# Data-Efficient Sleep Staging with Synthetic Time Series Pretraining

Niklas Grieger[1,2,3,*]    Siamak Mehrkanoon[2]    Stephan Bialonski[1,3,*]

[1]Department of Medical Engineering and Technomathematics, FH Aachen University of Applied Sciences, 52428 Jülich, Germany
[2]Department of Information and Computing Sciences, Utrecht University, Utrecht, The Netherlands
[3]Institute for Data-Driven Technologies, FH Aachen University of Applied Sciences, 52428 Jülich, Germany
[*]grieger@fh-aachen.de, bialonski@fh-aachen.de

## Abstract

Analyzing electroencephalographic (EEG) time series can be challenging, especially with deep neural networks, due to the large variability among human subjects and often small datasets. To address these challenges, various strategies, such as self-supervised learning, have been suggested, but they typically rely on extensive empirical datasets. Inspired by recent advances in computer vision, we propose a pretraining task termed "frequency pretraining" to pretrain a neural network for sleep staging by predicting the frequency content of randomly generated synthetic time series. Our experiments demonstrate that our method surpasses fully supervised learning in scenarios with limited data and few subjects, and matches its performance in regimes with many subjects. Furthermore, our results underline the relevance of frequency information for sleep stage scoring, while also demonstrating that deep neural networks utilize information beyond frequencies to enhance sleep staging performance, which is consistent with previous research. We anticipate that our approach will be advantageous across a broad spectrum of applications where EEG data is limited or derived from a small number of subjects, including the domain of brain-computer interfaces.

## 1 Introduction

Deep neural networks have achieved significant advances in analyzing electroencephalographic (EEG) time series [39], ranging from brain-computer interfaces [30] to the intricacies of sleep stage scoring [36, 15]. Such successes are attributed to the ability of deep neural networks, as universal function approximators, to learn properties (features) from patient data that are difficult for humans to conceptualize and define. However, training neural networks requires large and diverse datasets that capture the considerable variety between individual subjects and their medical conditions (subject heterogeneity). Creating such datasets is challenging due to the typically limited amount of data per subject (data scarcity) and diverse measurement protocols used in different clinics, which can introduce additional variability in the data. Furthermore, acquiring large datasets is often expensive, complicated, or even intractable due to strict privacy policies and ethical guidelines. This hinders the advancement of deep neural networks for widespread application in real-world medical settings.

Efforts to mitigate the scarcity of large datasets have primarily followed two paths: (1) the development of network architectures that incorporate constraints mirroring the data's intrinsic characteristics, such as symmetries [7], and (2) enhancing model performance with additional or cross-domain data to learn effective priors. Pertaining to the first path, a common feature in time series processing networks is the use of convolutional layers. These layers are designed to be translation-equivariant [18], which ensures that a temporal shift in the input only affects the output by the same shift. This characteristic enables consistent network responses to temporal patterns, regardless of their temporal location, while reducing the number of model parameters compared to architectures lacking such constraints. For the second path, a variety of strategies have been proposed to learn useful priors from data. One approach is data augmentation, in which time series are transformed while preserving their annotations (labels) to artificially expand the dataset [31, 23]. Deep neural networks trained on such augmented datasets implicitly learn to become invariant under these transformations, which can lead to better out-of-sample prediction performance. Another strategy is transfer learning [12], a two-step process in which neural networks are trained on one task using a large dataset (pretraining step [37]) and then adapted to learn the actual task of interest using another (usually much smaller) dataset (fine-tuning step). A variant of this idea is self-supervised learning [32, 2], which allows neural networks to be pretrained on large and heterogeneous datasets without explicitly labeled examples. Finally, generative models such as VAEs, GANs, and diffusion models can be used to sample new time series to extend existing datasets [22, 8, 46]. Such generative models approximate a data distribution and require large heterogeneous datasets for training. While all of these approaches have been demonstrated to be able to improve the performance of neural networks, they still rely on large empirical datasets for training.

Recent advances in computer vision have demonstrated that it is possible to learn effective priors exclusively from synthetic images, which has the potential to significantly reduce the need for large empirical datasets [3, 26]. Synthetic images for image

classification tasks were generated by simple random processes, such as iterated function systems to produce fractals [26] or random placement of geometric objects to cover an image canvas [3]. Deep neural networks pretrained on such data were demonstrated to learn useful priors for image classification tasks, yielding competitive performance comparable to pretraining on natural images on various benchmarks [26]. This remarkable finding highlights the potential of synthetic datasets that can be generated without much computational resources and, theoretically, in unlimited amounts.

Inspired by these advances, we hypothesize that pretraining exclusively on synthetic time series data generated from simple random processes can also yield effective priors for sleep staging. Given the importance of frequencies for sleep stage scoring and other EEG-based applications [5, 33], we introduce a pretraining method called "frequency pretraining" (FPT) that centers on generating synthetic time series data with specific frequency content. During pretraining, deep neural networks learn to accurately predict the frequencies present in these synthetic time series. Despite the deliberate simplicity of our synthetic data generation process and the inherent domain shift between synthetic and EEG data, we observe that FPT allows a deep neural network to detect sleep stages more accurately than fully supervised training when few samples (few-sample regime) or data from few subjects (few-subject regime) are available for fine-tuning. The success of our method underscores the essential role of frequency content in enabling neural networks to accurately and reliably discern sleep stages. We consider pretraining techniques leveraging synthetic data, like the one we propose, as a promising area of research, offering the potential to develop models in sleep medicine and neuroscience that are particularly suited for scenarios involving small datasets. To facilitate testing and further advancements, we make the source code of our method publicly available [19].

## Contributions

- Novel synthetic pretraining approach: We introduce "frequency pretraining" (FPT), a pretraining approach using synthetic time series with random frequency content, eliminating the need for empirical EEG data during pretraining.

- Demonstrated data efficiency: We demonstrate superior sleep staging performance of FPT in few-sample and few-subject regimes across three datasets, with comparable results to fully supervised methods when data is abundant.

- Analysis of frequency-based priors: We evaluate the role of frequency information in our pretraining task and how synthetic sample diversity affects fine-tuning.

- Comparison with self-supervised methods: We benchmark our approach against established self-supervised learning methods, demonstrating that synthetic data pretraining can achieve comparable results without requiring EEG data for pretraining.

## 2  Results

Our approach, illustrated in Figure 1 and detailed in Section 4, is based on a two-phase training process that combines pretraining on synthetic time series with fine-tuning on clinical sleep data. During the pretraining phase, we generate synthetic time series composed of sine waves with random frequencies drawn from predefined frequency ranges (frequency bins). These synthetic signals are then used to train a deep neural network whose feature extractor, $f$, learns to extract useful features, while the classifier, $c_p$, learns to predict the frequency bins from which the frequencies were drawn to generate the synthetic time series. After pretraining, the feature extractor is transferred to the fine-tuning phase, where it is applied to real EEG and electrooculography (EOG) data to classify sleep stages. In this phase, the model processes sequences of eleven consecutive sleep epochs, indexed $i - 5$ to $i + 5$, with the feature extractor producing features $h_i$ from each epoch. Another classifier $c_f$ then aggregates these features to predict the sleep stage (Wake, N1, N2, N3, REM) for the central sleep epoch $i$.

We evaluated our approach on three publicly available datasets, *DODO/H*, *Sleep-EDFx*, and *ISRUC*, which provided data from 276 subjects, including both healthy individuals and those with various medical conditions (see Section 4.2.1). For sleep staging performance, we tracked the Macro-F1 score, with higher values indicating better classification accuracy across sleep stages. During the pretraining phase, we assessed the model's ability to predict frequency bins using the Hamming metric and the accuracy (see Section 4).

### 2.1  Training Configurations

We compared the performance of pretrained models against the performance of non-pretrained models in scenarios with varying amounts of training data. In particular, we studied the performance of our approach in few-sample and few-subject regimes, where the greatest benefit was expected. Furthermore, we analyzed the priors that the model learned during pretraining and the role of frequency information in the learned features. Finally, we investigated whether these features could be further improved by fine-tuning the feature extractor. To enable these investigations, we created four training configurations.

