# PULSE: Privileged Knowledge Transfer from Electrodermal Activity to Low-Cost Sensors for Stress Monitoring

**Zihan Zhao**                                                                    ZIZ078@UCSD.EDU
*University of California San Diego, San Diego, CA, USA*

**Ning Yan**                                                          YAN.NINGYAN@FUTUREWEI.COM
*IC Lab, Futurewei Technologies Inc., Santa Clara, CA, USA*

**Masood Mortazavi**                                       MASOOD.MORTAZAVI@FUTUREWEI.COM
*IC Lab, Futurewei Technologies Inc., Santa Clara, CA, USA*

## Abstract

Electrodermal activity (EDA), the primary signal for stress detection, requires costly hardware often unavailable in real-world wearables. In this paper, we propose PULSE, a framework that utilizes EDA exclusively during self-supervised pretraining, while enabling inference without EDA but with more readily available modalities such as ECG, BVP, ACC, and TEMP. Our approach separates encoder outputs into *shared* and *private* embeddings. We align "shared" embeddings across modalities and fuse them into a modality-invariant representation. The "private" embeddings carry modality-specific information to support the reconstruction objective. Pretraining is followed by knowledge transfer where a frozen EDA teacher transfers sympathetic-arousal representations into student encoders. On WESAD, our method achieves strong stress-detection performance, showing that representations of privileged EDA can be transferred to low-cost sensors to improve accuracy while reducing hardware cost.

**Keywords:** Electrodermal activity, privileged information, knowledge transfer, knowledge distillation, multimodal learning, stress detection, wearable computing, masked autoencoder

**Data and Code Availability** All experiments were performed on the publicly available WESAD dataset (Schmidt et al., 2018), which provides synchronized physiological signals (including ECG, BVP, ACC, TEMP, and EDA) with stress annotations. We are not sharing code at submission; we will release code upon acceptance.

**Institutional Review Board (IRB)** This study analyzes a publicly available dataset (WESAD) and does not involve interaction with human subjects; therefore IRB approval is not required.

## 1. Introduction

Wearables have enabled continuous monitoring of human physiological state. *Electrodermal activity* (EDA), minute changes in skin conductance, has served as an index of mental workload, affect, and acute stress (Critchley, 2002; Boucsein, 2012; Schmidt et al., 2018; Sánchez-Reolid et al., 2025; Roos and Slavich, 2023).

**When the cornerstone is missing.** Despite its diagnostic power (Stržinar et al., 2023; Matton et al., 2023), EDA is absent from many commercial trackers due to hardware requirements (Ag/AgCl electrodes and a constant-current source) and its vulnerability to motion artifacts (Hossain et al., 2022). Consequently, longitudinal datasets and deployed devices often provide only ECG/PPG (BVP), inertial signals, and temperature; models trained with EDA cannot rely on the modality at test time.

**Privileged knowledge transfer as a remedy.** We adopt the Learning Using Privileged Information (LUPI) paradigm: during training, a teacher receives EDA and transfers its knowledge to a student that *never* sees EDA, enabling the student to operate without the EDA sensor at inference (Vapnik and Izmailov, 2015; Lopez-Paz et al., 2016; Markov and Matsui, 2016). This approach differs from missing-modality robustness methods that align modalities symmetrically (e.g., PhysioOmni, ADAPT) but do not exploit a strictly more informative teacher modality (Jiang et al., 2025; Mordacq et al., 2024; Ibtehaz and Mortazavi, 2024).

**Identified gap.** Work on LUPI and knowledge distillation (KD) has shown clear benefits in vision and speech (Vapnik and Izmailov, 2015; Lopez-Paz et al., 2016; Markov and Matsui, 2016), and multimodal wearables research typically pursues robustness to missing channels via symmetric alignment (Jiang et al., 2025; Mordacq et al., 2024; Liu et al., 2024; Ibtehaz and Mortazavi, 2024). However, to our knowledge, no prior study leverages *EDA as a teacher* in a LUPI setup for knowledge transfer to student encoders of cheaper signals (e.g., ECG/BVP/ACC/TEMP).

**Contributions.** We propose *PULSE: Privileged EDA knowledge transfer Using Low-cost SEnsors* for wearable stress monitoring:

1. **EDA teacher.** We freeze a self-supervised EDA encoder and treat its representations and hidden state embeddings as privileged signals during training; EDA is not used at inference.

2. **Multi-sensor student.** A student consumes only deployable channels (ECG, BVP, ACC, and TEMP) and uses a shared–private latent split; we align *only* the shared subspace to avoid over-constraining modality-specific "private" features which may be needed for reconstruction.

3. **Loss suite.** Training combines embedding alignment and intermediate (hidden-state) alignment for stability; optional reconstruction heads regularize the student.

4. **Evaluation.** We evaluate on WESAD (Schmidt et al., 2018) with ablations isolating the effect of privileged EDA and sensor subsets, demonstrating that EDA-privileged knowledge transfer improves performance without requiring an EDA sensor at test time.

## 2. Related Work

Early applications of learning using privileged information (LUPI) to affect modeling show that models trained with rich laboratory signals can generalize to in-the-wild settings where only a subset of modalities is available. Makantasis et al. propose using privileged physiological and telemetry cues during training while deploying with video-only inputs (Makantasis et al., 2021, 2024).

*EmotionKD* distils across heterogeneous biosignals (EEG and galvanic skin response/EDA) so that a single modality can operate at test time with improved emotion-recognition accuracy (Liu et al., 2023). This aligns with our goal of exploiting a stronger training-time signal to benefit deployable sensors.

Aslam et al. introduce an optimal-transport (OT) formulation to distil the *structural* "dark" knowledge from a multimodal teacher to a lightweight student for expression recognition (Aslam et al., 2024b). They subsequently extend this to *multi-teacher* PKD that aligns diverse teacher representations before distillation (Aslam et al., 2024a), and demonstrate additional gains by leveraging diversity in privileged teacher ensembles (Aslam et al., 2025).

Abbaspourazad et al. distill representational knowledge from a large PPG encoder into an accelerometer encoder trained on population-scale data, yielding generalist accelerometer foundation models that do not require PPG at inference (Abbaspourazad et al., 2025).

Compared to video- or face-centric PKD for affect modeling (Aslam et al., 2024b,a, 2025; Makantasis et al., 2024), our EDA-privileged knowledge transfer uses a frozen EDA teacher to transfer sympathetic-arousal information directly into deployable wrist/chest channels (ECG, BVP, ACC, TEMP). Unlike *EmotionKD* (Liu et al., 2023), we target stress inference with *no EDA sensor* at test time and align only the shared subspace to avoid over-constraining modality-specific features.

## 3. Dataset and Preprocessing

**Source.** We use the publicly available **WESAD** benchmark (Schmidt et al., 2018), which contains 15 subjects wearing a RespiBAN chest band and an Empatica E4 wristband while undergoing a 35-min protocol (baseline $\rightarrow$ social-stress $\rightarrow$ comedy). Signals recorded on the wrist include *electrodermal activity* (EDA, 4 Hz), *photoplethysmography* (BVP, 64 Hz), skin *temperature* (TEMP, 4 Hz), and tri-axial *acceleration* (ACC, 32 Hz); the chest strap provides *ECG*, respiration, EMG and ACC at 700 Hz.

**Resampling.** All channels are resampled to a uniform 64 Hz target rate: first-order polyphase filtering for high-frequency signals (ECG, BVP, ACC), linear interpolation for low-frequency channels (EDA, TEMP).

**Signal cleaning.**

- **EDA** : tonic drift removed with a first-order 0.05 Hz high-pass; phasic band retained with a 0.05–1 Hz Butterworth filter.

- **BVP** : band-pass 0.5–2 Hz to isolate pulse wave.

- **ECG** : band-pass 0.5–40 Hz.

- **Motion** : net acceleration is computed as $\|\mathrm{ACC}_{x,y,z}\|$ from the *wrist* accelerometer only; chest-band ACC is excluded to match wrist-only deployment hardware.

**Segmentation.** The cleaned, synchronized streams are cut into 60-second windows with a 0.25-second stride (96% overlap), mirroring the original WESAD baseline and prior EDA-only studies (Schmidt et al., 2018; Stržinar et al., 2023). Label 0 (*transient*) samples are discarded; only windows whose entire $60s$ lie in a single class $\{1 = \text{baseline}, 2 = \text{stress}, 3 = \text{amusement}\}$ are kept. Windows whose ECG or BVP standard deviation is lower than 0.02 (after $z$-scoring) are rejected to remove sensor drop-outs.

