

---

# Promoting cross-modal representations to improve multimodal foundation models for physiological signals

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many healthcare applications are inherently multimodal and involve multiple types  
2 of physiological signals. As sensors for measuring these signals become more  
3 ubiquitous, it is increasingly important to improve machine learning methods  
4 that consume multimodal healthcare data. Pretraining foundation models is a  
5 promising avenue for success. However, methods for developing foundation models  
6 in healthcare are still early in exploration and it is unclear which pretraining  
7 strategies are most effective given the diverse set of physiological signals collected.  
8 This is in part due to challenges of multimodal learning with health data: data  
9 across many patients is difficult to obtain and expensive, and there is a lot of inter-  
10 subject variability. Furthermore, modalities are often heterogeneously informative  
11 across the downstream tasks of interest. Here, we explore these challenges in  
12 the PhysioNet 2018 Challenge dataset collected across 1,985 patients. We used a  
13 masked autoencoding objective to pretrain a multimodal model on the dataset. We  
14 show that the model learns representations that can be linearly probed for a diverse  
15 set of downstream tasks. We hypothesize that cross-modal reconstruction objectives  
16 are important for the success of multimodal training as they encourages the model  
17 to combine information across modalities. We demonstrate that adding modality  
18 drop in the input space improves model performance across downstream tasks. We  
19 also show that late-fusion models pretrained with contrastive learning objectives  
20 are not as effective as across multiple tasks. Finally, we analyze the representations  
21 developed in the model. We show how attention weights become more cross-modal  
22 and temporally aligned as a result of our chosen pretraining strategy. The learned  
23 embeddings also become more distributed in terms of the modalities that each  
24 unit in the model encodes. Taken together, our work demonstrates the utility of  
25 multimodal foundation models with health data, even across diverse physiological  
26 data sources. We further argue how more explicit means of inducing cross-modality  
27 may be valuable additions to any multimodal pretraining strategy.

## 28 1 Introduction

29 Healthcare applications often involve integrating information across many modalities. For instance,  
30 to diagnose sleep disorders, physicians may evaluate neural, muscular, and respiratory signals [Ibáñez  
31 et al., 2018]. Adding to the complexity, the data used in healthcare spans a wide variety of formats  
32 (imaging data, time series, etc) and are collected from sensors placed on many different body locations  
33 [Acosta et al., 2022]. Many of these sensors for health data are becoming increasingly prevalent  
34 in everyday wearable devices [Jeong et al., 2018] [Wu and Luo, 2019] [Iqbal et al., 2021]. This  
35 technological advance is a promising opportunity for personalized healthcare and improving patient  
36 care. Thus, it is more and more important to leverage artificial intelligence to aid the interpretation of  
37 health data with heterogenous sensors.

38 In many settings, artificial intelligence has achieved unprecedented success in the development of  
39 multimodal foundation models [Jin et al., 2024, Bordes et al., 2024, Wadekar et al., 2024]. For  
40 instance, models can now integrate information across language, vision, audio, and video to solve  
41 complex tasks and perform human-like feats of reasoning [Radford et al., 2021, Alayrac et al., 2022,  
42 Wu et al., 2023, Lu et al., 2024, Mizrahi et al., 2024]. Multimodal foundation models are pretrained  
43 in a self-supervised manner on vast amounts of data to link information across modalities. The  
44 representations developed by these models are useful for tasks that require multimodal understanding.  
45 After pretraining, these models may be further trained on a downstream task or the representations  
46 they produce can be used as is. Pretraining strategies often outperform models trained from scratch on  
47 the same tasks and require less labeled data [Jin et al., 2024]. The success of multimodal foundation  
48 models in other domains suggests that similar advances can be achieved in healthcare settings.

49 There are further reasons to believe that health data in particular can benefit from foundation model  
50 strategies. Annotated data is limited in health data because clinical expertise is often necessary to  
51 create labels. Thus, the label efficiency of pretrained models is very useful in this setting. Furthermore,  
52 when considering wearable health devices, it becomes more important to develop models that are  
53 size-efficient. If a model pretrained on health data can successfully transfer its representations across  
54 many downstream tasks, this can greatly save on memory and runtime costs for wearable devices.

55 However, working with health data also introduces new types of challenges. Pretraining often  
56 consumes large amounts of unlabeled data, but patient privacy concerns limit the amount of large  
57 datasets available in this domain [Acosta et al., 2022, Shaik et al., 2023]. In addition, the cost  
58 associated with deploying many health sensors can make large-scale data collection prohibitively  
59 expensive [Acosta et al., 2022]. Thus, it becomes less clear whether pretraining can be as effective  
60 as it is in settings like natural language, where large corpora are more widely available. In health  
61 applications, it is also common for certain modalities to vary greatly in their informativeness for  
62 different downstream tasks [Krones et al., 2024]. This problem is exacerbated in wearable devices  
63 since different sensors may suffer from unequal amounts of noise, perhaps due to weaker contact or  
64 interference from other devices [Ates et al., 2022, Canali et al., 2022]. This poses a challenge for  
65 developing general purpose models that can be used for diverse tasks.

66 Here, we investigate these challenges by pretraining a multimodal model in a publicly available  
67 dataset with 1,985 patients. We are specifically concerned with time series data collected from  
68 physiological signals measured overnight from patients. Our contributions are the following:

- 69 • We explore the development of a multimodal foundation model in a dataset of diverse  
70 physiological signals: electroencephalography (EEG), electromyography (EMG), electrooculography (EOG), and electrocardiology (ECG). We demonstrate the strength of the  
71 learned representations in linear probe experiments on a disparate set of downstream tasks.  
72
- 73 • We show how explicitly enforcing cross-modal reconstruction in the pretraining objective  
74 improves the quality of the learned representations over standard multimodal MAE. We  
75 also show how late-fusion models pretrained with contrastive learning does not effectively  
76 transfer across multiple tasks.
- 77 • We analyze the learned representations to show that attention weights in the model be-  
78 come increasingly cross-modal under the pretraining objective we use. We also show that  
79 individual units in the model become more diversely tuned to the different modalities.

## 80 2 Related Work

81 Pretraining models with self-supervised objectives is a popular and effective strategy in machine  
82 learning [Ericsson et al., 2022, Gui et al., 2023]. After pretraining, the parameters of the model  
83 can be finetuned for some downstream task. Alternatively, another common approach is to freeze  
84 the pretrained model and train a lightweight readout head that uses the learned representations  
85 from the model to solve downstream tasks. This approach is especially attractive if efficiency in  
86 parameter tuning is a priority. In either case, a pretraining paradigm often outperforms training a  
87 model from scratch. Pretraining is especially useful if labeled data in the downstream task is limited  
88 as it provide a means for experimenters to define inductive biases on the model representations.  
89 Self-supervised strategies span several categories, including generative methods, contrastive learning,  
90 and autoencoding [Del Pup and Atzori, 2023, Gui et al., 2023]. We limit our discussion to the latter  
91 two in the context of multimodal pretraining.

