm vmd})$

**Supplemental Table A5-2.** List of different data modalities and data sources used to test the HAIM framework based on the MIMIC-IV-MM database. There are a total of four data modalities and eleven data sources. All data sources correspond to only one data modality. Thus, a model trained on a single data modality can have as little as 1 data source and many as 4 different data sources (of the same kind) as inputs. Double, triple and quadruple modality models can have a number of data sources ranging from [2 to 7], [3 to 9] and [4 to 11], respectively.

Feature Name	Missing %	Source	Handling
anchor_age	0.0	Demographics	N/A
gender_int	0.0	Demographics	N/A
ethnicity_int	0.0	Demographics	N/A
marital_status_int	0.0	Demographics	N/A
language_int	0.0	Demographics	N/A
insurance_int	0.0	Demographics	N/A
Foley Catheter	82.6	Procedure	Fill with 0
PICC Line	63.7	Procedure	Fill with 0
Intubation	75.3	Procedure	Fill with 0
Peritoneal Dialysis	99.7	Procedure	Fill with 0
Bronchoscopy	81.5	Procedure	Fill with 0
EEG	91.5	Procedure	Fill with 0
Dialysis - CRRT	93.1	Procedure	Fill with 0
Dialysis Catheter	88.9	Procedure	Fill with 0
Chest Tube Removed	93.1	Procedure	Fill with 0
Hemodialysis	92.9	Procedure	Fill with 0
Glucose	4.4	Lab	Fill with 0
Sodium	4.7	Lab	Fill with 0
Potassium	4.7	Lab	Fill with 0
Chloride	4.7	Lab	Fill with 0
Creatinine	4.7	Lab	Fill with 0
Urea Nitrogen	4.7	Lab	Fill with 0
Bicarbonate	4.7	Lab	Fill with 0
Anion Gap	4.7	Lab	Fill with 0
Hemoglobin	4.7	Lab	Fill with 0
Hematocrit	4.8	Lab	Fill with 0
Magnesium	5.4	Lab	Fill with 0
Platelet Count	9.8	Lab	Fill with 0

Feature Name	Missing %	Source	Handling
Phosphate	6.0	Lab	Fill with 0
White Blood Cells	4.9	Lab	Fill with 0
Calcium, Total	6.0	Lab	Fill with 0
MCH	4.9	Lab	Fill with 0
Red Blood Cells	4.9	Lab	Fill with 0
MCHC	4.9	Lab	Fill with 0
MCV	4.9	Lab	Fill with 0
RDW	4.9	Lab	Fill with 0
Neutrophils	36.9	Lab	Fill with 0
Vancomycin	60.0	Lab	Fill with 0
Heart Rate	19.5	Chart	Fill with 0
Non-Invasive Blood Pressure systolic	23.4	Chart	Fill with 0
Non-Invasive Blood Pressure diastolic	23.4	Chart	Fill with 0
Non-Invasive Blood Pressure mean	23.3	Chart	Fill with 0
Respiratory Rate	19.5	Chart	Fill with 0
$O_2$ saturation pulse oximetry	19.6	Chart	Fill with 0
GCS - Verbal Response	20.8	Chart	Fill with 0
GCS - Eye Opening	20.7	Chart	Fill with 0
GCS - Motor Response	20.8	Chart	Fill with 0
Electrocardiogram Notes	11.2	Notes	Empty String
Echocardiogram Notes	30.5	Notes	Empty String
Radiology Notes	0.1	Notes	Empty String

**Supplemental Table A5-3.** List of missing data percentages by individual variables and handling strategy. Individual variables (i.e., feature name) within key MIMIC-IV-MM data source groups are shown. The strategy for missing value handling used in our tests is as follows: 1) We exclude patients with no available X-rays from our selection cohort; 2) Time-series features are imputed with 0 if there is no measurement at any timestamp; 3) Text embeddings are generated from an empty string if there is no note available; 4) There were no missing values for demographics data.

