



Deposited via The University of York.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/225992/>

Version: Accepted Version

Article:

Xing, Zeyu, Mehmood, Owais and Smith, William Alfred Peter (2025) Unsupervised anomaly detection with a temporal continuation, confidence-aware VAE-GAN. Pattern Recognition. 111699. ISSN: 0031-3203

<https://doi.org/10.1016/j.patcog.2025.111699>

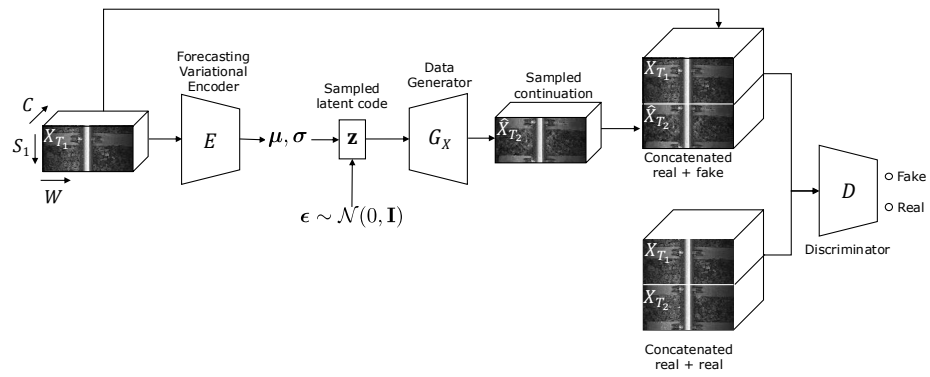
Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



- 1 Graphical Abstract
- 2 **Unsupervised anomaly detection with a temporal continuation, confidence-aware**
- 3 **VAE-GAN**
- 4 Zeyu Xing, Owais Mehmood, William A. P. Smith

5 **Highlights**

6 **Unsupervised anomaly detection with a temporal continuation, confidence-aware**
7 **VAE-GAN**

8 Zeyu Xing, Owais Mehmood, William A. P. Smith

- 9 • Propose an unsupervised, zero-shot anomaly detection method for spatiotempo-
10 ral signals, separating anomalies in predictable regions from unimportant stochas-
11 tic variations
- 12 • Our method is based on using a *forecasting VAE-GAN* to learn the space of plau-
13 sible continuations of a temporal sequence
- 14 • We make the model *confidence-aware* by also learning to predict the pointwise
15 confidence of the reconstruction, allowing us to separate structural from stochas-
16 tic uncertainty
- 17 • Achieve state-of-the-art performance on the ECG5000 [1, 2] and MIT-BIH [3]
18 time series anomaly detection datasets

Unsupervised anomaly detection with a temporal continuation, confidence-aware VAE-GAN

Zeyu Xing^a, Owais Mehmood^b, William A. P. Smith^a

^a*Department of Computer Science, University of York, York, UK*

^b*Omnicom Balfour Beatty, York, UK*

Abstract

We propose an unsupervised approach to anomaly detection in data with a temporal dimension. We adapt the VAE-GAN architecture to learn the proxy task of temporal sequence continuation. Rather than reconstructing the input, our variational decoder decodes to a forecast of the future sequence. In order to separate structural uncertainty (which our model can reconstruct by fitting to observed data) from stochastic uncertainty (which it cannot) we introduce an additional decoder that outputs the pointwise confidence of the prediction, after the optimal latent-variable has been found. We can use this for zero-shot anomaly detection, separating anomalies from stochastic variation that cannot be modelled, without any examples. This is important for domains in which anomalies are so rare that it is not possible or meaningful to train a supervised model. As an example of such a domain, we introduce a new dataset comprising linescan imagery of railway lines which we use to illustrate our methods. We also achieve state-of-the-art performance on the ECG5000 and MIT-BIH time series anomaly detection datasets. We make an implementation of our method available at <https://github.com/YorkXingZeyu/ECG-VAEGAN-Project>.

Keywords: time series anomaly detection, unsupervised anomaly detection, variational autoencoder, VAE-GAN

1. Introduction

Temporal sequential data arises across a whole host of data modalities from time series to video to audio. For such data, sequence continuation, completion, interpolation or reordering are emerging as promising proxy tasks for self-supervised feature

29 learning. The overarching premise is that, in order to reason about the future, ordering
30 or interpolation, it is necessary to learn a model that extracts not only low level fea-
31 tures but high level concepts as abstract as physical laws (for example, predicting that
32 a falling ball will bounce). Temporal sequences can be further subdivided into those
33 where the observations are overlapping and non-overlapping. Overlapping sequences
34 may observe the same part of the world at different times. For example, adjacent video
35 frames are likely to contain many of the same scene components. Such sequences can
36 be handled in a special way by explicitly modelling the relationship between the same
37 points at different times, for example using optical flow motion fields. This makes the
38 task of future prediction easier since it can, at least partly, be posed as motion prediction
39 of observed scene components.

40 In this paper, we focus on *non-overlapping* temporal sequences. Examples include
41 audio streams, time series data and linescan images from pushbroom cameras. We
42 propose a generative framework for self-supervised feature learning and anomaly de-
43 tection based on continuation of such sequences. We use a VAE-GAN [4] as our under-
44 lying architecture. The GAN discriminator component ensures that continuations are
45 natural and realistic, by encouraging them to follow the distribution of real complete
46 sequences. This avoids blurring multiple possible futures together. The VAE latent
47 variational variable model captures the stochasticity of future prediction. The distribu-
48 tion mean computed by our variational encoder can be seen as capturing the predictable
49 elements of the future which depend only on the observed portion of the data. The ran-
50 dom sampling process from the resulting latent distribution can be seen as exploring
51 possible futures. Within this space we expect to be able to reconstruct structural as-
52 pects of the actual future but not stochastic ones. For example (see Figure 5), if we
53 observe a section of a linescan image of a railway line, this constrains the positioning
54 of the next sleeper to a small range of possibilities (structural uncertainty) but the exact
55 configuration of the ballast stones cannot be meaningfully constrained (stochastic un-
56 certainty). We therefore augment the VAE-GAN model with an additional decoder that
57 predicts spatially varying confidence, i.e. the remaining pointwise similarity once the
58 optimal sample from the latent space has been found. Using the same example, we ex-
59 pect high confidence to be assigned to sleepers and low confidence to the ballast. Once

60 trained, our model learns an efficient encoder of the observed data that can be used as
61 a pretrained backbone for downstream tasks. However, the model can additionally be
62 used for *unsupervised* anomaly detection. Where a high confidence region cannot be
63 reconstructed accurately, we can assume the feature is anomalous. It is on this task that
64 we evaluate our proposed model.

65 While *supervised* anomaly detection methods provide state-of-the-art performance
66 in some domains, for some problems anomalies are so rare that a supervised approach
67 is not possible. For example, in rail surveying, we would like to detect anomalies that
68 have never been observed before. Severe anomalies such as cracks in the railhead are
69 so rare that only single examples may be observed over a period of decades. Posing this
70 as a supervised or weakly supervised problem leads to such severe class imbalance that
71 such approaches fail to learn any meaningful features. On the other hand, unsupervised
72 approaches can use the abundance of non-anomalous data to learn a rich model of
73 normal appearance and treat anomaly detection as the problem of detecting out-of-
74 distribution features. It is this problem setting that we address with the particularly
75 challenging case of also learning to ignore uninteresting stochastic variations.

76 Our contributions are as follows:

- 77 1. We propose an unsupervised, zero-shot anomaly detection method for spatiotem-
78 poral signals, separating anomalies in predictable regions from unimportant stochas-
79 tic variations;
- 80 2. Our method is based on using a *forecasting VAE-GAN* to learn the space of plau-
81 sible continuations of a temporal sequence;
- 82 3. We make the model *confidence-aware* by also learning to predict the pointwise
83 confidence of the reconstruction, allowing us to separate structural from stochas-
84 tic uncertainty in a self-supervised manner;
- 85 4. We achieve state-of-the-art performance on the ECG5000 [1, 2] and MIT-BIH [3]
86 time series anomaly detection datasets while also showing application to linescan
87 imagery on a new rail track surveying dataset.

88 While our method is general and could, in principle, be applied to temporal sequential
89 data from any domain, our evaluation focusses on linescan images and time series data

90 (specifically electrocardiogram traces).

91 **2. Related work**

92 *2.1. Self-supervised and generative models*

93 Most commonly, self-supervised learning refers to *feature learning* [5]. Here, self-
94 supervision is used for pretraining only, to discover useful representations of data that
95 are subsequently fine-tuned for other tasks. Examples of proxy tasks that have been
96 used for this purpose include predicting relative position of two regions in the same
97 image [6], colourisation [7], orientation prediction [8], video frame ordering [9], video
98 playback direction [10] and cycle-consistent point tracking [11]. Although these meth-
99 ods obviate the need for supervision, they only provide a route to *feature learning* -
100 i.e. they do not solve any useful task directly, just provide learnt features that can be
101 used for subsequent fine-tuning for a specific task. Another class of approaches use
102 generative models such as GANs. Here, a discriminator or critic provides a supervi-
103 sion signal from an unlabelled dataset while some component of the model learns to
104 extract useful features. Bi-directional GAN (BiGAN) [12] is a variant of a conven-
105 tional GAN in which an encoder is additionally learnt that maps from the data space
106 to the latent space. Our approach also learns to encode from data space to latent space
107 but this forms only the conditioning signal of our generator, like a conditional GAN
108 [13] and, rather than learning to reconstruct data from the latent space, we learn to pre-
109 dict temporal sequence continuations along with confidence in our prediction. Kingma
110 and Welling [14] introduced the Variational Autoencoder (VAE), a powerful generative
111 model that combines variational inference with autoencoders. This method has proven
112 effective in generating realistic data and learning latent representations. Similarly, He
113 et al. [15] introduced Masked Autoencoders (MAEs) which have demonstrated their
114 scalability and effectiveness in vision learning by leveraging masked signal modelling
115 to improve the representation learning capability of autoencoders. VAE-GANs have
116 been used for stochastic future video frame prediction [16], however we are the first to
117 tackle the problem of non-overlapping sequential data and to introduce estimation of
118 the spatially-varying confidence of the future prediction. Recent advancements in self-

119 supervised learning, such as [17] introduced the Joint-Embedding Predictive Architec-
120 ture (JEPA) which has further improved visual representation learning by predicting the
121 embedding of masked or missing portions of images. While VAEs have been widely
122 adopted for various applications, such as image generation and data compression, we
123 extend these concepts to tackle the problem of non-overlapping sequential data and
124 introduce the estimation of spatially-varying confidence for future predictions.