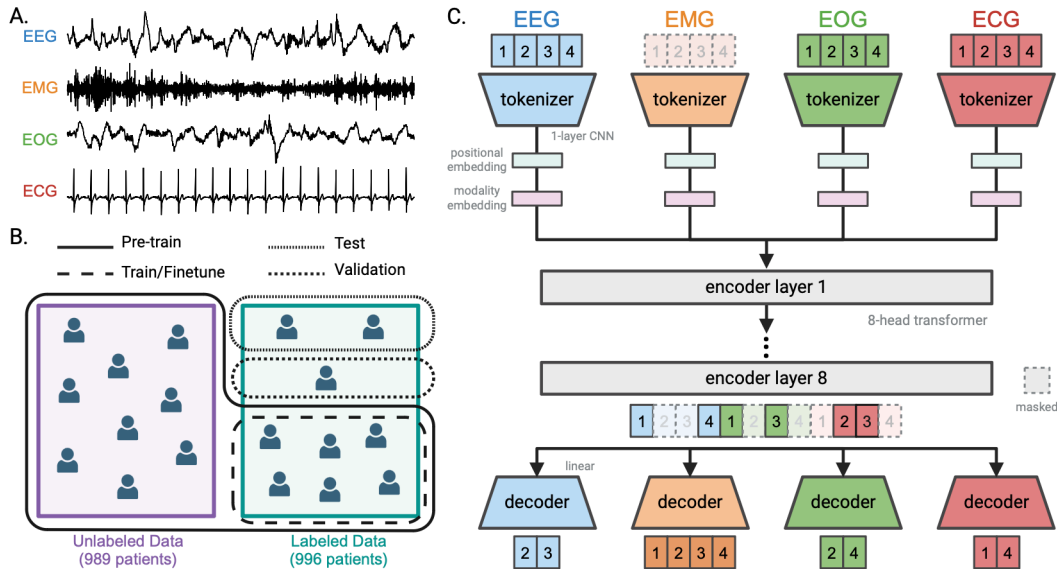


Figure 1: **A.** A 30-second sample from the training dataset. **B.** Data is split by patient identity for each part of the training procedure. The PhysioNet 2018 dataset consists of unlabeled data from 989 patients and labeled data from 996 patients, where each patient contributes 7.7 hours of data on average. The data for pretraining consists of all patients in the unlabeled dataset and 657 patients from the labeled dataset. The data for training and finetuning is drawn from the patients of the labeled dataset that were also used for pretraining. The data for the validation and test are drawn from the remaining patients of the labeled dataset not used for either pretraining or training. **C.** Diagram of the main pretraining strategy we use: multimodal masked autoencoding with modality drop in the input space. Tokenizers are modality-specific.

92 Contrastive learning is a self-supervised learning framework where models are optimized such that  
 93 representations of data in positive pairs become more similar while representations of data in negative  
 94 pairs become more dissimilar [Chen et al. 2020, Purushwalkam and Gupta 2020]. The definition of  
 95 positive and negative pairs is crucial to the success of these methods. One way to define these  
 96 pairs is to construct multiple “views” of a data sample through augmentations, like in the SimCLR  
 97 algorithm [Chen et al. 2020, Yuan et al. 2021]. Thus, a positive pair of data may be two different  
 98 augmentations of a single data sample (negative pairs would then be constructed across different  
 99 data samples). When working with multimodal data, another option is to consider each modality  
 100 as a distinct view of a data sample. In this case, positive pairs can be constructed by comparing  
 101 representations across modalities, as in CLIP-style pretraining [Radford et al. 2021, Yuan et al. 2021,  
 102 Zhang et al. 2022].

103 Masked autoencoding (MAE) is another popular pretraining strategy. In MAE, random patches of  
 104 the input are masked, and the model must use the remaining portions of the input to reconstruct  
 105 the masked portion [He et al. 2022]. This method has been extended to settings with multimodal  
 106 data, often to combine text data and vision data [Arici et al. 2021, Geng et al. 2022, Bachmann  
 107 et al. 2022, Zhao et al. 2023, Mizrahi et al. 2024]. To do so, these models combine data across  
 108 modalities early on so that representations are multimodally fused through layers of the model. The  
 109 joint embeddings are then used for the MAE task to reconstruct inputs across all modalities. This  
 110 structure inherently allows for the possibility of cross-modal reconstruction, as information from  
 111 one modality can be used to reconstruct another. MAE methods can be more compute- and size-  
 112 efficient due to the fused encoding structure used and the large amounts of data typically dropped out  
 113 as a result of the masking strategy [Bachmann et al. 2022, Mizrahi et al. 2024]. Of the above, the  
 114 method used for our model is most similar to MultiMAE introduced in [Bachmann et al. 2022].

115 Both of these pretraining strategies have been applied to physiological signals, although examples  
 116 are sparser than in other domains. We first discuss examples using contrastive learning strategies.  
 117 [Abbaspourazad et al. 2023] uses a large-scale Apple Watch dataset to classify demographics and  
 118 health information from two modalities: photoplethysmography (PPG) and ECG. [Thapa et al. 2024]

119 used sleep data collected across EEG, EMG, ECG, and EOG sensors for downstream sleep-related  
120 classification tasks. [Raghu et al. \[2022\]](#) uses cardiac and blood-related signals to predict mortality  
121 rate and pulmonary arterial pressure. Both [Abbaspourazad et al. \[2023\]](#) and [Raghu et al. \[2022\]](#) use  
122 a SimCLR-like strategy through data augmentations, while [Thapa et al. \[2024\]](#) uses a CLIP-like  
123 strategy and construct data pairs across modalities.

124 In comparison to contrastive methods, MAE pretraining is less common for multimodal physiological  
125 signals. [Mathew et al. \[2024\]](#) uses MAE-pretraining in a model for phonocardiogram (PCG) and  
126 ECG data. The data was collected from digital stethoscopes, and the model was finetuned to classify  
127 signatures of cardiovascular disease. The closest example to our work is from [Liu et al. \[2023\]](#), where  
128 a multimodal transformer model is pretrained on EEG, EMG, and EOG signals with a MultiMAE-like  
129 objective. However, this work was more limited in dataset size (100 patients in each pretraining  
130 dataset) and focused on one specific downstream task per pretrained model. In our work, we use a  
131 larger dataset with 1,985 patients and evaluate how well MultiMAE-pretrained models can perform  
132 on diverse downstream tasks. We later will make comparisons with contrastive methods as well.

133 A focus of our work is in encouraging cross-modal representation learning. This is inspired by works  
134 arguing that multimodal learning can be improved by optimizing for cross-modal reconstruction  
135 [\[Kleinman et al. 2023\]](#), [\[Hussen Abdelaziz et al. 2020\]](#), [\[Hazarika et al. 2022\]](#). While this objective is  
136 already present in the original MultiMAE algorithm, a simple way to further encourage cross-modal  
137 learning is to randomly drop modalities from the input [\[Hazarika et al. 2022\]](#), [\[Hussen Abdelaziz  
138 et al. 2020\]](#), [\[Arici et al. 2021\]](#), [\[Deldari et al. 2023\]](#). This pressures the model to learn relationships  
139 across modalities in order to satisfy the reconstruction task. In the health data field, modality dropout  
140 strategies have been used to improve performance in tasks with missing modalities or heterogeneous  
141 noise, but they are still limited in their use in a general pretraining strategy. Furthermore, analyses of  
142 how representations are shaped by multimodal fusion are largely unexplored. We investigate both  
143 these questions in this work.

## 144 3 Methods

### 145 3.1 Dataset

146 We use the publicly available PhysioNet 2018 Challenge dataset [\[Ghassemi et al. 2018\]](#). This dataset  
147 consists of physiological signals collected during overnight sleep from 1,985 subjects. On average,  
148 each subject contributes 7.7 hours of recording [\[Ghassemi et al. 2018\]](#). The dataset contains many  
149 sensors, but here we focus on EEG, EMG, EOG, and ECG recordings (Figure 1A). For EEG, we use  
150 only the F3-M2 differential pair for our main results. We note that the signals from these sensors  
151 show distinct characteristics and are not obviously related (Figure 1A).

