

EXPLANATION USING SIMULATION

Anonymous authors

Paper under double-blind review

ABSTRACT

In safety-critical domains, such as industrial systems, the lack of explainability in predictive ‘black-box’ machine learning models can hinder trust and adoption. Standard explainability techniques, while powerful, often require deep expertise in data analytics and machine learning and fail to align with the sequential, dynamic nature of data in these environments. In this paper, we propose a novel explainability framework that leverages reinforcement learning (RL) to support model predictions with visual explanations based on dynamical system simulation. By training RL agents to simulate events that require prediction, we use these agents’ critics to make classifications. Next, we employ the actors of the RL agents to simulate the potential future trajectories underlying these classifications, providing visual explanations that are more intuitive and align with the expertise of industrial domain experts. We demonstrate the applicability of this method through a case study involving monitoring a small industrial system for cyberattacks, showing how our framework generates actionable predictions that are supported with visual explanations. This approach aims to bridge the gap between advanced machine learning models and their real-world deployment in safety-critical environments.

1 INTRODUCTION

With the increasing volume of data generated in industrial environments and the growing complexity of managing these systems, deep learning has become a promising tool to automate or support decision-making for industrial domain experts (Wang et al., 2022). A key function of deep learning in these settings is predictive analytics—making classifications to forecast future events, which can assist operators in planning for critical decisions. Despite the potential, the adoption of deep learning in industrial environments faces a major challenge: explainability. While deep learning models have demonstrated high reliability and accuracy, their ‘black-box’ nature has put into question their trustworthiness and limited their adoption in industries that must adhere to stringent regulatory standards (Ahmed et al., 2022). In safety-critical industries, transparency in decision-making is often mandated by regulations.

Several methods, notably model-agnostic techniques like Local Interpretable Model-agnostic Explanations (LIME) (Ribeiro et al., 2016) and SHapley Additive exPlanations (SHAP) (Lundberg, 2017), have been utilized to provide a degree of explainability for deep learning models. However, these techniques require a level of expertise in data analytics that most domain experts lack, making them difficult for industrial practitioners to interpret and apply effectively. Instead, many of these domain experts rely on time-series analysis, which involves observing data sequences over time to understand system behavior, predict future trends, and make informed decisions (Fatima & Rahimi, 2024). As the volume of data grows, this manual analysis becomes increasingly challenging, creating a greater need for deep learning systems that can support decision-making while remaining interpretable.

In this paper, we propose a novel approach that generates and explains deep learning-based classifications for industrial dynamic systems by providing time-series explanations that align with the workflows and practices of domain experts. Our approach leverages reinforcement learning (RL) agents to model the sequences of events that lead to key events that domain experts need to predict. Next, we demonstrate the use of the RL agents’ state-action values or state values—generated by the agents’ deep neural network-based critics—to make predictive classifications, and the learned action policies to simulate the sequence of events that lead up to those outcomes and underlie the classifi-

054 cation. This simulation offers an visually intuitive, dynamic explanation of how the model reaches
 055 its decisions that is more aligned with the skills and practice of industrial domain experts. Conse-
 056 quently, this method stands to enhance trust in artificial intelligence (AI), facilitating its adoption in
 057 industries where safety, reliability, and transparency are of paramount importance.

