

Received 20 February 2025, accepted 28 March 2025, date of publication 1 April 2025, date of current version 11 April 2025.

Digital Object Identifier 10.1109/ACCESS.2025.3556957

## RESEARCH ARTICLE

# PFML: Self-Supervised Learning of Time-Series Data Without Representation Collapse

EINARI VAARAS<sup>1</sup>, MANU AIRAKSINEN<sup>2</sup>, AND OKKO RÄSÄNEN<sup>1</sup>, (Senior Member, IEEE)

<sup>1</sup>Signal Processing Research Centre, Tampere University, 33720 Tampere, Finland

<sup>2</sup>BABA Center, University of Helsinki, 00029 Helsinki, Finland

Corresponding author: Einari Vaaras (einari.vaaras@tuni.fi)

This work was supported in part by the Research Council of Finland under Grant 343498, and in part by the Sigrid Jusélius Foundation.

**ABSTRACT** Self-supervised learning (SSL) is a data-driven learning approach that utilizes the innate structure of the data to guide the learning process. In contrast to supervised learning, which depends on external labels, SSL utilizes the inherent characteristics of the data to produce its own supervisory signal. However, one frequent issue with SSL methods is representation collapse, where the model outputs a constant input-invariant feature representation. This issue hinders the potential application of SSL methods to new data modalities, as trying to avoid representation collapse wastes researchers' time and effort. This paper introduces a novel SSL algorithm for time-series data called Prediction of Functionals from Masked Latents (PFML). Instead of predicting masked input signals or their latent representations directly, PFML operates by predicting statistical functionals of the input signal corresponding to masked embeddings, given a sequence of unmasked embeddings. The algorithm is designed to avoid representation collapse, rendering it straightforwardly applicable to different time-series data domains, such as novel sensor modalities in clinical data. We demonstrate the effectiveness of PFML through complex, real-life classification tasks across three different data modalities: infant posture and movement classification from multi-sensor inertial measurement unit data, emotion recognition from speech data, and sleep stage classification from EEG data. The results show that PFML is superior to a conceptually similar SSL method and a contrastive learning-based SSL method. Additionally, PFML is on par with the current state-of-the-art SSL method, while also being conceptually simpler and without suffering from representation collapse. The code is freely available at <https://github.com/SPEECHCOG/PFML>.

**INDEX TERMS** EEG data, embedding masking, human activity recognition, multi-sensor inertial measurement unit data, representation collapse, self-supervised learning, sleep stage classification, speech emotion recognition, statistical functionals, time-series data.

## I. INTRODUCTION

Self-supervised learning (SSL) can be described as a data-driven learning paradigm where the training process is guided by the inherent structure of the data itself. Unlike supervised learning that relies on externally provided labels, SSL exploits the intrinsic properties of the data to generate its own supervisory signal [1], [2]. SSL enables the model to learn

The associate editor coordinating the review of this manuscript and approving it for publication was Inês Domingues<sup>1</sup>.

rich feature representations from large amounts of unlabeled data that can be used as a starting point for downstream tasks, either as such or by fine-tuning the feature extractor to be better suited for solving some specific task [2], [3]. Since typically there is an abundance of unlabeled data but a scarcity of labeled data, the use of SSL has been shown to reduce the need for large, manually annotated datasets [4], [5], [6]. In addition to SSL algorithms that have been developed for a single data modality, SSL algorithms that can be applied to multiple different data modalities have gained popularity

in recent years [2], [4], [7], [8], [9]. These methods and their extensions have shown great success in e.g. audio, image, and text data [4], [7], [8], [9], [10], [11], [12], [13], [14].

However, many SSL algorithms suffer from two issues: First, SSL algorithms are usually complex, with a plethora of hyperparameters that need careful tuning for the algorithm to work properly. This hinders the ability of SSL algorithms to be applied to new data domains, where the selection of these hyperparameters is not self-evident. For example, in contrastive learning-based SSL, the selection of positive and negative samples during training is essential for the algorithm to work properly. However, deciding which samples should be assigned to positive and negative categories is not always apparent [1], [15], [16]. As another example, determining the number of clusters for clustering-based SSL algorithms (such as Caron et al. [17] and Hsu et al. [18]) in a new data domain or task can be difficult. Examples of such domains could include, for instance, different types of medical time-series data (e.g. electroencephalography (EEG), electrocardiography (ECG), or electromyography (EMG) recordings) that come in various dataset sizes and from various recording configurations. Second, a common failure mode during SSL pre-training is representation collapse, where the model ends up outputting a constant, time-invariant feature representation. Representation collapse is very common in SSL pre-training [1], [19], [20], [21], and many SSL methods apply different countermeasures to tackle the problem (see Section III-A).

In the present study, we propose a new SSL algorithm for time-series data called Prediction of Functionals from Masked Latents (PFML). In PFML, the aim is to predict statistical functionals of the input signal corresponding to masked embeddings, given a sequence of unmasked embeddings. The overall methodological aim of our method is to have an SSL algorithm that would be as straightforward as possible to apply to various time-series data domains with minimal hyperparameter optimization, and without the risk of representation collapse. The contributions of the present study are as follows:

- 1) We propose a novel SSL algorithm for time-series data, PFML, that does not suffer from representation collapse, rendering the method straightforward to apply to new time-series data domains. To the best of our knowledge, PFML is the first work within the field of SSL for time-series data where the central idea of reconstructing statistical functionals is utilized.
- 2) We demonstrate the effectiveness of PFML using three different data modalities with complex, real-life classification tasks: infant posture and movement classification from multi-sensor inertial measurement unit (IMU) data, emotion recognition from speech data, and sleep stage classification from EEG data.
- 3) We show that PFML obtains superior results against both a conceptually similar SSL method and a contrastive learning-based SSL method. Additionally, PFML achieves results that are on par with the current

state-of-the-art data modality-agnostic SSL method, while also being conceptually simpler and without suffering from representation collapse.

## II. RELATED WORK

Most of the advances in SSL have focused on developing new, better-performing algorithms with some specific data modality in mind. For speech data, Baevski et al. [5] presented an SSL algorithm where the basic idea is to mask speech embeddings and then solve a contrastive task that is defined over a quantization of the embeddings which are simultaneously learned during the pre-training task. Hsu et al. [18] proposed that instead of solving a contrastive task, they predict cluster targets of masked embeddings. Furthermore, the SSL method by Chen et al. [22] also uses masking of embeddings, but the authors simulate noisy speech inputs and predict pseudo-labels of the original speech from the masked embeddings.

Similar to the advances in SSL for speech data, there have been significant developments in SSL for image data as well [6], [23], [24], [25], [26], [27], [28], [29], [30]. Grill et al. [26] presented an SSL method that uses two neural networks that learn from each other's representations of differently augmented views of the same image. He et al. [28] proposed masked autoencoders (MAE) that try to reconstruct masked patches of input images using an asymmetric encoder-decoder architecture. The SSL algorithm by Bao et al. [29] tokenizes images into visual tokens, followed by masking some image patches and then trying to recover the original tokens from the masked patches.

SSL has also excelled in natural language processing [31], [32], [33], [34]. Devlin et al. [31] introduced an SSL method which obtains bidirectional feature representations from unlabeled text by conditioning on both the left and right textual context. The method by Brown et al. [32] uses an autoregressive model which alternates dense and locally banded sparse attention patterns in their Transformer model. OpenAI [34] proposed an expanded version of the method by Brown et al. [32] by making the model not only larger, but also capable of handling image inputs in addition to text inputs.