152 Patient demographics such as age and gender were also recorded in the dataset. The physiological data  
153 comprises of unlabeled data from 989 patients and labeled data from 996 patients. In the labeled set,  
154 30-second contiguous windows were manually annotated by several certified sleep technologists into  
155 one of five sleep stages: wakefulness, stage 1, stage 2, stage 3, or rapid eye movement (REM). The  
156 same windows were also manually annotated for the presence of arousals (e.g. snores, vocalizations,  
157 respiratory effort, leg movement, etc.). Prior literature using this dataset mostly focus on the sleep  
158 staging task [\[Perslev et al. 2019\]](#), [\[Banville et al. 2021\]](#), [\[Phan et al. 2021\]](#), and comparisons to these  
159 works are discussed in the Appendix.

160 To prevent data leakage, data is split over patient identity. The pretraining dataset is comprised of  
161 all 989 patients in the unlabeled set and 657 patients in the labeled set. The training/finetuning dataset  
162 for downstream tasks is comprised of the 657 patients in the labeled set that were used for pretraining  
163 (that is, the training/finetuning dataset is a subset of the pretraining dataset). The validation set and  
164 test set are constructed from the remaining 117 and 219 patients of the labeled set, respectively,  
165 and are not seen in either pretraining or training/finetuning. The validation set is used to select  
166 hyperparameters of the model and the test set is used for the evaluation scores reported in the results.  
167 A visualization of these data splits are in Figure 1B.

168 The signals are preprocessed with an anti-aliasing bandpass FIR filter then downsampled from 200  
169 Hz to 100 Hz using decimation by 2. Specifically, EEG and EOG signals were filtered to 0.1-30 Hz  
170 [\[Feng et al. 2021\]](#), [\[Satapathy et al. 2024\]](#). EMG and ECG signals were filtered to 0.1-70 Hz [\[Burns](#)

Table 1: *Balanced accuracy with linear probe evaluation: unimodal vs multimodal.* All models are pretrained before the encoder is frozen and representations are linearly probed for each task. We show the test balancy accuracy for a random guess (“Random”), for models trained entirely from scratch (“Scratch”), and for models pretrained and then linearly probed for the task (“Pretrained”). Note that “Pretrained-All” is a multimodal model pretrained with MultiMAE and input modality drop. Mean score and standard deviation for the three tasks are shown in columns. We additionally define an aggregate score which gives the average score over all tasks, normalized by the corresponding chance performance value (a score of 0 would indicate no improvement from chance). 500 patients are used in the training set for task finetuning. 5 random seeds are used in each training/finetuning stage. Asterisks indicate the best-performing unimodal model for each task.

		Sleep	Age	Arousal	Aggregate
Random	–	0.2	0.5	0.5	0.0
Scratch	EEG	0.717 ± 0.003	0.641 ± 0.004	0.568 ± 0.093	1.0 ± 0.101
	EMG	0.461 ± 0.004	0.55 ± 0.006	0.538 ± 0.074	0.494 ± 0.076
	EOG	0.697 ± 0.006	0.626 ± 0.006	0.56 ± 0.082	0.952 ± 0.084
	ECG	0.279 ± 0.006	0.605 ± 0.022	0.516 ± 0.038	0.213 ± 0.025
	All	0.737 ± 0.003	0.626 ± 0.018	0.595 ± 0.013	1.042 ± 0.009
Pretrained	EEG	<b>0.745 ± 0.001*</b>	0.662 ± 0.001	0.604 ± 0.093	1.085 ± 0.106
	EMG	0.442 ± 0.001	0.615 ± 0.003	0.533 ± 0.048	0.502 ± 0.052
	EOG	0.727 ± 0.001	0.653 ± 0.003	0.636 ± 0.071*	1.07 ± 0.078
	ECG	0.339 ± 0.002	0.703 ± 0.002*	0.526 ± 0.04	0.385 ± 0.042
	All	0.744 ± 0.001	<b>0.719 ± 0.002</b>	<b>0.637 ± 0.081</b>	<b>1.144 ± 0.09</b>

171 et al. [2007] Feng et al. [2021] Satapathy et al. [2024]. All signals are then resampled to 100 Hz. We  
 172 use 30-second samples of data for pretraining and for the downstream classification tasks.

173 Three tasks are constructed from this dataset: (1) sleep scoring, (2) age classification, and (3) arousal  
 174 identification. Sleep scoring is a 5-way classification problem. Both arousal and age will be treated as  
 175 a binary classification problem. In the age classification task, we aim to identify whether a patient’s  
 176 age is under 55 (the mean age) or not.

### 177 3.2 Model architecture

178 Our model architecture is based off that of the vision transformer [Alexey [2020]]. Modality-specific  
 179 tokenizer layers are followed by fused encoding layers, so that multimodal information is fused early  
 180 on (Figure 1C). The input to the model is a 30 second time series from multiple sensors sampled at  
 181 100 Hz. We divide each time series into 30 chunks that are one second each. These chunks are then  
 182 fed to the tokenizer layers. Tokenizers are trained for each modality and consist of one convolutional  
 183 layer and one linear layer. Specifically, each signal chunk first passes through a 1D convolutional  
 184 layer (with 64 channels and kernel size of 21) before a max pooling operation. Then, a linear layer  
 185 projects each token into a 512-dimensional embedding space. This is followed by layer normalization  
 186 to ensure signals from all modalities have comparable scales. Given 1,985 patients with an average  
 187 of 7.7 hours of recording time each, the total dataset size is 1,834,140.

188 To summarize, the output of a tokenizer for one modality is 30 tokens with embedding dimension  
 189  $D = 512$ . Sinusoidal positional embeddings and a learnable modality embedding are then added to  
 190 each token. Finally, tokens across modalities are fused through concatenation.

191 This fused vector is then passed to the joint encoding layers, which is comprised of eight transformer  
 192 layers with multi-head self-attention [Vaswani [2017]] and normalization before attention layers  
 193 [Xiong et al. [2020]]. Each transformer layer has 8 heads, and each layer has a 10% dropout rate  
 194 during training over attention weights and projection weights.

### 195 3.3 Pretraining objectives

196 We use a multimodal masked autoencoding (MAE) objective similar to MultiMAE from [Bachmann  
 197 et al. [2022]]. As mentioned above, tokens across all modality tokenizers are fused via concatenation.  
 198 In MultiMAE, a fixed portion of these tokens are masked at uniform and dropped from the fused

199 vector. We use a 70% masking rate (see Appendix for how the mask rate was selected). The  
200 remaining unmasked tokens are passed into the encoder and processed. To prepare the input for the  
201 decoder layers, the tokens that are output from the encoder are then interleaved with learnable mask  
202 tokens. Values in the mask token are initialized from  $\mathcal{N}(0, 0.02)$  with truncation at  $[-2, 2]$ . These  
203 learnable mask tokens act as placeholders for the signal to be reconstructed (i.e., the dropped tokens).  
204 Mask tokens are inserted in the location of the previously dropped tokens. Positional information is  
205 preserved by adding the appropriate positional embedding to the newly interleaved mask tokens.

206 A decoder is trained for each modality to reconstruct the original signal. Each decoder consists  
207 of a cross-attention layer and a transformer layer before a linear projection. The input into each  
208 modality decoder is the subset of tokens from the encoder output that corresponds to that modality.  
209 Cross-modal reconstruction is enabled through the cross-attention layer, where the input is the query  
210 and the entire encoder output is passed as keys/values. The linear layer projects each token from the  
211 embedding dimension (512) to the original signal dimension (100). The loss is calculated only over  
212 the reconstructed signal corresponding to the dropped tokens.