#### A.6 Multitask Comparison 457

We implemented three methods using a universal problem setting of N tasks with feature dimension of d, with input features  $X = \{x_j\}_{j=1}^N$  where  $x_j \in \mathbb{R}^d$ . We do not include reshaping operations or the batch size dimension in the description to capture only the mathematical essence of the 458 459 460 implementations. 461

#### **Multi-head attention:** 462

463	Initialize linear transformation matrices:
464	– $W^Q, W^K, W^V \in \mathbb{R}^{d \times d}$ as query, key, value transformations
465	– $W^O \in \mathbb{R}^{d \times d}$ as the output transformation
466	- $H = 4$ as the number of heads
467	- $d_H = d/H$ as the dimension per head
468	• Apply linear transformation and projection on input features:
469	$- Q = XW^Q, K = XW^K, V = XW^V$
470	• Scaled dot-product to obtain attention weight ( $\sqrt{d_h}$ is used to stabilize gradient):
471	- $A = \operatorname{softmax}\left(\frac{QK^{\top}}{\sqrt{d_h}}\right)$
472	• Apply attention weights to obtain output:
473	$- O = (AV)W^O$
474	Cross-stitch:
475	• Initialize task interaction matrix:
476	- $\{T_{ij}\}_{i \neq j}^{1:N}$ where $T_{ij} \in \mathbb{R}^{2 \times 2}$
477	Apply interaction matrix:
478	- $z_{ij} = T_{ij} \cdot [x_i, x_j]  \forall (i, j)$
479	• Aggregation:
480	$- z_i = \sum_{j=1}^N z_{ij}  \forall i = 1, \dots, N$
481	• Output learned features:
482	$- \{z_i\}_{i=1}^N  \forall i = 1, \dots, N$
483	For <i>n</i> tasks, this requires $\frac{n(n-1)}{2}$ weight matrices of size $2 \times 2$ .
484	Multilinear relationship network (MRN):
195	Initialize linear transformation matrices:
486	- $\{T_{ij}\}_{i=1}^{N}$ where $T_{ij} \in \mathbb{R}^{d \times d}$
487	<ul> <li>Apply linear transformation and projection on input features:</li> </ul>
488	$- z_{ii} = T_{ii} \cdot [x_i, x_i]  \forall (i, j)$
489	• Aggregation:
100	$- \gamma - \sum^{N} \gamma \cdots \forall i - 1 $ N
450	• Output learned features:
491	$= \int z_i \sqrt{N}  \forall i = 1 \qquad N$
492	$- \left\{ z_{i} \right\}_{i=1}^{n}  \forall i = 1, \dots, i \forall$
493	For <i>n</i> tasks, this requires $\frac{\pi}{2}$ weight matrices of size $d \times d$ .
494	
495	Specifically, we conduct experiments in the original dataset on 10 different combinations of

of multi-4 tasks that comprehensively evaluate multitask strategies across all four types of machine learning 496 problem classes. The choice of diabetes and heart failure is arbitrary. 497

• Length of stay (regression), patient phenotyping (clustering) 498

- Length of stay (regression), thoracic testing (multiclass classification)
- Thoracic testing (multiclass classification), patient phenotyping (clustering)
- Diabetes (binary classification), length of stay (regression)
- Diabetes (binary classification), patient phenotyping (clustering)
- Diabetes (binary classification), thoracic testing (multiclass classification)
- Heart failure (binary classification), length of stay (regression)
- Heart failure (binary classification), patient phenotyping (clustering)
- Heart failure (binary classification), thoracic testing (multiclass classification)
- Diabetes (binary classification), Heart failure (binary classification)

<sup>508</sup> We observe that cross-task attention has a clear advantage in the majority of the cases across all three

509 strategies, with cross-stitch being a close competitor in these 2-tasks experiments (but with qualitative 510 disadvantages discussed below).

