

HiMAE: HIERARCHICAL MASKED AUTOENCODERS DISCOVER RESOLUTION-SPECIFIC STRUCTURE IN WEARABLE TIME SERIES DOWNSTREAM TASKS

Simon A. Lee^{1,2,*}, Cyrus Tanade¹, Hao Zhou¹, Juhyeon Lee¹, Megha Thukral¹, Minji Han¹, Md Sazzad Hissain Khan¹, Keum San Chun¹, Baiying Lu¹, Migyeong Gwak¹, Mehrab¹, Morshed, Viswam Nathan¹, Md Mahbubur Rahman¹, Li Zhu¹, Subramaniam Venkatraman¹, Sharanya Arcot Desai¹,

¹Digital Health Team, Samsung Research America

²Department of Computational Medicine, UCLA

* Work completed during AI Residency

simonlee711@g.ucla.edu

ABSTRACT

Wearable sensors provide abundant physiological time series observations, yet the resolution at which we should extract features for downstream tasks remain unclear. We hypothesize that temporal resolution is a fundamental axis of representation learning, with different clinical and behavioral outcomes relying on features at distinct scales. To test this resolution hypothesis, we introduce HiMAE (Hierarchical Masked Autoencoder), a self-supervised framework that combines masked autoencoding with a hierarchical convolutional encoder-decoder. HiMAE produces multi-resolution embeddings across its intermediate layers that enable systematic evaluation of which temporal scales carry predictive signal, transforming resolution from a hyperparameter into a probe for interpretability. Across classification and generative benchmarks, HiMAE consistently outperforms state-of-the-art foundation models that collapse scale, while being orders of magnitude smaller. Due to the convolution based design choices behind HiMAE, the model is also compact enough to run entirely on-device, achieving sub-millisecond inference on smartwatch-class CPUs for true edge inference. Together, these contributions position HiMAE as both an efficient self supervised learning method and a discovery tool for understanding how time resolution contributes to downstream task alignment.

1 INTRODUCTION

Wearable sensors have emerged as a primary modality for continuous health monitoring, providing access to rich physiological and behavioral signals in free-living settings (Erturk et al., 2025). Despite their ubiquity, the utility of wearable signals for machine learning in healthcare remains poorly understood. Unlike images (Dosovitskiy et al., 2021; Simonyan et al., 2014; Zhou et al., 2015; Petsiuk et al., 2018) or text (Brown et al., 2020; Li et al., 2016; Sundararajan et al., 2017; Arras et al., 2017), physiological time series rarely admit obvious visual cues that map cleanly to clinical outcomes, leaving open fundamental questions about which features carry predictive value. A particularly unresolved issue concerns temporal resolution: should models operate at a single universal resolution, or do different health outcomes depend on resolution-specific structure? Clinically actionable events can arise on second-level timescales, requiring representations that both capture fine-grained temporal patterns and support real-time inference under the computational constraints of wearable devices. We hypothesize that resolution is not a nuisance parameter but a fundamental axis of physiological representation learning. We refer to this as the *resolution hypothesis*, which posits that temporal granularity governs predictive performance in clinical and behavioral tasks. In this framing, “resolution” denotes the effective temporal context over which representations are formed—from fine-scale waveform morphology to coarse-scale dynamics spanning the whole sequence.

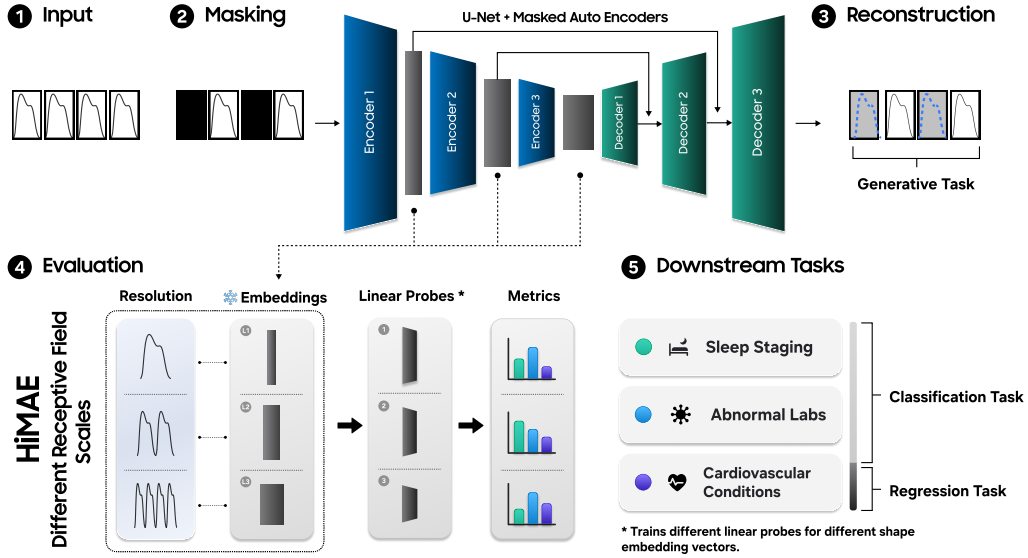


Figure 1: **HiMAE pre-training and evaluation pipeline.** (1) Physiological sequences are split into temporal patches. (2) Selected patches are masked randomly or contiguously. (3) A U-Net-style CNN encoder-decoder reconstructs missing values, with loss applied only to masked regions. (4) Multi-resolution embeddings feed linear probes for classification and regression benchmarking. (5) Three categorized task-lists are evaluated.

From an algorithmic perspective, much of the field defaults to transformer-based architectures (Vaswani et al., 2017), implicitly assuming that flexibility and capacity outweigh inductive bias. Yet wearable signals, while long in sequence length, are often generated by a few latent processes driven by biological mechanisms and captured through only a handful of sensor modalities. In this sense they are low-dimensional and highly structured. This raises the possibility that transformers may not only overfit but also obscure resolution-specific structure, rather than expose it. By contrast, hierarchical convolutional biases offer a natural mechanism for aligning architectures with the resolution hypothesis, capturing both local detail and long-range dependencies in a principled way. This motivates a re-examination of architectural design choices for self-supervised learning (SSL) on raw physiological time series.

In this work, we address these challenges by introducing *HiMAE* (Hierarchical Masked Autoencoder), a self-supervised pretraining framework for wearable time series that directly operationalizes the *resolution hypothesis* (Figure 1). HiMAE combines the masked autoencoding paradigm with 1D physiological signals by coupling patch-masking objectives (Wang et al., 2023) with a U-Net-inspired encoder-decoder (Ronneberger et al., 2015). Crucially, HiMAE produces multi-resolution embeddings, with each level of the hierarchy corresponding to a distinct temporal granularity. This design enables systematic interrogation of which resolutions carry predictive signal, while simultaneously yielding lightweight, efficient representations. Beyond its architectural advantages, HiMAE allows us to benchmark the resolution hypothesis across 14 classification. Our results reveal resolution-specific structure in wearable signals that is not readily identifiable by human experts, offering new insights into both representation learning and the interpretability of physiological time series in the time domain.

2 RELATED WORK

Self-Supervised Pretraining Objectives for Wearable Signals Wearable devices equipped with photoplethysmography (PPG), electrocardiography (ECG), and accelerometry generate long, multi-channel time series encoding diverse physiological and behavioral phenomena, including cardiovascular dynamics (Castaneda et al., 2018), activity patterns (Yuan et al., 2024; Xu et al., 2025), sleep cycles (Li et al., 2021; Thapa et al., 2024; Logacjov et al., 2025), and other latent processes. These data streams are abundant, and predominantly unlabeled, making them well suited for large-scale

self-supervised learning (Kaplan et al., 2020; Bommasani et al., 2021; Zhou et al., 2024; Liang et al., 2024).

SSL has become the dominant paradigm for wearable time-series representation learning, given the scarcity of labeled data and the ubiquity of unlabeled signals in free-living settings (Lee & Akamatsu, 2025). Among SSL strategies, masked autoencoding has emerged as a central approach, inspired by its success in vision (He et al., 2022; Vaid et al., 2023) and language modeling (Devlin et al., 2019). The method randomly occludes patches of the signal and tasks the model with reconstructing them, encouraging representations that capture latent physiological structure and temporal regularities (Zhang et al., 2022a; Kong et al., 2023). Recent large-scale efforts, most notably Google’s LSM series (Narayanswamy et al., 2024; Xu et al., 2025), rely heavily on masked autoencoding, establishing it as a pretraining standard for multi-modal wearable datasets. Yet despite its effectiveness for local pattern recovery, vanilla masked autoencoding often struggles to capture multi-resolution features unless coupled with explicitly hierarchical architectures.

In parallel, contrastive learning enforces invariance by pulling semantically similar samples together in latent space while pushing dissimilar ones apart (Schmitt & Kuljanin, 2008; Jaiswal et al., 2020). The central challenge for wearables is defining positive and negative pairs without labels. One solution is participant-level contrastive training, where samples from the same individual are positives and samples from different individuals are negatives, an approach adopted in Apple’s ECG and PPG foundation models (Abbaspourazad et al., 2023) and closely related to the SimCLR framework (Chen et al., 2020b). Other domain-specific innovations define pairs through physiological priors: PaPaGei leverages PPG morphology (Pillai et al., 2024), while SleepFM extends the paradigm across EEG, ECG, and EMG to enforce cross-modal consistency (Thapa et al., 2024). Additional embedding-level regularizers, such as differential entropy constraints (Jing et al., 2021; Abbaspourazad et al., 2023), further enrich learned representations. However, contrastive methods are highly sensitive to augmentation heuristics (which are rarely physiologically meaningful), computationally intensive, and limited in interpretability, providing little insight into which temporal structures are preserved.

HiMAE departs from both flat masked and contrastive approaches in two ways. First, instead of relying on a single-scale reconstruction or augmentation heuristics, HiMAE couples masked autoencoding with a hierarchical encoder–decoder that integrates information across resolutions, treating temporal scale as an explicit dimension of representation. Second, by extracting embeddings at multiple scales and probing them independently, HiMAE transforms SSL from a pretraining mechanism into a discovery tool: it directly tests which temporal resolutions carry predictive signal for downstream tasks. In doing so, HiMAE preserves the efficiency of masked autoencoding while introducing interpretability absent in contrastive or flat masked objectives.

Multi-scale Learning The emphasis on resolution awareness connects naturally to multi-scale learning, where modeling temporal signals across multiple granularities has emerged as a powerful inductive bias. In vision, multi-scale architectures such as pyramidal CNNs and hierarchical attention enable models to integrate fine-scale edges with coarse semantic structures, substantially improving recognition and generation in 2D (Wang et al., 2016; Yang et al., 2016; Liu et al., 2021a; Kusupati et al., 2024; Liu et al., 2024) and 3D (He et al., 2017; Ghadai et al., 2019; Zhang et al., 2022b).

In time series, multi-scale methods are fewer but increasingly influential. N-HiTS (Challu et al., 2022) improves long-horizon forecasting by allocating capacity across frequencies via hierarchical interpolation. Pyraformer (Liu et al., 2022) leverages pyramidal attention to capture dependencies over a tree of scales, while Scaleformer (Shabani et al., 2023) introduces iterative refinement across resolutions. Pathformer (Chen et al., 2024) further adapts pathways dynamically to match input-specific temporal dynamics.

Prior multi-scale methods typically rely on fixed hierarchies or task-specific refinement stages (e.g., for forecasting), which constrains their generality. While HiMAE also inherits inductive biases from convolutional design choices (e.g., step size, padding, kernel width), these parameters define receptive fields rather than dictate which scales are salient. By coupling self-supervised reconstruction with these fields, HiMAE induces a hierarchy of temporal embeddings that can be probed independently.

3 METHODS

Hierarchical Masked Autoencoders (HiMAE) HiMAE combines masked autoencoding (Baldi, 2012; He et al., 2022) with 1-D physiological time series by coupling a patch-masking objective with a U-Net-style convolutional encoder-decoder (Ronneberger et al., 2015). Given an input sequence $x \in \mathbb{R}^{C \times L}$, we partition it into $N = L/P$ non-overlapping patches of length P . A binary mask $m \in \{0, 1\}^N$ is sampled from a Bernoulli distribution with parameter r , indicating the masking ratio. Masked indices are selected uniformly at random without replacement, expanded to match temporal resolution as $m' \in \{0, 1\}^L$, and applied to the sequence, yielding $\tilde{x} = x \odot (1 - m')$. This masking procedure removes substantial context, forcing the model to infer higher-order dependencies. In addition to random masking, we also employ contiguous masking, in which adjacent patches are removed to mimic sensor dropout similar to recent protocols showing benefits (Xu et al., 2025). Both regimes are interleaved during pretraining to promote robustness across reconstruction settings.

Architecture The encoder f_θ is a hierarchical 1D CNN composed of residual convolutional blocks with stride-2 convolutions that downsample the temporal resolution by half at each stage, expanding the receptive field so that deeper layers capture long-range dependencies while shallow layers retain local detail. Each residual block consists of two convolutions with kernel size 5, batch normalization (Ioffe & Szegedy, 2015), and GELU activations (Hendrycks & Gimpel, 2023), along with a projection shortcut when input and output dimensions differ. The decoder g_ϕ mirrors this structure with transposed convolutions for upsampling and incorporates skip connections from encoder layers, concatenating intermediate features to restore fine-grained temporal structure. All convolutions are standard 1D operations defined over temporal windows, and striding handles subsampling directly. Intermediate activations use GELU, while the final layer applies a tanh nonlinearity so that outputs $\hat{x} \in \mathbb{R}^{C \times L}$ are bounded in $[-1, 1]$, matching the normalized input range.

We deliberately adopt a convolutional U-Net backbone rather than a transformer-based encoder for two reasons. First, physiological signals exhibit strong local dependencies governed by morphology (e.g., PPG waveform shape, ECG peaks), which are naturally modeled by finite receptive fields. Convolutions (O’Shea & Nash, 2015) encode this locality directly, whereas transformers must simulate it through restricted attention, often at higher parameter cost. Second, multi-resolution structure is intrinsic to physiology (e.g., heartbeats unfold over milliseconds, rhythms span seconds). A hierarchical CNN with skip connections provides an architectural bias toward such nested timescales, aligning directly with the resolution hypothesis and being orders of magnitude smaller than other proposed foundation models in this space (See Figure 2 for comparison). In contrast, transformers emphasize global mixing, which may obscure resolution-specific structure while consuming substantially more compute (Table 6). This rationale motivates HiMAE’s design as not only efficient but also inductively aligned with the temporal statistics of wearable signals.

Resolution Probes Multi-resolution embeddings extracted from different levels of the hierarchy are probed independently, with distinct linear classifiers trained per resolution (Alain & Bengio, 2018). This design enables us to systematically evaluate which temporal granularity carries predictive signal for downstream tasks, rather than collapsing embeddings into a single latent space. Finally, choices of patch length P and kernel size were guided by ablations (Appendix Section F.1), which confirmed that $P = 5$ and kernel size 5 yield the best balance between local fidelity and receptive field expansion when all other hyperparameters were fixed.

Wearable Foundation Models

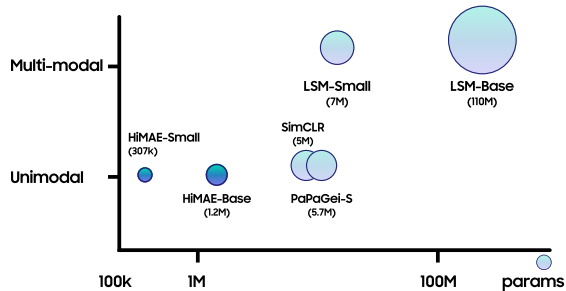


Figure 2: **HiMAE is lightweight**

Training minimizes a masked reconstruction loss restricted to occluded regions: $\mathcal{L}_{\text{MSE}}(\theta, \phi) = \frac{\|(\tilde{x} - x) \odot m'\|_2^2}{\sum_{t=1}^L m'_t}$, where m' ensures that gradients are only computed on masked segments. This objective estimates $p(x_{\mathcal{M}} | x_{\mathcal{O}})$, with \mathcal{M} and \mathcal{O} denoting masked and observed indices, preventing trivial copying of visible inputs and promoting temporally coherent, multi-scale representations.