059 2 RELATED WORK

061 **Explainability:** Post-hoc explainable AI (XAI) methods, designed to explain the predictions of
 062 trained black-box models, can be broadly categorized into *model-specific* and *model-agnostic* ap-
 063 proaches (Minh et al., 2022). Model-specific methods are tailored to a particular class of models,
 064 utilizing their internal structure to generate explanations. For instance, methods such as Layer-
 065 wise Relevance Propagation (LRP) (Montavon et al., 2019), DeepLIFT (Shrikumar et al., 2017),
 066 heatmaps (Payer et al., 2019), saliency maps (Adebayo et al., 2018), GradCAM (Selvaraju et al.,
 067 2020) backpropagate gradients from the output to the input layers in neural networks to visually
 068 emphasize regions in the input that have the greatest contribution to the model’s decision. Model-
 069 agnostic methods, on the other hand, do not rely on the internals of a specific model. Instead, they
 070 aim to provide explanations based purely on input-output behavior. SHAP (Lundberg, 2017) is one
 071 such method, which explains the contribution of each input feature to the prediction using game-
 072 theoretic principles. LIME (Ribeiro et al., 2016) approximates the black-box model locally using a
 073 simpler, interpretable model and uses the coefficients of this simpler model to explain the relative im-
 074 portance of each feature to the prediction. More recently, textual justification methods have emerged
 075 as an explainability approach (Shi et al., 2018; Sabol et al., 2020; Musto et al., 2021; Aminimehr
 076 et al., 2024; Hartmann et al., 2022). These methods generate natural language explanations of model
 077 predictions, aiming to make the reasoning accessible to general users. Musto et al. (2021) showed
 078 that text justification can generate user-friendly explanations that improve transparency and enhance
 079 user trust and satisfaction in the context of recommendation systems. A notable trend is that text-
 080 ual justification methods have been primarily proposed for applications such as recommendation
 081 systems, medical diagnosis, and stock prediction, highlighting the effectiveness of this explanation
 082 approach for domain experts in specialized fields. The motivation for our work aligns with that
 083 of textual justification, but offers a different medium of explanation tailored to domain experts in
 industrial domains: time-series plots.

084 **Explainable RL:** While sharing the same tools, the field of explainable RL (Wells & Bednarz, 2021;
 085 Heuillet et al., 2021) typically focuses on interpreting the reasoning, actions, and decision-making
 086 processes of RL agents after they have learned a task. In contrast, our approach does not aim to
 087 explain the internal decision-making of the agents themselves. Instead, we leverage the agents’
 088 learned behaviors to explain how specific predictions could unfold in the system. The goal is not to
 089 understand the agent’s reasoning but to use its actions to simulate potential future events, thereby
 090 providing a clearer explanation of the predicted outcomes.

091 **RL in Industry:** RL research is pervasive across industrial sectors, including power systems (Zhang
 092 et al., 2019), autonomous driving (Aradi, 2020), smart cities (Ullah et al., 2020), and manufactur-
 093 ing (Wang et al., 2021).

095 3 CONCEPT

096
 097 RL agents learn to make decisions by interacting with a dynamical system environment, aiming to
 098 maximize cumulative rewards (or minimize cumulative penalties). In an actor-critic RL architecture,
 099 the agent consists of two key components: the actor and the critic. The critic estimates a state value:
 100 the expected return (cumulative future rewards) from a given state S under a policy π , expressed as

$$101 V^\pi(S) = \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i R_{t+i} \mid S_t = S \right], \quad (1)$$

102 or, alternatively, a state-action value (or Q-value): the expected return starting from state S , taking
 103 action A , and then following policy π , expressed as

$$104 Q^\pi(S, A) = \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i R_{t+i} \mid S_t = S, A_t = A \right]. \quad (2)$$

108 $\gamma \in [0, 1)$ is a discount factor determining the importance of future rewards.

109
110 Through training, the critic learns to map the data, collected as observations, of the environment to
111 a numerical value that predicts the event the RL agent is trained to achieve. For instance, if an RL
112 agent is trained to inject malicious disturbances or manipulate machinery in an industrial setting, the
113 critic learns to map system data to a state value that predicts when these manipulations could lead to
114 equipment damage or system failure. Alternatively, in an autonomous driving simulation, where an
115 RL agents affect the external vehicle’s motion, the critic learns to anticipate risky driving situations
116 based on the car’s external environment data.

117 In our research, we extend the critic’s state value to generate both numerical and categorical values
118 that provide predictive classifications based on system data.