125 2.2. Temporal models

126 Generative modelling for self-supervised learning has been applied to a number
127 of different data modalities. For time series samples, self-supervision can learn the
128 underlying structural features of unlabeled time series by exploring the inter-sample
129 relationship and intra-time relationship of time series [18]. When dealing with audio
130 or speech data, it is often necessary to convert them into feature vectors [19]. Giri et
131 al. [20] use self-supervised learning to learn a compact representation of normal data
132 using self-supervised classification of metadata based on audio files, to detect anoma-
133 lies in sound data. At the same time, self-supervised pretraining for Automated Speech
134 Recognition (ASR) also makes great progress in processing audio data [21]. Later,
135 based on ASR and to supplement the ability to compare learning in self-supervision,
136 [22] proposed to co-learn the presentation from different models of speech and literacy
137 during pre-training.

138 For video data, the first is Arrow of Time, which will help tell whether a video is
139 running forward or backward [23]. Since video data cannot be captured simply through
140 a two-dimensional CNN, some researchers propose to use three-dimensional CNN to
141 solve space-Time cubic puzzles of videos [24]. Long short term memory (LSTM) net-
142 works tend to be used when processing such temporal data. [25] tried to use LSTM
143 to learn the representation of time series, using the encoder-decoder LSTM model to
144 rearrange the shuffled input sequence in the correct order. Tao et al. [26] propose the
145 pretext-Contrastive Learning (PCL) model on the basis of pretext-task and compari-
146 son learning and applied it to self-supervised video feature learning. Similarly, the
147 Video-based Temporal-Discriminative Learning (VTDL) framework is used to process
148 unlabelled video data [27]. For the video future prediction task, the purpose is to pre-

149 dict the future frame sequence or the future frame sequence feature. The idea is to
150 make predictions by parsing a given finite number of video frames [28]. For models,
151 the hope is that they can learn the dynamics of these known sequences of frames, the
152 more famous of which is the LSTM [29], and many methods have been proposed af-
153 terwards [30, 31, 32, 33, 34]. Many algorithms use LSTM to deal with time dynamic
154 problems in video [31, 33, 34]. These methods can be used in some self-supervised
155 feature learning tasks, and the advantage is that no manual labelling of data is required.
156 MCnet [34] has two encoders that learn the spatial features of the image and the tem-
157 poral dynamics of the video. They output temporal and spatial characteristics of the
158 data, which are fed into the decoder to predict future videos.

159 In the exploration of time series anomaly detection, many outstanding methods
160 have achieved remarkable results. Giannoulis et al. [35] presents Ditan, a deep-learning
161 domain-agnostic framework tailored for the detection and interpretation of anomalies
162 in multivariate time series data. The framework employs Convolutional Neural Net-
163 works (CNNs) to extract local features from time series data and LSTMs to capture
164 long-term dependencies in the sequences to identify temporal patterns and anomalies
165 across various datasets and applications, demonstrating its adaptability and effective-
166 ness in handling diverse time series anomaly detection tasks. Audibert et al. [36] ex-
167 plore the role of deep neural networks in multivariate time series anomaly detection.
168 The study utilizes deep learning techniques such as CNNs, Recurrent Neural Networks
169 (RNNs) and their variant LSTM networks, as well as Autoencoders. These models ef-
170 fectively capture complex temporal dependencies and patterns, significantly improving
171 the performance of anomaly detection in complex multivariate time series. Mokoena et
172 al. [37] address the challenge of explaining anomalies detected in time series data us-
173 ing a method called sequential explanations. It underscores the importance of not just
174 identifying anomalies but also understanding their underlying causes. The proposed
175 method provides detailed, step-by-step explanations for detected anomalies, enhancing
176 interpretability and aiding in more informed decision-making for time series anomaly
177 detection.

178 Pereira and Silveira [38] explore learning representations from healthcare time se-
179 ries data for unsupervised anomaly detection. The study utilizes Autoencoders to learn

180 normal data patterns and detect reconstruction errors, CNNs to extract local features
181 and patterns, and RNNs along with LSTM networks to capture temporal dependen-
182 cies. By combining these models, the research effectively extracts meaningful features
183 from complex healthcare time series data to achieve efficient unsupervised anomaly
184 detection.

185 2.3. Anomaly detection

186 Significant advancements have also been made in the exploration of unsupervised
187 and semi-supervised learning methods. Yang et al. [39] propose an unsupervised
188 anomaly detection and segmentation method by learning deep feature correspondence.
189 The method effectively detects and segments anomalies without the need for labeled
190 data, using deep neural networks to automatically extract relevant features from the in-
191 put data. The approach demonstrated superior performance across multiple real-world
192 datasets. Zhang et al. [40] introduces a novel deep anomaly detection method combin-
193 ing self-supervised learning and adversarial training. By employing Generative Adver-
194 sarial Networks (GANs), the model is able to self-supervise during the training process,
195 thereby improving the accuracy and robustness of anomaly detection. Experimental
196 results show that this method significantly enhances detection performance across var-
197 ious datasets. Zavrtnik et al. [41] presents a visual anomaly detection method based
198 on image inpainting, utilizing inpainting techniques to detect and localize anomalous
199 regions. By comparing the normal parts of an image with the inpainted version, the
200 method effectively identifies and marks anomalies. It demonstrated high efficiency and
201 accuracy in multiple visual detection tasks. Similar to our approach, Zhou et al. [42]
202 use an autoencoder but propose to use the latent representation itself as part of the
203 anomaly detection process. However this approach requires weak supervision whereas
204 ours is completely unsupervised.

205 Akcay et al. [43] introduce GANomaly, a semi-supervised anomaly detection method
206 using adversarial training. GANomaly employs a combination of a generator and dis-
207 criminator within a GAN framework to learn the underlying data distribution and iden-
208 tify anomalies. The method shows strong performance in various computer vision
209 tasks, such as image-based anomaly detection, by effectively learning to differentiate

210 between normal and abnormal data patterns. BeatGAN [44] also uses a generative
211 model for anomaly detection. Like GANomaly, the idea is to learn the distribution of
212 normal data and detect anomalies as hard-to-reconstruct data samples. However, unlike
213 our model, they model and reconstruct the whole signal whereas we learn to forecast
214 the continuation of a given signal segment. We believe this proxy task of temporal
215 continuation leads to a better model of the underlying features of the data.

216 Like in our work, Tang et al. [45] consider linescan image data. However, they
217 do not treat the data as temporal, instead working with fixed size images in which the
218 temporal dimension is a second spatial dimension. They tackle the supervised anomaly
219 detection problem for industrial inspection using a skip autoencoder and deep feature
220 extractor. The skip autoencoder captures multi-scale features by incorporating skip
221 connections, while the deep feature extractor enhances the representation of the input
222 data. This combination significantly improves the accuracy of anomaly detection in
223 industrial settings, demonstrating robustness in identifying defects in complex envi-
224 ronments.

225 Specifically related to anomaly detection in rail images, Liu et al. [46] present a ma-
226 chine vision-based method for inspecting rail fastener defects across multiple railways.
227 The approach utilizes image processing and deep learning techniques to automatically
228 detect and classify defects in rail fasteners, ensuring the safety and reliability of railway
229 infrastructure. The proposed method achieves high precision and efficiency, making it
230 suitable for large-scale railway maintenance applications.

231 Modern approaches to time series anomaly detection were recently surveyed by
232 Zamanzadeh et al. [47]. We conclude the literature review by mentioning the most re-
233 cent and relevant methods. Wang et al. [48] propose a VAE that is conditioned on both
234 global and local frequency features. This improves reconstruction normal data signif-
235 icantly. Kang and Kang [49] use a transformer to model temporal dependencies and
236 relationships among variables via self attention across these two dimensions. Miao et
237 al. [50] combine GAN losses with a transformer-based autoencoder while incorporat-
238 ing a contrastive loss into the discriminator which helps improve generalisation of the
239 normal model. CARLA [51] also uses a contrastive loss but proposes to inject anoma-
240 lies to create negative samples for contrastive learning. Kim et al. [52] consider the

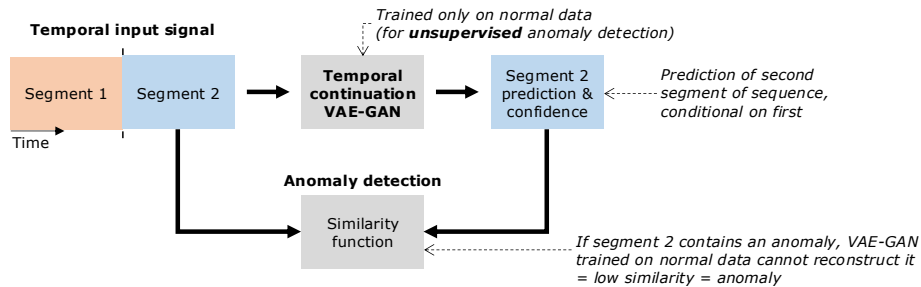


Figure 1: Overview of unsupervised anomaly detection method using a temporal continuation VAE-GAN. An input signal with a temporal dimension is divided into two segments. The temporal continuation VAE-GAN predicts the second segment, conditional on the first. This VAE-GAN is trained to learn the space of normal signals, including the subspace of plausible continuations and a pointwise confidence estimate to distinguish structural uncertainties (which we expect the model to be able to capture) from stochastic uncertainties (which we do not). If the second segment contains an anomaly, we do not expect the VAE-GAN to be able to accurately reconstruct it and this dissimilarity should be measurable and indicative of an anomaly. Since the VAE-GAN only needs to see normal data, this provides a means to perform unsupervised anomaly detection.

241 problem of test-time adaptation when a learnt normal model must deal with distribu-
 242 tional shift at test-time. Other generative architectures have also been considered. Zhou
 243 et al. [53] use normalising flows as a generative model for both anomaly detection and
 244 localisation. Yao et al. [54] use a diffusion model to remove anomalies. However, they
 245 propose to adapt the level of noise such that it is appropriate to the scale of anomaly.
 246 Dai et al. [55] also use a diffusion model but for generating synthetic anomalies without
 247 a prior. Finally, Liu et al. [56] tackle the problem of unsupervised anomaly detection
 248 in the context of continual learning. Here, the task is to incrementally learn different
 249 anomalies without forgetting those learnt earlier.

