

Predicting Machine Failures with Reinforcement Learning

DALTON SCHUTTE

Georgia Institute of Technology

Abstract

We investigate methods from statistical process control, deep learning, and reinforcement learning for predicting when two machines will fail. Each machine is fitted with sensors that report data at semi-regular intervals for each time step. Our results suggest that there is promise in reinforcement learning agents trained using proximal policy optimization. Our findings also suggest that there is value in daily generated T^2 control charts. The state of the art time series transformer model could not be trained so we do not make an assessment of its abilities, however, given the nature of the model it may be worth additional troubleshooting to evaluate its potential. Ultimately, the control charts are likely the best performance to cost option due to the lower compute, expertise, and data availability they require.

1 Introduction

As time progresses, machines gradually wear down and will fail. When a machine fails, it is likely to require maintenance or repair before it can return to its normal function. Depending on the application, failure may not be an issue if the process is low-stakes or has redundancies or backups or catastrophic if the process is of critical importance. In most cases, machine failure may result in lost revenue as production is halted or goods are damaged.

Knowing when a machine is likely to fail can reduce the impact on down- and upstream processes by allowing necessary preparations to be taken before the machine is stopped for maintenance. This is preferable to the alternative where the failure remains unaddressed until it is noticed. This can add to down times and makes preparing for the downtime rushed.

1.1 Background

1.1.1 Statistical Process Control

Statistical Process Control (SPC), or, Statistical Quality Control, is a set of techniques that allow one to track the results of a process over time [1]. The general procedure is:

1. Determine which time period will be used for establishing in-control markers
2. Gather observations from this period
3. Calculate the center line
4. Calculate any control limits (such as 2 standard deviations above/below the center line)

5. Plot new observations against these limits
6. Evaluate any patterns that emerge

There are different types of charts depending on what observations are available, the nature of the process, etc. Typically, a single characteristic is measured (e.g. weight for canned goods, length of extruded metal, etc.) per chart for simplicity. Tracking multiple variables is complex and requires special methodologies.

1.1.2 Deep Learning

Deep learning is a procedure where deep neural networks are trained on a dataset using an optimization algorithm in an iterative manner. There are many flavors of neural network, but two that are particularly well suited to sequential data are the Long-Short Term Memory (LSTM) network [2] and the transformer [3].

Transformer-based architectures have been at the forefront of NLP research in recent years. While they have demonstrated excellent performance on some generative language tasks, their ability to process long-term dependencies in sequential data and the ability to pre-train models once and fine-tune for multiple other tasks suggests they could be well suited to applications in time series analysis. Indeed there has been some research on this front. Recently Goswami et al. [4] collected a large amount of time series data to pre-train a family of models called MOMENT. This family of foundational time series models is pre-trained on time series from a variety of domains, with varying time horizons, and for various tasks.

1.1.3 Reinforcement Learning

Reinforcement learning (RL) is a special subset of ML/AI focused on choosing actions in some environment [5]. This process begins with some "agent" that will use some reinforcement learning algorithm to interact with an environment which, in turn, provides a reward signal that is used to help train the agent. The ultimate goal is for the agent to learn an optimal policy, π , that can be used to choose actions given a state. Perhaps one of the most famous applications of reinforcement learning is in games. Many board games, such as Go [6], have a huge number of discrete game states that one cannot simply find the best play with brute force.

Two of the, perhaps, most well-known techniques are Proximal Policy Optimization (PPO) [7] and Neural Guided Monte Carlo Tree Search (MCTS) [8]. The former was used to train the OpenAI 5, a collection of AI agents trained to play the video game Dota 2 [9] and the latter to train AlphaGo, which beat the then Go world champion in a set of 4-1.

1.2 Research Questions

The states for many machine sensors are not discrete, but rather, take on values in \mathbb{R} . The T² and, most, deep learning models handle this situation without issue. There are plenty of RL algorithms that can solve continuous control tasks [10, 11]. However, we wish to examine the efficacy of RL for stopping machines that are about to fail. We will seek to answer the following questions:

1. How well can MCTS and PPO perform at this task?
2. How do these algorithms compare to traditional methods from SPC?
3. How do these algorithms compare to models, such as MOMENT, designed for sequential data?
4. How do the RL methods compare to each other?

2 Methods

2.1 Data

We are using data from two machines, a blood refrigerator and a nitrogen generator [12].

A blood refrigerator is a machine designed to store blood safely at proper temperatures. The general composition of a unit is the compressor to pump coolant throughout, the condenser to remove heat from the collant and condense it, the evaporator which evaporates the coolant to cool the interior of the unit, an expansion valve to regulate the flow of coolant, and tubing to move the coolant throughout.

A nitrogen generator is a machine that separates nitrogen from other gases in compressed air. This unit typically includes an air compressor, carbon sieves for filtering, absorption vessels to collect nitrogen, and towers to increase production.

2.1.1 Analysis

The data includes a mix of binary and continuous variables with a binary target variable. The class imbalance is heavily in favor of the 'normal' class ($PW_0.5h=0$), which indicates the machines are operating within expectations the majority of the time, a desireable phenomenon. The ratio of normal to failed states is roughly 99:1 for the blood refrigerator and 49:1 for the nitrogen generator. In each dataset, there are several days where there is no failure state, meaning the machine operated without issue the entire day.

The measurement intervals are spaced roughly 1 minute apart for the nitrogen generator and about 34 seconds for the blood refrigerator. The label 'PW_0.5', using the notation and convention from the paper that provided this dataset [12], represents the Prediction Window (PW) of 30 minutes (0.5 hours). This means that in the 30 minute interval beginning at a time stamp t , the machine is operating outside of parameters for at least one of those time stamps. The other convention from that paper is the Reading Window over some interval. The Reading Window (RW) is the collection of sequential time stamps used as input to a model. In our experiments, for example, we use a RW of 20 minutes where we collect the minimum number of sequence of time steps necessary to span 20 minutes of time. In Pincioli et al. [12], they explore a wide combination of PW and RW but found that most of their models performed best with a 20 minute RW, which is why that is the window size we chose to use.

Figures 1 and 2 show some examples of the time series data in our dataset. Initially, the blood refrigerator dataset had 16 sensors providing output and the nitrogen generator had 7. We removed all the output from any sensor whose standard deviation was zero over the entire dataset. This left 12 variables for the blood refrigerator and 4 for the nitrogen generator.

The time series for the features in the blood refrigerator (figure 1) highlight that there is, intuitively, some change occurring in the machine that is leading to these sensor readings. Particularly, the product temperature base, evaporator temperature, and power supply have noticeably different patterns in the red failure region than in the period before.

	Blood Refrigerator		Nitrogen Generator	
	Train	Test	Train	Test
Days	25	27	29	8
Timesteps	60166	65763	40354	11162
Stops (1)	704	642	810	242

