

Investigation of Attention as a Predictor for Stream Data Anomalies

Authors

Affiliations

Abstract—This paper presents a novel approach to anomaly detection in stream data. Modern machine learning methods are often used for this task, by training Autoencoders and Transformers on reconstruction and prediction, and then assuming the reconstruction/prediction error as an indication of anomalies. These approaches only take into account the morphology of the stream, and are not aware of the different, often use case-specific and highly contextualized, aspects of what is considered as an anomaly. In this article we explore a novel idea that aims to involve the use case into the training of the anomaly detection model: Instead of aiming at a general-purpose anomaly detector, we aim at a general-purpose way to train a case-specific anomaly detector without direct supervision. The core idea is to extract anomaly indicators from the attention layers of a deep neural network trained on a use case-specific task. By removing the need to train on reconstruction or prediction, we are free to train on any supervision that is available and is characteristic of the use case. This alleviates the most pressing challenge of anomaly detection, namely that anomalies are vaguely defined and thus difficult to directly supervise, while simultaneously allowing use case-specific knowledge to be extracted from whatever other (non anomaly-related) supervision might be available. We evaluated our method using electroencephalogram (EEG) signals on a sleep-stage labeling task. The performance of the proposed approach was compared against reconstruction/prediction error and against EEG-specific noise detection methods. The results indicate that the proposed method is a promising task-agnostic tool for anomaly detection in streaming data.

Index Terms—AI, Machine learning, Transformers, Robustness, Time-Series, anomaly detection, noise estimation, outlier detection.

1. Introduction

The term *anomaly* is used in time series processing to indicate a wide range of deviations and encompasses not only unforeseen events in the observed system but also transient problems in the observation mechanism. Indicatively, across various domains, an anomaly in the observation data can be attributed to a cyberattack, a disease, a faulty subsystem or, on the other hand, to a sensor momentarily losing contact

or being heavily affected by noise in the environment [1]–[4].

This vague definition of what constitutes an anomaly brings with it the range of technical challenges that distinctively characterize the field: the long and complex temporal dependencies that need to be considered due to the priorly unknown duration of what constitutes an event, non-stationarity when aiming at correctly separating anomalies from data shifts due to long-term trends, high-dimensional data as there is no clear prior knowledge of what variables are needed. As conventional approaches often struggle to provide satisfactory results, or are limited by the need to tune their parameters for each specific application on a case-by-case basis, recent advancements in *deep learning* have gained significant attention in the research community [5]. And, given the vagueness of what constitutes ground truth, the most successful approaches are unsupervised or self-supervised: *forecasting-based* approaches train a neural network to predict the next datapoint in the sequence and *reconstruction-based* approaches train a neural network to reconstruct the sequence from a learned, compressed representation. In either case, the assumption is that a high error between the predicted/reconstructed sequence and the actual data is indicative of an abnormality.

However, as Mejri et al. [5] note, there is no general-purpose approach (or even paradigm) that works best across domains and applications and performance depends on the use case and the nature of the encountered anomalies; Deep learning methods appear in general superior, but conventional approaches (mostly PCA) can also in some cases perform. To address this, in this article we explore a novel idea that aims to involve the use case into the training of the anomaly detection model in a way that the system *adapts itself* to different use cases. In other words, instead of aiming at a general-purpose anomaly detector, we aim at a general-purpose way to train a case-specific anomaly detector without direct supervision.

The core of the idea is that anomaly detection is a pre-processing step for some sequence processing task. We operate under the assumption that there is supervision for this downstream task, although there is no supervision for what constitutes an anomaly. We then make the following research hypothesis: *Anomalies are the parts of the sequence that have the property that ignoring them gives superior performance despite the fact that decisions are made from fewer datapoints.*

One can easily see how this maps directly to specific instances of anomalies. For example, assuming a classification task on sensor data, a burst of environment noise looks the same regardless of the correct label during that period; a successful classifier would learn to ignore such bursts as the same pattern maps to any label. As another example, a successful model for predicting upcoming network traffic would learn to ignore some kinds of bursts (although they might indicate a cyberattack) and learn to emphasize ‘normal’ bursts (that indicate, for instance, the start of working hours).

Our contribution is the formulation of a methodology for leveraging the outputs of intermediate layers of deep neural networks in order to extract the level of significance the network places on the different parts of the sequence being learned. In the remainder of this article, we first provide the necessary background (Section 2) and then proceed to describe our methodology (Section 3), which we evaluate on an EEG dataset where we have supervision for both a downstream task (sleep stage prediction) *and* carefully curated anomaly annotations (Section 4). We then present and discuss the experimental results (Section 5) and conclude (Section 6).