250 3. Method

251 Our goal is to learn the space of normal variation of a temporal signal. We pose this
 252 in terms of estimating the temporal continuation of a given signal segment. However,
 253 rather than estimate a single point estimate, we predict the subspace of possible contin-
 254 uations. This provides a route to unsupervised anomaly detection since we can measure

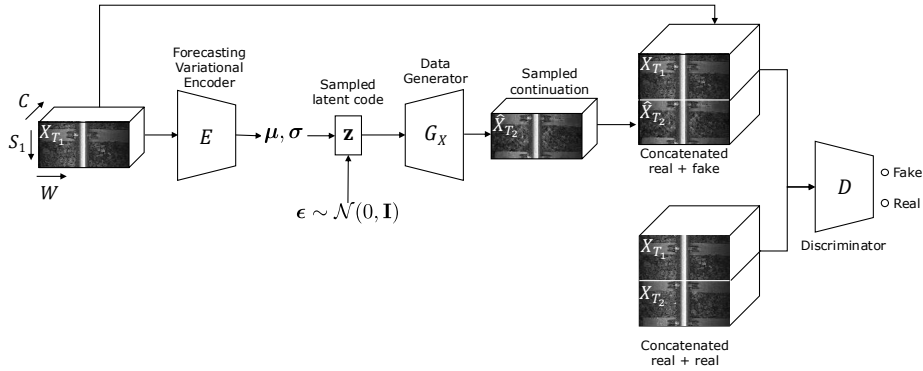


Figure 2: The temporal continuation VAE-GAN architecture. From the observed part of the time series X_{T_1} , the forecasting variational encoder E computes the parameters of a latent distribution, μ and σ . The data generator, G_X , decodes a sample from this latent distribution, \mathbf{z} , into a prediction of the following time-steps \hat{X}_{T_2} . The discriminator D is given real, $\text{cat}(X_{T_1}, X_{T_2})$, or fake, $\text{cat}(X_{T_1}, \hat{X}_{T_2})$, concatenated time series and seeks to distinguish them.

255 how the true continuation differs from the predicted subspace (see the overview in Fig-
 256 ure 1). Our underlying model is a temporal continuation VAE-GAN (see Section 3.2).
 257 This model learns to encode a given segment of signal to the subspace of possible con-
 258 tinuations, represented as a mean and variance of a latent representation. Sampling
 259 from this distribution and decoding provides a possible continuation. Our model also
 260 learns to predict a pointwise confidence map so that it learns in an unsupervised man-
 261 ner which regions of the continuation are predicted with high confidence (see Section
 262 3.3). It is in these regions that we expect to be able to reliably detect anomalies. The
 263 confidence map represents the predicted pointwise confidence of the continuation *after*
 264 *the optimal latent representation has been found*. This optimal representation is found
 265 in practice via a process of analysis-by-synthesis to fit the model (see Section 3.4). Our
 266 model is trained with several losses described in Section 3.6. Specifically, the objective
 267 is that the predicted subspace contains the true continuation observed in the training
 268 data and that the latent space is well-behaved (achieved using conventional VAE-GAN
 269 losses) but also that the subspace of continuations is diverse and not overfitted to the
 270 particular observed continuations.

271 Intuitively, our model allows us to answer the question: “Given the first part of a

272 temporal sequence, what possible continuations do we expect to see?” Then, given an
 273 actual continuation, we can ask: “How far does the actual continuation lie from the
 274 subspace of possible continuations that the model predicted?” Finally, our confidence
 275 map allows us to ask: “How confident is the model in its prediction at each output
 276 point?” Together, the answers to the second and third question allows us to detect
 277 anomalies when we see a features in a continuation that our model cannot predict yet
 278 our model is confident in the prediction of those features.

279 3.1. Problem statement

280 Consider a signal with zero or more spatial dimensions, one or more channels and
 281 a temporal dimension that is observed at S_1 evenly spaced time steps. We represent
 282 this observation by the tensor $X_{T_1} \in \mathbb{R}^{W \times C \times S_1}$, where C is the number of channels and
 283 the spatial dimension W may be expanded or dropped as appropriate to the particular
 284 signal. We are interested in the task of predicting the signal at the following S_2 time
 285 steps, i.e. predicting the tensor $X_{T_2} \in \mathbb{R}^{W \times C \times S_2}$ given X_{T_1} . Hat denotes an estimated
 286 quantity, e.g. \hat{X}_{T_2} is the prediction of the true X_{T_2} .

287 3.2. Temporal Continuation VAE-GAN

288 The first component of our model is a VAE-GAN, as shown in Figure 2. However,
 289 unlike a conventional VAE-GAN, we do not seek to autoencode, i.e. to reconstruct
 290 samples similar to the input. Instead, we decode to a continuation of the temporal se-
 291 quence. Therefore, the job of the encoder is to find latent distribution parameters that
 292 model the space of possible continuations. We do not use a ‘content’ (or ‘data’) loss
 293 that directly penalises differences between X_{T_2} and \hat{X}_{T_2} as in a VAE or autoencoder.
 294 Instead, we require only that the continuation is plausible (as measured by the discrim-
 295 inator) as in a GAN. The discriminator sees the concatenation of the observed part of
 296 the sequence and its predicted continuation and can therefore judge whether the contin-
 297 uation is plausible given the observation. The VAE-GAN part of our model comprises
 298 the following components.

299 *Forecasting Variational Encoder.* The forecasting variational encoder is a pair of func-
 300 tions $\mu, \sigma : \mathbb{R}^{C \times S_1 \times W} \rightarrow \mathbb{R}^d$ such that $\mu(X_{T_1}), \sigma(X_{T_1})$ provides the parameters of the

301 normal distribution corresponding to the embedding of X_{T_1} into a d -dimensional space.
 302 The mean, $\boldsymbol{\mu}(X_{T_1})$, of this distribution encodes the predictable aspects of X_{T_2} , while
 303 $\boldsymbol{\sigma}(X_{T_1})$ describes the shape of the distribution characterising the uncertain aspects.

304 *Data generator.* Unlike in a conventional VAE or VAE-GAN, our generator (or de-
 305 coder) does not seek to reconstruct the original input data. Instead, it predicts the
 306 temporal continuation of the input data. We call this our data generator to distinguish
 307 it from the confidence generator later. The data generator is a function $G_X : \mathbb{R}^d \rightarrow$
 308 $\mathbb{R}^{C \times S_2 \times W}$ such that $\hat{X}_{T_2} = G_X(\mathbf{z})$ is a prediction of X_{T_2} conditioned on latent vector
 309 $\mathbf{z}(X_{T_1}, \boldsymbol{\epsilon}) = \boldsymbol{\mu}(X_{T_1}) + \boldsymbol{\epsilon} \odot \boldsymbol{\sigma}(X_{T_1})$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ is random noise drawn from a
 310 normal distribution. The idea is that $\boldsymbol{\mu}$ should encode the predictable aspects of X_{T_2}
 311 while $\boldsymbol{\epsilon}$ provides a space in which to explore the structurally or stochastically uncertain
 312 aspects.

313 *Discriminator.* The discriminator is a function $D : \mathbb{R}^{C \times S_1 + S_2 \times W} \rightarrow [0, 1]$ that is given
 314 a concatenation of the observed X_{T_1} and either the true (X_{T_2}) or predicted (\hat{X}_{T_2}) contin-
 315 uation and returns the probability that the concatenated observation is drawn from the
 316 true data distribution. i.e. $D(\text{cat}(X_{T_1}, X_{T_2}))$ aims to predict whether $\text{cat}(X_{T_1}, X_{T_2})$ is real
 317 or fake, where cat concatenates tensors along the temporal dimension.

318 3.3. Confidence prediction and model fitting

319 We further augment our Temporal Continuation VAE-GAN with a means to predict
 320 confidence in the continuation for each spatiotemporal location. This is important for
 321 distinguishing between anomalous deviations from normality and stochastic variations
 322 that we do not expect the model to be able to reconstruct. We supervise the confidence
 323 prediction based on the actual error between the true continuation and the *best possible*
 324 *fit* of the model to the continuation. Concretely, when given access to the true X_{T_2} , we
 325 optimise the noise $\boldsymbol{\epsilon}$ to minimise the error between X_{T_2} and \hat{X}_{T_2} . The remaining error
 326 represents the inability of the model to explain all of X_{T_2} and we use this to supervise
 327 the confidence map.

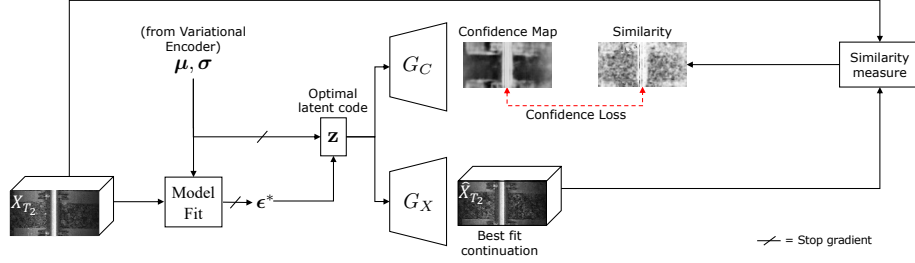


Figure 3: The confidence generator, G_C , predicts a confidence map from an optimal latent code. This should correspond to the spatial similarity between X_{T_2} and \hat{X}_{T_2} when the optimal ϵ^* has been found via a model fitting procedure (through which we do not propagate gradients) that minimises the reconstruction error.

328 *Confidence generator.* The confidence generator is a function $G_C : \mathbb{R}^d \rightarrow [0, 1]^{S_2 \times W}$
 329 such that $\mathbf{C}_m = G_C(\mathbf{z}(X_{T_1}, \epsilon^*))$ is a single channel confidence map of the same spa-
 330 tiotemporal dimension as X_{T_2} . Entries in \mathbf{C}_m represent the confidence (a value in the
 331 range $0 \dots 1$) of the prediction of X_{T_2} at the corresponding spatiotemporal location. The
 332 intention is that this confidence value reflects the similarity between X_{T_2} and \hat{X}_{T_2} when
 333 the latent vector with optimal ϵ is passed to G_X (see Model Fitting below). The defini-
 334 tion of similarity depends upon the choice of similarity measure used in the confidence
 335 loss.