**Normalization.** For each subject, $z$-score parameters are estimated from all baseline ($label = 1$) windows and applied to every channel. This preserves inter-subject differences during training and removes long-term drift within a recording.

**Leave-one-subject-out (LOSO) folds.** The pipeline yields on the order of $8 \times 10^3$ valid windows per subject ($3,840$ samples per window). We build fifteen LOSO folds: in fold $k$, subject $S_k$ is the *test* set and the remaining 14 subjects form the *training* set. Each fold is saved as a compressed `.npz` archive containing:

- `X_train`, `X_test` : input channels ($\mathrm{ECG}, \mathrm{BVP}, \mathrm{TEMP}, \text{net-ACC}_{\mathrm{w}}$);

- `Y_train`, `Y_test` : ground-truth EDA waveform;

- `L_train`, `L_test` : window-level labels;

- per-subject baseline statistics (mean, std) for denormalization.

This segmentation yields on the order of $10^5$ windows across all subjects, sufficient for self-supervised pretraining.

# 4. Methods

Our training pipeline is summarized in Figure 1. Each modality-specific `PhysioMAE` produces *shared* and *private* embeddings; the shared embeddings are aligned across modalities with a hinge loss and then averaged to form a single modality-invariant embedding. This averaged shared embedding, together with the private embeddings, drives reconstruction during pretraining. In the knowledge transfer stage, a frozen EDA `PhysioMAE` serves as the teacher and the student matches its representations. Finally, during fine-tuning and inference, only deployable sensors (ECG, BVP, ACC, TEMP) are used to produce predictions. No EDA is required at test time.

## 4.1. Pretraining setup

Before knowledge transfer, the cheap-sensor `PhysioMAE` models (ECG, BVP, ACC, TEMP) are pretrained jointly in a self-supervised manner, while the EDA `PhysioMAE` is pretrained separately. Inputs of length 3840 samples at 64 Hz are split into non-overlapping patches (size 96), embedded, and processed by a transformer encoder (embedding dimension 1024, depth 8, heads 8). A lightweight decoder (dimension 512, depth 4, heads 8) reconstructs the masked signal. Each cheap-sensor encoder outputs two embeddings: *shared* embeddings (modality-invariant, shown in colors) and *private* embeddings (modality-specific, shown in white in Figure 2). The private/shared embedding separation is implemented with a random private mask. The EDA encoder only outputs private embeddings, since EDA does not undergo alignment with other signals.

**Alignment loss.** Shared embeddings across modalities are aligned using a hinge loss with in-batch negatives. Let $\mathcal{M}$ be the set of input modalities and let $x_i$ denote the input from modality $i$ (with $i \in \mathcal{M}$). Let $\mathcal{P}$ denote the set of positive (matched) cross-modal pairs in the batch, i.e., $\mathcal{P} = \{(i,j) : x_i \text{ and } x_j\}$ are views of the same instance from different modalities. For a positive pair $(i,j) \in \mathcal{P}$, the in-batch negatives $\mathcal{N}(j)$ (resp. $\mathcal{N}(i)$) contain all other items from modality $j$ (resp. $i$) that correspond to different underlying instances than $x_i$ (resp. $x_j$).
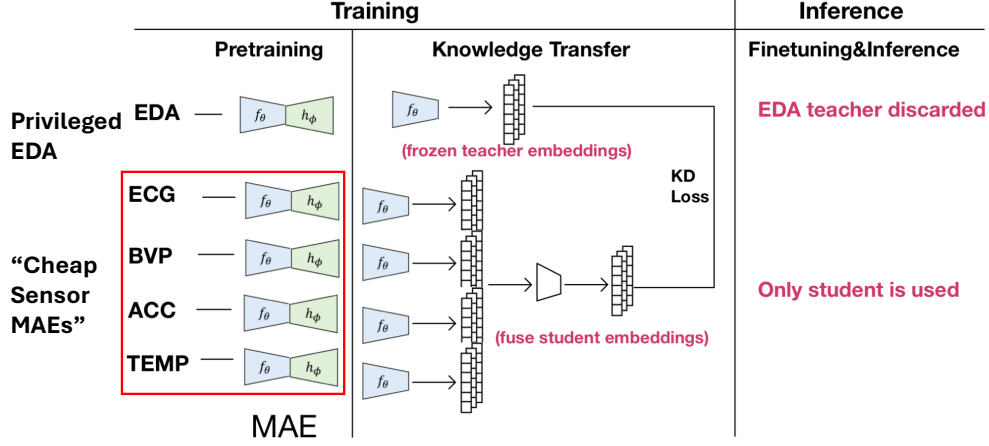
3

Figure 1: **The PULSE Framework.** Our framework uses privileged knowledge transfer from a frozen EDA encoder to students built on low-cost sensors. In the pretraining stage, student encoders learn modality-invariant shared embeddings alongside modality-specific private embeddings. Knowledge transfer is then achieved by aligning the students' shared embeddings with the privileged EDA teacher. Finally, during finetuning, the learned embeddings are used for supervised stress detection without requiring EDA at inference.
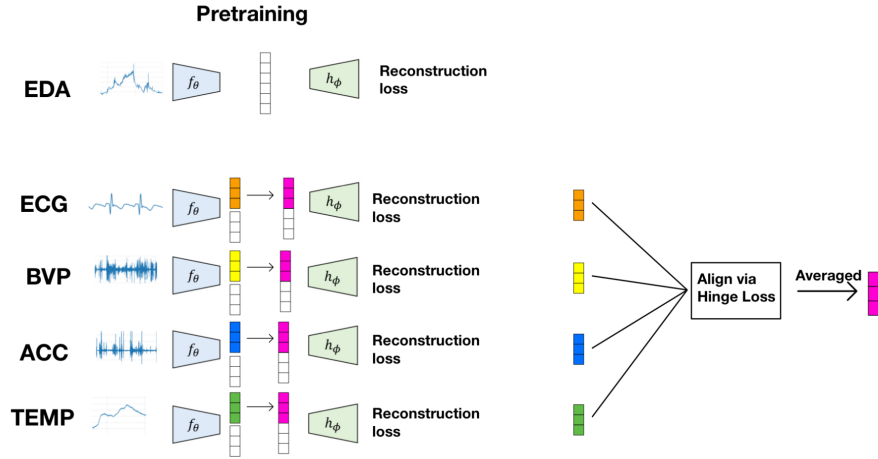


Figure 2: **Pretraining Setup.** In pretraining, each of the cheap sensor encoders outputs shared embeddings (colored boxes) and private embeddings (white boxes). The shared embeddings are aligned across modalities via a hinge loss objective, then averaged into a single shared embedding (magenta boxes). This averaged shared embedding, together with the private embeddings, is fed into the decoder for signal reconstruction. The EDA MAE is trained separately via only reconstruction loss.

We minimize alignment loss as follows.

$$\mathcal{L}_{\text{align}} = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \left[ \frac{1}{|\mathcal{N}(j)|} \sum_{j' \in \mathcal{N}(j)} \max(0, \right.$$
$$\cos\langle s_i, s_{j'}\rangle - \cos\langle s_i, s_j\rangle + \alpha)$$
$$+ \frac{1}{|\mathcal{N}(i)|} \sum_{i' \in \mathcal{N}(i)} \max(0, \tag{1}$$
$$\left. \cos\langle s_{i'}, s_j\rangle - \cos\langle s_i, s_j\rangle + \alpha) \right].$$

where $\alpha$ is the margin (default 0.2). This enforces the matched pair $(i,j)$ to be at least $\alpha$ more similar than any mismatched cross-modal pair in either direction.

**Reconstruction loss.** For reconstruction, the averaged shared embedding and private embeddings are passed through the decoder. With masked patches indexed by $\Omega$, the reconstruction loss is defined as follows.

$$\mathcal{L}_{\text{rec}} = \frac{1}{|\Omega|} \sum_{\omega \in \Omega} \| \underbrace{x_\omega}_{\text{original}} - \underbrace{\hat{x}_\omega}_{\text{reconstruction}} \|^2. \tag{2}$$

In PhysioMAE, 'MAE' denotes Masked AutoEncoder, not mean absolute error; our reconstruction loss is mean squared error (MSE).

**Total loss.** The overall pretraining objective combines both terms is defined as follows.

$$\mathcal{L}_{\text{pre}} = \lambda_{\text{align}} \mathcal{L}_{\text{align}} + \lambda_{\text{rec}} \mathcal{L}_{\text{rec}}. \tag{3}$$

The default weights are $\lambda_{\text{align}}$=1, and $\lambda_{\text{rec}}$=1.