Machine Learning	Cross-task	Multi-Head	Multilinear	Cross
Problem Class	(M3H)	Attention	<b>Relationship Network</b>	Stitch
Regression	0.567	0.562 (-0.88%)	0.431 (-23.99%)	0.565 (-0.35%)
Clustering	0.405	0.521 (+28.64%)	0.176 (-56.54%)	0.487 (+20.25%)
Multiclass	0.755	0.715 (-5.30%)	0.595 (-21.19%)	0.755 (+0%)
Binary (diabetes)	0.873	0.824 (-5.61%)	0.873 (+0%)	0.869 (-0.46%)
Binary (heart failure)	0.881	0.864 (-1.93%)	0.896 (+1.7%)	0.888 (+0.79%)

**Supplemental Table A6.** Comparison of machine learning problem classes across different models. The values represent performance metrics and percentage differences from the baseline (Cross-task M3H).

Beyond the quantitative advantage of the proposed cross-task framework, we would also like to emphasize the qualitative advantage of the chosen framework over existing methods:

- Interpretability: Available multimodal multi-task foundation models heavily rely on com-513 plex architectures, for example, with repeated use of multi-head attention mechanisms tens or 514 hundreds of times to achieve good performance guarantees. Even with known visualization 515 efforts to interpret these architectures, in practice, these attention weights are almost very 516 often not interpretable and non-sensible. This is why we opted for such a model structure 517 design. As reviewer 2 later correctly pointed out, the existing style of complex architecture 518 makes it very difficult to obtain clinician trust in hospital settings precisely because of such 519 lack of interpretability. Instead, in our case, we apply a single cross-task attention with one 520 single channel and a clean 2D attention weight to explicitly model how self-attention and 521 cross-attention interact. Such design allows for future analysis of interpretability a lot more 522 easily. 523
- Lightweight design for deployment: Existing architectures, such as Google's Med-PaLM 524 2 (released March 2023), contain 540 billion parameters and can be estimated usually 525 to need months to train with commercial-grade GPUs (such as Nvidia Volta V100) with 526 heavy RAM memory requirements (at least 1000GB if not parallelized). Although lighter-527 weight versions of these models exist, they remain in the billion-level parameters and pose a 528 significant implementation challenge for hospitals if they wish to host in-house models for 529 data privacy reasons. M3H, on the other hand, can be offered as a much lighter solution to 530 avoid these issues. 531

Similarly, all three of the compared multiclass methods require significantly more complex network
 structures. Multi-head model (in our case with 4 heads) requires 4 additional channels to integrate
 the data from separate heads; cross-stitch models would require significantly more weight matrices as
 the number of co-learned tasks increases, MRN models will require even more parameters, as they
 require a linear transformation of each combination of task pairs.

#### A.7 Loss Function Definition 537

- **Binary cross entropy loss (Binary classification):** 538
- Given  $x \in \mathbb{R}^d$  as an input feature of dimension  $d, y \in \{0,1\}^d$  as the binary outcomes,  $\hat{y} = \sigma(w^T x + b)$  is the predicted outcome from the M3H framework. Here  $\sigma(z) = \frac{1}{1+e^{-z}}$  is the 539 540 sigmoid function, w is the weight matrix, b is the bias vector, the loss function is defined as: 541  $l_{\text{binary}}(y, \hat{y}) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})).$ 542

#### Negative log-likelihood loss (Multiclass classification): 543

Given  $x \in \mathbb{R}^d$  as an input feature of dimension  $d, y \in \{1, 2, ..., K\}^d$  as the multiclass outcomes from K classes,  $\hat{y} = \sigma(w^T x + b)$  is the predicted outcome from the M3H framework. Here:  $\sigma(z) = z - \log\left(\sum_{k=1}^{K} e^{z_k}\right)$  is the log-softmax function, w is the weight matrix, b is the bias vector, 544 545 546 the loss function is defined as:  $l_{\text{multiclass}}(y, \hat{y}) = -\log(\hat{y})$ . 547