Table 1: Summary statistics of our datasets

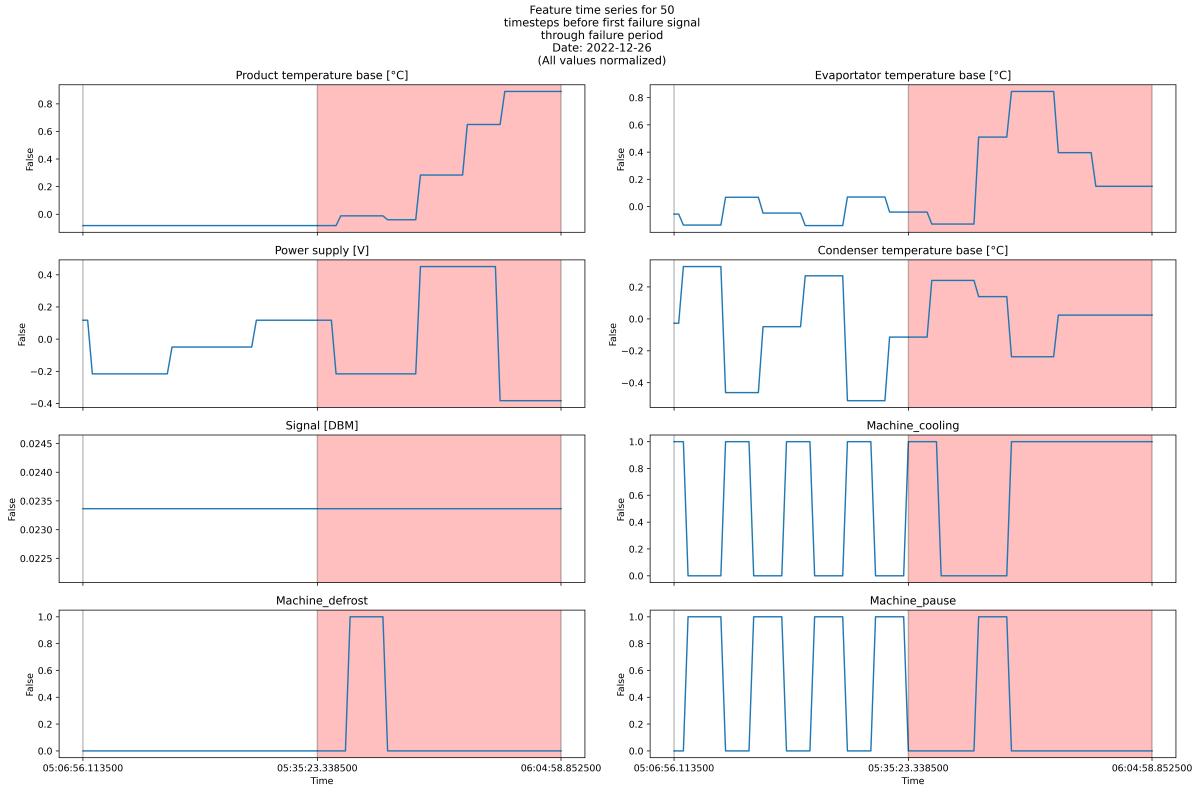


Figure 1: Plots of variables around the point of failure for the blood refrigerator. Red regions are where regions where the data label indicates failure (1)

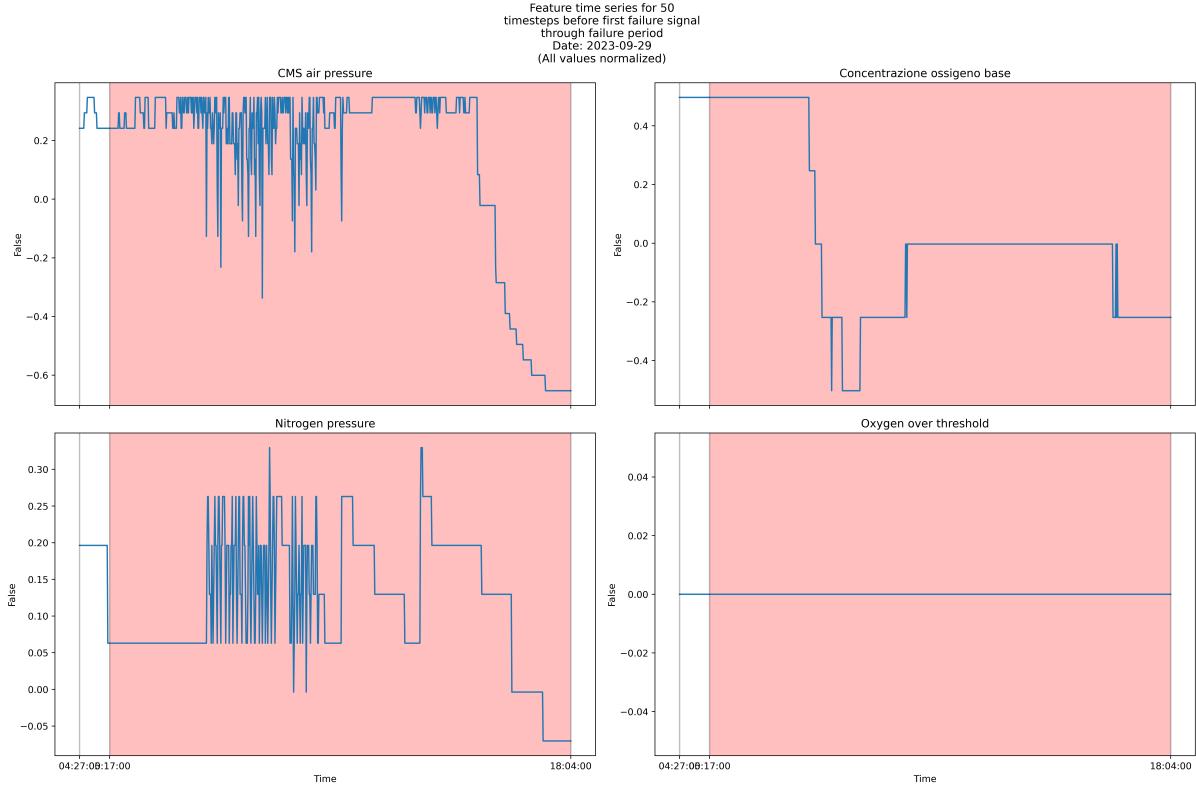


Figure 2: Plots of variables around the point of failure for the nitrogen generator. Red regions are where regions where the data label indicates failure (1)

The nitrogen generator (figure 2) shows an instance where the failure happens early in the observation period, only a few time steps into the beginning of the period. This poses a challenge for the methods as there is very little data for them to work from before having to determine if the machine is on a trajectory to failing. However, this may be the case in real operating conditions and so is an important case to keep in our dataset.

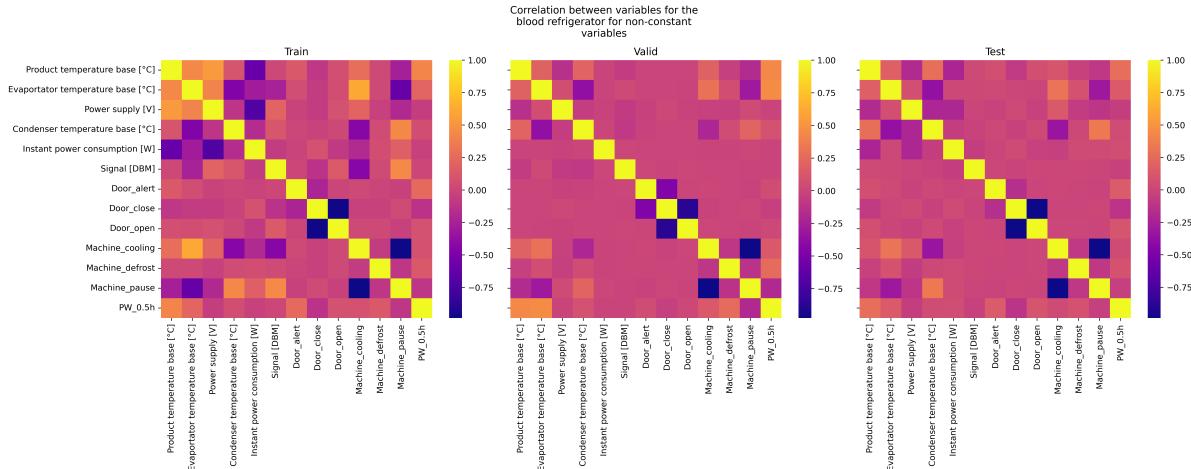


Figure 3: Correlation between variables in the blood refrigerator dataset splits. The label column is 'PW_0.5h'.

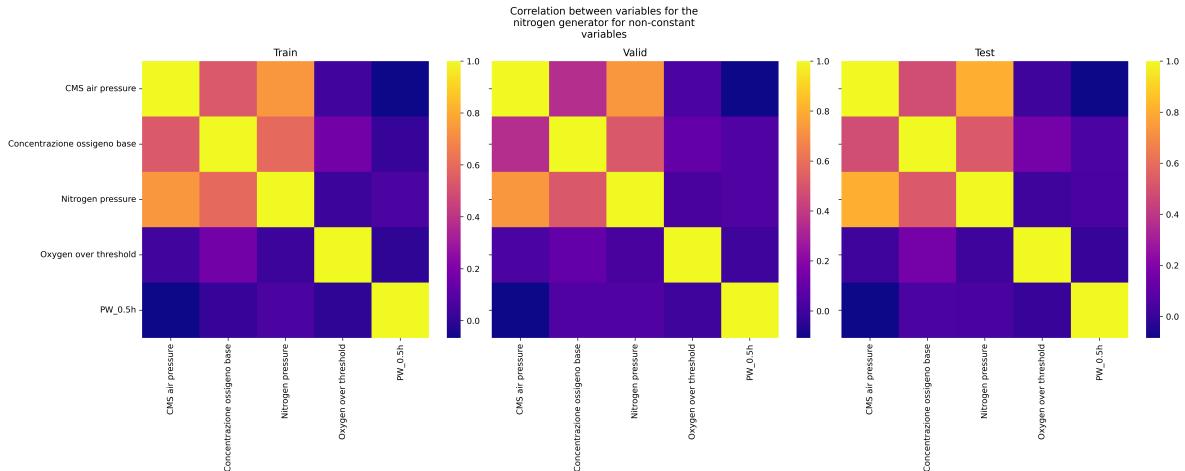


Figure 4: Correlation between variables in the nitrogen generator dataset splits. The label column is 'PW_0.5h'.