2. Background

One of the main frameworks in sequence processing are *recurrent neural networks* where the sequence is presented to the network one token at a time and the network maintains a *hidden state* which distills the information needed from past tokens to provide a context for the processing of the current token. One of the most successful recurrent architectures is the *Long Short-Term Memory (LSTM)* where trainable *gates* control the flow of information to and from the hidden state [6]. This allows LSTM to capture long temporal dependencies in a low-dimensional state representation [7], [8].

LSTMs are often combined with other techniques, such as convolution and encoder-decoder architectures. Encoder-decoder LSTMs, in particular, learn a compressed representation of the data. When training the decoder to reconstruct the original sequence back (which is known as *autoencoding*), the compressive encoding is trained to drop information that is not detrimental for the loss estimation. Anomalies (in the sense of patterns not encountered during training) can then be identified by higher reconstruction error [9], [10]. Reconstruction error over autoencoding has been used extensively for anomaly detection in various architecture besides LSTM [11], and its limitations are well-understood: Autoencoders struggle with high-dimensional data, as it is up to the system designer to find the layer widths (effectively, the level of compression) that drop the correct amount of information. Further, Autoencoders are limited with the respect to the kind of information they drop: Since this decision is driven by a loss that compares the reconstructed sequences against the input sequence, a commonly occurring pattern will be retained even when it is an anomaly in a given context.

Besides recurrency, the other major approach to sequence-processing is the *Transformer* architecture. Transformers are Autoencoders based on the *self-attention* mechanism [12]. Unlike LSTMs and other recurrent architectures, Transformers receive the complete sequence as input and model relationships across the entire sequence. This allows gradients to flow directly across the entire sequence, rather than being propagated step by step, making it easier to discover long-distance dependencies. At the core of the Transformer architecture lies the self-attention mechanism, which enables the model to assign different levels of importance to elements in the input sequence based on the value of *any other element in the sequence*. Specifically, Transformers learn three sets of weights which are applied to input of dimensionality d to get the *query* (Q) vector, the *key* (K) vector, and the *Value* (V) vector. The attention mechanism computes the *similarity score* QK^T which is a $d \times d$ matrix that determines the contribution of each element in the final representation of each other element. The similarity score is scaled and softmax’ed into a matrix of weights, which are applied to the value vector. This yields the representation $\text{softmax}(QK^T/\text{sqrt}(d)) \cdot V$ where each element contains information aggregated from the entire sequence, improving the model’s ability to detect long-range dependencies.

3. Research Methodology

3.1. Research Hypothesis

As discussed in the previous section, reconstruction and prediction errors are widely used as key indicators for anomaly detection in time-series data due to their intuitive appeal and straightforward implementation. However, these approaches do not take into account differences in the nature of what is considered an anomaly for each use case.

As we framed our work in a context where there is no anomaly supervision, it follows that it is also not possible to select training data that is not contaminated with anomalies. This means that the model may inadvertently learn to reconstruct anomalies. Additionally, reconstruction-based methods often struggle to detect contextual anomalies, where an observation may be anomalous only in specific temporal or multivariate contexts. For example, a high-temperature reading might be expected in the summer but anomalous in the winter, and reconstruction models may overlook such contextual nuances.

Similarly, prediction error, which measures deviations between predicted and actual values, can be susceptible to noise and non-stationarity in time-series data. In highly dynamic systems, normal variations may result in significant prediction errors, leading to false positives. Moreover, models relying on prediction error often assume that future patterns can be reliably forecasted based on past observations, an assumption that may not hold in volatile or chaotic systems.

On the other hand, the attention mechanism offers an alternative, currently unexplored, way to extract indications

TABLE 1: Methods under comparison.

Acronym	Architecture	Detection
<i>LSTM</i>	LSTM Autoencoder	Reconstruction error
<i>C-LSTM</i>	Convolutional LSTM Autoencoder	Reconstruction error
<i>AE_err</i>	Attention-based Autoencoder	Reconstruction error
<i>TP_err</i>	Transformer Predictor	Prediction error
<i>AE_att</i>	Attention-based Autoencoder	Attention
<i>TP_att</i>	Transformer Predictor	Attention
<i>PCA</i>	Principle component analysis	
<i>MNE</i>	IIR filter	

about what parts of the sequence are anomalies. The hypothesis is that the part of the sequence that receives the least attention while performing a task relevant to the use case, then is this part is an anomaly. The intuition is that neural networks will happily overfit the data when they are given enough parameters to do so. Such a network might not be good to actually perform the task, but is good at recognizing two ways in which a sub-sequence is anomalous: (a) it does not follow any pattern that was boosted (gradient-wise) by the training; (b) it follows a pattern that is inconsistent, it is sometimes associated with one class on the task and sometimes with another, so it alternates between being boosted and penalized by the loss function, again resulting in low attention.