Effective Global Context via Receptive Field Expansion While Transformers achieve global dependency modeling via $O(L^2)$ self-attention, the U-Net architecture in HiMAE approximates this behavior at $O(L)$ complexity through hierarchical spatial contraction. In a D -layer encoder, the effective receptive field (ERF) at layer d grows exponentially as $R_d = R_{d-1} + (k-1) \cdot \prod_{i=1}^{d-1} s_i$, where k is the kernel size and s is the stride. By the bottleneck, the ERF encompasses a significant portion of the input sequence L , allowing the model to capture “global” context without the quadratic memory overhead of an attention matrix. This hierarchical aggregation acts as a multi-scale proxy for global attention: deep layers integrate coarse, long-range context, while skip connections inject high-resolution local features back into the decoder. Consequently, HiMAE simulates the communicative benefits of attention through a series of local-to-global inductive biases, achieving competitive representation power at a fraction of the FLOPs required by vanilla ViT or Transformer-based autoencoders.

Pretraining and Evaluation Protocol PPG Sequences were sampled at $f_s = 100$ Hz over fixed windows of $T = 10$ s ($L = 1000$ timesteps). 10 second windows were selected due to clinically actionable events occurring in these time scales (ECG is collected at 10s intervals in clinical settings (Shuai et al., 2016; Elgendi, 2012)) and due to our interest in real-time monitoring on edge devices. Each signal was divided into non-overlapping patches of length $P = 5$ (200 patches total), and a masking ratio $r = 0.8$ was applied with patterns resampled per sequence and iteration to mitigate overfitting (we empirically tested this masking ratio in Appendix Section F.1 with similar observations made in (Narayanswamy et al., 2024)). The encoder architecture employed channel widths [16, 32, 64, 128], mirrored in the decoder. Optimization was performed with AdamW (Loshchilov & Hutter, 2019) ($\text{lr} = 10^{-3}$, weight decay = 10^{-3}) using a warmup–cosine schedule (10% linear warmup steps followed by cosine decay). Models trained up to 100k steps with batch size 2048 and early stopping triggered after 3 epochs without improvement similar to the protocols found in (Narayanswamy et al.). Data splits followed a 90/10 (train/validation) protocol across subjects, ensuring no identity overlap between pretraining and validation. Pretraining converged within 12 hours when distributing training across 4 Tesla T4 GPUs using PyTorch lightning (Paszke et al., 2019).

Pretraining datasets. We construct our pretraining corpus from approximately 80,000 hours of wearable green PPG signals, drawn from seven large-scale free world studies conducted at Samsung Research. These datasets include recordings from 47,644 participants across seven distinct wearable devices, capturing broad demographic, behavioral, and hardware variability in a noisy environment (See Appendix Section B for ethics considerations). Although our modeling framework is modality-agnostic and can extend to other physiological signals such as electrocardiograms (see Appendix F.2), we focus here on PPG due to its prevalence and the scale of available data (we lack the same order of magnitude of ECG compared to PPG because ECG is not passively collected). To ensure reliability, we apply a standardized preprocessing pipeline that retains only high-quality segments, filtering by a Signal Quality Index (SQI). The retained signals are further refined using a bandpass filter of 0.5–8 Hz (Christiano & Fitzgerald, 2003), consistent across all pretraining and evaluation studies, to isolate physiologically relevant dynamics. Finally, signals are normalized to the range $[-1, 1]$ to match the output range of the tanh activation function used in our models.

4 EXPERIMENTAL DESIGN

We follow the evaluation protocol of Narayanswamy et al. (2024) and extend it into a unified benchmark suite spanning generative, and classification, along with ablations to quantify how key architectural components interact with scaling. Across all experiments, our goal is not only to assess HiMAE’s efficiency and transferability, but also to test the *resolution hypothesis*: whether predictive signal concentrates at specific levels of the hierarchical embeddings. Further analysis and results are displayed in full in Appendix Section F.

Model scaling and generative reconstruction. We first study HiMAE’s scaling properties by measuring how reconstruction performance varies as a function of dataset size, number of participants, model capacity, and training compute capacity (batch size). For each axis, we systematically sub-sample or expand the relevant resource while holding others fixed, enabling us to isolate its contribution to representation quality. Scaling is assessed through mean squared error on masked reconstruction on a held out validation set, which provides a direct measure of how model capacity and data availability govern loss reduction. We also squeeze in ablations in this experiment to assess how removing skip connections, and removing the hierarchal design affect scaling.

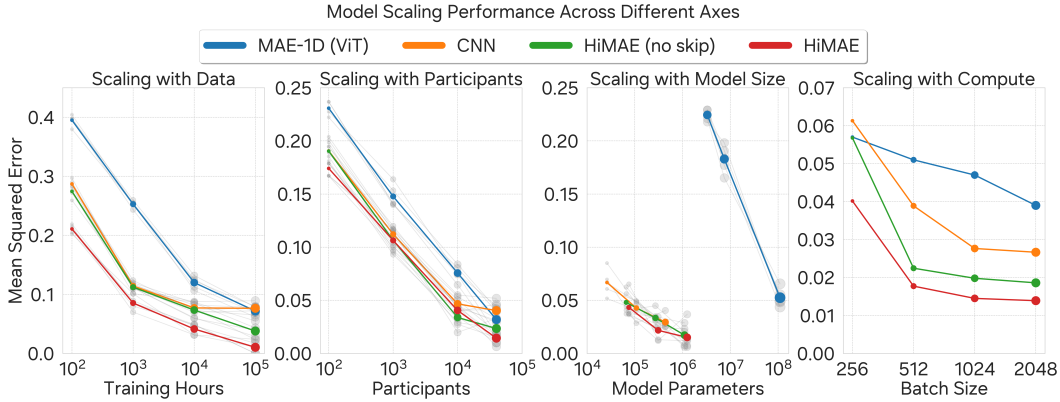


Figure 3: **HiMAE exhibits superior scaling across axes.** Mean squared error decreases most rapidly for HiMAE as data, participants, model size, and compute scale on a held out validation set. Ablations without skip connections confirm that both the hierarchical design and skip pathways are helpful for generative performance. Grey lines indicate multiple runs whereas colored lines are average performance.

To complement this aggregate view, we also evaluate generative performance under three increasingly challenging reconstruction regimes defined in the LSM papers (Narayanswamy et al.; Xu et al., 2025): (i) random imputation, where patches are masked at random uniformly; (ii) temporal interpolation, where contiguous spans are removed to simulate sensor dropout; and (iii) temporal extrapolation, where future spans are occluded and predictions must rely solely on past context. We compute the mean squared error (MSE) for these evaluations.

Classification To assess downstream transferability and adaptability, we benchmark HiMAE on 12 binary classification tasks drawn from labeled datasets fully disjoint from our pretraining sources. We organize these into three groups: cardiovascular outcomes, sleep staging, and abnormal laboratory prediction. Cardiovascular outcomes, provide the most established benchmarks, with well-documented links between PPG and clinical endpoints (Shabaan et al., 2020). These include hypertension detection, and arrhythmia-related events such as Premature Ventricular Contractions (PVCs) detections, typically identified via electrocardiograms (ECGs). Sleep staging is another task we include which is of high interest, given the demand for wearables to track fine-grained sleep states despite the temporal and physiological complexity of the task (Imtiaz, 2021; Thapa et al., 2024; Birrer et al., 2024). Laboratory predictions, on the other hand, serves as a discovery setting, testing whether PPG contains sufficient biomarker information to separate abnormal from healthy labs—an open question compared to patient-record benchmarks where such signals are more explicit (Kolo et al., 2024; McDermott et al., 2025). Together, these canonical and exploratory tasks form a spectrum that enables a comprehensive evaluation of representation quality across diverse digital health applications. All tasks are described in greater detail in Appendix Section D.

We evaluate HiMAE against two complementary classes of baselines. The first comprises established self-supervised representation learning methods for time series, including SimCLR (Chen et al., 2020b), DINO (Caron et al., 2021), Masked Siamese Networks (MSN) (Assran et al., 2022), and a ViT-based 1D masked autoencoder that follows the LSM training protocol of (Narayanswamy et al.). The second class consists of state-of-the-art time-series and wearable foundation models. This includes PaPaGei, a leading foundation model for PPG (Pillai et al., 2024), evaluated both using its publicly released Bell Labs checkpoint (PaPaGei-BL)¹ and a variant retrained on our pre-training corpus (PaPaGei-SRA). We additionally benchmark against Chronos (Ansari et al., 2024), a large-scale time-series foundation model, and the hierarchical Swin Transformer (Liu et al., 2021b), configured to match the LSM setting for controlled comparison. Further implementation details for all baselines are provided in Appendix E. All models are evaluated using a standard linear probing protocol, in which the pretrained encoder is frozen and a linear classifier is trained on top of the learned representations. Performance is reported using AUROC as the primary metric of discriminative ability. For every architecture, we expose the full sequence of embeddings along the temporal

¹<https://zenodo.org/records/13983110>

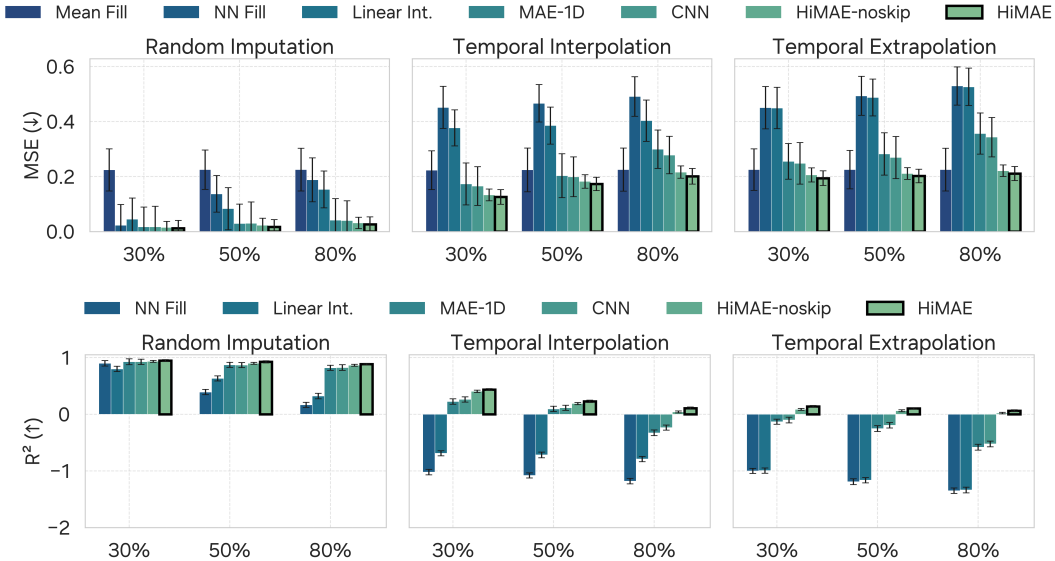


Figure 4: **Performance on generative benchmarks.** Mean squared error and r^2 for random imputation, temporal interpolation, and temporal extrapolation at varying missingness levels. Bold outline indicates best performing model.

dimension, rather than collapsing representations to a single summary token, ensuring that downstream probes retain access to resolution-specific information. This evaluation protocol allows us to assess whether pretraining yields representations that are both discriminative and transferable across tasks.

Resolution Hypothesis HiMAE produces embeddings at multiple temporal scales, and we probe each scale independently with linear classifiers. This allows us to test whether predictive information is concentrated at fine, intermediate, or coarse resolutions depending on the clinical endpoint. In this way, the classification tasks serve not only as benchmarks for transfer learning, but also as controlled tests of the resolution hypothesis (Receptive field lengths are described in Section C.1).

5 RESULTS

5.1 SCALING AND GENERATIVE BENCHMARK

Scaling: We first examine the scaling behavior in Figure 3 of HiMAE relative to baselines across data, participants, model parameters, and compute capacity (batch size). The overall scaling trends follow conventional expectations, error decreases monotonically with additional data, participants, or compute. However, scaling with model parameters reveals an interesting insight. HiMAE achieves substantially lower loss at smaller parameter capacities, while transformers only begin to close the gap once scaled to orders of magnitude more parameters (we chose transformer parameter count based on LSM’s original paper (Narayanswamy et al., 2024)). This difference reflects an inductive bias. Transformer which assume global receptive fields, appear to require considerably larger model capacity before capturing the local dynamics of the data. In contrast, HiMAE’s hierarchical convolutional structure exploits spatial and temporal locality efficiently, yielding superior performance at modest scales. This observation reinforces the importance of architectural priors in low-capacity regimes.

Generative: Turning to generative benchmarks, HiMAE consistently outperforms all baselines across random imputation, temporal interpolation, and temporal extrapolation tasks (Table 4). In terms of mean squared error, HiMAE achieves the lowest reconstruction error in every setting, including cases with heavy missingness. This advantage persists when evaluated with R^2 , where the mean-fill baseline serves as the reference. By achieving positive R^2 scores even in challenging extrapolation scenarios, HiMAE demonstrates reconstruction ability beyond naive heuristics (e.g., mean fill, nearest neighbor, or linear interpolation). Together, these results establish HiMAE as a strong generative model for missing data problems, with advantages that persist across scaling regimes and input corruption patterns.

Table 1: Linear probing classification performance comparison against baselines on different tasks. AUROC is reported in percent with 95% confidence intervals. The best performance is **bold**, the second best model is underscored. * denotes $p < 0.05$, ** denotes $p < 0.01$ from a two-sided z-test comparing HiMAE with the second-best model.

Model	#param (M)	Cardiovascular Conditions			Abnormal Blood Labs				Sleep Staging				
		Hyptn (lab)	Hyptn (free-living)	PVC	A1C	Hemoglobin	Platelets	Sodium	Potassium	Wake	Light	Deep	REM
SimCLR	5.0	53.4 (± 3.6)	53.7 (± 4.1)	51.7 (± 5.9)	60.9 (± 5.5)	50.8 (± 4.6)	44.4 (± 6.2)	58.6 (± 4.8)	67.0 (± 4.9)	56.6 (± 4.3)	52.7 (± 5.1)	67.0 (± 4.0)	51.0 (± 5.7)
DINO	6.5	51.7 (± 4.8)	52.2 (± 3.4)	47.0 (± 4.8)	58.9 (± 3.9)	49.6 (± 3.2)	42.9 (± 5.4)	56.9 (± 3.3)	64.5 (± 3.8)	55.3 (± 3.6)	55.2 (± 4.4)	68.8 (± 3.3)	46.0 (± 6.0)
MSN	2.5	45.2 (± 2.8)	55.2 (± 2.5)	62.9 (± 3.0)	62.9 (± 2.2)	52.1 (± 2.4)	45.9 (± 3.7)	60.4 (± 2.1)	69.5 (± 2.3)	57.8 (± 2.7)	50.3 (± 2.9)	65.3 (± 2.8)	56.0 (± 3.5)
MAE (ViT)	110.6	43.2 (± 7.0)	65.0 (± 5.5)	72.2 (± 7.0)	79.6 (± 6.5)	57.6 (± 4.9)	56.1 (± 5.8)	48.8 (± 6.1)	76.5 (± 6.8)	63.8 (± 5.3)	60.8 (± 5.8)	69.3 (± 6.6)	59.7 (± 6.2)
HiMAE	1.2	65.1** (± 1.7)	65.1 (± 1.6)	80.2* (± 1.4)	70.1 (± 2.0)	56.2 (± 1.3)	68.5* (± 1.8)	63.3* (± 1.9)	83.1 (± 1.5)	66.8 (± 1.8)	59.3 (± 2.1)	72.3 (± 1.4)	58.2 (± 2.2)

Table 2: Linear probing classification performance comparison against state-of-the-art wearable and time-series foundation models. AUROC is reported in percent with 95% confidence intervals. The best performance is **bold**, the second best model is underscored. * denotes $p < 0.05$, ** denotes $p < 0.01$ from a two-sided z-test comparing HiMAE with the second-best model.