Correlations Upon examination of the correlations, shown in figures 3 and 4, between variables for each dataset, there were no instances where any one variable was significantly correlated with the target variable, 'PW_0.5h'. We felt no need to drop additional features based on these observations.

Feature Distributions Figures 5 and 6 show the distributions of the continuous variables for both machines. The blood refrigerator has, good representation of most variables in the training set that appear in the validation and test splits.

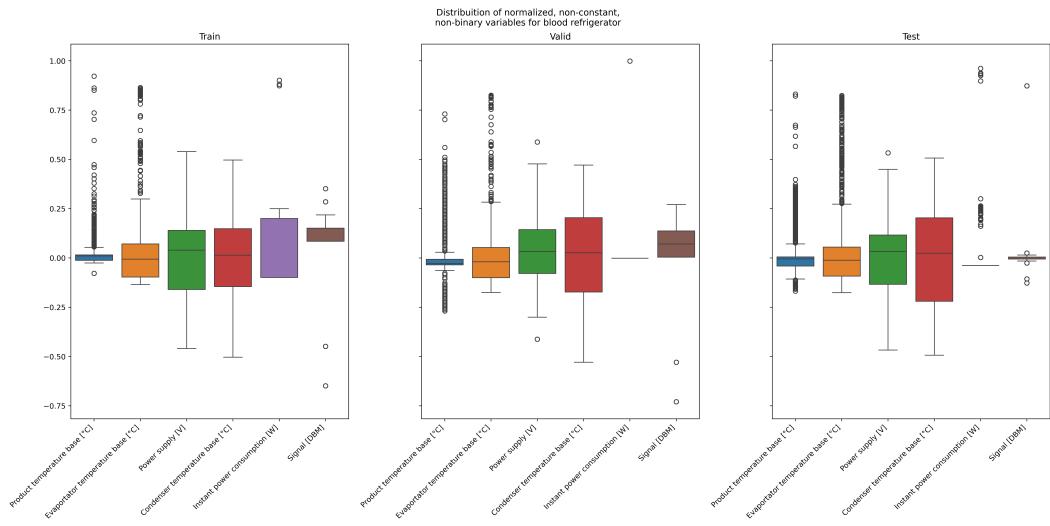


Figure 5: Distribution of normalized, continuous variables for the blood refrigerator dataset splits.

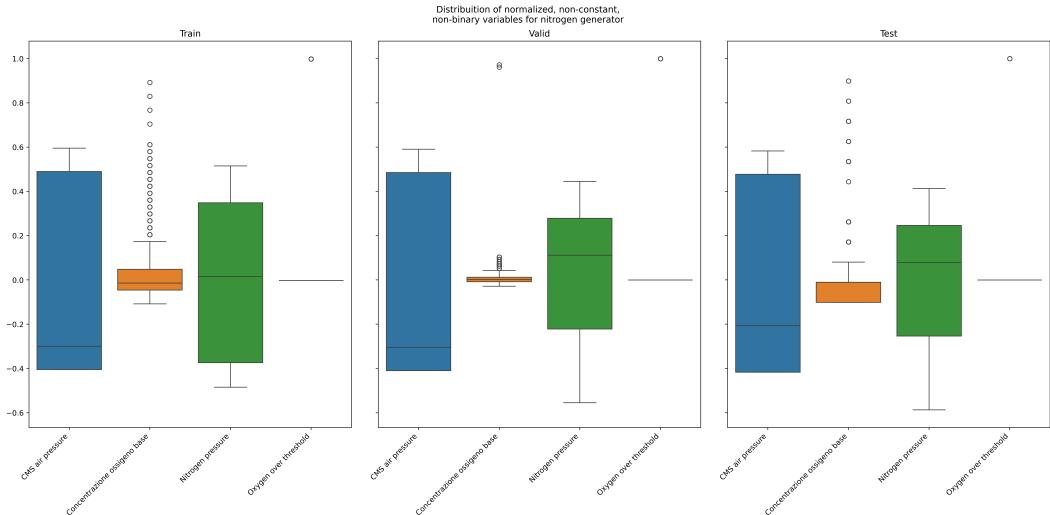


Figure 6: Distribution of normalized, continuous variables for the nitrogen generator dataset splits.

For additional detailed discussion and analysis of the datasets, the reader is referred to [12].

2.1.2 Processing

Pre-processing of the data was kept fairly minimal. As mentioned above, we decided to drop any columns with zero variance as these would not contribute any meaningful signal for any of our methods to learn from. This resulted in 12 columns for the blood refrigerator data and 4 columns for the nitrogen generator data. The remaining continuous variables were normalized with:

$$\hat{x}_{it} = \frac{x_{it} - \bar{x}_i}{\max x_i - \min x_i} \quad (1)$$

for variable i at time t .

2.2 Models

2.2.1 T² Control Chart

One method from SPC that is commonly used is the T²-Control Chart with Upper Control Limit (UCL). Given a set of n correlated characteristics X_i , assumed to follow a multivariate-normal distribution, the T² statistic is:

$$T^2 = m(\bar{x} - \mu)\Sigma^{-1}(\bar{x} - \mu) \quad (2)$$

Where m is the number of samples in the subgroup ($m = 1$ is referred to as individual control charts), \bar{x} is the mean of a sample of m vectors of measurements, Σ is the covariance matrix (assuming it is known), and μ is the in-control process mean (assuming it is known) [13]. For simplicity, it should be assumed that all discussion is with regards to individual T²-control charts unless otherwise stated.

The UCL is χ^2 distributed and is calculated at $\frac{\alpha}{2}$. In practice,

$$\text{UCL}_{\frac{\alpha}{2}} = \frac{(k-1)^2}{k} f\left(1 - \frac{\alpha}{2}; \frac{p}{2}, \frac{(k-p-1)}{2}\right) \quad (3)$$

where f is the beta distribution. The Lower Control Limit (LCL) is set to zero.

It is often the case that Σ and μ are not known. These can be estimated from the sample by:

$$\mu \approx \hat{\mu} = \frac{1}{k} \sum_{i=1}^k \hat{x}_i \quad (4)$$

where \hat{x}_i are the observations from the in-control process sample, and

$$\Sigma \approx \hat{\Sigma} = \begin{bmatrix} \hat{\sigma}_{11}^2 & \dots & \hat{\sigma}_{1n}^2 \\ \vdots & \ddots & \vdots \\ \hat{\sigma}_{n1}^2 & \dots & \hat{\sigma}_{nn}^2 \end{bmatrix} \quad (5)$$

where $\hat{\sigma}_{ij}^2$ is the variance between characteristics X_i and X_j from the in-control process sample.

To plot the i^{th} point on the individual chart, (1) becomes:

$$T_i^2 = (x_i - \hat{\mu})\hat{\Sigma}^{-1}(x_i - \hat{\mu}) \quad (6)$$

As new observations are made and new points added to the plot various signals may arise. Perhaps the most apparent is a breach of the UCL (a new point with a T^2 -statistic $>$ UCL). This could suggest that something unusual happened with that single instance. There are many other patterns that may arise. A continuous streak of 3-5 points beyond the 2 standard deviation mark, 7 or more beyond 1 standard deviation, 8-9 points alternating sides of the center line in a zig-zag pattern. All of these suggest different things about the underlying process and that either the process should be investigated or a new in-control sample taken to establish new charts.

We will be using the above technique largely as-is because it is a standard methodology and has demonstrated itself to be very effective. We will use the last 25%, 50%, 75%, and all of the training data to establish the control parameters. This will allow us to evaluate the effects of recency on the results. To determine when a "stop" signal should be raise, we have defined two patterns that accept a variety of parameters. When the conditions specified by the patterns are met, a 1 ("stop") is returned. Because this method is not trained in an iterative manner, we do not train or evaluate using the validation time series.

During the experiments, we found that the performance when using the training data to construct control charts for the test data yielded very poor results. So we included a set of results where the control limits are set using the first 150 or 70 time steps for the blood refrigerator and nitrogen generator, respectively, for each day in the test set. In other words, the control chart for a given day would be determined based on the sensor readings during the first part of its normal operation for a day.