#### Mean absolute error (Regression): 548

Given  $x \in \mathbb{R}^d$  as an input feature of dimension  $d, y \in \mathbb{R}^d$  as the regression outcomes,  $\hat{y} = w^T x + b$  is the predicted outcome from the M3H framework. Here w is the weight matrix, b is the bias vector, the loss function is defined as:  $l_{\text{regression}}(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^{N} |\hat{y}_i - y_i|$ . 549 550 551

#### Mean squared error (Clustering): 552

Given  $x \in \mathbb{R}^d$  as an encoder input of dimension  $d, \hat{x} \in \mathbb{R}^d$  as the decoder output of the same dimension, here w is the weight matrix, b is the bias vector, the loss function is defined as:  $l_{\text{clustering}}(x, \hat{x}) = \frac{1}{d} \sum_{i=1}^{d} (\hat{x}_i - x_i)^2$ . 553 554 555

#### **Contrastive Learning:** 556

This learning aims to project embeddings of different modalities into the same embedding space by 557 contrasting positive pairs (modalities from the same samples) and negative pairs (modalities from 558 dissimilar samples). In the M3H framework, because of the small dimension of the tabular features 559 (6) in comparison to the rest of the three modalities, we only apply contrastive learning among time 560 series, vision, and language data inputs. The formulation is as follows: 561

Given  $\hat{N}$  as the number of permutations between the N samples' three modalities (or N choose 562 2), and given  $E_i$  and  $E_j$  as pairs of embedding vectors from different modalities,  $y_i$  as the la-563 bel for the pair of (i, j), where they are either from the same sample (1) or different samples (0). We define a positive margin p = 0, and a negative margin n = 1. Specifically, for positive pairs, the loss is only computed if  $|E_i - E_j| > p$ , which aims to decrease positive 564 565 566 pairs' distance to 0, and for negative pairs, we only compute the loss when  $|E_i - E_j| < n$ , which aims to push the distance to be close to 1. The contrastive loss is computed as follows: 567 568

569 
$$L = \frac{1}{N} \sum_{i=1}^{N} \left( y_i \max\left(0, |E_i - E_j| - p\right)^2 + (1 - y_i) \max\left(0, n - |E_i - E_j|\right)^2 \right)$$

#### **NeurIPS Paper Checklist** 570

#### 1. Claims 571

Question: Do the main claims made in the abstract and introduction accurately reflect the 572 paper's contributions and scope? 573

Answer: [Yes] 574

Justification: We have provided detailed outline of how experiments are conducted and the 575 improvements of our model in comparison to the nominal single-task models both in the 576 introduction and abstract. We have also highlighted key findings and technical novelties 577 introduced in the later sections as well. 578

#### Guidelines: 579

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
  - · The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
  - · The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- 586 587 588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

580

581

582

583

584

585

# • It is fine to include aspirational goals as motivation as long as it is clear that these goals

are not attained by the paper.

### 2. Limitations

- Question: Does the paper discuss the limitations of the work performed by the authors?
- Answer: Yes

Justification: We have a limitation section at the end of the paper detailing the several possibilities for limitations, ranging from data, to inclusion of other tasks. We have also provided a practical implication section to reflect rigorously how the framework should be adopted in new data and system settings, and what are the potential remedies to deal of potential challenges.

- Guidelines:
  - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
  - The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
  - · The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
  - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by 618 reviewers as grounds for rejection, a worse outcome might be that reviewers discover 619 limitations that aren't acknowledged in the paper. The authors should use their best 620 judgment and recognize that individual actions in favor of transparency play an impor-621 tant role in developing norms that preserve the integrity of the community. Reviewers 622 will be specifically instructed to not penalize honesty concerning limitations. 623