Model	#param (M)	Cardiovascular Conditions			Abnormal Blood Labs				Sleep Staging				
		Hyptn (lab)	Hyptn (free-living)	PVC	A1C	Hemoglobin	Platelets	Sodium	Potassium	Wake	Light	Deep	REM
PaPaGei-BL	5.7	57.3 (± 4.7)	60.9 (± 4.1)	74.2 (± 6.4)	59.2 (± 5.8)	58.5 (± 5.2)	59.9 (± 4.9)	59.0 (± 4.3)	75.5 (± 5.5)	56.8 (± 4.9)	55.6 (± 5.0)	53.9 (± 4.5)	56.3 (± 5.7)
PaPaGei-SRA	5.7	59.3 (± 3.5)	<u>62.9</u> (± 3.7)	<u>75.2</u> (± 5.6)	<u>61.2</u> (± 4.1)	60.5 (± 2.7)	<u>61.9</u> (± 3.6)	<u>61.0</u> (± 3.3)	<u>77.5</u> (± 4.6)	56.8 (± 4.2)	57.6 (± 3.7)	55.9 (± 3.4)	<u>58.3</u> (± 5.1)
Swin-Transformer	110.6	58.3 (± 6.2)	61.9 (± 5.8)	74.2 (± 7.2)	58.2 (± 6.7)	<u>59.5</u> (± 8.0)	60.9 (± 7.1)	60.0 (± 5.6)	76.5 (± 6.8)	56.7 (± 5.4)	54.7 (± 6.1)	53.3 (± 7.5)	54.4 (± 5.3)
Chronos	200.0	67.3 (± 1.7)	59.9 (± 2.9)	65.7 (± 3.1)	58.2 (± 3.4)	53.6 (± 3.3)	60.9 (± 2.7)	63.3 (± 2.3)	63.5 (± 2.8)	<u>64.9</u> (± 2.7)	63.2 (± 2.1)	<u>72.2</u> (± 2.5)	57.3 (± 2.9)
HiMAE	1.2	<u>65.1</u> (± 1.7)	65.3 (± 1.6)	80.2 (± 1.4)	70.1** (± 2.0)	56.2 (± 1.3)	68.5** (± 1.8)	63.3 (± 1.9)	83.1* (± 1.5)	66.8 (± 1.8)	<u>59.2</u> (± 2.1)	72.3 (± 1.4)	58.5 (± 2.2)

Ablations: Ablation in Figure 3 and 4 further highlights the contributions of hierarchical design and skip connections in HiMAE. Removing either component results in increased error, indicating that both are crucial for effective representation learning. Nevertheless, even without these architectural elements, HiMAE variants remain competitive with larger transformer based models, underscoring the robustness of the approach. More importantly, the full model exhibits improved generalization across scaling axes (Appendix Section F.3), suggesting that the combination of hierarchy and skip connections facilitates better transfer as data and compute grow.

5.2 CLASSIFICATION BENCHMARKING

Classification In Tables 1 and 2, HiMAE consistently secures the majority of wins, frequently outperforming or matching models that are considerably larger. This is particularly striking given that prior work has typically relied on heavy architectures to reach similar levels of performance, highlighting HiMAE’s ability to capture a broad spectrum of physiological features with a compact design. These outcomes emphasize the model’s robustness when applied to structured, temporally dependent problems that demand sensitivity to subtle variations in wearable signals.

Taken together, these results position HiMAE as the most consistently strong performer across the benchmark suite. In cases where HiMAE does not place first it is only ~ 1 -2% behind the winning model. Crucially, this level of performance is achieved with a substantially smaller model than competing approaches, demonstrating a favorable tradeoff between efficiency and predictive power. Rather than excelling only in isolated cases, HiMAE delivers broad, cross-domain competitiveness, suggesting that compact models, when designed with the right inductive biases, can rival or even surpass far larger architectures.

5.3 RESOLUTION SPECIFIC CLINICAL INTERPRETABILITY

The resolution hypothesis predicts that different health outcomes depend on distinct temporal granularities. To test this, we analyze performance across HiMAE layers, where each layer corresponds to a progressively coarser resolution. Figure 5 reveals clear resolution-specific structure: individual downstream tasks achieve maximal AUROC at different layers, highlighted by the red boundaries.

This layer-task alignment underscores two key insights. First, temporal resolution is not a nuisance parameter but an axis of predictive structure: different outcomes are best represented at different scales (we show that collapsing an encoder decoder still has concordant results showing that our hierarchal model is not an artifact in Appendix Section F.4). Second, HiMAE naturally exposes this heterogeneity, functioning as a discovery tool for identifying the most informative resolution per task. This complements conventional interpretability methods (Amann et al., 2022; Xu et al., 2023; Lee et al., 2025) by shifting the focus from *which features* drive predictions to *which resolutions* matter. In doing so, HiMAE operationalizes the resolution hypothesis and provides insights to tasks where the resolution needed is not entirely clear.

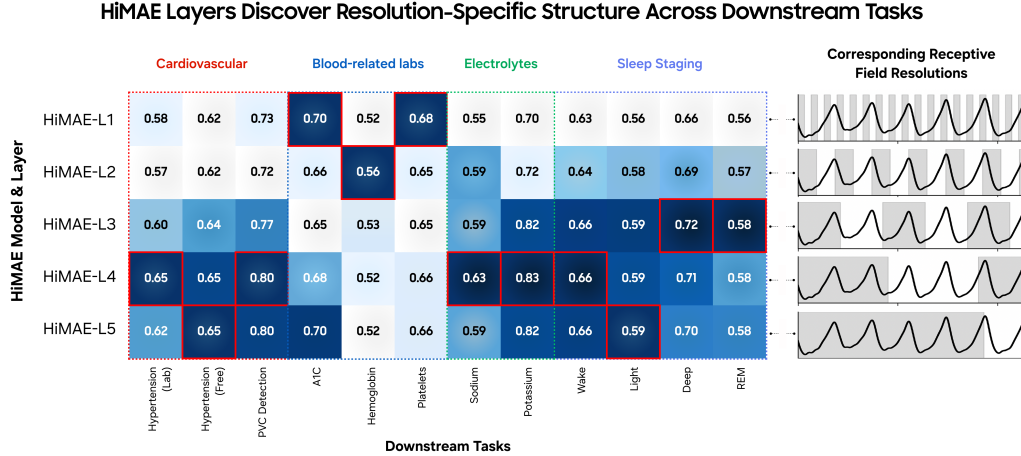
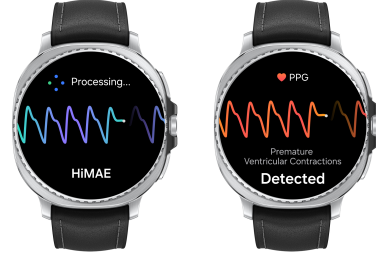


Figure 5: **HiMAE discovers task-specific structures for downstream tasks.** AUROC across layers shows that tasks rely on distinct temporal scales, highlighting HiMAE as a tool for discovering the most informative resolution in clinical machine learning.

Clinical Interpretation

The resolution-specific structure discovered by HiMAE carries clinical implications that resonate with existing literature or provide new clinical insights from a scientific discovery perspective. Cardiovascular and sleep-staging outcomes achieve maximal performance in later layers, which aligns with physiological understanding: cardiovascular (Zhou et al., 2021; Tang et al., 2025) and sleep dynamics (Patanaiik et al., 2018), which evolve gradually over longer time horizons. Capturing these slower patterns requires less temporal granularity, consistent with the notion that general trends, not transient spikes, dominate predictive structure in chronic or cyclic physiological processes.

In contrast, tasks involving laboratory measurements, are a more exploratory and scientific discovery setting where there isn’t much intuition on what resolution should reveal the most predictive signal. When looking at Figure 5, particularly the blood-related biomarkers, it exhibit optimal performance in earlier layers corresponding to finer temporal scales. These outcomes reflect inherently volatile physiological processes, where shifts in morphology can signal meaningful physiological change on lab measurements.



Model	Params (↓)	FLOPs (↓)	Memory (↓)	On-device Lat. (↓)
HiMAE	1.2M	0.0647 gFLOPs	4.8 MB	0.99 ms
Efficient-Net B-1	7.8M	0.70 gFLOPs	31.1 MB	1.42 ms
Swin-Transformer	110.6M	11.89 gFLOPs	423.8 MB	2.95 ms
MAE-1D	110.6M	15.94 gFLOPs	441.3 MB	3.36 ms

5.4 CASE STUDIES

Case Study 1: On-Device Benchmarking A central novelty of HiMAE is that it is, to our knowledge, the first SSL method compact enough to run entirely *on-watch*, rather than on phone-class hardware. We evaluate on-device PVC detection on smartwatch-class CPUs sampled at 100 Hz (Figure 6). HiMAE is exceptionally lightweight (1.2M parameters, 0.0647 gFLOPs, 4.8 MB) and achieves 0.99 ms latency per sample, equivalent to processing $\approx 1,010$ samples/s or ≈ 2.8 hours of signal per minute of wall time. By contrast it shows massive performance gains against transformer baselines, Swin-Transformer (110M parameters, 11.9 gFLOPs, 423 MB) and a MAE-1D (ViT) (110M, 15.9 gFLOPs, 441 MB). HiMAE also outperforms optimized models like Efficient-Net B1 (Tan & Le, 2020) providing context to the latency and compactness of our model. HiMAE is thus $\sim 3\text{--}4\times$ more efficient compared to trans-

Figure 6: **Model efficiency and on-device inference:** Sample on-device detections on Samsung Galaxy device. Size, compute cost, memory footprint, and CPU latency (ms per sample, batch size 2048) measured over a 10s sequence at 100Hz.

formers while fitting fully on-watch (without quantization (Jacob et al., 2017)), enabling continuous, private inference at the point of signal collection. *This prototype is strictly for research and is not deployed commercially.*

Case Study 2: HiMAE is adaptable in few shot settings

A central challenge in the wearable domain is that labels are scarce across tasks. Models that can adapt quickly from generic pretraining to specific detection tasks with limited supervision are therefore essential. Figure 7 illustrates this setting: HiMAE provides strong representations that can be adapted to diverse tasks such as PVC detection or hypertension monitoring with only a handful of labeled examples as reflected by the shape of the learning curves on the few-shot learning experiments. By reducing the supervision required to reach high performance, HiMAE enables new tasks to be supported on-device without the prohibitive cost of large curated datasets which help bolster its practical utility.

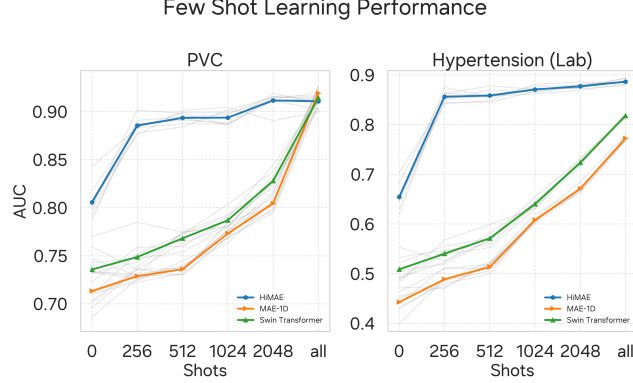


Figure 7: **Few-shot adaptation.** HiMAE adapts efficiently to new wearable tasks under sparse labels indicated by curve shape over transformer baselines.

6 DISCUSSION

Summary. HiMAE advances wearable self supervised methods along three dimensions: (i) its flexible architecture is expressly designed for multi-resolution mapping, enabling seamless adaptation across heterogeneous tasks, (ii) by aligning task-dependent resolutions with model representations, it not only optimizes predictive performance but also offers a window into the temporal organization of physiological biomarkers, and (iii) by design of the compactness, it achieves the first demonstration of true *on-watch* inference, running entirely within smartwatch-class constraints while matching or surpassing performance on far larger models. These results position HiMAE as an efficient representation learner but also as a framework for interrogating which temporal resolutions carry signal.

Resolution as a structural prior. Our findings validate the resolution hypothesis and suggest a shift in how representation learning on wearables should be conceptualized. This reframing implies that representation learning for physiological signals should expose, rather than collapse, scale-specific embeddings. The layer-wise AUROC profiles in Figure 5 show that predictive performance peaks at different levels of the hierarchy depending on the task, with fine-scale embeddings capturing short-lived physiological events and coarse-scale embeddings capturing slower behavioral phenomena. By revealing this heterogeneity, HiMAE provides empirical evidence that resolution-specific representations are essential for wearable health modeling.

From “on-device” to “on-watch.” HiMAE demonstrates that convolutional hierarchies can reduce model size by two orders of magnitude relative to transformer-based models, enabling the first instance of true *on-watch* inference. This moves the deployment frontier from phone-class to watch-class processors, where inference occurs exactly at the point of sensing. Beyond efficiency, this shift has consequences for privacy (data never leave the device) and for clinical viability (continuous real-time monitoring becomes feasible).

Limitations and Future Works While we focus on PPG, the principles underlying HiMAE generalize to multimodal settings. Physiological signals are inherently multi-scale across modalities (e.g., ECG beats, accelerometer motion cycles, EEG rhythms), and resolution-aware architectures could expose complementary temporal signatures across them. Another limitation of our work is we don’t handle sequences beyond 10 second windows which could unlock another breadth of tasks. Future works also warrants a clinical validation to the discoveries made by HiMAE which could be of significant interest to the health community.

LLM USAGE

A large language model (LLM) was used to assist in refining the phrasing and structure of the manuscript. Its role was limited to improving clarity, coherence, and readability of the text based on author-provided drafts. All scientific content, experimental design, and analysis were conceived, implemented, and verified by the authors.

ACKNOWLEDGMENTS

We thank Minji Han and Rachel Choi for their expertise in UX/UI design and for crafting the specialized visualizations not supported by standard Python libraries; their design contributions were essential to this work. We also thank Praveen Raja, Matthew Wiggins, and Mike Freedman for their invaluable feedback and insightful discussions throughout the project.

REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: a system for large-scale machine learning. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation, OSDI’16*, pp. 265–283, USA, 2016. USENIX Association. ISBN 9781931971331.
- 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.
- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes, 2018. URL <https://arxiv.org/abs/1610.01644>.
- Julia Amann, Dennis Vetter, Stig Nikolaj Blomberg, Helle Collatz Christensen, Megan Coffee, Sara Gerke, Thomas K Gilbert, Thilo Hagendorff, Sune Holm, Michelle Livne, et al. To explain or not to explain?—artificial intelligence explainability in clinical decision support systems. *PLOS Digital Health*, 1(2):e0000016, 2022.
- Ulzee An, Moonseong Jeong, Simon A Lee, Aditya Gorla, Yuzhe Yang, and Sriram Sankararaman. Raptor: Scalable train-free embeddings for 3d medical volumes leveraging pretrained 2d foundation models. *arXiv preprint arXiv:2507.08254*, 2025.
- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, Jasper Zschiegner, Danielle C. Maddix, Hao Wang, Michael W. Mahoney, Kari Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Yuyang Wang. Chronos: Learning the language of time series, 2024. URL <https://arxiv.org/abs/2403.07815>.
- Leila Arras, Grégoire Montavon, Klaus-Robert Müller, and Wojciech Samek. Explaining recurrent neural network predictions in sentiment analysis, 2017. URL <https://arxiv.org/abs/1706.07206>.
- Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Florian Bordes, Pascal Vincent, Armand Joulin, Mike Rabbat, and Nicolas Ballas. Masked siamese networks for label-efficient learning. In *European conference on computer vision*, pp. 456–473. Springer, 2022.
- Pierre Baldi. Autoencoders, unsupervised learning, and deep architectures. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pp. 37–49. JMLR Workshop and Conference Proceedings, 2012.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? . In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’21*, pp. 610–623, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445922. URL <https://doi.org/10.1145/3442188.3445922>.