119 The actor, on the other hand, dictates the actions taken by the agent. Since the actor learns to simulate
120 the sequence of actions that lead to the predicted event, we leverage the actor to provide explain-
121 ability. Specifically, the actor is used to simulate the sequence of events leading to a classification,
122 providing supporting visual explanation.

123 Figure 1 illustrates our conceptual framework. We train multiple RL agents, each designed to cause
124 a different event that requires prediction. Each agent’s critic processes observations from the pro-
125 duction system and feeds its state-action (or state) value to a function that maps it to a numerical
126 score representing the event’s impact. These mappings may include the probability of the event,
127 which in negative impact cases, translates to a measure of risk—commonly defined as the product
128 of impact and likelihood. The outcome with the highest impact (or risk) determines the classifica-
129 tion. When prompted for an explanation, the actor corresponding to the highest impact generates the
130 action sequence that leads to the predicted outcome. Next, the system’s simulated behavior under
131 the actor’s actions is visualized as a time-series plot, allowing the user to understand the system’s
132 future trajectory underlying the classification.

133 To demonstrate the advantages of our approach, particularly in safety-critical systems, we consider
134 the scenario of a small, localized electric microgrid targeted by cyberattackers. In this scenario, a
135 predictive analytics system must monitor the microgrid’s data and output a classification that reflects
136 the potential impact of cyberattacks manipulating a generator’s load-frequency control. Maintaining
137 a stable power frequency—at 60 Hz in North America—is critical for the proper functioning of
138 the grid. Deviations cause power flicker (experienced as lights rapidly dimming and brightening),
139 which can damage sensitive equipment by forcing it to operate outside of its designed-frequency
140 range. Power system equipment is expensive, and significant frequency anomalies are typically
141 detected by protection devices, which isolate the equipment to prevent their damage. Consequently,
142 by strategically compromising control, cyberattackers could potentially isolate generators and cause
143 power blackouts. The dangers posed by such cyberattacks have been the topic of numerous research
144 studies, surveyed in (Mohan et al., 2020).

145 We choose this scenario for several important reasons. First, power systems are increasingly the tar-
146 get of sophisticated cyberattacks, with substantial evidence pointing to nation-state actors involved
147 in cyber espionage (Hjortdal, 2011) and attacks on electric grids, as demonstrated in the attacks on
148 Ukraine’s grid in 2015 (Case, 2016) and 2016. Power systems are also massive in scale in complex-
149 ity, and stand to significantly benefit from incorporating deep learning-based methods for monitoring
150 and risk assessment. Additionally, power systems domain experts cannot be expected to have exten-
151 sive data analytics expertise. Second, the lack of data on cyberattacks makes it challenging to
152 predict threats from historical data. RL becomes a viable tool here, as it can generate simulated data
153 to anticipate the impact of unknown cyberattacks and help detect them when they occur. Finally, cy-
154 berattacks often involve subtle and strategic manipulations of the system. The critic in RL can help
155 uncover these complex failure modes, making it an effective approach for predicting and explaining
156 potential risks in critical infrastructure systems.

157 4 METHOD

158
159 Considering a dynamical system expressed as:

$$160 \dot{\mathbf{x}} = g(\mathbf{x}, \mathbf{u}, \mathbf{A}) \quad (3)$$

$$161 \mathbf{S} = h(\mathbf{x}, \mathbf{u}, \mathbf{A}) \quad (4)$$

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

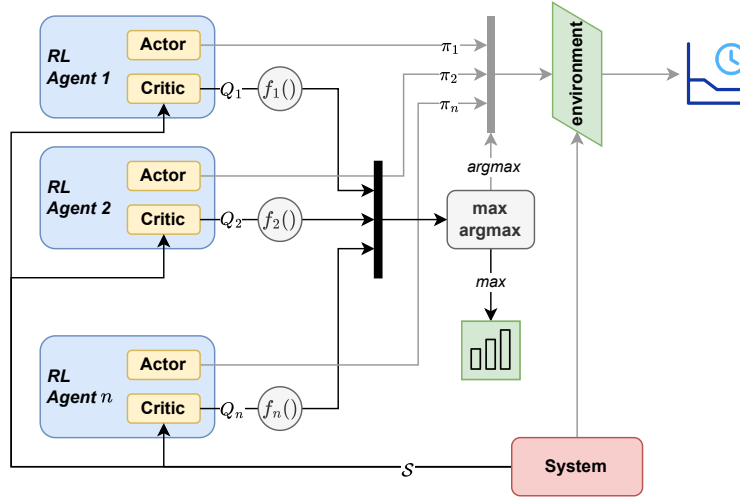


Figure 1: Concept illustrating n critics providing classifications and n actors supporting with explanations.

where $\mathbf{A} = \{A_1 \cup A_2 \cup \dots \cup A_n\}$ represents the set of actions that n RL agents can input into the system, and $\mathbf{S} = \{S_1 \cup S_2 \cup \dots \cup S_n\}$ denotes their respective observations. \mathbf{u} denotes inputs not generated by the RL agents, and \mathbf{x} represents the state of the dynamical system. Each RL agent is trained to achieve specific outcomes that domain experts need to predict. Depending on the use case, training can involve different paradigms, including single-agent or multi-agent reinforcement learning. This formulation allows for both single-agent and multi-agent RL setups.

Once each RL agents converges on an action policy, each agent i will have trained deep neural networks for both its actor and critic. The actor network π_i maps the agent's observations to actions:

$$\pi_i : S_i \rightarrow A_i \quad (5)$$

which selects actions that, when applied greedily, aim to maximize the agent's expected cumulative reward. The critic network estimates the value of the agent's state, either as a state-value function:

$$C_i : S_i \rightarrow V_i \quad (6)$$

or as a state-action value function (Q-value):

$$C_i : (S_i, A_i) \rightarrow Q_i \quad (7)$$

The critic's values are dependent on the reward function defined for each RL agent. These values are then mapped to a measure of predicted impact or risk using the following function:

$$f_i(V_i(S_i)) = P_{\text{likelihood},i} \cdot f_{\text{impact},i}(V_i(S_i)) \quad (8)$$

where $f_{\text{impact},i} \geq 0$ is a non-negative function that increases with the predicted impact of the outcome, and $P_{\text{likelihood},i}$ (optional) represents the probability of the outcome, which allows risk estimation. Alternatively, when using the Q-value, the impact or risk is computed as:

$$f_i(Q_i(S_i, A_i)) = P_{\text{likelihood},i} \cdot f_{\text{impact},i}(Q_i(S_i, A_i)) \quad (9)$$

At runtime, each agent's critic receives real-time observations from the production system, and the agent with the highest predicted impact or risk score is selected:

$$\{j, f_j(S_j)\} = \arg \max_i \{f_i(S_i)\}^{i \in \{1, \dots, n\}} \quad (10)$$

For better comprehension, the score can also be categorized based on predefined thresholds or classes. When an explanation is required for a classification, the policy π_j of the agent with the

highest impact or risk score is used to simulate the sequence of actions from time t up to a pre-defined horizon $t + T$ or until the outcome has been realized in the simulation, i.e.,

$$A_j = \pi_j(S_j) \quad (11)$$

$$A_i = \mathbf{0} \text{ for } i \in \{1, \dots, n\} \setminus j \quad (12)$$

The system state is then presented as a time-series plot, providing the domain expert with a visual explanation of the classification or predicted impact or risk score.

4.1 EXPERIMENT

We present a case study of a small microgrid, adapted from Mohamed & Kundur (2024). The microgrid is equipped with a security system that monitors the frequency and its rate-of-change (derivative), providing a classification that represents the vulnerability of the system to potential cyberattacks. When a high-urgency classification is detected, it signals the likelihood of an ongoing cyberattack that may be compromising the system. In such a scenario, a critical decision—such as isolating compromised communication channels—must be made to mitigate the attack. However, communication between system components is vital for the safety and stability of the microgrid, so this decision must be based on a highly trustworthy classification, triggered only under real threat to prevent further system damage.