The patterns we use are " n sequential breaches of the UCL or LCL" and " n total breaches of the UCL or LCL in a window of length t ". For the former case, if we raise a stop signal at the first instance of any window where all n observations are above or below (including a mix of above and below) the control limits. For the latter, we raise a stop signal at the first instance of a window of t time steps there are n breaches above or below the control limits.

2.2.2 Transformer

The transformer is an architecture with several components: a tokenizer, an embedding layer, transformer layers, a task-specific head [3]. The transformer layer is a collection of alternating multi-head attention layers and linear feed forward layers. These models are often pre-trained on large collections of text as foundation models then fine-tuned for specific tasks by adding a task-specific head that converts the representation from the last hidden layer of the transformer body into a useable output. It is possible to train transformers on more than just text. Many variations exist that operate on video, audio, time series, and combinations of the above (multi-modal).

A family of open-source foundational time series transformer models was recently released [4] that has been pre-trained on a rich collection of time series data and can be fine-tuned for specific tasks such as anomaly detection or classification. The models (MOMENT) were pre-trained with univariate time series data with a variety of time horizons. When given multivariate time series data, each variable is treated separately in its own channel.

The specific model we will be using has an architecture similar to the T5 encoder [14] and has approximately 385MM parameters. We will be fine-tuning a freshly initialized linear layer to perform the classification task.

For this model, we will extract rolling windows comprising 20 minutes (determined using the time stamps provided in the data) as the input sequence and the label for the time step immediately after the end of the window as the target for learning. MOMENT requires all input sequences are 512 tokens in length, so all inputs will be padded with zeros to meet the length requirement and masks will be generated so the loss is calculated only from features related to the input sequence. We chose 20

minutes since previous research found that was the window where the models typically performed best [12].

MOMENT will be fine tuned using the Adam optimizer [15] with a learning rate of 0.001 that uses a cosine annealing schedule to a minimum of 10^{-6} for binary classification. Weight decay is 0.1 and β_1, β_2 are 0.9 and 0.999, respectively. We use a batch size of 8. Fine-tuning is carried out over 5 epochs and the model is checkpointed at the end of each epoch if the validation loss decreased. Due to the high degree of class imbalance, we use normalized class weights in the loss function equal to:

$$[w_0, w_1] = [1/p_0, 1/p_1] * (1/p_0 + 1/p_1)^{-1} \quad (7)$$

2.2.3 Reinforcement Learning

Reinforcement learning is defined in the context of solving a Markov Decision Process (MDP). An MDP is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P})$ with a state space, action space, reward function, and a probability to transition from state to state given an action. A state in the state space $s_i \in \mathcal{S}$ is a representation of the environment at time i . In our case, each s is a vector in \mathbb{R}^n . The action space is comprised only of "Do nothing", "Stop", or, $0, 1 = \mathcal{A}$. Our reward function is defined so that the agent receives a positive reward if it returns 0 while the label is still 0 or 1 when the label is 1. The agent will receive a negative reward if it returns a 1 when the label is 0 and vice-versa. There is an additional penalty that is incurred for each time step that passes where the agent has not stopped the machine (returned a 0) but the label is 1, specifically:

$$\mathcal{R}(f, \dashv) = \begin{cases} a = 1 \wedge \text{PW_0.5h} = 1 & \rightarrow +1 \\ a = 1 \wedge \text{PW_0.5h} = 0 & \rightarrow -10 \\ a = 0 \wedge \text{PW_0.5h} = 1 & \rightarrow \sum_{i=0}^k i \text{ where } k \text{ is the time since the first stop signal} \\ a = 0 \wedge \text{PW_0.5h} = 0 & \rightarrow +0.01 \end{cases}$$

where $\text{PW_0.5h} = 0, 1$ indicates the label and a the action chosen by the agent at that step. The transition probabilities are defined by the temporal sequencing naturally present in the dataset and is deterministic in nature.

However, we observed extremely brittle policies by leaving the time series for each day in temporal order. So, we shuffle all of the days during each episode to prevent the agent from trying to learn dependencies day-to-day. While there is an argument to be made for the importance of day-to-day patterns, that is a problem that warrants its own paper. For simplicity, we will proceed only trying to have the agents learn how to act within each day.

For both algorithms, the hidden dimension for the neural networks is determined by the number of input features for the machine. PPO uses the first power of 2 larger than the input dimension and MCTS uses the second power of 2 larger than the input dimension. These were determined empirically over a limited space. For the blood refrigerator, with 12 variables, the hidden layers of the networks in the PPO agent have dimension 16 but dimension 32 for the MCTS agent. For the nitrogen generator the layers are 8 and 16, respectively.

Monte Carlo Tree Search The Monte Carlo Tree Search (MCTS) with neural guidance stores information in a tree structure with states as nodes and actions connecting from that node to the next state that it transitioned to as a result of that action. It gathers a number of paths by running simulations down the tree by selecting the action at each node that leads to the node with the largest expected value, starting at the root, and propagating values back up through the path based on a combination of the

expected value for that state and the number of visits to that node over the simulations. Some randomness is used to ensure some exploration can occur to increase the likelihood that the algorithm will find the optimal policy.

The MCTS algorithm has 4 steps:

1. Starting from the root, select child nodes until a leaf node is reached
2. If the leaf node is not a state that is equal to the process being stopped, create child nodes for each possible action
3. Perform a simulation by collecting paths from moving down the tree
4. Pass the result from this simulation (rollout) up through all nodes in their respective paths

Step 3 can be done randomly or with some sort of direction. This can be as simple as a greedy search (select the child node with the largest reward) or can use a more complex method such as a neural network [6,8]. The neural network takes the position in the tree and returns a vector of probabilities to select each action (the policy). These policy vectors are proportional to the exponential of the visit count for each node (i.e. $\pi \propto N(s, a)^{1/\tau}$, where N is the visit count and τ a scaling value). The neural network has two heads, one to return a policy and one to return a value for the node. Training is done by minimizing the cross entropy loss between the policy and the actions yielding the highest reward, plus, the mean squared error of the value of the current node and the discounted expectation of future returns. The ideal situation is for the network to learn to correctly choose an action that will keep moving the agent towards high-value states that yield, in aggregate, the largest cumulative reward.

We explored this method due to the unique nature of the policy optimization process and believed that it would provide an interesting and challenging implementation.

The network is trained using the Adam optimizer [15] with an initial learning rate of 0.0001 and uses cosine annealing to a minimum of 10^{-6} . Weight decay is set to 0.01 and β_1, β_2 are 0.9 and 0.999, respectively. The discount factor for future rewards is $\gamma = 0.99$. At the end of each time series for a day, end of an episode, or every 4th time step, a "learning" process is triggered. The process begins by collecting 100 simulations (paths) from the current tree using the current network. The agent samples 128 tuples from the 10,000 most recent experiences (s_t, a_t, r_t, s_{t+1}) to train the network. Each epoch is considered complete when the agent has progressed through all of the days' time series in the training set. We train the agent for 30 epochs and use early stopping if there is no improvement in the mean reward on the validation set after 5 epochs.

Proximal Policy Optimization The proximal policy optimization (PPO) [7] is an extention of their earlier trust-region policy optimization algorithm [16]. It is an on policy algorithm and uses two neural networks during training. The agent uses both a neural network to learn a policy, that is sampled from using a categorical distribution, to determine what action to take in a given state and a separate network to learn the expected value of a given state. We sample from a distribution determined by the policy to allow some amount of exploration to persist during training. This can help reduce the likelihood an agent gets stuck in a locally optimal policy and unable to find the globally optimal policy. The policy updates are clipped to prevent any single update from being too large and destabilizing the learning process.

The PPO algorithm contains some simplifications and improves the sample efficiency by using multiple policy updates per episode (or epoch). During each episode, the policy and value networks are updated k times, each time using a different mini-batch of recent experiences stored in the agent's memory. A simple queue is used to store the experiences in the order they occur up to a total of 10,000 experiences. As

new ones are added the oldest experiences will be removed from the memory to maintain 10,000 examples. We use 10 updates per learning instance with an experience replay batch size of 64 per update. The ϵ -clipping parameter is set to 0.2 (taken from the paper) and a learning rate of 0.0003 is used with the Adam [15] optimizer. The reward discounting factor is set to 0.999.

We chose to explore this algorithm since it has demonstrated excellent performance in a variety of continuous state spaces and is likely to demonstrate better convergence in this scenario. While the MCTS algorithm is very powerful, it tends to perform best when used either when there is an opponent the agent is playing against and/or when the state space is discrete.

