

# Detecting Abnormal Operations in Concentrated Solar Power Plants from Irregular Sequences of Thermal Images

Sukanya Patra  
University of Mons  
Mons, Belgium

sukanya.patra@umons.ac.be

Nicolas Sournac  
University of Mons  
Mons, Belgium

nicolas.sournac@student.umons.ac.be

Souhaib Ben Taieb  
University of Mons  
Mons, Belgium

souhaib.bentaieb@umons.ac.be

## ABSTRACT

Concentrated Solar Power (CSP) plants store energy by heating a storage medium with an array of mirrors that focus sunlight onto solar receivers atop a central tower. Operating at high temperatures these receivers face risks such as freezing, deformation, and corrosion, leading to operational failures, downtime, or costly equipment damage. We study the problem of anomaly detection (AD) in sequences of thermal images collected over a year from an operational CSP plant. These images are captured at irregular intervals ranging from one to five minutes throughout the day by infrared cameras mounted on solar receivers. Our goal is to develop a method to extract useful representations from high-dimensional thermal images for AD. It should be able to handle temporal features of the data, which include irregularity, temporal dependency between images and non-stationarity due to a strong daily seasonal pattern. The co-occurrence of low-temperature anomalies that resemble normal images from the start and the end of the operational cycle with high-temperature anomalies poses an additional challenge. We first evaluate state-of-the-art deep image-based AD methods, which have been shown to be effective in deriving meaningful image representations for the detection of anomalies. Then, we introduce a forecasting-based AD method that predicts future thermal images from past sequences and timestamps via a deep sequence model. This method effectively captures specific temporal data features and distinguishes between difficult-to-detect temperature-based anomalies. Our experiments demonstrate the effectiveness of our approach compared to multiple SOTA baselines across multiple evaluation metrics. We have also successfully deployed our solution on five months of unseen data, providing critical insights to our industry partner for the maintenance of the CSP plant. Our code<sup>1</sup> is publicly accessible. Additionally, as our dataset is confidential, we release a simulated dataset<sup>2</sup>.

## CCS CONCEPTS

• **Computing methodologies** → **Machine learning; Anomaly detection; Unsupervised learning.**

<sup>1</sup><https://github.com/sukanyapatra1997/ForecastAD>

<sup>2</sup><https://tinyurl.com/kdd2024Dataset>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

KDD '24, August 25–29, 2024, Barcelona, Spain

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 979-8-4007-0490-1/24/08

<https://doi.org/10.1145/3637528.3671623>

## KEYWORDS

Deep Image Anomaly Detection, Unsupervised Learning, Irregular Time-series, Non-stationarity, Concentrated Solar Power Plants

## ACM Reference Format:

Sukanya Patra, Nicolas Sournac, and Souhaib Ben Taieb. 2024. Detecting Abnormal Operations in Concentrated Solar Power Plants from Irregular Sequences of Thermal Images. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24)*, August 25–29, 2024, Barcelona, Spain. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3637528.3671623>

## 1 INTRODUCTION

The focus on renewable energies to counteract climate change has intensified recently. However, a critical challenge in adopting renewable energy sources is ensuring on-demand generation and dispatchability. A promising solution to this challenge is the integration of Thermal Energy Storage (TES) facilities, which temporarily store energy by heating or cooling a storage medium, such as water or molten salt. Concentrated Solar Power (CSP) plants effectively utilize TES for storing energy by heating the medium with an array of mirrors focused on solar receivers atop a central tower [40]. These solar receivers are composed of vertical heat exchanger tubes arranged in panel form, allowing the medium to flow through them.

Operating at extreme temperatures, these systems are prone to adverse effects, including the freezing of the medium (affecting a subset of vertical tubes with significantly higher temperatures), damage to heat-resistant coatings, and deformation and corrosion of the heat exchanger tubes. Therefore, meticulous monitoring of the process is crucial. Given the vast amount of data generated from multiple sensors, manually detecting abnormal behaviours becomes impractical. This necessitates an automated system capable of immediately identifying abnormal behaviours. The advantages of such a system are twofold: it ensures smooth operation and uninterrupted power generation by minimizing downtime, and it reduces the risk of further equipment damage by allowing for prompt failure responses. This approach also leads to an extended operational lifetime for the CSP plant.

In this paper, our goal is to develop a deep image-based anomaly detection (AD) [23, 27] method to identify abnormal behaviours in sequences of thermal images collected over a span of one year from an operational CSP plant. These images are captured at irregular intervals ranging from one to five minutes throughout the day by infrared cameras mounted on solar receivers. Our problem is related to data-driven Predictive Maintenance (PdM), where the state of equipment in industrial processes is monitored to predict future failures [31].

Specifically, we aim to develop an AD method capable of extracting useful representations from high-dimensional thermal images. It should be able to handle temporal features of the data, which include irregularity, temporal dependency between images and non-stationarity due to a strong daily seasonal pattern. An additional challenge is the coexistence of low-temperature anomalies that resemble low-temperature normal images from the start and the end of the operational cycle alongside high-temperature anomalies. This necessitates learning the current state of the operational cycle to correctly identify anomalous operations.

We first examine the performance of state-of-the-art (SOTA) deep AD methods that have been successful in extracting useful image representations for anomaly detection, such as CFlow [12], PatchCore [26], and DRÆM [38]. Our experiments confirm that neglecting the temporal features of the data leads to low accuracy, especially in distinguishing low-temperature normal samples from anomalies. Then, we explore a new forecasting-based AD method, ForecastAD, which predicts the image for a given future time based on a sequence of past observed images and their timestamps using a deep sequence model. ForecastAD extracts relevant representations from the high-dimensional images and captures the normal behaviour of the solar receivers, taking into account the temporal features of the data. An anomaly is then defined as a significant deviation from the learned normal behaviour. Our experiments demonstrate the effectiveness of ForecastAD compared to multiple SOTA baselines across various evaluation metrics. We have also successfully deployed our solution on five months of unseen data, providing critical insights for the maintenance of the CSP plant.

We first discuss related work in Section 2. Secondly, in Section 3, we present the case study on detecting anomalous behaviours in CSP plants. Then, we explain our forecasting-based approach, ForecastAD in Section 4. Finally, in Section 5, we discuss our empirical results before providing our concluding remarks in Section 6.

## 2 RELATED WORK

AD has been extensively studied over several decades [27]. The major AD methods can be broadly classified into four categories - density-based, reconstruction-based, classification-based approaches, and feature embedding-based methods.

**Density-based methods.** The density-based approach aims to estimate the probability distribution of normal data by assuming that normal samples are more likely to occur under the estimated distribution than anomalous samples. Traditional methods [4, 14, 24] fit a model to arbitrary data distribution but encounter challenges in higher-dimensional input spaces due to the curse of dimensionality. To overcome this, they are often applied to low-dimensional latent representations obtained using techniques like Autoencoder (AE) and Variational Autoencoder (VAE). Neural generative models, such as VAE and Generative Adversarial Networks (GANs), are deep learning-based methods that estimate a neural network's parameters to map a predetermined source distribution to the input data distribution. Recent AD methods such as CFlow [12] and FastFlow [37] further build on normalizing flows. However, studies demonstrated that normalizing flows often struggle to detect anomalies and assign them a higher likelihood [15, 20, 21].

**Reconstruction-based methods.** Reconstruction-based methods operate on the assumption that encoder-decoder models trained on normal samples will exhibit poor performance for anomalous samples. Common deep reconstruction models used include AE or VAE-based approaches, while advanced strategies involve reconstruction by memorized normality [11], model architecture adaptation [16] and partial/conditional reconstruction [22, 36]. Recent approach DRÆM [38] trains a discriminative network alongside the reconstruction network to localize anomalies without the need for any further post-processing steps. Generative models like GANs are also widely employed for anomaly detection, as the discriminator inherently calculates reconstruction loss for samples [39].

**One-class classification.** Anomaly detection can be approached as a one-class classification [13, 32] or single-class classification [8, 19] problem. Unlike density-based methods, classification-based techniques, such as One-Class Support Vector Machine (OC-SVM) [29], directly estimate a decision boundary to differentiate between normal and anomalous samples. However, this task can be challenging due to imbalanced datasets, where normal samples vastly outnumber anomalous samples. To address this, techniques like Support Vector Data Descriptor (SVDD) [32–34] derive a tight spherical bound. To enhance the expressivity of the classical models, deep learning models are used to learn the features from the data [10, 28].

**Feature Embedding-based methods** There are mainly two different types of feature embedding-based anomaly detection methods: memory bank [6, 17, 26], student-teacher [2, 41]. The main idea of memory bank methods is to extract features of nominal images and store them in a memory bank during the training phase. During the testing phase, the feature of a test image is used as a query to match the stored nominal features. The performance of the memory bank methods heavily depends on the completeness of the memory bank requiring a large number of nominal images. In the student-teacher approaches [7, 41], the student network learns to extract features of the nominal samples, similar to the teacher model. For anomalous images, the features extracted by the student network should be different from the teacher network.

## 3 A CASE STUDY ON DETECTING ANOMALOUS BEHAVIOURS IN CSP PLANTS

A CSP plant consists of two main components, namely: (i) the Thermal Solar Receiver and (ii) the Steam Generator. The Thermal Solar Receiver placed on top of a central tower in the plant acts as a solar furnace. On the ground surrounding the tower, an array of flat, movable mirrors called heliostats concentrate the sun rays on the solar receiver. The receiver consists of vertical heat exchanger tubes through which the heat transfer medium flows, absorbing the heat from the concentrated sun rays. Then, the absorbed thermal energy is utilized to generate superheated steam, which runs the Steam Generator for the production of energy. In this work, we focus on detecting anomalous behaviours of the Thermal Solar Receiver using data obtained from an operational plant.

CSP plants utilize high-capacity fluids like molten salts as the heat transfer medium, which are stored in TES facilities for future use. This allows for the on-demand generation of energy, making CSP plants a viable alternative to fossil fuel-based energy plants.

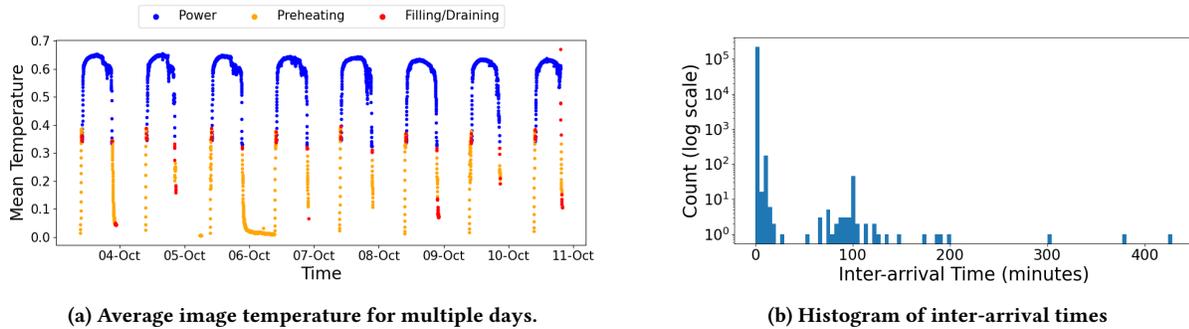


Figure 1: Visualisation of different properties of the data

However, due to operation in extreme temperatures, the solar thermal receivers are adversely impacted in several ways:

**[i] Blocked Tubes.** The molten salts passing through the heat exchanger tubes tend to freeze in localized zones when the temperature falls below a certain threshold, blocking them.

**[ii] Deformity.** The metal heat exchanger tubes in the receiver tend to expand due to the high temperatures. Uneven dilation of the tubes could eventually lead to deformity.