In this experiment, we train RL agents to simulate cyberattacks aimed at compromising the system and forcing it into blackout conditions. The RL agents are tasked with learning strategies that satisfy several key objectives:

1. The cyberattacks must introduce disturbances to the generator’s control system, manipulating its power output. An increase in power generation will lead to a rise in the microgrid’s frequency, while a reduction in power generation will decrease the frequency.
2. Both the frequency and its rate-of-change are measured in per-unit (pu), with frequency normalized by the nominal 60 Hz. Similarly, the rate-of-change is divided by 60 Hz for a consistent measurement unit. Protective devices within the microgrid are programmed to activate when the frequency or its rate-of-change deviates from safe thresholds, typically 3% pu for frequency deviation and 5% pu for the rate-of-change.
3. It is desirable for the attack to keep the frequency deviations minimal to avoid triggering countermeasures. The more subtle the deviations, the harder it is to detect the attack. Therefore, the RL agents should aim to manipulate the rate-of-change of the frequency rather than the absolute frequency itself to ensure stealthier attacks.

To meet these objectives, we design a reward function that guides the RL agent to force large deviation in the rate-of-change of frequency while keeping the frequency deviation minimal. The reward function is defined as follows:

$$R_t = \left(\frac{s_2}{5\%}\right)^2 \cdot \max\left(0, 1 - \left(\frac{s_1}{3\%}\right)^2\right) + 30\{|s_2| > 5\%\} - 5\{|s_1| > 3\%\} \quad (13)$$

In this reward formulation, s_1 denotes the frequency deviation from nominal, and s_2 is the rate-of-change of the frequency. The agents observations are $S = (s_1, s_2)$. The first term in equation 13 encourages the agent to increase the rate-of-change of frequency while minimizing the frequency deviation. The second and third terms provide additional incentive: a large positive reward for pushing the rate-of-change outside the safe operating range of 5%, and a large penalty for large frequency deviations. Episodes terminate when the conditions in the second and third terms are satisfied.

Appendices A and B provide further details on the state-space representation of the dynamical system and the architecture and hyperparameters of the Deep Deterministic Policy Gradient (DDPG) agent used in the experiments.

5 RESULTS

Figure 2a illustrates the action sequence injected into the system by a RL agent. The effect is shown in Figure 2b, where the agent induces a resonance in the system, causing the rate-of-change of