**Fully Supervised.** The Fully Supervised training configuration is similar to many existing deep learning approaches for sleep staging [36] and served as a baseline to compare our pretrained models against. In this configuration, we skipped the pretraining step and trained (fine-tuned) the feature extractor $f$ and classifier $c_f$ from scratch using sleep staging data.
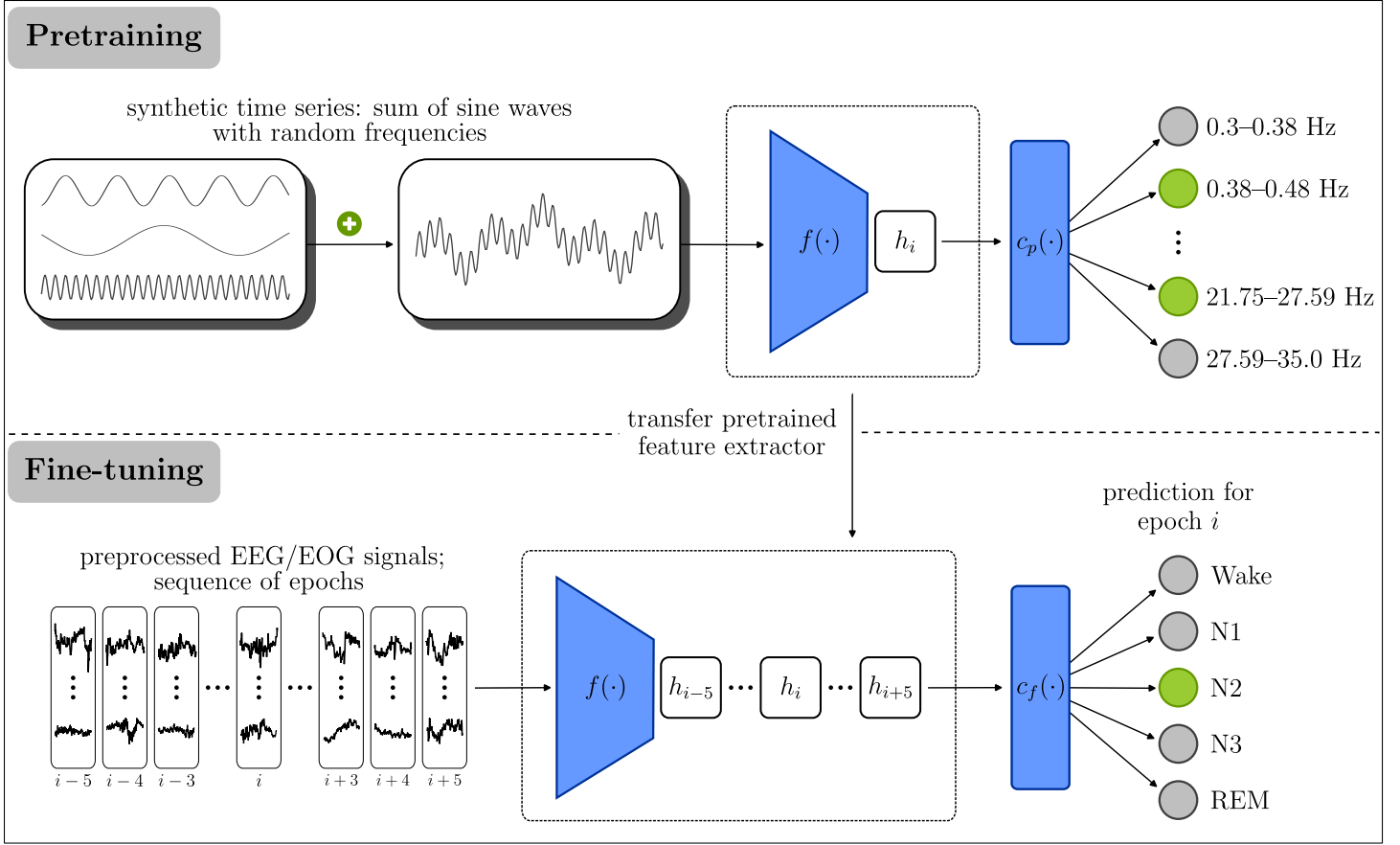
Figure 1: The training process consists of a *pretraining* and a *fine-tuning* phase. In the pretraining phase ("frequency pretraining"), we generated synthetic time series signals by summing sine waves with random frequencies. These synthetic signals were used to train a deep neural network consisting of a feature extractor $f$ and a classifier $c_p$ to predict the frequencies present in the signals (multi-label classification problem). In the fine-tuning phase, the pretrained feature extractor $f$ produced features $h_i$ from individual EEG and EOG epochs. The features of a sequence of these epochs were then used by a classifier $c_f$ that was trained to predict the sleep stage of the middle epoch in the sequence (multi-class classification problem).

**Fixed Feature Extractor.** We employed the Fixed Feature Extractor configuration to investigate the relevance of the features generated by the pretrained feature extractor for sleep staging. After pretraining the feature extractor $f$ on synthetic data, we kept its model weights and BatchNorm statistics fixed and only fine-tuned the sleep staging classifier $c_f$ on sleep data.

**Fine-Tuned Feature Extractor.** With this training configuration, we studied (i) how model performance changes when the pretrained feature extractor is allowed to change during fine-tuning and (ii) whether the priors learned during pretraining can prevent overfitting in few-sample or few-subject regimes. As in the previous configuration, we first pretrained the feature extractor $f$ on synthetic data, but then fine-tuned both the feature extractor and the classifier $c_f$ on sleep data without keeping any model weights fixed. Consequently, this configuration is similar to the Fully Supervised configuration with the key distinction that the feature extractor is initialized with pretrained weights.

**Untrained Feature Extractor.** The Untrained Feature Extractor configuration served as a baseline to study whether our pretraining scheme produces priors that are superior to random weights for sleep staging. We randomly initialized the feature extractor $f$ using Kaiming normal initialization [24] and then kept its weights fixed while fine-tuning the classifier $c_f$. This approach mirrors the Fixed Feature Extractor configuration, but with a random feature extractor instead of a pretrained one.

## 2.2 Data Efficiency

To assess the data efficiency of our pretraining method, we compare the performance of the different training configurations when fine-tuned with a reduced amount of training data (low-data regime) or the full training data (high-data regime). We investigated the low-data regime by first pretraining models with the full synthetic dataset. Depending on the training configuration, we then fine-tuned the pretrained or randomly initialized models with the data of 50 randomly sampled sleep staging samples from one subject. The sampling procedure selected sleep staging data without class stratification, and each sleep stage was represented at least once in the reduced datasets. In the high-data regime, we followed the same procedure
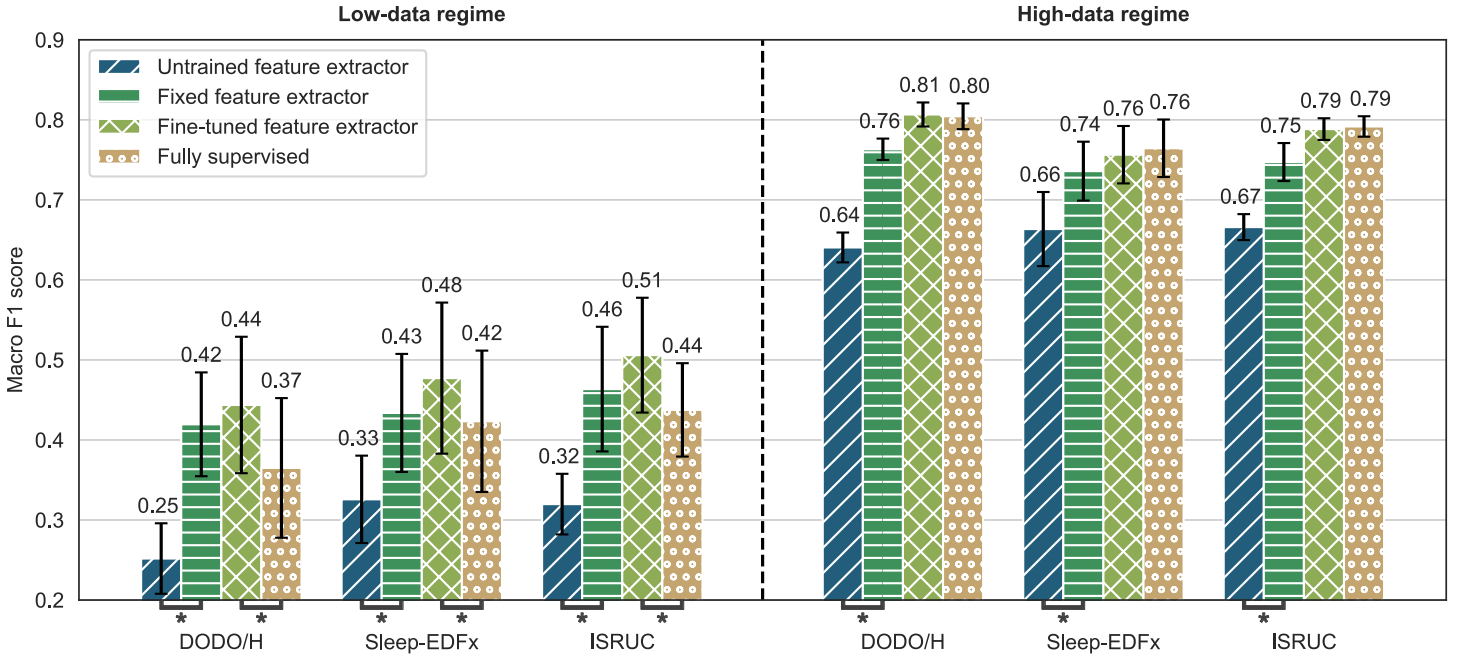
Figure 2: Average macro F1 scores achieved by the different training configurations when fine-tuned with a small set of training data (left side) or a large set (full set) of training data (right side) from one of the three datasets DODO/H, Sleep-EDFx, and ISRUC. The bars indicate the mean of the macro F1 scores measured on the test sets and averaged over 15 trainings (3 repetitions of a 5-fold cross-validation). Error bars show the standard deviation of the macro F1 scores, and * symbols indicate significant differences between Untrained and Fixed Feature Extractor or between Fine-Tuned Feature Extractor and Fully-Supervised training configuration ($p < 0.01$) according to a one-sided Wilcoxon signed-rank test.

but fine-tuned models with the full training data instead of a reduced amount of data. We repeated all experiments three times in a 5-fold cross-validation scheme, resulting in 15 training runs for each configuration and dataset. This approach allowed us to estimate the spread of the macro F1 scores when using different model initializations, sleep staging samples, and data folds for training and testing (see Section 4).