**[iii] Stress and Metal Fatigue.** The metal tubes in the receiver undergo expansion when exposed to high temperatures during regular operation and contraction when the operation ends. Such repeated changes lead to metal fatigue. Additionally, the pressure generated from the flowing molten salts exerts stress on the tubes.

**[iv] Corrosion.** Due to the interaction of the metal with the molten salt flowing through the tubes, it tends to deteriorate over time. These reactions are further accelerated due to the high temperatures in the receiver.

Hence, CSP plants require close monitoring to guarantee seamless operation and continuous power generation. Achieving this requires the analysis of data collected by numerous sensors installed on the Solar Receiver. Yet, the vast volume of data generated renders manual inspection unfeasible, thus underscoring the need for an automated, data-driven monitoring system.

### 3.1 Data Description

The Thermal Solar Receiver is composed of several panels, each featuring vertical heat exchanger tubes. These tubes allow the heat transfer medium to flow through, effectively absorbing heat from the concentrated sunlight. Infrared (IR) cameras, strategically positioned around the solar receiver, capture the surface temperature, producing thermal images with dimensions of  $184 \times 608$ . These images are captured approximately every one to five minutes, with each image’s timestamp recorded. During normal operations, the temperature of the heat transfer medium gradually increases as it traverses the vertical tubes from one end of the panel to the other, a direct result of absorbing heat from the concentrated sunlight. Consequently, the surface temperature patterns recorded by the IR cameras are anticipated to exhibit a smooth gradient, aligning with the medium’s flow direction. Our dataset covers a year of operational data without ground truth labels for the images (normal or

abnormal), making it an unsupervised anomaly detection problem. Note that throughout this work, we provide normalized images from the dataset for the sake of confidentiality.

Operation in CSP plants occurs across three distinct phases as depicted in Figure 1a: (i) *Preheating*, (ii) *Filling/Draining*, and (iii) *Power*. The molten salt used in CSP plants freezes when the temperature drops below a certain threshold. To avoid this, solar panels are initially heated during the *Preheating* phase. Then, in the subsequent *Filling/Draining* phase, molten salt is circulated within the panels. The *Power* phase initiates as the molten salt absorbs heat from sunlight, facilitating power generation. As operations conclude, the molten salt is drained from the panels during the *Filling/Draining* phase. Consequently, the panels commence cooling down, transitioning back to the *Preheating*. Our work focuses solely on the *Power* phase, as it is crucial for power generation and susceptible to damage from prolonged exposure to high temperatures.

**Data characteristics and modelling challenges.** Through extensive data analysis, we identified the following additional challenges, which are essential for modelling the solar receiver data:

**[i] Non-stationarity.** Figure 1a presents the average surface temperature across a specific week, highlighting temporal variations in the mean image temperature and demonstrating a clear pattern of daily seasonality in the data.

**[ii] Irregular sampling.** The images were captured at irregular time intervals, as illustrated in Figure 1b, which depicts the distribution of inter-arrival times. Additionally, the dataset lacks data for the extended periods when the plant was not operational.

**[iii] Temporal dependence.** The images exhibit a strong temporal dependence, influenced significantly by weather conditions.

**[iv] High dimensionality.** Anomalous characteristics often stay hidden and unnoticed due to data sparsity in high-dimensional spaces. Identifying features that capture the essential high-order, non-linear interactions needed for AD is thus challenging.

**[v] Large volume of data.** The dataset comprises images captured throughout a year of operation at approximately one to five-minute intervals, leading to a vast volume of data.

**[vi] Unlabeled data.** Our dataset lacks ground truth labels for the images, whether they are normal or abnormal, classifying our task as an unsupervised anomaly detection problem.

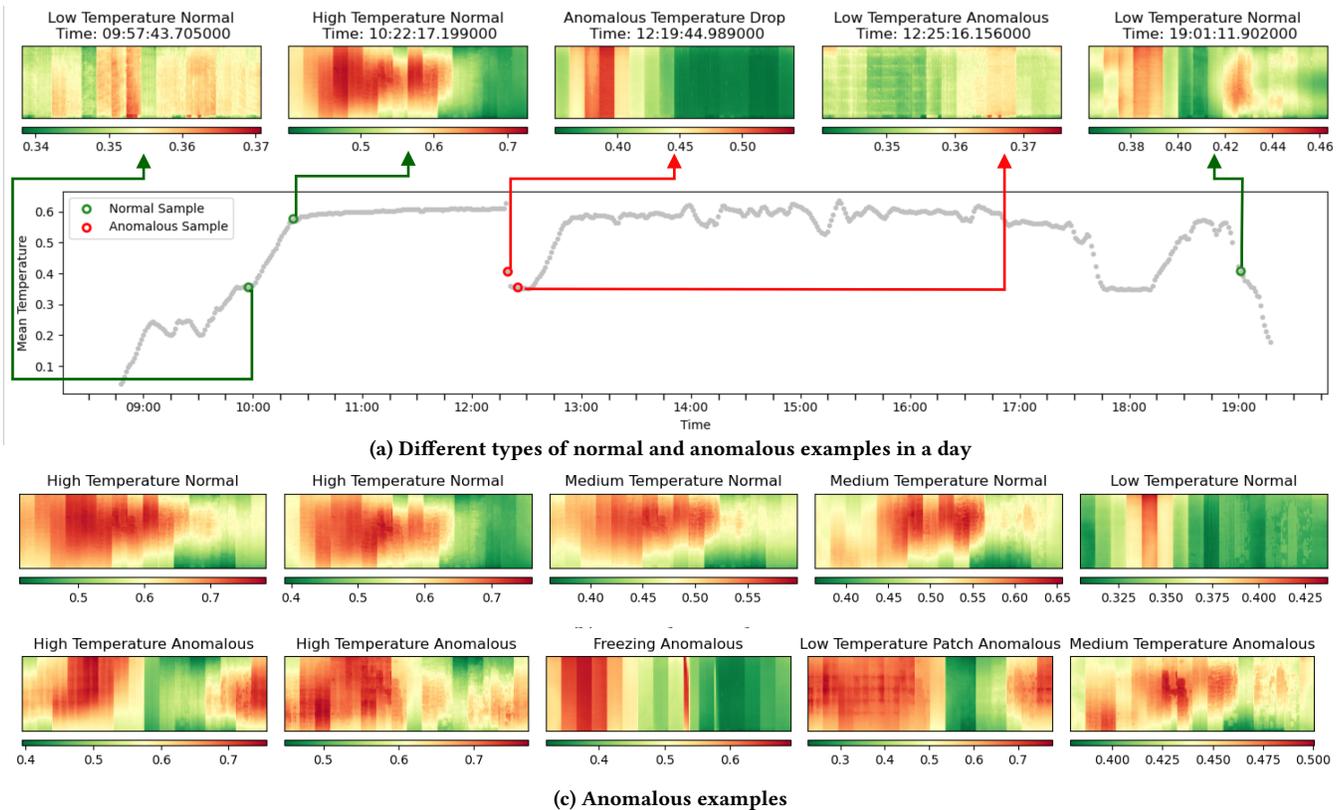


Figure 2: Examples of different types of normal and anomalous images

### 3.2 Data Labelling

To effectively assess the performance of various AD methods, we have labelled a subset of data from the CSP plant. This endeavour is notably complex due to the plant’s operation across multiple phases, each characterized by unique temperature ranges. Consequently, this diversity leads to a range of normal and anomalous sample types, as depicted in Figure 2. The challenge of identifying anomalies through the plant’s operational phases is evident in Figure 2a. Notably, normal images with low temperatures at the operation’s start and end (the left-most and right-most images in Figure 2a) closely resemble low-temperature anomalies (the second image from the right in Figure 2a). The distinction between these samples relies heavily on context.

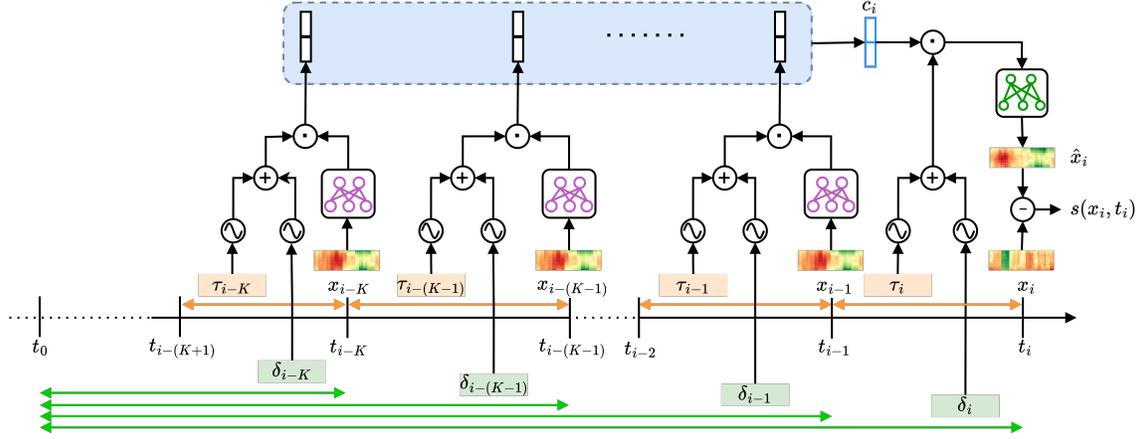
Moreover, the variable nature of anomalies adds a layer of complexity to the labelling process. Our approach to this challenge is informed by a deep understanding of the CSP plant’s operations and expert insights from the field. We categorize the *Power* phase into three distinct segments: (i) Starting (S), where the solar receiver’s mean temperature begins to rise; (ii) Middle (M), where it reaches and maintains its peak; and (iii) Ending (E), where it declines as the day concludes. In our preprocessing, we exclude days with significantly few samples or with a consistently low temperature throughout the M segment, likely indicative of sensor or system failures. For the S and E segments, samples showing a consistent temperature increase ( $> 5^{\circ}\text{C}$ ) or decrease ( $< -5^{\circ}\text{C}$ ), respectively, are deemed normal, whereas those displaying contrary trends are marked as anomalous. In the M segment, we apply the following

four rules for labelling:

**[R1.] Difference between consecutive images.** During the M segment of the *Power* phase, we expect a stable temperature. Significant deviations from the preceding observation indicate an abnormality. To detect such anomalies, we first compute the pixel-wise squared differences between every two consecutive images. For each pair, we select the 95th percentile of these pixel-wise differences as our *score*. Samples are then labelled as anomalous if their score exceeds the 99.9th percentile of the scores for all samples in the dataset

**[R2.] Difference from average daily temperature.** Samples with average temperatures that significantly deviate from the daily average temperature are labelled as anomalous. To identify these anomalies, we first compute the daily mean temperature. Then, we calculate the difference between each image’s average temperature and the mean temperature of the corresponding day, which serves as our *score*. Finally, samples are labelled as anomalous if their score falls below the 1st percentile of the distribution of scores across all samples in the dataset.

**[R3.] Difference with specific daily normal samples.** Rules R1 and R2 are limited to the detection of low-temperature anomalous samples. To address this, we select the first five images from the M segment of *Power* phase of each day to serve as a set of templates for that day. We then employ a similar methodology as in Rule R1, but instead of comparing an image to just the prior image, we compute the mean difference between the image and all



**Figure 3: Illustration of the end-to-end architecture of ForecastAD.** The model is trained to forecast the next image in the sequence given a context embedding  $c_i$  of  $K$  prior data points obtained using a sequence-to-sequence model. For  $(x_{i-k}, t_{i-k}, y_{i-k}) \in \mathcal{D}$  in the context, we sum the embeddings of inter-arrival time  $\tau_{i-k}$  and interval since the start of operation  $\delta_{i-k}$  and concatenate it with the image embedding. The anomaly score  $s(x_i, t_i)$  is computed as the difference between the forecasted and original image.

five templates of the corresponding day. Applying this rule allows us to obtain sets of high-temperature normal and abnormal samples, along with a diverse set of low-temperature abnormal samples.