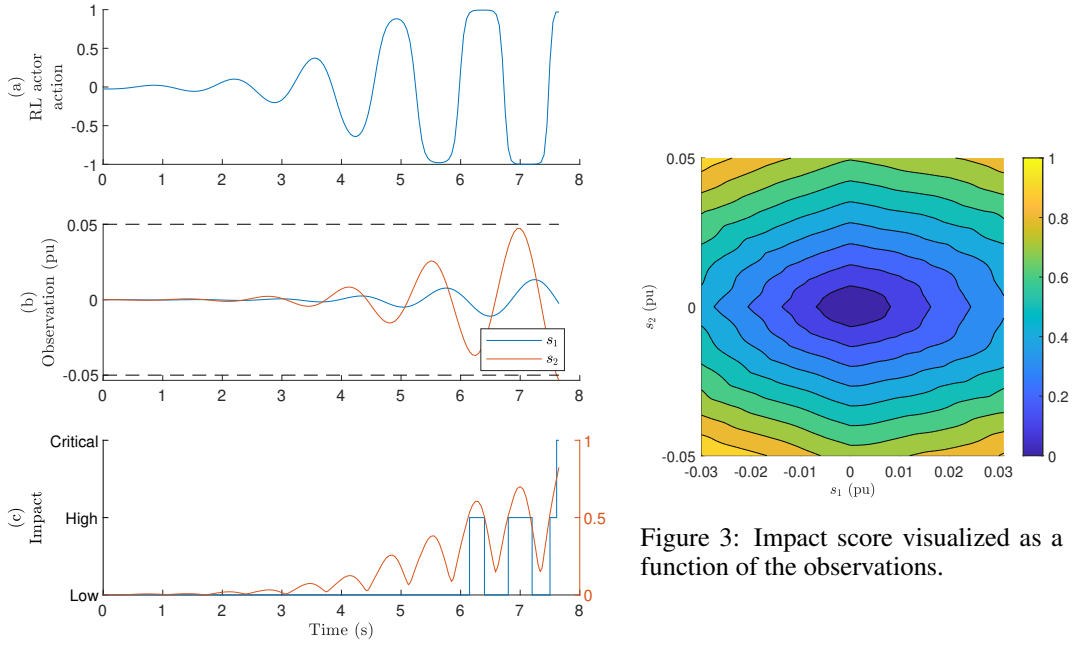


Figure 2: (a) action sequence generated by RL agent; (b) effect on system observations; (c) impact score and category.

frequency (depicted in red) to progressively increase until it breaches the lower threshold, marked by the lower dashed horizontal line. During training, this RL agent consistently favored driving the rate-of-change below the lower safety threshold. To include contrasting behavior, we trained a second RL agent to manipulate the system such that the rate-of-change of frequency exceeds the upper safety threshold.

After training, we observed that the Q-values of the agents decrease (i.e., become more negative) as the system approaches failure. This behavior occurs because, near failure, the agent has fewer remaining timesteps to collect potential rewards, whereas, further from failure, the agent can accumulate rewards over a larger time window. Consequently, the Q-value is larger when the system is far from failure. Based on this, we defined the function mapping the Q-value to the impact score as follows:

$$f_1(Q(S, \pi(S))) = f_2(Q(S, \pi(S))) = \frac{e^{-Q(S, \pi(S))} - 1}{M} \quad (14)$$

where M normalizes the impact corresponding to the Q-value near system failure. Specifically, M is determined as:

$$M = e^{-\min\{Q(S, \pi(S))\}} - 1 \quad (15)$$

The forms of $f_1(\cdot)$ and $f_2(\cdot)$ ensure that the impact score is non-negative and increases as the system approaches failure. For simplicity, we do not weigh $f_1(\cdot)$ or $f_2(\cdot)$ over one another, nor do we consider likelihoods of events, i.e., $P_{\text{likelihood}, i} = 1, i \in \{1, 2\}$.

Figure 3 visualizes the impact as a function of the system’s observations. The impact \mathcal{I} is calculated as:

$$\mathcal{I} = \max\{f_1(Q_1(S, \pi(S))), f_2(Q_2(S, \pi(S)))\} \quad (16)$$

Given the low-dimensional observation space, this visualization effectively captures the relationship between the system’s state and the risk of cyberattacks. Intuitively, the plot shows that the impact

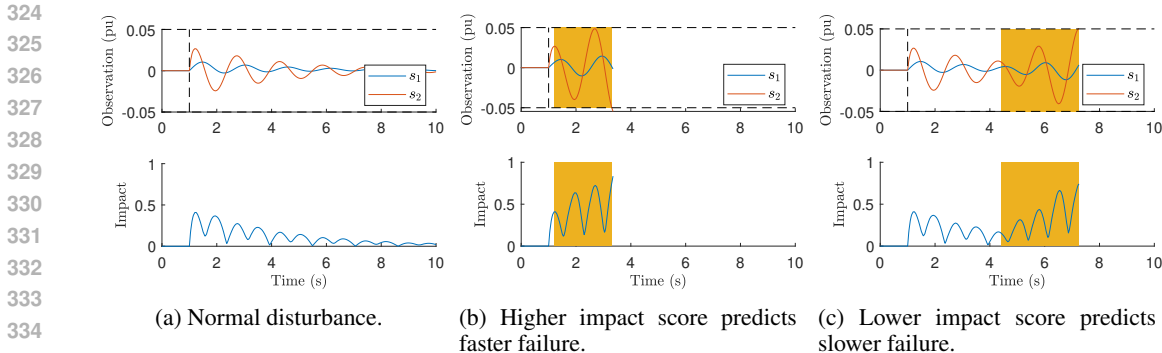


Figure 4: Explaining real-time classifications made by the RL agents’ critics via simulations. The yellow zones represent the simulated portions.