213 To encourage additional cross-modal interactions, we also use input modality drop during pretraining  
214 (Figure 1B) [Hazarika et al., 2022, Hussen Abdelaziz et al., 2020, Arici et al., 2021, Deldari et al.,  
215 2023]. In every batch, one randomly chosen modality is completely dropped on top of the typical  
216 MultiMAE uniform masking over tokens.

217 In later experiments we will make comparisons with contrastive learning objectives, resulting in  
218 modifications to the pretraining loss and the model architecture. In this case, the model will be  
219 converted to a late fusion structure that is typical for models trained with contrastive objectives.  
220 Further details can be found in the corresponding results section (§4.3) and Appendix F.

### 221 3.4 Finetuning

222 We are most interested in understanding how well representations learned by the pretrained model can  
223 transfer to multiple tasks. As such, after pretraining, we discard the decoder and freeze the encoder.  
224 The output of the encoder is layer normalized and average pooled over the token dimension. This  
225 512-dimensional vector is then passed to a linear classification head. A classification head is trained  
226 for each of the downstream tasks with weighted cross entropy loss to account for class imbalance.  
227 This is most relevant for the arousal detection task, where arousal events are extremely rare (2.7%  
228 of data samples). Although not the focus of this paper, we also conduct full finetuning experiments  
229 where both encoder and classifier parameters are trained (Appendix E).

### 230 3.5 Optimization

231 Models are pretrained for 2000 epochs or until a fixed compute time of 10 days is exceeded. For  
232 finetuning, models are trained for 200 epochs. Learning rates were scheduled with 10 epochs of  
233 linear warmup to  $1 \times 10^{-4}$  and cosine annealing thereafter [Loshchilov and Hutter, 2016]. We used  
234 the AdamW optimizer [Loshchilov and Hutter, 2017]. The model checkpoint chosen for evaluation  
235 was from the epoch where the lowest validation error was achieved, except in the case of pretraining  
236 the MultiMAE model with input modality drop. In this case, the validation error was quite noisy and  
237 the most recent checkpoint was chosen instead. Additional details can be found in the Appendix.

## 238 4 Experiments

### 239 4.1 A pretrained multimodal model develops representations that support a diverse set of 240 tasks in the PhysioNet18 dataset.

241 We first assess the extent to which pretraining and multimodal learning benefits downstream task  
242 performance in this dataset. We evaluate performance on the three tasks when both unimodal and  
243 multimodal models are trained from scratch. The balanced accuracy achieved by these models on  
244 the test set is shown in Table 1 (“Scratch” rows). In addition to the three tasks, we also define an  
245 aggregate score to highlight models that perform well across all tasks. The aggregate score is defined  
246 as  $\frac{1}{N} \sum_i \frac{s_i - r_i}{r_i}$ , where  $s_i$  is the average test score on task  $i$ ,  $r_i$  is the chance level performance for  
247 task  $i$ , and  $N = 3$  is the total number of tasks. Scores are measured using balanced accuracy. We  
248 see that the multimodal model performs overall better than any of the unimodal models (compare

Table 2: *Linear probe evaluation on all three tasks, comparing multimodal pretraining strategies.* All models are pretrained before the encoder is frozen and representations are linearly probed for each task. Mean test score and standard deviation for the three tasks are shown in columns. Aggregate score is defined as in Table 1. 500 patients are used in the training set for task finetuning. 5 random seeds are used in each training/finetuning stage.

Pretraining Strategy	Sleep		Age	
	Balanced Acc.	Cohen Kappa	Balanced Acc.	AUROC
Contrastive CLIP-style (LOO)	0.708 ± 0.0004	0.572 ± 0.001	0.643 ± 0.004	0.705 ± 0.006
Contrastive CLIP-style (Pairwise)	0.703 ± 0.001	0.559 ± 0.001	0.646 ± 0.0003	0.698 ± 0.0003
Contrastive SimCLR-style	0.656 ± 0.001	0.52 ± 0.001	0.624 ± 0.009	0.673 ± 0.015
MultiMAE Only	0.734 ± 0.001	0.618 ± 0.001	0.684 ± 0.001	0.758 ± 0.001
MultiMAE + Input Mod. Drop	<b>0.744 ± 0.001</b>	<b>0.63 ± 0.002</b>	<b>0.719 ± 0.002</b>	<b>0.785 ± 0.002</b>

Pretraining Strategy	Arousal		Aggregate Score
	Balanced Acc.	AUROC	
Contrastive CLIP-style (LOO)	<b>0.71 ± 0.002</b>	<b>0.776 ± 0.001</b>	1.082 ± 0.005
Contrastive CLIP-style (Pairwise)	0.708 ± 0.002	0.772 ± 0.001	1.075 ± 0.002
Contrastive SimCLR-style	0.585 ± 0.048	0.616 ± 0.070	0.900 ± 0.040
MultiMAE Only	0.604 ± 0.089	0.638 ± 0.136	1.082 ± 0.062
MultiMAE + Input Mod. Drop	0.637 ± 0.081	0.677 ± 0.128	<b>1.144 ± 0.058</b>

249 aggregate scores), although its performance on the sleep classification task slightly lags behind the  
 250 unimodal EEG model.

251 We next examine the benefits of pretraining the model and transferring the learned representations to  
 252 each of the downstream tasks. We first pretrain the unimodal models with masked autoencoding. The  
 253 test scores for these models are shown in Table 1 as well (“Pretrained” rows). Pretraining seems to  
 254 benefit all models, whether unimodal or multimodal. Interestingly, the pretrained unimodal models  
 255 reveal that a different modality is most informative for each task: EEG is more effective for sleep  
 256 staging, ECG for age classification, and EOG for arousal classification.

257 We then pretrain a multimodal model with MultiMAE and input modality drop. We evaluate this  
 258 model on the downstream tasks (“Pretrained, All” in Table 1). The multimodal model outperforms  
 259 the unimodal model in age classification and arousal classification, and performs very similarly to  
 260 EEG in the sleep staging task (Table 1). We find that the multimodal model performs well in all tasks  
 261 and achieves a higher aggregate score, despite the imbalance in modality dominance across tasks.

262 Notably, the improvement in aggregate score obtained by the multimodal model is greater when  
 263 training data is more limited (Appendix D). Given full-finetuning, though, the differences across  
 264 models are more minimal (Appendix E).

## 265 4.2 Adding input modality drop to MAE pretraining is important for downstream task 266 performance.

267 We chose our particular pretraining strategy with the hypothesis that encouraging multimodal fusion  
 268 improves performance in the downstream tasks. We investigate whether this is the case by first testing  
 269 the importance of using input modality drop (which theoretically should result in more cross-modal  
 270 learning). We compare task performance to that of a standard MultiMAE strategy, which does  
 271 not include input modality drop (Table 2). We see that removing input modality drop causes a  
 272 performance drop in downstream tasks (compare “MAE Only” to “MAE + Input Mod. Drop” in  
 273 Table 2). In fact, without input modality drop, MultiMAE underperforms the most informative  
 274 unimodal models across all tasks. Overall, dropping modalities in the input appears to be a simple  
 275 and effective means to increase performance over standard MultiMAE.