- Vera Birrer, Mohamed Elgendi, Olivier Lambercy, and Carlo Menon. Evaluating reliability in wearable devices for sleep staging. *NPJ Digital Medicine*, 7(1):74, 2024.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Lucía Bouza, Aurélie Bugeau, and Loïc Lannelongue. How to estimate carbon footprint when training deep learning models? a guide and review. *Environmental Research Communications*, 5(11):115014, November 2023. ISSN 2515-7620. doi: 10.1088/2515-7620/acf81b. URL <http://dx.doi.org/10.1088/2515-7620/acf81b>.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Ebubekir Buber and DIRI Banu. Performance analysis and cpu vs gpu comparison for deep learning. In *2018 6th International Conference on Control Engineering & Information Technology (CEIT)*, pp. 1–6. IEEE, 2018.
- Stephen Butterworth et al. On the theory of filter amplifiers. *Wireless Engineer*, 7(6):536–541, 1930.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers, 2021. URL <https://arxiv.org/abs/2104.14294>.
- Denisse Castaneda, Aibhlin Esparza, Mohammad Ghamari, Cinna Soltanpur, and Homer Nazeran. A review on wearable photoplethysmography sensors and their potential future applications in health care. *International journal of biosensors & bioelectronics*, 4(4):195, 2018.
- Yong-Mei Cha, Glenn K Lee, Kyle W Klarich, and Martha Grogan. Premature ventricular contraction-induced cardiomyopathy: a treatable condition. *Circulation: Arrhythmia and Electrophysiology*, 5(1):229–236, 2012.
- Cristian Challu, Kin G. Olivares, Boris N. Oreshkin, Federico Garza, Max Mergenthaler-Canseco, and Artur Dubrawski. N-hits: Neural hierarchical interpolation for time series forecasting, 2022. URL <https://arxiv.org/abs/2201.12886>.
- Irene Y Chen, Emma Pierson, Sherri Rose, Shalmali Joshi, Kadija Ferryman, and Marzyeh Ghassemi. Ethical machine learning in healthcare. *Annual review of biomedical data science*, 4(1): 123–144, 2021.
- Peng Chen, Yingying Zhang, Yunyao Cheng, Yang Shu, Yihang Wang, Qingsong Wen, Bin Yang, and Chenjuan Guo. Pathformer: Multi-scale transformers with adaptive pathways for time series forecasting. *arXiv preprint arXiv:2402.05956*, 2024.
- 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, 2020a.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020b. URL <https://arxiv.org/abs/2002.05709>.
- Lawrence J Christiano and Terry J Fitzgerald. The band pass filter. *International economic review*, 44(2):435–465, 2003.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423/>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL <https://arxiv.org/abs/2010.11929>.
- Mohamed Elgendi. On the analysis of fingertip photoplethysmogram signals. *Current cardiology reviews*, 8(1):14–25, 2012.
- Eray Erturk, Fahad Kamran, Salar Abbaspourazad, Sean Jewell, Harsh Sharma, Yujie Li, Sinead Williamson, Nicholas J Foti, and Joseph Futoma. Beyond sensor data: Foundation models of behavioral data from wearables improve health predictions. *arXiv preprint arXiv:2507.00191*, 2025.
- Hany Ferdinando, Matti Huotari, and Teemu Myllylä. Photoplethysmography signal analysis to assess obesity, age group and hypertension. In *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pp. 5572–5575. IEEE, 2019.
- Chloë FitzGerald and Samia Hurst. Implicit bias in healthcare professionals: a systematic review. *BMC medical ethics*, 18(1):19, 2017.
- Sambit Ghadai, Xian Yeow Lee, Aditya Balu, Soumik Sarkar, and Adarsh Krishnamurthy. Multi-level 3d cnn for learning multi-scale spatial features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 0–0, 2019.
- Thomas D Giles, Bradford C Berk, Henry R Black, Jay N Cohn, John B Kostis, Joseph L Izzo Jr, and Michael A Weber. Expanding the definition and classification of hypertension. *The Journal of Clinical Hypertension*, 7(9):505–512, 2005.
- Thomas D Giles, Barry J Materson, Jay N Cohn, and John B Kostis. Definition and classification of hypertension: an update. *The journal of clinical hypertension*, 11(11):611–614, 2009.
- Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220, 2000.
- Awni Hannun, Jagrit Digani, Angelos Katharopoulos, and Ronan Collobert. MLX: Efficient and flexible machine learning on apple silicon, 2023. URL <https://github.com/ml-explore>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. URL <https://arxiv.org/abs/1512.03385>.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- Mingyi He, Bo Li, and Huahui Chen. Multi-scale 3d deep convolutional neural network for hyperspectral image classification. In *2017 IEEE International Conference on Image Processing (ICIP)*, pp. 3904–3908. IEEE, 2017.
- Yuting He, Fuxiang Huang, Xinrui Jiang, Yuxiang Nie, Minghao Wang, Jiguang Wang, and Hao Chen. Foundation model for advancing healthcare: challenges, opportunities and future directions. *IEEE Reviews in Biomedical Engineering*, 2024.

- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023. URL <https://arxiv.org/abs/1606.08415>.
- Syed Anas Imtiaz. A systematic review of sensing technologies for wearable sleep staging. *Sensors*, 21(5):1562, 2021.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift, 2015. URL <https://arxiv.org/abs/1502.03167>.
- Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference, 2017. URL <https://arxiv.org/abs/1712.05877>.
- Ashish Jaiswal, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia Makedon. A survey on contrastive self-supervised learning. *Technologies*, 9(1):2, 2020.
- Li Jing, Pascal Vincent, Yann LeCun, and Yuandong Tian. Understanding dimensional collapse in contrastive self-supervised learning. *arXiv*, 2021. doi: 10.48550/arxiv.2110.09348. URL <https://arxiv.org/abs/2110.09348>.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Yasin Kaya and Hüseyin Pehlivan. Classification of premature ventricular contraction in eeg. *International Journal of Advanced Computer Science and Applications*, 6(7), 2015.
- Aleksia Kolo, Chao Pang, Edward Choi, Ethan Steinberg, Hyewon Jeong, Jack Gallifant, Jason A Fries, Jeffrey N Chiang, Jungwoo Oh, Justin Xu, et al. Meds decentralized, extensible validation (meds-dev) benchmark: Establishing reproducibility and comparability in ml for health. *Machine Learning For Health Conference*, 2024.
- Lingjing Kong, Martin Q Ma, Guangyi Chen, Eric P Xing, Yuejie Chi, Louis-Philippe Morency, and Kun Zhang. Understanding masked autoencoders via hierarchical latent variable models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7918–7928, 2023.
- Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham Kakade, Prateek Jain, and Ali Farhadi. Matryoshka representation learning, 2024. URL <https://arxiv.org/abs/2205.13147>.
- Simon A Lee and Kai Akamatsu. Foundation models for physiological signals: Opportunities and challenges. 2025.
- Simon A Lee, Sujay Jain, Alex Chen, Kyoka Ono, Arabdha Biswas, Ákos Rudas, Jennifer Fang, and Jeffrey N Chiang. Clinical decision support using pseudo-notes from multiple streams of ehr data. *npj Digital Medicine*, 8(1):394, 2025.
- Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. Visualizing and understanding neural models in NLP. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 681–691, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1082. URL <https://aclanthology.org/N16-1082/>.
- Qiao Li, Qichen Li, Ayse S Cakmak, Giulia Da Poian, Donald L Bliwise, Viola Vaccarino, Amit J Shah, and Gari D Clifford. Transfer learning from eeg to ppg for improved sleep staging from wrist-worn wearables. *Physiological measurement*, 42(4):044004, 2021.
- Yuxuan Liang, Haomin Wen, Yuqi Nie, Yushan Jiang, Ming Jin, Dongjin Song, Shirui Pan, and Qingsong Wen. Foundation models for time series analysis: A tutorial and survey. In *Proceedings of the 30th ACM SIGKDD conference on knowledge discovery and data mining*, pp. 6555–6565, 2024.

- Anne-Laure Ligozat, Julien Lefevre, Aurélie Bugeau, and Jacques Combaz. Unraveling the hidden environmental impacts of ai solutions for environment life cycle assessment of ai solutions. *Sustainability*, 14(9):5172, April 2022. ISSN 2071-1050. doi: 10.3390/su14095172. URL <http://dx.doi.org/10.3390/su14095172>.
- Yihan Lin, Zhirong Bella Yu, and Simon Lee. A case study exploring the current landscape of synthetic medical record generation with commercial llms. *arXiv preprint arXiv:2504.14657*, 2025.
- Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=0EXmFzUn5I>.
- Yun Liu, Yu-Huan Wu, Guolei Sun, Le Zhang, Ajad Chhatkuli, and Luc Van Gool. Vision transformers with hierarchical attention. *Machine intelligence research*, 21(4):670–683, 2024.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021a.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows, 2021b. URL <https://arxiv.org/abs/2103.14030>.
- Aleksej Logacjov, Kerstin Bach, and Paul Jarle Mork. Long-term self-supervised learning for accelerometer-based sleep–wake recognition. *Engineering Applications of Artificial Intelligence*, 141:109758, 2025.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2017. URL <https://arxiv.org/abs/1711.05101>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- Jianwen Luo, Kui Ying, and Jing Bai. Savitzky–golay smoothing and differentiation filter for even number data. *Signal processing*, 85(7):1429–1434, 2005.
- Connor MacIsaac, Macros Nguyen, Alexander Uy, Tianmin Kong, and Ava Hedayatipour. A programmable gain calibration method to mitigate skin tone bias in ppg sensors. *Biosensors*, 15(7):423, 2025.
- Melissa D McCradden, Shalmali Joshi, James A Anderson, Mjaye Mazwi, Anna Goldenberg, and Randi Zlotnik Shaul. Patient safety and quality improvement: Ethical principles for a regulatory approach to bias in healthcare machine learning. *Journal of the American Medical Informatics Association*, 27(12):2024–2027, 2020.
- Matthew BA McDermott, Justin Xu, Teya S Bergamaschi, Hyewon Jeong, Simon A Lee, Nassim Oufattole, Patrick Rockenschaub, Kamilė Stankevičiūtė, Ethan Steinberg, Jimeng Sun, et al. Meds: Building models and tools in a reproducible health ai ecosystem. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 6243–6244, 2025.
- Andrew C Miller, Joseph Futoma, Salar Abbaspourazad, Christina Heinze-Deml, Saba Emrani, Ian Shapiro, and Guillermo Sapiro. A wearable-based aging clock associates with disease and behavior. *Nature communications*, 16(1):9264, 2025.
- Girish Narayanswamy, Xin Liu, Kumar Ayush, Yuzhe Yang, Xuhai Xu, Jake Garrison, Shyam A Tailor, Jacob Sunshine, Yun Liu, Tim Althoff, et al. Scaling wearable foundation models. In *The Thirteenth International Conference on Learning Representations*.
- 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.

- Keiron O'Shea and Ryan Nash. An introduction to convolutional neural networks, 2015. URL <https://arxiv.org/abs/1511.08458>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019. URL <https://arxiv.org/abs/1912.01703>.
- Amiya Patanaik, Ju Lynn Ong, Joshua J Gooley, Sonia Ancoli-Israel, and Michael WL Chee. An end-to-end framework for real-time automatic sleep stage classification. *Sleep*, 41(5):zsy041, 2018.
- Ignacio Perez-Pozuelo, Dimitris Spathis, Jordan Gifford-Moore, Jessica Morley, and Josh Cowsls. Digital phenotyping and sensitive health data: Implications for data governance. *Journal of the American Medical Informatics Association*, 28(9):2002–2008, 2021.
- Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models, 2018. URL <https://arxiv.org/abs/1806.07421>.
- Arvind Pillai, Dimitris Spathis, Fahim Kawsar, and Mohammad Malekzadeh. Papagei: Open foundation models for optical physiological signals. *arXiv preprint arXiv:2410.20542*, 2024.
- Arvind Pillai, Dimitris Spathis, Fahim Kawsar, and Mohammad Malekzadeh. Papagei: Open foundation models for optical physiological signals, 2025. URL <https://arxiv.org/abs/2410.20542>.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Neal Schmitt and Goran Kuljanin. Measurement invariance: Review of practice and implications. *Human resource management review*, 18(4):210–222, 2008.
- Muhammad Shabaan, Kaleem Arshid, Muhammad Yaqub, Feng Jinchao, M Sultan Zia, Giridhar Reddy Bojja, Muazzam Iftikhar, Usman Ghani, Loknath Sai Ambati, and Rizwan Munir. Survey: smartphone-based assessment of cardiovascular diseases using ecg and ppg analysis. *BMC medical informatics and decision making*, 20(1):177, 2020.
- Mohammad Amin Shabani, Amir H. Abdi, Lili Meng, and Tristan Sylvain. Scaleformer: Iterative multi-scale refining transformers for time series forecasting. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=sCrn11CtjoE>.
- Wei Shuai, Xi-xing Wang, Kui Hong, Qiang Peng, Ju-xiang Li, Ping Li, Jing Chen, Xiao-shu Cheng, and Hai Su. Is 10-second electrocardiogram recording enough for accurately estimating heart rate in atrial fibrillation. *International journal of cardiology*, 215:175–178, 2016.
- Gerald Simonneau, Nazzareno Galiè, Lewis J Rubin, David Langleben, Werner Seeger, Guido Domenighetti, Simon Gibbs, Didier Lebrec, Rudolf Speich, Maurice Beghetti, et al. Clinical classification of pulmonary hypertension. *Journal of the American College of Cardiology*, 43(12S):S5–S12, 2004.
- G  rald Simonneau, Ivan M Robbins, Maurice Beghetti, Richard N Channick, Marion Delcroix, Christopher P Denton, C Gregory Elliott, Sean P Gaine, Mark T Gladwin, Zhi-Cheng Jing, et al. Updated clinical classification of pulmonary hypertension. *Journal of the American college of cardiology*, 54(1_Supplement_S):S43–S54, 2009.
- Gerald Simonneau, Michael A Gatzoulis, Ian Adatia, David Celermajer, Chris Denton, Ardeschir Ghofrani, Miguel Angel Gomez Sanchez, R Krishna Kumar, Michael Landzberg, Roberto F Machado, et al. Updated clinical classification of pulmonary hypertension. *Journal of the American College of Cardiology*, 62(25S):D34–D41, 2013.

- Gérald Simonneau, David Montani, David S Celermajer, Christopher P Denton, Michael A Gatzoulis, Michael Krowka, Paul G Williams, and Rogerio Souza. Haemodynamic definitions and updated clinical classification of pulmonary hypertension. *European respiratory journal*, 53(1), 2019.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps, 2014. URL <https://arxiv.org/abs/1312.6034>.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks, 2020. URL <https://arxiv.org/abs/1905.11946>.
- Rui Tang, Jaejin An, Brandon K Bellows, Andrew E Moran, and Yiyi Zhang. Trends in hypertension prevalence, awareness, treatment, and control among us young adults, 2003–2023. *American Journal of Hypertension*, pp. hpaf044, 2025.
- Benjamin A Teplitzky, Michael McRoberts, and Hamid Ghanbari. Deep learning for comprehensive ecg annotation. *Heart rhythm*, 17(5):881–888, 2020.
- Rahul Thapa, Bryan He, Magnus Ruud Kjaer, Hyatt Moore Iv, Gauri Ganjoo, Emmanuel Mignot, and James Zou. Sleepfm: Multi-modal representation learning for sleep across brain activity, ecg and respiratory signals. In *International Conference on Machine Learning*, pp. 48019–48037. PMLR, 2024.
- Anja Thieme, Aditya Nori, Marzyeh Ghassemi, Rishi Bommasani, Tariq Osman Andersen, and Ewa Luger. Foundation models in healthcare: Opportunities, risks & strategies forward. In *Extended abstracts of the 2023 CHI conference on human factors in computing systems*, pp. 1–4, 2023.
- Akhil Vaid, Joy Jiang, Ashwin Sawant, Stamatios Lerakis, Edgar Argulian, Yuri Ahuja, Joshua Lampert, Alexander Charney, Hayit Greenspan, Jagat Narula, et al. A foundational vision transformer improves diagnostic performance for electrocardiograms. *NPJ Digital Medicine*, 6(1):108, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Haochen Wang, Kaiyou Song, Junsong Fan, Yuxi Wang, Jin Xie, and Zhaoxiang Zhang. Hard patches mining for masked image modeling, 2023. URL <https://arxiv.org/abs/2304.05919>.
- Ke Wang, Jiamu Yang, Ayush Shetty, and Jessilyn Dunn. Dreamt: Dataset for real-time sleep stage estimation using multisensor wearable technology. *PhysioNet* <https://doi.org/10.13026/62AN-CB28>, 2024.
- Peng Wang, Yuanzhouhan Cao, Chunhua Shen, Lingqiao Liu, and Heng Tao Shen. Temporal pyramid pooling-based convolutional neural network for action recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(12):2613–2622, 2016.
- Lukasz Wesolowski, Bilge Acun, Valentin Andrei, Adnan Aziz, Gisle Dankel, Christopher Gregg, Xiaoqiao Meng, Cyril Meurillon, Denis Sheahan, Lei Tian, et al. Datacenter-scale analysis and optimization of gpu machine learning workloads. *IEEE Micro*, 41(5):101–112, 2021.
- Michael Wornow, Yizhe Xu, Rahul Thapa, Birju Patel, Ethan Steinberg, Scott Fleming, Michael A Pfeffer, Jason Fries, and Nigam H Shah. The shaky foundations of large language models and foundation models for electronic health records. *npj digital medicine*, 6(1):135, 2023.
- 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.