**[R4.] Freezing statistics.** To identify the anomalous samples with characteristics such as freezing and low-temperature patches, we compute row-wise and column-wise differences within each image. First, we calculate the maximum value of the element-wise differences between two consecutive rows, which we term the *horizontal score*. Next, we compute the element-wise differences between consecutive columns and apply a Sobel filter [30] to detect vertical edges. The mean value of the elements detected by the Sobel filter across all columns is referred to as the *vertical score*. An image is labelled as anomalous if either the horizontal or the vertical score exceeds a predefined threshold.

Given the labelling rules, we first apply them to obtain an initial set of labels. Then, in collaboration with domain experts, we conduct a visual analysis of the labelled thermal images. For the visual analysis, we also perform clustering on the labelled samples and inspect the cluster centres in addition to analyzing each image individually. This thorough inspection leads to subsequent refinements of the labelled set, enhancing the reliability of the labels for accurately assessing anomaly detection methods.

### 3.3 General problem formulation

Consider a dataset  $\mathcal{D} = \{(x_i, t_i, y_i)\}_{i=1}^n$  consisting of  $n = 16,917$  triplets. Each  $x_i \in \mathcal{X} = \mathbb{R}_+^d$  corresponds to a thermal image with dimension  $d = H \times W$ , where the height  $H$  is 184 and the width  $W$  is 608. These images were captured at times  $t_i \in \mathbb{R}_+$ , and each  $y_i \in \{0, 1\}$  denotes the corresponding label, with 0 representing the normal class and 1 representing the anomalous class.

Let  $\mathcal{D}_N$ ,  $\mathcal{D}_V$  and  $\mathcal{D}_T$  denote disjoint training, validation and test sets, respectively, with  $\mathcal{D}_N \cup \mathcal{D}_V \cup \mathcal{D}_T = \mathcal{D}$ .  $\mathcal{D}_N$  is exclusively composed of normal samples, i.e.,  $y_i = 0$  for all  $(x_i, t_i, y_i) \in \mathcal{D}_N$ .  $\mathcal{D}_V$  and  $\mathcal{D}_T$  include both normal and anomalous samples.

Using the training set  $\mathcal{D}_N$ , the AD methods aim to learn a scoring function  $s(\cdot, \cdot) : \mathbb{R}_+^d \times \mathbb{R}_+ \rightarrow \mathbb{R}$  that assigns an anomaly score  $s(x, t)$  to any given point  $(x, t)$ . By using a threshold  $\lambda \in \mathbb{R}$ , this anomaly score can then be converted into a predicted label  $\hat{y}$  as follows:

$$\hat{y} = \begin{cases} 1, & \text{if } s(x, t) \geq \lambda; \\ 0, & \text{if } s(x, t) < \lambda. \end{cases} \quad (1)$$

## 4 A FORECASTING-BASED AD MODEL

We present a new forecasting-based AD method, denoted ForecastAD, to detect anomalous operations in the Thermal Solar Receiver of a CSP plant from irregular sequences of thermal images. The proposed method builds a forecasting model to reconstruct the thermal images using past observations as context. Images that are hard to reconstruct are considered anomalous. For a given image, our procedure can be summarized in the following steps: (i) extract feature embeddings for that image (Section 4.1), (ii) use the previous  $K$  images as *context* and encode them using a deep sequence model (Section 4.2), and (iii) using the context, reconstruct the image with a decoder forecasting model (Section 4.3), then assign an anomaly score based on the reconstruction error between the original and predicted image. We provide an overview of the architecture of ForecastAD in Figure 3 and summarize it in Algorithm 1.

### 4.1 Image Encoder

Using the training data,  $\mathcal{D}_N$ , we pre-train an encoder network to capture the inherent structure of our dataset’s images. The image encoder, denoted by  $\phi_e(\cdot; W_e) : \mathcal{X} \rightarrow \mathcal{Z}$ , transforms images from the high-dimensional input space  $\mathcal{X}$  to a compact latent space  $\mathcal{Z} = \mathbb{R}^{d'}$ , significantly reducing dimensionality where  $d' \ll d$ . We use an autoencoder framework for image reconstruction, with a decoder network  $\phi_d(\cdot; W_d) : \mathcal{Z} \rightarrow \mathcal{X}$  to project images from the latent space  $\mathcal{Z}$  back to the original input space  $\mathcal{X}$ . The autoencoder is given by  $\phi = \phi_e \circ \phi_d$ , with  $\circ$  indicating function composition.

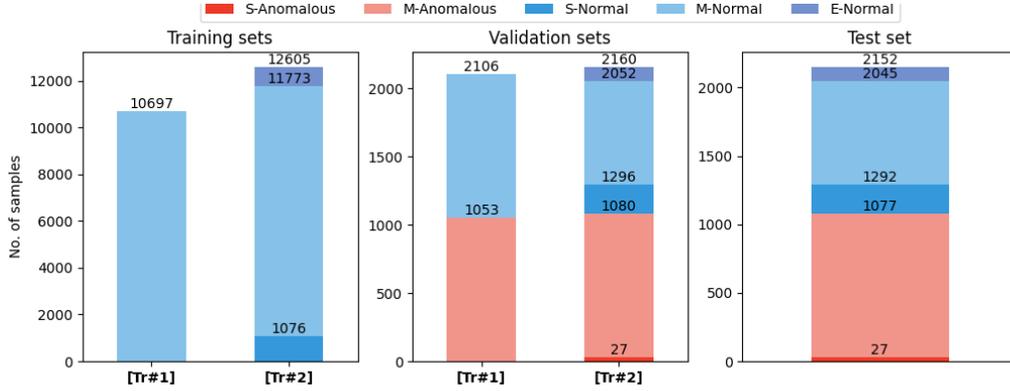


Figure 4: Dataset split for two different training setups and the test set

Given the high dimensionality of the input image, we opt for a multi-layer deep convolutional network as the image encoder, exploiting its effectiveness in extracting meaningful representations directly from the data [5]. We calculate the reconstruction loss between original data points  $x_i$  and their reconstructions  $\hat{x}_i = \phi(x_i)$  as:

$$\mathcal{L}_{\text{pre-train}} = \frac{1}{|\mathcal{D}_N|} \sum_{i=1}^{|\mathcal{D}_N|} \|x_i - \hat{x}_i\|_F^2, \quad (2)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm.

## 4.2 Context Encoder

To handle the irregular inter-arrival times  $\tau_i = t_i - t_{i-1}$  between successive ( $i$ )-th and ( $i-1$ )-th images, our deep sequence model incorporates both the image sequences and their associated irregular inter-arrival times. We employ a sinusoidal encoding  $\psi_i = f_{\text{sin}}(\tau_i)$ , inspired by the positional encoding technique in transformer models [35]. This method aligns with strategies used in Neural Temporal Point Processes [9].

In addition to the inter-arrival times between consecutive images, we also embed the relative time since the start of the operation  $t_0$ , i.e.,  $\delta_i = t_i - t_0$ , which provides information about the position of an image  $x_i$  within the operational cycle. Such temporal context helps in detecting challenging temperature-based anomalies, as it helps distinguish between low-temperature anomalies occurring mid-cycle and low-temperature normal images at the start of the operation. We use the same sinusoidal encoding for the interval  $\delta_i$  as  $\Psi_i = f_{\text{sin}}(\delta_i)$ . The sum of the two time embeddings  $\psi_i$  and  $\Psi_i$  is combined with the image embedding to obtain the final embedding  $\hat{z}_i = [z_i \oplus (\psi_i + \Psi_i)]$ , where  $z_i = \phi_e(x_i)$  represents the image embedding and  $\oplus$  denotes the concatenation operator.

We compactly encode the embeddings of the  $K$  samples preceding the image at timestep  $t_i$  into a fixed-dimensional vector  $c_i$ , termed the context vector for the  $i$ -th image. This can be accomplished with a deep sequence model. In our implementation, we opt for an LSTM  $\varphi(\cdot; W_c)$ , parameterized by  $W_c$ . For a given context sequence  $C_i = \{\hat{z}_{i-K}, \dots, \hat{z}_{i-1}\}$ , the hidden state is recursively updated from previous states as  $c_i = \varphi(c_{i-1}, \hat{z}_{i-1}; W_c)$ , starting from a random state.

## 4.3 Image Decoder

To predict the  $i$ -th image, we use  $c_i$ , the past context encoding,  $\psi_i$ , the embedding of the next inter-arrival time  $\tau_i$ , as well as  $\Psi_i$ , the embedding of the time duration since the start of the operation  $\delta_i$ . Specifically, we compute  $\hat{x}_i = \phi_d([c_i \oplus (\psi_i + \Psi_i)]; W_d)$  where  $\phi_d(\cdot; W_d)$  is the decoder network. Note that the decoder network is pre-trained along with the image encoder using the image reconstruction task on the images from the training set  $\mathcal{D}_N$ . The prediction error is computed as the Frobenius norm difference between the original and the forecasted image. The total training loss is obtained by averaging the prediction errors over all the training examples:

$$\mathcal{L}_{\text{train}} = \frac{1}{|\mathcal{D}_N|} \sum_{i=1}^{|\mathcal{D}_N|} \|x_i - \hat{x}_i\|_F^2 \quad (3)$$

Finally, the anomaly score of a new point  $(x, y, t)$  is defined as the associated prediction error, i.e.  $s(x, y, t) = \|x - \hat{x}\|_F^2$ . Algorithm 1 summarizes the different steps of our ForecastAD method.

## 5 EXPERIMENTS

### 5.1 Baselines

We first compare ForecastAD against simple methods, which detect anomalies based on statistical features extracted from the images. These features include the corresponding time of day, as well as the mean, maximum, and standard deviation of the temperature, to distinguish between normal and abnormal samples. We also evaluate against deep image-based AD methods, namely, autoencoder, FastFlow [37], PatchCore [26], PaDiM [6], CFlow [12], DRÆM [38], and Reverse Distillation [7]. Deep methods have been shown to be more effective than shallow ones for image AD [27], leveraging the deep neural networks' capability to extract representative features through multiple layers of abstraction.