We pretrain each cheap-sensor `PhysioMAE` jointly (EDA separately) for 300 epochs with Adam ($lr = 10^{-4}$, batch size = 128).

## 4.2. Knowledge Transfer.

We freeze an EDA encoder (`PhysioMAE`) pretrained with reconstruction loss as the *teacher* model and distill it into four *student* `PhysioMAE` encoders for ECG, BVP, ACC, and TEMP. All five models share the same ViT-style configuration (signal length 3840 samples at 64 Hz; default patch length 96; encoder embed dimension 1024; depth 8; heads 8; decoder dimension 512, depth 4, heads 8).

**Transfer heads and fusion.** We attach a lightweight transfer head that (i) layer-normalizes student embeddings, (ii) projects each student's hidden embeddings into *shared* and *private* subspaces via per-modality linear layers, and (iii) *fuses* the shared

embeddings across modalities using a single linear fusion layer initialized to the exact average of modalities. The teacher side applies only LayerNorm (no teacher projector). The fusion layer can be frozen for a warm start and optionally unfrozen at a chosen epoch. By default, the fusion layer is unfrozen from the beginning.

**What is transferred.** We *do not* use logits or labeled supervision during knowledge transfer. Instead, we align (a) *hidden tokens* at all student layers against teacher layers and (b) a *final pooled embedding* constructed from the last-layer embeddings. Hidden-embedding knowledge transfer proceeds as follows: (1) take per-block tokens from each student (excluding CLS), (2) project to shared/private spaces, (3) map each student layer to a teacher layer, (4) fuse students' shared tokens across modalities, and (5) minimize a token-wise cosine similarity loss. The final-embedding knowledge transfer fuses last-layer shared embeddings, mean-pools over time, and matches (cosine) to the teacher's last-layer embeddings.

Formally, for each matched layer $\ell$ we optimize

$$\mathcal{L}_{\text{hid}} = \frac{1}{|\mathcal{L}|} \sum_{\ell \in \mathcal{L}} \left( 1 - \cos\langle \underbrace{\texttt{Fuse}(\{S_m^\ell\})}_{\text{students' shared}}, \underbrace{T^\ell}_{\text{teacher}} \rangle \right), \tag{4}$$

and for the final embedding

$$\mathcal{L}_{\text{emb}} = 1 - \cos\langle \text{mean}_t \, \texttt{Fuse}(\{S_m^{\text{final}}\}), \, \text{mean}_t \, T^{\text{final}} \rangle. \tag{5}$$

Function $\text{mean}_t$ stands for averaging over the time dimension.

**Optional regularizers.** There are two optional terms can be toggled: (i) a *decorrelation* penalty between each student's shared and private embeddings using a squared cross-covariance, and (ii) a *reconstruction* loss from each student's MAE decoder with a masking ratio (default 0.5) applied only when enabled. These are added as $\lambda_{\text{perp}} \cdot \mathcal{L}_{\text{perp}}$ and $\lambda_{\text{rec}} \cdot \mathcal{L}_{\text{rec}}$. We discovered that the *reconstruction* loss is essential for avoiding collapse during empirically study.

**Total loss.** The total loss is defined as

$$\mathcal{L} = \lambda_{\text{hid}} \mathcal{L}_{\text{hid}} + \lambda_{\text{emb}} \mathcal{L}_{\text{emb}} + \lambda_{\text{rec}} \mathcal{L}_{\text{rec}} + \lambda_{\text{perp}} \mathcal{L}_{\text{perp}}. \tag{6}$$

The default weights are $\lambda_{\text{hid}}$=1, $\lambda_{\text{emb}}$=1, $\lambda_{\text{rec}}$=0.1, $\lambda_{\text{perp}}$=0.

**Optimization for knowledge transfer.** We train on WESAD windows using a 90/10 train/validation split. We set batch size = 128 and use Adam ($lr = 10^{-4}$) for 100 epochs. An optional cosine scheduler can be enabled. The fusion layer is unfrozen by

default; if frozen for warm start, it can be unfrozen at a specified epoch. Transfer head parameters are included in the optimizer to allow unfreezing.

### 4.3. Finetuning and Inference.

At test time, EDA is absent. The student stack consumes only deployable modalities (ECG, BVP, ACC, TEMP), forms the fused shared representation, and produces predictions via a lightweight classifier; the frozen EDA teacher and all EDA paths are unused. During finetuning, the encoders of each signal modality's `PhysioMAE` are frozen. Only the lightweight classifier is allowed to update.

**Overfitting control and head design.** During finetuning we observed that the densely sampled and highly overlapping windows can lead to rapid overfitting, especially in leave-one-subject-out evaluation. To mitigate this, we *uniformly subsample* the training windows by a factor of 1/40, which (i) reduces correlation between adjacent windows, (ii) acts as an effective regularizer, and (iii) lowers training variance across subjects. In parallel, we use a deliberately small *2-layer MLP* classification head (hidden dimension = 4) on top of the fused embedding, which further limits the model capacity and improves generalization. Unless noted otherwise, encoders are frozen and only the 2-layer MLP head is trained.

Finetuning proceeds as: (1) build per-modality embeddings, (2) mean pool over time, (3) fuse (averaging shared tokens or concatenating shared and private embeddings), and (4) classify with the 2-layer MLP head. We train with the 1/40 downsampling regime and report AUROC/AUPRC as primary metrics; accuracy at a fixed threshold is included as a secondary summary.

**Optimization and schedules for finetuning.** We train each fold for 300 epochs with Adam ($lr = 10^{-3}$), batch size 128, with a cosine learning rate scheduler. Model selection saves the checkpoint with the highest *validation AUPRC*; ties are broken by *validation accuracy* computed at the threshold that maximizes accuracy on the validation set. We do not use early stopping.

## 5. Experiments & Results

We evaluate four different training configurations, as shown in Table 1. This also corresponds to the ablation study (IDs **A**–**D**). In addition, we also evaluate the pretrained EDA-only `PhysioMAE` (ID **E**). We report the aggregate metrics in Table 2, where each row maps to one of the configurations in Table 1. No-EDA and Full-sensor baselines use the same pretraining pipeline, then finetune (none trained from scratch).

**Metric choice (threshold independence).** We treat threshold–independent ranking metrics, *AUROC* and *AUPRC*, as the primary indicators of model quality. These metrics summarize performance across all possible decision thresholds, avoiding the instability that arises from selecting a single operating point.

In contrast, *accuracy* is threshold–dependent. In leave-one-subject-out settings, per-subject score calibration and prevalence can vary markedly, so any fixed threshold (even the one tuned on validation) can yield large swings in accuracy across test subjects. Subject-specific thresholding would artificially inflate results by using test information. For these reasons, we report AUROC/AUPRC as the main metrics and include accuracy only as a secondary summary at a fixed threshold.

**Comparisons to prior WESAD work.** Our objective is to *remove* EDA at inference phase while maintaining EDA as privileged information during training phase. We found no suitable like-for-like baseline in prior WESAD results that uses a comparable *multimodal* input set *excluding* EDA; most reports are either ECG-only or include EDA, and they also differ in class definitions and evaluation protocols (e.g., LOSO vs. subject-dependent). For completeness, representative results are summarized in Appendix A.