336 3.4. Model fitting

337 Suppose we are given both an observed X_{T_1} and the true continuation X_{T_2} . We
 338 want to find the best representation within our model of this observation, i.e. to fit the
 339 model. This entails finding the optimal ϵ^* such that $G_X(\mathbf{z}(X_{T_1}, \epsilon^*))$ best fits the true
 340 continuation X_{T_2} . We solve the analysis-by-synthesis optimisation problem:

$$\epsilon^* = \arg \min_{\epsilon} \|G_X(\mu(X_{T_1}) + \epsilon \odot \sigma(X_{T_1})) - X_{T_2}\|_1. \quad (1)$$

341 This seeks to minimise the L_1 difference between actual and synthesised X_{T_2} . We use
 342 this optimal model fit to compute the similarities that are used to train the confidence
 343 generator. Specifically, we solve the optimisation problem using gradient descent for a
 344 fixed number of iterations within the outer training loop.

345 *3.5. Training the confidence generator*

346 The confidence generator is trained using model fitting as shown in Figure 3. The
347 model fitting process is used as an oracle that provides the optimal latent code cor-
348 responding to the best fit continuation. The difference between the best fit and true
349 continuations is determined using a data-specific similarity measure. The confidence
350 generator is supervised to predict confidence maps that are close to the true similarity.
351 We do not propagate gradients from the confidence generator or through the model fit-
352 ting optimisation process into the variational encoder. So the confidence generator can
353 either be trained independently of the temporal continuation VAE-GAN or in parallel
354 with it.

355 *3.6. Losses*

356 The goal of our Temporal Continuation VAE-GAN is to learn the space of plausible
357 continuations, conditioned on the observation X_{T_1} . Training only with a reconstruction
358 or content loss as in a VAE encourages overfitting and collapse of the latent space to
359 predict only the true continuation without any diversity. Instead, we use a discriminator
360 and GAN loss to ensure that all continuations are plausible and a diversity loss to ensure
361 the latent distributions capture meaningful and significant variation. Our assumption is
362 that, if both of these are satisfied, then the true continuation lies somewhere in the latent
363 space. In addition, as in a VAE we impose a prior regularisation loss on the predicted
364 distributions using the KL divergence. This encourages a well-behaved latent space in
365 which the model fitting optimisation can smoothly converge to a good solution. We
366 now describe the various losses used during training.

367 *Generator loss.* We use a binary cross entropy loss for the discriminator. Although
368 other GAN losses could be used (we experimented with the WGAN) we found this
369 simple loss to work well for our applications. To update the generator, we compute this
370 with inverted labels, i.e. seek to maximise the probability of being real for a batch of N
371 fake images:

$$\mathcal{L}_{\text{gen}} = - \sum_{i=1}^N \log D(\text{cat}[X_{T_1}^i, G_X(\mathbf{z}(X_{T_1}^i, \epsilon))]). \quad (2)$$

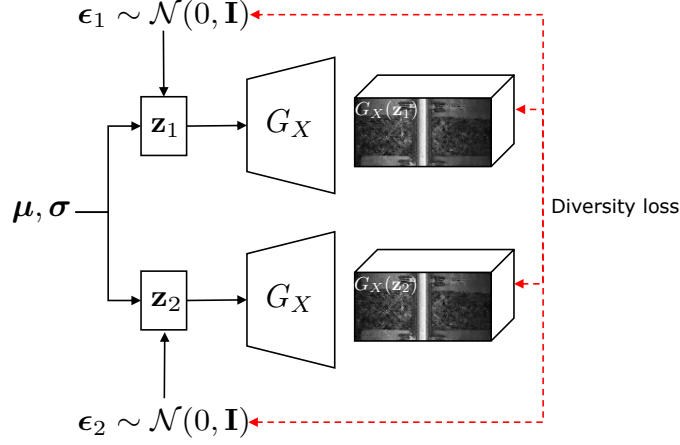


Figure 4: Each time the forecasting variational encoder estimates the latent distribution, we draw two different samples. The diversity loss encourages that, when the samples are further apart, so should the decoded continuations be further apart.

372 *Discriminator loss.* To update the discriminator, we compute binary cross entropy loss
 373 for a batch of correctly labelled fake and real images:

$$\mathcal{L}_{\text{dis}} = - \sum_{i=1}^N \log D(\text{cat}[X_{T_1}^i, X_{T_2}^i]) + \log(1 - D(\text{cat}[X_{T_1}^i, G_X(\mathbf{z}(X_{T_1}^i, \epsilon))])). \quad (3)$$

374 *Prior loss.* We use the KL divergence as a prior loss to encourage the latent distribution
 375 for every input to be close to a standard normal distribution:

$$\mathcal{L}_{\text{KL}} = \sum_{i=1}^N \sum_{j=1}^d \mu_j(X_{T_1}^i)^2 + \sigma_j(X_{T_1}^i)^2 - \log \sigma_j(X_{T_1}^i) - 1 \quad (4)$$

376 *Diversity loss.* This loss encourages diversity within the latent space, i.e. that a large
 377 change in \mathbf{z} should correspond to a large change in \hat{X}_{T_2} . We ensure this using the
 378 diversity loss proposed by [57]:

$$\mathcal{L}_{\text{diversity}} = \sum_{i=1}^N \frac{\|\mathbf{z}(X_{T_1}^i, \epsilon_1) - \mathbf{z}(X_{T_1}^i, \epsilon_2)\|_1}{\|G_X(\mathbf{z}(X_{T_1}^i, \epsilon_1)) - G_X(\mathbf{z}(X_{T_1}^i, \epsilon_2))\|_1 + \epsilon}, \quad (5)$$

379 where ϵ is a small constant for numerical stability, and ϵ_1 and ϵ_2 are two random sam-
 380 ples. See Figure 4. This loss was previously used in the context of GANs and has not
 381 been used in the context of VAE-GANs or temporal continuation previously.

Algorithm 1 Training our proposed model

- 1: **while** not converged **do**
- 2: Sample random batch of real images $(X_{T_1}^1, \dots, X_{T_1}^B)$
- 3: ***# Phase 1: encourage G_X and E to produce***
- 4: ***# images that are more realistic and diverse***
- 5: Generate two continuations for each real image:

$$G_X(\mathbf{z}(X_{T_1}^i, \epsilon_1)) \text{ and } G_X(\mathbf{z}(X_{T_1}^i, \epsilon_2))$$

- 6: Compute \mathcal{L}_{gen} , $\mathcal{L}_{\text{diversity}}$ and \mathcal{L}_{KL} and backpropagate into G_X and G_E
 - 7: ***# Phase 2: encourage G_C to predict confidence***
 - 8: ***# consistent with similarity using optimal \mathbf{z}***
 - 9: Fit the model to the current target images by solving (1)
 - 10: Compute $\mathcal{L}_{\text{confidence}}$ and backpropagate into G_C
 - 11: ***# Phase 3: improve discriminator D***
 - 12: ***# to better detect fake images***
 - 13: Compute \mathcal{L}_{dis} and backpropagate into discriminator D
 - 14: Take gradient descent step
 - 15: Zero gradients
 - 16: **end while**
-

382 *Confidence loss.* The confidence loss measures the L_1 error between the confidence
383 map predicted by G_C and the true similarity map $\mathbf{S} \in [0, 1]^{S_2 \times W}$:

$$\mathcal{L}_{\text{confidence}} = \sum_{i=1}^N \|G_C(\mathbf{z}(X_{T_1}^i, \epsilon^*) - \mathbf{S}^i\|_1 \quad (6)$$

384 where $\mathbf{S}^i = s(G_X(\mathbf{z}(X_{T_1}^i, \epsilon^*), X_{T_2}^i))$ is computed according to some similarity function
385 $s : \mathbb{R}^{S_2 \times W} \times \mathbb{R}^{S_2 \times W} \rightarrow [0, 1]^{S_2 \times W}$. There are many ways we might choose to define
386 similarity depending on the nature of the data. We specify what was used for each
387 dataset below.

388 *3.7. Implementation*

389 Each iteration of our training pipeline comprises three phases, as shown in Algo-
390 rithm 1. In each phase, losses that relate to different components of our model are
391 calculated and backpropagated before a gradient descent step is taken on the accumu-
392 lated gradients. We use the RMSProp optimiser. We implement our generators and
393 discriminator as convolutional neural networks, though this choice is orthogonal to our
394 overall idea and any architecture (such as a transformer) could be used. We follow the
395 DCGAN architecture for each component, adapting filter sizes to accommodate spatial
396 input size of $S_1 \times W$ for the encoder, $S_2 \times W$ for the generator and $S_1 + S_2 \times W$ for the
397 discriminator. Our generators use batchnorm and ReLU activation with tanh activation
398 at the output layer while our discriminator uses batchnorm, LeakyReLU activation and
399 sigmoid activation for the output.

400 *3.8. Unsupervised anomaly detection*

401 Assuming that our model has been trained only on normal data (i.e. excluding
402 anomalies) then, given real observation X_{T_1} and its true continuation X_{T_2} , we can
403 use our model to assess whether X_{T_2} contains any anomalies. The difference between
404 $\hat{X}_{T_2} = G_X(\mathbf{z}(X_{T_1}), \epsilon^*)$ and X_{T_2} indicates which parts of the real continuation were diffi-
405 cult for the model to reconstruct. However, we know that our prediction will only be
406 reliable in non-stochastic regions of the continuation, i.e. where the model is confident.
407 We can therefore produce an anomaly map, \mathbf{A} , that scales errors by their corresponding
408 confidence:

$$\mathbf{A} = G_C(\mathbf{z}(X_{T_1}, \epsilon^*)) \odot e(G_X(\mathbf{z}(X_{T_1}, \epsilon^*)), X_{T_2}), \quad (7)$$

409 where $e(\cdot, \cdot)$ is a data-specific error function. Large values in this map, indicate regions
410 where the model is confident in its prediction but the prediction is very different to the
411 data - i.e. an anomaly. For anomaly detection, we threshold these anomaly maps, count
412 the number of anomalous-labelled points and then threshold the count to classify data
413 as anomalous.