### 5.2 Experimental Setup

**Network architectures and hyperparameters.** Except for the autoencoder, the baselines follow the implementation from Anomalib [1], which is a widely used library for benchmarking AD methods. Based on our experiments, we opted for a Deep Convolutional Autoencoder (DCAE) with a latent dimension of  $d' = 128$ . Detailed

**Table 1: Anomaly detection performance. Style: best in bold and second best using underline**

Train Setting	Model	AUROC (%)			AUPR (%)		
		[Ts#1]	[Ts#2]	[Ts#3]	[Ts#1]	[Ts#2]	[Ts#3]
[Tr#1]	Autoencoder	98.05 ( $\pm$ 0.74)	46.43 ( $\pm$ 1.61)	87.87 ( $\pm$ 0.26)	98.50 ( $\pm$ 0.54)	6.62 ( $\pm$ 0.19)	81.46 ( $\pm$ 0.55)
	CFlow [12]	94.68 ( $\pm$ 1.26)	39.99 ( $\pm$ 2.33)	82.91 ( $\pm$ 1.08)	96.28 ( $\pm$ 0.92)	5.94 ( $\pm$ 0.24)	76.11 ( $\pm$ 1.04)
	DRÆM [38]	97.70 ( $\pm$ 0.77)	40.48 ( $\pm$ 2.15)	87.38 ( $\pm$ 0.61)	97.97 ( $\pm$ 0.92)	6.13 ( $\pm$ 0.27)	82.05 ( $\pm$ 0.57)
	FastFlow [37]	99.83 ( $\pm$ 0.03)	47.32 ( $\pm$ 0.29)	<b>91.36 (<math>\pm</math> 0.25)</b>	99.87 ( $\pm$ 0.02)	<b>9.42 (<math>\pm</math> 1.18)</b>	86.02 ( $\pm$ 0.52)
	PaDiM [6]	99.85 ( $\pm$ 0.02)	49.86 ( $\pm$ 0.47)	91.23 ( $\pm$ 0.10)	<b>99.89 (<math>\pm</math> 0.01)</b>	7.73 ( $\pm$ 0.18)	<b>87.25 (<math>\pm</math> 0.21)</b>
	PatchCore [26]	99.23 ( $\pm$ 0.08)	<b>50.58 (<math>\pm</math> 0.37)</b>	89.04 ( $\pm$ 0.30)	99.43 ( $\pm$ 0.05)	7.25 ( $\pm$ 0.08)	83.41 ( $\pm$ 0.27)
	Reverse Distillation [7]	93.88 ( $\pm$ 1.13)	41.31 ( $\pm$ 2.19)	84.61 ( $\pm$ 1.54)	95.47 ( $\pm$ 0.79)	6.08 ( $\pm$ 0.30)	80.70 ( $\pm$ 1.37)
	ForecastAD	<b>99.86 (<math>\pm</math> 0.05)</b>	46.22 ( $\pm$ 1.06)	89.89 ( $\pm$ 0.35)	<b>99.89 (<math>\pm</math> 0.04)</b>	6.57 ( $\pm$ 0.09)	85.75 ( $\pm$ 0.65)
[Tr#2]	Autoencoder	96.67 ( $\pm$ 0.77)	45.92 ( $\pm$ 2.47)	85.45 ( $\pm$ 1.18)	96.91 ( $\pm$ 0.93)	6.69 ( $\pm$ 0.35)	78.61 ( $\pm$ 1.30)
	CFlow [12]	84.91 ( $\pm$ 2.72)	42.90 ( $\pm$ 2.71)	77.38 ( $\pm$ 2.98)	88.18 ( $\pm$ 2.02)	6.51 ( $\pm$ 0.39)	74.80 ( $\pm$ 3.24)
	DRÆM [38]	93.52 ( $\pm$ 0.52)	40.51 ( $\pm$ 1.33)	85.71 ( $\pm$ 0.78)	94.56 ( $\pm$ 0.44)	7.62 ( $\pm$ 1.01)	83.36 ( $\pm$ 1.08)
	FastFlow [37]	92.38 ( $\pm$ 0.72)	52.51 ( $\pm$ 1.09)	89.92 ( $\pm$ 0.68)	93.46 ( $\pm$ 0.60)	8.87 ( $\pm$ 0.46)	88.76 ( $\pm$ 0.56)
	PaDiM [6]	95.99 ( $\pm$ 0.37)	58.14 ( $\pm$ 1.00)	<u>92.28 (<math>\pm</math> 0.32)</u>	96.77 ( $\pm$ 0.32)	<u>11.50 (<math>\pm</math> 0.86)</u>	<u>90.73 (<math>\pm</math> 0.42)</u>
	PatchCore [26]	<b>96.78 (<math>\pm</math> 0.57)</b>	<u>60.15 (<math>\pm</math> 1.82)</u>	91.38 ( $\pm$ 0.42)	<b>97.57 (<math>\pm</math> 0.37)</b>	9.77 ( $\pm$ 0.79)	88.92 ( $\pm$ 0.72)
	Reverse Distillation [7]	87.19 ( $\pm$ 0.99)	57.22 ( $\pm$ 5.77)	84.04 ( $\pm$ 1.64)	86.09 ( $\pm$ 1.34)	10.64 ( $\pm$ 1.55)	78.83 ( $\pm$ 2.66)
	ForecastAD	94.78 ( $\pm$ 1.09)	<b>85.81 (<math>\pm</math> 1.23)</b>	<b>92.53 (<math>\pm</math> 0.81)</b>	96.92 ( $\pm$ 0.57)	<b>28.73 (<math>\pm</math> 1.70)</b>	<b>92.97 (<math>\pm</math> 0.36)</b>

**Algorithm 1:** Training process of ForecastAD

**Require:** Training dataset  $\mathcal{D}_N$ , Sinusoidal encoder  $f_{\sin}$   
 Image encoder  $\phi_e$ , Image decoder  $\phi_d$ , Number of epochs  $e$   
 Learning rate  $\eta$ , Context length  $K$

```

1 for (epoch = 1, 2, ..., e) and  $((x_i, t_i, y_i) \in \mathcal{D}_N)$  do
2   Initialize context embedding:
3    $c_i \leftarrow \text{random}()$ 
4   for  $k = K, K - 1, \dots, 1$  do
5     Calculate the time embeddings:
6      $\psi_{i-k} \leftarrow f_{\sin}(\tau_{i-k})$ 
7      $\Psi_{i-k} \leftarrow f_{\sin}(\delta_{i-k})$ 
8     Create joint embedding:
9      $\hat{z}_{i-k} \leftarrow [\phi_e(x_{i-k}; W_e) \oplus (\psi_{i-k} + \Psi_{i-k})]$ 
10    Update context embedding
11     $c_i \leftarrow \varphi(c_i, \hat{z}_{i-k}; W_c)$ 
12  end
13  Encode target time embeddings:
14   $\psi_i \leftarrow f_{\sin}(\tau_i)$ ,  $\Psi_i \leftarrow f_{\sin}(\delta_i)$ 
15  Predict the next data point:
16   $\hat{x}_i \leftarrow \phi_d([c_i \oplus (\psi_i + \Psi_i)]; W_d)$ 
17  Update the model parameters  $W_e$ ,  $W_d$  and  $W_c$  by
    minimising the loss  $\mathcal{L}_{\text{train}}$  (Eq. 3)
18 end
Return:  $\phi_e(\cdot; W_e)$ ,  $\phi_d(\cdot; W_d)$ ,  $\varphi(\cdot; W_c)$ 

```

architectural specifications are provided in Appendix A. The image encoder employed in ForecastAD mirrors the structure of the downsampling branch in DCAE. In ForecastAD, we adopt a 4-layer LSTM network with a hidden dimension of 128 to serve as the context encoder  $\varphi$ . For time encoding, the sinusoidal embedding has a dimension of 16. We adhere to the hyperparameters mentioned by the authors for the baseline methods. For ForecastAD, we use MSE and train using an Adam optimizer with a learning rate of 0.001 and weight decay of 0.00001. We use a pre-processing step for all the experiments where the images in the dataset are resized

to  $256 \times 256$  to be compatible with the baselines. Unless otherwise specified, we use a sequence length of  $K = 30$ .

**Dataset.** Our labelled dataset comprises days, which are segmented into training, validation, and test sets. Days featuring exclusively normal samples are allocated across these three sets, while those with anomalous samples are included in both the validation and test sets. To underscore the challenges presented by low-temperature samples, we adopt two training setups: (i) [Tr#1], incorporating training and validation samples solely from the M phase, and (ii) [Tr#2], comprising training and validation samples from the S, M, and E phases. Importantly, the test set in both scenarios consists of samples spanning the S, M, and E phases. The distribution of normal and anomalous samples across S, M, and E phases for these setups is depicted in Figure 4. For ForecastAD, we generate a sequence for each data point by selecting  $K$  preceding samples. If there are less than  $K$  prior samples, we duplicate the corresponding day’s first data point to form a  $K$ -length sequence. Lastly, for the first data point captured each day, we set the  $\tau$  and  $\delta$  to a small positive value  $\epsilon = 1e - 5$ .

**Model evaluation.** We evaluate the models based on the Area under the Receiver Operating Characteristics curve (AUROC) and the Area under the Precision-Recall curve (AUPR). To highlight the effectiveness of each model in distinguishing between low-temperature normal and anomalous behaviours, we utilize three test setups containing: (i) test samples in M [Ts#1], (ii) test samples in S-E [Ts#2], and (iii) test samples in S-M-E [Ts#3]. For the experiments below, we report mean over 5 runs along with one standard error.

### 5.3 Results and Discussion

We summarize the results over five runs for different training setups in Table 1. For the training setup [Tr#1], ForecastAD provides competitive results when compared to image-based SOTA models as measured by both AUROC and AUPR metrics over the test samples in [Ts#3]. Additionally, we observe good performance for image-based SOTA approaches in [Ts#1]. This performance can be attributed to the training exclusively on samples from M, which predominantly fall within the high-temperature region where temporal context is less critical. However, since the models are not

trained on low-temperature normal samples from the start and end of the operational cycle, their performance in [Ts#2] naturally declines. Specifically, the AUPR score is significantly low in [Ts#2], as the models tend to assign very high anomaly scores to most low-temperature samples found in S-E.

Considering the setup [Tr#2], we observe a drop in performance for the SOTA methods compared to [Tr#1]. This observation can be attributed to the fact that when the model is exposed to a limited number of low-temperature normal samples, it struggles to learn from them. Instead, these samples act as contamination, diminishing performance over the high-temperature samples in [Ts#1]. Additionally, baselines fail to distinguish between low-temperature normal and anomalous samples in [Ts#2], as they do not incorporate temporal features. ForecastAD significantly outperforms all baselines by approx. 25% in [Ts#2], while maintaining competitive performance across all test samples in [Ts#3]. In Appendix D, we provide an extended version of Table 1.

## 5.4 Ablation Study

**Importance of time-embedding and pre-training.** Table 2 shows the results of an ablation study to understand the importance of  $\tau$  and  $\delta$  in ForecastAD. In all configurations, we always keep the image encoding as part of the input. Firstly, we observe the lowest AUROC and AUPR scores in [Ts#2] when the context has only the encodings of  $K$ -prior images. It emphasizes the need to address the challenge posed by irregular sequences and co-occurrence of low-temperature normal and anomalous samples. Then, on considering either  $\tau$  or  $\delta$ , we observe a significant improvement in [Ts#2]. Furthermore, incorporating both  $\tau$  and  $\delta$  yields the best performance, highlighting that both  $\tau$  and  $\delta$  are necessary for reliable detection of anomalies. Lastly, we also empirically validate the impact of pre-training the image encoder and decoder using the image reconstruction task. Using the pre-trained models offers substantial enhancements in performance when compared to a randomly-initialized backbone.