**Main observations.** From Table 2, symmetric alignment (**B**) yields only a *minor* improvement over the no–EDA baseline (**A**) when EDA is removed at test time. This indicates that enforcing cross-modal consistency alone does not transfer enough of the EDA-specific information to materially influence downstream stress classification once the privileged channel is absent. In contrast, PULSE (**C**) delivers a *larger and consistent* gain in all evaluation metrics relative to (**A**) and (**B**). Freezing the EDA encoder as a teacher supplies a stable target that shapes the students' shared representation toward EDA's sympathetic-arousal code, which persists even when only cheap sensors are used at inference.

| ID | Name / Description | Train-time inputs | Teacher? | Test-time inputs | Purpose |
|----|----|----|----|----|----|
| **A** | No-EDA baseline | Cheap sensors only (ECG, BVP, ACC, TEMP) | – | Cheap sensors only | Shows what is achievable if EDA is never used. |
| **B** | Symmetric alignment | Cheap sensors + EDA; all encoders updated jointly with masked-reconstruction / alignment losses (no frozen branch) | – | Cheap sensors only | Tests the benefit of standard cross-modal alignment without privileged knowledge transfer. |
| **C** | **PULSE** | Cheap sensors + EDA; EDA encoder *frozen* → students match its embeddings | Yes (frozen EDA teacher) | Cheap sensors only | Measures the additional lift from LUPI-style privileged knowledge transfer. |
| **D** | Full-sensor baseline | Same as **B** | – | Cheap sensors + EDA | Provides an empirical baseline if the EDA sensor is kept at deployment. |
| **E** | EDA-only baseline (teacher, frozen) | Pretrain EDA MAE; finetune classifier with the EDA encoder *frozen* | – | EDA only | Demonstrates that the teacher learns useful, non-collapsed features and isolates the value of privileged transfer vs. direct EDA at inference. |

Table 1: Training/evaluation configurations.

| ID / Name | Test-time inputs | AUROC | AUPRC | Accuracy (%) |
|----|----|----|----|----|
| A/ No–EDA Baseline | Cheap sensors | $0.963 \pm 0.050$ | $0.937 \pm 0.101$ | $91.64 \pm 6.61$ |
| B/ Symmetric Alignment | Cheap sensors | $0.972 \pm 0.031$ | $0.944 \pm 0.061$ | $88.83 \pm 6.24$ |
| C/ **PULSE** | Cheap sensors | $\mathbf{0.994 \pm 0.011}$ | $\mathbf{0.988 \pm 0.022}$ | $\mathbf{96.08 \pm 4.52}$ |
| D/ Full-Sensor Baseline | Cheap sensors + EDA | $0.983 \pm 0.028$ | $0.963 \pm 0.048$ | $90.74 \pm 5.58$ |
| E/ EDA MAE (teacher, frozen at finetune) | EDA | $0.962 \pm 0.067$ | $0.924 \pm 0.122$ | $87.20 \pm 17.69$ |

Table 2: Comparison of training strategies with and without privileged EDA knowledge. (Reported means and standard deviations are computed across folds.)

Table 3: Three-class stress classification (baseline / stress / amusement) under LOSO. Mean $\pm$ sd across 15 subjects. PULSE uses no EDA at inference.

| Model | AUROC | AUPRC | Accuracy (%) |
|----|----|----|----|
| no-EDA baseline | $0.891 \pm 0.090$ | $0.810 \pm 0.130$ | $71.62 \pm 10.27$ |
| symmetric alignment | $0.811 \pm 0.103$ | $0.710 \pm 0.126$ | $69.12 \pm 8.82$ |
| **PULSE** | $\mathbf{0.956 \pm 0.058}$ | $\mathbf{0.894 \pm 0.115}$ | $\mathbf{85.40 \pm 7.18}$ |
| full-sensor baseline | $0.812 \pm 0.128$ | $0.703 \pm 0.124$ | $73.38 \pm 7.89$ |

Contrary to the expectation that full-sensor training yields best performance, we observe that PULSE (**C**) slightly *exceeds* the joint full-sensor model (**D**) on the threshold-independent metrics (AUROC/AUPRC). This outcome is plausible for several reasons: (i) the frozen EDA teacher injects a strong prior and acts as a data-dependent regularizer, constraining students to a modality-invariant subspace that generalizes better on small, cross-subject datasets; (ii) joint optimization in **D** can overfit to idiosyncratic EDA artifacts and/or converge to suboptimal equilibria across branches, whereas **C** benefits from a stationary target and a smoother optimization landscape; (iii) knowledge transfer effectively denoises and compresses the EDA signal into the cheap-sensor representation, while **D** consumes raw EDA at test time and may rely on spurious cues. Taken together, these effects allow **C** to outperform **D** despite not using EDA at inference.

### 5.1. Results on 3-class classification

To further evaluate the generality of PULSE beyond binary stress detection, we extended our experiments to a **three-class classification** setup distinguishing *baseline*, *stress*, and *amusement* conditions. The same leave-one-subject-out (LOSO) protocol and model configurations were used as in the binary case. Metrics are macro-averaged AUROC and AUPRC, with accuracy reported as the mean $\pm$ standard deviation across 15 held-out subjects in Table 3.
**Findings.** PULSE achieves the best performance across all metrics while operating without EDA at inference, confirming that privileged EDA knowledge effectively transfers to the deployable modalities. The gap between PULSE and the no-EDA baseline widens in the multiclass setting, suggesting that fine-grained distinctions between emotional states benefit disproportionately from privileged supervision. These results reinforce that the transferred sympathetic-arousal structure from EDA enriches representations of low-cost physiological signals, improving robustness and separability across multiple affective states.

## 6. Conclusion & Future Work

We present PULSE, an EDA-privileged knowledge transfer framework that separates modality-invariant *shared* and modality-specific *private* embeddings, aligns the shared space with a hinge objective, and uses a frozen EDA encoder as a teacher to transfer EDA-specific dynamics into cheap sensors. On WE-SAD, PULSE achieves the strongest performance, outperforming both the no–EDA baseline and symmetric alignment without a teacher; notably, it can even exceed the joint full-sensor model, consistent with the regularizing effect of a stationary teacher and task-aligned representation compression. These gains are obtained while removing the need for EDA at inference.

In this paper, our experiments focus on WE-SAD with leave-one-subject-out evaluation. We also conducted experiments where ECG serves as the privileged teacher and transfers knowledge into the remaining deployable modalities (BVP, ACC, and TEMP) and found that this setup also yields a gain in threshold-independent metrics (details are shown in Appendix B.4).

In future work, we plan to enhance by: (i) expanding to additional datasets, devices, and stressors to quantify cross-domain generalization, (ii) running systematic ablations on teacher quality, hinge margin, and shared/private capacity, and (iii) conducting a comprehensive calibration and deployment analysis (operating points, thresholds). These studies will strengthen the generalization and robustness claims of the approach.

Moreover, we plan to pursue several directions to further advance this work:

- Evaluate on additional multimodal datasets to assess external validity, starting with *SWELL-KW* Koldijk et al. (2014): (i) matched-modality evaluation (ECG student at test time; EDA used only as a teacher during training) and (ii) cross-dataset transfer (*train on WESAD → test on SWELL-KW* and the reverse), leveraging shared physiology (ECG, EDA).

- Explore richer knowledge transfer objectives (e.g., multi-layer feature matching, CKA/MMD, mutual-information–based losses), adaptive margins, and temperature schedules.

- Study capacity/control: shared/private dimensionality, projector designs, and multi-teacher ensembles (e.g., respiration/EEG as auxiliary teachers).

- Add frequency-domain branches that ingest STFT/CWT features per modality and fuse them with time-domain embeddings (e.g., via concatenation or cross-attention) to capture

rhythmic autonomic patterns and improve noise robustness.

- Personalization and domain adaptation (subject-wise adapters), plus uncertainty estimation to select deployment thresholds post hoc while reporting AUROC/AUPRC.

Overall, PULSE provides a simple, effective path to compress the informative EDA signal into cheaper modalities, improving stress recognition while reducing hardware requirements; we anticipate further gains from stronger teachers, better objectives, and broader validation.

# References

Salar Abbaspourazad, Anshuman Mishra, Joseph Futoma, Andrew C. Miller, and Ian Shapiro. Wearable accelerometer foundation models for health via knowledge distillation. *arXiv preprint arXiv:2412.11276*, 2025. doi: 10.48550/arXiv. 2412.11276. URL https://arxiv.org/abs/2412. 11276.

Muhammad Haseeb Aslam, Marco Pedersoli, Alessandro Lameiras Koerich, and Eric Granger. Multi teacher privileged knowledge distillation for multimodal expression recognition. *arXiv preprint arXiv:2408.09035*, 2024a. URL https://arxiv.org/abs/2408.09035.

Muhammad Haseeb Aslam, Muhammad Osama Zeeshan, Soufiane Belharbi, Marco Pedersoli, Alessandro L. Koerich, Simon Bacon, and Eric Granger. Distilling privileged multimodal information for expression recognition using optimal transport. *arXiv preprint arXiv:2401.15489*, 2024b. URL https://arxiv.org/abs/2401.15489.

Muhammad Haseeb Aslam, Marco Pedersoli, Alessandro L. Koerich, and Eric Granger. Leveraging diversity for privileged multi-teacher knowledge distillation for facial expression recognition. SSRN preprint 5270694, 2025. URL https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5270694.