## 276 4.3 Late fusion models with contrastive learning objectives are more variable in performance.

277 Multimodal fusion is additionally encouraged in the MultiMAE model through the early fusion  
 278 architecture, where representations across modalities are mixed early in the network. We next make  
 279 comparisons to models with a late fusion structure where representations across modalities are not

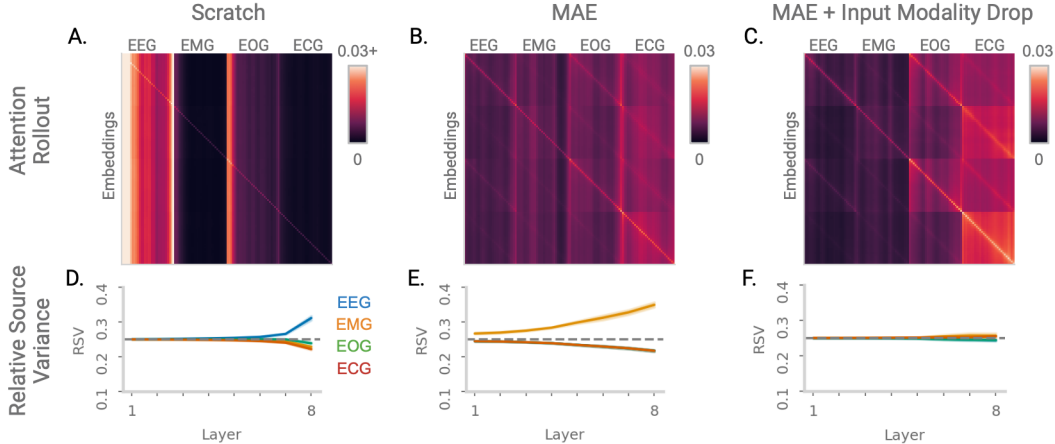


Figure 2: Measures of modality fusion across model representations. **A.** Attention rollout from tokens in the embeddings to tokens in the input. Here, the model is trained from scratch on sleep staging. Values are capped at 0.03 for comparisons with (BC). **B.** As in (A), but for the model pretrained with MAE. **C.** As in (A), but for the model pretrained with MAE and input modality drop. **D.** Relative source variance (RSV) of units across layers of the model in (A) to each of the four modalities. 95% confidence intervals shown, over 512 units in each embedding vector. **EF.** As in (D), but for the models in (B) and (C), respectively.

280 mixed except in the decoders for downstream tasks (Appendix F). To do so, we pretrain late fusion  
 281 models with contrastive learning, a common pretraining objective for these types of model.

282 We first test SimCLR-style multiview contrastive learning [Chen et al., 2020, Purushwalkam and  
 283 Gupta, 2020], with particular inspiration from [Raghu et al., 2022]. We randomly generate aug-  
 284 mentations for all input data samples (using the same signal augmentations from [Raghu et al.,  
 285 2022]). Positive pairs are defined as representations from adjacent time windows. We find that the  
 286 SimCLR-style model underperform standard MultiMAE in all tasks (Table 2). This may indicate that  
 287 defining desired relationships between of modality embeddings is important for the performance of a  
 288 contrastive learning model.

We next test CLIP-style pretraining to assess the benefits of using modality contrast in the contrastive  
 learning loss [Radford et al., 2021, Yuan et al., 2021, Zhang et al., 2022]. We will use two objectives  
 defined in [Thapa et al., 2024], a previous work in physiological signals that inspired our approach  
 here. [Thapa et al., 2024] defined a pairwise loss and a leave-one-out (LOO) loss:

$$l_{ijk}^{pair} = -\log \frac{\exp(\text{sim}(x_k^i, x_k^j)) * \tau}{\sum_{m=1}^N \exp(\text{sim}(x_k^i, x_m^j)) * \tau} \quad l_{ik}^{LOO} = -\log \frac{\exp(\text{sim}(x_k^i, \bar{x}_k^{\neq i})) * \tau}{\sum_{m=1}^N \exp(\text{sim}(x_k^i, \bar{x}_m^{\neq i})) * \tau}$$

289 for modalities  $i$  and  $j$ , sample  $k$ , temperature  $\tau$ , and modality embedding  $x$ .  $N$  is the total number of  
 290 samples, and  $\bar{x}_k^{\neq i}$  is the average of representations that are not modality  $i$  given data sample  $k$ . We find  
 291 that both CLIP-style models underperform even standard MultiMAE in sleep and age classification  
 292 (Table 2). Surprisingly, the contrastive model does extremely well in arousal classification. However,  
 293 in terms of aggregate performance, using MultiMAE with input modality drop is still most effective  
 294 out of the strategies we tested.

295 Despite these results, we speculate that developing new formulations of contrastive learning may  
 296 improve task performance. These methods are highly sensitive to the choice of positive and negative  
 297 pairs. It may be that contrastive methods in multimodal biosignals require domain-specific design to  
 298 reach their full potential.

#### 299 4.4 MAE + input modality drop encourages cross-modal fusion in attention weights and 300 model representations.

301 Finally, we wanted to understand whether our intuition about cross-modal fusion was indeed reflected  
 302 in the representations developed by the model. We first examine the attention weights of the model to



303 understand how much each output token from the encoder is influenced by input tokens from each  
 304 modality. We use a method called attention rollout [Abnar and Zuidema, 2020]. Attention rollout  
 305 accounts for the effects of the residual layers by defining the attention at layer  $l$  as a sum of the raw  
 306 attention weights and the identity matrix:  $A_l = 0.5W_l + 0.5I$  where  $W_l = \text{softmax}(Q_lK_l^T)$ . Thus,  
 307 to obtain the attention of the output embedding to the inputs, attention weights are rolled out across  
 308 model layers:  $A_L * A_{L-1} * \dots * A_2 * A_1$ , for  $L$  layers in the encoder.

309 We plot the results of attention rollout first for a multimodal model trained from scratch on sleep-  
 310 scoring (Figure 2A). The attention matrix develops strong vertical bands, indicating that model  
 311 embeddings attend to specific tokens without any context-specificity. In this case, EEG and EOG  
 312 tokens are most dominant. We next plot the attention matrix for a model pretrained with MultiMAE  
 313 and MultiMAE with input modality drop (Figure 2BC). The attention weights are more evenly spread  
 314 across the matrix, indicating greater cross-modal attention, although some sparse vertical bands can  
 315 still be observed. This can be interpreted as greater context-specificity in the attention weights. We  
 316 also observe an additional effect from both MultiMAE models where attention weights become more  
 317 temporally aligned. That is, tokens largely attend to other tokens that occurred around the same  
 318 window of time (Figure 2BC), an effect that is also visible when examining the raw attention matrices  
 319  $W_l$  (Appendix).

320 Although attention rollout allowed us to better understand the benefits of MultiMAE pretraining, it is  
 321 unclear how input modality drop affects representations. To further investigate this, we next analyze  
 322 individual embedding units in the model to see how tuned they are to different modalities. We use  
 323 relative source variance (RSV), which quantifies the variance in the activity of a unit due to a particular  
 324 input modality [Kleinman et al., 2023]. As an example, assume we want to calculate RSV due to  
 325 EEG. First, let  $x_{EEG} \sim X_{EEG}$  be a sample of EEG data from the dataset (with similar notation for  
 326 all other modalities). The source variance of a unit  $a$  due to EEG when all other modalities  $j$  are  
 327 fixed at samples  $x_j$  is defined as

$$SV_a(X_{EEG}, x_{EMG}, x_{EOG}, x_{ECG}) = \text{Var}(f(X_{EEG}, X_{EMG} = x_{EMG}, X_{EOG} = x_{EOG}, X_{ECG} = x_{ECG})_a)$$

328 where  $f$  gives the output embedding from the encoder, averaged over tokens. Symmetrically, source  
 329 variance can also be defined for the other modalities. Taking the softmax over these source variances  
 330 for a unit  $a$  gives the relative source variance of  $a$ . Thus if unit  $a$  is uniformly tuned to all input  
 331 sources, it would have a RSV value of 0.25 for each modality.