2.3 Evaluation

The trained model or agent is used to generate predictions on the held-out test set. We collect the recall, precision, and F_1 (in both strict and relaxed settings) as well as the mean-time-from-event (MTFE). The MTFE measures the mean mistiming of the method, mathematically:

$$MTFE = \frac{1}{K} \sum_k (t_k m - t_k e) \quad (8)$$

The strict version of the classification metrics is determined where the model or agent returns the stop signal only on the very first occurrence of the stop signal in that time series or, in the case where a time series does not stop, if the agent never raises the stop signal. The relaxed variation considers any stop signal returned by the model at any stop label in the time series as correct.

The precision, recall, and F_1 are calculated in the usual manner. The only difference arises in the use of strict or lenient evaluation protocols.

Each method is trained on the training split then evaluated on the held-out test split. The control chart and transformer return predictions for each time step while the RL agents stop making predictions for a given day until they return a stop action or reached the end of the day's time series. The predictions are calculated on a per-day basis for each test split (e.g. the 27 days for the blood refrigerator test set are each given a single prediction).

The combination of traditional classification metrics allows us to evaluate the daily performance while the MTFE allows us to evaluate the intra-day performance. Meaning we can make more nuanced comparisons between techniques by having an understanding of how they perform at the macro and micro levels. Methods that have excellent precision and recall but high MTFE would suggest, for example, that while across days the method is accurately stopping the machine or allowing it to run without interruption, stopping the machine far too early might translate to more wasted person-hours evaluating and examining the machine before returning it to operation only for it to actually fail later in the day and require additional maintenance.

3 Results

3.1 T^2 Control Charts

3.1.1 Blood Refrigerator

Tables 2-5 report the MTFE, F_1 , Precision, and Recall scores for the T^2 control charts that were determined based on the last p proportion of the training split. Table 6 includes the same metrics but in the case where the first 150 time steps were used to

determine the control chart parameters for the day. The patterns are described in the first column of each table and translate to "XperYata" or "XseqAt α " where $\alpha = 0.05$, X is the number of breaches, and Y is the length of the window.

Overall, there are not many differences between the results for the lenient and strict evaluation protocols for all tables. The approximate spacing between time intervals is 34 seconds. The values that most closely mimic the 20 minute RW with 30 minute PW would be the "20per40at0.05" and "40seqAt0.05" rows from each table.

	Proportion of Training Split							
	0.25		0.50		0.75		1.00	
	Lenient	Strict	Lenient	Strict	Lenient	Strict	Lenient	Strict
5per10at0.05	-328	-328	-249	-249	-421	-421	-422	-422
10per20at0.05	-398	-398	-243	-243	-422	-422	-422	-422
20per40at0.05	-400	-400	-232	-232	-422	-422	-422	-422
30per60at0.05	-400	-400	-213	-213	-422	-422	-422	-422
40per80at0.05	-400	-400	-185	-185	-422	-422	-422	-422
60per120at0.05	-391	-391	-192	-192	-422	-422	-422	-422
120per240at0.05	-365	-365	-110	-110	-422	-422	-422	-422
5seqAt0.05	-326	-326	-245	-245	-421	-421	-421	-421
10seqAt0.05	-392	-392	-177	-177	-421	-421	-421	-421
20seqAt0.05	-343	-343	9	9	-421	-421	-418	-418
30seqAt0.05	-316	-316	197	198	-397	-397	-383	-383
40seqAt0.05	-299	-299	107	110	-397	-397	-383	-383
60seqAt0.05	-305	-305	188	193	-377	-377	-236*	-234*
120seqAt0.05	-284	-284	<i>-27*</i>	<i>-22*</i>	-298	-298	34	42
Mean	-353	-353	-98	-97	-406	-406	-370	-369

Table 2: Mean Time From Event for the T^2 Control Charts trained on varying portions of the training split. Best absolute result per column in bold, best without surpassing the start of failure in italics with *

Across all of tables 2-5, the best results are when only the most recent 50% of the training split is used to determine the charts. The weakest results are when 75% of the training split is used. It also tends to be the case that a higher number of breaches per window produces stronger results.

	Proportion of Training Split							
	0.25		0.50		0.75		1.00	
	Lenient	Strict	Lenient	Strict	Lenient	Strict	Lenient	Strict
5per10at0.05	0.00	0.00	0.20	0.20	0.00	0.00	0.00	0.00
10per20at0.05	0.11	0.11	0.29	0.29	0.00	0.00	0.00	0.00
20per40at0.05	0.11	0.11	0.29	0.29	0.00	0.00	0.00	0.00
30per60at0.05	0.11	0.11	0.36	0.36	0.00	0.00	0.00	0.00
40per80at0.05	0.11	0.11	0.36	0.36	0.00	0.00	0.00	0.00
60per120at0.05	0.11	0.11	0.36	0.36	0.00	0.00	0.00	0.00
120per240at0.05	0.11	0.11	0.36	0.36	0.00	0.00	0.00	0.00
5seqAt0.05	0.00	0.00	0.20	0.20	0.00	0.00	0.00	0.00
10seqAt0.05	0.11	0.11	0.29	0.29	0.00	0.00	0.00	0.00
20seqAt0.05	0.11	0.11	0.36	0.36	0.00	0.00	0.00	0.00
30seqAt0.05	0.11	0.11	0.48	0.42	0.00	0.00	0.00	0.00
40seqAt0.05	0.11	0.11	0.59	0.54	0.00	0.00	0.00	0.00
60seqAt0.05	0.20	0.20	0.67	0.62	0.00	0.00	0.11	0.11
120seqAt0.05	0.29	0.29	0.82	0.79	0.10	0.10	0.27	0.27
Mean	0.11	0.11	0.40	0.39	0.01	0.01	0.03	0.03

Table 3: F₁ scores for T² Control Charts on varying portions of the training split. Best in bold.

	Proportion of Training Split							
	0.25		0.50		0.75		1.00	
	Lenient	Strict	Lenient	Strict	Lenient	Strict	Lenient	Strict
5per10at0.05	0	0	1	1	0	0	0	0
10per20at0.05	1	1	1	1	0	0	0	0
20per40at0.05	1	1	1	1	0	0	0	0
30per60at0.05	1	1	1	1	0	0	0	0
40per80at0.05	1	1	1	1	0	0	0	0
60per120at0.05	1	1	1	1	0	0	0	0
120per240at0.05	1	1	1	1	0	0	0	0
5seqAt0.05	0	0	1	1	0	0	0	0
10seqAt0.05	1	1	1	1	0	0	0	0
20seqAt0.05	1	1	1	1	0	0	0	0
30seqAt0.05	1	1	1	1	0	0	0	0
40seqAt0.05	1	1	1	1	0	0	0	0
60seqAt0.05	1	1	1	1	0	0	1	1
120seqAt0.05	1	1	1	1	1	1	1	1
Mean	0.86	0.86	1.00	1.00	0.07	0.07	0.14	0.14

Table 4: Precision scores. Best in bold.

	Proportion of Training Split							
	0.25		0.50		0.75		1.00	
	Lenient	Strict	Lenient	Strict	Lenient	Strict	Lenient	Strict
5per10at0.05	0.00	0.00	0.11	0.11	0.00	0.00	0.00	0.00
10per20at0.05	0.06	0.06	0.17	0.17	0.00	0.00	0.00	0.00
20per40at0.05	0.06	0.06	0.17	0.17	0.00	0.00	0.00	0.00
30per60at0.05	0.06	0.06	0.22	0.22	0.00	0.00	0.00	0.00
40per80at0.05	0.06	0.06	0.22	0.22	0.00	0.00	0.00	0.00
60per120at0.05	0.06	0.06	0.22	0.22	0.00	0.00	0.00	0.00
120per240at0.05	0.06	0.06	0.22	0.22	0.00	0.00	0.00	0.00
5seqAt0.05	0.00	0.00	0.11	0.11	0.00	0.00	0.00	0.00
10seqAt0.05	0.06	0.06	0.17	0.17	0.00	0.00	0.00	0.00
20seqAt0.05	0.06	0.06	0.22	0.22	0.00	0.00	0.00	0.00
30seqAt0.05	0.06	0.06	0.32	0.26	0.00	0.00	0.00	0.00
40seqAt0.05	0.06	0.06	0.42	0.37	0.00	0.00	0.00	0.00
60seqAt0.05	0.11	0.11	0.50	0.45	0.00	0.00	0.06	0.06
120seqAt0.05	0.17	0.17	0.70	0.65	0.05	0.05	0.16	0.16

Table 5: Recall scores. Best in bold.