624	3.	Theory Assumptions and Proofs
625 626		Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?
627		Answer: [NA].
628		Justification: Our paper is focused on model architecture design and thus does not have
629		theoretical results.
630		Guidelines:
631		• The answer NA means that the paper does not include theoretical results.
632 633		• All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
634		• All assumptions should be clearly stated or referenced in the statement of any theorems.
635		• The proofs can either appear in the main paper or the supplemental material, but if
636		they appear in the supplemental material, the authors are encouraged to provide a short
637		proof sketch to provide intuition.
638 639		• Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
640		• Theorems and Lemmas that the proof relies upon should be properly referenced.
641	4.	Experimental Result Reproducibility
642		Question: Does the paper fully disclose all the information needed to reproduce the main ex-
643		perimental results of the paper to the extent that it affects the main claims and/or conclusions
644		of the paper (regardless of whether the code and data are provided or not)?
645		Answer: [Yes]
646		Justification: The paper clearly outlined the architectures, parameters, data cohort availability,
647 648		processing steps to ensure that all necessary data is needed for the reproduction of the computational results.
649		Guidelines:
650		• The answer NA means that the paper does not include experiments.
651		• If the paper includes experiments, a No answer to this question will not be perceived
652		well by the reviewers: Making the paper reproducible is important, regardless of
653		whether the code and data are provided or not.
654 655		• If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable
656		• Depending on the contribution reproducibility can be accomplished in various ways
657		For example, if the contribution is a novel architecture, describing the architecture fully
658		might suffice, or if the contribution is a specific model and empirical evaluation, it may
659		be necessary to either make it possible for others to replicate the model with the same
660		dataset, or provide access to the model. In general. releasing code and data is often
661		one good way to accomplish this, but reproducibility can also be provided via detailed
662		instructions for now to replicate the results, access to a nosted model (e.g., in the case
663 664		appropriate to the research performed
665		• While NeurIPS does not require releasing code, the conference does require all submis-
666		sions to provide some reasonable avenue for reproducibility, which may depend on the
667		nature of the contribution. For example
668		(a) If the contribution is primarily a new algorithm, the paper should make it clear how
669		to reproduce that algorithm.
670		(b) If the contribution is primarily a new model architecture, the paper should describe
671		the architecture clearly and fully.
672		(c) If the contribution is a new model (e.g., a large language model), then there should
673		the model (e.g. with an open-source dataset or instructions for how to construct
675		the dataset).

676 677 678 679 680			(d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.
681		5.	Open access to data and code
682 683 684 685			Question: Does the paper provide open access to the data and code, with sufficient instruc- tions to faithfully reproduce the main experimental results, as described in supplemental material?
606			Institution: Unfortuantely due to privacy reasons we are unable to release the code. But
687			readers are encouraged to contact the authors to its access.
688			Guidelines:
689 690 691 692 693 694 695	•		<ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).</li> </ul>
696 697 698	i.		• The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
699 700 701			<ul> <li>The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li> <li>The authors should provide scripts to reproduce all experimental results for the new</li> </ul>
702 703			proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
704 705			• At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
706 707			• Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
708		6.	Experimental Setting/Details
709 710 711			Question: Does the paper specify all the training and test details (e.g., data splits, hyper- parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?
712			Answer: [Yes]
713 714 715			Justification: The paper outlined all the necessary details for where to obtain the data cohort, how to preprocess the dataset, how to split the train, validation, and test sets, the hyperparameters chosen in the paper in order to reproduce all the results.
716			Guidelines:
717			• The answer NA means that the paper does not include experiments.
718			• The experimental setting should be presented in the core of the paper to a level of detail
719			that is necessary to appreciate the results and make sense of them.
720 721			• The full details can be provided either with the code, in appendix, or as supplemental material.
722		7	Experiment Statistical Significance
723			Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?
725			Answer: [Ves]
726			Justification: Error bars and confidence intervals are provided for the main computational
727 728			results in the supplemental figure. Details on how these results are computed are also included in the methods.