**Table 2: Ablation of time-embedding and pre-training.**

Pre-train	$\tau$	$\delta$	AUROC (%)			AUPR (%)		
			[Ts#1]	[Ts#2]	[Ts#3]	[Ts#1]	[Ts#2]	[Ts#3]
-	✓	✓	94.60 ( $\pm 1.60$ )	75.30 ( $\pm 5.89$ )	90.58 ( $\pm 1.00$ )	96.78 ( $\pm 0.81$ )	23.54 ( $\pm 4.17$ )	89.93 ( $\pm 1.29$ )
-	✓	-	<b>97.12 (<math>\pm 0.44</math>)</b>	72.94 ( $\pm 7.15$ )	<b>92.74 (<math>\pm 1.26</math>)</b>	98.06 ( $\pm 0.29$ )	21.32 ( $\pm 4.29$ )	91.41 ( $\pm 1.88$ )
✓	-	-	94.59 ( $\pm 0.93$ )	<b>84.08 (<math>\pm 3.83</math>)</b>	92.49 ( $\pm 0.79$ )	96.56 ( $\pm 0.49$ )	<b>28.83 (<math>\pm 4.17</math>)</b>	<b>92.67 (<math>\pm 0.57</math>)</b>
✓	✓	-	92.71 ( $\pm 1.32$ )	82.71 ( $\pm 3.09$ )	91.12 ( $\pm 1.09$ )	95.58 ( $\pm 0.95$ )	26.31 ( $\pm 3.60$ )	92.15 ( $\pm 0.96$ )
✓	✓	✓	<b>94.78 (<math>\pm 1.09</math>)</b>	<b>85.81 (<math>\pm 1.23</math>)</b>	<b>92.53 (<math>\pm 0.81</math>)</b>	<b>96.92 (<math>\pm 0.57</math>)</b>	<b>28.73 (<math>\pm 1.70</math>)</b>	<b>92.97 (<math>\pm 0.36</math>)</b>

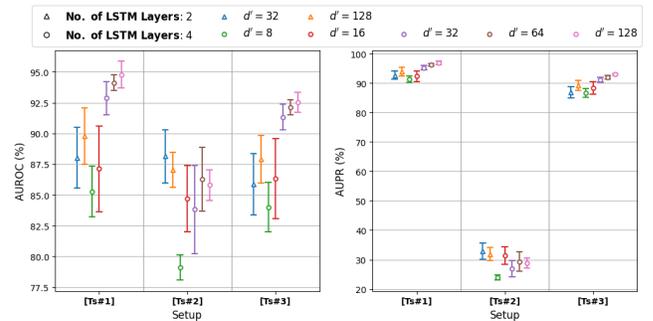
**Effect of context length ( $K$ ).** We report the AD performance of ForecastAD with varying context lengths  $K$  in Table 3. For context lengths  $K \leq 20$ , we do not observe any correlation between performance and context length. However, larger sequence lengths of 30 or 40 yield better performance. To limit computational demands, we did not consider larger sequence lengths and chose a sequence length of 30 for all our experiments.

**Effect of different architecture.** In Figure 5, we analyze the effect of the number of layers in LSTM and the latent dimension  $d'$  on the AUROC and AUPR scores. Firstly, we observe that larger latent dimensions lead to higher scores in most cases, regardless of the number of layers in LSTM. Secondly, ForecastAD performs better on [Ts#1] and [Ts#3] with a 4-layer LSTM, while a 2-layer LSTM yields better results on [Ts#2]. Based on this empirical observation,

**Table 3: Ablation of  $K$**

$K$	AUROC (%)			AUPR (%)		
	[Ts#1]	[Ts#2]	[Ts#3]	[Ts#1]	[Ts#2]	[Ts#3]
1	88.85 ( $\pm 2.55$ )	78.26 ( $\pm 1.86$ )	83.64 ( $\pm 1.72$ )	92.68 ( $\pm 1.49$ )	23.14 ( $\pm 2.63$ )	83.22 ( $\pm 1.01$ )
5	91.21 ( $\pm 0.94$ )	<b>87.83 (<math>\pm 1.45</math>)</b>	89.25 ( $\pm 0.81$ )	94.44 ( $\pm 0.51$ )	<b>32.24 (<math>\pm 2.15</math>)</b>	89.82 ( $\pm 0.49$ )
10	94.02 ( $\pm 1.81$ )	78.13 ( $\pm 3.91$ )	89.82 ( $\pm 0.77$ )	96.40 ( $\pm 0.93$ )	23.05 ( $\pm 2.07$ )	89.29 ( $\pm 0.86$ )
20	92.64 ( $\pm 1.21$ )	83.22 ( $\pm 2.90$ )	90.46 ( $\pm 1.31$ )	95.65 ( $\pm 0.60$ )	27.20 ( $\pm 3.45$ )	91.39 ( $\pm 0.93$ )
30	<b>94.78 (<math>\pm 1.09</math>)</b>	85.81 ( $\pm 1.23$ )	<b>92.53 (<math>\pm 0.81</math>)</b>	<b>96.92 (<math>\pm 0.57</math>)</b>	28.73 ( $\pm 1.70$ )	<b>92.97 (<math>\pm 0.36</math>)</b>
40	92.66 ( $\pm 1.46$ )	85.87 ( $\pm 0.92$ )	91.13 ( $\pm 1.26$ )	95.59 ( $\pm 0.73$ )	29.24 ( $\pm 1.49$ )	91.88 ( $\pm 0.83$ )
50	93.44 ( $\pm 0.72$ )	87.29 ( $\pm 1.77$ )	92.09 ( $\pm 0.45$ )	96.06 ( $\pm 0.34$ )	31.51 ( $\pm 1.52$ )	92.62 ( $\pm 0.31$ )
60	93.36 ( $\pm 0.75$ )	81.03 ( $\pm 4.25$ )	91.03 ( $\pm 1.51$ )	95.95 ( $\pm 0.44$ )	28.46 ( $\pm 3.83$ )	91.45 ( $\pm 1.58$ )

we chose a 4-layer LSTM with latent dimension  $d' = 128$ , which has the highest scores in [Ts#1] and [Ts#3] while having a comparable performance with the best configuration on [Ts#2].



**Figure 5: Ablation of different architectures**

## 5.5 Interpretability of ForecastAD

Interpretability of deep learning models is critical for high-risk applications to enhance transparency and trustworthiness. Therefore, we extract anomaly maps from ForecastAD corresponding to each image during inference. Recall that ForecastAD is trained with pixel-wise regression loss, and thus, the anomaly map can be computed as the difference between the original and forecasted images. Based on recent works on IAD [6, 26], we smoothed the anomaly maps using a Gaussian filter and normalized it using the minimum and maximum anomaly scores for the normal samples in the validation set. In Figure 6, we show the anomaly maps of 4 normal and 4 anomalous test samples, along with the image for reference. It can be seen that for specific types of anomalies, such as freezing, where we observe high-temperature streaks, ForecastAD assign high anomaly scores to those regions. Therefore, it aids the interpretability of the results from ForecastAD. To further enhance the understanding, the anomaly maps can be complemented by plots of mean temperature to show the sudden drops or rises in temperature resulting in the samples being anomalous.

## 5.6 Simulated Dataset

We have prepared a simulated dataset to ensure reproducibility and validation of the results. We use a variational autoencoder to generate the data. Additional details about the data generation are deferred to Appendix B. The distribution of normal and anomalous samples across S, M, and E phases for these setups is depicted in Figure 7. We have also compared our method to the baselines on the simulated dataset. The results are reported in Table 4 for training setup [Tr#2] and test setup [Ts#3], which are the main focus of our work. The results highlight the effectiveness of ForecastAD,

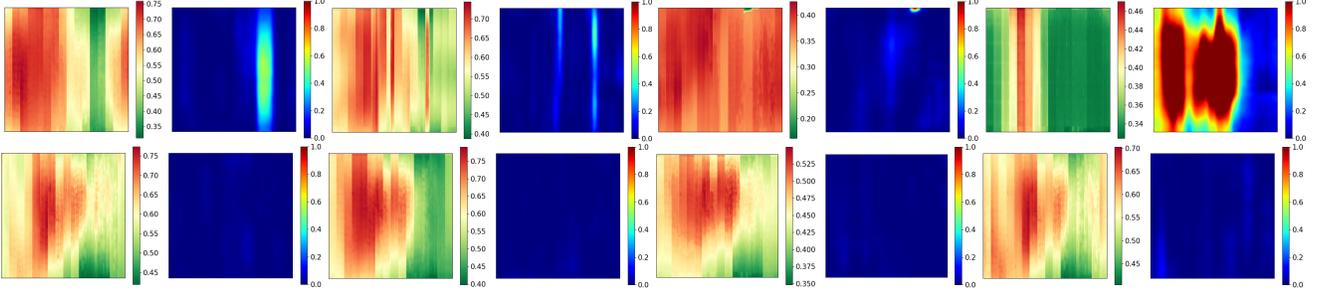


Figure 6: Examples of anomaly maps for anomalous (top) and normal (bottom) images.

similar to our results on the original dataset as reported in Table 1 of the paper. Furthermore, in Appendix B, we provide qualitative evidence to support the validity of the simulated dataset by visualizing generated and original images for a random set of timestamps.

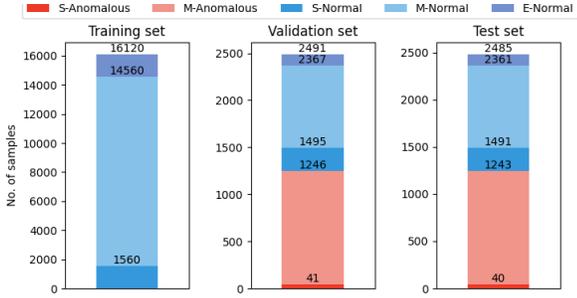


Figure 7: Data split for simulated dataset

Table 4: Anomaly detection performance on simulated data.

Model	AUROC (%)			AUPR (%)		
	[Ts#1]	[Ts#2]	[Ts#3]	[Ts#1]	[Ts#2]	[Ts#3]
Autoencoder	87.97 ( $\pm 4.08$ )	66.34 ( $\pm 2.49$ )	82.00 ( $\pm 1.58$ )	94.04 ( $\pm 1.99$ )	24.72 ( $\pm 6.51$ )	83.46 ( $\pm 1.61$ )
CFlow [12]	83.42 ( $\pm 2.97$ )	51.32 ( $\pm 4.16$ )	70.30 ( $\pm 2.67$ )	90.67 ( $\pm 1.97$ )	10.46 ( $\pm 0.80$ )	69.42 ( $\pm 2.14$ )
DR,EM [38]	98.11 ( $\pm 0.81$ )	61.89 ( $\pm 5.32$ )	89.02 ( $\pm 0.81$ )	99.02 ( $\pm 0.40$ )	25.90 ( $\pm 4.32$ )	88.52 ( $\pm 0.75$ )
FastFlow [37]	97.24 ( $\pm 0.54$ )	52.23 ( $\pm 3.63$ )	87.98 ( $\pm 0.67$ )	98.43 ( $\pm 0.26$ )	9.49 ( $\pm 0.63$ )	87.76 ( $\pm 0.93$ )
PaDiM [6]	97.93 ( $\pm 0.56$ )	56.04 ( $\pm 0.42$ )	88.97 ( $\pm 0.44$ )	98.76 ( $\pm 0.31$ )	9.89 ( $\pm 0.07$ )	88.25 ( $\pm 0.25$ )
PatchCore [26]	98.28 ( $\pm 0.29$ )	66.42 ( $\pm 1.96$ )	92.31 ( $\pm 0.31$ )	98.81 ( $\pm 0.20$ )	21.57 ( $\pm 2.79$ )	92.28 ( $\pm 0.26$ )
Reverse Distillation [7]	75.80 ( $\pm 5.53$ )	57.59 ( $\pm 4.25$ )	65.60 ( $\pm 4.60$ )	86.23 ( $\pm 3.10$ )	12.14 ( $\pm 1.95$ )	63.36 ( $\pm 3.35$ )
ForecastAD	<b>98.73 (<math>\pm 0.64</math>)</b>	<b>98.61 (<math>\pm 0.43</math>)</b>	<b>97.84 (<math>\pm 0.65</math>)</b>	<b>99.30 (<math>\pm 0.33</math>)</b>	<b>88.03 (<math>\pm 3.18</math>)</b>	<b>97.96 (<math>\pm 0.54</math>)</b>

## 5.7 Deployment

We have tested ForecastAD over five months of data from an operational CSP plant. A freshly labelled dataset was curated by initially applying a predefined set of labelling rules, followed by a meticulous review and cleanup of the dataset with guidance from domain experts. The performance metrics of ForecastAD on this labelled set containing 8373 normal and 1,321 abnormal samples are detailed in Table 5. Furthermore, the deployment results are broken down per month over different operating stages. Please note that for some months, we could not compute the performance metrics as there are no anomalous samples present in the dataset. Such cases are marked as “–” in the table. It is important to note that there is variability in this data, such as different stages of operations (starting, ending, and middle) and varying external weather conditions. We can observe that ForecastAD is fairly robust in the detection of anomalies over this period. The actionable insights derived from ForecastAD contribute to the strategic maintenance planning of the CSP plant, thereby enhancing the durability of its equipment.