332 We first measure the RSV values of embedding units in the model trained from scratch on the sleep  
 333 staging task (Figure 2D). We find that representations become more tuned for EEG in the later layers  
 334 of the model. This is likely because EEG is more informative for the task and thus the decoder places  
 335 greater emphasis on EEG over the other modalities. We next measure the RSV values for the MAE-  
 336 pretrained model (Figure 2E). Interestingly, we see that across layers, units in the model become  
 337 increasingly tuned to EMG input. This is likely because the model struggles most to reconstruct EMG  
 338 (Appendix) and thus places greater representation weight onto that modality. In contrast, the model  
 339 trained with MAE and modality drop is equally tuned to all modalities across all layers (Figure 2F).

## 340 5 Limitations and Discussion

341 We have shown the strength of a foundation model-style approach using physiological data with  
 342 a diverse set of downstream tasks. We compare a variety of approaches and argue that explicitly  
 343 incorporating objectives that promote cross-modal reconstruction greatly improves representation  
 344 quality for solving downstream tasks. Specifically, we find that incorporating input modality drop  
 345 is a simple, yet especially effective strategy. We note that making comparisons with other datasets  
 346 would be additionally informative, especially since multimodal fusion strategies are often dependent  
 347 on the dataset and task at hand [Ma et al., 2022]. In addition, developing a large range of downstream  
 348 tasks will provide better insights into the strengths of different pretraining strategies and help identify  
 349 those that are especially useful for general purpose training.

## 350 References

- 351 Salar Abbaspourazad, Oussama Elachqar, Andrew C Miller, Saba Emrani, Udhyakumar Nallasamy,  
352 and Ian Shapiro. Large-scale training of foundation models for wearable biosignals. *arXiv preprint*  
353 *arXiv:2312.05409*, 2023.
- 354 Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. *arXiv preprint*  
355 *arXiv:2005.00928*, 2020.
- 356 Julián N Acosta, Guido J Falcone, Pranav Rajpurkar, and Eric J Topol. Multimodal biomedical ai.  
357 *Nature Medicine*, 28(9):1773–1784, 2022.
- 358 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel  
359 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language  
360 model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736,  
361 2022.
- 362 Dosovitskiy Alexey. An image is worth 16x16 words: Transformers for image recognition at scale.  
363 *arXiv preprint arXiv: 2010.11929*, 2020.
- 364 Tarik Arici, Mehmet Saygin Seyfioglu, Tal Neiman, Yi Xu, Son Train, Trishul Chilimbi, Belinda  
365 Zeng, and Ismail Tutar. Mlim: Vision-and-language model pre-training with masked language and  
366 image modeling. *arXiv preprint arXiv:2109.12178*, 2021.
- 367 H Ceren Ates, Peter Q Nguyen, Laura Gonzalez-Macia, Eden Morales-Narváez, Firat Güder, James J  
368 Collins, and Can Dincer. End-to-end design of wearable sensors. *Nature Reviews Materials*, 7(11):  
369 887–907, 2022.
- 370 Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimae: Multi-modal multi-  
371 task masked autoencoders. In *European Conference on Computer Vision*, pages 348–367. Springer,  
372 2022.
- 373 Hubert Banville, Omar Chehab, Aapo Hyvärinen, Denis-Alexander Engemann, and Alexandre  
374 Gramfort. Uncovering the structure of clinical eeg signals with self-supervised learning. *Journal*  
375 *of Neural Engineering*, 18(4):046020, 2021.
- 376 Florian Bordes, Richard Yuanzhe Pang, Anurag Ajay, Alexander C Li, Adrien Bardes, Suzanne  
377 Petryk, Oscar Mañas, Zhiqiu Lin, Anas Mahmoud, Bargav Jayaraman, et al. An introduction to  
378 vision-language modeling. *arXiv preprint arXiv:2405.17247*, 2024.
- 379 Joseph W Burns, Flavia B Consens, Roderick J Little, Karen J Angell, Sid Gilman, and Ronald D  
380 Chervin. Emg variance during polysomnography as an assessment for rem sleep behavior disorder.  
381 *Sleep*, 30(12):1771–1778, 2007.
- 382 Stefano Canali, Viola Schiaffonati, and Andrea Aliverti. Challenges and recommendations for  
383 wearable devices in digital health: Data quality, interoperability, health equity, fairness. *PLOS*  
384 *Digital Health*, 1(10):e0000104, 2022.
- 385 Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum  
386 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- 387 Federico Del Pup and Manfredo Atzori. Applications of self-supervised learning to biomedical  
388 signals: A survey. *IEEE Access*, 2023.
- 389 Shohreh Deldari, Dimitris Spathis, Mohammad Malekzadeh, Fahim Kawsar, Flora Salim, and Akhil  
390 Mathur. Latent masking for multimodal self-supervised learning in health timeseries. *arXiv*  
391 *preprint arXiv:2307.16847*, 2023.
- 392 Linus Ericsson, Henry Gouk, Chen Change Loy, and Timothy M Hospedales. Self-supervised  
393 representation learning: Introduction, advances, and challenges. *IEEE Signal Processing Magazine*,  
394 39(3):42–62, 2022.
- 395 LX Feng, X Li, HY Wang, WY Zheng, YQ Zhang, DR Gao, and MQ Wang. Automatic sleep staging  
396 algorithm based on time attention mechanism. *front hum neurosci* 15: 692054, 2021.