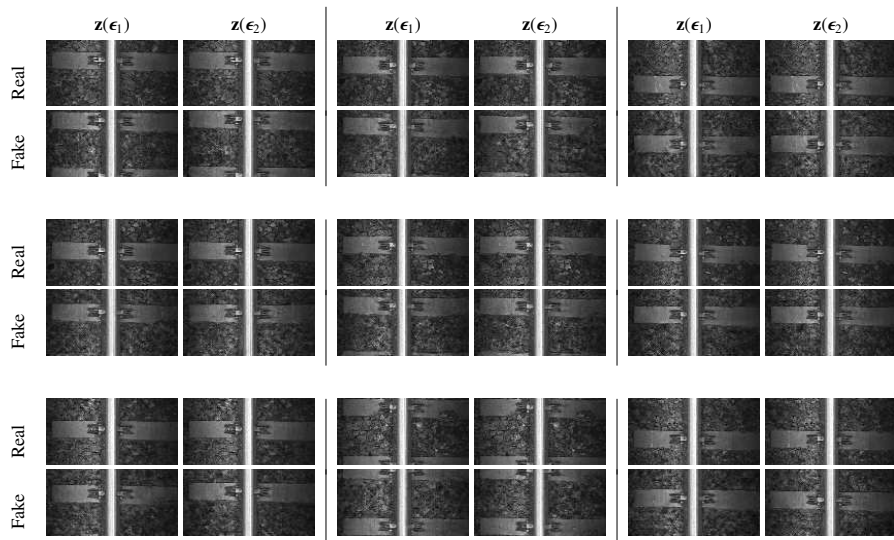


Figure 5: Illustration of quality of temporal continuation and diversity. Each 2×2 block of images shows the same X_{T_1} (observed real image) in the top row and two different \hat{X}_{T_2} (fake images) in the bottom row, produced by two different random samples from the latent space. Both should provide plausible continuations of the real image while also showing diversity between the two samples either in stochastic elements (such as the ballast in the background) or structural elements (such as the precise positioning of the sleepers or clips).

414 4. Results

415 4.1. Datasets

416 We provide experimental results on three different datasets across two modalities to
 417 demonstrate the performance of our method. Testing our approach on other modalities
 418 of data such as audio or time series from a source other than ECG is left to future work.

419 We use a dataset of grayscale ($C = 1$) railtrack images in order to qualitatively
 420 evaluate the behaviour of our model. This is captured with a linescan camera mounted
 421 on the underside of a track inspection car, the vision system illuminates the track with
 422 a series of LED wire lights and gets images of the track and its surroundings as the car
 423 moves along it at speeds of up to 125mph. The resolution of the original images is $W =$
 424 $2,048$, $H = 15,000$ where H corresponds to the temporal dimension. We take crops of
 425 size 2048×2048 via sliding a window vertically along the images with a step size of

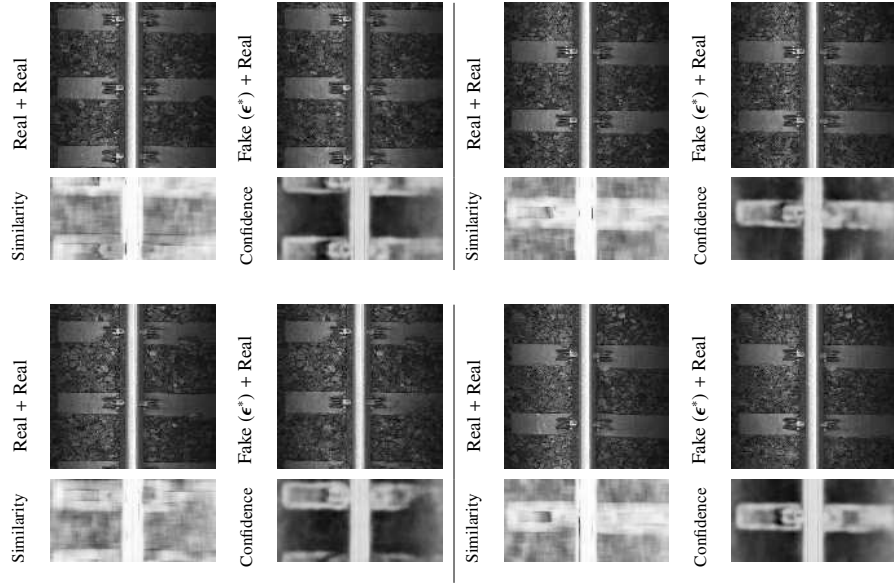


Figure 6: Model fitting and confidence prediction. For each example (comprising two rows and two columns), the first row of the first column shows an observed image X_{T_1} and its true continuation X_{T_2} . The first row of the second column shows the observed image X_{T_1} and its predicted continuation \hat{X}_{T_2} using the optimal ϵ^* after fitting the model to the observed X_{T_2} . The similarity between X_{T_2} and \hat{X}_{T_2} (according to the structural similarity index) is shown in the second row of the first column while the estimated confidence, having seen only X_{T_1} is shown in the second row of the second column.

426 100 pixels. We then downsample the images to size 128×128 for $W = 128$ and split
427 equally into size $S_1 = S_2 = 64$. From 20 linescan images, this leads to a dataset of
428 10k 128×128 images. It is assumed that there are no anomalies within this training set
429 and we use no labels. For this dataset, we use the structural similarity [58] to supervise
430 confidence maps, i.e. $s(x, y) = \text{SSIM}(x, y)$. This measures similarity over a local region
431 at each point, rather than only pixel-wise similarity. This is helpful in reflecting low
432 confidence in stochastic regions. For anomaly detection we use negated similarity as
433 our error measure, i.e. $e(x, y) = 1 - s(x, y)$.

434 We provide quantitative evaluation on the ECG5000 benchmark. This is a time
435 series anomaly detection benchmark. It forms part of the UCR time series archive [1]
436 and comprises 5,000 electrocardiogram (ECG) single channel ($C = 1$) traces from the
437 dataset originally collected by [2]. Each trace consists of a total of $W = 140$ uniform

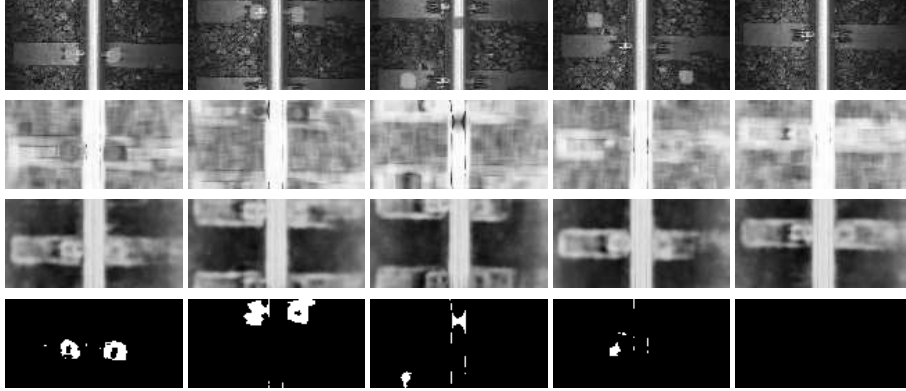


Figure 7: Anomaly detection on synthetic examples. From top to bottom: real image with synthetic anomaly, similarity $s(G_X(\mathbf{z}(X_{T_1}, \epsilon^*)), X_{T_2})$, confidence map and thresholded error map. The first three examples show anomalies on the clips, sleeper and rail, the fourth shows an anomaly on the ballast and the fifth no anomaly.

438 time steps corresponding to one heartbeat from a patient with congestive heart failure.
 439 We use $S_1 = 76$ time steps for the observed portion and $S_2 = 64$ time steps for the
 440 predicted portion. We supervise the confidence generator with L1 error, i.e. $s(x, y) =$
 441 $\text{abs}(x - y)$, hence our confidence generator is actually predicting error. This means that
 442 when we perform anomaly detection we can directly use the scaled “confidence” value
 443 as error: $e(x, y) = w \cdot s(x, y)$, where w is a scalar weight parameter.

444 We also use the MIT-BIH Arrhythmia Database [3] for quantitative evaluation. This
 445 is a widely used reference dataset for ECG signal analysis. This data set contains
 446 multichannel ECG recordings with detailed annotations for each heartbeat, identifying
 447 beat types and rhythm information. Each recording is sampled at 360 Hz and represents
 448 a complete ECG trace of a patient. We use timesteps $S_1 = S_2 = 64$ and again use L1
 449 error to supervise the correspondence generator.

450 4.2. Qualitative analysis

451 We begin by providing qualitative analysis of the behaviour of our model on the
 452 railway dataset.

453 In Figure 5 we show the ability of the model to generate plausible and diverse
 454 samples. For each example, we take a single real $X_{T_1}^i$ and generate fake continuations
 455 $G_X(\mathbf{z}(X_{T_1}^i, \epsilon_1))$ and $G_X(\mathbf{z}(X_{T_1}^i, \epsilon_2))$ where ϵ_1 and ϵ_2 are two different random samples.

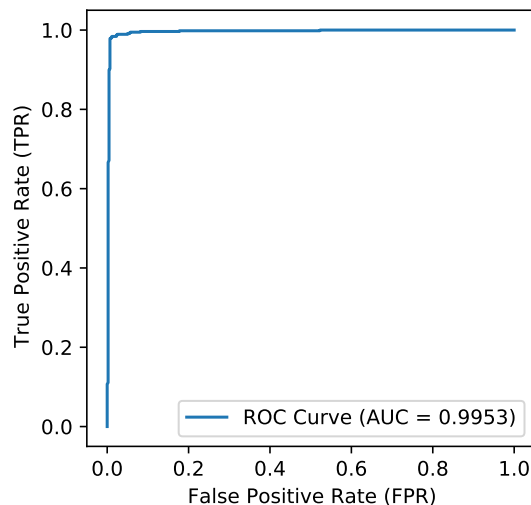
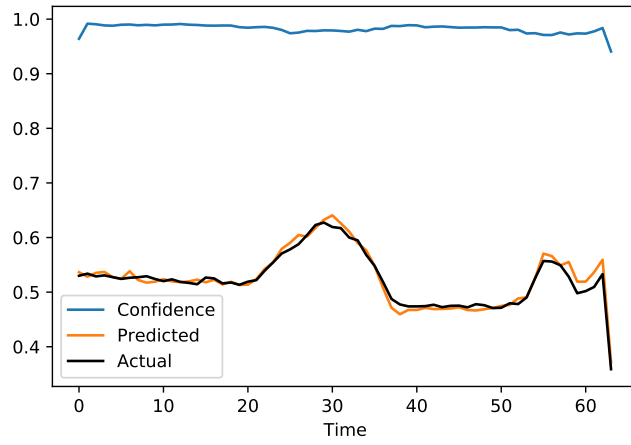


Figure 8: ROC curve for the ECG5000 dataset [1, 2].