Table 5: Deployment performance

Month	AUROC (%)			AUPR (%)		
	[Ts#1]	[Ts#2]	[Ts#3]	[Ts#1]	[Ts#2]	[Ts#3]
1	0.89	0.71	0.89	0.66	0.17	0.62
2	0.86	0.94	0.85	0.82	0.70	0.79
3	0.96	0.93	0.95	0.91	0.19	0.87
4	0.91	–	0.91	0.75	–	0.63
5	0.85	–	0.85	0.60	–	0.57
Overall	0.88	0.81	0.88	0.72	0.25	0.69

## 6 CONCLUSION

We address the problem of anomaly detection in irregular sequences of thermal images collected from IR cameras in an operational CSP plant. Extensive analysis of our dataset reveals distinctive temporal characteristics, setting it apart from established AD industrial image benchmark datasets like MVTec [3]. We empirically demonstrate that image-based SOTA AD methods underperform, especially when context is critical for anomaly detection. We also introduce a forecasting-based AD method, ForecastAD, that predicts future thermal images from past sequences and timestamps using a deep sequence model. This method effectively captures specific temporal data features and distinguishes between difficult-to-detect temperature-based anomalies. Experimental results demonstrate the effectiveness of ForecastAD, outperforming existing SOTA methods as measured by AUROC and AUPR. Notably, ForecastAD exhibits significant enhancements in detecting anomalous behaviours, particularly among low-temperature samples. Furthermore, ForecastAD has been successfully deployed, providing critical insights for the maintenance of the CSP plant to our industry partner. For future work, we aim to further study the role of context and sequence lengths in anomaly detection performance. We also aim to extend our model to be more robust to distribution shifts inherent in industrial processes, notably by considering probabilistic forecasting models.

## ACKNOWLEDGMENTS

This work is supported by the research project “Federated Learning and Augmented Reality for Advanced Control Centers”. We thank Thibault GEORGES and Adrien FARINELLE from John Cockerill for helping us understand the dataset along with the associated abnormal behaviours.

## REFERENCES

- [1] Samet Akcay, Dick Ameln, Ashwin Vaidya, Barath Lakshmanan, Nilesh Ahuja, and Utku Genc. 2022. Anomalib: A Deep Learning Library for Anomaly Detection. In *2022 IEEE International Conference on Image Processing (ICIP)*. IEEE, 1706–1710. <https://doi.org/10.1109/ICIP46576.2022.9897283>
- [2] Kilian Batzner, Lars Heckler, and Rebecca König. 2024. EfficientAD: Accurate Visual Anomaly Detection at Millisecond-Level Latencies. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. 128–138. <https://arxiv.org/abs/2303.14535v2>
- [3] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. 2019. MVTEC ad-A comprehensive real-world dataset for unsupervised anomaly detection. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (2019), 9584–9592. <https://doi.org/10.1109/CVPR.2019.00982>
- [4] C. M. Bishop. 1994. Novelty detection and neural network validation. *IEE Proceedings: Vision, Image and Signal Processing* 141, 4 (8 1994), 217–222. <https://doi.org/10.1049/IP-VIS:19941330>
- [5] Christopher M Bishop and Nasser M Nasrabadi. 2006. *Pattern recognition and machine learning*. Vol. 4. Springer.
- [6] Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. 2021. PaDiM: a Patch Distribution Modeling Framework for Anomaly Detection and Localization. In *International Conference on Pattern Recognition*, Vol. 12664 LNCS. Springer Science and Business Media Deutschland GmbH, 475–489. <https://doi.org/10.48550/arxiv.2011.08785>
- [7] Hanqiu Deng and Xingyu Li. 2022. Anomaly detection via reverse distillation from one-class embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9737–9746.
- [8] Ran El-Yaniv and Mordechai Nisenson. 2006. Optimal Single-Class Classification Strategies. In *Advances in Neural Information Processing Systems*, B Schölkopf, J Platt, and T Hoffman (Eds.), Vol. 19. MIT Press. <https://proceedings.neurips.cc/paper/2006/file/ae1d2c2d957a01dcb3f3b39685cdeb4fa-Paper.pdf>
- [9] Joseph Enguehard, Babylon Health, Dan Busbridge, Adam Bozson, Claire Woodcock, and Nils Hammerla. 2020. Neural Temporal Point Processes For Modelling Electronic Health Records. *Proceedings of Machine Learning Research* (7 2020), 85–113. <https://arxiv.org/abs/2007.13794v2>
- [10] Sarah M. Erfani, Sutharshan Rajasegarar, Shanika Karunasekera, and Christopher Leckie. 2016. High-dimensional and large-scale anomaly detection using a linear one-class SVM with deep learning. *Pattern Recognition* 58 (10 2016), 121–134. <https://doi.org/10.1016/j.patcog.2016.03.028>
- [11] Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton Van Den Hengel. 2019. Memorizing Normality to Detect Anomaly: Memory-augmented Deep Autoencoder for Unsupervised Anomaly Detection. *Proceedings of the IEEE International Conference on Computer Vision* 2019-October (4 2019), 1705–1714. <https://doi.org/10.48550/arxiv.1904.02639>
- [12] Denis Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. 2021. CFLOW-AD: Real-Time Unsupervised Anomaly Detection with Localization via Conditional Normalizing Flows. *Proceedings - 2022 IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2022* (7 2021), 1819–1828. <https://doi.org/10.1109/WACV51458.2022.00188>
- [13] Shehroz S. Khan and Michael G. Madden. 2014. One-class classification: taxonomy of study and review of techniques. *The Knowledge Engineering Review* 29, 3 (2014), 345–374. <https://doi.org/10.1017/S026988891300043X>
- [14] Andreas Kind, Marc Stoecklin, and Xenofontas Dimitropoulos. 2009. Histogram-based traffic anomaly detection. *IEEE Transactions on Network and Service Management* 6, 2 (6 2009), 110–121. <https://doi.org/10.1109/TNSM.2009.090604>
- [15] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. 2020. Why Normalizing Flows Fail to Detect Out-of-Distribution Data. *Advances in Neural Information Processing Systems* 2020-December (6 2020). <https://arxiv.org/abs/2006.08545v1>
- [16] Chieh-Hsin Lai, Dongmian Zou, and Gilad Lerman. 2019. Robust Subspace Recovery Layer for Unsupervised Anomaly Detection. *arXiv preprint* (3 2019). <https://doi.org/10.48550/arxiv.1904.00152>
- [17] Sungwook Lee, Seunghyun Lee, and Byung Cheol Song. 2022. CFA: Coupled-hypersphere-based Feature Adaptation for Target-Oriented Anomaly Localization. *IEEE Access* 10 (6 2022), 78446–78454. <https://doi.org/10.1109/ACCESS.2022.3193699>
- [18] Leland McInnes, John Healy, and James Melville. 2018. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. [arXiv:1802.03426](https://arxiv.org/abs/1802.03426) [stat.ML]
- [19] T. Minter. 1975. Single-Class Classification. *LARS Symposia* (1 1975). [https://docs.lib.purdue.edu/lars\\_symp/54](https://docs.lib.purdue.edu/lars_symp/54)
- [20] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. 2018. Do Deep Generative Models Know What They Don't Know? *7th International Conference on Learning Representations, ICLR 2019* (10 2018). <https://doi.org/10.48550/arxiv.1810.09136>
- [21] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, and Balaji Lakshminarayanan. 2019. Detecting Out-of-Distribution Inputs to Deep Generative Models Using Typicality. In *4th Workshop on Bayesian Deep Learning (NeurIPS 2019)*. <https://arxiv.org/abs/1906.02994v2>
- [22] Duc Tam Nguyen, Zhongyu Lou, Michael Klar, and Thomas Brox. 2019. Anomaly Detection With Multiple-Hypotheses Predictions. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, 4800–4809. <https://proceedings.mlr.press/v97/nguyen19b.html>
- [23] Guansong Pang, Chunhua Shen, Longbing Cao, and Anton van den Hengel. 2020. Deep Learning for Anomaly Detection: A Review. *Comput. Surveys* 54, 2 (7 2020). <https://doi.org/10.1145/3439950>
- [24] Emanuel Parzen. 1962. On Estimation of a Probability Density Function and Mode. *The Annals of Mathematical Statistics* 33, 3 (1962), 1065–1076. <http://www.jstor.org/stable/2237880>
- [25] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury Google, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf Xamla, Edward Yang, Zach Devito, Martin Raison Naba, Alykhan Tejani, Sasank Chilamkurthy, Qure Ai, Benoit Steiner, Lu Fang Facebook, Junjie Bai Facebook, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *Advances in Neural Information Processing Systems* 32 (2019).
- [26] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. 2022. Towards Total Recall in Industrial Anomaly Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14318–14328. <https://doi.org/10.48550/arxiv.2106.08265>
- [27] Lukas Ruff, Jacob R. Kauffmann, Robert A. Vandermeulen, Gregoire Montavon, Wojciech Samek, Marius Kloft, Thomas G. Dietterich, and Klaus Robert Müller. 2021. A Unifying Review of Deep and Shallow Anomaly Detection. *Proc. IEEE* 109, 5 (5 2021), 756–795. <https://doi.org/10.1109/JPROC.2021.3052449>
- [28] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. 2018. Deep One-Class Classification. In *International Conference on Machine Learning*. PMLR, 4393–4402. <https://proceedings.mlr.press/v80/ruff18a.html>
- [29] Bernhard Schölkopf, John C. Platt, John Shawe-Taylor, Alex J. Smola, and Robert C. Williamson. 2001. Estimating the Support of a High-Dimensional Distribution. *Neural Computation* 13, 7 (2001), 1443–1471. <https://doi.org/10.1162/089976601750264965>
- [30] I Sobel. 1968. An isotropic 3 × 3 image gradient operator, presentation at Stanford Artificial Intelligence Project (SAIL).
- [31] Wuqin Tang, Qiang Yang, Kuixiang Xiong, and Wenjun Yan. 2020. Deep learning based automatic defect identification of photovoltaic module using electroluminescence images. *Solar Energy* 201 (5 2020), 453–460. <https://doi.org/10.1016/j.solener.2020.03.049>
- [32] David Tax. 2001. *One-Class Classification; Concept-Learning In The Absence Of Counter-Examples*. Ph.D. Dissertation. Technische Universiteit Delft, Delft.
- [33] David M.J. Tax and Robert P.W. Duin. 1999. Support vector domain description. *Pattern Recognition Letters* 20, 11-13 (11 1999), 1191–1199. [https://doi.org/10.1016/S0167-8655\(99\)00087-2](https://doi.org/10.1016/S0167-8655(99)00087-2)
- [34] David M.J. Tax and Robert P.W. Duin. 2004. Support Vector Data Description. *Machine Learning* 54, 1 (2004), 45–66. <https://doi.org/10.1023/B:MACH.0000008084.60811.49>
- [35] Ashish Vaswani, Google Brain, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *Advances in Neural Information Processing Systems*, Vol. 30. 5998–6008.
- [36] Xudong Yan, Huaidong Zhang, Xuemiao Xu, Xiaowei Hu, and Pheng-Ann Heng. 2021. Learning Semantic Context from Normal Samples for Unsupervised Anomaly Detection. *Proceedings of the AAAI Conference on Artificial Intelligence* 35, 4 (5 2021), 3110–3118. <https://ojs.aaai.org/index.php/AAAI/article/view/16420>
- [37] Jiawei Yu, Ye Zheng, Xiang Wang, Wei Li, Yushuang Wu, Rui Zhao, and Liwei Wu. 2021. FastFlow: Unsupervised Anomaly Detection and Localization via 2D Normalizing Flows. *arXiv preprint arXiv:2111.07677* (11 2021). <https://doi.org/10.48550/arxiv.2111.07677>
- [38] Vitjan Zavrtnik, Matej Kristan, and Danijel Škočaj. 2021. DRÆM - A discriminatively trained reconstruction embedding for surface anomaly detection. *Proceedings of the IEEE International Conference on Computer Vision* (8 2021), 8310–8319. <https://doi.org/10.1109/ICCV48922.2021.00822>
- [39] Houssam Zenati, Chuan Sheng Foo, Bruno Lecouat, Gaurav Manek, and Vijay Ramaseshan Chandrasekhar. 2018. Efficient GAN-Based Anomaly Detection. *arXiv preprint* (2 2018). <https://doi.org/10.48550/arxiv.1802.06222>
- [40] H. L. Zhang, J. Baeyens, J. Degrève, and G. Caceres. 2013. Concentrated solar power plants: Review and design methodology. *Renewable and Sustainable Energy Reviews* 22 (6 2013), 466–481. <https://doi.org/10.1016/j.rser.2013.01.032>
- [41] Jie Zhang, Masanori Suganuma, and Takayuki Okatani. 2023. Contextual Affinity Distillation for Image Anomaly Detection. *arXiv preprint arXiv:2307.03101* (7 2023). <https://arxiv.org/abs/2307.03101v1>

## A NETWORK ARCHITECTURE

ForecastAD is implemented using the PyTorch framework [25]. The architecture of the encoder  $\phi_e$  and  $\phi_d$  is presented in Tables 6 and 7, respectively. The input dimension of the encoder network is  $(3 \times 256 \times 256)$ . The period of the sinusoidal encoding is 1000.