To investigate this hypothesis, we experimented with several methods for anomaly and noise detection ranging from attention-based methods to state-of-the-art machine learning approaches, as well as conventional anomaly detection methods. These are listed in Table 1 and described in more detail in the remainder of this section.

3.2. Reconstruction and Prediction Error

The *LSTM Autoencoder (LSTM)* uses a sequence-to-sequence architecture with LSTM layers for both encoding and decoding. The encoder compresses the input into a fixed-size latent representation by processing the input sequence and retaining the final hidden state of the LSTM. This is achieved using an LSTM layer, followed by a dropout layer for regularization, and a fully-connected layer for dimensionality reduction. The decoder then processes this compressed representation using the reverse architecture (FC, dropout, and LSTM) to reconstruct the original data.

In this architecture, fully connected (FC) layers process the final hidden state of the LSTM (the last timestep in the sequence) as a single, comprehensive representation of the entire input sequence. This output is then projected into a lower-dimensional latent space through a linear transformation. While this approach is computationally efficient, it assumes that the LSTM’s final hidden state sufficiently captures all relevant temporal dependencies. As a result, it can struggle to retain fine-grained temporal details, especially for data like EEG signals, where localized patterns are crucial.

The *Convolutional LSTM Autoencoder (C-LSTM)* enhances feature extraction by combining LSTM layers with

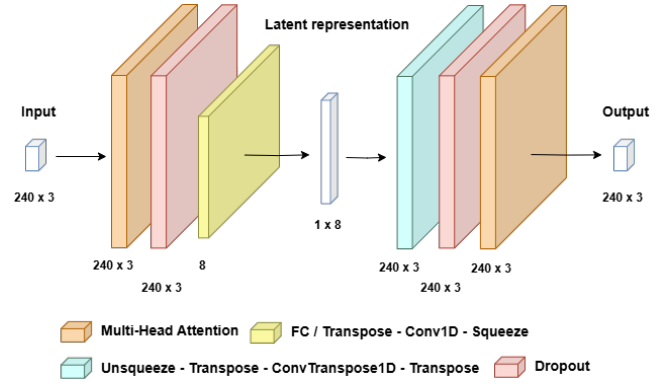


Figure 1: AE_err/AE_att architecture.

convolutional layers. The key difference between the LSTM and C-LSTM architectures lies in how dimensionality reduction is achieved in the encoder and decoder: the first uses fully connected (FC) layers, while the latter uses Conv1D layers. These operate directly on the sequence of hidden states produced by the LSTM. By applying a kernel across the temporal dimension, they extract localized patterns and dependencies within the sequence. This convolutional operation integrates information from multiple timesteps, creating a more nuanced and structured representation. After convolution, the output is reduced to a lower-dimensional latent space, where temporal features are preserved and compactly encoded. This method emphasizes localized temporal dynamics while reducing dimensionality.

The *Attention-based Autoencoder (AE_err)* combines convolutional layers with an attention mechanism to reconstruct the input data, again using the convolutional layers to capture local features while replacing LSTM with attention to capture long-term patterns.

Fig. 1 illustrates this architecture. The input to the encoder consists of a sequence of 240 measurements from two simultaneous channels, along with a same-size sequence of time representation. The input data passes through the encoder and is compressed into a latent representation of size 1×8 , where the first dimension represents a single compressed time step, and the second dimension corresponds to eight learned features. This is achieved by applying multi-head attention to the input data, followed by a dropout layer for regularization, and a convolutional layer for dimensionality reduction. The decoder then processes this compressed representation using the reverse architecture (transposed convolution, dropout, and attention) to reconstruct the original data, capturing both trends and amplitudes.

In all three autoencoding systems, anomaly detection is based on the assumption that the compressed latent representation preserves only recurring, periodic fluctuations and trends. Atypical spikes, often caused by noise, do not cause enough loss to be worth the space to represent them (in terms of nodes in the latent representation). Therefore, the level of anomaly is estimated as the reconstruction error, the difference between the original input sequence and the sequence decoded from the latent representation.

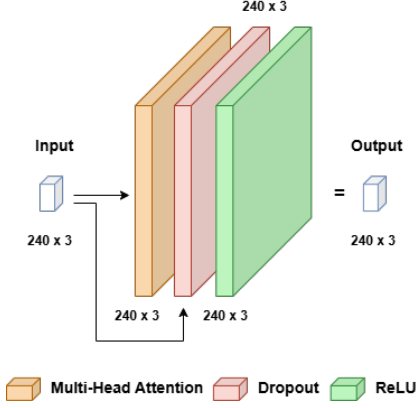


Figure 2: TP_err/TP_att architecture.