- Qian Xu, Wenzhao Xie, Bolin Liao, Chao Hu, Lu Qin, Zhengzijin Yang, Huan Xiong, Yi Lyu, Yue Zhou, and Aijing Luo. Interpretability of clinical decision support systems based on artificial intelligence from technological and medical perspective: A systematic review. *Journal of healthcare engineering*, 2023(1):9919269, 2023.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical attention networks for document classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pp. 1480–1489, 2016.
- 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 person-days of wearable data. *NPJ digital medicine*, 7(1):91, 2024.
- Qi Zhang, Yifei Wang, and Yisen Wang. How mask matters: Towards theoretical understandings of masked autoencoders. *Advances in Neural Information Processing Systems*, 35:27127–27139, 2022a.
- Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Point-m2ae: multi-scale masked autoencoders for hierarchical point cloud pre-training. *Advances in neural information processing systems*, 35:27061–27074, 2022b.
- Bin Zhou, Rodrigo M Carrillo-Larco, Goodarz Danaei, Leanne M Riley, Christopher J Paciorek, Gretchen A Stevens, Edward W Gregg, James E Bennett, Bethlehem Solomon, Rosie K Singleton, et al. Worldwide trends in hypertension prevalence and progress in treatment and control from 1990 to 2019: a pooled analysis of 1201 population-representative studies with 104 million participants. *The lancet*, 398(10304):957–980, 2021.
- Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization, 2015. URL <https://arxiv.org/abs/1512.04150>.
- Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *International Journal of Machine Learning and Cybernetics*, pp. 1–65, 2024.

APPENDIX

A AUTHOR CONTRIBUTION

We attribute proper credit to the following authors for the development of this project

Table 3: Overview of author contributions.

Author	Concept	Experiment Design	Coding	Analysis	Writing	Visualization	Project Mgmt.	Discussion	Resources
Simon Lee	✓	✓	✓	✓	✓	✓	✓	✓	
Cyrus Tanade			✓	✓	✓		✓	✓	✓
Hao Zhou			✓	✓		✓		✓	
Juhyeon Lee				✓	✓			✓	✓
Megha Thurkal				✓				✓	✓
Minji Han						✓		✓	✓
Baiying Liu				✓				✓	
Keum San Chun								✓	✓
Migyeok Gwak								✓	✓
Mehrab Bin Morshed								✓	
Viswam Nathan								✓	
Mahbubur Rahman								✓	✓
Li Zhu								✓	
Sharanya Desai		✓			✓		✓	✓	✓

B ETHICS CONSIDERATIONS

B.1 DATA PRIVACY AND CONSENT

Wearable signals capture sensitive physiological and behavioral information (Erturk et al., 2025). Our study relies on both publicly available and proprietary (company-owned) datasets that have been carefully vetted. These datasets include transparent disclosure of data usage, explicit opt-in mechanisms, and the option for participants to withdraw (Perez-Pozuelo et al., 2021). Across the seven datasets used in this study, we obtained written consent (via paper or digital waivers) that clearly informed participants that their data may be used for commercial research purposes.

B.2 BIAS AND REPRESENTATIVENESS

Physiological signals vary across age, gender, ethnicity, health status, and socioeconomic context, yet most existing datasets underrepresent key populations (FitzGerald & Hurst, 2017; McCradden et al., 2020; Chen et al., 2021). Such underrepresentation risks embedding biases into foundation models, leading to inequitable performance in downstream applications. Mitigation requires deliberate corpus curation, bias auditing, and systematic evaluation across diverse cohorts. In this study, we sought to mitigate bias by incorporating a pre-training corpus drawn from a wide range of wearable devices, collected across multiple regions of the world and over many years.

B.3 CLINICAL IMPLICATIONS

Wearable foundation models are not substitutes for medical judgment. Their predictions require regulatory approval and clinical validation before integration into healthcare practice. Without safeguards, model misinterpretation could lead to misdiagnosis or inappropriate treatment. Development should involve clinical collaborators, real-world evaluations, and explicit positioning of models as decision-support rather than diagnostic systems. In our group, ongoing collaborations aim to evaluate where our foundation model performs well and how it may assist in forming clinical insights. We emphasize that no definitive clinical conclusions should be drawn from this work.

B.4 ENVIRONMENTAL IMPACT

Training generative models entails substantial computational and environmental costs (Ligozat et al., 2022; Bender et al., 2021; Bouza et al., 2023). To minimize our footprint, we limited redundant runs, and reused checkpoints to avoid unnecessary GPU usage. All experiments were conducted on data-center GPUs with efficient cooling systems and renewable energy credits to reduce carbon intensity. We emphasize that transparent reporting of compute usage and bounding resource allocation are necessary steps toward sustainable machine learning research.

C REPRODUCIBILITY STATEMENT

Table 4: HiMAE architecture components.

Encoder-Decoder			
Layer	Output Shape	EncoderConvBlock	DecoderSkipBlock
Input	[B, 1, T]		
EncoderConvBlock(1→16)	[B, 16, $T/2$]	Layer	Layer
EncoderConvBlock(16→32)	[B, 32, $T/4$]	Conv1d ($k = 5, s=2, p=2$)	ConvTranspose1d ($k = 5, s=2, p=2, op=1$)
EncoderConvBlock(32→64)	[B, 64, $T/8$]	BatchNorm	Concat skip connection
EncoderConvBlock(64→128)	[B, 128, $T/16$]	GELU	Conv1d ($k = 5, s=1, p=2$)
EncoderConvBlock(128→256)	[B, 256, $T/32$]	Conv1d ($k = 5, s=1, p=2$)	BatchNorm
DecoderSkipBlock(256→128)	[B, 128, $T/16$]	BatchNorm	GELU
DecoderSkipBlock(128→64)	[B, 64, $T/8$]	Conv1d ($k = 1, s=2$) + BN	Conv1d ($k = 5, s=1, p=2$)
DecoderSkipBlock(64→32)	[B, 32, $T/4$]	GELU	BatchNorm
DecoderSkipBlock(32→16)	[B, 16, $T/2$]		GELU
Final Deconv (16→1)	[B, 1, T]		
Tanh	[B, 1, T]		

Due to restrictions around data licensing and industry policies, we are unable to release the full source code associated with HiMAE. However, To mitigate this limitation, we provide a simplified code base in this <https://github.com/Simonlee711/HiMAE> as well as complete details of the model architecture, layer configurations, and hyperparameters in Table 4. This includes all encoder, decoder, and skip connection blocks, along with kernel sizes, strides, padding, activation functions, and normalization layers. Together, these descriptions and codebases are sufficient to re-implement the model faithfully in any modern deep learning framework (Paszke et al., 2019; Abadi et al., 2016; Bradbury et al., 2018; Hannun et al., 2023). In addition, we report all training settings (e.g., optimizer, learning rate schedule, and batch size) in the Appendix Section E to further support reproducibility. Our goal is to ensure that, while the exact implementation cannot be shared, independent researchers can replicate the methodology and validate the findings presented in this work.

C.1 TEMPORAL RESOLUTION AS AN EXPLICIT SCALE AXIS

HiMAE’s encoder implements a structured mapping from depth to temporal scale. Each encoder block halves the temporal resolution while increasing the span of input samples contributing to each feature, yielding a hierarchy of representations indexed by effective temporal support. This makes temporal resolution an explicit axis of representation, rather than an emergent byproduct of depth.

Concretely, the encoder is composed of b convolutional blocks, each reducing the sequence length by a factor of two. As a result, the representation at depth b operates on a grid of resolution $T/2^b$. At the same time, each block aggregates information over an increasingly large window of the input signal. Because kernel size is fixed across layers, the cumulative temporal support of encoder features grows exponentially with depth, scaling as

$$R_b = \Theta(2^b),$$

up to architecture-dependent constants. Thus, encoder depth simultaneously controls both the granularity at which the signal is represented and the temporal extent over which features are computed.

Table 5 instantiates this scale hierarchy for the HiMAE encoder. Shallow layers operate at high temporal resolution with receptive fields spanning only a few tens of samples, capturing fine-scale waveform morphology. Intermediate layers aggregate information over 10^1 – 10^2 samples, corresponding to sub-second temporal structure such as beat-to-beat variability. The deepest layers integrate over several hundred samples, encoding longer-range physiological dynamics across multiple cardiac cycles.

This explicit scale stratification is central to masked autoencoding on physiological signals. Because masking removes contiguous temporal regions, successful reconstruction requires contextual information at a scale comparable to the masked interval. Features whose receptive fields are too small lack sufficient context, while features whose receptive fields are too large oversmooth across distinct physiological events. HiMAE’s exponential scale ladder ensures that intermediate encoder depths naturally align with the characteristic temporal extent of masked regions, concentrating learning signal at those resolutions.

Table 5: Temporal resolution and cumulative receptive field through the encoder. T denotes the input length in samples. R_ℓ is the receptive field after layer ℓ .

Layer	Kernel k	Stride s	Output length	R_ℓ
Enc1-conv1	5	2	$T/2$	5
Enc1-conv2	5	1	$T/2$	13
Enc2-conv1	5	2	$T/4$	21
Enc2-conv2	5	1	$T/4$	37
Enc3-conv1	5	2	$T/8$	53
Enc3-conv2	5	1	$T/8$	85
Enc4-conv1	5	2	$T/16$	117
Enc4-conv2	5	1	$T/16$	181
Enc5-conv1	5	2	$T/32$	245
Enc5-conv2	5	1	$T/32$	373

Viewed this way, encoder depth in HiMAE should not be interpreted as a measure of abstraction alone, but as an index over temporal scales. This perspective explains why linear probes trained on intermediate layers often outperform both shallower and deeper representations: they correspond to resolutions at which physiological structure is most predictive.

D DATASETS

D.1 AQUITION AND APPROVAL

All data analyzed in this study were collected under informed consent, with participants explicitly agreeing for their wearable-derived signals to be used in health-related research. The consent language stated that data could be used for developing new health features and algorithms and for inclusion in scientific publications. In particular, participants were informed that health and wellness data such as steps, heart rate, sleep, and photoplethysmography (PPG) signals could contribute to findings aimed at advancing general knowledge of health and science. No data used in this study included personally identifying information such as names or email addresses. We attach a portion of the protocols defined in our user data agreements below:

The use of these de-identified data for data usage was reviewed and classified as exempt. In addition, because the supporting records constitute case histories and document exposure to devices, we complied with the recordkeeping requirements in 21 CFR § 812.140(a)(3), including obtaining written digital consent and dated information. Participants could withdraw at any time; such withdrawals were documented in the case history, and data collected up to the point of withdrawal were retained and used for the investigation in accordance with the consent and applicable regulations.

For downstream evaluations, we relied on a combination of institutional review board (IRB)-approved datasets and publicly available resources. For instance, the PVC detection task used paired PPG and ECG recordings to derive annotations of premature ventricular contractions, with ECG-based labels verified both algorithmically and manually. The hypertension classification tasks were drawn from the My Heart Lab Study collected in a lab Setting (ID NCT04314947) and My BP Lab (Clinical Trials ID 19-27169) studies collected in a free-world setting, both of which collected wrist-based PPG alongside reference blood pressure measurements under IRB-approved protocols. Sleep staging was evaluated using the DREAMT dataset, which combines PPG with gold-standard polysomnography annotations in individuals with and without diagnosed sleep disorders. Finally, a range of abnormal lab test prediction tasks were derived from the Tulane University dataset (ID 20242033), linking PPG from Samsung devices with clinical laboratory values for biomarkers (More details in Appendix Section D).

Across all studies, participants consented to data collection through mobile platforms that supported eligibility screening and enrollment, provided full informed consent, and enabled seamless integration of Samsung devices for continuous signal acquisition. Where appropriate, participants also reported medical histories or completed questionnaires through these platforms. All data were de-identified and stored in accordance with the approved study protocols, ensuring compliance with ethical and regulatory standards.

This layered consent and governance framework ensures that the data underpinning our pretraining and evaluation tasks are both ethically sourced and scientifically robust, supporting the broader goal of advancing health monitoring through consumer wearables.

D.2 PRE-TRAINING DATASETS

Table 6: **Demographic Characteristics of the Study Population.** Distributions are shown by biological sex, age group, racial identity, and BMI category ($N = 47,644$).

Category	Subgroup	N	% of total
Sex	Male	36,990	77.6
	Female	10,532	22.1
	Another gender identity	122	0.3
Age	18–29	12,019	25.2
	30–49	27,207	57.2
	50–64	7,067	14.8
	65+	1,351	2.8
Race	White	31,029	65.2
	Asian or Pacific Islander	7,630	16.0
	Black or African American	3,414	7.2
	American Indian or Alaskan Native	592	1.2
	Another race	4,979	10.4
BMI	Underweight (< 18.5)	823	1.7
	Normal weight (18.5 – 24.9)	13,626	28.5
	Overweight (25 – 29.9)	16,634	34.9
	Obese I (30 – 34.9)	8,745	18.5
	Obese II and III (≥ 35)	7,816	16.4

D.2.1 DEVICE DISTRIBUTION

The distribution of participants and data availability highlights both the diversity of collection devices and the heterogeneity of study contributions (Figure 8). At the device level, participation is primarily sourced from Watch Active 2, Watch 3, Watch Active, each contributing a lot of participants, while older models such as the Galaxy Gear S3 are represented by fewer users. This heterogeneity in devices provide us with a realistic and diverse set of raw wearable signals that can help us build generalizable foundation models. The presence of entries labeled as “NA” further reflects the mixture of collection devices and the occasional incompleteness of metadata. *We note that the devices used in our study are provided by two distributors limiting its generalizability and causing potential biases due to not having access to other consumer wearable devices.*

D.2.2 PARTICIPANT COUNTS

In terms of study based segmentation, the dataset contains a handful of large-scale cohort studies, leading to diverse representation (Figure 8). Efforts were made to ensure representation across studies of varying sizes. This underscores the necessity of leveraging the vast scale of high-volume cohorts while simultaneously preserving the heterogeneity introduced by smaller studies, since both dimensions are essential for building foundation models that truly capture the variability and complexity of one-dimensional PPG signal modeling. Our data was collected across 4 countries (USA, South Korea, Brazil, Bangladesh) and the demographics are highlighted in Table 6. Note that missing demographics were imputed via KNN based on average PPG segments which have shown to recover these people specific attributes (Miller et al., 2025; MacIsaac et al., 2025; Ferdinando et al., 2019).

D.2.3 PRE-PROCESSING PIPELINE

We segment raw PPG signals into fixed-length 10 s windows and apply a lightweight quality-control pipeline to remove motion artifacts and non-physiologic segments. Each window is first standardized to remove scale and offset differences across devices and recording conditions. Windows with extreme amplitude fluctuations, indicative of motion bursts or sensor saturation, are identified using a simple distributional check and either trimmed to remove outliers or discarded if the signal remains unstable. This step prioritizes precision over recall to ensure high-quality pretraining data.

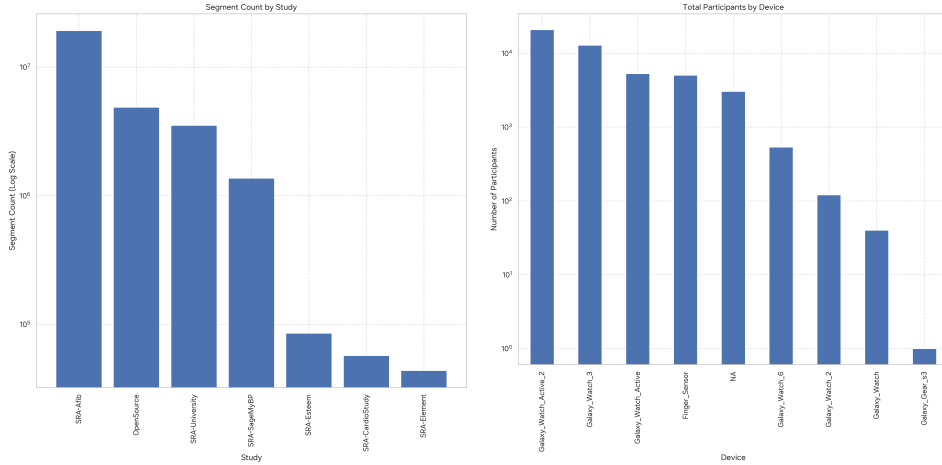


Figure 8: **Segment Count by Study.** This bar chart shows the number of data segments collected for each study, with the y-axis on a logarithmic scale to account for the large differences in segment counts.