**Table 6: Encoder architecture**

Layer
Conv2d-1(3, 32, 5, padding=2, stride=2)
BatchNorm2d-1(32, eps=1e-04)
MaxPool2d-1(2,2)
Conv2d-2(32, 64, 5, padding=2, stride=2)
BatchNorm2d-2(64, eps=1e-04)
MaxPool2d-2(2,2)
Conv2d-3(64, 128, 5, padding=2, stride=2)
BatchNorm2d-3(128, eps=1e-04)
MaxPool2d-3(2,2)
Conv2d-4(128, 128, 5, padding=2, stride=2)
BatchNorm2d-4(128, eps=1e-04)
MaxPool2d-4(2,2)
Linear-1(128, 128, padding=2, stride=2)
BatchNorm2d-1(128, eps=1e-04)

**Table 7: Decoder architecture**

Layer
ConvTranspose2d-1(128, 64, 5, padding=2)
BatchNorm2d-1(64, eps=1e-04)
ConvTranspose2d-2(64, 64, 5, padding=2)
BatchNorm2d-2(64, eps=1e-04)
ConvTranspose2d-3(64, 32, 5, padding=2)
BatchNorm2d-3(32, eps=1e-04)
ConvTranspose2d-4(32, 3, 5, padding=2)

## B DATA GENERATION

We generate a public dataset mirroring the properties of our private dataset using a variational autoencoder (VAE). For daily sequence generation, we condition the VAE on the time of day, class (Positive, Negative, Unlabelled), and phase (Start, Middle, End). A standard convolutional neural network (CNN) with three convolution blocks with an increasing number of feature maps is used to encode images, and a multilayer perceptron (MLP) is used to encode conditioning variables. Then, a two-layer LSTM network produces the context embedding. Finally, a deconvolutional CNN, mirroring the encoder’s architecture, reconstructs the original image from a latent vector sampled from the latent distribution. We use a multivariate Gaussian distribution with a diagonal covariance matrix as our latent distribution. We train the VAE over 20 epochs using Adam optimizer, a learning rate of 1e-4, and a batch size of 16. Figure 8

allows us to compare the generated images with the original images from the private dataset. From the figure, we can observe that the VAE can generate normal images and a diverse set of anomalies which are visually similar to images in the original dataset.

## C SENSITIVITY TO DATA LABELLING

Some images in the deployment set are wrongly labelled as normal. Such inconsistencies result in a distribution shift between the labelled training and deployment sets, leading to significant degradation of model performance. Thus, we cleaned the deployment set by calculating the distance of each labelled normal image  $x_i$  in the deployment set to the images in the training set  $\mathcal{D}_N$ . For this, we first obtain the context embedding  $c_i$  for each image  $x_i$  using the context encoder of ForecastAD. We then compute the distance between context embeddings as  $\xi_i = \min_{j=1, \dots, |\mathcal{D}_N|} \|c_i - c_j\|_2$ . Images from the deployment set for which the distance  $d_i$  exceeds a predetermined threshold are removed. Figure 9 shows the UMAP projection [18] of the context embeddings corresponding to the samples in the training set  $\mathcal{D}_N$  and the samples from the deployment set that are removed during the cleaning process.

## D ADDITIONAL RESULTS

**Simple baselines.** We also consider four simple baselines which compute thresholds based on the features extracted from the samples in the validation set to detect anomalies. The extracted features can be seen as the anomaly scores. In this study, we consider the following features:

- Time of day: the time when the image is recorded.
- Negative mean: negative of the average pixel value of the image.
- Negative max: negative of the maximum pixel value of the image.
- Negative Std: negative of the standard deviation of the pixel value in the image.

**Threshold selection.** Given the anomaly scores assigned to the labelled validation samples, we explore two different threshold selection approaches:

- **F1-score.** We compute F1 scores by considering various thresholds. Given the F1 scores, we select the optimal threshold  $\lambda_f$  as the one that corresponds to the maximum F1 score.
- **G-Mean.** Similar to  $\lambda_f$ , we calculate the G-Mean value by applying multiple thresholds. The optimal threshold  $\lambda_g$  corresponds to the threshold where we obtain the highest G-Mean value. G-Mean is the geometric mean of specificity and recall.

$$\text{G-Mean} = \sqrt{\text{Specificity} \cdot \text{Recall}} = \sqrt{(1 - \text{FPR}) \cdot \text{TPR}}$$

Based on the selected threshold, Accuracy and F1-score are reported for [Tr#1] and [Tr#2] over three test setups in Table 8.

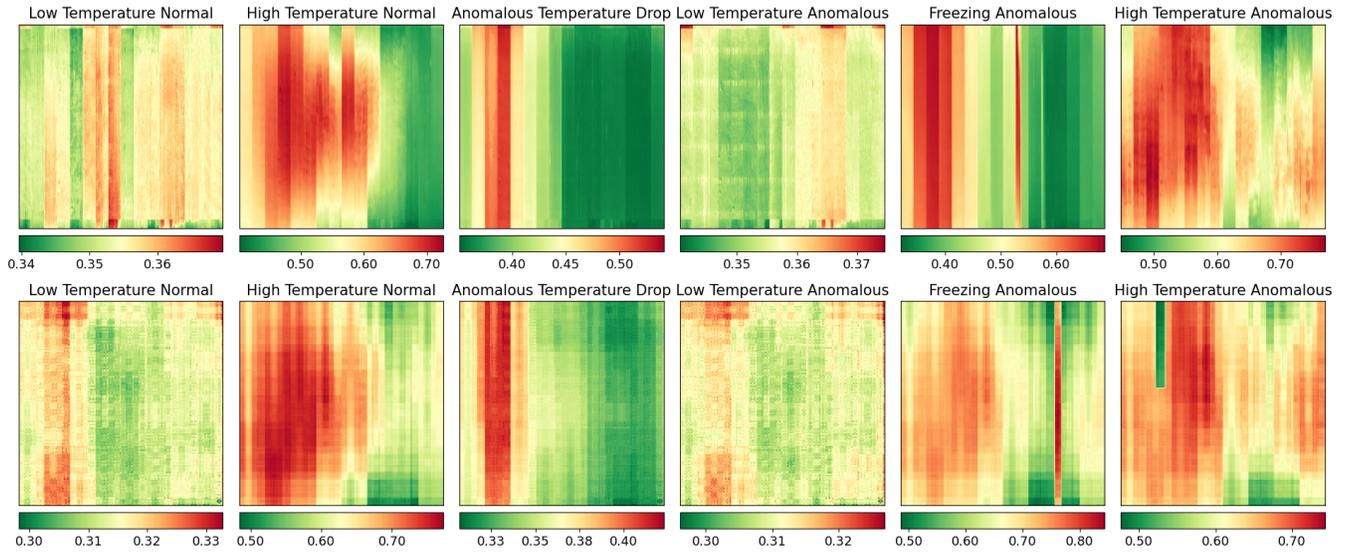


Figure 8: Examples of different types of images in Original (top) and simulated (bottom) dataset

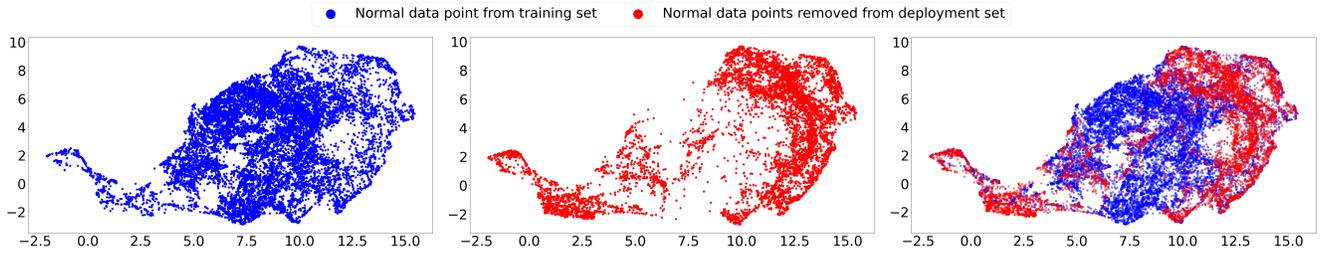


Figure 9: UMAP plot of training set normal images (left), removed deployment normal images (middle) and the both (right).