Behnam Behinaein, Anubhav Bhatti, Dirk Rodenburg, Paul Hungler, and Ali Etemad. A transformer architecture for stress detection from ecg. In *2021 International Symposium on Wearable Computers*, UbiComp '21, page 132–134. ACM, September 2021. doi: 10.1145/3460421. 3480427. URL http://dx.doi.org/10.1145/3460421.3480427.

Wolfram Boucsein. *Electrodermal Activity*. Springer, New York, NY, 2nd edition, 2012.

Hugo D. Critchley. Electrodermal responses: what happens in the brain. *The Neuroscientist*, 8(2): 132–142, 2002. doi: 10.1177/107385840200800209.

Sayandeep Ghosh, SeongKi Kim, Muhammad Fazal Ijaz, Pawan Kumar Singh, and Mufti Mahmud. Classification of mental stress from wearable physiological sensors using image-encoding-based deep neural network. *Biosensors*, 12(12), 2022. ISSN

2079-6374. doi: 10.3390/bios12121153. URL https://www.mdpi.com/2079-6374/12/12/1153.

Md-Billal Hossain, Hugo F. Posada-Quintero, Youngsun Kong, Riley McNaboe, and Ki H. Chon. Automatic motion artifact detection in electrodermal activity data using machine learning. *Biomed. Signal Process. Control*, 71:103176, 2022.

Nabil Ibtehaz and Masood S. Mortazavi. Modally reduced representation learning of multi-lead ECG signals through simultaneous alignment and reconstruction. In *ICLR 2024 Workshop on Learning from Time Series For Health*, 2024. URL https://openreview.net/forum?id=pav9rgzr1c.

Wei-Bang Jiang, Xi Fu, Yi Ding, and Cuntai Guan. PhysioOmni: Towards robust multimodal physiological foundation models with arbitrary missing modalities. *arXiv preprint arXiv:2504.19596*, 2025.

Saskia Koldijk, Maya Sappelli, Suzan Verberne, Mark A. Neerincx, and Wessel Kraaij. The swell knowledge work dataset for stress and user modeling research. In *Proceedings of the 16th International Conference on Multimodal Interaction*, ICMI '14, page 291–298, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450328852. doi: 10.1145/2663204.2663257. URL https://doi.org/10.1145/2663204.2663257.

Ran Liu, Ellen Zippi, Hadi Pour Ansari, Chris Sandino, Jingping Nie, Hanlin Goh, Erdrin Azemi, and Ali Moin. Frequency-aware masked autoencoders for multimodal pretraining on biosignals. In *ICLR Workshop*, 2024. URL https://arxiv.org/abs/2309.05927.

Yucheng Liu, Ziyu Jia, and Haichao Wang. Emotionkd: A cross-modal knowledge distillation framework for emotion recognition based on physiological signals. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*. ACM, 2023. doi: 10.1145/3581783.3612277. URL https://dl.acm.org/doi/10.1145/3581783.3612277.

David Lopez-Paz, Léon Bottou, Bernhard Schölkopf, and Vladimir Vapnik. Unifying distillation and privileged information. In *International Conference on Learning Representations (ICLR)*, 2016. URL http://leon.bottou.org/papers/lopez-paz-2016.

Konstantinos Makantasis, David Melhart, Antonios Liapis, and Georgios N. Yannakakis. Privileged information for modeling affect in the wild. *arXiv preprint arXiv:2107.10552*, 2021. URL https://arxiv.org/abs/2107.10552. Appeared at the 2021 9th Int. Conf. on Affective Computing and Intelligent Interaction (ACII).

Konstantinos Makantasis, Kosmas Pinitas, Antonios Liapis, and Georgios N. Yannakakis. From the lab to the wild: Affect modeling via privileged information. *IEEE Transactions on Affective Computing*, 2024. doi: 10.1109/TAFFC.2023.3265072. URL https://arxiv.org/abs/2305.10919. Early Access 2023.

Konstantin Markov and Tomoko Matsui. Robust speech recognition using generalized distillation framework. In *Proceedings of Interspeech*, pages 2364–2368, 2016. doi: 10.21437/Interspeech.2016-852.

Katie Matton, Robert Lewis, John Guttag, and Rosalind Picard. Contrastive learning of electrodermal activity representations for stress detection. In *Proceedings of Machine Learning Research (PMLR)*, volume 209, pages 411–425, 2023.

Julie Mordacq, Léo Milecki, Maria Vakalopoulou, et al. ADAPT: Multimodal learning for detecting physiological changes under missing modalities. *arXiv preprint arXiv:2407.03836*, 2024.

Eric Oliver and Sagnik Dakshit. Cross-modality investigation on wesad stress classification, 2025. URL https://arxiv.org/abs/2502.18733.

Pooja Prajod and Elisabeth André. On the generalizability of ecg-based stress detection models. In *2022 21st IEEE International Conference on Machine Learning and Applications (ICMLA)*, page 549–554. IEEE, December 2022. doi: 10.1109/icmla55696.2022.00090. URL http://dx.doi.org/10.1109/ICMLA55696.2022.00090.

L. G. Roos and G. M. Slavich. Wearable technologies for health research: Opportunities, limitations, and practical considerations. *Brain, Behavior, and Immunity*, 113:444–452, October 2023. doi: 10.1016/j.bbi.2023.08.008.

Roberto Sánchez-Reolid, María T. López, and Antonio Fernández-Caballero. Machine learning for stress detection from electrodermal activity: A scoping review. *Sensors*, 2025. to appear.

P. Schmidt, A. Reiss, R. Dürichen, C. Marberger, and K. Van Laerhoven. Introducing wesad: A multimodal dataset for wearable stress and affect detection. In *Proc. 20th ACM Int. Conf. on Multimodal Interaction (ICMI '18)*, pages 400–408, Boulder, CO, 2018.

Ghanapriya Singh, Orchid Chetia Phukan, Rinki Gupta, and Anand Nayyar. Hybrid deep learning model for wearable sensor-based stress recognition for internet of medical things (iomt) system. *International Journal of Communication Systems*, 37 (3):e5657, 2024. doi: https://doi.org/10.1002/dac. 5657. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/dac.5657.

Žiga Stržinar, Araceli Sanchis, Agapito Ledezma, Oscar Sipele, Boštjan Pregelj, and Igor Škrjanc. Stress detection using frequency spectrum analysis of wrist-measured electrodermal activity. *Sensors*, 23(2):963, 2023. doi: 10.3390/s23020963.

Vladimir Vapnik and Rauf Izmailov. Learning using privileged information: Similarity control and knowledge transfer. *Journal of Machine Learning Research*, 16:2023–2049, 2015. URL https://jmlr.org/papers/v16/vapnik15b.html.

Shiqi Yang, Yuan Gao, Yao Zhu, Li Zhang, Qinlan Xie, Xuesong Lu, Fang Wang, and Zhushanying Zhang. A deep learning approach to stress recognition through multimodal physiological signal image transformation. *Scientific Reports*, 15, 07 2025. doi: 10.1038/s41598-025-01228-3.

## Appendix A. Comparison with some existing work

We show comparison of our results with some existing work in Table 4. Most WESAD reports either (i) rely on a single sensor (often ECG) or (ii) include EDA at inference, and many use task definitions and splits that differ from ours (e.g., multi-class vs. binary; subject-dependent vs. LOSO). Consequently, these numbers provide useful context but are not like-for-like baselines for our setting where EDA is absent at test time.

## Appendix B. Hyperparameter and Ablation Studies

We report ablations on (i) modality dropout during pretraining, (ii) private-mask ratio controlling shared vs. private capacity, and (iii) hinge-loss margin used for shared-embedding alignment. Unless noted, metrics are mean $\pm$ std over folds.

### B.1. Modality Dropout in Finetuning

We assess the sensitivity of PULSE to missing test-time modalities using the configurations in Table 5. These results are obtained by removing signal modalities during finetuning and inference. Relative to the full setting (ECG+BVP+ACC+TEMP), removing TEMP yields small changes. Removing ACC in addition (ECG+BVP) produces the largest additional drop. Going to ECG only leads to a further but modest decline.

Variability increases as modalities are removed (e.g., AUPRC std: $0.033 \rightarrow 0.036 \rightarrow 0.088 \rightarrow 0.117$), indicating reduced stability under stronger input constraints. Overall, PULSE degrades gracefully and maintains strong performance even with ECG only, consistent with the goal of transferring EDA-informed arousal structure to cheaper sensors used at inference.