More recently, SSL literature has seen a growing number of work towards SSL algorithms capable of running the pre-training task on multiple different data modalities. The authors of van den Oord et al. [4] developed an SSL approach that predicts future embeddings based on previous context using contrastive learning. They showed that their method was able to learn useful feature representations for audio, image, text, and reinforcement learning in 3D environments. The SSL method by Akbari et al. [7] also uses contrastive learning, but their method simultaneously takes audio, video, and text data as input and creates multimodal feature representations. These features were shown to work well with multiple different downstream tasks, i.e. video action recognition, audio event classification, image classification, and text-to-video retrieval. Baevski et al. [8] proposed data2vec, an SSL method for audio, image, and text data. In their approach, the model tries to predict masked latent features of an older

version of itself that are both normalized and averaged over multiple Transformer layers. Their results in downstream tasks demonstrate the effectiveness of the method in all three data modalities. Yue et al. [35] presented a universal framework, TS2Vec, for learning time-series representations at different semantic levels using hierarchical contrastive learning over augmented context views. Their approach enables robust contextual representations for each timestamp and allows for simple aggregation to obtain representations of arbitrary sub-sequences. They demonstrated that their method obtained state-of-the-art performance in time-series classification, forecasting, and anomaly detection tasks on multiple datasets, most of which were small in terms of the number of training samples. Notably, TS2Vec achieved these results by adding very simple classifiers after the pre-trained model, which has fixed constraints on the model architecture to model the hierarchical representations. Wang et al. [9] proposed an SSL method that performs prediction of masked tokens in a unified manner on images, texts, and image-text pairs. Their experiments showed that their method achieves state-of-the-art performance on various vision and vision-language tasks.

For modality agnostic SSL algorithms, objective functions play a crucial role in guiding the learning process. These functions can be broadly categorized into three types: instance discrimination, clustering, and masked prediction. Instance discrimination aims to distinguish between different instances of data, thereby encouraging the model to learn unique features for each instance and enhancing the discriminative power of the learned representations. Contrastive learning methods, such as [4], [5], [7], [35], [36], and [37], are an example of instance discrimination-based SSL methods. Clustering, on the other hand, groups similar instances together in the feature space, fostering the model to learn common features among instances belonging to the same group. Methods like [17], [18], and [38] are examples of clustering-based SSL methods. Lastly, masked prediction involves the task of predicting masked parts of the input data based on the unmasked parts, thereby encouraging the model to learn contextual relationships within the data. Examples of such SSL methods include [8], [28], [31], [39], [40].

### III. METHOD

#### A. MOTIVATION

One key issue with many SSL methods is the problem of *representation collapse*, where the model outputs a constant, input-invariant feature representation, leading to a trivial solution of the pre-training task [1], [20]. This considerably slows down the development process for novel data domains and/or tasks due to the necessity of operating in uncertainty, when it is not clear whether the representation collapse is caused by an ill-posed task or by the SSL algorithm. To avoid this, SSL methods have taken several different countermeasures: Baeviski et al. [5] use the same target representations in their contrastive learning task in a dual

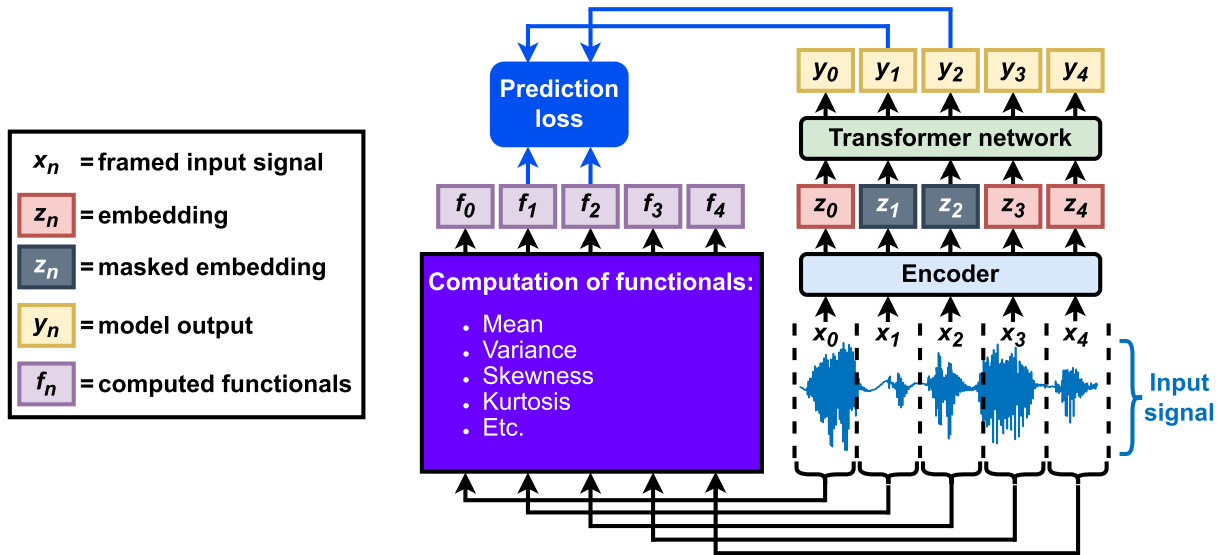
manner, i.e. both as a positive and a negative example. Grill et al. [26] both add an additional predictor to their training regime and use a moving average of their so-called online neural network to avoid representation collapse. Bardes et al. [41] add a regularization term to their loss function that both maintains variance of the embeddings and decorrelates each pair of variables. In data2vec [8], the authors tackle representation collapse by carefully selecting their model hyperparameters and promoting target representation variance through feature normalization. Also, in the code implementation of data2vec,<sup>1</sup> pre-training is stopped if the variance of either model predictions or training targets falls below a predefined threshold.

Intuitively, given a trivial task, the model does not learn useful feature representations during pre-training. In contrast, if the learning objective is too complicated, the model fails to converge to a useful solution. For time-series data, i.e. a waveform (e.g. audio) or a set of waveforms (e.g. multi-channel EEG), trying to reconstruct masked parts of the input signal given the unmasked parts of the signal (as in e.g. MAE [28]) is a very complex task. This is due to the fact that a time-series signal can have large temporal variation even between short periods of time. While joint learning of *a priori* unspecified latent representations and their prediction allows discarding of this irrelevant variation (as in, e.g., [4] or [5]), the problem requires learning algorithms that become susceptible to representation collapse and/or may require careful tuning of the training process. We hypothesize that for SSL pre-training with time-series data, a model would learn more useful features for downstream tasks if the complex setting of MAE would be alleviated slightly. Hence, we propose PFML, a novel SSL algorithm for time-series data. Our method builds on the concept of MAE and reduces the complexity of the pre-training task of MAE in two ways:

- 1) Instead of aiming to reconstruct the input signal, the model tries to predict a set of statistical functionals computed from the input signal.
- 2) Instead of masking the input signal directly, PFML borrows the idea of e.g. wav2vec 2.0 [5] and data2vec [8] and masks the embeddings created by the encoder model.

Regarding point (1), by making the model predict statistical functionals of masked latent features instead of predicting the input signal  $\mathbf{x}$  itself, we relieve the model from the complex task of modelling the high-dimensional distribution of  $\mathbf{x}$  in detail. We validate this argument of generating better features for downstream tasks by reducing the computational complexity of the pre-training task in Section IV, where we compare our proposed method against MAE. In theory, the set of statistical functionals can be chosen so that the desired and deterministically calculated statistical properties of the data, and thereby their variance, are preserved in the target features. Furthermore, regarding point (2), we show in our experiments in Section IV that it is more beneficial during pre-training to

<sup>1</sup><https://github.com/facebookresearch/fairseq/tree/main/examples/data2vec>



**FIGURE 1.** An overview of the PFML pre-training pipeline. Note that in the figure, the input signal has only a single channel, whereas PFML can also be applied to multi-channel time-series data.

mask the latent features instead of masking the input directly. This further alleviates the complexity of the learning task in particular for the encoder module.