Table 8: Extended Anomaly Detection Performance. Style: best and second best

Train Setting	Model	AUROC (%)			Accuracy (%)						F1-score (%)					
		[Ts#1]	[Ts#2]	[Ts#3]	$\lambda_f$	[Ts#1]	$\lambda_g$	$\lambda_f$	[Ts#2]	$\lambda_g$	$\lambda_f$	[Ts#3]	$\lambda_g$	$\lambda_f$	[Ts#3]	$\lambda_g$
[Tr#1]	Time of day	86.99 (± 0.00)	41.44 (± 0.00)	79.97 (± 0.00)	83.03 (± 0.00)	82.86 (± 0.00)	61.60 (± 0.00)	61.60 (± 0.00)	79.55 (± 0.00)	79.41 (± 0.00)	85.34 (± 0.00)	85.17 (± 0.00)	0.00 (± 0.00)	0.00 (± 0.00)	80.20 (± 0.00)	80.02 (± 0.00)
	Negative Mean	81.67 (± 0.00)	46.73 (± 0.00)	71.51 (± 0.00)	70.49 (± 0.00)	76.43 (± 0.00)	10.89 (± 0.00)	13.47 (± 0.00)	60.83 (± 0.00)	66.22 (± 0.00)	74.40 (± 0.00)	76.91 (± 0.00)	14.33 (± 0.00)	12.72 (± 0.00)	65.46 (± 0.00)	66.76 (± 0.00)
	Negative STD	79.31 (± 0.00)	40.76 (± 0.00)	70.35 (± 0.00)	73.60 (± 0.00)	73.60 (± 0.00)	34.10 (± 0.00)	34.10 (± 0.00)	67.19 (± 0.00)	67.19 (± 0.00)	72.49 (± 0.00)	72.49 (± 0.00)	14.18 (± 0.00)	14.18 (± 0.00)	64.66 (± 0.00)	64.66 (± 0.00)
	Negative Max	77.61 (± 0.00)	44.85 (± 0.00)	68.64 (± 0.00)	75.10 (± 0.00)	75.37 (± 0.00)	14.33 (± 0.00)	14.90 (± 0.00)	65.24 (± 0.00)	65.57 (± 0.00)	75.29 (± 0.00)	75.47 (± 0.00)	11.80 (± 0.00)	11.87 (± 0.00)	65.31 (± 0.00)	65.49 (± 0.00)
	Autoencoder	98.05 (± 0.74)	46.43 (± 1.61)	87.87 (± 0.26)	94.50 (± 1.22)	94.49 (± 1.22)	36.96 (± 3.72)	37.19 (± 3.68)	85.17 (± 0.46)	85.20 (± 0.47)	95.27 (± 1.06)	95.26 (± 1.06)	12.58 (± 0.64)	12.62 (± 0.63)	86.45 (± 0.54)	86.46 (± 0.55)
	CFlow [12]	94.68 (± 1.26)	39.99 (± 2.33)	82.91 (± 1.08)	87.07 (± 1.67)	86.98 (± 1.77)	31.69 (± 1.70)	32.84 (± 1.40)	78.09 (± 1.45)	78.20 (± 1.47)	88.50 (± 1.50)	88.31 (± 1.64)	12.75 (± 0.30)	12.55 (± 0.43)	79.53 (± 1.37)	79.42 (± 1.49)
	Deep SVDD (one-class) [28]	52.76 (± 8.71)	49.22 (± 2.66)	51.85 (± 6.15)	61.35 (± 1.95)	58.34 (± 5.62)	20.00 (± 7.87)	54.10 (± 8.04)	54.65 (± 2.86)	57.65 (± 3.52)	70.79 (± 1.77)	54.97 (± 6.57)	13.79 (± 0.87)	13.33 (± 1.53)	64.12 (± 1.60)	50.54 (± 5.49)
	Deep SVDD (soft-boundary) [28]	30.22 (± 6.77)	49.68 (± 4.26)	35.45 (± 3.84)	58.18 (± 0.04)	43.04 (± 5.65)	7.79 (± 0.06)	68.42 (± 9.38)	50.01 (± 0.03)	47.16 (± 3.99)	73.56 (± 0.03)	33.14 (± 7.34)	14.37 (± 0.01)	11.66 (± 1.73)	66.67 (± 0.02)	31.14 (± 6.59)
	DREM [38]	97.70 (± 0.77)	40.48 (± 2.15)	87.38 (± 0.61)	91.70 (± 1.08)	91.81 (± 1.05)	30.49 (± 3.85)	30.95 (± 3.68)	81.78 (± 1.39)	81.94 (± 1.35)	92.97 (± 0.83)	93.04 (± 0.81)	12.14 (± 0.92)	11.94 (± 1.07)	83.66 (± 0.89)	83.75 (± 0.87)
	FastFlow [37]	95.83 (± 0.03)	47.32 (± 0.29)	91.36 (± 0.25)	97.39 (± 0.29)	97.38 (± 0.29)	42.12 (± 1.27)	42.18 (± 1.29)	88.43 (± 0.38)	88.43 (± 0.38)	97.72 (± 0.26)	97.71 (± 0.26)	13.22 (± 0.26)	13.24 (± 0.25)	89.17 (± 0.35)	89.17 (± 0.35)
	PaDiM [6]	99.85 (± 0.02)	49.86 (± 0.47)	91.23 (± 0.10)	96.92 (± 0.72)	96.45 (± 0.58)	43.44 (± 1.91)	44.76 (± 1.28)	88.24 (± 0.34)	88.07 (± 0.33)	97.28 (± 0.65)	96.86 (± 0.53)	13.76 (± 0.29)	13.46 (± 0.10)	88.92 (± 0.43)	88.66 (± 0.39)
	PatchCore [26]	99.23 (± 0.08)	50.58 (± 0.37)	89.04 (± 0.30)	95.50 (± 0.25)	95.52 (± 0.25)	31.29 (± 1.54)	31.46 (± 1.65)	85.08 (± 0.41)	85.13 (± 0.43)	96.21 (± 0.22)	96.23 (± 0.21)	15.16 (± 0.16)	15.31 (± 0.19)	86.77 (± 0.34)	86.81 (± 0.35)
Reverse Distillation [7]	93.88 (± 1.13)	41.31 (± 2.19)	84.61 (± 1.54)	87.01 (± 1.68)	86.66 (± 1.52)	35.01 (± 2.90)	39.03 (± 2.51)	78.58 (± 1.64)	78.93 (± 1.58)	89.27 (± 1.44)	88.66 (± 1.31)	12.91 (± 0.36)	12.81 (± 0.74)	81.17 (± 1.43)	80.86 (± 1.40)	
ForecastAD	99.86 (± 0.05)	46.22 (± 1.06)	89.89 (± 0.35)	97.73 (± 0.34)	97.65 (± 0.27)	36.10 (± 1.19)	36.62 (± 1.35)	87.73 (± 0.29)	87.75 (± 0.26)	98.04 (± 0.29)	97.97 (± 0.24)	14.02 (± 0.50)	14.13 (± 0.50)	88.76 (± 0.25)	88.75 (± 0.22)	
[Tr#2]	Time of day	86.99 (± 0.00)	41.44 (± 0.00)	79.97 (± 0.00)	83.03 (± 0.00)	82.86 (± 0.00)	61.60 (± 0.00)	61.60 (± 0.00)	79.55 (± 0.00)	79.41 (± 0.00)	85.34 (± 0.00)	85.17 (± 0.00)	0.00 (± 0.00)	0.00 (± 0.00)	80.20 (± 0.00)	80.02 (± 0.00)
	Negative Mean	81.67 (± 0.00)	46.73 (± 0.00)	71.51 (± 0.00)	71.16 (± 0.00)	77.45 (± 0.00)	9.46 (± 0.00)	43.27 (± 0.00)	61.15 (± 0.00)	71.89 (± 0.00)	76.17 (± 0.00)	75.96 (± 0.00)	14.59 (± 0.00)	13.91 (± 0.00)	67.24 (± 0.00)	68.54 (± 0.00)
	Negative STD	79.31 (± 0.00)	40.76 (± 0.00)	70.35 (± 0.00)	66.00 (± 0.00)	73.60 (± 0.00)	9.17 (± 0.00)	45.56 (± 0.00)	56.78 (± 0.00)	69.05 (± 0.00)	74.45 (± 0.00)	70.98 (± 0.00)	12.67 (± 0.00)	12.84 (± 0.00)	66.33 (± 0.00)	64.16 (± 0.00)
	Negative Max	77.61 (± 0.00)	44.85 (± 0.00)	68.64 (± 0.00)	78.37 (± 0.00)	74.71 (± 0.00)	33.81 (± 0.00)	49.86 (± 0.00)	71.14 (± 0.00)	70.68 (± 0.00)	77.19 (± 0.00)	72.26 (± 0.00)	14.13 (± 0.00)	14.63 (± 0.00)	68.62 (± 0.00)	65.87 (± 0.00)
	Autoencoder	96.67 (± 0.77)	45.92 (± 2.47)	85.45 (± 1.18)	93.09 (± 1.21)	93.16 (± 1.40)	26.02 (± 3.26)	28.83 (± 2.91)	82.21 (± 1.27)	82.72 (± 1.39)	94.15 (± 1.02)	94.10 (± 1.22)	13.99 (± 0.36)	13.69 (± 0.52)	84.26 (± 1.01)	84.40 (± 1.22)
	CFlow [12]	84.91 (± 2.72)	42.90 (± 2.71)	77.38 (± 2.98)	76.78 (± 3.03)	77.15 (± 2.29)	29.86 (± 2.77)	39.77 (± 4.33)	69.17 (± 2.77)	71.09 (± 2.36)	81.93 (± 2.03)	80.16 (± 1.93)	12.83 (± 0.93)	12.35 (± 1.22)	74.43 (± 1.83)	73.19 (± 1.94)
	Deep SVDD (one-class) [28]	45.93 (± 5.07)	52.24 (± 2.24)	46.56 (± 4.80)	58.11 (± 0.07)	54.04 (± 4.88)	7.85 (± 0.91)	61.60 (± 2.26)	49.96 (± 0.06)	55.27 (± 4.34)	73.51 (± 0.06)	47.40 (± 4.76)	14.38 (± 0.01)	12.36 (± 1.67)	66.63 (± 0.06)	44.36 (± 4.31)
	Deep SVDD (soft-boundary) [28]	29.21 (± 15.05)	52.48 (± 2.83)	35.44 (± 11.72)	62.34 (± 4.10)	38.41 (± 11.14)	17.65 (± 9.91)	65.04 (± 4.93)	55.09 (± 5.05)	42.73 (± 9.28)	74.76 (± 1.16)	46.42 (± 8.16)	14.06 (± 0.30)	14.09 (± 0.89)	68.12 (± 1.42)	44.14 (± 7.40)
	DREM [38]	93.52 (± 0.52)	40.51 (± 1.33)	85.71 (± 0.78)	85.00 (± 1.41)	87.14 (± 0.98)	28.54 (± 4.62)	35.70 (± 3.60)	75.85 (± 1.66)	78.80 (± 0.98)	88.07 (± 1.00)	88.96 (± 0.88)	11.84 (± 0.36)	11.40 (± 0.36)	79.62 (± 1.04)	80.67 (± 0.85)
	FastFlow [37]	92.38 (± 0.72)	52.51 (± 1.09)	89.92 (± 0.68)	88.04 (± 0.78)	87.83 (± 0.87)	56.62 (± 3.06)	59.03 (± 2.49)	82.95 (± 1.02)	83.16 (± 1.05)	90.15 (± 0.61)	89.86 (± 0.68)	13.39 (± 0.46)	13.14 (± 0.38)	84.49 (± 0.78)	84.45 (± 0.84)
	PaDiM [6]	95.99 (± 0.37)	58.14 (± 1.00)	92.28 (± 0.32)	88.77 (± 0.26)	87.58 (± 0.26)	46.25 (± 3.26)	65.73 (± 1.90)	81.88 (± 0.48)	84.03 (± 0.38)	90.92 (± 0.24)	88.87 (± 0.27)	14.81 (± 0.39)	17.14 (± 1.42)	84.09 (± 0.21)	84.07 (± 0.27)
	PatchCore [26]	96.78 (± 0.57)	60.15 (± 1.82)	91.38 (± 0.42)	90.67 (± 1.24)	90.34 (± 1.30)	55.82 (± 4.20)	56.85 (± 4.14)	85.02 (± 0.61)	84.91 (± 0.63)	91.84 (± 1.20)	91.49 (± 1.26)	17.86 (± 0.88)	17.86 (± 0.97)	85.71 (± 0.79)	85.49 (± 0.83)
Reverse Distillation [7]	87.19 (± 0.99)	57.22 (± 5.77)	84.04 (± 1.64)	84.39 (± 1.15)	82.48 (± 1.14)	44.41 (± 6.44)	57.02 (± 5.75)	77.91 (± 1.16)	78.36 (± 0.97)	87.32 (± 0.91)	84.92 (± 1.12)	14.91 (± 1.50)	15.50 (± 1.41)	80.60 (± 0.71)	79.53 (± 0.91)	
ForecastAD	94.78 (± 1.09)	85.81 (± 1.23)	92.53 (± 0.81)	88.82 (± 1.33)	88.75 (± 1.28)	76.68 (± 2.13)	78.40 (± 1.36)	86.85 (± 0.89)	87.07 (± 0.98)	90.03 (± 1.24)	89.88 (± 1.19)	35.40 (± 1.59)	36.29 (± 1.40)	86.83 (± 1.01)	86.89 (± 1.06)	