For windows that pass amplitude screening, we assess temporal regularity by measuring short-range autocorrelation. Physiologically plausible PPG signals exhibit quasi-periodic structure; windows with highly irregular or unstable periodicity are rejected, as these patterns typically arise from motion or sensor decoupling. We additionally enforce a minimum number of cycles to eliminate degenerate or truncated traces.

Surviving windows are band-pass filtered to the cardiac frequency range to remove baseline drift and high-frequency noise while preserving pulse morphology. Signal quality is then evaluated via template matching against a canonical PPG waveform. We compute a per-window quality score that jointly reflects the fraction of the signal that matches the template and the strength of that match, penalizing cases where apparent agreement is driven by only a small portion of the window. Windows that fail this final morphology check are excluded.

This filtering is applied at scale across the corpus, retaining only windows that are clean, periodic, and morphologically consistent. The resulting pretraining set emphasizes physiologically meaningful PPG signals across devices and sampling rates, substantially reducing motion artifacts without relying on labels, heuristics tied to specific hardware, or subject-level metadata.

D.3 DOWNSTREAM EVALUATION DATA

We evaluate HiMAE across diverse downstream tasks to assess the generality of wearable PPG representations. Rather than assuming a fixed mapping between PPG and outcomes, we exploit HiMAE’s ability to learn hierarchical temporal features and adaptively resolve signal segments at scales most informative for prediction. This design allows us to probe the representational value of optical physiological signals across clinically and behaviorally relevant applications.

D.3.1 PVC DETECTION

Table 7: Stratified 80/20 Train/Test splits for PVC tasks (with per-task totals).

Task	Split	Negative	Positive	Total
PVC Detection	train	369987 (91.8%)	32832 (8.2%)	402819
	test	69880 (89.7%)	8019 (10.3%)	77899
	totals	439767 (91.4%)	40950 (8.6%)	480717

Premature Ventricular Contractions (PVCs) (Number Breakdowns in Table 7) are abnormal beats arising in the ventricles (Cha et al., 2012; Kaya & Pehlivan, 2015). We use paired PPG–ECG data, with ECG annotations generated using BeatLogic (Teplitzky et al., 2020) and manually verified.

PPG inputs are 10s non-overlapping wrist segments, pre-processed with a Savitzky–Golay filter (Luo et al., 2005), a 0.5–4.0 Hz bandpass, normalization to $[-1, 1]$, and exclusion of segments with motion artifacts or disruptions > 1 s. This task evaluates whether ubiquitous PPG can approximate arrhythmia detection typically restricted to ECG.

D.3.2 HYPERTENSION CLASSIFICATION

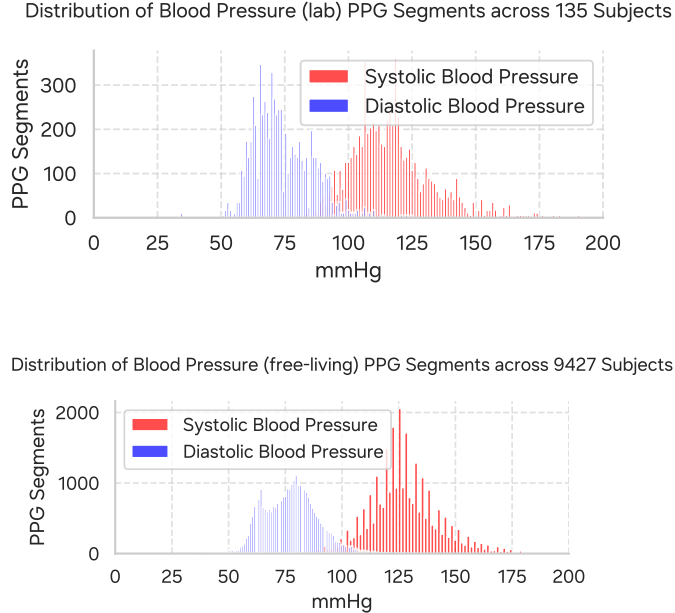


Figure 9: **Blood Pressure Distribution:** The distribution of Blood Pressure Values (mmHG) across the lab and free-living studies. We define hypertension as systolic over 130 and diastolic over 80 to generate binary outcomes.

Hypertension classification (Number Breakdowns in Figure 9) relies on cuff-based references (Simonneau et al., 2004; Giles et al., 2005; 2009; Simonneau et al., 2009; 2013; 2019). Subjects within ± 8 mmHg of the diagnostic cutoff are excluded to reduce label noise, with remaining individuals labeled hypertensive or normotensive. Each 10s PPG segment undergoes Savitzky–Golay smoothing, 0.5–4.0 Hz bandpass filtering, normalization to $[-1, 1]$, and artifact removal. Unlike PVC detection, which is event-based, this task leverages PPG morphology and temporal dynamics to reflect vascular state. These evaluations contain both hypertension data collected in a naturalistic free world environment and within a controlled lab environment for both the hypertensive and blood pressure regression tasks.

D.3.3 SLEEP STAGING

Table 8: Stratified 80/20 Train/Test splits for Sleep Staging.

Task	Split	Wake	Light	Deep	REM	Total
Sleep Staging (4-class)	train	44829 (23.9%)	115932 (61.8%)	6696 (3.6%)	20214 (10.8%)	187671
	test	11298 (23.6%)	30153 (63.1%)	1416 (3.0%)	4881 (10.2%)	47748
	totals	56127 (23.8%)	146085 (61.9%)	8112 (3.4%)	25095 (10.6%)	235419

Sleep staging (Number Breakdowns in Tables 8) is evaluated on the DREAMT dataset (Wang et al., 2024) hosted on PhysioNet (Goldberger et al., 2000), which includes overnight wristband data with simultaneous PSG. Annotations follow AASM standards into wake, REM, NREM1, NREM2, and NREM3, excluding missing and preparation segments. PPG is bandpass filtered (0.5–12 Hz) (Butterworth et al., 1930), segmented into 10s windows, and normalized to zero mean and unit variance. Performance is measured with five-fold subject-independent cross-validation. This task examines whether PPG encodes temporal patterns sufficient for sleep stage classification. *We note that sleep*

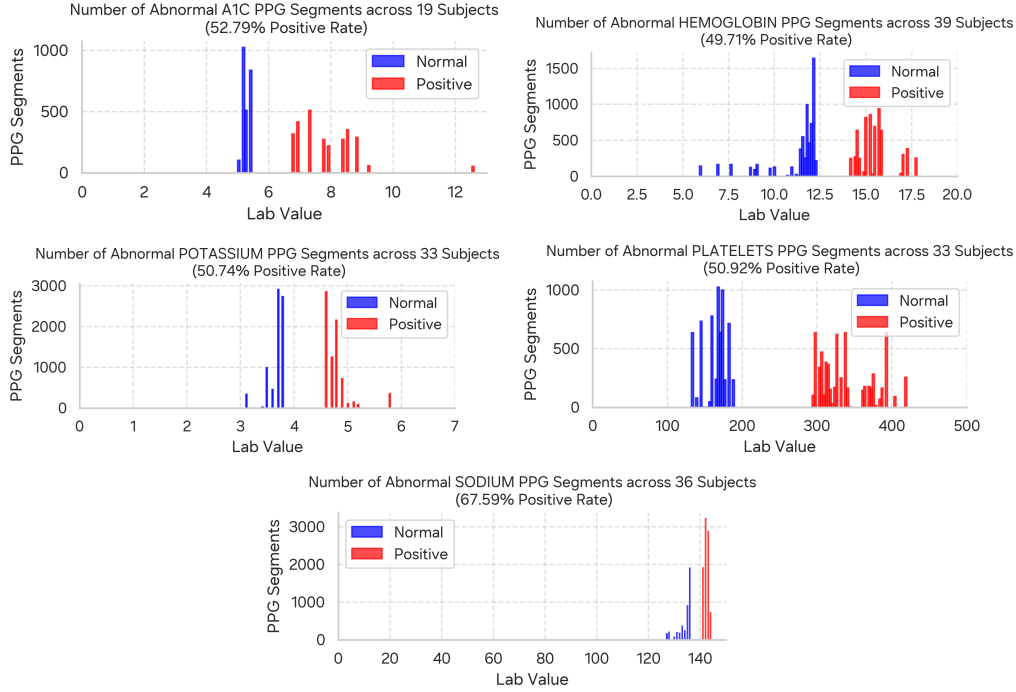


Figure 10: **Abnormal Labs Distribution:** The number of PPG segments for abnormal labs separated based on lab-specific cutoffs. We define an abnormal lab as falling above the 75th percentile of values and a normal lab as falling within the 25th percentile.

staging has canonically been designed by leveraging the whole sleep cycle but we are assessing the ability to monitor real time sleep staging from much shorter PPG segments.

D.3.4 ABNORMAL LAB TESTS

For abnormal lab test prediction, we use Samsung Galaxy Watch PPG collected at Tulane University paired with clinical laboratory results. Each test is framed as a binary classification task: outcomes are labeled negative if within the 25th percentile of lab values and the positive labels are anything above the 75th percentile (Figure 10). All other labels are excluded. Preprocessing matches other tasks. Targets include A1C, hemoglobin, platelets, potassium, and sodium, each selected for established clinical relevance. This task extends evaluation beyond cardiovascular and behavioral endpoints to systemic markers of metabolic, and hematologic health. *We note that it is unclear whether PPG can predict abnormal from healthy lab values based on the PPG alone. Despite this, Tulane university presents us with an opportunity to discover if PPG signal can provide digital signatures making this an exploratory task in our benchmark.*

E BASELINES AND MODEL CONFIGURATION

Self Supervised Pre-trained methods have become a dominant paradigm for health and wellness to study a variety of applications (Wornow et al., 2023; Thieme et al., 2023; He et al., 2024; An et al., 2025; Lin et al., 2025). Foundation models for one-dimensional signals are predominantly repurposed from architectures designed for vision, with adaptations that reinterpret temporal structure as a flattened analogue of spatial correlation. In this section we highlight our baseline models and model configurations

E.1 BASELINES

MAE-1D We introduce a Masked Autoencoder that mimics the protocol of LSM (Narayanswamy et al., 2024). This model introduces a large-scale foundation model trained on multimodal wearable sensor data but we adapt it to 1D PPG Signal. Specifically, it adopts a vision transformer architecture trained via masked autoencoding with random masking. In our work, we do not replicate the full multimodal design; instead, we adapt and constrain the model to a unimodal setting for fair comparison and due to lack of open source code.

Swin-Transformer (Liu et al., 2021a) is a hierarchical Transformer that forms multi-scale representations by restricting self-attention to non-overlapping windows and alternating partitions with a shifted-window scheme, which enables cross-window communication while keeping computation near-linear in sequence length. We use this baseline as this is a direct comparison and counterpart to our proposed hierarchical HiMAE model. For wearable sensing, we adopt a 1D adaptation that tokenizes temporal patches and applies windowed attention along time, capturing both fine-grained waveform morphology and longer-range dependencies.

Masked Siamese Networks (MSN) (Assran et al., 2022) learn label-efficient representations by combining masked signal modeling with Siamese-style contrastive objectives. Instead of relying on class labels, MSN masks portions of the input and enforces consistency between augmented views. Architecturally, it employs a Vision Transformer encoder shared across views, while leveraging a predictor network to stabilize training. The key idea is to couple self-distillation with masked reconstruction to reduce sample complexity.

DINO (Caron et al., 2021) is a self-supervised framework that leverages knowledge distillation without labels. Using a teacher-student setup, the student network is trained to match the output distribution of the teacher under different data augmentations. Both networks are 1D-ViTs, and the method induces cluster-like emergent properties in the learned embedding space, enabling strong transfer performance without explicit contrastive pairs or handcrafted pretext tasks.

SimCLR (Chen et al., 2020b) establishes contrastive learning as a competitive self-supervised paradigm. The core idea is to maximize agreement between augmented views of the same signal in a latent space while pushing apart representations of different images. This is implemented using a ResNET encoder (He et al., 2015), a projection head, and a contrastive loss (NT-Xent (Chen et al., 2020a)).

PaPaGei (Pillai et al., 2024) is a domain-specific foundation model designed for optical physiological sensing, particularly photoplethysmography (PPG). It adapts ResNET-style CNN architectures to learn robust, generalizable representations from large-scale optical physiological datasets. PaPaGei releases both model weights and datasets to support reproducibility and broader adoption in physiological signal analysis. In our work, we used their source code to benchmark their method by pre-training on our volume of data to ensure fair comparison.

E.2 HYPERPARAMETERS FOR HiMAE AND BASELINES

To ensure a fair comparison across models, we aligned the training setup as closely as possible to the original implementations while maintaining consistency in optimizer choice and scheduling. All the methods trained from scratch (HiMAE, MAE-1D, Swin-Transformer, MSN, DINO, SimCLR) were trained under identical optimization regimes, while PaPaGei follows its released open source training protocol. Table 9 summarizes the key hyperparameters for all models.

Table 9: Hyperparameter Configurations for Different Models

Configuration	HiMAE	MAE-1D	Swin-Transformer	MSN	DINO	SimCLR	PaPaGei
Training Steps				50000			15000
Warmup Steps				2500			—
Optimizer			AdamW (Loshchilov & Hutter (2017))				
Opt. momentum $[\beta_1, \beta_2]$	[0.9, 0.95]	[0.9, 0.95]	[0.9, 0.95]	[0.9, 0.99]	[0.9, 0.99]	[0.9, 0.99]	—
Base learning rate	0.001	0.005	0.005	0.001	0.004	0.001	0.0001
Batch size			2048				256
Weight decay			0.0001				—
Gradient clipping	1.0	1.0	1.0	3.0	3.0	3.0	—
Dropout			0.0				—
Learning rate schedule			Linear Warmup & Cosine Decay				
Loss Function			Mean Squared Error	Cross Entropy		Contrastive Loss	
Data resolution			1 (signal) - 100 Hz (Sampling Rate) \times 10 (seconds)				
Augmentation			Flip, Time-Warping, Noise				

E.3 LAYER WISE ANALYSIS

Layer Wise Analysis

```

1 def layerwise_probe(model, dataloader, labels, device):
2     """
3     For all architectures we use the full sequence embedding across
4     the temporal dimension, without collapsing to a single summary token,
5     to ensure that downstream probes have access to resolution-specific
6     information.
7     """
8     model.to(device).eval()
9     hooks, acts = [], {}
10
11     # capture encoder layer outputs
12     for i, layer in enumerate(model.encoder_layers):
13         hooks.append(layer.register_forward_hook(lambda m, x, y, i=i: acts
14             .setdefault(i, y.detach().cpu()))))
15
16     Xs = {i: [] for i in range(len(model.encoder_layers))}
17     ys = []
18     for xb in dataloader:
19         xb = xb.to(device)
20         with torch.no_grad():
21             _ = model(xb.transpose(1, 2))
22         for i, a in acts.items():
23             Xs[i].append(a.flatten(1).numpy()) # flatten embedding to
24             # preserve time
25         ys.append(labels[: xb.size(0)])
26         labels = labels[xb.size(0):]
27
28     results = {}
29     for i, feats in Xs.items():
30         X = np.concatenate(feats)
31         y = np.concatenate(ys)
32         X = StandardScaler().fit_transform(X)
33         clf = LogisticRegression(max_iter=1000)
34         auc = cross_val_score(clf, X, y, cv=StratifiedKFold(5), scoring="
35             roc_auc").mean()
36         results[f"layer_{i}"] = auc
37
38     for h in hooks: h.remove()
39     return results

```

The layer-wise analysis examines how temporal resolution in the learned representations aligns with performance on downstream tasks. For each encoder block, we extract and flatten the full sequence embedding to preserve temporal detail, allowing probes to access features at different levels of abstraction. By training logistic regression classifiers on these embeddings, we can assess which layers best capture task-relevant temporal patterns.