### B.2. Modality Dropout in Pretraining

From Table 6 we can see that pretraining with all four cheap sensors is best; dropping modalities degrades AUPRC and accuracy most sharply, indicating reduced transfer of privileged EDA structure. Compared with Table 5, removing modalities during pretraining degrades performance more than removing

Table 4: Selected WESAD results from prior work, augmented with our results. Metrics are reported as in the original papers and are *not* directly comparable across differing task setups (binary vs. 3/4-class) and evaluation protocols.

| Paper (Year) | Test-time sensors | Eval protocol | Accuracy (%) | AUROC |
|---|---|---|---|---|
| **Ours (PULSE, 2025)** | ECG + BVP + ACC + TEMP | LOSO, binary | 93.97±5.77 | 0.989 ± 0.017 |
| **Ours (Symmetric alignment, 2025)** | ECG + BVP + ACC + TEMP | LOSO, binary | 88.83±6.24 | 0.972 ± 0.031 |
| **Ours (No-EDA baseline, 2025)** | ECG + BVP + ACC + TEMP | LOSO, binary | 91.64±6.61 | 0.963 ± 0.050 |
| **Ours (Full-sensor baseline, 2025)** | ECG + BVP + ACC + TEMP + EDA | LOSO, binary | 90.74±5.58 | 0.983 ± 0.028 |
| Prajod and André (2022) | ECG only | LOSO, binary | 90.8 | N/A |
| Behinaein et al. (2021) | ECG only | LOSO, binary | 91.1 | N/A |
| Oliver and Dakshit (2025) | ECG / BVP / EDA / RESP / TEMP (single-modality) | Subject-dependent (random 85:15) | up to 99.95 (validation; multiclass) | N/A |
| Yang et al. (2025) | EDA + ACC + BVP + RESP + TEMP | LOSO, 3-class | 90.96 | N/A |
| Singh et al. (2024) | ECG + EDA + RESP (+TEMP/ACC) | LOSO, 3-class | 90.45 | N/A |
| Ghosh et al. (2022) | EDA + ECG + ACC + RESP + TEMP | LOSO, 4-class | 94.77 | N/A |

| Name | Test-time inputs | AUROC | AUPRC | Accuracy (%) |
|---|---|---|---|---|
| PULSE (full) | ECG + BVP + ACC + TEMP | **0.989 ± 0.017** | **0.977 ± 0.033** | **93.97 ± 5.77** |
| PULSE (–TEMP) | ECG + BVP + ACC | 0.985 ± 0.021 | 0.972 ± 0.036 | 90.85 ± 6.53 |
| PULSE (–TEMP, –ACC) | ECG + BVP | 0.966 ± 0.072 | 0.948 ± 0.088 | 87.74 ± 8.53 |
| PULSE (ECG only) | ECG | 0.960 ± 0.082 | 0.936 ± 0.117 | 86.22 ± 10.37 |

Table 5: Ablations on test-time input availability for PULSE. (Reported means and standard deviations are computed across folds.)

Table 6: Effect of removing modalities during pretraining. Full-sensor pretraining yields the strongest downstream performance.

| Setup | AUROC | AUPRC | Accuracy |
|---|---|---|---|
| PULSE (ECG + BVP + ACC + TEMP) | **0.989 ± 0.017** | **0.977 ± 0.033** | **93.97 ± 5.77 %** |
| ECG + BVP + ACC | 0.949 ± 0.084 | 0.923 ± 0.117 | 86.12 ± 8.22 % |
| ECG + BVP | 0.973 ± 0.048 | 0.942 ± 0.112 | 90.44 ± 9.92 % |
| Only ECG | 0.960 ± 0.077 | 0.927 ± 0.137 | 75.74 ± 10.81 % |

them at test time, underscoring the benefits of multimodal representation learning over single-modality training.

### B.3. Hidden-state matching ablations

To assess the effect of hidden-state supervision during knowledge distillation, we performed ablations over which layers of the student were matched to the frozen EDA teacher. Earlier experiments used mid-level layers (3, 5, 7) based on empirical stability, but we now set **all-layer matching** as the default configuration. Table 7 summarizes the results.

We compare three configurations: (i) removing hidden-state matching entirely ($\mathcal{L}_{\text{hid}} = 0$), equivalent to final-representation–only transfer; (ii) matching only intermediate layers (3, 5, 7) as in the original setup; and (iii) matching all hidden layers. For reference, the no-EDA baseline undergoes the same pretraining and finetuning pipeline but without knowledge transfer.

Removing hidden-state matching degrades all metrics relative to both PULSE and the no-EDA baseline, indicating that intermediate feature supervision is critical for transferring sympathetic-arousal structure. Matching only a subset of layers recovers much of the lost performance, but supervising *all* layers yields the highest AUROC, AUPRC, and accuracy. This suggests that multi-depth guidance stabilizes optimization and enables richer cross-modal transfer, supporting the adoption of all-layer matching as the new default configuration.

### B.4. ECG as privileged teacher

To evaluate whether privileged knowledge transfer in PULSE depends specifically on EDA or can generalize to other physiological modalities, we conducted an ablation using **ECG as the privileged teacher**. In this setup, the pretrained ECG encoder is frozen and used to distill knowledge into the remaining deployable modalities (BVP, ACC, and TEMP) through the same shared/private alignment and hidden-state matching objectives. Table 8 reports the results under the LOSO protocol.

Using ECG as the teacher yields moderate improvements in AUROC and AUPRC compared to the no-KD counterpart, while accuracy remains comparable within variance. This demonstrates that the PULSE framework is *modality-agnostic*: privileged distillation benefits are not limited to EDA, though the magnitude of gain depends on the informativeness of the teacher modality. EDA remains the strongest privileged signal for stress-related supervision due to its direct coupling with sympathetic arousal, yet ECG-based distillation still provides measurable transfer, confirming that PULSE can flexibly leverage alternative privileged channels when EDA is unavailable.

### B.5. Private Mask Ratio

Table 9 shows that both insufficient private capacity (ratio 0; no private embeddings to encode modality-specific features) and excessive private capacity (ratio 1; no shared embeddings for alignment and knowledge transfer) hurt performance. A *balanced* split between shared and private capacity (0.5) most effectively supports privileged EDA transfer, while a smaller private fraction (0.2) trades a slight drop in AUROC/AUPRC for the highest accuracy.

### B.6. Hinge-Loss Margin

Table 10 indicates that a modest margin (0.2) provides the strongest signal for aligning student shared embeddings with the EDA teacher; margins that are too small (0) or too large ($\geq 0.6$) underperform.

Table 7: Effect of hidden-state matching during KD.

| Model | AUROC | AUPRC | Accuracy (%) |
|---|---|---|---|
| PULSE (match layers 3,5,7) | $0.989 \pm 0.017$ | $0.977 \pm 0.033$ | $93.97 \pm 5.77$ |
| PULSE ($\mathcal{L}_{\text{hid}}$=0; final-only) | $0.953 \pm 0.096$ | $0.922 \pm 0.155$ | $90.95 \pm 9.73$ |
| **PULSE (match all layers)** | $\mathbf{0.994 \pm 0.011}$ | $\mathbf{0.988 \pm 0.022}$ | $\mathbf{96.08 \pm 4.52}$ |
| no-EDA baseline (no KD) | $0.963 \pm 0.050$ | $0.937 \pm 0.101$ | $91.64 \pm 6.61$ |

Table 8: ECG as frozen teacher distilled to BVP/ACC/TEMP students.

| Setup | AUROC | AUPRC | Accuracy (%) |
|---|---|---|---|
| ECG teacher $\rightarrow$ (BVP, ACC, TEMP) | $\mathbf{0.956 \pm 0.067}$ | $\mathbf{0.930 \pm 0.085}$ | $86.63 \pm 5.86$ |
| (BVP, ACC, TEMP) w/o KD | $0.949 \pm 0.058$ | $0.902 \pm 0.102$ | $\mathbf{87.06 \pm 8.77}$ |

## B.7. Adaptive fusion vs. static averaging

We further examined whether learning per-modality fusion weights could improve performance compared to static averaging of modality embeddings. In the **adaptive fusion** variant, a small trainable gating module predicts fusion weights for each modality during both knowledge distillation and finetuning, initialized to uniform values and updated end-to-end. The default PULSE configuration instead uses a **static average** of shared embeddings across modalities. Results are summarized in Table 11. Note that for this experiment, we used the setup where only layers 3, 5, and 7 are matched during knowledge transfer (The first row in Table 7.)

