LEVERAGING SHARED PROTOTYPES FOR A MULTIMODAL PULSE MOTION FOUNDATION MODEL

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

033

034

036

039

040

041

042

043

044

045

046

047

048

049

050

051

052

Paper under double-blind review

ABSTRACT

Modeling multi-modal time-series data is critical for capturing system-level dynamics, particularly in biosignals where modalities such as ECG, PPG, EDA, and accelerometry provide complementary perspectives on interconnected physiological processes. While recent self-supervised learning (SSL) advances have improved unimodal representation learning, existing multi-modal approaches often rely on CLIP-style contrastive objectives that overfit to easily aligned features and misclassify valid cross-modal relationships as negatives, resulting in fragmented and non-generalizable embeddings. To overcome these limitations, we propose ProtoMM, a novel SSL framework that introduces a shared prototype dictionary to anchor heterogeneous modalities in a common embedding space. By clustering representations around shared prototypes rather than explicit negative sampling, our method captures complementary information across modalities and provides a coherent "common language" for physiological signals. In this work, we focus on developing a Pulse Motion foundation model with ProtoMM and demonstrate that our approach outperforms contrastive-only and prior multimodal SSL methods, achieving state-of-the-art performance while offering improved interpretability of learned features.

1 Introduction

Digital biomarkers (for stress, physical activity, sleep, etc.) obtained from wearable sensors, such as smart watches and smartphones, provide unprecedented opportunities to give individuals novel insights into their states of health and wellness throughout their daily life, along with new tools for managing their health-related behaviors (Rehg et al., 2017). In order to realize this potential, however, it is critical to develop effective models for *multi-modal* time series biosignal data, so that complementary sensing modalities can be leveraged to overcome the ambiguities and noise that are inherent in wearable signals collected in the field environment.

Recently, there has been substantial progress in developing unimodal Foundation Models (FMs) which are pre-trained using large datasets on modalities such as accelerometry (Xu et al., 2024b; Yuan et al., 2024), ECG (Abbaspourazad et al., 2023; McKeen et al., 2024), and PPG (Saha et al., 2025; Pillai et al., 2024). These models have demonstrated effective generalization to downstream tasks and have established new benchmarks for performance. Building on these successes, recent works have focused on the challenge of how to align multiple signal modalities in pretraining multimodal FMs, often using CLIP-style contrast objectives that pull temporally aligned signals together (Thapa et al., 2024; Deldari et al., 2022; 2024; Zhang et al., 2025). A key challenge is to ensure that the resulting multimodal embedding captures both between-modality information (i.e., features shared across modalities, such as the features of the cardiac cycle that are present in all cardiovascular signals such as ECG, PPG, and ICG) and within-modality information (i.e., features that are unique to a single modality, such as the signatures of kinematic motion that are present in accelerometry). When modalities are highly complementary, such as PPG and accelerometry, there is a danger that the alignment process could emphasize between-modality features at the expense of within-modality features. This is because cardiovascular activity (as captured in the PPG signal) is only indirectly connected to kinematic movement (as captured via accelerometry), e.g. through the increase in cardiac activity which accompanies strenuous exercise. Emphasizing signal alignment could inadvertently discard information which is unique to each modality and critical for downstream tasks such as stress detection.

This paper introduces **ProtoMM**, a novel self-supervised alignment strategy for pre-training a pulse motion foundation models with PPG and accelerometry signal modalities. The goal of ProtoMM is to test the hypothesis that alignment of complementary signal modalities can be facilitated via a prototype-based approach, in which biosignals are discretized into prototype vectors as part of the embedding process. We hypothesize that these protypes will be effective in encoding within-modality features and preserving them during the embedding process. ProtoMM achieves this goal by first creating multiple augmented views from each modality. Embeddings from all views, derived using different augmentations of the same modality or from different modalities entirely, are then projected onto a shared set of prototype vectors. The model is trained using a Multimodal Prototype Prediction loss, where the prototypes probabilities of one view must predict the prototype assignment of another. By enforcing this consistency across all pairs of views, the prototypes become discrete, learnable anchors for the shared latent space for both within-modality and between-modality information. We develop and test ProtoMM for the task of joint multi-modal modeling of PPG and accelerometry data. This is beneficial because PPG can be used to detect the physiological stress response through cardiovascular changes (Jahanjoo et al., 2024), while integrating accelerometer data is essential to disambiguate between responses due to physical activity versus responses caused by psychological stress (Sevil et al., 2020; Sun et al., 2012).

We validate the potential of ProtoMM via thorough experimental evaluation that first demonstrate how ProtoMM captures within and between modality information, next demonstrate how the explicit modeling of latent states via discrete prototypes is particularly useful in our multimodal setting. With this, ProtoMM achieves achieves superior performance against state-of-the-art multimodal self-supervised learning methodologies, and we can qualitatively validate that our prototypes capture morphological similarities with higher-level semantic information. We will release our model weights and a codebase with the full training methodology, architecture, and reproducible evaluation code, upon acceptance. The main contributions of this work are:

- 1. We introduce **ProtoMM**, a novel multimodal self-supervised framework for pulse motion foundation model, that resolves a key limitation of existing alignment methods. By using a shared set of prototypes and a swapped prediction objective, our model is designed to capture both within-modality (unique) and between-modality (shared) information.
- 2. We conduct extensive experiments by evaluating its transferability on three downstream datasets across six distinct tasks. We obtain superior performance, showing that ProtoMM consistently outperforms leading multimodal and unimodal baselines.
- We demonstrate that the explicit nature of our prototype-based learning leads to improved interpretability. Through qualitative analysis, we show that individual prototypes can learn to represent specific, semantically meaningful physiological and behavioral states.

2 Related Work

In recent years, learning useful representations from unlabeled sensor data has become the predominant paradigm, leveraging the ease of wearable sensors to record large quantities of data in naturalistic conditions (Bycroft et al., 2018). Popular approaches include future prediction (Narayanswamy et al., 2024; Haresamudram et al., 2021), contrasting between randomly augmented segments (Tang et al., 2020; Haresamudram et al., 2022), probabilistic transformation prediction (Saeed et al., 2019; Yuan et al., 2024), and reconstruction of randomly masked data (Haresamudram et al., 2020; Narayanswamy et al., 2024; Xu et al., 2025; Miao et al., 2024).