In the low-data regime, we observed that models pretrained with synthetic data outperformed models trained from scratch in sleep stage classification (see Figure 2). The performance gap between pretrained and non-pretrained models was most pronounced when comparing the fine-tuned feature extractor to the fully supervised configuration, with the former achieving average macro F1 scores that were 0.06–0.07 higher than those of the latter across datasets. When removing the fine-tuning of the pretrained features (fixed feature extractor configuration), our method still yielded macro F1 scores that were 0.01–0.05 higher than those of the fully supervised configuration. Comparing the fixed feature extractor to the untrained feature extractor configuration further highlights the importance of the learned features. Pretraining the feature extractor with synthetic data improved the macro F1 scores by 0.10–0.17 compared to a random initialization. While the macro F1 scores of the training configurations varied between datasets, the general trends observed in the low-data regime were consistent across all three datasets.

In the high-data regime, our pretrained models were on par with fully supervised trained models ($p < 0.05$ for a paired TOST test with a margin of $\pm 0.01$ on the DODO/H and ISRUC datasets) and achieved competitive performance in sleep stage classification (see Figure 2). The fine-tuned feature extractor configuration achieved average macro F1 scores of 0.76–0.81 across datasets, comparable to the macro F1 scores of 0.76–0.80 achieved by the fully supervised configuration. As in the low-data regime, we observed that fine-tuning the pretrained features was beneficial, as the macro F1 scores achieved by the fine-tuned feature extractor configuration were 0.02–0.05 higher than those of the fixed feature extractor configuration. The performance gap between the untrained feature extractor and the fixed feature extractor remained substantial even in the high-data regime, with the former achieving average macro F1 scores that were 0.08–0.12 lower than those of the latter.

## 2.3 Impact of Subject Diversity and Number of Training Samples

We further explored the data efficiency of our pretraining method by investigating how model performance is affected by subject diversity and sample volume (i.e., number of samples) in the training data used for fine-tuning. To study the effect of subject diversity independent of sample volume, we separately varied the number of subjects and the number of samples in the training data. The number of subjects was randomly sampled as $n_{subj} \in \{1, 2, 3, 4, 5\}$, and the number of samples was randomly sampled from those subjects as $n_{samples} \in \{50, 130, 340, 900, \text{all}\}$ ("all" indicates that all available training samples were used). Our sampling strategy selected sleep staging data without class stratification, and each sleep stage was

represented at least once in the reduced datasets. For each parameter combination, we trained the fully supervised and fine-tuned feature extractor configurations on all three datasets with three repetitions of 5-fold cross-validation.

In our results, we observed that both the fully supervised and the fine-tuned feature extractor configurations benefited from increased subject diversity, even when the total number of training samples was held constant (see rows in Figure 3a–f). Similarly, both configurations benefited from an increased number of training samples when the number of subjects was held constant (see columns in Figure 3a–f). For the considered parameter ranges, the impact of reduced subject diversity and reduced sample volume on model performance appeared to be comparable.

Similar to the observations made in Section 2.2, the fine-tuned feature extractor configuration achieved better macro F1 scores than the fully supervised configuration in the few-sample regime. The performance gap between the two configurations was most evident when the number of samples was limited to 50, with bootstrap estimates of the mean differences in macro F1 scores ranging from 0.05 to 0.09 across datasets and subject numbers (see Figure 3g–i). These differences in macro F1 scores decreased as the number of samples increased. When training on all available training samples, the fine-tuned feature extractor configuration achieving comparable or slightly better performance than the fully supervised configuration. Interestingly, the performance gap between the two configurations with all available training samples was most pronounced for the DODO/H dataset. The fine-tuned feature extractor configuration achieved average macro F1 scores that were 0.02–0.06 higher than those of the fully supervised configuration. For the Sleep-EDFx and ISRUC datasets, the performance differences between the two configurations with all available training samples were minimal at 0.00–0.01.

Depending on the dataset, the fine-tuned feature extractor configuration showed varying degrees of improvement over the fully supervised configuration in the few-subject regime. When training with data of a single subject, the fine-tuned feature extractor configuration achieved average macro F1 scores that were 0.06–0.11 higher across sample numbers than those of the fully supervised configuration for the DODO/H dataset (see rows in Figure 3g). This performance gap decreased as the number of subjects increased, with mean bootstrapped differences in macro F1 scores between 0.02 and 0.09 for five subjects. In contrast, varying the number of subjects had less impact on the performance gap between the fine-tuned feature extractor and the fully supervised configurations for the Sleep-EDFx and ISRUC datasets (see rows in Figure 3h–i). The improvements achieved by the fine-tuned feature extractor configuration when trained with one subject (0.01–0.05 for the Sleep-EDFx dataset and 0.00–0.07 for the ISRUC dataset) were comparable to those achieved when trained with five subjects (0.00–0.06 for the Sleep-EDFx dataset and 0.01–0.06 for the ISRUC dataset).

## 2.4   Priors Towards Frequency Information

To get a better understanding of the pretraining process and the priors learned by the model, we recorded several metrics during pretraining. The recorded loss function converged to a low value, indicating that the model has learned effectively (see Figure 4a). At the same time, the Hamming metric reached a high value of 0.9 on the validation data (see Figure 4b). This value can be interpreted as the model predicting the frequency bins that were used to create the synthetic signals with an accuracy of 90%. The model was particularly proficient in predicting higher frequencies starting from 2.5 Hz (accuracies $> 90\%$; see Figure 4c). Lower frequencies, especially those below 1 Hz, were predicted with lower accuracy (75–85%).

We hypothesize that the differences in prediction accuracy across frequency bins are not due to the pretrained model being unable to predict lower frequencies. Instead, we believe this discrepancy arises from the varying width of the frequency bins (see x-axis of Figure 4c) as the bins for lower frequencies were narrower than those for higher frequencies due to their logarithmic scaling. While we applied a logarithmic binning scheme to increase our model's focus on the lower frequencies important for identifying slow wave sleep (N2 and N3 sleep), the narrow low frequency bands increase the difficulty of the pretraining task leading to decreased accuracies. Preliminary exploration of this trade-off between model focus and task difficulty showed that the current binning scheme slightly outperforms a linear binning scheme in the fixed feature extractor setup.

Interestingly, the ability of the pretrained feature extractor to extract useful features for sleep staging was strongly influenced by the diversity of the synthetic pretraining data (see Figure 4d). We investigated this influence of sample diversity by pretraining models with varying amounts of different synthetic samples $n_{synthetic} \in \{1, 10, 100, 10^3, 10^4, 10^5, 10^6\}$. To isolate the effect of sample diversity from the effect that the number of training steps has on model performance, we kept the number of gradient updates per training epoch constant by under- or oversampling the synthetic data to 100,000 samples. The pretrained models were then trained for sleep staging using data from one random subject from the DODO/H dataset and evaluated on the validation data of each cross-validation fold. As in the previous experiments, we performed three repetitions of a 5-fold cross-validation scheme for each number of synthetic samples.