729		Guidelines:
730		• The answer NA means that the paper does not include experiments.
731		• The authors should answer "Yes" if the results are accompanied by error bars, confi-
732		dence intervals, or statistical significance tests, at least for the experiments that support
733		the main claims of the paper.
734		• The factors of variability that the error bars are capturing should be clearly stated (for
735		example, train/test split, initialization, random drawing of some parameter, or overall
736		run with given experimental conditions).
737		• The method for calculating the error bars should be explained (closed form formula, call to a library function, bestetran, etc.)
738		• The assumptions made should be given (e.g. Normally distributed errors)
739		<ul> <li>The assumptions made should be given (e.g., Normany distributed errors).</li> <li>It should be clear whether the error her is the standard deviation or the standard error.</li> </ul>
740		of the mean
742		• It is OK to report 1-sigma error bars, but one should state it. The authors should
743		preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
744		of Normality of errors is not verified.
745		• For asymmetric distributions, the authors should be careful not to show in tables or
746		figures symmetric error bars that would yield results that are out of range (e.g. negative
747		error rates).
748		• If error bars are reported in tables or plots, The authors should explain in the text how
749		they were calculated and reference the corresponding figures or tables in the text.
750	8.	Experiments Compute Resources
751		Question: For each experiment, does the paper provide sufficient information on the com-
752		puter resources (type of compute workers, memory, time of execution) needed to reproduce the amerimenta?
/53		
754		Answer: [Yes],
755		Justification: The paper's section on practical implementation details the computational
756		resources needed to run the model, as well as an estimated time to run on a local computer. This is also accompanied by the operating system details
758		Guidelines:
		• The answer NA means that the nemer does not include auroriments
759		• The answer NA means that the paper does not include experiments.
760 761		• The paper should indicate the type of compute workers CPO or GPO, internal cluster, or cloud provider, including relevant memory and storage.
762		• The paper should provide the amount of compute required for each of the individual
763		experimental runs as well as estimate the total compute.
764		• The paper should disclose whether the full research project required more compute
765		than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper)
766	0	Code Of Ethics
767	9.	
768		Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
770		Answer: [Yes],
771		Justification: The content of this paper complies with NeuRIPS code of ethics, and was
772		conducted with the hopes to advance our understanding of medicine to better patient care.
773		Guidelines:
774		• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
775		• If the authors answer No, they should explain the special circumstances that require a
776		deviation from the Code of Ethics.
777		• The authors should make sure to preserve anonymity (e.g., if there is a special consid-
778		eration due to laws or regulations in their jurisdiction).
779	10.	Broader Impacts

780 781		Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?
782		Answer: [Yes]
783		Justification: In the implication to practice section, we discuss thoroughly both the negative
784		and positive implications of the introduced model and its impact if applied to hospital
785		systems.
786		Guidelines:
707		• The answer NA means that there is no conjuted impact of the work performed
/8/		• The answer two means that there is no societal impact of the work performed.
788 789		• If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
790		• Examples of negative societal impacts include potential malicious or unintended uses
791		(e.g., disinformation, generating fake profiles, surveillance), fairness considerations
792		(e.g., deployment of technologies that could make decisions that unfairly impact specific
793		groups), privacy considerations, and security considerations.
794		• The conference expects that many papers will be foundational research and not tied
795		to particular applications, let alone deployments. However, if there is a direct path to
796		any negative applications, the authors should point it out. For example, it is legitimate
797		to point out that an improvement in the quality of generative models could be used to
798		generate deepfakes for disinformation. On the other hand, it is not needed to point out
799		that a generic algorithm for ontimizing neural networks could enable people to train
800		models that generate Deepfakes faster.
001		• The authors should consider possible harms that could arise when the technology is
801		being used as intended and functioning correctly harms that could arise when the
002 902		technology is being used as intended but gives incorrect results and harms following
003		from (intentional or unintentional) misuse of the technology
804		If the second of uninclutional finishese of the technology.
805		• If there are negative societal impacts, the authors could also discuss possible mitigation
806		strategies (e.g., gated release of models, providing defenses in addition to attacks,
807		mechanisms for monitoring misuse, mechanisms to monitor now a system learns from
808	11	Seference and accessionity of ML).
809	11.	Saleguards
810		Question: Does the paper describe safeguards that have been put in place for responsible
811		release of data or models that have a high risk for misuse (e.g., pretrained language models,
812		image generators, or scraped datasets)?
813		Answer: [Yes]
814		Justification: In the practical implementation section, we discuss how the model should
815		be validated and evaluated prior to its implementation. This includes safeguards against
816		for example data perturbations to ensure that the model does not negatively impact patient
817		outcome predictions.
818		Guidelines:
819		• The answer NA means that the paper poses no such risks.
820		• Released models that have a high risk for misuse or dual-use should be released with
821		necessary safeguards to allow for controlled use of the model. for example by requiring
822		that users adhere to usage guidelines or restrictions to access the model or implementing
823		safety filters.
824		• Datasets that have been scraped from the Internet could nose safety risks. The authors
825		should describe how they avoided releasing unsafe images
000		• We recognize that providing effective sofequends is shallonging, and many papers do
826		• we recognize that providing enective sateguards is chantenging, and many papers do not require this, but we appourage outpose to take this into account and make a bast
827		for the fort
628		
829	12.	Licenses for existing assets
830 831		Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and