The *Transformer Predictor (TP_err)* method utilizes a Transformer that uses multi-head attention for time-series forecasting. It consists of a sequence-to-sequence architecture, with multi-head attention layers in both the encoder and decoder components. The encoder captures temporal dependencies in the input sequence, while the decoder predicts the next sequence based on these encoded features.

Fig. 2 illustrates this architecture. The input consists of a sequence of 240 measurements from two simultaneous channels, along with a same-size sequence of time representation. This data passes through the embedding layer, which applies multi-head attention to capture key temporal patterns. The output is then processed with dropout regularization and passed through a LeakyReLU activation function to introduce non-linearity.

The TP_err method estimates noise through prediction error, i.e. the difference between the model's predicted output and observed values. When the model predicts the next time step in the sequence, a significant difference between the predicted and actual values suggests that the input data may be noisy.

3.3. Attention-Based Detection

The exact same architectures as in AE_err and TP_err above are also used for the respective attention-based detection methods *Attention-based Autoencoder (AbAE_att)* and *Transformer Predictor (TP_att)*, except that now anomaly is estimated by the attention weights.

The key idea is that when a part of the sequence receives a low attention weight, it suggests that the information at that moment is likely noisy and, therefore, unimportant for the task.

3.4. Conventional Approaches

Alongside machine learning methods, conventional techniques provide a reliable alternative for detecting noise in time-series data and are commonly utilized. Among the various conventional methods proposed by many researchers

over the years, a statistical method and a signal-processing technique have been considered for this work.

Principal Component Analysis (PCA) is a linear dimensionality reduction technique. It analyzes a data table where observations are described by quantitative dependent inter-correlated variables. With PCA, the important information is extracted from the table and represented in a new set of orthogonal variables, called principal components [13]. The similarity patterns of the observations and the variables can then be depicted as points on maps [14].

Multinomial Noise Exponential filtering (MNE) is a method used to remove unwanted frequencies from time-series data, typically EEG signals. It applies a bandpass filter to the data, allowing signals within a specific frequency range to pass through while reducing frequencies outside this range.

4. Experimental Setup

4.1. Test case description

This study employs a case analysis to validate the efficacy of the proposed noise estimation methods. The focus is detecting noise within electroencephalographic (EEG) signals, recorded via headbands utilized during slumber to monitor various stages of the user's sleep cycle.

One of the main concerns when dealing with electroencephalographic signals (EEG) is assuring that clean data with a high signal-to-noise ratio is recorded. The EEG signal amplitude is in the microvolts range, and it is easily contaminated with noise, known as artifacts, which need to be filtered from the neural processes to keep the valuable information needed for different applications.

In this domain, an artifact is denoted as any component of the EEG signal not directly produced by human brain activity, making the system register noise that contaminates the neural EEG data. The ability to recognize these artifacts is the first step in removing them. EEG artifacts can be classified depending on their origin, which can be physiological or external to the human body (technical/non-physiological). Fig. 3 shows some examples of EEG signals contaminated with noise.

In this figure, pink-colored areas are typical of low-frequency noise (0.2-4Hz), which is mainly due to perspiration originating from small drops of sweat produced by the skin glands, which cause changes in the electrical baseline of the electrodes. Brown-colored and green-colored areas are typical of high-frequency noise (30-45Hz) that may originate from electrical activity produced by the muscles when they are contracted, like, for example, muscle tension in the jaw or forehead that can take place when clenching or frowning, respectively. Orange-colored areas are typical of high-amplitude noise that may be due to temporary failures in contact between the EEG sensor and the scalp produced by touching the sensor or by spontaneous changes in electrode-skin contact.

In this analysis, EEG signals collected from a headband with two EEG channels synchronized with medical-grade

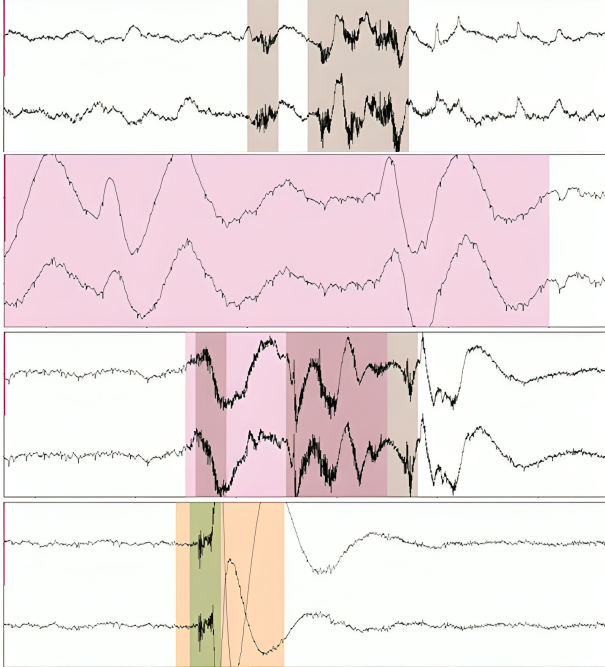


Figure 3: Examples of the temporal visualization of EEG signals and coloring of noisy segments.