Both the fixed feature extractor and the fine-tuned feature extractor configurations required a high level of diversity in the synthetic pretraining data to achieve good sleep staging performance (see Figure 4d). When pretrained with only one synthetic sample, the performance of the two training configurations with pretraining differed only slightly from the performance of the configurations without pretraining (untrained feature extractor and fully supervised configurations). As the number of synthetic samples was increased to more than 100, the performance of the pretrained models improved substantially until reaching a plateau at around 10,000 samples. We hypothesize that this plateau was reached because the model had learned all the relevant features from the pretraining task, and additional samples did not provide any further benefit. The simplicity
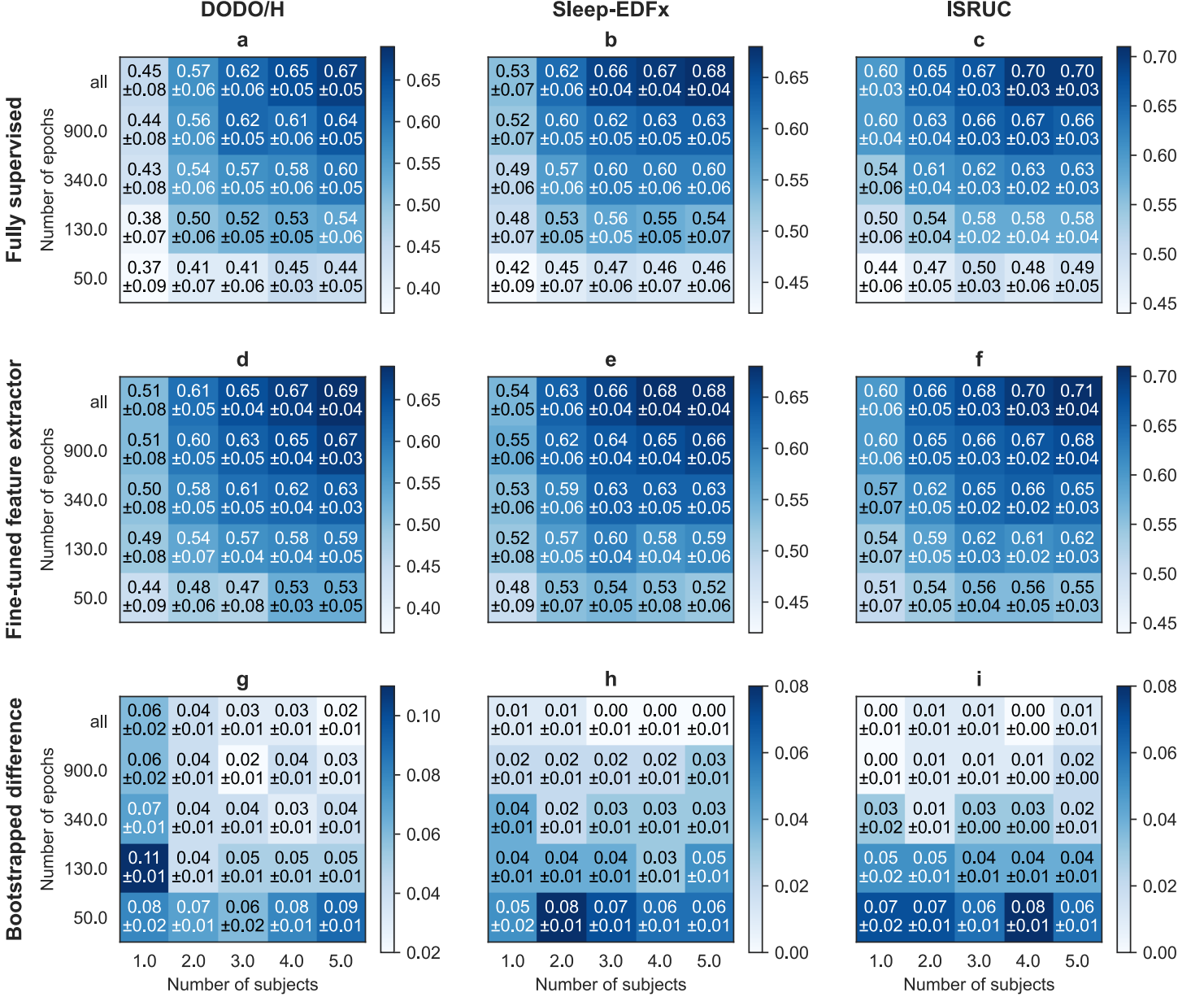
Figure 3: Detailed comparison of the fully supervised and fine-tuned feature extractor configurations in few-sample and few-subject regimes for each of the three datasets DODO/H (panels **a**,**d**,**g**), Sleep-EDFx (matrices **b**,**e**,**h**), and IS-RUC (panels **c**,**f**,**i**). Panels (**a**–**c**) and Panels (**d**–**f**) show the macro F1 scores achieved by the fully supervised and the fine-tuned feature extractor configurations, respectively. The displayed scores are the mean and the standard deviation of the macro F1 scores measured on the test set over 15 trainings (3 repetitions of a 5-fold cross-validation). Panels (**g**–**i**) show the average macro F1 scores achieved by the fine-tuned feature extractor configuration minus the average macro F1 scores achieved by the fully supervised configuration. These differences were calculated using a bootstrapping approach with 10,000 bootstrap samples. In this bootstrapping approach, we first paired the 15 macro F1 scores available for both configurations for each matrix entry based on the used seed and fold. We then calculated the difference between each pair of macro F1 scores. For each bootstrap sample, we sampled 15 of these differences with replacement and calculated their average. Finally, we display the mean and the standard deviation of these bootstrap samples.
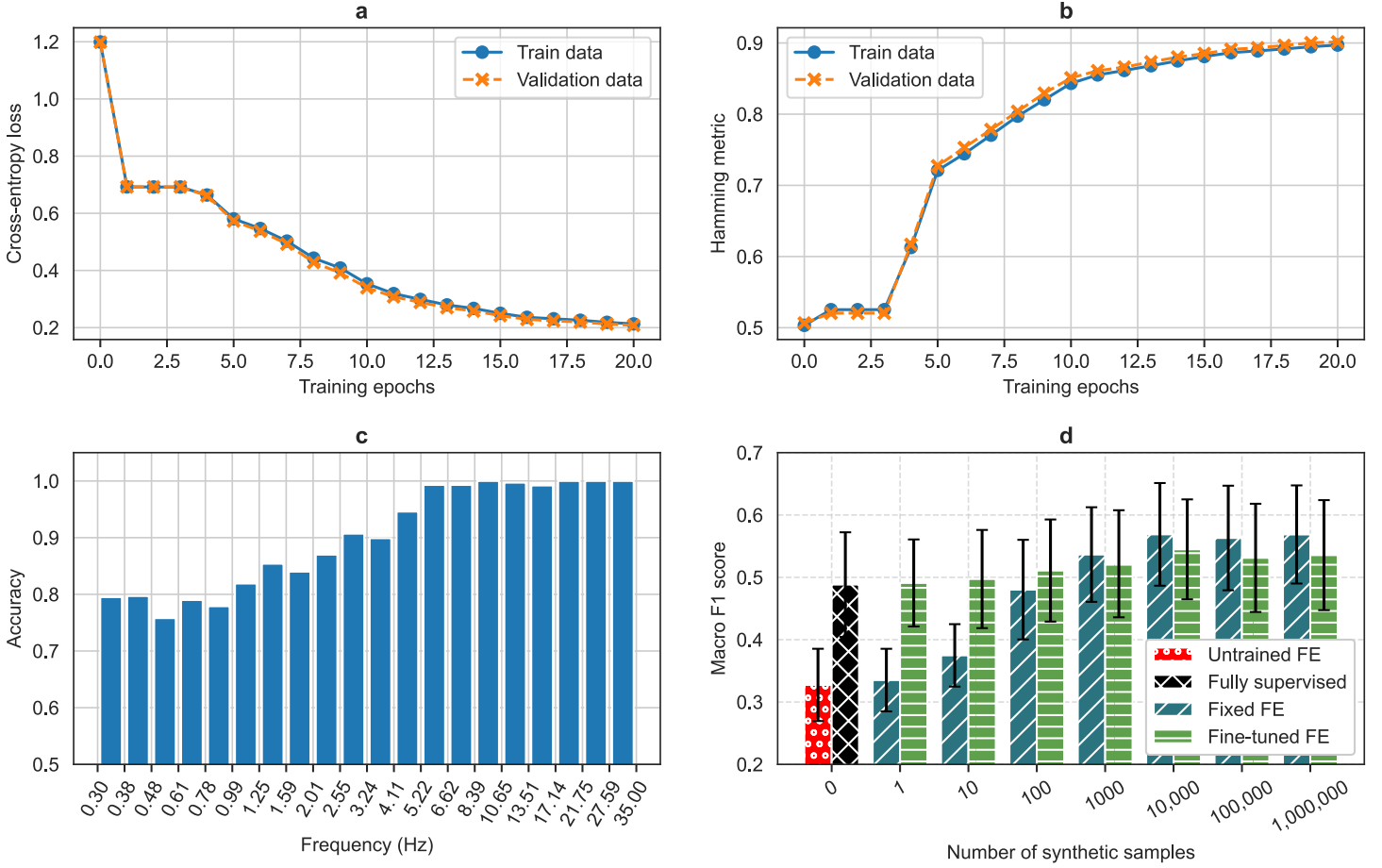
Figure 4: Analysis of the "frequency pretraining" task. Panel (**a**) shows the development of the loss function during a single pretraining run for both the training and validation data. The following two panels quantify the accuracy of a pretrained model to predict the frequencies of the synthetic signals. The hamming metric in panel (**b**) measures the overall accuracy summarized across all frequency bins, while panel (**c**) shows the accuracy for each frequency bin separately. Panel (**d**) illustrates how the diversity of the synthetic pretraining data affects the performance of the pretrained feature extractor (FE) after fine-tuning for sleep staging on a single subject from the DODO/H dataset. Each bar indicates the mean of the macro F1 scores on the validation data averaged over 3 repetitions of a 5-fold cross-validation, while the error bars show the standard deviation of the macro F1 scores.

of the pretraining task could also explain the negligible performance differences between the pretrained configurations with less than 100 synthetic samples and the training configurations without pretraining. For such few synthetic samples, the model may have memorized the synthetic data rather than learn general features useful for sleep staging.