Contrastive Representation Learning for Time-Series Data: Prior work has demonstrated that contrasting randomly transformed windows of sensor data is highly effective, e.g., SimCLR (Tang et al., 2020). Similarly, motif-based positive pair generation for contrastive training has shown great promise (Xu et al., 2024a;b). Using data from a single modality, these methods essentially learn within-modality information.

Alternatively, approaches like ColloSSL (Jain et al., 2022), CroSSL (Deldari et al., 2024), and CO-COA (Deldari et al., 2022) mine positives and negatives *across sensors and modalities*. They evaluate using diverse modalities, including accelerometers, gyroscopes, ECG, EMG, and EDA. As such, a critical drawback is the non-trivial nature of mining the pairs. More recently, aligning sensor data with natural language descriptions has emerged as an effective option, essentially adopting the

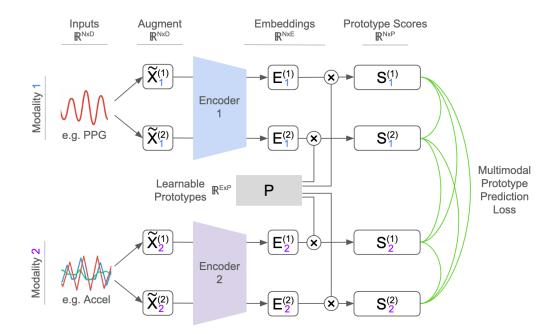


Figure 1: ProtoMM processes augmented segments from multiple modalities (i.e. PPG and Accelerometry) through dedicated encoders to produce the embeddings. The embeddings are then projected onto a shared set of prototype vectors, and the model is trained with a Multimodal Prototype Prediction Loss (\mathcal{L}_{MPP}) that learns to capture both within- and between modality information without relying on negative sampling.

CLIP framework (Radford et al., 2021) for time-series data. Methods such as IMU2CLIP (Moon et al., 2022), Ts2Act (Xia et al., 2024), and SensorLM (Zhang et al., 2025) have demonstrated the capabilities of such modeling. As such, a majority of these methods model only the common *between-modality information*. Our approach–ProtoMM–models both within- and between-modality information. Further, it sidesteps the challenges associated with mining positive/negative pairs by enforcing consistency between cluster assignments of augmented inputs.

Prototype-Based Representation Learning: Instead of instance-based discrimination used in approaches like SimCLR, SwAV (Caron et al., 2020) jointly clusters the data using prototype vectors, and enforces consistency between soft cluster assignments produced by different augmented inputs. VQ-VAE (Van Den Oord et al., 2017) also employs such vectors, and combines an autoencoder with vector quantization in order to perform online clustering with hard assignment. This setup has been extended to pose data as well (Zhang et al., 2023; Wang et al., 2024). Instead of the autoencoder, VQ-Wav2vec (Baevski et al., 2019; 2020) and VQ-CPC (Haresamudram et al., 2024) use CPC (Oord et al., 2018) as the base. As such, vector quantization based methods perform hard assignments of the prototype vectors. Our work builds on SwAV, which performs soft cluster assignments, leading to richer expressivity as there can be semantic overlap within the prototype vectors.

3 Protomm: Methodology and Design

We introduce ProtoMM, a multimodal self-supervised learning framework for pre-training pulse motion foundation model, which learns semantically meaningful representations by aligning multimodal time-series data to a shared set of prototypes. Our approach generalizes the swapped assignment prediction mechanism, originally proposed in the unimodal, two-view setting by SwAV (Caron et al., 2020), to a more general setting involving an arbitrary number of modalities as well as views per modality. By enforcing both within and between-modal consistency, the model learns features that are robust to augmentations while being coherent across different sensor modalities.

3.1 PROBLEM SETUP AND NOTATION

We define a multimodal time-series as $\mathbf{X}_t = \{\mathbf{X}_{t,1}, \mathbf{X}_{t,2}, \dots, \mathbf{X}_{t,M}\}$, where M is the number of sensor modalities and $\mathbf{X}_{t,m} \in \mathbb{R}^{T_m \times C_m}$ represents the data for the m-th modality. The subscript

t indicates that all modal data are temporally aligned to the same window; we omit it for brevity when the context is clear. The dimensions T_m and C_m allow for varying sampling frequencies and channel sizes across modalities.

A data augmentation module, $\mathcal{A}(\cdot)$, is used to generate A distinct views for each modality, creating an augmented set $\{\tilde{\mathbf{X}}_m^{(a)}\}$ for all modalities $m \in \{1,\ldots,M\}$ and views $a \in \{1,\ldots,A\}$. Each augmented view $\tilde{\mathbf{X}}_m^{(a)}$ is then passed through a modality-specific encoder $E_m(\cdot)$ to produce a normalized embedding $\mathbf{E}_m^{(a)} = E_m(\tilde{\mathbf{X}}_m^{(a)})$, where $\mathbf{E}_m^{(a)} \in \mathbb{R}^E$. The resulting set of embeddings can then be aligned using a shared prototype space.

3.2 SHARED PROTOTYPE SPACE

To align representations across modalities, we introduce a set of P trainable prototype vectors, organized as columns in a matrix $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_P] \in \mathbb{R}^{E \times P}$. This prototype space is shared across all modalities and training instances, serving as a common set of semantic anchors.

Each embedding $E_m^{(a)}$ is projected onto the prototypes yielding similarity scores $S \in \mathbb{R}^P$. To convert these scores into a soft assignment, we use two distinct transformations. First, we compute a probability distribution $\mathbf{U}_m^{(a)}$ using a softmax function with a temperature parameter τ . Second, we compute an assignment target $\mathbf{V}_m^{(a)}$ using the Sinkhorn-Knopp algorithm (Cuturi, 2013), which enforces an equipartition constraint to prevent mode collapse by ensuring all prototypes are utilized equally across a batch.

$$\boldsymbol{S}_{m}^{(a)} = \boldsymbol{E}_{m}^{(a)} \boldsymbol{P} \tag{1}$$

$$U_m^{(a)} = \operatorname{Softmax}(S_m^{(a)}/\tau)$$
 (2)

$$V_m^{(a)} = \operatorname{Sinkhorn}(S_m^{(a)}) \tag{3}$$

Now, the probability vector **U** and an assignment target **V** can be aligned via a cross entropy loss:

$$\ell(\boldsymbol{z}_t, \boldsymbol{q}_s) = -\boldsymbol{V} \cdot \log \boldsymbol{U} \tag{4}$$

3.3 Multimodal Prototype Prediction Loss

The key intuition of our method is that any pair of views originating from the same underlying system, regardless of modality, should be able to predict each other's prototype assignment. This is enforced through a swapped prediction loss comprising two components.