Adaptive fusion offered no consistent improvement and slightly degraded performance across AUROC, AUPRC, and accuracy. We attribute this to the high inter-subject variability under LOSO evaluation: fusion weights learned on training subjects do not necessarily generalize to unseen participants. Static averaging, by contrast, enforces a more stable modality combination and avoids overfitting to subject-specific patterns. Consequently, we retain static fusion as the default configuration for all reported results.

## B.8. Effect of alignment loss

To quantify the contribution of the cross-modal alignment objective, we conducted an ablation removing the alignment loss $\mathcal{L}_{\text{align}}$ during pretraining. In this setting, each modality's encoder is trained solely via its reconstruction objective (MAE) without enforcing consistency across shared embeddings. Table 12 summarizes the results for the no-EDA baseline.

Removing the alignment loss leads to a sharp drop in both ranking metrics and accuracy, confirming that explicit cross-modal regularization is essential for learning a coherent shared latent space. Without this constraint, modalities drift apart during pretraining, producing representations that are less transferable and less robust for downstream stress classification. The results validate $\mathcal{L}_{\text{align}}$ as a key component in maintaining cross-modal consistency and enabling effective privileged knowledge transfer.

## Appendix C. Why hinge (vs. anchored contrastive)

During pretraining we first tried an **anchored contrastive** objective: each modality's shared embeddings were aligned to ECG (the "anchor") via cosine similarity. In practice, the alignment loss collapsed to 0 within the first few iterations and then remained nearly constant, providing little training signal thereafter, as shown in Figure 3. Although we did not complete a full performance evaluation under this objective, this early collapse motivated switching to a **hinge-based alignment** with a small positive margin, which maintained non-trivial gradients throughout training, as we have shown in the appendix.

## Appendix D. Evidence that reconstruction prevents collapse

When we train without reconstruction (KD-only), the shared embeddings collapse toward a constant

Table 9: Varying the private-mask ratio that reserves capacity for modality-specific features. A balanced split (0.5) maximizes AUROC/AUPRC, while a slightly smaller private region (0.2) peaks accuracy. Removing knowledge transfer (ratio 1) hurts all metrics.

| Private mask ratio | AUROC | AUPRC | Accuracy |
|---|---|---|---|
| 0 | $0.959 \pm 0.062$ | $0.927 \pm 0.100$ | $70.41 \pm 8.60\,\%$ |
| 0.2 | $0.980 \pm 0.044$ | $0.966 \pm 0.076$ | $\mathbf{94.46 \pm 7.78}\,\%$ |
| 0.4 | $0.919 \pm 0.116$ | $0.894 \pm 0.148$ | $79.51 \pm 13.92\,\%$ |
| 0.5 (default) | $\mathbf{0.989 \pm 0.017}$ | $\mathbf{0.977 \pm 0.033}$ | $93.97 \pm 5.77\,\%$ |
| 0.6 | $0.967 \pm 0.072$ | $0.956 \pm 0.084$ | $92.62 \pm 9.89\,\%$ |
| 0.8 | $0.975 \pm 0.047$ | $0.958 \pm 0.082$ | $93.76 \pm 5.96\,\%$ |
| 1 (no knowledge transfer) | $0.945 \pm 0.078$ | $0.906 \pm 0.127$ | $70.41 \pm 8.60\,\%$ |

Table 10: Effect of hinge-loss margin for shared-embedding alignment. A small positive margin (0.2) is optimal; too large a margin degrades performance, likely by over-separating positive pairs and harming transfer.

| Margin | AUROC | AUPRC | Accuracy |
|---|---|---|---|
| 0.0 | $0.972 \pm 0.039$ | $0.954 \pm 0.060$ | $70.41 \pm 8.60\,\%$ |
| 0.2 (default) | $\mathbf{0.989 \pm 0.017}$ | $\mathbf{0.977 \pm 0.033}$ | $\mathbf{93.97 \pm 5.77}\,\%$ |
| 0.4 | $0.977 \pm 0.051$ | $0.966 \pm 0.071$ | $89.08 \pm 13.39\,\%$ |
| 0.6 | $0.983 \pm 0.052$ | $0.975 \pm 0.070$ | $77.57 \pm 8.85\,\%$ |
| 0.8 | $0.957 \pm 0.077$ | $0.938 \pm 0.112$ | $70.41 \pm 8.60\,\%$ |
| 1.0 | $0.969 \pm 0.070$ | $0.933 \pm 0.160$ | $76.43 \pm 8.69\,\%$ |

vector across modalities: the mean pairwise cosine $\approx 1.0$ for ECG/BVP/ACC (TEMP 0.998), and the feature variance is near-zero, indicating almost no dispersion. After adding the reconstruction objective, feature variance returns by several orders of magnitude, and mean pairwise cosine drops to 0.027–0.137, showing that samples are no longer trivially aligned, as shown in Figure 4. This pattern holds consistently across modalities, supporting our claim that reconstruction acts as an information-preserving regularizer: it counteracts the KD objective's tendency to minimize alignment by shrinking representation variance, thereby preventing representational collapse and yielding meaningful geometry for downstream LOSO generalization.

## Appendix E. Training Dynamics

In Figures 5, 6, and 7, we show the convergence behavior in different stages of training (pretraining, knowledge transfer, and finetuning).

Table 11: Fusion strategy during finetuning.

| Model | AUROC | AUPRC | Accuracy (%) |
|---|---|---|---|
| PULSE (static fusion; default) | **0.989 ± 0.017** | **0.977 ± 0.033** | **93.97 ± 5.77** |
| PULSE (adaptive fusion) | 0.967 ± 0.046 | 0.951 ± 0.064 | 90.91 ± 7.35 |

Table 12: Ablation on alignment loss $\mathcal{L}_{\text{align}}$ during pretraining. Removing alignment substantially degrades performance, indicating that shared-space regularization is critical for multimodal representation learning.

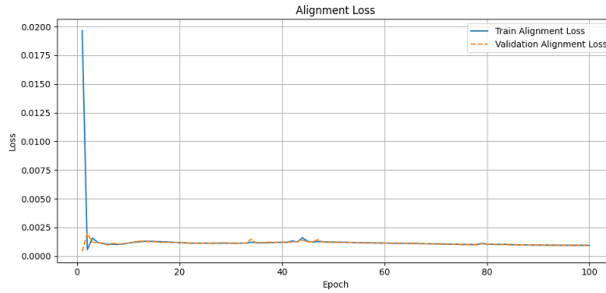| Model | AUROC | AUPRC | Accuracy (%) |
|---|---|---|---|
| no-EDA baseline (with alignment loss) | **0.963 ± 0.050** | **0.937 ± 0.101** | **91.64 ± 6.61** |
| no-EDA baseline (without alignment loss) | 0.945 ± 0.078 | 0.906 ± 0.127 | 70.41 ± 8.60 |



Figure 3: Learning curve showing anchored contrastive loss in early experiments

| Modality | mean_pairwise_cosine | mean_feature_variance |
|---|---|---|
| ECG | 1.000 | 2.57e–05 |
| BVP | 1.000 | 6.26e–05 |
| ACC | 1.000 | 7.53e–05 |
| TEMP | 0.998 | 7.96e–04 |
| | Added reconstruction | |
| Modality | mean_pairwise_cosine | mean_feature_variance |
| ECG | 0.070 | 2.52e–02 |
| BVP | 0.0645 | 1.51e–02 |
| ACC | 0.137 | 9.20e–02 |
| TEMP | 0.0273 | 2.06e–03 |