		Mean Time From Event	Precision	Recall	F1
5per10at0.05	Lenient	-421.37	1.00	0.74	0.85
	Strict	-421.37	1.00	0.74	0.85
10per20at0.05	Lenient	-420.52	1.00	0.74	0.85
	Strict	-420.52	1.00	0.74	0.85
20per40at0.05	Lenient	-410.81	1.00	0.74	0.85
	Strict	-410.81	1.00	0.74	0.85
30per60at0.05	Lenient	-413.67	1.00	0.74	0.85
	Strict	-413.67	1.00	0.74	0.85
40per80at0.05	Lenient	-402.89	1.00	0.74	0.85
	Strict	-402.89	1.00	0.74	0.85
60per120at0.05	Lenient	-402.11	1.00	0.74	0.85
	Strict	-402.11	1.00	0.74	0.85
120per240at0.05	Lenient	-377.67	1.00	0.74	0.85
	Strict	-377.67	1.00	0.74	0.85
5seqAt0.05	Lenient	-421.19	1.00	0.74	0.85
	Strict	-421.19	1.00	0.74	0.85
10seqAt0.05	Lenient	-419.74	1.00	0.74	0.85
	Strict	-419.74	1.00	0.74	0.85
20seqAt0.05	Lenient	-393.89	1.00	0.74	0.85
	Strict	-393.89	1.00	0.74	0.85
30seqAt0.05	Lenient	-380.89	1.00	0.74	0.85
	Strict	-380.89	1.00	0.74	0.85
40seqAt0.05	Lenient	-392.04	1.00	0.74	0.85
	Strict	-392.04	1.00	0.74	0.85
60seqAt0.05	Lenient	-417.85	1.00	0.79	0.88
	Strict	-417.85	1.00	0.79	0.88
120seqAt0.05	Lenient	-421.52	1.00	0.84	0.91
	Strict	-421.52	1.00	0.84	0.91

Table 6: Results when using the first 150 time steps of each day. Best for each metric in bold.

3.1.2 Nitrogen Generator

Tables 7 and 8 report the MTFE, F1, Precision, and Recall scores for the T^2 control charts that were determined based on the last p proportion of the training split. The precision, recall, and F_1 results were identical for all patterns and all splits so the results in table 8 represent all of those metrics. During the experiments to use the first portion of the day, we were unable to find an interval of time that would result in an invertible covariance matrix, thus the statistics could not be computed. The patterns are described in the first column of each table and translate to "XperYata α " or "XseqAt α " where $\alpha = 0.05$, X is the number of breaches, and Y is the length of the window.

Overall, there are not many differences between the results for the lenient and strict evaluation protocols for all tables. The approximate spacing between time intervals is 1 minute. The values that most closely mimic the 20 minute RW with 30 minute PW would be the "30seqAt0.05" row from each table.

	Proportion of Training Split							
	0.25		0.5		0.75		1	
	Lenient	Strict	Lenient	Strict	Lenient	Strict	Lenient	Strict
5per10at0.05	-764	-764	-764	-764	-791	-791	-791	-791
10per20at0.05	-755	-755	-764	-764	-791	-791	-791	-791
20per40at0.05	-757	-757	-763	-763	-791	-791	-791	-791
30per60at0.05	-758	-758	-765	-765	-791	-791	-791	-791
40per80at0.05	-756	-756	-766	-766	-791	-791	-791	-791
60per120at0.05	-683	-683	-767	-767	-791	-791	-791	-791
120per240at0.05	-691	-691	-691	-691	-791	-791	-791	-791
5seqAt0.05	-755	-755	-762	-762	-791	-791	-791	-791
10seqAt0.05	-753	-753	-762	-762	-791	-791	-791	-791
20seqAt0.05	-676	-676	-759	-759	-791	-791	-791	-791
30seqAt0.05	-676	-676	-759	-759	-791	-791	-791	-791
40seqAt0.05	-676	-676	-676	-676	-791	-791	-791	-791
60seqAt0.05	-676	-676	-676	-676	-791	-791	-791	-791
120seqAt0.05	-676	-676	-676	-676	-791	-791	-791	-791
Mean	-718	-718	-739	-739	-791	-791	-791	-791

Table 7: Mean Time From Event for the T^2 Control Charts trained on varying portions of the training split. Best absolute result per column in bold, best without surpassing the start of failure in italics with *

	Proportion of Training Split							
	0.25		0.5		0.75		1	
	Lenient	Strict	Lenient	Strict	Lenient	Strict	Lenient	Strict
5per10at0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
10per20at0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
20per40at0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
30per60at0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
40per80at0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
60per120at0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
120per240at0.05	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00
5seqAt0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
10seqAt0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
20seqAt0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
30seqAt0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
40seqAt0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
60seqAt0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
120seqAt0.05	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00
Mean	0.14	0.14	0.00	0.00	0.00	0.00	0.00	0.00

Table 8: F_1 scores for T^2 Control Charts on varying portions of the training split. The Precision and Recall score tables are identical, so we will simply reference these numbers. Best in bold.

		Mean Daily Test Reward	MTFE	Precision	Recall	F1
MCTS	Blood Refrigerator	-7.78	-391	0.47	0.41	0.44
	Nitrogen Generator	-10	-735	0.00	0.00	0.00
PPO	Blood Refrigerator	-0.03	-9	0.91	0.59	0.71
	Nitrogen Generator	-1.01	-107	1.00	0.66	0.80

Table 9: Results from evaluating the RL agents on the held out test set for each machine. The mean daily test reward is calculated by determining the total reward for each day then averaging the days together.

3.2 Transformer

During training a number of issues were encountered that prevented a useable, fine-tuned MOMENT model from being used. Since we were attempting to perform time series classification for each 20 minute reading window, we tried to use the package provided by the research group*. There was a bug in calling the classification architecture specifically and mention of this was found in the issues of their GitHub repo. To circumvent this, we implemented the necessary functionality from scratch in our implementation based on the code in their repo that would be used for classification tasks. We could successfully instantiate a model and complete forward passes and gradient backpropagations.

We then encountered issues during fine-tuning. Over epochs the loss would not decrease. After checking that gradients were being generated and pushed through the model at each step, tensor computations were being properly executed, and that the logits and labels were as expected, we proceeded to attempt a hyperparameter search. Given the size of the model and the limitations of the local compute (NVIDIA 4070 TI Super 16GB), we opted for a modest hyperparameter space that searched over various learning rates, weight decays, gradient clipping sizes, and dropout probabilities. Still, no improvements were noticed.

At this point we decided to modify the class weights we were passing to the loss function to address class imbalance. As shown in equation 7, we used the normalized inverse ratio to give more importance to the loss for the minority class. We experimented with modifications to this that increased and decreased the ratio, to no avail. We attempted this along side under- and over-sampling methods. Finally, we tried freezing the encoder and embedder in various combinations and using LoRA [17] and PiSSA [18] parameter efficient fine-tuning methods in a last-ditch effort to train the model.

Since we were unable to obtain a tuned model, we will not report results for this as they would be uninformative.

3.3 Reinforcement Learning

Table 9 shows the results after training the RL agents. The reinforcement learning algorithms have very different results. The MCTS didn't perform as strongly as the PPO agent. The PPO agent had the lowest MTFE that was before the first failure

*<https://github.com/moment-timeseries-foundation-model/moment>

label in the time series for each day. The RL algorithms were trained on a single NVIDIA L4 Tensor Core GPU with 16GB memory and 4 vCPUs.

4 Discussion

4.1 Findings

4.1.1 Blood Refrigerator

The use of control charts parameterized on various amounts of the training data provides mixed results. For most patterns and most proportions, the MTFE is often far in advance of when the failure signals occur in the labels. For example, the mean MTFE for the best proportion (0.5) was -97 time steps, in the strict case, which corresponds to stopping over 45 minutes early. This is not particularly bad, especially in the case where a delicate substance such as blood is concerned, but for other machines this can translate to lost revenue or production. This is much better though when we consider the case where control charts are calculated daily and the mean MTFE for the best pattern was -378 time steps, around 214 minutes early. Regardless of machine, this is very imprecise. The traditional metrics may have improved when switching from 50% of the training data to daily calculations, but that has its trade-offs. The daily charts are less likely to stop the machines on a daily basis, but when it does stop them it is likely to be 3-4 hours earlier than necessary. Again, for blood refrigerators this may not be terrible, but for machine in manufacturing settings these early signals may lead to diagnostic results that suggest the machine is fine, only for it to fail later in the day requiring it to be shutdown again.