First, the within-modality loss, $\mathcal{L}_{within-mod}$, enforces consistency between all pairs of augmentations within a modality:

$$\mathcal{L}_{\text{within-mod}} = \sum_{m=1}^{M} \sum_{a=1}^{A} \sum_{b=1, b \neq a}^{A} \ell\left(U_m^{(a)}, V_m^{(b)}\right)$$
 (5)

Second, the between-modality loss, $\mathcal{L}_{between-mod}$, enforces consistency between all pairs from different modalities:

$$\mathcal{L}_{\text{between-mod}} = \sum_{m=1}^{M} \sum_{n=1}^{M} \sum_{n=1}^{A} \sum_{k=1}^{A} \ell\left(\boldsymbol{U}_{m}^{(a)}, \boldsymbol{V}_{n}^{(b)}\right)$$
(6)

Our final objective, the Multimodal Prototype Prediction Loss (\mathcal{L}_{MSP}), is a linear combination of these two losses, normalized for stability:

$$\mathcal{L}_{MPP} = \frac{1}{A \times M} (\alpha \mathcal{L}_{within-mod} + (1 - \alpha) \mathcal{L}_{between-mod}), \tag{7}$$

The hyperparameter $\alpha \in [0,1]$ balances the contribution of the two objectives. Setting $\alpha=1$ reduces the objective to independent, unimodal swapped prediction on each modality found in SwAV (Caron et al., 2020). In this work, we set $\alpha=0.5$ to equally weight both sources of learning.

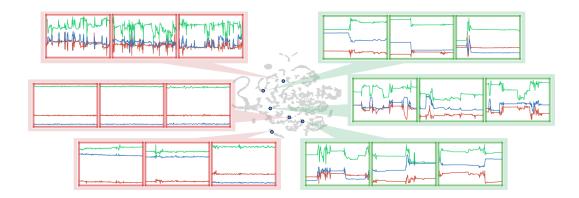


Figure 2: t-SNE of learned prototypes (gray), with k-means centroids (blue) and their top-three nearest **accelerometer** time-series. Panel borders denote ground-truth labels (green = Unstressed, red = Stressed). Each centroid captures a distinct motion motif, from active oscillatory bursts (top left) to sedentary plateaus (top right).

4 EXPERIMENTAL DESIGN

In Section 4.1, we first detail our large-scale pre-training dataset, along with the architecture and training settings. Then, in Section 4.2, we describe the unimodal and multimodal baselines that we compare against. Finally, in Section 4.3, we discuss our downstream datasets, their associated tasks, as well as our evaluation procedure.

4.1 Pre-training Setup

Our pre-training data comprises the initial 10 days of data from the large-scale Mobile Open Observation of Daily Stressors (MOODS) study Neupane et al. (2024). It contains 794,872 synchronized 30-second segments of accelerometer and PPG signals from 122 participants (39 men, 77 women, 6 non-binary; mean age 38±13 years). The data was collected "in-the-wild" from a wrist-worn WearOS smartwatch as participants went about their daily lives without specific instructions. Further details on the study design are available in Neupane et al. (2024).

To ensure fair comparison, all models and modalities use the same encoder architecture: a 1D ResNet-26 with a kernel size of 11, a stride of 2, and a final embedding dimension of 512. The only architectural difference is the number of input channels for the initial convolution layer: 3 for accelerometer data, 1 for PPG, and 4 for the early-fusion models that concatenate inputs. In Table 1, the "Input" column denotes early-fusion models with concatenated inputs as P+A and multimodal models with separate encoders as P|A. We also use a consistent set of augmentations: for the accelerometer, we use an empirically-established set comprising additive Gaussian noise, scaling, 3D rotation, negation, time reversal, channel shuffle, segment shuffle, and time warping (Tang et al., 2020; Xu et al., 2024b). For the single-channel PPG signal, we use the same set, omitting the inapplicable rotation and channel shuffling augmentations.

We utilize the Adam optimizer (Kingma, 2014) with a learning rate of 10^{-5} , no weight decay, and batch size of 256. Models are trained for 100 epochs or 96 hours (whichever comes first) on an NVIDIA L40S GPU. Both training and validation loss are calculated per batch but aggregated at epoch level, and the checkpoint with lowest validation loss is used for all downstream experiments.

4.2 BENCHMARKS

We compare ProtoMM against a comprehensive suite of multimodal and unimodal self-supervised baselines. Unless otherwise specified, all baselines are trained from scratch using the identical pre-training setup described above. Shared or modality-specific projection heads are incorporated following each baseline's original specifications. The baselines are as follows:

- **CLIP** (Radford et al., 2021; Thapa et al., 2024): Employs a contrastive loss to align the representations of temporally-paired between-modal signals.
- **COCOA** (Deldari et al., 2022): A multimodal contrastive method that aligns representations across modalities. The positive pairs are obtained from other temporally aligned sensors, whereas the negatives are mined from the same sensor, but from temporally misaligned data.
- **CroSSL** (Deldari et al., 2024): It stacks embeddings from modality-specific encoders, and performs random masking in the embedding space. A between-modal aggregator is utilized to obtain global embeddings, and the training is performed using the VICReg loss
- FOCAL (Liu et al., 2023): A multimodal contrastive method that factorized representations into orthogonal between and within-modal subspaces, while enforcing temporal locality. Training in the frequency domain uses temporal and spectral augmentations to form within-modal positives, and temporally aligned sensors to form between-modal positives.
- SLIP (Mu et al., 2022): A straightforward multimodal method that combines SimCLR and CLIP
 objectives to learn within and between-modality interactions, respectively.
- SimCLR (Chen et al., 2020): A unimodal contrastive learning framework that utilizes augmented versions of input data. This is the unimodal version of SLIP.
- **ProtoMM Within-Mod**: An unimodal ablation of our method with $\alpha = 1$ that only performs within-modal. This is equivalent to applying SwAV (Caron et al., 2020).