F ADDITIONAL RESULTS

F.1 MODEL CONFIGURATIONS ABLATIONS

We conducted a comprehensive ablation study of HiMAE on a 100 Hz dataset comprising ten million segments (roughly 30k hours). The experiments systematically varied architecture and hyperparameters to understand their effect on reconstruction quality (Extrapolation task from our generative benchmark in tables where it is not explicitly stated as previously done in (Narayanswamy et al., 2024)), with multiple independent training runs averaged to reduce variance from stochastic initialization and data sampling. Unless otherwise noted, all training employed AdamW with a learning rate of 3×10^{-4} , cosine decay scheduling, and a batch size of 512.

Architecture. We evaluated HiMAE alongside CNN baselines across increasing network depths, defined by the sequence of hidden channel dimensions [16, 32, 64], [16, 32, 64, 128], and [16, 32, 64, 128, 256]. Table 10 lists the parameter counts, showing a modest growth for HiMAE compared to CNN baselines, with the skip-connected HiMAE exhibiting slightly higher capacity than its no-skip variant.

Table 10: Model Parameters (in K or M)

Model <i>Depth</i>	HiMAE-tiny [16,32,64]	HiMAE-small [16,32,64,128]	HiMAE-Base [16,32,64,128,256]
CNN	26.2 K	108 K	437 K
HiMAE-no skip	66.1 K	271 K	1.10 M
HiMAE	75.3 K	309 K	1.25 M

The impact of network depth on mean absolute error (MAE) and mean squared error (MSE) is summarized in Table 11. Increasing depth consistently reduced both MAE and MSE for HiMAE, with the deepest configuration yielding the lowest reconstruction error. Skip connections were critical, as HiMAE consistently outperformed its no-skip variant across all depths.

Table 11: MAE and MSE for Different Network Depths

Model <i>Depth</i>	HiMAE-tiny [16,32,64]		HiMAE-small [16,32,64,128]		HiMAE-Base [16,32,64,128,256]	
	MAE ↓	MSE ↓	MAE ↓	MSE ↓	MAE ↓	MSE ↓
CNN	0.405,2	0.234,5	0.417,7	0.249,1	0.400,8	0.231,5
HiMAE-noskip	0.403,1	0.236,5	0.400,6	0.246,5	0.397,5	0.233,9
HiMAE	0.400,8	0.230,9	0.389,2	0.223,2	0.3827	0.2210

Patch Size. We varied the spatial-temporal patch sizes over 1, 5, 10, and 20. The results in Table 13 indicate that 5 provided the best trade-off between local resolution and generative performance. Smaller patches increased flexibility but slightly degraded performance due to reduced receptive field per token, while overly large patches caused loss of fine-grained structure.

Table 13: Model Performance for Different Patch Sizes

Model	1		5		10		20	
	MAE ↓	MSE ↓	MAE ↓	MSE ↓	MAE ↓	MSE ↓	MAE ↓	MSE ↓
CNN	0.414,0	0.239,1	0.400,8	0.231,5	0.412,2	0.244,9	0.427,4	0.261,3
HiMAE-noskip	0.406,9	0.239,8	0.397,6	0.233,9	0.403,7	0.246,2	0.419,5	0.262,9
HiMAE	0.389,9	0.226,8	0.3827	0.2210	0.386,1	0.231,2	0.403,9	0.247,9

Convolution Kernel Size. Kernel size was varied over {1, 5, 10, 20}. Table 14 shows that 5 yielded the lowest errors across all models, suggesting moderate receptive fields match the temporal and

spatial scales of our data. Very small kernels restricted context aggregation, while very large kernels oversmoothed latent features.

Table 14: Model Performance Across Convolution Kernel Sizes

Model	1		5		10		20	
	MAE ↓	MSE ↓	MAE ↓	MSE ↓	MAE ↓	MSE ↓	MAE ↓	MSE ↓
CNN	0.416,2	0.241,3	0.401,0	0.230,9	0.410,3	0.241,8	0.424,1	0.257,6
HiMAE-noskip	0.409,0	0.242,7	0.395,9	0.233,1	0.403,2	0.244,0	0.420,8	0.259,1
HiMAE	0.392,1	0.228,3	0.3821	0.2206	0.388,5	0.231,6	0.404,7	0.248,5

Stride. We evaluated stride values of 2, 4, and 8 (Table 15). Smaller strides yielded the best performance, particularly for HiMAE, by preserving high temporal resolution in early feature maps. Performance degraded monotonically with stride increases.

Table 15: Model Performance Across Stride Values

Model	2		4		8	
	MAE ↓	MSE ↓	MAE ↓	MSE ↓	MAE ↓	MSE ↓
CNN	0.4016	0.2312	0.413,9	0.244,5	0.431,8	0.267,8
HiMAE-noskip	0.3976	0.2334	0.409,8	0.247,1	0.427,2	0.270,2
HiMAE	0.3829	0.2209	0.392,8	0.232,5	0.410,3	0.250,4

Masking Ratio. Finally, we explored the effect of varying the latent masking ratio in the masked autoencoding objective for generative tasks, with ratios from 0.5 to 0.9. As shown in Table 16, interpolation and extrapolation both improved when increasing the ratio up to 0.8, after which performance degraded for interpolation and collapsed for extrapolation.

Table 16: MAE and MSE for HiMAE Across Different Masking Ratios Evaluated on Generative Tasks

HiMAE Masking Ratio	Temporal Interpolation		Temporal Extrapolation	
	MAE ↓	MSE ↓	MAE ↓	MSE ↓
0.5	0.397,2	0.229,2	0.407,7	0.251,9
0.6	0.388,9	0.222,3	0.397,5	0.229,4
0.7	0.384,8	0.220,7	0.396,3	0.227,8
0.8	0.3796	0.2183	0.3879	0.2217
0.9	0.381,8	0.221,9	0.288,1	0.221,6

Final Selection. These controlled experiments informed the final HiMAE configuration: the deepest architecture [16, 32, 64, 128, 256] with skip connections, patch size 5, kernel size 5, stride 2, and a masking ratio of 0.8, which jointly achieved the best trade-off between reconstruction fidelity and parameter efficiency.

F.2 ECG PRE-TRAINING

HiMAE attains the lowest masked-reconstruction error on ECG (Table 17), indicating that its hierarchical masking and reconstruction inductive biases capture reconstruction capacity beyond PPG. MAE-1D (ViT) is a close second, while the ablated HiMAE and CNN trail, reinforcing that the full HiMAE design transfers effectively to the ECG domain.

Table 17: Masked-reconstruction loss on ECG masked auto encoding task.

Model	MSE (↓)
HiMAE	0.148
MAE-1D (ViT)	0.162
HiMAE (no skip)	0.184
CNN	0.207

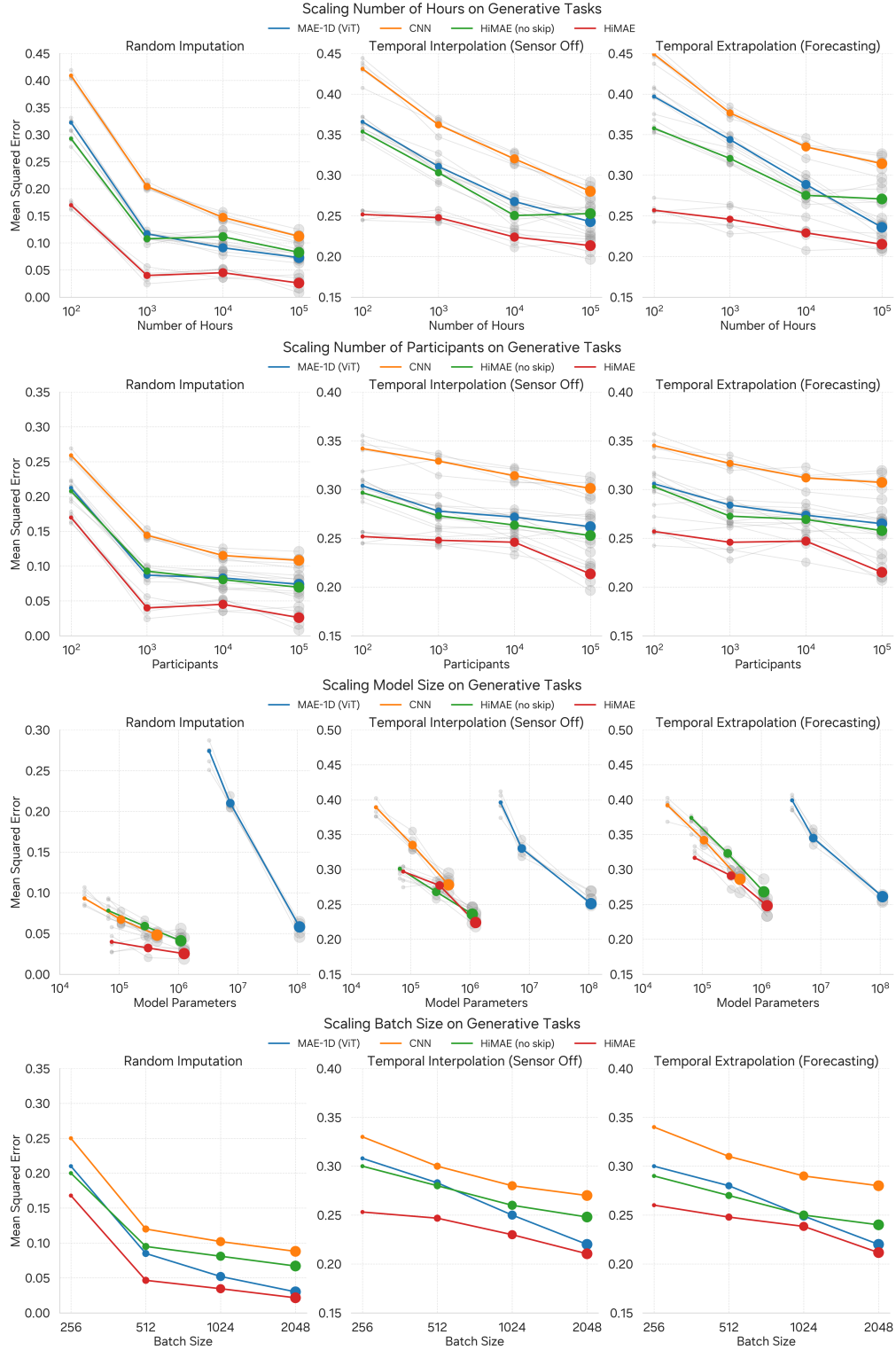


Figure 11: **Scaling Experiments on Generative Tasks:** Evaluation on the three generative tasks. HiMAE consistently outperforms all model at our scale of data

F.3 SCALING RESULTS FOR GENERATIVE TASKS

Scaling analysis. We evaluate HiMAE’s reconstruction error under participant, recording hour, batch size, and model size scaling, following the regimes of Narayanswamy et al. (2024); Xu et al. (2025): random imputation, temporal interpolation, and temporal extrapolation. Across all settings HiMAE follows clean scaling law trends (Kaplan et al., 2020) and maintains a margin over MAE-1D (ViT) and CNN baselines.

The most pronounced effect is model size. At small capacities HiMAE achieves lower error than much larger transformer baselines, highlighting the advantage of hierarchical inductive bias over sheer parameter count. MAE-1D only begins to close the gap at orders of magnitude more parameters. The transformer could surpass our HiMAE model when given a larger capacity but this again highlights the effectiveness of the inductive bias that we are conveying. Participant, hour, and batch size scaling follow canonical patterns. More participants and longer recordings steadily reduce error, with HiMAE continuing to improve where baselines saturate, especially on interpolation and extrapolation.

F.4 HIERARCHICAL CONCORDANCE

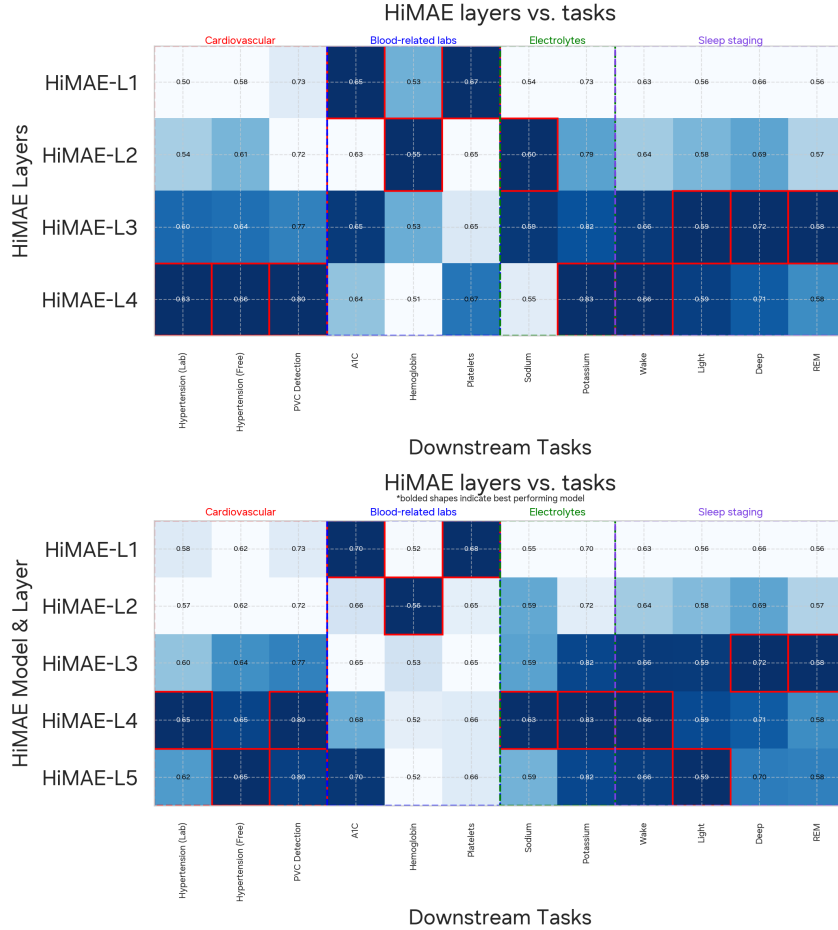


Figure 12: **HiMAE layer concordance across encoder depths.** Heatmaps compare downstream AUROC when probing HiMAE at 4 layers (top) versus 5 layers (bottom). Despite the removal of an encoder–decoder stage, the resolution–task alignment remains highly concordant: tasks such as PVC detection and hypertension consistently peak at similar layers, while sleep staging benefits from coarser representations. Minor deviations appear in intermediate layers, but the overall hierarchy of predictive resolutions is preserved, indicating robustness of the resolution hypothesis to architectural depth.

Layer concordance across depths. We further assess the stability of the resolution hypothesis by comparing HiMAE trained with four versus five encoder-decoder stages (Figure 12). The resulting heatmaps reveal that the alignment between downstream tasks and temporal resolutions is largely preserved across depths. Cardiovascular endpoints such as PVC detection and hypertension consistently achieve their best performance at finer layers, while blood related labs benefits from coarser layers. Although minor fluctuations appear in intermediate levels, the overall hierarchy of predictive resolutions is concordant. This suggests that the resolution-task mapping uncovered by HiMAE is not an artifact of architectural depth, but a robust property of the representations themselves.

F.5 TRANSFORMER LAYER INTERPRETABILITY

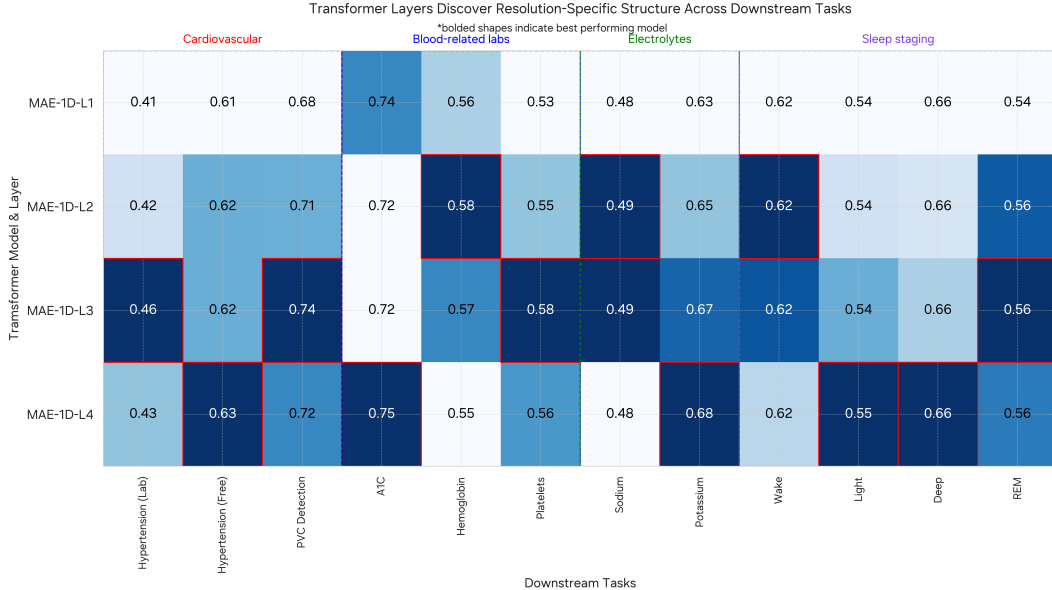


Figure 13: **Transformer based Layer Wise Analysis:** Unlike HiMAE, MAE-1D exhibits non-monotonic trends and lacks concordance with HiMAE’s internal representation hierarchy. We observe that performance typically peaks in intermediate and task-specific later layers, though this pattern carries important nuance.