**B. PREDICTION OF FUNCTIONALS FROM MASKED LATENTS**

Figure 1 depicts an overview of the PFML pre-training pipeline. First, a single- or multi-channel signal  $\mathbf{x}$  is framed into a sequence of short-term frames  $\{\mathbf{x}_0, \mathbf{x}_1, \dots\}$ ,  $\mathbf{x}_n = \{x_t, x_{t+1}, \dots, x_{t+N-1}\}$ , of  $N$  samples each. Then, a set of  $m$  functionals,  $\mathcal{F} = \{F_0, F_1, \dots, F_{m-1}\}$ , is computed for each frame  $\mathbf{x}_n$  to produce corresponding functional values  $\mathbf{f}_n = \{F_0(\mathbf{x}_n), F_1(\mathbf{x}_n), \dots, F_{m-1}(\mathbf{x}_n)\}$ . Here, functionals are defined as mathematical operations which map a time series of arbitrary length into a single value, such as the mean or variance of the signal. The frames  $\mathbf{x}_n$  are also fed to an encoder model, which converts the framed signals into embeddings  $\mathbf{z}_n$ . Some of these embeddings are masked randomly at time steps  $M$  (for example,  $M \in \{1, 2\}$  in Figure 1), after which all  $\mathbf{z}_n$  are used as an input for a Transformer-based model to obtain outputs  $\mathbf{y}_n$ . Finally, a prediction loss is computed between the outputs of masked time steps  $\mathbf{y}_M$  and their functional counterparts  $\mathbf{f}_M$ . As a result, PFML pre-training optimizes the prediction of functionals of input signal frames corresponding to the masked embeddings, given the unmasked embeddings from the temporal context of these frames.

In PFML, predicting only one or a few functionals of a framed signal can be a trivial task, and will most probably lead to learning feature representations that are not very useful for downstream tasks. However, as the number of functionals that each describe some property of the framed signal grows, a more accurate description of the signal can be obtained (see e.g. McDermott & Simoncelli [42] for representing perceptual properties of sound textures with functionals). Therefore, as the number of different functionals grows, the PFML

algorithm is getting closer to predicting all of the nuances of the input signal.

Let us assume the following in PFML pre-training:

- Assumption 1: There is temporal variability across the frames  $\mathbf{x}_n$ . This assumption is reasonable as real-world data typically exhibits temporal variability.
- Assumption 2: Given Assumption 1, a set of non-trivial functionals  $\mathcal{F}$  computed from  $\mathbf{x}_n$  also contains variance across the frames. This follows naturally since non-constant functionals derived from variable data also exhibit variability.

Under these assumptions, as the model is trying to predict the computed functionals  $\mathbf{f}_n$  given the embeddings  $\mathbf{z}_n$ , good model predictions  $\mathbf{y}_n$  that lead to low prediction loss values also inherently contain variance. On the contrary, if  $\mathbf{y}_n$  were to contain zero variance across the frames while  $\mathbf{f}_n$  contains variance, the prediction loss would be high. Consequently, PFML pre-training does not converge to collapsed feature representations, as long as Assumptions 1 and 2 hold true. For a more detailed formulation, see Appendix A in the supplementary material. Empirical results (see Section V) support this theoretical claim, showing that PFML maintains variance in predictions across various datasets.

In the present study, we selected 11 mathematical operations as our set of functionals: mean, variance, skewness, kurtosis, minimum value, maximum value, zero-crossing rate (ZCR), and the mean, variance, skewness, and kurtosis of the autocorrelation function (ACF). The ZCR for a signal  $\mathbf{x} = \{x_0, x_1, \dots, x_{N-1}\}$  is defined as

$$\text{ZCR}(\mathbf{x}) = \frac{1}{N-1} \sum_{k=1}^{N-1} |\text{sgn}(x_k) - \text{sgn}(x_{k-1})| \quad (1)$$

where  $\text{sgn}$  denotes the sign function [43]. The ACF for a signal  $\mathbf{x}$  at lag  $\tau$  is defined as

$$\text{ACF}(\mathbf{x}, \tau) = \frac{1}{(N - \tau)\sigma^2} \sum_{k=0}^{N-\tau-1} (x_{k+\tau} - \mu)(x_k - \mu) \quad (2)$$

where  $\tau < N$ ,  $\mu$  is the mean of  $\mathbf{x}$ , and  $\sigma^2$  is the variance of  $\mathbf{x}$  [43]. Note that Equation 2 returns a vector of measurements when applied to all lags  $\tau < N$ .

For masking the embeddings, in each training and validation minibatch we randomly select frames with a probability of  $p_m$  to be mask starting indices, and we mask the embedding of that frame and  $l_m - 1$  subsequent frames, resulting in a minimum mask length of  $l_m$  frames. We replace each embedding that is selected for masking with a vector of ones. Masks can overlap, enabling longer mask spans than  $l_m$  frames (especially with high  $p_m$ ). Furthermore, we also define that each training and validation sequence needs to have at least one mask starting index during PFML pre-training.

Note that the PFML pre-training process is not restricted to any specific type of neural networks. In the present study, we used convolutional neural networks (CNNs) as our encoder model, and  $T$  Transformer encoder blocks as the temporal model. However, any type of encoder could be used for PFML, as long as the encoder can convert time-series data into a sequence of embeddings. Furthermore, other temporal models, such as conformer-based models [44] or bidirectional recurrent neural networks [45], [46], could also be used for PFML, as long as the model is able to take contextual information into account.

## IV. EXPERIMENTS

We evaluate our PFML method using three different datasets of time-series data with complex classification tasks: infant posture and movement classification from multi-sensor IMU data, emotion recognition from speech data, and sleep stage classification from EEG data. For each dataset, we first run SSL pre-training with unlabeled data using PFML, after which we fine-tune our models for downstream classification tasks using labeled data. We compare PFML against four different baselines: MAE [28], data2vec [8], TS2Vec [35], and not using pre-training at all. We selected MAE for our experiments since it is conceptually very similar to PFML, and we chose data2vec since it is the current state-of-the-art data-modality-agnostic SSL method. We also included the conceptually different TS2Vec in the present experiments due to its reported state-of-the-art performance on multiple different datasets of single- and multi-channel time-series data. In order to make the prediction of functionals directly comparable with predicting the input signal, we use a slightly modified version of MAE where we mask embeddings instead of masking inputs.

This section is organized as follows: First, Section IV-A gives general-level information on the pre-training and fine-tuning experiments regarding PFML, MAE, and data2vec methods that utilize identical neural network architecture for the end-use system. This is followed by

Sections IV-B, IV-C, and IV-D, which give modality-specific experimental details for IMU, speech, and EEG data, respectively. Finally, Section IV-E explains the key differences between TS2Vec-related experiments and the experiments described in Section IV-A.