increases as the rate-of-change of frequency approaches the 5% pu thresholds. We can further categorize the impact into three levels, which can be used for system alerts or state classification:

- $I \in [0, 0.5)$: Low; system is safe, attack unlikely
- $I \in [0.5, 0.7)$: High; potential risk to the system, possible attack
- $I \in [0.7, 1]$: Critical; imminent system failure, active attack

Notably, the ‘Critical’ category alert to s_2 —the rate-of-change of frequency—nearing the critical 5% pu threshold, signaling a high likelihood of system failure due to a cyberattack. Hence, action must be taken to mitigate attacks when the system is in ‘High’ category.

To illustrate, Figure 2c shows the evolution of the numerical impact value (in red) during the cyber-attack, alongside the corresponding categorical impact levels (in blue). As the attack progresses, the impact value continues to rise. The attack starts at 0 seconds. Approximately 6 seconds after, the classification spikes to level 2, indicating a high likelihood of an ongoing cyberattack. Just before the 8-second mark, the classification briefly spikes to level 3, signifying imminent system failure right before the attack succeeds in failing the system.

To further illustrate the use of the classification system and its explainability during normal operation, Figure 4a shows the effect of a normal system disturbance on the observations s_1 and s_2 . During this disturbance, the impact increases but remains within level 1, indicating that the system is safe. The impact value reaches its peak at 1.2 seconds, signaling a moment of relatively heightened vulnerability. The impact value suggests that a cyberattack initiated precisely at this point in time would pose a greater risk compared to one initiated later. We use our proposed method to verify and explain this. Figure 4b shows the actor’s explanation of an attack that begins exactly at the moment when the impact value peaks, and continues through the period highlighted by the yellow zone. For comparison, Figure 4c illustrates the effects of an attack that starts later. By examining these figures side by side, it becomes evident that the system fails more quickly when the attack is launched at the instance of peak impact in Figure 4a. This demonstrates the system’s heightened susceptibility to failure at that specific point in time, validating the predictive and interpretable utility of the proposed method.

6 DISCUSSION

Accessibility of explanation through simulation: One of the key strengths of this approach is its ability to offer explanations for classifications and predictions in an intuitive and accessible format, particularly for industrial domain experts. By leveraging simulations, this method aligns well with the existing expertise of these professionals, who are accustomed to monitoring time-series plots, events, and sequences. The simulations create a visual narrative that makes understanding the underlying dynamics more straightforward, thereby facilitating the integration of machine learning insights into industrial systems.

Data generation in the absence of historical data: Another advantage of this method is its capability to generate synthetic data for events that lack historical precedent. By employing multiple RL

agents, the system explores and predicts a diverse range of events, including those that are subtle or previously unknown. Moreover, this approach enables domain experts to iteratively interact with the RL environment, allowing them to introduce new test cases or expand the number of events requiring prediction. This adaptability is key to refining the predictive model over time, as experts can augment RL agents and environments with additional details as necessary.

6.1 CHALLENGES

Expertise required for RL design: Although our approach alleviates the need for deep technical expertise in explainability methods, we acknowledge the complexity involved in designing and training RL agents. This presents a potential barrier, as creating effective RL models demands specialized knowledge. However, the trade-off is that this complexity is primarily concentrated in the RL development phase, which occurs infrequently. In contrast, the resulting system, once in place, simplifies interpretation and understanding, a task that will be required far more often by a broader audience of users.