## 2.5 Comparison to Self-Supervised Methods

We compared our frequency pretraining approach (FPT) with two popular self-supervised learning (SSL) methods, Sim-CLR [9] and VICReg [4], using the same model architecture as in FPT. The projection head (SimCLR) and expander (VICReg) followed their original designs [9, 4]. Different from the conceptually simple multi-label classification task in FPT, both SSL methods follow the objective of contracting representations of data-augmented "positive" samples while keeping representations of dissimilar "negative" samples far apart (see original works for details [9, 4]). Positive samples were generated using standard data augmentations: amplitude scaling (random factor 0.5–2), Gaussian noise injection ($\sigma = 0.05$), random temporal masking (10 segments of 1.5–3 s), time shifting (up to $\pm1.5$ s), and time warping (random factor 0.67–1.5). We evaluated two pretraining configurations: pretraining on (i) synthetic data (as in FPT) and (ii) EEG recordings from the ISRUC and Sleep-EDFx datasets. Fine-tuning was performed on DODO/H datasets under low-data (50 random epochs from one random subject) and high-data (all available training data) regimes, respectively.

When pretrained on synthetic data, both SimCLR and VICReg achieved performance comparable to the FPT method, with FPT slightly outperforming the SSL methods in the high-data regime when the feature extractor was fine-tuned (average MF1 scores of 0.80 versus 0.81, see Table 1). Pretraining on EEG data yielded only modest improvements, with both SSL methods showing slightly better performance in the low-data regime when the feature extractor was fixed (difference in average MF1 scores of 0.02) and in the high-data regime when fine-tuning the feature extractor (difference in average MF1

scores of 0.01). Compared to the EEG-pretrained SimCLR and VICReg models in the high-data regime, FPT showed similar performance (average MF1 scores of 0.76 and 0.81) but, unlike the two SSL methods, did not require EEG data for pretraining. In the low-data regime, the SSL methods performed comparable to the FPT method when fine-tuning the feature extractor and achieved slightly higher scores than FPT when keeping the feature extractor fixed (average MF1 scores of 0.44 and 0.45 versus 0.42). None of the differences in MF1 scores between the SSL methods and FPT were significant according to a two-sided Wilcoxon signed-rank test ($p$-values $> 0.05$).

Table 1: Comparison of different pretraining methods when fine-tuned with 50 random sleep epochs of one subject (low-data regime) or the full training set (high-data regime) of the DODO/H dataset. During fine-tuning, the feature extractor was either kept fixed (FE fixed) or updated together with the classifier (FE fine-tuned). We pretrained SimCLR and VICReg on the same synthetic data used for the frequency pretraining task, in addition to the traditional approach of pretraining on EEG data (i.e., ISRUC and Sleep-EDFx datasets). MF1 scores were calculated as the average of three repetitions of a five-fold cross validation (standard deviation in parentheses).

| Pretraining Method | Avg. MF1 Low-Data Regime | | Avg. MF1 High-Data Regime | |
| --- | --- | --- | --- | --- |
| | FE Fixed | FE Fine-Tuned | FE Fixed | FE Fine-Tuned |
| SimCLR, pretrained with synth. data | 0.42 (0.08) | 0.44 (0.09) | 0.76 (0.02) | 0.80 (0.02) |
| SimCLR, pretrained with EEG data | 0.44 (0.08) | 0.43 (0.09) | 0.76 (0.02) | 0.81 (0.02) |
| VICReg, pretrained with synth. data | 0.43 (0.07) | 0.44 (0.08) | 0.76 (0.02) | 0.80 (0.02) |
| VICReg, pretrained with EEG data | 0.45 (0.08) | 0.45 (0.09) | 0.76 (0.02) | 0.81 (0.01) |
| Frequency pretraining | 0.42 (0.06) | 0.44 (0.09) | 0.76 (0.01) | 0.81 (0.02) |

# 3  Discussion

In this work, we propose a novel pretraining scheme for EEG time series data that leverages synthetic data generated by a simple random process. We specialized the hyperparameters of the pretraining task to the typical frequency range and distribution of sleep EEG signals and demonstrated the effectiveness of this task for sleep stage classification. Due to the availability of several open sleep staging datasets [35, 47], we were able to fully control the amount and diversity of the training data, which allowed us to study the impact of our method in different data regimes. We hypothesize that our pretraining scheme could be particularly beneficial in few-sample and few-subject regimes, which we argue could benefit greatly from the priors towards frequency information that a model learns during pretraining.

Our results confirm the effectiveness of our pretraining scheme, particularly in few-sample and few-subject regimes. Pretrained models outperformed non-pretrained models when fine-tuned with a reduced number of subjects or training samples (see rows and columns in Figure 3g–i, respectively). The performance gap between pretrained and non-pretrained models was most pronounced in the few-sample regime, where pretrained models consistently achieved improvements over non-pretrained models across multiple datasets. In the few-subject regime, this performance gap was not as consistent across datasets, with our pretraining method showing the most substantial improvements in the DODO/H dataset. Our findings support observations made in the field of Self-Supervised Learning (SSL) that pretrained models generally have better data efficiency than fully supervised ones [2, 13]. In contrast to SSL methods, however, our pretraining scheme improves data efficiency without requiring empirical data, while achieving comparable or only slightly reduced performance (see Table 1). Interestingly, we observed that two data-augmentation-based SSL methods performed well even when pretrained with synthetic data instead of EEG data (see Table 1), suggesting that SSL approaches are promising, yet computationally intensive, alternatives to the FPT method. Further exploration of the use of SSL methods for pretraining on synthetic data is warranted and seems promising. Generating synthetic data can be much cheaper and more cost-effective than collecting EEG data, as generating 10,000 samples of the synthetic data used in this study required only $8.17 \pm 0.37$ s of a single consumer-grade CPU core (Lenovo Legion S7 16IAH7 Laptop with an Intel® Core™ i7-12700H CPU).

We hypothesize that the potential of synthetic data stems from the priors that the model learns during pretraining. These priors could prevent overfitting to a small number of training samples, particularly those from minority classes (e.g., N1 sleep), or subject-specific features, which is especially problematic in situations with very little training data. As expected, we observed that all of our training configurations improved with a larger training dataset (see Figure 2). This aligns with the prevalent view in the literature that deep learning models for sleep staging need substantial amounts of diverse data to perform well [36, 1, 15, 14]. When trained with the full training data, pretrained models performed comparably to fully supervised models (see Figure 2), achieving macro F1 scores similar to those of other deep learning approaches for sleep staging [36, 16]. In conclusion, our pretraining method was most beneficial in situations with limited training data, where it outperformed models trained from scratch, but had less impact in situations with large amounts of training data.

We further observed that, while the frequency content of a signal is crucial for sleep staging, deep neural networks extract additional information from the data that exceeds the frequency domain. The importance of the frequency content of a signal for sleep staging is demonstrated by the high macro F1 scores achieved by the fixed feature extractor configuration

and the substantial performance improvements it achieved over untrained feature extractors (see Figure 2). We attribute the performance gap between the two training configurations to the priors learned during pretraining the feature extractor. These priors biased the model to extract frequency information from the data, which it achieved with high accuracy after pretraining (see Figure 4b,c). Our finding is consistent with previous studies that reported frequency-based features to be important for sleep staging [33]. When the feature extractor of our model was allowed to be fine-tuned after pretraining, model performance increased (see Figure 2). We hypothesize that this increase in the macro F1 score is due to the feature extractor learning to extract information beyond the frequency content of the signal during fine-tuning. This hypothesis is in line with the AASM annotation guidelines [5], which consider several frequency-unrelated features essential for sleep staging. These features include time-domain information, which is important for spindles as well as amplitudes and specific patterns, such as k-Complexes [15]. Recent studies that applied feature engineering approaches to sleep staging further support our hypothesis by including additional features from the time-domain in the models used [43, 10]. In their work, Vallat and Walker analyzed the most important features for their model and found that time-related features, such as the time elapsed from the beginning of a recording, were among the top 20 most important features [43]. Although it remains unclear what additional information our pretrained models learn during fine-tuning, our method offers a promising avenue for future research into the interpretability of deep neural networks for sleep staging.