### A. PFML, MAE, AND DATA2VEC EXPERIMENTS

For PFML pre-training, our models consist of a modality-specific frame-level encoder (detailed in Sections IV-B, IV-C, and IV-D for IMU, speech, and EEG data, respectively) and a Transformer network consisting of  $T$  Transformer encoder blocks. Between the encoder and Transformer networks there is a CNN-based relative positional encoder followed by a Gaussian error linear unit (GELU) [47] activation and layer normalization [48]. We frame our input signals before feeding the data into an encoder model, and we compute functionals from these frames as our training targets. For multi-channel data, we compute functionals separately for each channel. The functionals are then z-score normalized across the entire pre-training dataset. For computational efficiency, we pre-compute the functionals of each signal frame before the pre-training process. After the Transformer encoder blocks, we add a linear projection to convert the Transformer outputs into predicted functionals. After pre-training, this linear projection is discarded. Pre-training is run until validation loss convergence, and we use the model with the lowest validation loss as our pre-trained model. Starting from an initial learning rate, we gradually reduce the learning rate during model training with a reduction factor of 0.5 based on the plateauing of the validation loss.

We pre-train our models using MAE and data2vec in a similar manner as for PFML, and we use the same model architecture for all three pre-training algorithms. MAE pre-training is run in a similar manner as PFML pre-training, with the only exception of predicting the input signal frames instead of functionals. For data2vec pre-training, we used the instance-normalized [49] and averaged outputs of each feed-forward part of all Transformer encoder blocks as our training targets. If we observed that a representation collapse occurred during data2vec pre-training, we restarted the pre-training process. For further details on the data2vec algorithm, see Baevski et al. [8]. We used mean squared error loss for all pre-training processes except for PFML with speech data, where we found L1 loss to work better.

We fine-tune our pre-trained models in two stages. In the first stage, two randomly initialized fully-connected GELU layers followed by a softmax function are added after the Transformer model. Then, these layers are fine-tuned separately as the weights of the encoder and Transformer are frozen. In the second stage, the entire model is fine-tuned with the same hyperparameters as in the first fine-tuning stage with one exception: The learning rate  $\eta$  is linearly increased from  $0.001 \cdot \eta$  to  $\eta$  during a warm-up period of 20 training epochs, followed by reduction by a factor of 0.5 based on validation loss plateauing. We use weighted categorical cross-entropy

loss by weighting the loss of each sample by its class' inverse frequency.

We also test the linear separability of the features learned by our pre-trained models. In this case, we only add one linear layer followed by a softmax function after the Transformer model, and we fine-tune this single layer while the weights of the encoder and Transformer are frozen. As a baseline, we perform the same linear evaluation for a randomly initialized model without any pre-training.

In order to demonstrate the superiority of PFML against the state-of-the-art SSL method for multiple data modalities, `data2vec`, in terms of representation collapse, we ran each of the pre-training algorithms of the present experiments 10 times using the best hyperparameter combinations for each SSL method and for each data modality. We defined representation collapse to have occurred if the variance of either the embeddings or model outputs fell below 0.01 for 10 consecutive pre-training epochs, during which the validation loss was decreasing. In our preliminary experiments, we found that this condition was a good indicator of an upcoming representation collapse: A systematic decrease in the variance of a model's embeddings or outputs indicates impending representation collapse in SSL methods where the model can invent its own training targets.

For pre-training, we use the RAdam [50], AdamW [51], and Adam [52] optimizers for PFML, MAE, and `data2vec`, respectively. For fine-tuning, we use the Adam optimizer. For the "no pre-training" condition, we simply omit pre-training, the first fine-tuning stage, and the learning rate warm-up period of the second fine-tuning stage. To demonstrate fair comparison, we carefully select the pre-training and fine-tuning hyperparameters for each SSL method separately in order to minimize the number of representation collapses during SSL pre-training, and to maximize the fine-tuning performance. For a complete list of pre-training and fine-tuning hyperparameters, refer to Appendix B in the supplementary material. We used an NVIDIA Tesla V100 GPU to train our models, and we implemented the code using PyTorch version 1.13.1. Our implementation is publicly available on GitHub.<sup>2</sup>

### B. INFANT POSTURE AND MOVEMENT CLASSIFICATION

For infant posture and movement classification, we use the multi-sensor IMU data from Airaksinen et al. [53]. The data contains 24-channel signals from infants (three gyroscope and three accelerometer channels, four limbs) with a sampling rate of 52 Hz. We window the signals into 120-sample frames (approx. 2.3 seconds) with 50% overlap. For further details about the dataset, see Airaksinen et al. [53].

For model pre-training, we use a 387-hour subset of a 1333-hour dataset of unlabeled IMU data [54]. This subset contains infant free-form play that has been automatically screened for signal quality, and it was selected for the present experiments based on the findings of Vaaras et al. [54], where using the

subset yielded the best results in terms of SSL pre-training. This subset contains 4669 sequences of 260 consecutive frames, each corresponding to five minutes of data. As the encoder, we use the same four-layer CNN-based encoder architecture as in Airaksinen et al. [53] with three minor modifications that were found to improve training efficiency and system performance when replicating the experiments of Airaksinen et al. [53]: We added layer normalization after the last two convolutions to make the pre-training process more stable, the kernel size of the second convolutional layer of the CNN encoder was changed from [4, 5] to [3, 5], and the originally temporally asymmetrical padding was set to [1, 2] to make it symmetric. The pre-training data is randomly split into a training and validation set in a ratio of 80:20 sequences, and we input 260-frame sequences into the model.

For fine-tuning our pre-trained models, we use a 29-hour (91,449 frames) labeled dataset of IMU data [53] (41 recordings and distinct participants) for two separate tasks: posture classification and movement classification. The data contains nine annotated movement categories (still, roll left/right, pivot left/right, proto/elementary/fluent movement, transition) and seven annotated posture categories (prone, supine, left/right side, crawl posture, sitting, standing) for each 2.3-second frame. For model training, we use all annotated data, but we only use the frames in which all annotators agreed on the label for model testing. We train our models separately for both classification tasks using the so-called iterative annotation refinement labels from Airaksinen et al. [55].

Model fine-tuning is run using recording-level 10-fold cross-validation on the 41 distinct recordings of the labeled dataset. We split each training fold into separate training and validation sets in a ratio of 80:20 recordings. Due to large class imbalances in the labeled dataset of IMU data (see Airaksinen et al. [53]), the unweighted average F1 score (UAF1) is used as our performance metric. We use UAF1 on the validation set as the training criterion, and we select the best-performing model based on validation set UAF1 score. We use random sensor dropout ( $p = 0.3$ ) for data augmentation during model fine-tuning. The final UAF1 score of fine-tuning is computed from an aggregate confusion matrix across all test folds. For further details regarding the pre-training and fine-tuning hyperparameters, see Appendix B in the supplementary material.

### C. SPEECH EMOTION RECOGNITION

We use the 56-hour subset of Finnish speech of the NICU-A corpus [56] for our speech emotion recognition experiments. This subset contains 129,007 utterances with a sampling rate of 16 kHz, of which 5198 and 345 belong to annotated training and testing sets, respectively. Each annotated utterance in NICU-A contains binary labels for emotional valence (positive/non-positive) and arousal (high/low). We window each speech signal into 30-ms frames with a 20-ms overlap. Each sequence is z-score normalized, and we zero-pad or truncate each normalized sequence into 3-second segments

<sup>2</sup><https://github.com/SPEECHCOG/PFML>

(301 frames). See Vaaras et al. [56] for further details on NICU-A.

For model pre-training, we use all 129,007 utterances, and we input 301-frame sequences to our model. We use a four-layer CNN encoder (slightly modified audio encoder of van den Oord et al. [4]) with output channels [128, 128, 128, 128], kernel sizes [10, 8, 4, 4], strides of [5, 4, 2, 2], and paddings of [3, 2, 1, 1]. Each layer is followed by layer normalization, a GELU nonlinearity, and dropout. The last CNN layer is followed by average pooling with a kernel size of 6 before dropout. The pre-training utterances are randomly split into a training and validation set in a ratio of 80:20 sequences.