4.3 DOWNSTREAM EVALUATION SETUP

For evaluation, we use three datasets that contain synchronous PPG and accelerometer data, along with annotations: MOODS (stress detection, activity recognition) (Neupane et al., 2024), WESAD (stress detection) (Schmidt et al., 2018), and PPG-DaLiA (activity recognition, instantaneous heart rate prediction) (Reiss et al., 2019). All signals are resampled to 50 Hz to match the sampling frequency in pre-training dataset. After pre-training is complete, we freeze the model encoders and train a linear probe on the downstream tasks. For multimodal models, we generally will concatenate the embeddings from the PPG and accelerometer encoders to form the final representation. This is indicated by P+A in the "Out" column in our Results Table 1 and 2. For unimodal baselines, we use the single embedding directly. This is indicated by a P in the Out column in our Results Table 1.

- MOODS contains PPG and accelerometer signals collected using a Fossil Sport (Version 4) smartwatch from 122 participants. For stress detection (Stressed vs. Unstressed), we split the dataset into 1-minute windows following previous works (Mishra et al., 2018; Toshnazarov et al., 2024); whereas for binary activity recognition (Stationary vs. Non-stationary), we use 20-second samples for downstream evaluations.
- WESAD uses a wrist-worn Empatica E4 (McCarthy et al., 2016) to record accelerometer data (at 32 Hz) and blood volume pulse (at 64 Hz) from 15 participants. The stress detection task includes four sessions: Stress, Baseline, Amusement, and Meditation. For the binary classification task, following prior work (Schmidt et al., 2018; Dahal et al., 2023; Lange et al., 2024), we drop Meditation, merge Baseline and Amusement into Non-stress, while retaining the original Stress sessions. Each session is segmented into non-overlapping 1-minute windows for downstream evaluations.
- PPG-DaLiA contains accelerometer data (at 32 Hz) and blood volume pulse (at 64 Hz), collected using an Empatiaca E4 (McCarthy et al., 2016) worn on the wrist and ECG (700 Hz) from chestworn RespiBAN (biosignalsplux, 2019) to offer the ground truth of heart rate prediction. Fifteen subjects followed a semi-structured daily life protocol comprising eight distinct activities (Sitting, Ascending and Descending stairs, Table soccer, Cycling, Driving, Lunch break, Walking, and Working), with transient segments between activities annotated as an additional class. For model input, we segment all signals into 8-second windows with a 2-second sliding step.

We report macro F1-score and accuracy for classification tasks (stress, activity) and mean absolute error (MAE) and \mathbb{R}^2 for the regression task (heart rate prediction).

5 RESULTS AND DISCUSSION

In this section we discuss four key findings. First, we clearly demonstrate that prototype-based approach improves upon the established contrastive learning setup. Next, we contrast ProtoMM against a full suite of multimodal benchmarks and demonstrate state-of-the-art performance. Then,

Table 1: Linear-probe evaluation of unimodal and multimodal baselines. ProtoMM achieves best overall performance which indicates that grounding the alignment in a shared, discrete prototype space is a more effective mechanism for learning both shared and unique features simultaneously.

				MO	ODS			WE	SAD		PPG-DaLiA			
			Stress (2)		Activity (2)		Stress (4)		Stress (2)		Activity (9)		HR (R)	
Model	In	Out	↑F1	↑Acc	↑F1	↑Acc	↑F1	↑Acc	↑F1	↑Acc	↑F1	↑Acc	↓MAE	$\uparrow R^2$
SimCLR	A	A	0.477	0.598	0.861	0.925	0.557	0.629	0.793	0.836	0.559	0.560	15.37	0.101
SimCLR	P	P	0.527	0.607	0.655	0.857	0.349	0.517	0.692	0.699	0.302	0.399	8.14	0.670
PMM WMod	Α	Α	0.464	0.604	0.857	0.924	0.533	0.584	0.675	0.726	0.586	0.577	15.45	0.097
PMM WMod	P	P	0.470	<u>0.611</u>	0.595	0.848	0.522	0.607	0.847	0.877	0.285	0.397	8.29	0.670
SimCLR	P+A	P+A	0.488	0.601	0.849	0.919	0.445	0.596	0.721	0.753	0.542	0.536	9.91	0.534
PMM WMod	P+A	P+A	0.467	0.600	0.836	0.913	0.486	0.562	0.753	0.795	0.543	0.547	15.07	0.155
CLIP	P A	P+A	0.524	0.578	0.869	0.930	0.496	0.618	0.910	0.918	0.640	0.624	9.37	0.609
COCOA	P A	P+A	0.520	0.579	0.858	0.924	0.508	0.607	0.778	0.808	0.620	0.601	9.43	0.615
CroSSL	P A	P+A	0.461	0.622	0.751	0.882	0.431	0.494	0.783	0.808	0.553	0.540	11.47	0.464
FOCAL	P A	P+A	0.514	0.590	0.846	0.918	0.632	0.708	0.866	0.890	0.595	0.599	12.17	0.410
SLIP	P A	P+A	0.524	0.586	0.867	0.929	0.500	0.652	0.848	0.863	0.633	0.625	8.89	0.634
ProtoMM	P A	P+A	0.532	0.591	0.872	0.932	0.623	0.719	0.910	0.918	0.656	0.638	8.74	0.648

KEY: Encoder (In)put, Final Embedding (Out)put; (A)ccel, (P)PG; + designates concatenation, | designates separate encoders

we demonstrate how ProtoMM performs well because it is able to effectively integrate within- and between-modality information. Finally, we visualize the learned prototype space, and through examples show how the prototypes capture semantic meaning and specific morphologies.

Prototypes Explicitely Improve Performance. SLIP serves as the prototype-free analogue to ProtoMM. We utilize identical architecture and augmentation set for both models, making them equivalent up to the final embedding layer that produces $E_m^{(a)}$. Then, ProtoMM projects these embeddings onto a learned set of prototypes before applying a swapped prediction loss across between and within-modality pairs, but SLIP directly applies a contrastive (NT-Xent) loss to the embeddings themselves across all between- and within-modality pairs. Consequently, this comparison enables us to quantify the contribution of the prototype mechanism.

The results in Table 1 confirm the advantages of our prototype-based method. It demonstrates superior performance over its direct prototype-free analogue, SLIP, on every metric across all six downstream tasks. The performance gains are particularly pronounced on the WESAD dataset, where ProtoMM improves the F1-score for 4-class stress detection by a significant 24.6% (0.623 vs. 0.500) and for binary stress detection by over 7% (0.910 vs. 0.848).