- 397 Xinyang Geng, Hao Liu, Lisa Lee, Dale Schuurmans, Sergey Levine, and Pieter Abbeel. Multimodal  
398 masked autoencoders learn transferable representations. *arXiv preprint arXiv:2205.14204*, 2022.
- 399 Mohammad M Ghassemi, Benjamin E Moody, Li-Wei H Lehman, Christopher Song, Qiao Li, Haoqi  
400 Sun, Roger G Mark, M Brandon Westover, and Gari D Clifford. You snooze, you win: the  
401 physionet/computing in cardiology challenge 2018. In *2018 Computing in Cardiology Conference*  
402 *(CinC)*, volume 45, pages 1–4. IEEE, 2018.
- 403 Jie Gui, Tuo Chen, Jing Zhang, Qiong Cao, Zhenan Sun, Hao Luo, and Dacheng Tao. A survey on self-  
404 supervised learning: Algorithms, applications, and future trends. *arXiv preprint arXiv:2301.05712*,  
405 2023.
- 406 Devamanyu Hazarika, Yingting Li, Bo Cheng, Shuai Zhao, Roger Zimmermann, and Soujanya Poria.  
407 Analyzing modality robustness in multimodal sentiment analysis. *arXiv preprint arXiv:2205.15465*,  
408 2022.
- 409 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
410 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer*  
411 *vision and pattern recognition*, pages 16000–16009, 2022.
- 412 Ahmed Hussen Abdelaziz, Barry-John Theobald, Paul Dixon, Reinhard Knothe, Nicholas Apostoloff,  
413 and Sachin Kajareker. Modality dropout for improved performance-driven talking faces. In  
414 *Proceedings of the 2020 International Conference on Multimodal Interaction*, pages 378–386,  
415 2020.
- 416 Vanessa Ibáñez, Josep Silva, and Omar Cauli. A survey on sleep assessment methods. *PeerJ*, 6:  
417 e4849, 2018.
- 418 Sheikh MA Iqbal, Imadeldin Mahgoub, E Du, Mary Ann Leavitt, and Waseem Asghar. Advances in  
419 healthcare wearable devices. *NPJ Flexible Electronics*, 5(1):9, 2021.
- 420 In Cheol Jeong, David Bychkov, and Peter C Searson. Wearable devices for precision medicine and  
421 health state monitoring. *IEEE Transactions on Biomedical Engineering*, 66(5):1242–1258, 2018.
- 422 Yizhang Jin, Jian Li, Yexin Liu, Tianjun Gu, Kai Wu, Zhengkai Jiang, Muyang He, Bo Zhao, Xin  
423 Tan, Zhenye Gan, et al. Efficient multimodal large language models: A survey. *arXiv preprint*  
424 *arXiv:2405.10739*, 2024.
- 425 Michael Kleinman, Alessandro Achille, and Stefano Soatto. Critical learning periods for multisensory  
426 integration in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
427 *and Pattern Recognition*, pages 24296–24305, 2023.
- 428 Felix Krones, Umar Marikkar, Guy Parsons, Adam Szmul, and Adam Mahdi. Review of multimodal  
429 machine learning approaches in healthcare. *arXiv preprint arXiv:2402.02460*, 2024.
- 430 Ran Liu, Ellen L Zippi, Hadi Pouransari, Chris Sandino, Jingping Nie, Hanlin Goh, Erdrin Azemi,  
431 and Ali Moin. Frequency-aware masked autoencoders for multimodal pretraining on biosignals.  
432 *arXiv preprint arXiv:2309.05927*, 2023.
- 433 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv*  
434 *preprint arXiv:1608.03983*, 2016.
- 435 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
436 *arXiv:1711.05101*, 2017.
- 437 Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem,  
438 and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision  
439 language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
440 *Pattern Recognition*, pages 26439–26455, 2024.
- 441 Mengmeng Ma, Jian Ren, Long Zhao, Davide Testuggine, and Xi Peng. Are multimodal transformers  
442 robust to missing modality? In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
443 *Pattern Recognition*, pages 18177–18186, 2022.

- 444 George Mathew, Daniel Barbosa, John Prince, and Subramaniam Venkatraman. Foundation models  
445 for cardiovascular disease detection via biosignals from digital stethoscopes. 2024.
- 446 David Mizrahi, Roman Bachmann, Oguzhan Kar, Teresa Yeo, Mingfei Gao, Afshin Dehghan, and  
447 Amir Zamir. 4m: Massively multimodal masked modeling. *Advances in Neural Information*  
448 *Processing Systems*, 36, 2024.
- 449 Mathias Perslev, Michael Jensen, Sune Darkner, Poul Jørgen Jennum, and Christian Igel. U-time: A  
450 fully convolutional network for time series segmentation applied to sleep staging. *Advances in*  
451 *Neural Information Processing Systems*, 32, 2019.
- 452 Huy Phan and Kaare Mikkelsen. Automatic sleep staging of eeg signals: recent development,  
453 challenges, and future directions. *Physiological Measurement*, 43(4):04TR01, 2022.
- 454 Huy Phan, Oliver Y Chén, Minh C Tran, Philipp Koch, Alfred Mertins, and Maarten De Vos.  
455 Xsleepnet: Multi-view sequential model for automatic sleep staging. *IEEE Transactions on Pattern*  
456 *Analysis and Machine Intelligence*, 44(9):5903–5915, 2021.
- 457 Senthil Purushwalkam and Abhinav Gupta. Demystifying contrastive self-supervised learning:  
458 Invariances, augmentations and dataset biases. *Advances in Neural Information Processing*  
459 *Systems*, 33:3407–3418, 2020.
- 460 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
461 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
462 models from natural language supervision. In *International conference on machine learning*, pages  
463 8748–8763. PMLR, 2021.
- 464 Aniruddh Raghu, Payal Chandak, Ridwan Alam, John Guttag, and Collin Stultz. Contrastive pre-  
465 training for multimodal medical time series. In *NeurIPS 2022 Workshop on Learning from Time*  
466 *Series for Health*, 2022.
- 467 Santosh Kumar Satapathy, Biswajit Brahma, Baidyanath Panda, Paolo Barsocchi, and Akash Kumar  
468 Bhoi. Machine learning-empowered sleep staging classification using multi-modality signals.  
469 *BMC Medical Informatics and Decision Making*, 24(1):119, 2024.
- 470 Thanveer Shaik, Xiaohui Tao, Lin Li, Haoran Xie, and Juan D Velásquez. A survey of multimodal  
471 information fusion for smart healthcare: Mapping the journey from data to wisdom. *Information*  
472 *Fusion*, page 102040, 2023.
- 473 Rahul Thapa, Bryan He, Magnus Ruud Kjaer, Hyatt Moore IV, Gauri Ganjoo, Emmanuel Mignot,  
474 and James Y Zou. Sleepfm: Multi-modal representation learning for sleep across ecg, eeg and  
475 respiratory signals. In *AAAI 2024 Spring Symposium on Clinical Foundation Models*, 2024.
- 476 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- 477 Shakti N Wadekar, Abhishek Chaurasia, Aman Chadha, and Eugenio Culurciello. The evolution of  
478 multimodal model architectures. *arXiv preprint arXiv:2405.17927*, 2024.
- 479 Min Wu and Jake Luo. Wearable technology applications in healthcare: a literature review. *Online J.*  
480 *Nurs. Inform.*, 23(3), 2019.
- 481 Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal  
482 llm. *arXiv preprint arXiv:2309.05519*, 2023.
- 483 Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang,  
484 Yanyan Lan, Liwei Wang, and Tiejian Liu. On layer normalization in the transformer architecture.  
485 In *International Conference on Machine Learning*, pages 10524–10533. PMLR, 2020.
- 486 Xin Yuan, Zhe Lin, Jason Kuen, Jianming Zhang, Yilin Wang, Michael Maire, Ajinkya Kale, and  
487 Baldo Faieta. Multimodal contrastive training for visual representation learning. In *Proceedings of*  
488 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6995–7004, 2021.
- 489 Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Con-  
490 trastive learning of medical visual representations from paired images and text. In *Machine*  
491 *Learning for Healthcare Conference*, pages 2–25. PMLR, 2022.

492 Zijia Zhao, Longteng Guo, Xingjian He, Shuai Shao, Zehuan Yuan, and Jing Liu. Mamo: Fine-grained  
493 vision-language representations learning with masked multimodal modeling. In *Proceedings of the*  
494 *46th International ACM SIGIR Conference on Research and Development in Information Retrieval*,  
495 pages 1528–1538, 2023.

496 **NeurIPS Paper Checklist**

497 **1. Claims**

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

500 Answer: [Yes]

501 Justification: Yes, the abstract and introduction match the claims made in the paper.

502 Guidelines:

- 503 • The answer NA means that the abstract and introduction do not include the claims  
504 made in the paper.
- 505 • The abstract and/or introduction should clearly state the claims made, including the  
506 contributions made in the paper and important assumptions and limitations. A No or  
507 NA answer to this question will not be perceived well by the reviewers.
- 508 • The claims made should match theoretical and experimental results, and reflect how  
509 much the results can be expected to generalize to other settings.
- 510 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
511 are not attained by the paper.