We fine-tune and test our models separately for both classification tasks (valence/arousal) using the labeled 5198- and 345-utterance training and testing sets, respectively. The training set is randomly split into a training and validation set in a ratio of 80:20 utterances, and we select the best-performing model of the fine-tuning process based on the unweighted average recall (UAR) performance score on the validation set. This model is then used to compute the UAR performance score of the test set. See Appendix B in the supplementary material for further details regarding the pre-training and fine-tuning hyperparameters.

#### D. SLEEP STAGE CLASSIFICATION

For sleep stage classification, we use the pre-processed expanded Sleep-EDF Database [57], [58] from a study by Eldele et al. [59]. The dataset contains 30-second segments of the Fpz-Cz channel with a sampling rate of 100 Hz, comprising a total of 195,479 segments of EEG data. Each 30-second segment belongs to one of five annotated categories: wake, rapid eye movement (REM), non-REM stage 1, non-REM stage 2, or non-REM stages 3 and 4 combined. We z-score normalize each 30-second segment, and we window each segment into 4-second frames with 2 seconds of overlap, resulting into 14 frames for each segment.

We pre-train our models using all 195,479 EEG segments. We use the 14-frame sequences as our input for a three-layer CNN encoder with output channels [128, 128, 128, 128], kernel sizes [10, 8, 4], strides of [5, 5, 3], and paddings of [3, 2, 1]. Each convolution is followed by layer normalization, a GELU nonlinearity, and dropout. The third CNN layer is followed by average pooling with a kernel size of 5 before dropout. We randomly split the EEG segments for pre-training into a training and validation set in a ratio of 80:20 segments.

We fine-tune our models for sleep stage classification using 10-fold cross-validation at the test subject-level on the 78 test subjects of the dataset. Each training fold is split into training and validation sets at the test subject-level in a ratio of 80:20 test subjects. Similar to Sec. IV-B, we use the validation UAF1 score as our training criterion, and the testing UAF1 score is computed from an aggregate confusion matrix across all test folds. For further details on the training hyperparameters, see Appendix B in the supplementary material.

#### E. TS2VEC EXPERIMENTS

We also conducted a similar set of experiments as described in Section IV-A for the TS2Vec method. TS2Vec stands out from the other three SSL methods used in the present experiments due to its hierarchical contrastive learning algorithm, which is integrated into the pre-defined CNN encoder architecture. Therefore, TS2Vec does not allow the use of the same modality-specific encoders used with PFML, MAE, and data2vec in our experiments. Additionally, it is not possible to pre-train the classification model (i.e. Transformer) in TS2Vec. Consequently, for all three data modalities (IMU, speech, and EEG data), we use the pre-defined TS2Vec encoder, and we do not pre-train the Transformer network. In order to compare the linear separability of the SSL features, we add one linear layer directly after the TS2Vec pre-trained encoder, and we fine-tune this single linear layer while the weights of the encoder are frozen.

For TS2Vec pre-training, we used the original TS2Vec implementation<sup>3</sup> with default hyperparameter settings, except for two modifications: First, we observed that the pre-defined number of minibatch updates was insufficient for our pre-training experiments. We found it more effective to run TS2Vec pre-training similarly to PFML, MAE, and data2vec, i.e. until validation loss convergence, and then use the model with the lowest validation loss as the final pre-trained model. Second, we adjusted the output dimensionality of the TS2Vec encoder to match the modality-specific input dimensionalities of the Transformer (see Table 5 of Appendix B in the supplementary material).

For fine-tuning the TS2Vec pre-trained models, we added a randomly initialized Transformer model and two randomly initialized fully-connected GELU layers followed by a softmax function after the TS2Vec pre-trained encoder. The models were then fine-tuned similarly to the process described in Section IV-A, except that in the first fine-tuning stage, both the Transformer and the GELU layers were fine-tuned while the weights of the encoder were frozen. The fine-tuning hyperparameters were the same as those used for the other pre-training algorithms (see Table 6 of Appendix B in the supplementary material).

To demonstrate fair comparison, we also conducted three additional experiments using TS2Vec: First, similar to Section IV-A, we tested the “no pre-training” condition of the neural network architecture combining the TS2Vec encoder and Transformer classifier by fine-tuning the entire model at once using a randomly initialized model. Second, we tested the linear separability of a randomly initialized TS2Vec encoder by adding a single linear layer after the encoder and fine-tuning only this layer. Third, since the authors of the TS2Vec paper [35] used a support vector machine (SVM) classifier with a radial basis function kernel for the SSL features in classification tasks, we adopted the same procedure for our experimental pipeline. All of the TS2Vec-related results are presented in Appendix F in the supplementary material.

<sup>3</sup><https://github.com/zhihanyue/ts2vec>

**TABLE 1.** Downstream task fine-tuning results for PFML, data2vec, MAE, TS2Vec, and not using pre-training at all for the five different classification tasks across the three different data modalities (IMU, speech, and EEG data).

	Multi-sensor IMU data (infant motility assessment)		Speech data (speech emotion recognition)		EEG data (sleep stage classification)
	Movement	Posture	Valence	Arousal	Sleep stage
No pre-training	80.6	94.9	68.2	65.5	69.1
MAE	81.0	95.6	69.9	68.1	70.5
data2vec	<b>81.9</b>	<b>95.8</b>	<b>70.7</b>	68.5	69.8
TS2Vec <sup>a</sup>	73.4	93.4	69.0	66.4	63.0
PFML (ours)	81.8	95.7	<b>70.7</b>	<b>68.6</b>	<b>71.2</b>
	UAF1 (%)		UAR (%)		UAF1 (%)

**TABLE 2.** Linear evaluation results for PFML, data2vec, MAE, TS2Vec, and a randomly initialized model.

	Multi-sensor IMU data (infant motility assessment)		Speech data (speech emotion recognition)		EEG data (sleep stage classification)
	Movement	Posture	Valence	Arousal	Sleep stage
Random initialization	10.8	45.9	51.4	50.8	22.8
MAE	39.9	87.2	60.9	58.8	56.0
data2vec	41.7	87.1	61.8	<b>59.3</b>	53.9
TS2Vec <sup>b</sup>	36.6	83.1	<b>64.9</b>	54.3	<b>57.0</b>
PFML (ours)	<b>43.8</b>	<b>87.4</b>	61.6	59.2	56.3
	UAF1 (%)		UAR (%)		UAF1 (%)

**TABLE 3.** Frequency of representation collapse across 10 runs of PFML, data2vec, MAE, and TS2Vec for each tested data modality.

	Multi-sensor IMU data	Speech data	EEG data
MAE	<b>0/10</b>	<b>0/10</b>	1/10
data2vec	9/10	8/10	8/10
TS2Vec	<b>0/10</b>	<b>0/10</b>	<b>0/10</b>
PFML (ours)	<b>0/10</b>	<b>0/10</b>	<b>0/10</b>

## V. RESULTS

Table 1 presents the fine-tuning results of the comparison of our PFML method against MAE, data2vec, TS2Vec, and not using pre-training at all. Across all three data modalities and five classification tasks, the end-use results show that PFML outperformed MAE and TS2Vec, and achieved results that are on par with data2vec. Using pre-training with any SSL method other than TS2Vec provided superior results as opposed to not using pre-training at all. Furthermore, for the classification of posture from IMU data, which is considered an easy task [53], there were only minor differences in performance between all SSL methods other than TS2Vec. The comparably weak results for TS2Vec are most probably due to the fact that TS2Vec cannot utilize modality-specific encoders, and also since the Transformer model is not pre-trained during TS2Vec pre-training. Using SVM models instead of Transformers in TS2Vec classification, as in the original article of Yue et al. [35], did not improve the results (see Appendix F in the supplementary material). In sleep stage classification from EEG data, both MAE and PFML outperformed data2vec by a large margin. Also, the comparison between PFML and MAE showcases that it is more beneficial to predict functionals than to predict the input signal.