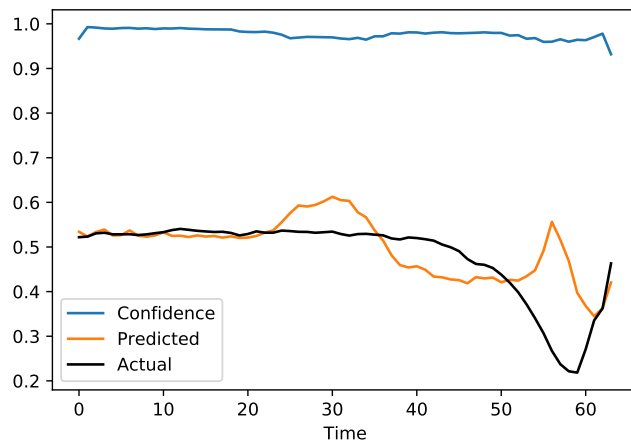
456 The generator is able to create images with the right structure (e.g. spacing between
 457 sleepers) and detail while using different random samples leads to slight changes in
 458 stochastic and structural elements. This illustrates that our model not only learns to
 459 decode the latent space to a plausible continuation but also that it learns a subspace of
 460 variation for the possible continuations.

461 In Figure 6 we illustrate fitting our model to observed data. Given a real observed
 462 X_{T_1} , we optimise ϵ in order to minimise the error to the real observed X_{T_2} by solving
 463 the optimisation problem in (1). Note that this successfully adjusts structural elements
 464 of the fake image such that the main features align well. The similarity maps show
 465 which regions are reconstructed accurately (white means perfect local similarity). The
 466 predicted confidence map shows the model prediction of which regions in the image
 467 the model will be able to generalise to well. This includes the rail itself, the sleeper and
 468 certain elements of the clamp while it has low confidence for the stochastic background
 469 as expected. We emphasise that this separation of learnable structural uncertainty from
 470 unlearnable stochastic uncertainty is learnt without supervision. The structural ele-
 471 ments are effectively ‘detected’ by the fact that they can be reliably modelled.

472 To qualitatively evaluate anomaly detection, we manually painted anomalies onto



(a)



(b)

Figure 9: ECG traces for a normal (a) and abnormal (b) heartbeat.

473 the rail, sleeper, clamp and background ballast. In Figure 7 we show qualitative exam-
 474 ples of anomaly detection on these images. In the top row, our synthetic anomalies are
 475 visible as gray blobs. In the second row, the raw similarity between the reconstructed
 476 and observed images does show low similarity in the anomaly regions but also in the
 477 stochastic parts of the image. In the third row, the confidence map predicted by our
 478 model allows suppression of dissimilarity in regions of low confidence. The resulting

Source	[S]upervised/ [U]nsupervised	AUC	Acc	F1
Ours	U	0.9953	0.9860	0.9875
Pereira and Silveira [38]	S	0.9836	0.9843	0.9844
	U	0.9819	0.9596	0.9522
Lei et al. [59]	S	0.9100	-	-
Karim et al. [60]	S	-	0.9496	-
Malhotra et al. [61]	S	-	0.9340	-
Liu et al. [62]	U	-	-	0.8084

Table 1: Quantitative anomaly detection results on the ECG5000 dataset.

479 anomaly maps in the bottom row detect only badly reconstructed regions in areas of
480 high confidence. This is crucial to limit false positives.

481 4.3. Quantitative evaluation

482 Although the ECG5000 dataset was originally used for five-way classification (nor-
483 mal plus four abnormalities), this dataset is now widely used for time series anomaly
484 evaluation (normal versus any abnormality). We follow [38] and divide the dataset
485 randomly into 80% training and 20% testing. We train the model using only the nor-
486 mal portion of the training set (i.e. excluding anomalies), comprising 2,359 traces in
487 total. We then test on the whole of the test set which comprises 1,000 traces in to-
488 tal, 560 of which are normal and 440 of which are anomalous. Note that we operate
489 in an unsupervised setting: we never see abnormal traces at training time. We show
490 our ROC curve in Figure 8 and quantitative results in Table 1. Our approach outper-
491 forms all previous unsupervised methods and even outperforms the best supervised
492 method on both area under curve and accuracy. In Figure 9 we show qualitative results
493 for a normal (a) and anomalous (b) trace. We plot the true X_{T_2} (black), best-fit con-
494 tinuation $\hat{X}_{T_1} = G_X(\mathbf{z}(X_{T_1}, \epsilon^*))$ (orange) and predicted error (i.e. one minus predicted
495 confidence). The model can fit the normal trace well but cannot explain the anoma-
496 lous trace. In other words, conditional on the first observed segment, the anomalous

Model	[S]upervised/ [U]nsupervised	AUC	F1 (%)	Accuracy (%)	Recall (%)	Precision (%)
Stacked LSTM [63]	U		81.0	-	87.0	82.0
LSTM with MLP [64]	S		87.0	95.0	75.0	-
VAE [65]	U		76.6	87.8	-	-
Transformer [66]	U	0.93	92.3	89.5	98.2	87.1
Our work	U	0.93	93.2	90.1	95.1	91.4

Table 2: Quantitative anomaly detection results on the MIT-BIH dataset.

second segment does not lie within the subspace forecast by our model. The predicted error is relatively flat though increases sharply towards the end where there is often a lot of variability in the training data. The fact that our model knows its prediction in this region is unreliable means differences here can be ignored - i.e. they cannot be confidently labelled as anomalies.

For the MIT-BIH dataset [3] we again divide into training and testing sets, where the training set consists only of normal beats, and the testing set included both normal and abnormal beats. We follow standard practices for the preprocessing and error evaluation for this dataset [66]. Segmentation of the continuous traces into training and testing samples was based on the annotated R-peak positions, with each heartbeat segment spanning from one R-peak to the next. To ensure consistent segment length, signals from both channels were resampled to a fixed length of 128 time steps per segment. Additionally, to reduce the impact of amplitude variations, the signal data from all channels were normalized using the 3rd and 97th percentiles as the range for scaling. Specifically, spectral error was employed as the core metric to measure the difference between predicted and actual signals. The spectral error was calculated by performing a Fast Fourier Transform (FFT) on the residuals of the first channel and applying our confidence-weighted scaling. As shown in Table 2, our method achieves good performance in terms of F1 score, recall, and precision, particularly in the unsupervised anomaly detection task. Compared to existing unsupervised methods and even some supervised approaches, our method set a new benchmark for the MIT-BIH dataset.

Dataset	Condition	F1 (%)	AUC (%)	Accuracy (%)
ECG5000	Baseline	98.3	99.56	98.1
	No confidence map	97.9	99.56	97.7
	No optimization of ϵ	98.3	99.4	98.1
	No diversity loss	97.6	99.5	97.3
MIT-BIH	Baseline	93.2	92.6	90.1
	No confidence map	93.2	92.6	90.1
	No optimization of ϵ	89.9	86.9	84.8
	No diversity loss	89.6	84.3	86.1

Table 3: Ablation experiments Results on ECG5000 and MIT-BIH Datasets.

519 *4.4. Ablation study*

520 Finally, we conducted an ablation study of the three key ingredients of our method
521 using the ECG5000 and MIT-BIH datasets. Specifically, we evaluate the impact of:
522 confidence map, model fitting (i.e. optimization of ϵ), and diversity loss. The results,
523 summarized in Table 3, highlight the varying impact of these components on model per-
524 formance. In the ECG5000 data set, incorporating the confidence map led to notable
525 improvements in the F1 score and accuracy, while the AUC remained consistent. How-
526 ever, on the MIT-BIH dataset, the inclusion of the confidence map did not show any
527 measurable effect, as all metrics remained identical with or without it. This suggests
528 that the utility of the confidence map may vary depending on the dataset characteristics
529 or noise levels. The optimization of the latent variable ϵ exhibited a more consistent
530 influence. On ECG5000, it slightly enhanced AUC while maintaining F1 and accu-
531 racy scores. On MIT-BIH, the impact was more pronounced, with F1 score, AUC and
532 accuracy all improving substantially. These results indicate that ϵ -optimization signif-
533 icantly contributes to the model’s ability to generalize, particularly on datasets with
534 diverse and complex patterns like MIT-BIH. The diversity loss consistently improved
535 model performance on both datasets. This demonstrates the robustness of diversity
536 loss in improving the detection of anomalous samples and reducing overfitting across

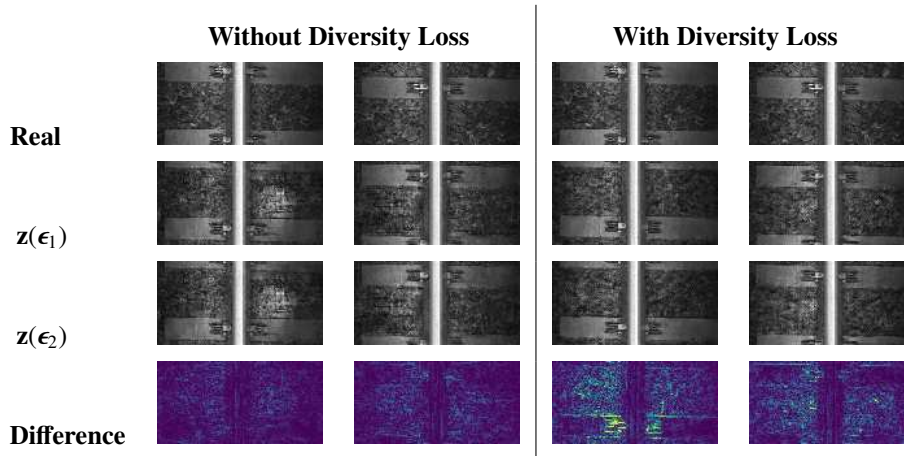


Figure 10: Impact of diversity loss: In the first row we show an observed real X_{T_1} . In the second and third rows we show two possible continuations in which we use different random noise samples ϵ_1 and ϵ_2 . In the fourth row we show a heat map visualisation of the absolute difference between the second and third rows (same colour map scale used for all four images). In the first two columns the results are an ablation in which diversity loss is not used during training. In the last two columns diversity loss is used during training.