The RL agents both had MTFEs closer to the start of the failure period, the PPO agent in particular averaged 9 time steps (4.5 minutes) ahead of failure, than the control chart approaches which would help catch the machine prior to failing without sacrificing as much operating time. However, the lower recall for both agents means that they may miss impending failure signals and the machine will proceed to fail without notice anyway. The precision for both agents is quite high though meaning it is less likely to raise false positive stops and therefore lead to less unnecessary down time than the control charting approaches. Some of the control chart patterns did also show high precision and recall, however, they also tended to have much larger MTFEs. The MCTS agent didn't perform as strongly as the PPO agent and this is likely due to the nature of the algorithm and the nature of the sensor data. By adopting the MCTS method to better operate on continuous state spaces, the performance may improve. However, as it stands, the PPO agent would certainly be the preferred approach.

Figure 7 shows some example T^2 statistics and control charts for the blood refrigerator on four different days. These charts depict those calculated using the first 150 time steps for each day. Figure 7a is a good example of the limitations of the chosen control chart patterns and how design of the stopping patterns is of critical importance to making the control charts effective. The PPO agent performed very well and stopped the machine close to the failure region. This would translate to less unnecessary downtime and person-hours spent servicing the machine. Figure 7b demonstrates one weakness of the T^2 control charts, the dependence on being able to calculate a reasonable inverse covariance matrix from the sample data. As a result of calculating the inverse covariance matrix, many of the statistics are very large. The PPO agent however does not suffer from this same limitation and is able to provide a more appropriately timed stopping signal, this suggests that the agent has generalized somewhat effectively from the training data. Figure 7c is an example where nearly all methods, except for the MCTS agent, stop the machine when it was operating within parameters. The 60seqAt0.05 and 120seqAt0.05 patterns were the only two that did not erroneously stop the machine on that day. Finally, figure 7d shows a machine

that is producing less varied T^2 statistics over the operational period of the machine. However, many of the patterns either stop the machine very early into the period or fail to stop it all together. The PPO agent provides a reasonably timed stop signal in this instance.

In response to our research questions, we can conclude:

1. PPO has the capacity to perform well for predicting failure for a blood refrigerator. MCTS seems less capable but further research is suggested before making a firm conclusion.
2. PPO shows less brittleness than the patterns we chose. With appropriate pattern definition the control charts may be more competitive. However, the overhead for implementing PPO is considerably greater compared to the control charts.
3. We cannot draw conclusions on this item since we could not successfully train the MOMENT model.
4. The PPO algorithm was able to converge and solve the problem more effectively than the MCTS algorithm for this machine.

4.1.2 Nitrogen Generator

As mentioned above, we were unable to find windows that would produce daily control charts for the nitrogen generator. Overall, the performance for all methods on the nitrogen generator was notably worse than for the blood refrigerator. This is likely due to there only being four variables remaining for the generator after removing variables with zero variance so any one variable has a much stronger effect on the difference between subsequent T^2 values than the relative effect of one variable for the blood refrigerator with 12 total variables. The charts all show processes that tend to "hug" one of the control limits. Typically, this would be interpreted as a fundamental shift having occurred in the machine that would require recalculation of the control limits. As shown in figure 8, we only observed one day from the test split where the control chart patterns didn't all stop the machine almost immediately. The T^2 statistics are still hugging the bottom control limit but the patterns with longer window or sequence lengths are more robust to this situation though still stop well ahead of the failure region.

The RL agents both struggled more with this machine as well as compared to the blood refrigerator. The time to converge for both algorithms was notably longer for both algorithms. This may be due to the few variables and their weak correlations with the label or the ease with which, even a small, neural network can over-fit. While the PPO agent did perform much better than the MCTS in this scenario, the performance is still not ideal and likely to generate very early or erroneous stop signals.

With regards to our research questions:

1. PPO seems to be able to perform acceptably. MCTS failed to converge entirely.
2. PPO shows some robustness to the trends observed in the control charts. The control charts would need considerable effort put into pattern selection to improve their performance or consideration of alternative control chart methods entirely.
3. We were unable to train MOMENT and cannot make conclusions on this question.
4. The PPO method can converge and perform reasonably well compared to the MCTS algorithm. However, both methods would want to have more comprehensive tuning before being considered for a deployment setting.

4.2 Limitations

Some of the limitations of this work include the lack of hyperparameter tuning for the transformer and reinforcement learning models. RL models, in particular, are rather finicky and have been found to be sensitive to random seed. Because they exhibit such a strong dependence on initial conditions and randomness, it would be ideal to include the random seed as a hyperparameter to be tuned during training. However, this comes at increased cost. Even on the data center class NVIDIA L4 GPU, training the PPO algorithm until it converged took an average of 6 hours per hyperparameter configuration and the MCTS algorithm varied from 2-8 hours. While the transformer wasn't able to be trained, a full training cycle over 5 epochs was expected to take around 3 hours on the NVIDIA 4070 TI Super. Even modest hyperparameter spaces would require days of compute and distribution across many machines to efficiently train these models. If the use case warranted such a sophisticated model, then the effort, and expense, may be worth the effort.

The dataset that is available is very good and a great resource, but still rather small for some of these models. For example, even though the blood refrigerator had over 60,000 time steps in the training split, it amounted to only 25 unique training days. This means the RL agents would see the same 25 days for each episode which is not very diverse when one considers the infinite state space possible for 12 variables over an entire day. If an individual works at an organization that owns the machines and has the sensors constantly delivering data, this issue can be mitigated as they will have access to considerably more data.

As mentioned, we only evaluated the methods on an intra-day basis. Some machines may be expected to run continuously 24/7 and would require a different approach to modelling.

After the issues experienced with the MOMENT model, further investigation of their training data revealed that they did use sensor time series in the dataset. However, it also included time series from biomedical phenomena, finance, traffic, etc. It may be necessary to have more sensor data and use a more aggressive fine-tuning approach to overcome the differences between these other types of time series data and our own. The model was also trained using single-channel (univariate) time series rather than multi-channel which may have contributed to the difficulties we experienced. Additional exploration is necessary to make this approach feasible. However, we believe that it would be a worthwhile effort considering the results shared in the MOMENT paper and the general abilities of foundational transformer models at large.

4.3 Future Work

Continuing to explore these threads would be a worthwhile endeavor if automated methods for machine stopping are required. Depending on the sophistication required or the criticality of proper stopping of the machine, some avenues may be more worthwhile than others.

Further research in the control charting methods would be to explore other multivariate control charting methods for this problem. It would also be pertinent to do a thorough comparison of various methods for determining the length of the baseline control period. Our results indicate that using some amount of the most recent sensor data may be effective in some cases but other cases may warrant daily recalculation of charts. Finding patterns that correspond to particular failure scenarios in the machine would also be an interesting direction.

Continuing to explore the use of RL algorithms would necessitate more work on the computational infrastructure side of things. Aside from the many developments in RL algorithms that have been published since the methods we used, finding useful hyperparameter combinations is very intense task. Additionally, interesting research could be done in increasing the sample efficiency to make the most of the limited dataset.

This may be achievable with simple tuning of the PPO parameters or with some modifications to the algorithm. It would also be interesting to attempt to train an agent that uses the data from various machines. This may improve the generalizability of the agent, which could alleviate some of the brittleness that trained deep RL algorithms can exhibit.

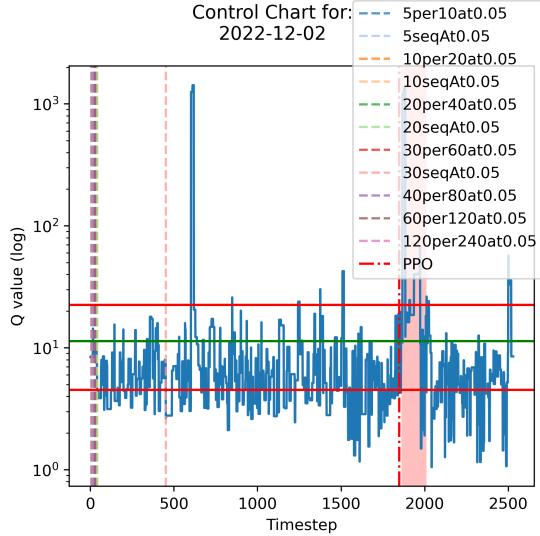
5 Conclusion

We have shown that classic techniques from statistical process control and more recent advances in reinforcement learning can both provide effective methods for detecting machine failure from multivariate time series data. Each comes with its own set of considerations and which is the best fit for a particular use case will strongly depend on a number of factors unique to each organization and machine. However, control charts can provide a good initial technique that can serve both as an initial method to use while other methods are explored and as a baseline to compare other techniques against. Due to its efficacy, low computation cost, and ease of implementation, T^2 control charts are a good starting point for machine monitoring.