832 properly respected?

833	Answer: [Yes]
834 835	Justification: The data and models used for this paper are properly introduced, elaborated, and cited by the paper.
836	Guidelines:
837	• The answer NA means that the paper does not use existing assets.
838	• The authors should cite the original paper that produced the code package or dataset.
830	• The authors should state which version of the asset is used and if possible include a
840	URL.
841	• The name of the license (e.g., CC-BY 4.0) should be included for each asset.
842	• For scraped data from a particular source (e.g., website), the convright and terms of
843	service of that source should be provided.
844	• If assets are released, the license, copyright information, and terms of use in the
845	package should be provided. For popular datasets, paperswithcode.com/datasets
846	has curated licenses for some datasets. Their licensing guide can help determine the
847	license of a dataset.
848 849	• For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
850	• If this information is not available online, the authors are encouraged to reach out to
851	the asset's creators.
852	13. New Assets
853	Question: Are new assets introduced in the paper well documented and is the documentation
854	provided alongside the assets?
855	Answer: [Yes].
856	Justification: The provided new model has been thoroughly outlined by the paper with
857	detailed instructions on its training parameters and architectures, data used, as well as
858	evaluations.
859	Guidelines:
860	• The answer NA means that the paper does not release new assets.
861	• Researchers should communicate the details of the dataset/code/model as part of their
862	submissions via structured templates. This includes details about training, license,
863	limitations, etc.
864 865	• The paper should discuss whether and how consent was obtained from people whose asset is used.
866 867	• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
868	14. Crowdsourcing and Research with Human Subjects
869	Question: For crowdsourcing experiments and research with human subjects, does the paper
870	include the full text of instructions given to participants and screenshots, if applicable, as
871	well as details about compensation (if any)?
872	Answer: [Yes].
873	Justification: This study is conducted using a licensed, but publicly available dataset of
874	human subjects. The details of this cohort has been thoroughly discussed in the patient
875	cohort section.
876	Guidelines:
877	• The answer NA means that the paper does not involve crowdsourcing nor research with
878	human subjects.
879	• Including this information in the supplemental material is fine, but if the main contribu-
880	tion of the paper involves numan subjects, then as much detail as possible should be included in the main paper.
881	According to the NeurIDS Code of Ethics, workers involved in data collection, exercice
883	• According to the Neuris Code of Eulics, workers involved in data conection, curation, or other labor should be paid at least the minimum wage in the country of the data
	collector

885 886	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects
887 888 889 890		Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?
891		Answer: [NA].
892		Justification: IRB approval was not needed due to the use of a public national database.
893		Guidelines:
894 895		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
896 897 898		• Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
899 900 901		• We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
902 903		• For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.