537 datasets.

538 In summary, while the impact of the confidence map appears dataset-dependent,
 539 the optimization of ϵ and the diversity loss consistently enhance model performance. ϵ -
 540 optimization is particularly effective in improving generalization on complex datasets,
 541 and the diversity loss contributes significantly to anomaly detection and classification
 542 robustness. These findings validate the necessity of these components in achieving
 543 state-of-the-art results on anomaly detection tasks.

544 Finally, in Figure 10 we show a qualitative illustration of the impact of the diversity
 545 loss on the railway image dataset. Without diversity loss we can see that the model
 546 has learnt limited dependence on ϵ (the two generated images are very similar and the
 547 difference map shows little change anywhere). With the diversity loss, more variation
 548 for different ϵ is evident and this is visible in the difference map. This shows both
 549 structural variation (around the clamps) and stochastic variation in the ballast. Interest-
 550 ingly, we also observe that diversity loss improves the model more generally. Without
 551 diversity loss the generations exhibit artefacts which are not present with it. This is

552 because diversity loss avoids overfitting to the particular continuation observed in the
553 training data and encourages learning a smooth subspace of plausible continuations.

554 **5. Conclusions**

555 We have shown that we can learn a generative model for stochastic continuation of
556 non-overlapping temporal sequences. Our unsupervised method automatically learns
557 which parts of the continuation are predictable and which are stochastic. This provides
558 a route to unsupervised anomaly detection.

559 From a practical perspective, deploying our approach in real-world settings for a
560 new data domain entails only the following steps. First, choose an appropriate architec-
561 ture for the encoder and data/confidence generators. Any off-the-shelf architecture that
562 is widely used for the data modality could be used here. Second, choose an appropriate
563 similarity or dissimilarity measure to supervise the confidence generator. Again, any
564 standard metric such as SSIM for images or MSE for time series signals could be used.
565 Finally, adjust the key hyperparameters of latent space dimension and segment sizes
566 (often $S_1 = S_2$ will prove the best choice, giving an equal balance between the size
567 of input data and model prediction). In terms of computational cost, in the simplest
568 case our method only requires a forward pass through the encoder and generator net-
569 works and evaluation of the anomaly map metric. For slightly improved performance,
570 optional iterative optimisation of *epsilon* requires a fixed number of gradient descent
571 steps.

572 There are a number of limitations to our approach. First, from a practical imple-
573 mentation perspective, our use of a convolutional encoder and decoder potentially lim-
574 its modelling of long-range dependencies. This is not a limitation of the method itself,
575 but rather the chosen architecture. This could be resolved by using a transformer so
576 that relationships between signal or image patches at distant spatio/temporal locations
577 could be captured. Nevertheless, the fact that our implementation is still competitive
578 with transformer-based architectures (see Table 2) shows the benefit of our method
579 regardless of architectural choices. Second, finding the optimal latent parameter via
580 optimisation requires in-network iterative optimisation which is more expensive than a

581 simple forward pass. While it is possible to use only the encoded mean latent variable
582 without optimisation of the additional noise parameter (results in third rows of ablation
583 study in Table 3) this does reduce performance. It amounts to assuming that the con-
584 tinuation is the most likely without considering the contents of the actual continuation.
585 In the context of anomaly detection, this is not optimal. Finally, while our model learns
586 the distribution of temporal continuation, its application to anomaly detection requires
587 selection of a similarity function in order to distinguish normal from abnormal. It is
588 likely that performance would be improved if this function could itself be learnt. How-
589 ever, this would require supervision in the form of example anomalies which is not
590 always available depending on the problem.

591 In future, we would like to explore using the trained encoder as a pre-trained back-
592 bone for downstream tasks. The encoder has learned to embed sufficient information
593 about a given time series segment to predict the following segment. We believe that
594 this means it would likely perform well when fine tuned for other tasks such as classi-
595 fication or object detection. Secondly, we would also like to explore the use of other
596 architectures such as transformers which naturally handle sequential data and so may
597 perform well for time series data. Finally, we would like to test whether our method
598 generalises to other temporal data modalities such as video and audio.

599 **References**

- 600 [1] H. A. Dau, A. Bagnall, K. Kamgar, C.-C. M. Yeh, Y. Zhu, S. Gharghabi, C. A.
601 Ratanamahatana, E. Keogh, The ucr time series archive, *IEEE/CAA Journal of*
602 *Automatica Sinica* 6 (6) (2019) 1293–1305.
- 603 [2] A. L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G.
604 Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, H. E. Stanley, Physiobank, phys-
605 iotoolkit, and physionet: components of a new research resource for complex
606 physiologic signals, *circulation* 101 (23) (2000) e215–e220.
- 607 [3] G. B. Moody, R. G. Mark, The impact of the mit-bih arrhythmia database, *IEEE*
608 *engineering in medicine and biology magazine* 20 (3) (2001) 45–50.

- 609 [4] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, O. Winther, Autoencoding be-
610 yond pixels using a learned similarity metric, in: International conference on
611 machine learning, PMLR, 2016, pp. 1558–1566.
- 612 [5] L. Jing, Y. Tian, Self-supervised visual feature learning with deep neural net-
613 works: A survey, *IEEE Transactions on Pattern Analysis and Machine Intelli-*
614 *gence* 43 (11) (2021) 4037–4058. doi:10.1109/TPAMI.2020.2992393.
- 615 [6] C. Doersch, A. Gupta, A. A. Efros, Unsupervised visual representation learning
616 by context prediction, in: *Proc. ICCV*, 2015, pp. 1422–1430.
- 617 [7] R. Zhang, P. Isola, A. A. Efros, Colorful image colorization, in: *Proc. ECCV*,
618 2016, pp. 649–666.
- 619 [8] S. Gidaris, P. Singh, N. Komodakis, Unsupervised representation learning by pre-
620 dicting image rotations, in: *Proc. ICLR*, 2018.
- 621 [9] I. Misra, C. L. Zitnick, M. Hebert, Shuffle and learn: unsupervised learning using
622 temporal order verification, in: *Proc. ECCV*, 2016, pp. 527–544.
- 623 [10] D. Wei, J. J. Lim, A. Zisserman, W. T. Freeman, Learning and using the arrow of
624 time, in: *Proc. CVPR*, 2018, pp. 8052–8060.
- 625 [11] X. Wang, A. Jabri, A. A. Efros, Learning correspondence from the cycle-
626 consistency of time, in: *Proc. CVPR*, 2019.
- 627 [12] S. Hitawala, Comparative study on generative adversarial networks, arXiv
628 preprint arXiv:1801.04271 (2018).
- 629 [13] M. Mirza, S. Osindero, Conditional generative adversarial nets, arXiv preprint
630 arXiv:1411.1784 (2014).
- 631 [14] D. P. Kingma, M. Welling, Auto-encoding variational bayes, arXiv preprint
632 arXiv:1312.6114 (2013).
- 633 [15] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, R. Girshick, Masked autoencoders are
634 scalable vision learners, in: *Proceedings of the IEEE/CVF conference on com-*
635 *puter vision and pattern recognition*, 2022, pp. 16000–16009.

- 636 [16] A. X. Lee, R. Zhang, F. Ebert, P. Abbeel, C. Finn, S. Levine, Stochastic adversarial
637 video prediction, arXiv preprint arXiv:1804.01523 (2018).
- 638 [17] M. Assran, Q. Duval, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, Y. LeCun,
639 N. Ballas, Self-supervised learning from images with a joint-embedding predic-
640 tive architecture, in: Proceedings of the IEEE/CVF Conference on Computer Vi-
641 sion and Pattern Recognition, 2023, pp. 15619–15629.
- 642 [18] H. Fan, F. Zhang, Y. Gao, Self-supervised time series representation learning by
643 inter-intra relational reasoning, arXiv preprint arXiv:2011.13548 (2020).
- 644 [19] S. Liu, A. Mallol-Ragolta, E. Parada-Cabeleiro, K. Qian, X. Jing, A. Kathan,
645 B. Hu, B. W. Schuller, Audio self-supervised learning: A survey, arXiv preprint
646 arXiv:2203.01205 (2022).
- 647 [20] R. Giri, S. V. Tenneti, F. Cheng, K. Helwani, U. Isik, A. Krishnaswamy, Self-
648 supervised classification for detecting anomalous sounds., in: DCASE, 2020, pp.
649 46–50.
- 650 [21] X. Chang, T. Maekaku, P. Guo, J. Shi, Y.-J. Lu, A. S. Subramanian, T. Wang, S.-w.
651 Yang, Y. Tsao, H.-y. Lee, et al., An exploration of self-supervised pretrained rep-
652 resentations for end-to-end speech recognition, in: 2021 IEEE Automatic Speech
653 Recognition and Understanding Workshop (ASRU), IEEE, 2021, pp. 228–235.
- 654 [22] Z. Chen, Y. Zhang, A. Rosenberg, B. Ramabhadran, G. Wang, P. Moreno, Inject-
655 ing text in self-supervised speech pretraining, in: 2021 IEEE Automatic Speech
656 Recognition and Understanding Workshop (ASRU), IEEE, 2021, pp. 251–258.
- 657 [23] L. C. Pickup, Z. Pan, D. Wei, Y. Shih, C. Zhang, A. Zisserman, B. Scholkopf,
658 W. T. Freeman, Seeing the arrow of time, in: Proceedings of the IEEE Conference
659 on Computer Vision and Pattern Recognition, 2014, pp. 2035–2042.
- 660 [24] D. Kim, D. Cho, I. S. Kweon, Self-supervised video representation learning with
661 space-time cubic puzzles, in: Proceedings of the AAAI conference on artificial
662 intelligence, Vol. 33, 2019, pp. 8545–8552.