Table 2 shows the results of the linear evaluation experiments. Similar to the results of Table 1, PFML outperformed MAE and was comparable to data2vec when using the pre-trained models as feature extractors for linear classifiers. Again, both MAE and PFML outperformed data2vec by a large margin in sleep stage classification from EEG data. In valence classification from speech data and sleep stage classification from EEG data, the TS2Vec method excelled. However, for other classification tasks, the SSL features provided by TS2Vec were the weakest performance-wise. In the case of using a randomly initialized model as a feature extractor for linear classifiers, the classification accuracy was at chance-level in all cases except when classifying posture for IMU data.

The results of representation collapse experiments are shown in Table 3. As can be seen from the results, it is very common for representation collapse to occur with data2vec across all data modalities. On the contrary, the results indicate that PFML, MAE, and TS2Vec do not suffer from representation collapse: PFML and TS2Vec did not experience representation collapses at all, and MAE had a representation collapse only once. Furthermore, we attribute this single representation collapse of MAE to bad luck in model weight initialization, as in this particular case the model loss started diverging from the beginning of the pre-training process. The results showcase that methods like PFML, MAE, and TS2Vec, whose training targets inherently contain variance, are less prone to representation collapse compared to methods like data2vec that learn their own prediction targets.

<sup>a</sup> TS2Vec uses a different encoder compared to other experiments, since the TS2Vec algorithm is built into the pre-defined CNN encoder architecture.

<sup>b</sup> In TS2Vec, the linear layer is added directly after the pre-trained encoder.

**TABLE 4.** A listing of pros and cons for each of the SSL methods used in the present experiments. For further details on these results, refer to Appendix D in the supplementary material.

	End-use performance	Robust to representation collapse	Expected pre-training time	Flexibility of neural network architecture	Linear separability of SSL features	Number of hyperparameters
TS2Vec	–	+	±	–	–	±
MAE	±	+	+	+	±	+
data2vec	+	–	–	+	±	–
PFML (ours)	+	+	+	+	+	±

## VI. ADDITIONAL HYPERPARAMETER EXPERIMENTS

In order to demonstrate that it is more beneficial during pre-training to mask the latent features instead of masking the input directly, we ran PFML pre-training for all three datasets twice: either by masking the inputs or by masking the embeddings. Subsequently, we fine-tuned our models for all five classification tasks, and the results are shown in Table 7 of Appendix C in the supplementary material. As can be observed from the results, it is more beneficial for downstream tasks if we alleviate the complexity of the pre-training task for the encoder by masking the embeddings instead of masking the inputs. The only exception was with EEG data, where it did not make a difference whether inputs or embeddings were masked.

For each data modality, we also experimented with different configurations of masking probability  $p_m$  and the length of the masks  $l_m$ . We ran PFML pre-training using different configurations of  $p_m$  and  $l_m$ , and then we fine-tuned the pre-trained models. For IMU and speech data, we only experimented with one classification task each, namely classification of movement from IMU data and classification of valence from speech data. The results for different configurations of  $p_m$  and  $l_m$  for IMU, speech, and EEG data are shown in Appendix C in the supplementary material in Tables 8, 9, and 10, respectively. For IMU data, the differences between different masking strategies are rather small, whereas for speech and EEG data the selection of masking hyperparameters has a notable effect on fine-tuning performance.

We also experimented with the effect of discarding some of the functionals in PFML pre-training for IMU data. After pre-training, we fine-tuned our model for movement classification, and the results are presented in Table 11 of Appendix C in the supplementary material. The results indicate that using the full set of 11 functionals during PFML pre-training provides the best outcome. As the number of discarded functionals increases, the prediction task becomes simpler and the training targets are able to capture less information of the input signal frames, leading to worse fine-tuning performance.

Finally, we tested different mask types for PFML pre-training using IMU data. We either replaced the masked embeddings with a fixed vector of zeros, ones, random Gaussian noise (as in e.g. Baevski et al. [11]), or a learnable mask token (as in e.g. Baevski et al. [8]). After PFML

pre-training using the four different mask types, we fine-tuned the pre-trained models for movement classification. Table 12 of Appendix C in the supplementary material presents the comparison results for different mask types. As can be observed, the choice between a mask of ones or random Gaussian noise does not have a notable impact on the performance. However, using a learnable mask token yielded slightly worse results than a vector of ones or random Gaussian noise, and a vector of zeros yielded the worst results. We observed that either using a vector of ones, random Gaussian noise, or learnable mask tokens for masking the embeddings promoted embedding variance, whereas using a vector of zeros provided a smaller level of variance for the embedding representations during pre-training. This lower level of variance for embeddings might potentially hinder the fine-tuning process, resulting into a lower performance in downstream tasks.

## VII. CONCLUSION

In this paper, we presented PFML, a novel SSL algorithm for time-series data that avoids the common SSL issue of representation collapse. PFML operates by predicting statistical functionals of the input signal corresponding to masked embeddings, given a sequence of unmasked embeddings. We demonstrated the effectiveness of PFML using five different classification tasks across three different data modalities: infant posture and movement classification from multi-sensor IMU data, emotion recognition from speech data, and sleep stage classification from EEG data. Our results show that PFML is superior to both a conceptually similar SSL method, MAE, and a contrastive learning-based SSL method, TS2Vec. Our results also show that PFML is on par with the current state-of-the-art data modality agnostic SSL method, data2vec, while being conceptually simpler and without suffering from representation collapse. The pros and cons for each of the SSL methods used in the present study are shown in Table 4. The fact that PFML matches the performance of data2vec while also avoiding the issue of representation collapse renders PFML more straightforward to apply to new time-series data domains, such as in the case of clinical time-series data. The present work may also be extended to other domains than time-series data, such as images where functionals could be computed of, e.g., image patches.

## A. LIMITATIONS

We selected the present set of 11 functionals for their effectiveness across the three data modalities used in the present study, aiming for potential generalizability and a robust starting point to other data domains and downstream tasks. However, carefully selecting the number and type of functionals specifically for different modalities may lead to better results than presented here. Also, we did not include data augmentation in our pre-training processes to save computational time for PFML pre-training, as we wanted to pre-compute the functionals before the model training. As shown in e.g. [1], [6], [26], [28], data augmentation during pre-training may lead to improved performance on downstream tasks. Nonetheless, performing masking for randomly sampled frames is already a form of data augmentation in itself. Furthermore, other model architectures besides CNN-based encoders or Transformer encoder blocks could also be used, and this may improve PFML pre-training performance. Lastly, we acknowledge that typically SSL pre-training is run with very large minibatch sizes using multiple GPUs, and the results of the present experiments might improve with larger minibatch sizes. However, to promote reproducibility and encourage other researchers to try PFML, we deliberately pre-trained our models using relatively small minibatches so that the pre-training processes could be run on a single GPU with 16 GB of VRAM. As detailed in Appendix E in the supplementary material, our method used only a moderate amount of computational resources.