Interestingly, performance gains remain isolated to the multimodal setting. In unimodal settings, SimCLR consistently outperforms its prototype-based counterpart, ProtoMM W-Mod. This aligns with prior work in unimodal time-series self-supervision (Meng et al., 2023). We hypothesize this is for two reasons: first, prototypes may act as a shared, discretized vocabulary that provides a common language to translate between disparate data streams like PPG and accelerometry. Second, as a negative-free method, ProtoMM avoids the "false negative" problem of contrastive learning. This issue is actively explored in the multimodal vision-language alignment domain (Byun et al., 2024; Chun, 2025), suggesting its potential relevance for multimodal time-series alignment as well.

ProtoMM Achieves State-of-the-art Performance. Table 1 presents comparisons against all baseline methods, which show that ProtoMM achieves SOTA results. It achieves the best overall performance in 4/6 downstream tasks and outperforms all other multimodal methods in 5/6 tasks.

The results show a clear trend where multimodal models outperform their unimodal counterparts, particularly on complex classification tasks. This highlights the value of alignment: one modality provides essential context that the other lacks. For example, heart rate information from PPG can help disambiguate activities with similar motion profiles from accelerometry, such as ascending versus descending stairs classes in PPG DaLiA activity classification or how research has shown that accelerometry and PPG signals can be used together for stress prediction (Sevil et al., 2020; Wu et al., 2015). The one exception to this trend is in the HR regression tasks. This highlights a key nuance in multimodal learning. Heart Rate is a metric that can be derived directly from the raw PPG waveform, and accelerometry signal gives little to no information on the heart rate. Therefore, the

Table 2: ProtoMM achieves the best performance at α =0.5, showing that the unified MPP objective successfully captures both within- and between-modality information.

					MO	ODS			WES	SAD		PPG-DaLiA			
				Stres	ss (2)	Activity (2)		Stress (4)		Stress (2)		Activity (9)		HR (R)	
α	Model	In	Out	↑F1	↑Acc	↑F1	†Acc	↑F1	†Acc	↑F1	↑Acc	↑F1	↑Acc	↓MAE	$\uparrow R^2$
.5	ProtoMM	P A	P+A	0.532	0.591	0.872	0.932	0.623	0.719	0.910	0.918	0.656	0.638	8.74	0.648
0	PMM BMod	P A	P+A	0.461	0.604	0.719	0.874	0.578	0.663	0.783	0.808	0.496	0.502	10.45	0.536
1	PMM WMod	P A	P+A	0.498	0.592	0.855	0.923	0.612	0.685	0.858	0.890	0.618	0.603	9.14	0.611

KEY: Encoder (In)put, Final Embedding (Out)put; (A)ccel, (P)PG; + designates concatenation, | designates separate encoders

Table 3: Integrating multimodal information into a unimodal embedding improves performance. ProtoMM rows show the performance of an unimodal embedding that was pre-trained multimodally (with both PPG and Accelerometer), while ProtoMM Within-Mod rows show the performance of an embedding pre-trained in isolation on just that single sensor.

					MO	ODS			WE	SAD		PPG-DaLiA			
				Stress (2)		Acitivity (2)		Stress (4)		Stress (2)		Activity (9)		HR (R)	
α	Model	In	Out	↑F1	↑Acc	↑F1	↑Acc	↑F1	↑Acc	↑F1	↑Acc	↑F1	↑Acc	↓MAE	$\uparrow R^2$
.5	ProtoMM	P A	A	0.496	0.597	0.872	0.931	0.620	0.663	0.837	0.863	0.615	0.597	15.0	0.152
1	PMM WMod	A	A	0.464	0.604	0.857	0.924	0.533	0.584	0.675	0.726	0.586	0.577	15.5	0.097
	ProtoMM PMM WMod				0.613 0.611										0.686 0.670

KEY: Encoder (In)put, Final Embedding (Out)put; (A)ccel, (P)PG; + designates concatenation, | designates separate encoders

inclusion of accelerometry within the embedding model only serves to further obfuscate the final embedding, such that multimodal models generally do worse. However, ProtoMM is able to more effectively disentangle the modalities, to better preserve within-modal PPG-specific information, achieving the best multimodal performance.

Out of the multimodal baselines, CLIP, COCA, and CroSSL focusing on modeling only the between-modal information, whereas SLIP, FOCAL and our ProtoMM method are the ones that explicitly model both between and within-modal information. However, interestingly, the models that model both do not uniformly outperform the models that model only between-modal information. Despite SLIP augmenting the CLIP objective with an additional within-modal loss, their overall performance is comparable as the 2nd best models after ProtoMM. This suggests that standard contrastive losses may struggle to balance the two objectives, potentially over-emphasizing the more difficult between-modal alignment. ProtoMM's success indicates that grounding the alignment in a shared, discrete prototype space is a more effective mechanism for learning both shared and unique features simultaneously.

Finally, the table shows the early fusion models that concatenate the raw signals before encoder input perform poorly. Within Table 1, they are marked as P+A in the In column. This finding demonstrates that despite both being biosignals, a naive concatenation of PPG and accelerometry is insufficient, and modality-specific encoders are essential for learning how to capture PPG-specific and Accelerometry-specific features from the raw sensor data.

Balancing within- and between-modal objectives. The core advantage of ProtoMM lies in its ability to effectively simultaneously and effectively capture between- and within-modality information through a unified objective, \mathcal{L}_{MPP} . We validate this design choice by modifying the loss weighting parameter to be $\alpha=1$, such that only within-modality ($\mathcal{L}_{within-mod}$) information is learned or to be $\alpha=0$, such that only between-modality ($\mathcal{L}_{between-mod}$) information is learned.

Table 2 shows that the optimal performance is achieved at $\alpha=0.5$, demonstrating that successful multimodal integration requires explicitly modeling both the unique contributions of each modality and their synergistic interactions. The balanced objective enables ProtoMM to learn representations that are simultaneously invariant to modality-specific augmentations while being semantically aligned across modalities.