EEG devices are considered. The recorded signal frequency for each channel of the headband is 128Hz. These signals are grouped into consecutive 30-second segments, and each segment is further annotated by three experts, into one of five sleep stages: Wake, N1, N2, N3, and REM. These stages correspond to specific brain activity patterns, such as slow eye movements, sleep spindles, and other characteristic waveforms. The goal of the current analysis is to detect the presence of noise, if any, at any of these segments. In total, 56 recordings were examined. Each recording refers to the EEG signals acquired by the headband of a user during a night-long sleep. These signals are grouped into 30-second consecutive segments.

As a ground-truth method to evaluate the results of the proposed noise estimation methods, the original estimation method used by the headband manufacturer is considered. This method consists of task-specific algorithms to automatically estimate noise and evaluate the overall quality of EEG signals, identifying the artifacts in recordings made using the wearable textile headband.¹

4.2. Data Pre-processing

Data normalization is required for the machine learning methods to ensure robust model training and reliable outcomes. It is performed using statistics derived from the

1. The device manufacturers and data providers are co-authors, so further details are withheld to preserve anonymity. More details about the device and the data (including data access) will be made available in the camera-ready.

TABLE 2: Absolute number and percentage of segments with artifacts in each channel and recording session.

RecID	A	B	C	D
HB1	0	57 (5%)	27 (3%)	9 (1%)
HB2	151 (14%)	56 (5%)	26 (3%)	9 (1%)

entire training dataset instead of relying on per-batch calculations. The median and *interquartile range (IQR)* are used instead of the mean and standard deviation, respectively, to improve robustness against outliers. The median measures central tendency and is resistant to extreme values, and IQR is a measure of statistical dispersion that represents the spread of the middle 50% of a dataset and reduces the influence of outliers by focusing on the range within which the central portion of the data lies.

PCA and MNE require raw data that preserves the characteristics of the original signal, so no normalization is applied for these methods.

To train the machine learning models, the available recordings were randomly split into training (43 recordings), validation (3 recordings), and testing subsets (10 recordings). Among the testing recordings, six had zero or one noisy segments are not reported here as it makes little sense to compare methods on having identified a single datapoint. For the record, only PCA identified the noisy segment in only one of the segments.

Table 2 gives the absolute and relative noise density for the remaining four recordings, arranged in descending order. Specifically, the table gives the count of 30-second segments within which at least one artifact is annotated as noise in the ground-truth labeling. Note that recordings are *not* of equal length.

What is noteworthy is that Recording A has considerably different noise rates between the two channels. Since the signals from the two channels are heavily correlated, this presents an opportunity for the machine learning methods to demonstrate learning a comparative model that exploits the fact that the noisy parts of HB2 are not correlated to their corresponding parts in HB1. However this is a difficult theory to construct, since it requires ‘discovering’ correlation first.

4.3. Training Process and Hyper-parameters

The implementations of the methods in Table 1 that we used in our experiments are publicly available. Our implementation of the machine learning methods have been reposit in Zenodo.² PCA is implemented using the Frouros³ open-source Python library for drift detection. MNE is used as implemented and configured in the ‘MNE

2. To preserve anonymity the Zenodo record is not published but left as a draft; reviewers can access it at <https://zenodo.org/records/14842198> (must use this link to include the access token)

3. Documentation and source code are available at <https://github.com/IFCA-Advanced-Computing/frouros>.

Tools’ Python package for exploring, visualizing, and analyzing human neurophysiological data.⁴ MNE is a bandpass filter with pre-configured cut-offs specifically targeting EEG data.

To process the data with the *AE* and *TP* methods, we split each 30-second segment of EEG data (3840 samples per channel) into 16 smaller chunks of 240 samples each. This point is moot for *LSTM*, *C-LSTM*, *PCA* and *MNE*.

The Autoencoder methods are trained on sequence reconstruction. The Predictor methods are trained on predicting sleep stage. To balance the autoencoders’ focus on both the amplitude and trends of the time-series data, we define a custom loss function, called *BlendedLoss*. This function combines the median and mean of the powered absolute differences between the predicted values (\hat{x}) and the target values (x):

$$\text{Loss} = (1 - b) \cdot \text{median}(|\hat{x} - x|^p) + b \cdot \text{mean}(|\hat{x} - x|^p)$$

where p is the power parameter that controls the sensitivity of the loss to differences and b the *blend factor* that controls the trade-off between learning the overall trend (median error) and closely following local patterns (mean error).