There are several opportunities for future work that could build upon our findings. One promising direction is to explore the pretraining task in more detail, for example, by investigating the synthetic data generation process and the impact of changing the used frequency range. Similar to previous work in the vision domain [3], it could also be promising to investigate which structural properties of synthetic time series are important for sleep staging. This could be achieved by defining new pretraining tasks that are based on different data generation processes that incorporate more complex structures like desynchronized phases across channels, noise, or polymorphic amplitude variations. Exploring such data generation processes may lead to a better understanding of what constitutes "natural" EEG time series and what information, besides the frequency content, is essential for sleep staging. In addition, we suggest exploring models with greater capacity and less inductive bias than the CNN-based architecture used in this work, such as transformer models [45, 6], which we expect to benefit even more from our pretraining method. Pretraining such models with synthetic data may alleviate their need for large amounts of training data [11]. Another avenue for future research is to investigate whether our pretraining method is beneficial for specific cohorts of subjects, such as patients with a specific disorder or specific age groups. Although we did investigate datasets with different demographics in this work, we did not perform detailed analyses of the impact of these demographics on model performance. Finally, it could be insightful to compare our approach with a broader range of SSL methods [32] and data augmentation strategies that employ synthetic EEG generators [31, 22]. To enable such comparisons and to facilitate future research in this direction, we make our code available online [19].

Our method presents a novel solution to address important issues that affect current deep learning models in the EEG time series domain, without requiring large amounts of patient data. We expect our approach to be advantageous in various applications where EEG data is scarce or derived from a limited number of subjects, such as brain–computer interfaces [30] or neurological disorder detection [44].

# 4    Materials and Methods

We trained deep neural networks in two phases: a pretraining phase based on synthetic data and a fine-tuning phase based on clinical sleep staging data. The training process is illustrated in Figure 1.

## 4.1    Phase 1: Pretraining with Synthetic Data

In the pretraining phase, we trained a deep neural network to predict the frequencies of randomly-generated synthetic time series ("frequency pretraining") [20]. The deep neural network consisted of a feature extractor to produce features and a classifier to predict which frequencies are present in the synthetic time series. Instead of training the model to predict the exact frequencies, we defined frequency ranges (frequency bins) and let the model identify the bins containing frequencies present in the signal (therefore turning the problem into a multi-label classification task).

### 4.1.1    Synthetic Data Generation

The synthetic samples used for pretraining consisted of a time series signal $s(t) \in \mathbb{R}^{3 \times 3000}$, featuring 3 channels of 30 s sampled at 100 Hz, and an associated label vector $l$. The synthetic time series were created by summing sine waves with randomly sampled frequencies, while the label vector encoded from which frequency ranges (frequency bins) the frequencies were sampled. Since we wanted to test our method in the context of sleep staging, we followed the American Academy of Sleep Medicine (AASM) guidelines and focused on the spectral band between 0.3–35 Hz. This range was divided into 20 frequency bins on a base-2 logarithmic scale, resulting in a label vector $l \in \{0, 1\}^{20}$ with binary entries for each bin that indicated whether frequencies in a given bin were present (1) or absent (0) in the time series. To generate a synthetic time series, we first randomly selected a subset of the 20 frequency bins, setting the corresponding entries in the label vector $l$ to 1 (each bin was selected with a probability of 50%). Next, we independently

sampled frequencies from the selected bins for each channel and generated sine waves with these frequencies. The sine waves were then shifted by a random phase sampled from a uniform distribution between 0 and $2\pi$, summed, and normalized to create the synthetic signal. Therefore, one channel $s_c(t)$ of a synthetic signal $s(t) = (s_1(t), s_2(t), s_3(t))$ is defined as

$$s_c(t) = \frac{\tilde{s}_c(t) - \mu_c}{\sigma_c}, \quad \tilde{s}_c(t) = \sum_i^{n_f} I_i \, \sin(2\pi \, t f_{c,i} + \phi_i), \quad I_i \in \mathcal{U}\{0,1\}, \tag{1}$$

where $t$ denotes time (in seconds), $n_f$ denotes the number of frequency bins, $I_i$ denotes the indicator function encoding whether a frequency within a bin $i$ is present (1) or not (0), $f_{c,i}$ denotes the frequencies of the sine waves (in Hertz), $\phi_i$ denotes their phase (in radians), and $\mu_c$ and $\sigma_c$ denote mean and standard deviation of $\tilde{s}_c$. Note that the phases $\phi_i$ were sampled for each frequency bin independently but remained the same across channels to keep the data generation process simple, while the frequencies $f_{c,i}$ were sampled independently for each bin and channel.

For each training run, we generated 101,000 synthetic samples, using 100,000 samples for training and the remaining 1000 samples as a validation set to track various metrics. The number of synthetic samples, frequency bins, and the logarithmic scale that defines the boundaries between the bins were determined through preliminary experiments where we explored a wide range of hyperparameters.

### 4.1.2 Model

We based our model on the TinySleepNet architecture, a conceptually simple deep neural network for sleep staging that has previously demonstrated competitive results [41]. This architecture consists of a convolutional feature extractor that extracts features from individual sleep epochs and a classifier that aggregates these feature across multiple epochs to perform sleep staging. The feature extractor consisted of four convolutional layers with 128 filters each, a kernel size of 50, 8, 8, and 8, respectively, and a stride of 25, 1, 1, and 1, respectively. Following each convolutional layer was a batch normalization layer and a ReLU activation function. The outputs of the first and last convolutional layer were passed through max pooling layers that reduced the temporal dimension of the feature maps by a factor of 8 and 4, respectively. After both max pooling layers, we applied dropout with a dropout rate of 0.5.

The feature extractor was followed by a classifier, which consisted of a dense layer with 80 neurons and ReLU activation function as well as a dense layer with 20 neurons and sigmoid activation function (i.e., one output for each frequency bin). We defined a threshold of 0.5 to determine whether a frequency from a frequency bin was present in the signal (output greater than 0.5) or not (output less than or equal to 0.5).

### 4.1.3 Pretraining for Frequency Prediction

We pretrained models for 20 training epochs using the Adam optimizer [29], a fixed learning rate of $10^{-4}$, a batch size of 64, and a binary cross-entropy loss. This loss function is commonly employed for multi-label classification problems with one-hot encoded labels. It is defined as

$$L_{bce} = -\frac{1}{N}\frac{1}{n_f}\sum_{i=1}^{N}\sum_{j=1}^{n_f} y_{i,j} \cdot \log(\hat{y}_{i,j}) + (1 - y_{i,j}) \cdot \log(1 - \hat{y}_{i,j}), \tag{2}$$

where $N$ is the number of samples, $n_f$ is the number of frequency bins, $y_{i,j}$ is the true one-hot encoded label of sample $i$ and frequency bin $j$, and $\hat{y}_{i,j}$ is the predicted one-hot encoded label of sample $i$ and frequency bin $j$.

In addition to the loss function, we recorded the hamming metric on the training and validation data after each training epoch. This metric is derived from the hamming loss, a common metric for multi-label classification problems [40]. Specifically, the hamming metric tracks the fraction of correctly predicted frequency bins and is defined as

$$H = \frac{1}{N}\frac{1}{n_f}\sum_{i=1}^{N}\sum_{j=1}^{n_f} I_{\{y_{i,j}=\hat{y}_{i,j}\}}, \tag{3}$$

where $I$ is the indicator function, which assumes the value 1 if the true one-hot encoded label $y_{i,j}$ of sample $i$ and frequency bin $j$ is equal to the predicted one-hot encoded label $\hat{y}_{i,j}$ of sample $i$ and frequency bin $j$.

## 4.2 Phase 2: Fine-Tuning with Sleep Staging Data

In the fine-tuning phase, we used the pretrained feature extractor of the pretraining phase for sleep stage classification on EEG and EOG data. Sleep stage classification is a multi-class classification problem where a short segment of EEG and EOG data (sleep epoch) is assigned to one of five sleep stages (Wake, N1, N2, N3, REM) based on the underlying brain activity. In our approach, the feature extractor of our model produced features from individual sleep epochs of EEG and EOG data. The extracted features were then concatenated over a sequence of multiple epochs and aggregated by a classifier to predict the sleep stage of the middle epoch in the input sequence (see Figure 1). By incorporating the surrounding sleep epochs into the classifier's input, we provided additional temporal context, potentially enhancing the accuracy of sleep stage prediction.

#### 4.2.1 Sleep Staging Data

To test our approach, we used four publicly available datasets encompassing a diverse range of subjects, including both healthy individuals and those with various medical conditions such as sleep apnea, REM sleep behavior disorder, and affective disorders: DODO [21], DODH [21], Sleep-EDFx [27, 34], and ISRUC [28]. The DODO and DODH datasets contain 55 recordings from subjects diagnosed with obstructive sleep apnea (OSA) and 25 recordings from healthy subjects, respectively. We combined the two datasets into a single dataset with 80 recordings, which we refer to as DODO/H. Each recording stems from a different subject and was split into 30-s non-overlapping windows (sleep epochs), which were annotated with sleep stages (Wake, N1, N2, N3, REM) by five sleep experts following the AASM guidelines [5]. We aggregated these five expert annotations into a consensus label for each sleep epoch using majority voting, with ties resolved by selecting the sleep stage determined by the most reliable scorer (defined as the scorer with the highest average agreement with all other scorers) [21]. Epochs that the consensus of the sleep experts annotated as artifacts were removed from the datasets. In our experiments, we focused on three EEG and EOG derivations available in both datasets and recommended by the AASM guidelines [5]: C3-M2, F3-M2, EOG1. All channels were sampled at 250 Hz.