Figure 4: Effect of adding reconstruction during KD. Without reconstruction, shared embeddings collapse (cosine $\approx 1.0$, near-zero variance). Adding reconstruction restores variance by several orders of magnitude, preventing collapse and enabling meaningful shared geometry.
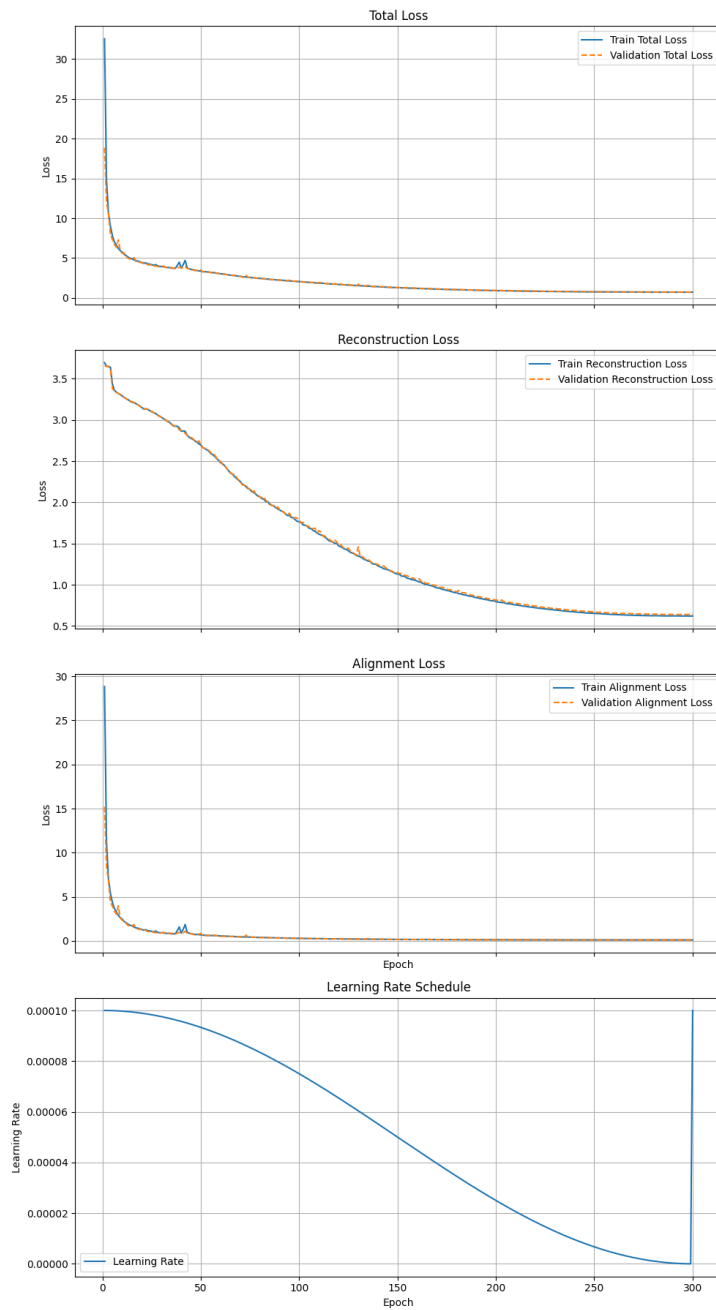
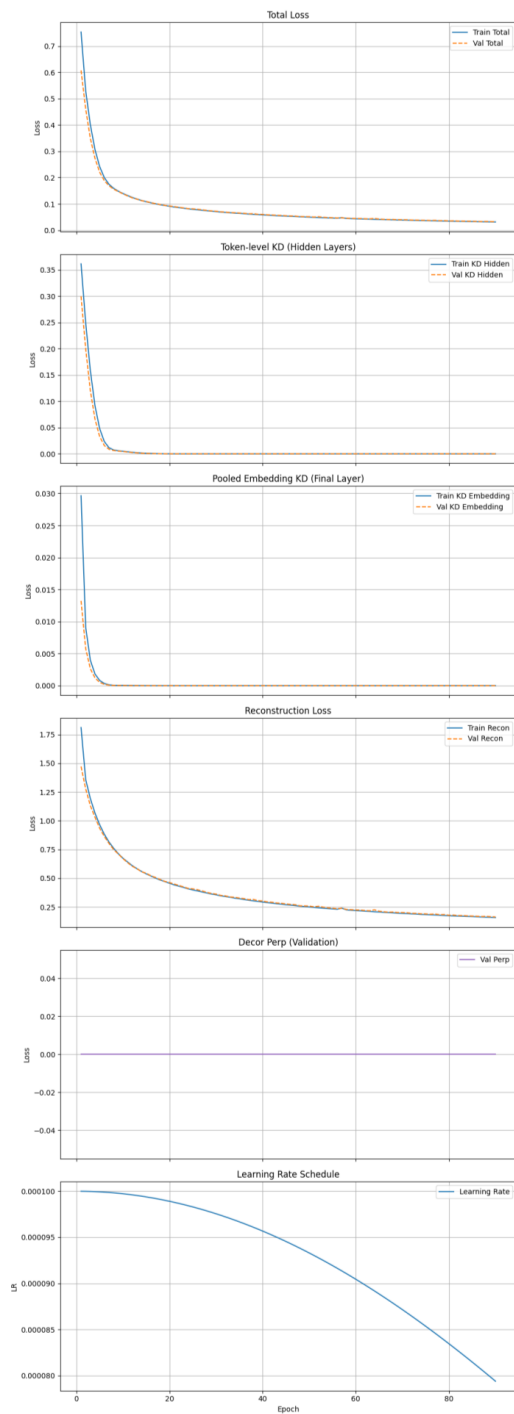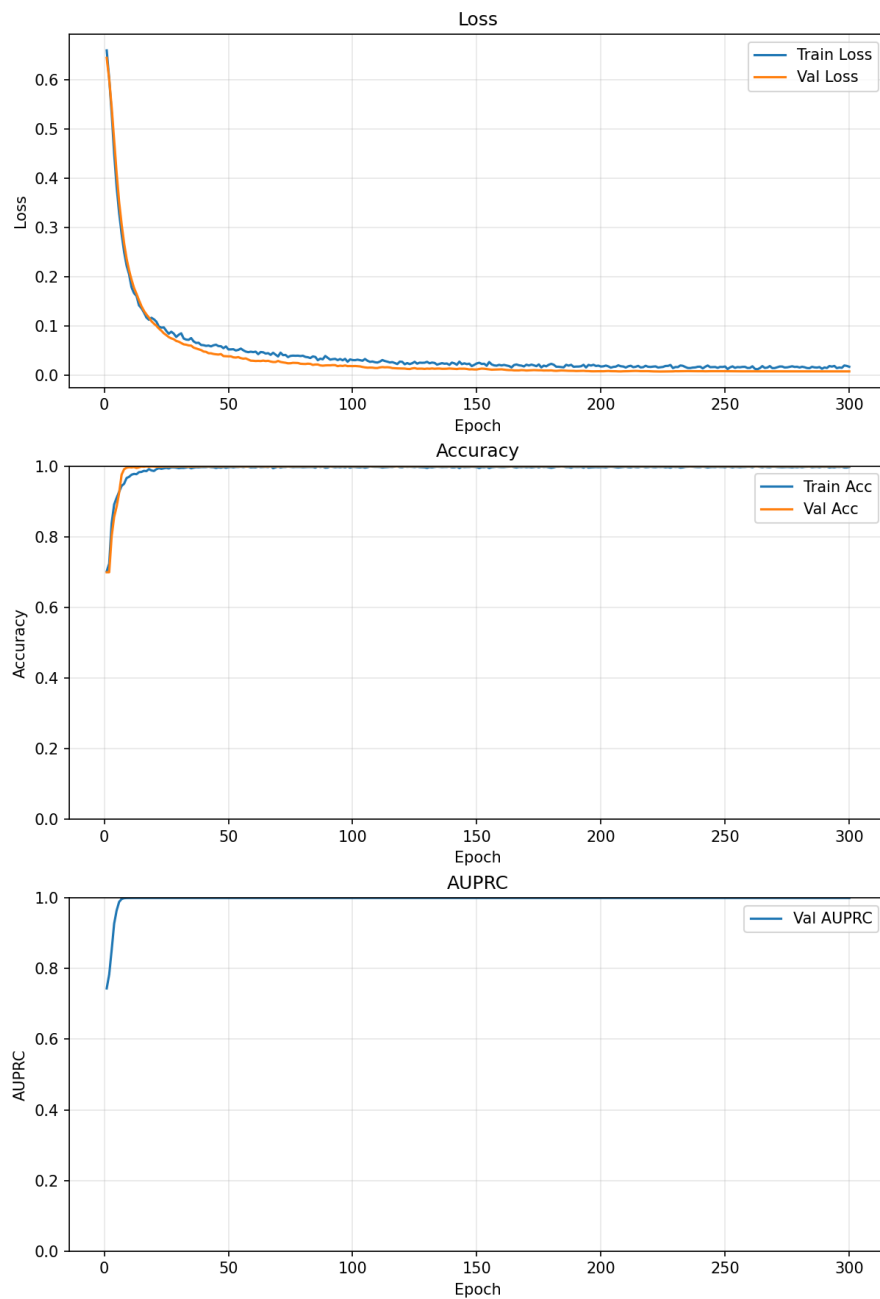Figure 5: Convergence behavior in pretraining

Figure 6: Convergence behavior in knowledge transfer

Figure 7: Convergence behavior in finetuning