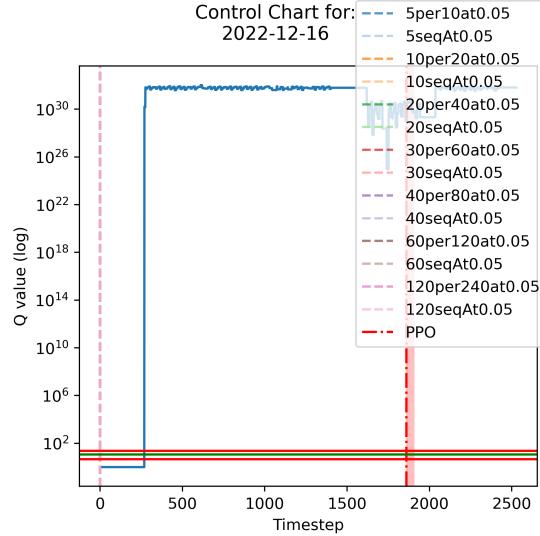
References

- [1] William F. Guthrie. NIST/SEMATECH e-Handbook of Statistical Methods (NIST Handbook 151), 2020.
- [2] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, November 1997.
- [3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need, August 2023. arXiv:1706.03762 [cs].
- [4] Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski. MOMENT: A Family of Open Time-series Foundation Models, May 2024.
- [5] Richard S. Sutton and Andrew G. Barto. *Reinforcement learning: an introduction*. Adaptive computation and machine learning. MIT Press, Cambridge, Mass, 1998.
- [6] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George Van Den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354–359, October 2017.
- [7] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347, 2017.
- [8] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, January 2016.

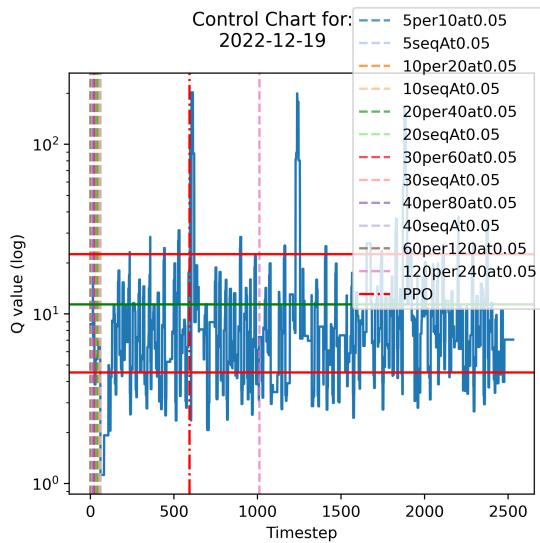
- [9] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Christopher Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub W. Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning. *ArXiv*, abs/1912.06680, 2019.
- [10] Vitchyr H. Pong, Shixiang Shane Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models: Model-free deep rl for model-based control. *ArXiv*, abs/1802.09081, 2018.
- [11] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Manfred Otto Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *CoRR*, abs/1509.02971, 2015.
- [12] Nicolò Oreste Pincioli Vago, Francesca Forbicini, and Piero Fraternali. Predicting Machine Failures from Multivariate Time Series: An Industrial Case Study. *Machines*, 12(6):357, May 2024.
- [13] John Lawson. *An introduction to acceptance sampling and SPC with R*. CRC Press, Boca Raton London New York, first edition edition, 2021.
- [14] Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2019.
- [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [16] John Schulman, Sergey Levine, P. Abbeel, Michael I. Jordan, and Philipp Moritz. Trust region policy optimization. *ArXiv*, abs/1502.05477, 2015.
- [17] J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685, 2021.
- [18] Fanxu Meng, Zhaohui Wang, and Muhan Zhang. Pissa: Principal singular values and singular vectors adaptation of large language models, 2024.



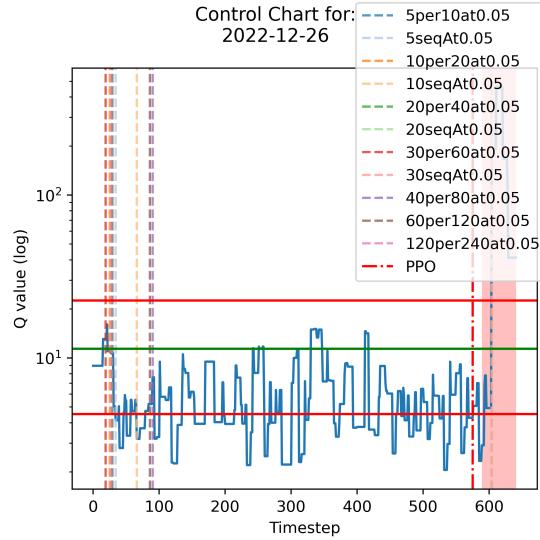
(a) An example where nearly all of the methods stopped the process prior to the failure region. The PPO agent stopped the process closest to the region, with the MCTS agent failing to stop. The control chart patterns all stopped very early into operation.



(b) A process where the T^2 statistics became unbounded resulting in premature stopping by all of the control chart patterns. The MCTS algorithm failed to stop the process. The PPO agent did stop the process closer to the region.

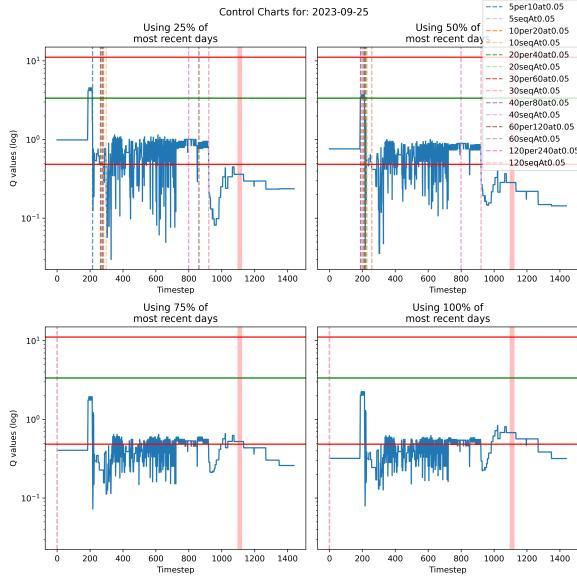


(c) A day where the machine never entered into a failure state. Notably, the PPO agent incorrectly stopped the machine while some of the control chart patterns did not.

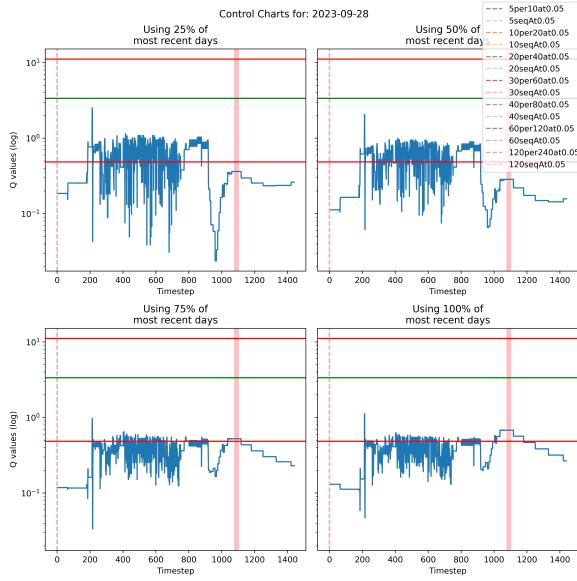


(d) An example where the MCTS algorithm stopped the process nearest the failure region. The PPO agent stopped the process near the region and many of the control charts stopped it very early.

Figure 7: Plots of the T^2 statistics and T^2 control chart elements for the blood refrigerator. Vertical lines indicate where the various methods raised a stop signal for the process. Red regions are where failures were reported in the labels. The absence of a red region means the machine never entered into a failure state. The absence of a method in the legend means the method did not raise a stop signal.



(a) The only example where the patterns didn't stop the machine almost immediately.



(b) A representative example of the other days.

Figure 8: T^2 control charts for the nitrogen generator calculated using proportions of the training split.