## B. BROADER IMPACTS

Since the main goal of PFML is to make the algorithm straightforwardly applicable to different time-series data domains, our method makes it easier to apply SSL pre-training for time-series data without complex tuning of hyperparameters or the need to profoundly understand the target data domain. As an example, properties of different medical time-series data, such as those obtained with EEG, ECG, or EMG, can be dependent on the clinical environment, the specific measurement equipment and setup, or clinical population being measured [60]. This limits the applicability of “universal” pre-trained models predominant in computer vision and speech technology. In a similar manner, various industrial sensor setups, such as those for system monitoring and predictive maintenance (accelerometers, magnetometers etc.), can result in data unique to a particular environment or machine type. In these cases, the use of PFML pre-training can be practical, since applying modality-specific SSL algorithms or fine-tuning pre-trained models from other data modalities might not generalize well to novel time-series data domains. Hence, PFML may promote the use of machine learning as an assisting tool in e.g. clinical healthcare or other limited-data domains. However, as with all classifiers, machine-learning models trained using PFML might make errors. Incorrect model-based decisions, such as incorrect diagnoses, may be detrimental in some cases. Lastly, any bias, private information, or harmful content in the pre-training data can,

in theory, be reflected to the feature representations that are learned by PFML.

## ACKNOWLEDGMENT

The authors would like to thank Tampere Center for Scientific Computing for the computational resources used in this study. The author Einari Vaaras would like to thank the Finnish Foundation for Technology Promotion for the encouragement grants and the Nokia Foundation for the Nokia Scholarship.

## REFERENCES

- [1] R. Balestriero, M. Ibrahim, V. Sobal, A. Morcos, S. Shekhar, T. Goldstein, F. Bordes, A. Bardes, G. Mialon, Y. Tian, A. Schwarzschild, A. Gordon Wilson, J. Geiping, Q. Garrido, P. Fernandez, A. Bar, H. Pirsaviash, Y. LeCun, and M. Goldblum, “A cookbook of self-supervised learning,” 2023, *arXiv:2304.12210*.
- [2] J. Gui, T. Chen, J. Zhang, Q. Cao, Z. Sun, H. Luo, and D. Tao, “A survey on self-supervised learning: Algorithms, applications, and future trends,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 46, no. 12, pp. 9052–9071, Dec. 2024.
- [3] D. Erhan, Y. Bengio, A. Courville, P.-A. Manzagol, P. Vincent, and S. Bengio, “Why does unsupervised pre-training help deep learning?” *J. Mach. Learn. Res.*, vol. 11, no. 19, pp. 625–660, Mar. 2010.
- [4] A. van den Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” 2018, *arXiv:1807.03748*.
- [5] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, “Wav2vec 2.0: A framework for self-supervised learning of speech representations,” in *Proc. NeurIPS*, Jan. 2020, pp. 12449–12460.
- [6] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, “A simple framework for contrastive learning of visual representations,” in *Proc. ICML*, Jan. 2020, pp. 1597–1607.
- [7] H. Akbari, L. Yuan, R. Qian, W.-H. Chuang, S. Chang, Y. Cui, and B. Gong, “VATT: Transformers for multimodal self-supervised learning from raw video, audio and text,” in *Proc. NeurIPS*, Jan. 2021, pp. 24206–24221.
- [8] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, and M. Auli, “Data2vec: A general framework for self-supervised learning in speech, vision and language,” in *Proc. ICML*, Jan. 2022, pp. 1298–1312.
- [9] W. Wang, H. Bao, L. Dong, J. Bjorck, Z. Peng, Q. Liu, K. Aggarwal, O. Khan Mohammed, S. Singhal, S. Som, and F. Wei, “Image as a foreign language: BEIT pretraining for vision and vision-language tasks,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2023, pp. 19175–19186.
- [10] O. J. Hénaff, A. Srinivas, J. D. Fauw, A. Razavi, C. Doersch, S. M. A. Eslami, and A. V. D. Oord, “Data-efficient image recognition with contrastive predictive coding,” in *Proc. ICML*, Jan. 2019, pp. 4182–4192.
- [11] A. Baevski, A. Babu, W.-N. Hsu, and M. Auli, “Efficient self-supervised learning with contextualized target representations for vision, speech and language,” in *Proc. ICML*, Jan. 2022, pp. 1416–1429.
- [12] J. W. Yoon, S. M. Kim, and N. S. Kim, “MCR-Data2vec 2.0: Improving self-supervised speech pre-training via model-level consistency regularization,” in *Proc. Interspeech*, Aug. 2023, pp. 2833–2837.
- [13] Q.-S. Zhu, L. Zhou, J. Zhang, S.-J. Liu, Y.-C. Hu, and L.-R. Dai, “Robust Data2 VEC: Noise-robust speech representation learning for ASR by combining regression and improved contrastive learning,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Jun. 2023, pp. 1–5.
- [14] J. Lian, A. Baevski, W.-N. Hsu, and M. Auli, “Av-Data2 Vec: Self-supervised learning of audio-visual speech representations with contextualized target representations,” in *Proc. IEEE Autom. Speech Recognit. Understand. Workshop (ASRU)*, Dec. 2023, pp. 1–8.
- [15] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus, “Hard negative mixing for contrastive learning,” in *Proc. NeurIPS*, Jan. 2020, pp. 21798–21809.
- [16] J. A. Robinson, C.-Y. Chuang, S. Sra, and S. Jegelka, “Contrastive learning with hard negative samples,” in *Proc. ICLR*, Jan. 2020, pp. 1–28.
- [17] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, “Unsupervised learning of visual features by contrasting cluster assignments,” in *Proc. NeurIPS*, Jan. 2020, p. 9912.