The Sleep-EDFx dataset is provided on the PhysioNet platform [17] and contains 197 recordings from two studies [27, 34], of which we used the 153 recordings of the "Sleep Cassette" study [34]. The recordings stem from 78 healthy subjects with two recordings per subject, except for subjects 13, 36, and 52, who have only one recording each due to technical issues. The recordings were split into 30-s sleep epochs and annotated with sleep stages (Wake, N1, N2, N3, N4, REM) according to the 1968 Rechtschaffen and Kales manual [38] by a single expert per recording. To be consistent with the annotations of the other datasets, we combined the N3 and N4 stages into a single N3 stage and removed data annotated as artifacts or movement. Furthermore, we removed excess daytime data from the recordings by cropping them to begin 30 min before the first annotated non-Wake stage and end 30 min after the last annotated non-Wake stage. All available EEG and EOG channels (Fpz-Cz, Pz-Oz, EOG horizontal) were sampled at 100 Hz.

The ISRUC dataset [28] contains 126 recordings and is divided into three subgroups. Subgroup one contains 100 recordings from 100 subjects suffering from various disorders affecting sleep (including OSA, REM disorder, affective disorder, snoring), subgroup two contains 16 recordings from 8 subjects diagnosed with OSA and snoring (2 recordings per subject), and subgroup three contains 10 recordings from 10 healthy subjects. Due to technical issues, we excluded two recordings (recordings 8 and 100) from the first subgroup. All recordings were split into 30-s sleep epochs and annotated with sleep stages by two experts following the AASM guidelines [5]. We only used the annotations of the first expert and focused on the EEG and EOG channels C3-A2, F3-A2, and LOC-A2 sampled at 200 Hz.

To prepare the data for training, we filtered the signals between 0.3 and 35 Hz with an 8th-order zero-phase Butterworth filter and downsampled the filtered signals to 100 Hz using polyphase filtering if necessary. We then normalized the amplitudes of each individual sleep epoch by subtracting the median and dividing the result by the interquartile range of its amplitude distribution. Finally, we clipped all signal amplitudes above 20 or below $-20$ to minimize outliers.

After preprocessing, the three datasets (DODO/H, Sleep-EDFx, ISRUC) were split subject-wise into five folds for cross-validation. Note that we created a separate cross-validation split for each dataset, resulting in three different 5-fold cross-validation configurations. When splitting the DODO/H and the ISRUC datasets, we ensured that the recordings within each split were balanced in regard to the sub-datasets or subgroups. In addition to splitting the data into training and test folds, we reserved approximately 10% of the data from each fold as validation data for hyperparameter tuning and early stopping (i.e., one recording of DODO and DODH each, two subjects of Sleep-EDFx, and two recordings from the first subgroup of ISRUC for each fold). We validated our models on the validation data of the four training folds. The validation data of the test fold was kept separate and not involved in validation or testing to prevent data leakage.

#### 4.2.2 Model

Since the sleep epochs used for fine-tuning had the same format as the synthetic data generated for pretraining (i.e., 3 channels with 30 s of data sampled at 100 Hz), we kept the architecture of the feature extractor unchanged. However, we replaced the classifier with a different architecture to predict sleep stages instead of frequency bins. The classifier was inspired by the TinySleepNet architecture [41] and consisted of a bidirectional Long-Short Term Memory (LSTM) [25] layer with a hidden size of 128, a dropout layer with a dropout rate of 0.5, and a dense layer with 5 neurons and softmax activation function.

#### 4.2.3 Fine-Tuning for Sleep Stage Classification

We fine-tuned the models with training samples consisting of 11 sleep epochs that provided the classifier with context information from 5 epochs before and after the epoch to be classified. To ensure that the model could use the same aggregation process for all sleep epochs of a recording, we padded the first and last five sequences of each recording with zeros to the full sequence length. Fine-tuning was performed using the Adam optimizer [29], a fixed learning rate of $10^{-4}$, weight decay of $10^{-3}$, a batch size of 32, and the categorical cross-entropy loss. To prevent overfitting, we limited each training run to a maximum of 50 training epochs and stopped training early if the macro F1 score on the validation data did not improve for 10 training epochs (early stopping). Additionally, we clipped all gradients with a maximum norm greater than 5.0 to prevent exploding gradients.

We tracked model performance by recording the macro F1 score on the training and validation data after each training epoch. The macro F1 score is calculated as the average F1 score [42] across all sleep stages and is defined as

$$\overline{F_1} = \frac{1}{k} \sum_{i=1}^{k} 2 \cdot \frac{p_i \cdot r_i}{p_i + r_i}, \tag{4}$$

where $k$ is the number of sleep stages, $p_i$ is the precision of sleep stage $i$, and $r_i$ is the recall of sleep stage $i$.

In all four training configurations (Fully Supervised, Fixed Feature Extractor, Fine-Tuned Feature Extractor, and Untrained Feature Extractor), we performed training using the cross-validation scheme described in Section 4.2.1. A single training run in this cross-validation scheme included both pretraining and fine-tuning where applicable. If not specified otherwise, all reported results were obtained on the test folds of the cross-validation splits, while the small part of the training data reserved for validation was used as early stopping set during fine-tuning.

To ensure comparability between different training and data configurations, we seeded the model initialization, the synthetic data generation during pretraining, and the data subsampling when fine-tuning with a limited amount of training data. This strategy reduced the influence of random initialization and data sampling on comparisons between training configurations. In particular, seeding ensured that models were fine-tuned with the same data across all training configurations in the few-sample and few-subject regimes. Furthermore, when training with a reduced amount of data, we maintained a constant number of gradient updates per training epoch by duplicating each training sample $\lfloor \frac{N}{N_{red}} \rfloor$ times, where $N$ is the number of samples in the full training data and $N_{red}$ is the number of samples after reduction.

Experiments were performed on an NVIDIA DGX A100 workstation with eight NVIDIA A100 GPUs and a Dell workstation with a single NVIDIA RTX A6000 GPU.

## Author Contributions

## Funding

## Data Availability

The DODO and DODH datasets analyzed in the current study are described in Guillot et al. [21] and are publicly available (see `https://github.com/Dreem-Organization/dreem-learning-open` (accessed on 30 March 2023) for how to access the datasets). The Sleep-EDFx dataset is described in Kemp et al. [27] and Mourtazaev et al. [34] and is publicly available at `https://physionet.org/content/sleep-edfx/1.0.0/` (accessed on 30 March 2023). The ISRUC dataset is described in Khalighi et al. [28] and is publicly available at `https://sleeptight.isr.uc.pt/` (accessed on 13 April 2023). The source code necessary to reproduce our results is available on GitHub at `https://github.com/dslaborg/frequency-pretraining/tree/full-paper-version` (accessed on 29 August 2025).

## Acknowledgments

## Conflicts of Interest

The authors declare no conflicts of interest.

## References

[1] D. Alvarez-Estevez. Challenges of applying automated polysomnography scoring at scale. *Sleep Med. Clin.*, 18(3):277–292, Sept. 2023.

[2] H. Banville, O. Chehab, A. Hyvärinen, D.-A. Engemann, and A. Gramfort. Uncovering the structure of clinical EEG signals with self-supervised learning. *J. Neural Eng.*, 18:046020, Mar. 2021.

[3] M. Baradad, J. Wulff, T. Wang, P. Isola, and A. Torralba. Learning to see by looking at noise. In M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, editors, *Annu. Conf. on Neural Information Processing Systems, NeurIPS, Virtual, 6–14 December*, pages 2556–2569, Red Hook, NY, USA, Dec. 2021. Curran Associates Inc.

[4] A. Bardes, J. Ponce, and Y. LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *The Tenth Int. Conf. on Learning Representations, ICLR, Virtual, 25–29 April*. OpenReview.net, 2022.

[5] R. B. Berry, R. Brooks, C. E. Gamaldo, S. M. Harding, R. M. Lloyd, C. L. Marcus, and B. V. Vaughn. *The AASM manual for the scoring of sleep and associated events: Rules, terminology and technical specifications, Version 2.6*. American Academy of Sleep Medicine, Darien, Illinois, 2020.

[6] G. Brandmayr, M. M. Hartmann, F. Fürbass, G. Matz, M. Samwald, T. Kluge, and G. Dorffner. Relational local electroencephalography representations for sleep scoring. *Neural Networks*, 154:310–322, 2022.

[7] M. M. Bronstein, J. Bruna, T. Cohen, and P. Velickovic. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *CoRR*, abs/2104.13478, 2021.