- 663 [25] S. Saha, F. Bovolo, L. Bruzzone, Change detection in image time-series using
664 unsupervised lstm, *IEEE Geoscience and Remote Sensing Letters* (2020).
- 665 [26] L. Tao, X. Wang, T. Yamasaki, Self supervised video representation using pretext-
666 contrastive learning, *arXiv preprint arXiv:2010.15464* 2 (2020) 2.
- 667 [27] J. Wang, Y. Lin, A. J. Ma, P. C. Yuen, Self-supervised temporal discrimina-
668 tive learning for video representation learning, *arXiv preprint arXiv:2008.02129*
669 (2020).
- 670 [28] J. Kaur, S. Das, Future frame prediction of a video sequence, *arXiv preprint*
671 *arXiv:2009.01689* (2020).
- 672 [29] Y. Wang, L. Jiang, M.-H. Yang, L.-J. Li, M. Long, L. Fei-Fei, Eidetic 3d lstm: A
673 model for video prediction and beyond, in: *International conference on learning*
674 *representations*, 2018.
- 675 [30] S. Oprea, P. Martinez-Gonzalez, A. Garcia-Garcia, J. A. Castro-Vargas, S. Orts-
676 Escolano, J. Garcia-Rodriguez, A. Argyros, A review on deep learning techniques
677 for video prediction, *IEEE Transactions on Pattern Analysis and Machine Intelli-*
678 *gence* (2020).
- 679 [31] L. Jiang, M. Xu, Z. Wang, Predicting video saliency with object-to-motion cnn
680 and two-layer convolutional lstm, *arXiv preprint arXiv:1709.06316* (2017).
- 681 [32] W. Lotter, G. Kreiman, D. Cox, Deep predictive coding networks for video pre-
682 diction and unsupervised learning, *arXiv preprint arXiv:1605.08104* (2016).
- 683 [33] C. Finn, I. Goodfellow, S. Levine, Unsupervised learning for physical interaction
684 through video prediction, *Advances in neural information processing systems* 29
685 (2016).
- 686 [34] R. Villegas, J. Yang, S. Hong, X. Lin, H. Lee, Decomposing motion and content
687 for natural video sequence prediction, *arXiv preprint arXiv:1706.08033* (2017).

- 688 [35] M. Giannoulis, A. Harris, V. Barra, Ditan: A deep-learning domain agnostic
689 framework for detection and interpretation of temporally-based multivariate
690 anomalies, *Pattern Recognition* 143 (2023) 109814.
- 691 [36] J. Audibert, P. Michiardi, F. Guyard, S. Marti, M. A. Zuluaga, Do deep neural
692 networks contribute to multivariate time series anomaly detection?, *Pattern*
693 *Recognition* 132 (2022) 108945.
- 694 [37] T. Mokoena, T. Celik, V. Marivate, Why is this an anomaly? explaining anomalies
695 using sequential explanations, *Pattern Recognition* 121 (2022) 108227.
- 696 [38] J. Pereira, M. Silveira, Learning representations from healthcare time series data
697 for unsupervised anomaly detection, in: 2019 IEEE international conference on
698 big data and smart computing (BigComp), IEEE, 2019, pp. 1–7.
- 699 [39] J. Yang, Y. Shi, Z. Qi, Learning deep feature correspondence for unsupervised
700 anomaly detection and segmentation, *Pattern Recognition* 132 (2022) 108874.
- 701 [40] X. Zhang, J. Mu, X. Zhang, H. Liu, L. Zong, Y. Li, Deep anomaly detection with
702 self-supervised learning and adversarial training, *Pattern Recognition* 121 (2022)
703 108234.
- 704 [41] V. Zavrtanik, M. Kristan, D. Skočaj, Reconstruction by inpainting for visual
705 anomaly detection, *Pattern Recognition* 112 (2021) 107706.
- 706 [42] Y. Zhou, X. Song, Y. Zhang, F. Liu, C. Zhu, L. Liu, Feature encoding with autoen-
707 coders for weakly supervised anomaly detection, *IEEE Transactions on Neural*
708 *Networks and Learning Systems* 33 (6) (2021) 2454–2465.
- 709 [43] S. Akcay, A. Atapour-Abarghouei, T. P. Breckon, Ganomaly: Semi-supervised
710 anomaly detection via adversarial training, in: *Computer Vision–ACCV 2018:*
711 *14th Asian Conference on Computer Vision*, Perth, Australia, December 2–6,
712 2018, Revised Selected Papers, Part III 14, Springer, 2019, pp. 622–637.
- 713 [44] B. Zhou, S. Liu, B. Hooi, X. Cheng, J. Ye, Beatgan: Anomalous rhythm detection
714 using adversarially generated time series., in: *IJCAI*, Vol. 2019, 2019, pp. 4433–
715 4439.

- 716 [45] T.-W. Tang, H. Hsu, W.-R. Huang, K.-M. Li, Industrial anomaly detection with
717 skip autoencoder and deep feature extractor, *Sensors* 22 (23) (2022) 9327.
- 718 [46] J. Liu, Y. Huang, S. Wang, X. Zhao, Q. Zou, X. Zhang, Rail fastener defect
719 inspection method for multi railways based on machine vision, *Railway Sciences*
720 1 (2) (2022) 210–223.
- 721 [47] Z. Zamanzadeh Darban, G. I. Webb, S. Pan, C. Aggarwal, M. Salehi, Deep learn-
722 ing for time series anomaly detection: A survey, *ACM Computing Surveys* 57 (1)
723 (2024) 1–42.
- 724 [48] Z. Wang, C. Pei, M. Ma, X. Wang, Z. Li, D. Pei, S. Rajmohan, D. Zhang, Q. Lin,
725 H. Zhang, et al., Revisiting VAE for unsupervised time series anomaly detection:
726 A frequency perspective, in: *Proceedings of the ACM Web Conference 2024*,
727 2024, pp. 3096–3105.
- 728 [49] H. Kang, P. Kang, Transformer-based multivariate time series anomaly detection
729 using inter-variable attention mechanism, *Knowledge-Based Systems* 290 (2024)
730 111507.
- 731 [50] J. Miao, H. Tao, H. Xie, J. Sun, J. Cao, Reconstruction-based anomaly detection
732 for multivariate time series using contrastive generative adversarial networks, *In-*
733 *formation Processing & Management* 61 (1) (2024) 103569.
- 734 [51] Z. Z. Darban, G. I. Webb, S. Pan, C. C. Aggarwal, M. Salehi, CARLA: Self-
735 supervised contrastive representation learning for time series anomaly detection,
736 *Pattern Recognition* 157 (2025) 110874.
- 737 [52] D. Kim, S. Park, J. Choo, When model meets new normals: test-time adaptation
738 for unsupervised time-series anomaly detection, in: *Proceedings of the AAAI*
739 *conference on artificial intelligence*, Vol. 38, 2024, pp. 13113–13121.
- 740 [53] Y. Zhou, X. Xu, J. Song, F. Shen, H. T. Shen, Msflow: Multiscale flow-based
741 framework for unsupervised anomaly detection, *IEEE Transactions on Neural*
742 *Networks and Learning Systems* (2024).

- 743 [54] H. Yao, M. Liu, Z. Yin, Z. Yan, X. Hong, W. Zuo, GLAD: towards better re-
744 construction with global and local adaptive diffusion models for unsupervised
745 anomaly detection, in: European Conference on Computer Vision, Springer,
746 2024, pp. 1–17.
- 747 [55] S. Dai, Y. Wu, X. Li, X. Xue, Generating and reweighting dense contrastive pat-
748 terns for unsupervised anomaly detection, in: Proceedings of the AAAI Confer-
749 ence on Artificial Intelligence, Vol. 38, 2024, pp. 1454–1462.
- 750 [56] J. Liu, K. Wu, Q. Nie, Y. Chen, B.-B. Gao, Y. Liu, J. Wang, C. Wang, F. Zheng,
751 Unsupervised continual anomaly detection with contrastively-learned prompt, in:
752 Proceedings of the AAAI conference on artificial intelligence, Vol. 38, 2024, pp.
753 3639–3647.
- 754 [57] H. Liu, Z. Wan, W. Huang, Y. Song, X. Han, J. Liao, PD-GAN: Probabilistic
755 diverse gan for image inpainting, in: Proceedings of the IEEE/CVF Conference
756 on Computer Vision and Pattern Recognition, 2021, pp. 9371–9381.
- 757 [58] D. Brunet, E. R. Vrscay, Z. Wang, On the mathematical properties of the structural
758 similarity index, *IEEE Transactions on Image Processing* 21 (4) (2011) 1488–
759 1499.
- 760 [59] Q. Lei, J. Yi, R. Vaculin, L. Wu, I. S. Dhillon, Similarity preserving representation
761 learning for time series clustering, in: Proceedings of the 28th International Joint
762 Conference on Artificial Intelligence, 2019, pp. 2845–2851.
- 763 [60] F. Karim, S. Majumdar, H. Darabi, S. Chen, Lstm fully convolutional networks
764 for time series classification, *IEEE access* 6 (2017) 1662–1669.
- 765 [61] P. Malhotra, A. Ramakrishnan, G. Anand, L. Vig, P. Agarwal, G. Shroff,
766 Lstm-based encoder-decoder for multi-sensor anomaly detection, *arXiv preprint*
767 *arXiv:1607.00148* (2016).
- 768 [62] Y. Liu, J. Chen, S. Wu, Z. Liu, H. Chao, Incremental fuzzy c medoids clustering
769 of time series data using dynamic time warping distance, *Plos one* 13 (5) (2018)
770 e0197499.

- 771 [63] M. Thill, S. Däubener, W. Konen, T. Bäck, P. Barancikova, M. Holena, T. Horvat,
772 M. Pleva, R. Rosa, Anomaly detection in electrocardiogram readings with stacked
773 lstm networks., in: ITAT, 2019, pp. 17–25.
- 774 [64] G. Sivapalan, K. K. Nundy, S. Dev, B. Cardiff, D. John, Annet: A lightweight
775 neural network for ECG anomaly detection in iot edge sensors, IEEE Transactions
776 on Biomedical Circuits and Systems 16 (1) (2022) 24–35.
- 777 [65] P. Matias, D. Folgado, H. Gamboa, A. V. Carreiro, Robust anomaly detection
778 in time series through variational autoencoders and a local similarity score., in:
779 Biosignals, 2021, pp. 91–102.
- 780 [66] A. Alamr, A. Artoli, Unsupervised transformer-based anomaly detection in ECG
781 signals, Algorithms 16 (3) (2023) 152.