- [18] W.-N. Hsu, B. Bolte, Y. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 29, pp. 3451–3460, 2021.
- [19] T. Hua, W. Wang, Z. Xue, S. Ren, Y. Wang, and H. Zhao, "On feature decorrelation in self-supervised learning," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2021, pp. 9578–9588.
- [20] J. Li, P. Vincent, Y. LeCun, and Y. Tian, "Understanding dimensional collapse in contrastive self-supervised learning," in *Proc. ICLR*, Jan. 2021, pp. 1–17.
- [21] Q. Garrido, R. Balestrieri, L. Najman, and Y. LeCun, "RankMe: Assessing the downstream performance of pretrained self-supervised representations by their rank," in *Proc. ICML*, Jan. 2022, pp. 10929–10974.
- [22] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao, J. Wu, L. Zhou, S. Ren, Y. Qian, Y. Qian, J. Wu, M. Zeng, X. Yu, and F. Wei, "WavLM: Large-scale self-supervised pre-training for full stack speech processing," *IEEE J. Sel. Topics Signal Process.*, vol. 16, no. 6, pp. 1505–1518, Oct. 2022.
- [23] H.-Y. Lee, J.-B. Huang, M. Singh, and M.-H. Yang, "Unsupervised representation learning by sorting sequences," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 667–676.
- [24] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *Proc. ICLR*, Jan. 2018, pp. 1–16.
- [25] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proc. ECCV*, Jan. 2018, pp. 139–156.
- [26] J.-B. Grill, F. Strub, F. Althé, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. Á. Pires, Z. D. Guo, M. G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko, "Bootstrap your own latent: A new approach to self-supervised learning," in *Proc. NeurIPS*, Jan. 2020, pp. 21271–21284.
- [27] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *Proc. ICML*, Jan. 2021, pp. 8748–8763.
- [28] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2022, pp. 15979–15988.
- [29] H. Bao, L. Dong, S. Piao, and F. Wei, "BEiT: BERT pre-training of image transformers," in *Proc. ICLR*, Jan. 2021, pp. 1–18.
- [30] M. Oquab et al., "DINOv2: Learning robust visual features without supervision," in *Proc. Trans. Mach. Learn. Res.*, Jan. 2023, pp. 1–32.
- [31] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, Jan. 2018, pp. 4171–4186.
- [32] T. B. Brown et al., "Language models are few-shot learners," in *Proc. NeurIPS*, Jan. 2020, pp. 1877–1901.
- [33] Y. Tay, M. Dehghani, V. Q. Tran, X. García, J. Lee, X. Wang, H. W. Chung, D. Bahri, T. Schuster, H. Zheng, D. Zhou, N. Houlsby, and D. Metzler, "UL2: Unifying language learning paradigms," in *Proc. ICLR*, Jan. 2022, pp. 1–39.
- [34] OpenAI et al., "GPT-4 technical report," 2023, *arXiv:2303.08774*.
- [35] Z. Yue, Y. Wang, J. Duan, T. Yang, C. Huang, Y. Tong, and B. Xu, "TS2 Vec: Towards universal representation of time series," in *Proc. AAAI*, vol. 36, Jun. 2022, pp. 8980–8987.
- [36] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 9726–9735.
- [37] E. Pizzi, S. D. Roy, S. N. Ravindra, P. Goyal, and M. Douze, "A self-supervised descriptor for image copy detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2022, pp. 14512–14522.
- [38] Y. M. Asano, C. Rupprecht, and A. Vedaldi, "Self-labelling via simultaneous clustering and representation learning," in *Proc. ICLR*, Jan. 2019, pp. 1–22.
- [39] W. Wang, Q. Tang, and K. Livescu, "Unsupervised pre-training of bidirectional speech encoders via masked reconstruction," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2020, pp. 6889–6893.
- [40] Z. Xie, Z. Zhang, Y. Cao, Y. Lin, J. Bao, Z. Yao, Q. Dai, and H. Hu, "SimMIM: A simple framework for masked image modeling," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2022, pp. 9643–9653.
- [41] A. Bardes, J. Ponce, and Y. LeCun, "VICReg: Variance-invariance-covariance regularization for self-supervised learning," in *Proc. ICLR*, Jan. 2021, pp. 1–23.
- [42] J. H. McDermott and E. P. Simoncelli, "Sound texture perception via statistics of the auditory periphery: Evidence from sound synthesis," *Neuron*, vol. 71, no. 5, pp. 926–940, Sep. 2011.
- [43] L. R. Rabiner and R. W. Schafer, "Introduction to digital speech processing," *Found. Trends Signal Process.*, vol. 1, nos. 1–2, pp. 1–194, 2007.
- [44] A. Gulati, J. Qin, C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, "Conformer: Convolution-augmented transformer for speech recognition," in *Proc. Interspeech*, Oct. 2020, pp. 5036–5040.
- [45] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [46] K. Cho, B. V. Merriënboer, Ç. Gulçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder–decoder for statistical machine translation," in *Proc. EMNLP*, Jan. 2014, pp. 1724–1734.
- [47] D. Hendrycks and K. Gimpel, "Gaussian error linear units (GELUs)," 2016, *arXiv:1606.08415*.
- [48] J. Lei Ba, J. Ryan Kiros, and G. E. Hinton, "Layer normalization," 2016, *arXiv:1607.06450*.
- [49] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," 2016, *arXiv:1607.08022*.
- [50] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, "On the variance of the adaptive learning rate and beyond," in *Proc. ICLR*, Aug. 2019, pp. 11–14.
- [51] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *Proc. ICLR*, Jan. 2017, pp. 1–19.
- [52] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, Dec. 2014, pp. 1–15.
- [53] M. Airaksinen, A. Gallen, A. Kivi, P. Vijayakrishnan, T. Häyriinen, E. Ilén, O. Räsänen, L. M. Haataja, and S. Vanhatalo, "Intelligent wearable allows out-of-the-lab tracking of developing motor abilities in infants," *Commun. Med.*, vol. 2, no. 1, p. 69, Jun. 2022.
- [54] E. Vaaras, M. Airaksinen, S. Vanhatalo, and O. Räsänen, "Evaluation of self-supervised pre-training for automatic infant movement classification using wearable movement sensors," in *Proc. 45th Annu. Int. Conf. IEEE Eng. Med. Biol. Soc. (EMBC)*, Jul. 2023, pp. 1–6.
- [55] M. Airaksinen, O. Räsänen, E. Ilén, T. Häyriinen, A. Kivi, V. Marchi, A. Gallen, S. Blom, A. Varhe, N. Kaartinen, L. Haataja, and S. Vanhatalo, "Automatic posture and movement tracking of infants with wearable movement sensors," *Sci. Rep.*, vol. 10, no. 1, p. 169, Jan. 2020.
- [56] E. Vaaras, S. Ahlqvist-Björkroth, K. Drossos, L. Lehtonen, and O. Räsänen, "Development of a speech emotion recognizer for large-scale child-centered audio recordings from a hospital environment," *Speech Commun.*, vol. 148, pp. 9–22, Mar. 2023.
- [57] B. Kemp, A. H. Zwiderman, B. Tuk, H. A. C. Kamphuisen, and J. J. L. Obery, "Analysis of a sleep-dependent neuronal feedback loop: The slow-wave microcontinuity of the EEG," *IEEE Trans. Biomed. Eng.*, vol. 47, no. 9, pp. 1185–1194, Sep. 2000.
- [58] A. L. Goldberger, L. A. N. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley, "PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for complex physiologic signals," *Circulation*, vol. 101, no. 23, pp. e215–e220, Jun. 2000.
- [59] E. Eldele, Z. Chen, C. Liu, M. Wu, C.-K. Kwok, X. Li, and C. Guan, "An attention-based deep learning approach for sleep stage classification with single-channel EEG," *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 29, pp. 809–818, 2021.
- [60] D. S. Watson, J. Krutzinna, I. N. Bruce, C. E. Griffiths, I. B. McInnes, M. R. Barnes, and L. Floridi, "Clinical applications of machine learning algorithms: Beyond the black box," *BMJ*, vol. 364, p. 1886, Mar. 2019.



**EINARI VAARAS** was born in Finland, in 1996. He received the B.Sc.(Tech.) and M.Sc.(Tech.) degrees in electrical engineering from Tampere University, Tampere, Finland, in 2019 and 2021, respectively, where he is currently pursuing the D.Sc.(Tech.) degree. His research interests include representation learning, active learning, paralinguistic speech processing, and applying machine-learning models for clinical healthcare.



**MANU AIRAKSINEN** was born in Finland, in 1986. He received the M.Sc.(Tech.) and D.Sc.(Tech.) degrees in speech and language technology from Aalto University, Espoo, Finland, in 2012 and 2018, respectively. He is currently a Senior Research Engineer with the BABA Center, University of Helsinki, and Helsinki University Hospital, Helsinki, Finland. His current research interests include wearable technology and machine learning applied to clinical and developmental science contexts. He has published over 45 peer-reviewed journal articles and conference papers on the related topics.



**OKKO RÄSÄNEN** (Senior Member, IEEE) was born in Finland, in 1984. He received the M.Sc.(Tech.) degree in language technology from Helsinki University of Technology, Espoo, Finland, in 2007, and the D.Sc.(Tech.) degree in language technology from Aalto University, Espoo, in 2013. In 2015, he was a Visiting Researcher with the Language and Cognition Laboratory, Stanford University, Stanford, CA, USA. He is currently a Professor with the Signal Processing Research Centre, Tampere University, Finland. He also holds the Title of the Docent from the School of Electrical Engineering, Aalto University, in the area of spoken language processing. His research interests include computational modeling of language acquisition, cognitive aspects of language processing, and speech processing and machine learning in general. He has published more than 110 journal articles and peer-reviewed conference papers on related topics.

...