Figure 13 outline a deep conceptual contrast between HiMAE’s layer wise interpretability from Figure 14 and Appendix Section F.4 and Transformer based encoders like MAE-1D (ViT) in how their internal representations evolve and why their layerwise “probing” behaviors differ.

In a U-Net with convolution operators, each layer has a localized receptive field that gradually expands with depth. Convolution is a local operator, so early layers capture fine-grained spatial details (edges, textures), mid layers combine local motifs into parts, and deeper layers encode semantic abstractions or whole objects. Skip connections reintroduce lost resolution but don’t globalize information. This creates a hierarchical representation: the notion of “scale” is physically encoded in the architecture, with clear separations between low-level and high-level representations. When you probe features layer by layer, you observe clean transitions.

Transformers, on the other hand, start with global receptive fields from the very first layer because self-attention mixes information across all positions in one step. Every token can, in principle, interact with every other token regardless of spatial proximity. Depth therefore does not correspond to expanding spatial scale but instead refines representations through repeated global mixing and specialization. Each layer develops different attention patterns and specialized circuits (like induction heads or copy-suppression heads), rather than encoding progressively larger spatial features. Depth adds precision and compression, not hierarchical abstraction.

This makes transformer probing results very different. Since there is no strict bottom-to-top feature pyramid, probes do not show a monotonic increase in semantic abstraction. Instead, they show nonmonotonic behavior: mid layers often peak in semantic information, and late layers compress

features toward task-specific prediction spaces. Representations are distributed, overlapping, and non-hierarchical—each layer contributes globally but in subtly different ways.

Mathematically, this difference arises because convolution enforces spatial locality and translation equivariance via weight sharing and limited receptive fields, while self-attention defines a global kernel $A = \text{softmax}(QK^\top / \sqrt{d})$ that mixes all spatial positions. Hence, the “hierarchy” in CNNs is an emergent geometric property of the convolution operator, whereas in transformers the representational geometry is depth-wise iterative refinement within a globally connected graph.

F.6 T-SNE VISUALIZATION AND REPRESENTATION INTERPRETABILITY

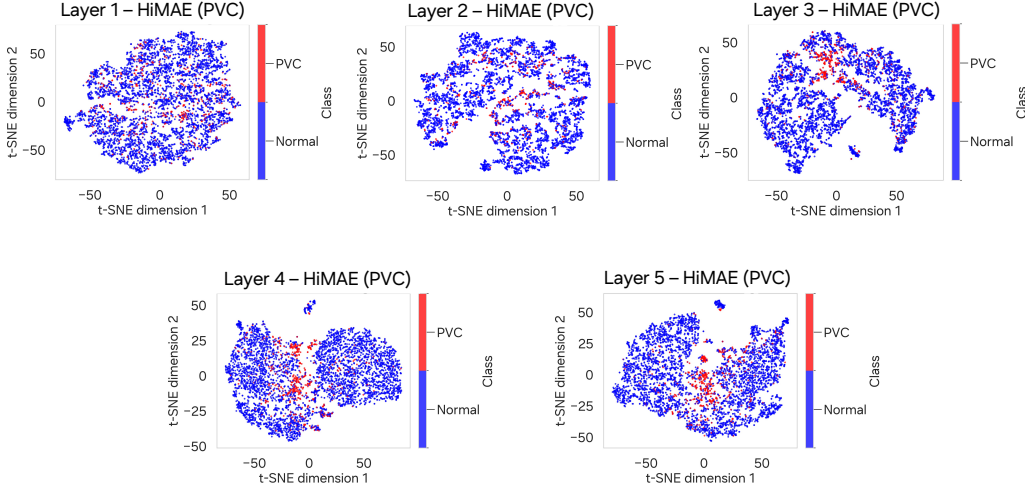


Figure 14: **t-SNE Visualization and Representation Interpretability.** We explore how the layer wise TSNE evolve on the PVC Classification task to help us understand how the representations are organized as they traverse the multiple encoders.

To characterize how HiMAE’s hierarchical representations evolve across depth, we visualize t-SNE projections of encoder embeddings from each layer on the PVC detection task in Figure 14. The corresponding AUROC scores (L1: 0.73, L2: 0.72, L3: 0.77, L4: 0.80, L5: 0.80; Figure 5) closely mirror the progressive emergence of class separability observed across layers.

Representations from early layers (L1–L2) form diffuse clusters with limited separation between normal and PVC segments, indicating that these layers predominantly encode local waveform morphology. At intermediate depth (L3), distinct PVC clusters begin to emerge, coinciding with the first substantial increase in AUROC and marking a transition toward rhythm-level abstraction. Deeper layers (L4–L5) exhibit compact, well-separated clusters, consistent with representations that integrate longer temporal context and capture higher-level cardiac dynamics relevant for arrhythmia discrimination.

Together, these visualizations provide an interpretable view of HiMAE’s hierarchical inductive bias: representations progressively abstract temporal information from fine-grained morphology to broader physiological context. The layer-wise evolution of t-SNE structure offers empirical support for the resolution hypothesis, suggesting that higher layers encode slower, more discriminative temporal processes that ultimately drive clinical performance.

G ON-DEVICE EXPERIMENTS

G.1 EXPERIMENTAL PROTOCOL

We evaluated the on-device performance of HiMAE using a Samsung Galaxy Watch 8 running Wear OS. All experiments were performed natively on-device to capture realistic latency and throughput characteristics under wearable hardware constraints. The model was deployed using PyTorch Mobile with TorchScript conversion to minimize runtime overhead and ensure compatibility with ARM-based computation. The device is powered by the Exynos W1000 chipset, featuring a 5-core CPU (1× Cortex-A78, 3× Cortex-A55, 1× Cortex-M55) fabricated on a 3 nm GAA process, and equipped with 2 GB of LPDDR5 memory.

Inference was performed at a fixed batch size of 1, corresponding to a 10-second physiological signal window sampled at 100 Hz. To ensure measurement stability, we used 20 warm-up runs followed by 100 timed inference passes. Latency was defined as the mean per-sample forward-pass time, with additional reporting of median and 95th percentile values to capture tail latencies. Throughput was defined as the total number of samples processed per second over a 10-second rolling interval.

For timing measurements, we used CUDA event synchronization on GPU and Python’s high-resolution wall-clock timers on CPU. All computations were executed using float32 precision. The benchmarking routines are provided in the accompanying code listings, which include throughput and latency measurement functions.

As a point of reference, we additionally report datacenter-grade inference metrics obtained using an NVIDIA T4 GPU to contextualize the mobile device performance. Although the T4 operates at higher power, modern mobile GPUs (e.g., Qualcomm Adreno 750) demonstrate comparable inference efficiency per watt (Buber & Banu, 2018; Wesolowski et al., 2021), validating the relevance of on-device inference as a proxy for real-world deployment on consumer hardware.

Throughput Code

```

1  def measure_throughput(model, dummy_input, device, num_seconds=10):
2      """Measures inference throughput (samples/sec) with a fixed batch
3          size of N."""
4
5      model.eval()
6      batch_size = N
7      dummy_input = dummy_input.to(device)
8      model.to(device)
9
10     with torch.no_grad():
11         for _ in range(10):
12             _ = model(dummy_input)
13         if device.type == 'cuda':
14             torch.cuda.synchronize()
15
16     num_inferences = 0
17     start_time = time.time()
18     with torch.no_grad():
19         while time.time() - start_time < num_seconds:
20             _ = model(dummy_input)
21             if device.type == 'cuda':
22                 torch.cuda.synchronize()
23             num_inferences += 1
24
25     total_time = time.time() - start_time
26     throughput = num_inferences * batch_size / total_time
27     return {"Throughput_samples_per_sec": throughput}

```

Latency Code (Batch Size = 1)

```

1  def measure_inference_time_bs1(model, dummy_input, device):
2      """Measures inference latency (mean, median, 95th percentile) with
3          a fixed batch size of 1."""
4      if dummy_input.shape[0] != 1:
5          raise ValueError("Input batch size must be 1 for this
6              function. Got {dummy_input.shape[0]}")
7
8      model.eval()
9      dummy_input = dummy_input.to(device)
10     model.to(device)
11
12     with torch.no_grad():
13         for _ in range(warmup_runs):
14             _ = model(dummy_input)
15     if device.type == 'cuda':
16         torch.cuda.synchronize()
17
18     timings = []
19     with torch.no_grad():
20         for _ in range(num_runs):
21             if device.type == 'cuda':
22                 start_event = torch.cuda.Event(enable_timing=True)
23                 end_event = torch.cuda.Event(enable_timing=True)
24                 start_event.record()
25                 _ = model(dummy_input)
26                 end_event.record()
27                 torch.cuda.synchronize()
28                 elapsed_time_ms = start_event.elapsed_time(end_event)
29             else:
30                 start_time = time.time()
31                 _ = model(dummy_input)
32                 end_time = time.time()
33                 elapsed_time_ms = (end_time - start_time) * 1000
34             timings.append(elapsed_time_ms)
35
36     mean_latency = np.mean(timings)
37     median_latency = np.median(timings)
38     std_latency = np.std(timings)
39     p95_latency = np.percentile(timings, 95)
40     return {
41         "Mean_Latency_ms": mean_latency,
42         "Median_Latency_ms": median_latency,
43         "P95_Latency_ms": p95_latency,
44     }

```

G.2 INFERENCE EFFICIENCY

We benchmarked the inference efficiency of our proposed HiMAE against the transformer baseline (MAE-1D), measuring three aspects: model footprint and computational complexity in terms of parameters, memory, and FLOPs per 10-second input window at 100 Hz (Table 18); latency, defined as mean per-sample forward-pass time at batch size 1; and throughput, defined as the maximum number of samples processed per second (Table 19).

Results Despite being more than two orders of magnitude smaller in parameter count, the HiMAE consistently outperforms the transformer baseline across all efficiency metrics. Between Efficient-Net (Tan & Le, 2020), it remains marginally better which is encouraging due to the optimizations designed in this model.

Model footprint: HiMAE reduces parameters from 110M to 0.31M ($\sim 355\times$ fewer), FLOPs from 15.94G to 0.0647G ($\sim 246\times$ fewer), and memory from 441.3MB to 3.6MB ($\sim 123\times$ smaller). These reductions highlight that computational savings scale with the compactness of the model, without loss of representational capacity for the task.

Latency: HiMAE achieves substantially faster per-sample inference. On GPU, latency drops from 0.80ms to 0.039ms ($\sim 20\times$ faster), while on CPU it falls from 3.93ms to 0.99ms ($\sim 4\times$ faster). The reduction in latency follows directly from the smaller computational footprint, reflecting a consistent efficiency advantage.

Throughput: These improvements translate into higher throughput across hardware. On GPU, throughput increases from 1.24k to 25.8k samples/s ($\sim 21\times$ higher), while CPU throughput rises from 0.255k to 1.2k samples/s ($\sim 5\times$ higher). These results confirm that computational gains extend beyond memory and FLOPs, yielding end-to-end speedups at inference time.

In summary, HiMAE achieves a favorable tradeoff between compactness and efficiency, providing lower FLOPs, smaller memory footprint, and faster inference despite its reduced model size. It also outperforms Efficient-Net B1 which was specially designed and optimized for performance and compactness giving a comparison and context to our models performance.

Model	Params	FLOPs	Memory
HiMAE	1.2M	0.0647 gFLOPS	4.8 MB
Efficient-Net	7.8M	0.70 gFLOPS	31.1 MB
Swin-Transformer	110.6M	11.89 gFLOPS	423.8 MB
MAE-1D	110.6M	15.94 gFLOPS	441.3 MB

Table 18: **HiMAE is lightweight and efficient:** Model size and compute cost comparison between HiMAE and MAE-1D. FLOPs measured per forward pass on a 10s sequence at 100Hz.

Model	GPU Lat.	GPU Thr.	CPU Lat.	CPU Thr.
HiMAE	0.039 ms	25.8k/s	0.99 ms	1.2k/s
Efficient-Net	0.082 ms	12.2k/s	1.42 ms	0.704k/s
Swin-Transformer	0.704 ms	1.42k/s	2.95 ms	0.456k/s
MAE-1D	0.80 ms	1.24k/s	3.36 ms	0.298k/s

Table 19: **Inference Performance:** Latency (ms per sample, batch size 2048) and throughput (samples/sec) measured over 10 s windows.

H FREQUENTLY ASKED QUESTIONS

What are the main conclusions from this work? We demonstrate that convolutional architectures benefit from inductive biases that remain advantageous for PPG signals. On our pre-training data, our model consistently outperforms alternative baselines. Furthermore, scaling experiments across model sizes reveal that brute-force scaling of generic architectures is possible, but less effective: our model achieves stronger performance and scales more gracefully due to a better initialization and inductive structure relative to other models. In addition to this inductive bias and compact design, our contributions are two fold in the sense that our model demonstrates the first on-device model which does not require phone level processors to run inference.

Is your pre-training dataset large enough? Our pre-training corpus was collected internally and is of comparable scale to recent public benchmarks such as PaPaGei and Apple’s datasets. In terms of magnitude, we position our dataset as PaPaGei (Pillai et al., 2025) < Ours < Apple (Abbaspourazad et al., 2023) < Google (Narayanswamy et al., 2024). Thus, while not the largest available, our dataset size is sufficiently large to validate the approach and lies within the range of accepted practice for self supervised learning wearable models.

Why do you model at 10-second windows? We deliberately adopt 10s windows sampled at 100Hz to balance physiological coverage with on-device feasibility. Many clinically actionable events, such as arrhythmic beats or premature ventricular contractions, unfold on the order of seconds and require rapid detection to enable continuous monitoring and real-time feedback. Shorter windows would impair the model’s ability to capture meaningful temporal context, while much longer windows would hinder low-latency inference on watch-class hardware. By constraining the receptive field to 10s, HiMAE preserves second-level resolution while remaining efficient enough to process signals continuously under the hardware limits of edge devices. Additionally, 10-second window are a standard protocol that are adopted in the clinical setting where ECG for example is collected and interpreted at 10 second segments (Shuai et al., 2016).

What are the advantages of smaller models? From a research perspective, smaller models foster inclusivity by reducing reliance on brute-force scaling of transformer-based architectures that only industry-scale labs can realistically afford. From a deployment standpoint, compact models enable on-device inference on constrained hardware such as wearables. This dual benefit—lower research barriers and wider deployment potential—underscores the importance of investigating architectures that remain competitive at modest scale.

How large is too large to deploy on a smart watch? In principle, models up to approximately 50MB can be stored and executed on modern smart watches or larger models can be quantized (Jacob et al., 2017). In practice, however, latency and energy considerations suggest that models exceeding roughly 10MB may already hinder real-time inference and limit commercial viability. Additionally quantization does not do due diligence to the original model and some level of the model’s performance is lost. While smartphones relax these constraints, our contribution highlights that the proposed model remains sufficiently compact to fit within the computational and storage budgets of wearable devices such as watches, thereby supporting direct on-device deployment.

Can PPG predict abnormal laboratory results? We frame this as a binary classification task, testing whether photoplethysmography signal encodes biomarkers that separate “normal” from “abnormal” lab classes. Our investigation probes whether learned PPG representations capture biomarker signatures correlated with out-of-range labs, using lightweight classifiers on frozen embeddings with strict temporal alignment. Preliminary evidence suggests discriminative signal above chance, but these findings are designed to be exploratory and not clinically actionable.