Between-modal Knowledge Transfer to the Other Modality. We would like to further explore how effective ProtoMM integrates Between-modal information by training ProtoMM normally, to capture both between-modal and within-modal information, and then instead of concatenating the

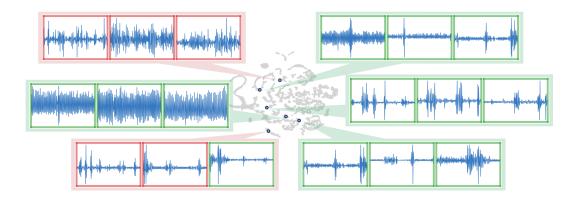


Figure 3: t-SNE of learned prototypes (gray), with k-means centroids (blue) and their top-three nearest **PPG** time-series. Panel borders denote ground-truth labels (green = Unstressed, red = Stressed). Each centroid captures a distinct pattern, from waveforms with high amplitude and variance (middle left) to those with a steady baseline and spiky variances (bottom right).

outputs of each modality's embedding for evaluation, instead just using one modality's embedding at a time, show by the 1st and 3rd rows in Table 3. This will help us investigate whether the modality's trained embedding captures more information than if it was trained independently. Therefore, the baseline is training each modality independently with a unimodal ProtoMM Within-Mod ($\alpha=1$) with only one modality as input, shown by the 2nd and 4th rows in Table 3.

As shown in Table 3, both encoders (i.e., for the accelerometer and PPG) outperform their unimodal ProtoMM Within-Mod counterparts on nearly all tasks, despite being evaluated without the full multimodal embedding. This indicates that the shared prototype space encourages each encoder to develop representations that are semantically aligned with the broader physiological context across modalities, not just its own signal characteristics, which results in more robust and informative features even when deployed as a unimodal embedding.

Interpreting Prototypes. To demonstrate that the explicit nature of our prototype-based learning leads to improved interpretability, we conduct a qualitative analysis on the WESAD stress detection dataset. We cluster the learned prototypes with k-means clustering (k=15), then for a given learned prototype centroid, we identify the accelerometer and PPG segments with highest cosine similarity in the representation space. Figure 2 and 3 shows a t-SNE projection of all prototypes, with representative examples highlighted alongside their top-3 nearest neighbors.

These neighbors exhibit both label consistency (stressed vs. unstressed) and coherent temporal dynamics that correspond to distinct physiological patterns (e.g., stable low-amplitude patterns, sharp oscillations, near-static segments with abrupt changes, strong trapezoid-like movements). This suggests that the shared prototype vectors are structured in the representation space in clusters that exhibit both label consistency and coherent physiological patterns. Rather than operating as a black box, ProtoMM's prototypes function as semantically meaningful anchors that capture both morphological signal characteristics and higher-level physiological contexts, offering a tangible insights into the latent physiological state.

6 Conclusion

In this paper, we presented ProtoMM, a prototype-based multimodal framework for self-supervised learning on time-series data to pre-train a pulse motion foundation model. By leveraging a shared prototype space, ProtoMM aligns embeddings from different sensor modalities and uncovers common latent physiological states. The learned prototypes act as cluster centers, structuring the representation space into coherent groups that reflect both morphological similarities and higher-level semantic consistency. Importantly, this explicit prototype space enhances interpretability, as individual prototypes correspond to meaningful physiological patterns. Comprehensive experiments across three datasets and six downstream tasks demonstrate that ProtoMM consistently outperforms twelve state-of-the-art baselines. Beyond contrastive or masking-based approaches, ProtoMM introduces a prototype-based swapped prediction objective, offering a new perspective on representation learning for multimodal time-series.

7 REPRODUCIBILITY

Our Methods section in Section 3 presents the model and training setup, and Experiments section in Section 4 describes the evaluation protocol and study design. Hyperparameters for every benchmark appear in Appendix A.1. We will release our model weights and a codebase with the full training methodology, architecture, and reproducible evaluation code, upon acceptance. The PPG-DaLiA and WESAD datasets are publicly available, and Section 4.3 explains how we curate and preprocess each one for our tasks.

8 ETHICS

Our paper develops models using physiological signals, with the goal of improving personal health. We acknowledge the associated risks, including privacy issues and the possibility of widening health disparities, as these models enable more detailed characterization of patients. Without effective regulation, patients may have limited control over their data, raising concerns about the upholding of autonomy, a core principle of medical ethics. Our study uses de-identified data from IRB-approved protocols, ensuring no participant identification information is included in our analysis. Nevertheless, we believe our work contributes positively to the field by advancing personalized health recommendations, which can enhance care quality, patient safety, and overall well-being. Additionally, we acknowledge the use of LLMs to assist in editing and polishing the writing for this submission, specifically, to edit phrasing and to clarify the framing of ideas in a manner that reflect the authors' original intent.

REFERENCES

- Salar Abbaspourazad, Oussama Elachqar, Andrew C Miller, Saba Emrani, Udhyakumar Nallasamy, and Ian Shapiro. Large-scale training of foundation models for wearable biosignals. *arXiv* preprint arXiv:2312.05409, 2023.
- Alexei Baevski, Steffen Schneider, and Michael Auli. vq-wav2vec: Self-supervised learning of discrete speech representations. *arXiv preprint arXiv:1910.05453*, 2019.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.
- biosignalsplux, 2019. URL https://bio-medical.com/media/support/ biosignalsplux_explorer_user_manual_v.1.0.pdf.
- Clare Bycroft, Colin Freeman, Desislava Petkova, Gavin Band, Lloyd T Elliott, Kevin Sharp, Allan Motyer, Damjan Vukcevic, Olivier Delaneau, Jared O'Connell, et al. The uk biobank resource with deep phenotyping and genomic data. *Nature*, 562(7726):203–209, 2018.
- Jaeseok Byun, Dohoon Kim, and Taesup Moon. Mafa: Managing false negatives for vision-language pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27314–27324, 2024.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Sanghyuk Chun. Multiplicity is an inevitable and inherent challenge in multimodal learning. *arXiv* preprint arXiv:2505.19614, 2025.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013.

Kamana Dahal, Brian Bogue-Jimenez, and Ana Doblas. Global stress detection framework combining a reduced set of hrv features and random forest model. *Sensors*, 23(11):5220, 2023.