Revisiting our research hypothesis, we expect methods trained on predicting sleep stage to outperform methods trained on autoencoding since they have access to task-specific labels. More generally our methods are ordered as follows in terms of task-specific knowledge they have access to:

- 1) *MNE* is specifically, expertly designed to detect noise in EEG signals and has proven to be very effective on this task. *PCA* is a general-purpose drift detection method, but for the purposes of this paper has been manually calibrated on the training data to look for ‘drifts’ that have the expected duration before returning to the previous data distribution.
- 2) *AE_att* and *TP_att* implement our hypothesis that the attention mechanism can exploit supervision unrelated to anomaly detection to automatically extract task-specific knowledge of anomalies.
- 3) *LSTM*, *C-LSTM*, *AE_err* and *TP_err* are general-purpose autoencoders with minimal access to task-specific knowledge in the form of a task-specific loss function.

This ordering reflects our prior expectation regarding their relative performance. *None of the above (neither machine learning nor manual calibration) has ever had any access to noise annotations in the training data, but only to sleep-stage annotations.*

For the training configuration, we set the batch size to 512 and trained the models for a maximum of 1000 epochs with patience of 30 epochs, meaning that if the validation loss does not improve for 30 consecutive epochs, training will stop. We set the learning rate to 1e-4 and use the Adam optimizer to adjust the model weights. Additionally, we implement a ReduceLROnPlateau scheduler to adjust the

learning rate dynamically based on the validation loss. If the validation loss plateaus, the scheduler reduces the learning rate to push the model to fit finer patterns.

5. Results and Discussion

We prepared three validation/testing setups. In the first setup validation and testing scores are calculated on the original 30-sec segments where each segments is annotated as ‘noise’ if it includes at least one sample marked as ‘noise’ in the ground-truth annotation. In the 5min and 10min-window setups segments evaluation was performed on 5min and 10min windows. If a method reports noise on any part of the window and the ground truth also annotates as noisy any part of the window, the window counts as a true positive. Obviously, these are easier tasks than the original 30-sec segmentation on two grounds: (a) positives get less sparse so task becomes less of a needle-in-a-haystack problem, and (b) learners have access to a longer context for the same number of input token, which know to be an important factor in sequence processing [15].

Table 3 presents F-score, Precision, and Recall metrics for the evaluated methods on all three evaluation setups. Since all methods give a numerical estimation and not a binary decision, a threshold was established. The threshold was calculated on the training data, separately for each method, as the threshold that makes 1% of the training data come out as noise. This is the only piece of prior domain knowledge shared by all learners. F-score is calculated for $\beta = 2$, in order place more weight on recall than precision. This is due to the fact that the dataset is very unbalanced and noise instances are rare, even in the 10min-window setup. F-score is adapted accordingly to reflect that recall is harder to achieve than precision.

The first observation is that by comparing the three sub-tables we can immediately see that the three windows behave as expected, with results improving with longer windows. By comparing horizontally from right to left we also see that results improve as noise gets denser.⁵ We cannot speculate on whether the top-to-bottom improvement is due to the longer context or due to denser positives. But either way, we have clear indications about how task difficulty scales across the table.

What is directly relevant for our hypothesis are the comparisons between (a) attention-based anomaly detection (*AE_att* and *TP_att*) against the error-based anomaly detection from the same models (*AE_err* and *TP_err*), and (b) models trained on the sleep stage task (*TP*) against autoencoded models (*AE*). *TP_att* generally outperforms all four *TP/AE_err/att* combinations, although autoencoding proved better on HB1-C (primarily *LSTM* and secondarily *AE_att*) with the gap becoming more pronounced on the denser setups. Although the conditions present in HB1-C are worth investigating, the results show that the ideas presented here are validated and promising.

4. Available from <https://mne.tools> and Pypi.

5. Note that HB1-A is out-of-order in this respect, as it has no noise.

TABLE 3: F-score, precision and recall on for different segment length. The F-score is calculated for $\beta = 2$, emphasising recall over precision, as it appropriate in our setting. Precision and recall, respectively, are given in parenthesis; When absent, it denotes that they are within three percentile points of the F2 score. Empty cells denote that the respective method did not mark any segment as noise, which implies that recall is 0% and precision is undefined.

(a) Evaluation on 30sec windows.