512 **2. Limitations**

513 Question: Does the paper discuss the limitations of the work performed by the authors?

514 Answer: [Yes]

515 Justification: Yes, limitations are explicitly mentioned in the discussion section. Throughout  
516 the results, we also discuss weaknesses in our approach and ways they can be improved.

517 Guidelines:

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

544 **3. Theory Assumptions and Proofs**

545 Question: For each theoretical result, does the paper provide the full set of assumptions and  
546 a complete (and correct) proof?

547 Answer: [NA]

548 Justification: We do not have theoretical results in this paper.

549 Guidelines:

- 550 • The answer NA means that the paper does not include theoretical results.
- 551 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 552 referenced.
- 553 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 554 • The proofs can either appear in the main paper or the supplemental material, but if
- 555 they appear in the supplemental material, the authors are encouraged to provide a short
- 556 proof sketch to provide intuition.
- 557 • Inversely, any informal proof provided in the core of the paper should be complemented
- 558 by formal proofs provided in appendix or supplemental material.
- 559 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 560 4. Experimental Result Reproducibility

561 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

562 perimental results of the paper to the extent that it affects the main claims and/or conclusions

563 of the paper (regardless of whether the code and data are provided or not)?

564 Answer: [Yes]

565 Justification: Yes, details of the implementation needed to reproduce the paper are mentioned

566 in the methods and the appendix.

567 Guidelines:

- 568 • The answer NA means that the paper does not include experiments.
- 569 • If the paper includes experiments, a No answer to this question will not be perceived
- 570 well by the reviewers: Making the paper reproducible is important, regardless of
- 571 whether the code and data are provided or not.
- 572 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 573 to make their results reproducible or verifiable.
- 574 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 575 For example, if the contribution is a novel architecture, describing the architecture fully
- 576 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 577 be necessary to either make it possible for others to replicate the model with the same
- 578 dataset, or provide access to the model. In general, releasing code and data is often
- 579 one good way to accomplish this, but reproducibility can also be provided via detailed
- 580 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 581 of a large language model), releasing of a model checkpoint, or other means that are
- 582 appropriate to the research performed.
- 583 • While NeurIPS does not require releasing code, the conference does require all submis-
- 584 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 585 nature of the contribution. For example
- 586 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 587 to reproduce that algorithm.
- 588 (b) If the contribution is primarily a new model architecture, the paper should describe
- 589 the architecture clearly and fully.
- 590 (c) If the contribution is a new model (e.g., a large language model), then there should
- 591 either be a way to access this model for reproducing the results or a way to reproduce
- 592 the model (e.g., with an open-source dataset or instructions for how to construct
- 593 the dataset).
- 594 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 595 authors are welcome to describe the particular way they provide for reproducibility.
- 596 In the case of closed-source models, it may be that access to the model is limited in
- 597 some way (e.g., to registered users), but it should be possible for other researchers
- 598 to have some path to reproducing or verifying the results.

#### 599 5. Open access to data and code

600 Question: Does the paper provide open access to the data and code, with sufficient instruc-

601 tions to faithfully reproduce the main experimental results, as described in supplemental

602 material?

603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654

Answer: [Yes]

Justification: The data is an open source dataset that has been used across many works. The code will be made publicly available after the double-blind review process.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 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).
- 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.
- 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.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

**6. Experimental Setting/Details**

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?

Answer: [Yes]

Justification: Yes, these details are provided in the methods and appendix. Furthermore, publicly provided code will have all training details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

**7. Experiment Statistical Significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Yes, error bars, number of samples, and the type of error bar statistic are all provided with each figure.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)



- 655 • The assumptions made should be given (e.g., Normally distributed errors).
- 656 • It should be clear whether the error bar is the standard deviation or the standard error
- 657 of the mean.
- 658 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 659 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 660 of Normality of errors is not verified.
- 661 • For asymmetric distributions, the authors should be careful not to show in tables or
- 662 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 663 error rates).
- 664 • If error bars are reported in tables or plots, The authors should explain in the text how
- 665 they were calculated and reference the corresponding figures or tables in the text.

## 666 8. Experiments Compute Resources

667 Question: For each experiment, does the paper provide sufficient information on the com-  
668 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
669 the experiments?

670 Answer: [Yes]

671 Justification: Yes, the appendix provides thorough detail of compute resources used and the  
672 time experiments took.

673 Guidelines:

- 674 • The answer NA means that the paper does not include experiments.
- 675 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
676 or cloud provider, including relevant memory and storage.
- 677 • The paper should provide the amount of compute required for each of the individual  
678 experimental runs as well as estimate the total compute.
- 679 • The paper should disclose whether the full research project required more compute  
680 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
681 didn't make it into the paper).

## 682 9. Code Of Ethics

683 Question: Does the research conducted in the paper conform, in every respect, with the  
684 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

685 Answer: [Yes]

686 Justification: Yes, we have reviewed the code of ethics and our paper conforms with this  
687 code.

688 Guidelines:

- 689 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 690 • If the authors answer No, they should explain the special circumstances that require a  
691 deviation from the Code of Ethics.
- 692 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
693 eration due to laws or regulations in their jurisdiction).

## 694 10. Broader Impacts

695 Question: Does the paper discuss both potential positive societal impacts and negative  
696 societal impacts of the work performed?

697 Answer: [Yes]

698 Justification: Yes, we have a healthcare-centered motivation which we discuss thoroughly  
699 in the introduction.

700 Guidelines:

- 701 • The answer NA means that there is no societal impact of the work performed.
- 702 • If the authors answer NA or No, they should explain why their work has no societal  
703 impact or why the paper does not address societal impact.

- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 723 11. Safeguards

724 Question: Does the paper describe safeguards that have been put in place for responsible  
725 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
726 image generators, or scraped datasets)?

727 Answer: [NA]

728 Justification: The paper does not pose a risk in terms of data or model misuse. The data is  
729 already openly available, and the model is centered around interpreting this data.

730 Guidelines:

- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 741 12. Licenses for existing assets

742 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
743 the paper, properly credited and are the license and terms of use explicitly mentioned and  
744 properly respected?

745 Answer: [Yes]

746 Justification: Yes, we give credit to existing data and code frameworks wherever relevant.

747 Guidelines:

- 748
- 749
- 750
- 751
- 752
- 753
- 754
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 755
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 756
- 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.
- 757
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 758
- 759
- 760
- 761
- 762

### 763 13. **New Assets**

764 Question: Are new assets introduced in the paper well documented and is the documentation  
765 provided alongside the assets?

766 Answer: [NA]

767 Justification: We do not introduce new assets.

768 Guidelines:

- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 769
- 770
- 771
- 772
- 773
- 774
- 775
- 776

### 777 14. **Crowdsourcing and Research with Human Subjects**

778 Question: For crowdsourcing experiments and research with human subjects, does the paper  
779 include the full text of instructions given to participants and screenshots, if applicable, as  
780 well as details about compensation (if any)?

781 Answer: [NA]

782 Justification: We did not collect new data with human subjects.

783 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 784
- 785
- 786
- 787
- 788
- 789
- 790
- 791

### 792 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 793 Subjects**

794 Question: Does the paper describe potential risks incurred by study participants, whether  
795 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
796 approvals (or an equivalent approval/review based on the requirements of your country or  
797 institution) were obtained?

798 Answer: [NA]

799 Justification: We did not collect new data with human subjects.

800 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - 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.
- 801
- 802
- 803
- 804
- 805

806  
807  
808  
809  
810

- 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.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.