[8] F. P. Carrle, Y. Hollenbenders, and A. Reichenbach. Generation of synthetic EEG data for training algorithms supporting the diagnosis of major depressive disorder. *Front. Neurosci.*, 17:1219133, 2023.

[9] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton. A simple framework for contrastive learning of visual representations. In *Proc. 37th Int. Conf. on Machine Learning, ICML 2020, Virtual, 13–18 July*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 2020.

[10] M. J. V. D. Donckt, J. V. D. Donckt, E. Deprost, N. Vandenbussche, M. Rademaker, G. Vandewiele, and S. V. Hoecke. Do not sleep on traditional machine learning: Simple and interpretable techniques are competitive to deep learning for sleep scoring. *Biomed. Signal Process. Control.*, 81:104429, 2023.

[11] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Int. Conf. on Learning Representations, ICLR, Vienna, Austria, 3–7 May*, 2021.

[12] A. Ebbehoj, M. Ø. Thunbo, O. E. Andersen, M. V. Glindtvad, and A. Hulman. Transfer learning for non-image data in clinical research: A scoping review. *PLOS Digital Health*, 1(2):e0000014, Feb. 2022.

[13] E. Eldele, M. Ragab, Z. Chen, M. Wu, C.-K. Kwoh, and X. Li. Self-supervised learning for label-efficient sleep stage classification: A comprehensive evaluation. *IEEE T. Neur. Sys. Reh.*, 31:1333–1342, 2023.

[14] L. Fiorillo, G. Monachino, J. van der Meer, M. Pesce, J. D. Warncke, M. H. Schmidt, C. L. A. Bassetti, A. Tzovara, P. Favaro, and F. D. Faraci. U-Sleep's resilience to AASM guidelines. *npj Digit. Medicine*, 6:33, 2023.

[15] L. Fiorillo, A. Puiatti, M. Papandrea, P.-L. Ratti, P. Favaro, C. Roth, P. Bargiotas, C. L. Bassetti, and F. D. Faraci. Automated sleep scoring: A review of the latest approaches. *Sleep Med. Rev.*, 48:101204, 2019.

[16] M. Gaiduk, A. Serrano Alarcón, R. Seepold, and N. Martínez Madrid. Current status and prospects of automatic sleep stages scoring: Review. *Biomed. Eng. Lett.*, 13(3):247–272, July 2023.

[17] 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*, 101(23), June 2000.

[18] I. J. Goodfellow, Y. Bengio, and A. C. Courville. *Deep Learning*. Adaptive computation and machine learning. MIT Press, Cambridge, Massachusetts, 2016.

[19] N. Grieger. Source code of the model presented in Grieger et al., "Data-efficient sleep staging with synthetic time series pretraining". https://github.com/dslaborg/frequency-pretraining, 2024.

[20] N. Grieger, S. Mehrkanoon, and S. Bialonski. Pretraining sleep staging models without patient data. In *Int. Conf. on Learning Representations (ICLR), Workshop on Learning from Time Series For Health, Vienna, Austria, 7–11 May*, Vienna, Austria, 2024.

[21] A. Guillot, F. Sauvet, E. H. During, and V. Thorey. Dreem open datasets: Multi-scored sleep datasets to compare human and automated sleep staging. *IEEE T. Neur. Sys. Reh.*, 28(9):1955–1965, Sept. 2020.

[22] A. G. Habashi, A. M. Azab, S. Eldawlatly, and G. M. Aly. Generative adversarial networks in EEG analysis: An overview. *J. NeuroEng. Rehabil.*, 20(1):40, Apr. 2023.

[23] C. He, J. Liu, Y. Zhu, and W. Du. Data augmentation for deep neural networks model in EEG classification task: A review. *Front. Hum. Neurosci.*, 15:765525, Dec. 2021.

[24] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In *2015 IEEE Int. Conf. on Computer Vision, ICCV 2015, 7–13 December*, pages 1026–1034, Santiago, Chile, Dec. 2015. IEEE Computer Society.

[25] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Comput.*, 9:1735–1780, Nov. 1997.

[26] H. Kataoka, K. Okayasu, A. Matsumoto, E. Yamagata, R. Yamada, N. Inoue, A. Nakamura, and Y. Satoh. Pre-training without natural images. *Int. J. Comput. Vis.*, 130(4):990–1007, 2022.

[27] B. Kemp, A. H. Zwinderman, B. Tuk, H. A. C. Kamphuisen, and J. J. L. Oberye. Analysis of a sleep-dependent neuronal feedback loop: The slow-wave microcontinuity of the EEG. *IEEE Trans. Biomed. Eng.*, 47(9):1185–1194, 2000.

[28] S. Khalighi, T. Sousa, J. M. dos Santos, and U. Nunes. ISRUC-Sleep: A comprehensive public dataset for sleep researchers. *Comput. Methods Programs Biomed.*, 124:180–192, 2016.

[29] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and Y. LeCun, editors, *3rd Int. Conf. Learning Representations, ICLR*, San Diego, CA, USA, May 2015.

[30] W. Ko, E. Jeon, S. Jeong, J. Phyo, and H.-I. Suk. A survey on deep learning-based short/zero-calibration approaches for EEG-based brain-computer interfaces. *Front. Hum. Neurosci.*, 15:643386, May 2021.

[31] E. Lashgari, D. Liang, and U. Maoz. Data augmentation for deep-learning-based electroencephalography. *J. Neurosci. Meth.*, 346:108885, Dec. 2020.

[32] Z. Liu, A. Alavi, M. Li, and X. Zhang. Self-supervised contrastive learning for medical time series: A systematic review. *Sensors*, 23(9):4221, 2023.

[33] S. Motamedi-Fakhr, M. Moshrefi-Torbati, M. Hill, C. M. Hill, and P. R. White. Signal processing techniques applied to human sleep EEG signals — A review. *Biomed. Signal Process. Control.*, 10:21–33, 2014.

[34] M. S. Mourtazaev, B. Kemp, A. H. Zwinderman, and H. A. C. Kamphuisen. Age and gender affect different characteristics of slow waves in the sleep EEG. *Sleep*, 18(7):557–564, Sept. 1995.

[35] I. Obeid and J. Picone. The Temple University hospital EEG data corpus. *Front. Neurosci.*, 10:196, May 2016.

[36] H. Phan and K. Mikkelsen. Automatic sleep staging of EEG signals: Recent development, challenges, and future directions. *Physiol. Meas.*, 43(4):04TR01, Apr. 2022.

[37] Y. Qiu, F. Lin, W. Chen, and M. Xu. Pre-training in medical data: A survey. *Mach. Intell. Res.*, 20(2):147–179, Feb. 2023.

[38] A. Rechtschaffen, A. Kales, R. Berger, W. Dement, A. Jacobson, L. Johnson, M. Jouvet, L. Monroe, I. Oswald, H. Roffwarg, B. Roth, and R. Walter. *A Manual of Standardized Terminology, Techniques and Scoring System for Sleep Stages of Human Subjects.* Public Health Service, U.S. Government Printing Office, Washington, D.C., 1968.

[39] Y. Roy, H. Banville, I. Albuquerque, A. Gramfort, T. H. Falk, and J. Faubert. Deep learning-based electroencephalography analysis: A systematic review. *J. Neural Eng.*, 16:051001, 2019.

[40] M. Sokolova and G. Lapalme. A systematic analysis of performance measures for classification tasks. *Inf. Process. Manag.*, 45(4):427–437, 2009.

[41] A. Supratak and Y. Guo. TinySleepNet: An efficient deep learning model for sleep stage scoring based on raw single-channel EEG. In *42nd Annual Int. Conf. of the IEEE Engineering in Medicine & Biology Society, EMBC 2020*, pages 641–644, Montreal, QC, Canada, July 2020. IEEE.

[42] A. Tharwat. Classification assessment methods. *Appl. Comput. Inform.*, 17(1):168–192, 2021.

[43] R. Vallat and M. P. Walker. An open-source, high-performance tool for automated sleep staging. *eLife*, 10:e70092, Oct. 2021.

[44] H. van Dijk, G. van Wingen, D. Denys, S. Olbrich, R. van Ruth, and M. Arns. The two decades brainclinics research archive for insights in neurophysiology (TDBRAIN) database. *Sci. Data*, 9(1):333, June 2022.

[45] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Annu. Conf. Neural Information Processing Systems, NeurIPS, 4–9 December*, pages 5998–6008, Red Hook, NY, USA, 2017. Curran Associates Inc.

[46] Z. You, Y. Guo, X. Zhang, and Y. Zhao. Virtual electroencephalogram acquisition: A review on electroencephalogram generative methods. *Sensors*, 25(10):3178, May 2025.

[47] G.-Q. Zhang, L. Cui, R. Mueller, S. Tao, M. Kim, M. Rueschman, S. Mariani, D. Mobley, and S. Redline. The national sleep research resource: Towards a sleep data commons. *J. Am. Med. Inform. Assn.*, 25(10):1351–1358, May 2018.