- Shohreh Deldari, Hao Xue, Aaqib Saeed, Daniel V Smith, and Flora D Salim. Cocoa: Cross modality contrastive learning for sensor data. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 6(3):1–28, 2022.
- Shohreh Deldari, Dimitris Spathis, Mohammad Malekzadeh, Fahim Kawsar, Flora D Salim, and Akhil Mathur. Crossl: Cross-modal self-supervised learning for time-series through latent masking. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pp. 152–160, 2024.
- Harish Haresamudram, Apoorva Beedu, Varun Agrawal, Patrick L Grady, Irfan Essa, Judy Hoffman, and Thomas Plötz. Masked reconstruction based self-supervision for human activity recognition. In *Proceedings of the 2020 ACM International Symposium on Wearable Computers*, pp. 45–49, 2020.
- Harish Haresamudram, Irfan Essa, and Thomas Plötz. Contrastive predictive coding for human activity recognition. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 5(2):1–26, 2021.
- Harish Haresamudram, Irfan Essa, and Thomas Plötz. Assessing the state of self-supervised human activity recognition using wearables. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 6 (3), sep 2022. doi: 10.1145/3550299. URL https://doi.org/10.1145/3550299.
- Harish Haresamudram, Irfan Essa, and Thomas Ploetz. Towards learning discrete representations via self-supervision for wearables-based human activity recognition. *Sensors*, 24(4):1238, 2024.
- Anice Jahanjoo, Nima TaheriNejad, and Amin Aminifar. High-accuracy stress detection using wristworn ppg sensors. In 2024 IEEE International Symposium on Circuits and Systems (ISCAS), pp. 1–5, 2024. doi: 10.1109/ISCAS58744.2024.10558012.
- Yash Jain, Chi Ian Tang, Chulhong Min, Fahim Kawsar, and Akhil Mathur. Collossl: Collaborative self-supervised learning for human activity recognition. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 6(1):1–28, 2022.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lucas Lange, Nils Wenzlitschke, and Erhard Rahm. Generating synthetic health sensor data for privacy-preserving wearable stress detection. *Sensors*, 24(10):3052, 2024.
- Shengzhong Liu, Tomoyoshi Kimura, Dongxin Liu, Ruijie Wang, Jinyang Li, Suhas Diggavi, Mani Srivastava, and Tarek Abdelzaher. Focal: Contrastive learning for multimodal time-series sensing signals in factorized orthogonal latent space. *Advances in Neural Information Processing Systems*, 36:47309–47338, 2023.
- Cameron McCarthy, Nikhilesh Pradhan, Calum Redpath, and Andy Adler. Validation of the empatica e4 wristband. In 2016 IEEE EMBS international student conference (ISC), pp. 1–4. IEEE, 2016.
- Kaden McKeen, Laura Oliva, Sameer Masood, Augustin Toma, Barry Rubin, and Bo Wang. Ecg-fm: An open electrocardiogram foundation model, 2024. URL https://arxiv.org/abs/2408.05178.
- Qianwen Meng, Hangwei Qian, Yong Liu, Lizhen Cui, Yonghui Xu, and Zhiqi Shen. Mhccl: masked hierarchical cluster-wise contrastive learning for multivariate time series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9153–9161, 2023.
- Shenghuan Miao, Ling Chen, and Rong Hu. Spatial-temporal masked autoencoder for multi-device wearable human activity recognition. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 7(4):1–25, 2024.

Varun Mishra, Tian Hao, Si Sun, Kimberly N Walter, Marion J Ball, Ching-Hua Chen, and Xinxin Zhu. Investigating the role of context in perceived stress detection in the wild. In *Proceedings of the 2018 ACM International Joint Conference and 2018 International Symposium on Pervasive and Ubiquitous Computing and Wearable Computers*, pp. 1708–1716, 2018.

- Seungwhan Moon, Andrea Madotto, Zhaojiang Lin, Alireza Dirafzoon, Aparajita Saraf, Amy Bearman, and Babak Damavandi. Imu2clip: Multimodal contrastive learning for imu motion sensors from egocentric videos and text. *arXiv preprint arXiv:2210.14395*, 2022.
- Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. In *European conference on computer vision*, pp. 529–544. Springer, 2022.
- Girish Narayanswamy, Xin Liu, Kumar Ayush, Yuzhe Yang, Xuhai Xu, Shun Liao, Jake Garrison, Shyam Tailor, Jake Sunshine, Yun Liu, et al. Scaling wearable foundation models. *arXiv preprint arXiv:2410.13638*, 2024.
- Sameer Neupane, Mithun Saha, Nasir Ali, Timothy Hnat, Shahin Alan Samiei, Anandatirtha Nandugudi, David M Almeida, and Santosh Kumar. Momentary stressor logging and reflective visualizations: Implications for stress management with wearables. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pp. 1–19, 2024.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Arvind Pillai, Dimitris Spathis, Fahim Kawsar, and Mohammad Malekzadeh. Papagei: Open foundation models for optical physiological signals. *arXiv preprint arXiv:2410.20542*, 2024.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL https://arxiv.org/abs/2103.00020.
- James M Rehg, Susan A Murphy, and Santosh Kumar. *Mobile Health: Sensors, Analytic Methods, and Applications*. Springer, 2017. doi: 10.1007/978-3-319-51394-2.
- Attila Reiss, Ina Indlekofer, Philip Schmidt, and Kristof Van Laerhoven. Deep ppg: Large-scale heart rate estimation with convolutional neural networks. *Sensors*, 19(14):3079, 2019.
- Aaqib Saeed, Tanir Ozcelebi, and Johan Lukkien. Multi-task self-supervised learning for human activity detection. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 3(2):1–30, 2019.
- Mithun Saha, Maxwell A Xu, Wanting Mao, Sameer Neupane, James M Rehg, and Santosh Kumar. Pulse-ppg: An open-source field-trained ppg foundation model for wearable applications across lab and field settings. *arXiv preprint arXiv:2502.01108*, 2025.
- Philip Schmidt, Attila Reiss, Robert Duerichen, Claus Marberger, and Kristof Van Laerhoven. Introducing wesad, a multimodal dataset for wearable stress and affect detection. In *Proceedings of the 20th ACM international conference on multimodal interaction*, pp. 400–408, 2018.
- Mert Sevil, Mudassir Rashid, Mohammad Reza Askari, Zacharie Maloney, Iman Hajizadeh, and Ali Cinar. Detection and characterization of physical activity and psychological stress from wristband data. *Signals*, 1(2):188–208, 2020.
- Feng-Tso Sun, Cynthia Kuo, Heng-Tze Cheng, Senaka Buthpitiya, Patricia Collins, and Martin Griss. Activity-aware mental stress detection using physiological sensors. In *Mobile Computing, Applications, and Services: Second International ICST Conference, MobiCASE 2010, Santa Clara, CA, USA, October 25-28, 2010, Revised Selected Papers 2*, pp. 282–301. Springer, 2012.
- Chi Ian Tang, Ignacio Perez-Pozuelo, Dimitris Spathis, and Cecilia Mascolo. Exploring contrastive learning in human activity recognition for healthcare. *arXiv preprint arXiv:2011.11542*, 2020.