RECID	HB1				HB2			
	A	B	C	D	A	B	C	D
LSTM		2%			2% (6%,1%)	4% (7%,4%)	4% (3%,4%)	
C-LSTM			4% (11%,4%)					
AE_ERR								
TP_ERR								
AE_ATT		4%			1% (20%,1%)			
TP_ATT		2%			15% (13%,15%)	10% (5%,14%)	4% (2%,8%)	3% (1%,11%)
PCA		88% (100%,85%)		14% (100%,11%)	6% (41%,5%)	2% (1%,4%)	60% (24%,96%)	26% (67%,22%)
MNE					2% (40%,1%)			

(b) Evaluation on 5min windows.

RECID	HB1				HB2			
	A	B	C	D	A	B	C	D
LSTM		13% (6%,18%)			19% (12%,22%)	22% (8%,38%)	11% (3%,35%)	
C-LSTM			21% (10%,30%)		1% (2%,1%)			4% (1%,11%)
AE_ERR		16% (27%,14%)			4% (10%,3%)			
TP_ERR								
AE_ATT		9% (3%,18%)	10% (3%,22%)		7% (18%,6%)			
TP_ATT		14% (6%,19%)			42% (14%,84%)	26% (7%,100%)	16% (4%,100%)	5% (1%,100%)
PCA			83% (77%,85%)	87% (80%,89%)	21% (25%,21%)	10% (3%,27%)	49% (17%,96%)	87% (80%,89%)
MNE				96% (82%,100%)	15% (38%,13%)			14% (100%,11%)

(c) Evaluation on 10min windows.

RECID	HB1				HB2			
	A	B	C	D	A	B	C	D
LSTM		19% (7%,37%)	5% (2%,11%)		28% (12%,42%)	25% (7%,64%)	9% (2%,42%)	
C-LSTM			34% (11%,67%)		1% (1%,1%)			21% (5%,100%)
AE_ERR		31% (30%,32%)			6% (8%,5%)			
TP_ERR								
AE_ATT		12% (3%,35%)	26% (7%,89%)		9% (14%,9%)			
TP_ATT		21% (7%,39%)	2%		41% (13%,90%)	22% (5%,100%)	14% (3%,100%)	5% (1%,100%)
PCA			78% (58%,85%)	96% (82%,100%)	36% (26%,40%)	19% (5%,63%)	44% (14%,96%)	96% (82%,100%)
MNE				96% (82%,100%)	25% (35%,23%)			96% (82%,100%)

What is also noteworthy is that TP_att outperforms the task-calibrated conventional methods on HB2-A and HB2-B. Regarding HB2-A it is worth investigating whether the transformer has ‘discovered’ the lack of correlation between HB1-A and HB2-A on noisy segments. Regarding HB2-B, a possible explanation is that TP_att is more robust to noise density considerably above the hard-wired 1%: it suffers the unavoidable precision loss just as all methods, but it gets a comparatively higher F-score through perfect recall.

6. Conclusions and Future Work

This paper investigates the extraction of anomaly indicators from *intermediate* layers of a deep neural network

trained on a sequence processing task. By removing the need to train on a sequence reconstruction/prediction task, we are free to train on whatever task-specific supervision is available. This alleviates the most pressing challenge of anomaly detection, namely that anomalies are vaguely defined and thus difficult to directly supervise, while simultaneously allowing task-specific knowledge to be extracted from whatever other (non anomaly-related) supervision might be available.

Our ideas are positively validated on EEG signals sourced from wearable sleep-monitoring devices. This dataset has the great advantage of having both task supervision (sleep stage detection) and noise supervision, both expertly annotated by the device manufacturer. Our current

findings already indicate that the attention layer of a Transformer trained on the sleep stage task exhibits promising performance when used as an anomaly detector.

Subsequent steps must focus in two directions: exploring the current experimental setup at more depth and expanding the breadth of the results with diverse datasets and tasks. Regarding the first direction, we mostly plan to investigate how the components of the attention layer (Q, K, and V) interact in order to understand whether a subset of these components is more tightly fitted on recognizing patterns that should be ignored. Another investigation would be to understand if the performance on HB2-A really indicates that the Transformer has ‘discovered’ the lack of correlation between HB1-A and HB2-A on noisy segments or is incidental. Despite the inherent difficulty in trying to understand the inner workings of neural networks even for single-layer attention, as the the field of *explainable machine learning* evolves the required tools are starting to materialize.

Regarding the second direction, the most immediate goal is to devise a second task on the same data besides sleep stage detection. This would allow us to experiment with training on different tasks and jointly training on multiple tasks, which we expect to give insights about what tasks are more appropriate as drivers for anomaly detection. A less directly accessible goal is to test our idea on different use cases and domains. Although both sequence processing and anomaly detection are ubiquitous, the existence of simultaneous task/noise supervision is less so. But we are confident that the exposure gained through publishing our EEG results will create an interest in compiling and publishing such datasets.

References

- [1] N. Laptev, S. Amizadeh, and I. Flint, “Generic and Scalable Framework for Automated Time-series Anomaly Detection,” in *Proc. of the 21th ACM SIGKDD Intl Conf. on Knowledge Discovery and Data Mining (KDD '15)*, Aug. 2015. [Online] DOI: [10.1145/2783258.2788611](https://doi.org/10.1145/2783258.2788611)
- [2] A. Geiger, D. Liu, S. Alnegheimish, K. Veeramachaneni, and A. Cuesta-Infante, “TadGAN: Time series anomaly detection using Generative Adversarial Networks,” in *2020 IEEE Intl Conference on Big Data (Big Data)*, Dec. 2020. [Online] DOI: [10.1109/bigdata50022.2020.9378139](https://doi.org/10.1109/bigdata50022.2020.9378139)
- [3] Z. Ji, J. Gong, and J. Feng, “A Novel Deep Learning Approach for Anomaly Detection of Time Series Data,” *Scientific Programming*, vol. 2021, pp. 1–11, Jul. 2021. [Online] DOI: [10.1155/2021/6636270](https://doi.org/10.1155/2021/6636270)
- [4] Y. Zhou, A. M. Al-Ahmari, Z. Li, N. Wu, and H. Ren, “Anomaly detection via a combination model in time series data,” *Applied Intelligence*, vol. 51, no. 7, Jan. 2021. [Online] DOI: [10.1007/s10489-020-02041-3](https://doi.org/10.1007/s10489-020-02041-3)
- [5] N. Mejri, L. Lopez-Fuentes, K. Roy, P. Chernakov, E. Ghorbel, and D. Aouada, “Unsupervised anomaly detection in time-series: An extensive evaluation and analysis of state-of-the-art methods,” *Expert Systems with Applications*, vol. 256, 2024. [Online] DOI: [10.1016/j.eswa.2024.124922](https://doi.org/10.1016/j.eswa.2024.124922)
- [6] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, 1997. [Online] DOI: [10.1162/neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735)
- [7] H. Hojjati, T. K. K. Ho, and N. Armanfard, “Self-supervised anomaly detection in computer vision and beyond: A survey and outlook,” *Neural Networks*, vol. 172, 2024. [Online] DOI: [10.1016/j.neunet.2024.106106](https://doi.org/10.1016/j.neunet.2024.106106)
- [8] K. U. Duja, I. A. Khan, and M. Alsuhaibani, “Video surveillance anomaly detection: A review on deep learning benchmarks,” *IEEE Access*, vol. 12, 2024. [Online] DOI: [10.1109/ACCESS.2024.3491868](https://doi.org/10.1109/ACCESS.2024.3491868)
- [9] A. Sgueglia, A. Di Sorbo, C. A. Visaggio, and G. Canfora, “A systematic literature review of iot time series anomaly detection solutions,” *Future Generation Computer Systems*, vol. 134, pp. 170–186, 2022. [Online] DOI: [10.1016/j.future.2022.04.005](https://doi.org/10.1016/j.future.2022.04.005)
- [10] L. Basora, X. Olive, and T. Dubot, “Recent advances in anomaly detection methods applied to aviation,” *Aerospace*, vol. 6, no. 11, 2019. [Online] DOI: [10.3390/aerospace6110117](https://doi.org/10.3390/aerospace6110117)
- [11] Z. Zamanzadeh Darban, G. I. Webb, S. Pan, C. Aggarwal, and M. Salehi, “Deep learning for time series anomaly detection: A survey,” *ACM Comput. Surv.*, vol. 57, no. 1, Oct. 2024. [Online] DOI: [10.1145/3691338](https://doi.org/10.1145/3691338)
- [12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems 30*, 2017. [Online]. Available: <http://papers.nips.cc/paper/7181-attention-is-all-you-need>
- [13] R. Bro and A. K. Smilde, “Principal component analysis,” *Analytical Methods*, vol. 6, no. 9, pp. 2812–2831, 2014. [Online] DOI: [10.1039/c3ay41907j](https://doi.org/10.1039/c3ay41907j)
- [14] H. Abdi and L. Williams, “Principal component analysis,” *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 2, pp. 433–459, 2010. [Online] DOI: [10.1002/wics.101](https://doi.org/10.1002/wics.101)
- [15] M.-C. Lee, J.-C. Lin, and E. G. Gran, “How far should we look back to achieve effective real-time time-series anomaly detection?” in *Advanced Information Networking and Applications*. Springer, 2021. [Online] DOI: [10.1007/978-3-030-75100-5_13](https://doi.org/10.1007/978-3-030-75100-5_13)