- Rahul Thapa, Bryan He, Magnus Ruud Kjaer, Hyatt Moore IV, Gauri Ganjoo, Emmanuel Mignot, and James Y. Zou. SleepFM: Multi-modal representation learning for sleep across ECG, EEG and respiratory signals. In AAAI 2024 Spring Symposium on Clinical Foundation Models, 2024. URL https://openreview.net/forum?id=cDXtscWCKC.
- Kobiljon Toshnazarov, Uichin Lee, Byung Hyung Kim, Varun Mishra, Lismer Andres Caceres Najarro, and Youngtae Noh. Sosw: Stress sensing with off-the-shelf smartwatches in the wild. *IEEE Internet of Things Journal*, 2024.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- Chongyang Wang, Yuan Feng, Lingxiao Zhong, Siyi Zhu, Chi Zhang, Siqi Zheng, Chen Liang, Yuntao Wang, Chengqi He, Chun Yu, et al. Ubiphysio: Support daily functioning, fitness, and rehabilitation with action understanding and feedback in natural language. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 8(1):1–27, 2024.
- Min Wu, Hong Cao, Hai-Long Nguyen, Karl Surmacz, and Caroline Hargrove. Modeling perceived stress via hrv and accelerometer sensor streams. In 2015 37th annual international conference of the IEEE engineering in medicine and biology society (EMBC), pp. 1625–1628. IEEE, 2015.
- Kang Xia, Wenzhong Li, Shiwei Gan, and Sanglu Lu. Ts2act: Few-shot human activity sensing with cross-modal co-learning. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 7(4):1–22, 2024.
- Maxwell Xu, Alexander Moreno, Hui Wei, Benjamin Marlin, and James Matthew Rehg. Rebar: Retrieval-based reconstruction for time-series contrastive learning. In *The Twelfth International Conference on Learning Representations*, 2024a.
- Maxwell A Xu, Jaya Narain, Gregory Darnell, Haraldur Hallgrimsson, Hyewon Jeong, Darren Forde, Richard Fineman, Karthik J Raghuram, James M Rehg, and Shirley Ren. Relcon: Relative contrastive learning for a motion foundation model for wearable data. *arXiv preprint arXiv:2411.18822*, 2024b.
- Maxwell A Xu, Girish Narayanswamy, Kumar Ayush, Dimitris Spathis, Shun Liao, Shyam A Tailor, Ahmed Metwally, A Ali Heydari, Yuwei Zhang, Jake Garrison, et al. Lsm-2: Learning from incomplete wearable sensor data. *arXiv preprint arXiv:2506.05321*, 2025.
- Hang Yuan, Shing Chan, Andrew P Creagh, Catherine Tong, Aidan Acquah, David A Clifton, and Aiden Doherty. Self-supervised learning for human activity recognition using 700,000 persondays of wearable data. *NPJ digital medicine*, 7(1):91, 2024.
- Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Yong Zhang, Hongwei Zhao, Hongtao Lu, Xi Shen, and Ying Shan. Generating human motion from textual descriptions with discrete representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 14730–14740, 2023.
- Yuwei Zhang, Kumar Ayush, Siyuan Qiao, A Ali Heydari, Girish Narayanswamy, Maxwell A Xu, Ahmed A Metwally, Shawn Xu, Jake Garrison, Xuhai Xu, et al. Sensorlm: Learning the language of wearable sensors. *arXiv preprint arXiv:2506.09108*, 2025.

A APPENDIX

A.1 BENCHMARK IMPLEMENTATIONS

Our code will be made publicly available upon acceptance. All benchmark implementations adhere to the experimental setup described in Section 4, utilizing identical 1D ResNet-26 encoders with global max temporal pooling to produce 512-dimensional embeddings. Each training sample is transformed into two stochastic views, and augmentations are sampled with equal probability. Training hyperparameters remain consistent across all methods as described in Section 4.1. Below we detail the implementation-specific parameters for each baseline.

• **CLIP** (Radford et al., 2021; Thapa et al., 2024): Pairwise contrastive alignment between temporally paired modalities. We employ modality-specific single-layer projection heads with same input and output dimensions (512). The temperature parameter is initialized to 1.0.

- COCOA (Deldari et al., 2022): We reimplement the original TensorFlow code (https://github.com/cruiseresearchgroup/COCOA) in PyTorch while maintaining architectural fidelity. Key parameters include temperature=0.1, scale_loss=1/32, and lambd=3.9e-3, window = 100. Modality-specific projectors consist of a flattening layer followed by a linear projection to maintain dimensional consistency during training. Due to the method's specific design, global pooling is disabled during training.
- CroSSL (Deldari et al., 2024): We reimplement the original TensorFlow code (https://github.com/Nokia-Bell-Labs/CroSSL) in PyTorch, which follows the original spatial masking approach with coverage=0.9. Embeddings from modality-specific encoders are processed through a shared projector network consists of several linear layers with ReLU activations in between, resulting in a final output dimension (proj_size) of 32.
- FOCAL (Liu et al., 2023): We adapt the official implementation (https://github.com/tomoyoshki/focal) to our codebase while preserving the core methodology. Modality-specific projectors consist of two linear layers with ReLU activation. Key parameters include temperature=0.5, sequence_length=4, and loss weights: shared contrastive=1, private contrastive=1, orthogonal=3, rank=5. Augmentations include standard temporal transformations followed by frequency-domain phase shifting.
- **SLIP** (Mu et al., 2022): We implement both within- and between-modality contrastive losses using the standard CLIP formulation with temperature initialized to 1.0.
- SimCLR (Chen et al., 2020): Unimodal contrastive learning with two augmented views per sample. temperature=0.1 is used for the contrastive loss.

All implementations maintain fairness through consistent encoder architecture, data preprocessing, and evaluation protocols. Differences arise only in method-specific components as detailed above, ensuring meaningful comparison while preserving each approach's distinctive characteristics.