The importance of historical data: While our focus has been on the use of RL in generating predictions through simulation, we recognize the value of historical data in classification algorithms. Historical data provides critical insights into patterns that can inform decision-making. Although our method primarily emphasizes online RL learning in simulated environments, there is potential to use offline RL (Levine et al., 2020), which leverages previously collected datasets. A hybrid approach, combining both historical data and simulated environments, could further enhance the predictive models.

Simulation environments: A key limitation of relying on online RL is that it necessitates the availability of accurate dynamical models that can be incorporated into the RL environment. For many industrial systems, models form the foundation of system design and are readily available. Furthermore, with the rapid growth of digital twins and simulation software, domain experts will increasingly have access to models that closely mimic real-world environments. This progression will help bridge the gap between virtual simulations and real-world applications, making the adoption of our method more feasible.

7 CONCLUSION

In this paper, we proposed a novel direction for explaining deep learning model classifications through time-series simulations. We developed a real-time predictive system for dynamical environments, where multiple reinforcement learning (RL) agents are employed to make predictions and support these predictions with simulations that visualize the forecasted system behavior underlying the classification. The motivation behind this approach is to offer an alternative method for explainability in contexts where domain experts are more accustomed to analyzing time-series plots to understand system behavior and anticipate future outcomes.

This work paves the way for further exploration of RL-based explainability, particularly in settings where traditional post-hoc XAI methods fall short. We encourage future research to apply this framework to a range of industrial case studies, explore hybrid RL approaches that combine offline and online learning, and investigate the method in diverse RL neural network architectures to further enhance the potential of this method.

REFERENCES

- Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *Advances in neural information processing systems*, 31, 2018.
- Imran Ahmed, Gwanggil Jeon, and Francesco Piccialli. From artificial intelligence to explainable artificial intelligence in industry 4.0: a survey on what, how, and where. *IEEE Transactions on Industrial Informatics*, 18(8):5031–5042, 2022.
- Amirhossein Aminimehr, Pouya Khani, Amirali Molaei, Amirmohammad Kazemeini, and Erik Cambria. Tbxplain: A text-based explanation method for scene classification models with the

- 432 statistical prediction correction. In *Proceedings of the Conference on Governance, Understanding*
 433 *and Integration of Data for Effective and Responsible AI*, pp. 54–60, 2024.
- 434
- 435 Szilárd Aradi. Survey of deep reinforcement learning for motion planning of autonomous vehicles.
 436 *IEEE Transactions on Intelligent Transportation Systems*, 23(2):740–759, 2020.
- 437
- 438 Defense Use Case. Analysis of the cyber attack on the ukrainian power grid. *Electricity information*
 439 *sharing and analysis center (E-ISAC)*, 388(1-29):3, 2016.
- 440
- 441 Syeda Sitara Wishal Fatima and Afshin Rahimi. A review of time-series forecasting algorithms for
 442 industrial manufacturing systems. *Machines*, 12(6):380, 2024.
- 443
- 444 Mareike Hartmann, Han Du, Nils Feldhus, Ivana Kruijff-Korbayová, and Daniel Sonntag. Xaines:
 445 Explaining ai with narratives. *KI-Künstliche Intelligenz*, 36(3):287–296, 2022.
- 446
- 447 Alexandre Heuillet, Fabien Couthouis, and Natalia Díaz-Rodríguez. Explainability in deep rein-
 448 forcement learning. *Knowledge-Based Systems*, 214:106685, 2021.
- 449
- 450 Magnus Hjortdal. China’s use of cyber warfare: Espionage meets strategic deterrence. *Journal of*
 451 *Strategic Security*, 4(2):1–24, 2011.
- 452
- 453 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tuto-
 454 rial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 455
- 456 Scott Lundberg. A unified approach to interpreting model predictions. *arXiv preprint*
 457 *arXiv:1705.07874*, 2017.
- 458
- 459 Dang Minh, H Xiang Wang, Y Fen Li, and Tan N Nguyen. Explainable artificial intelligence: a
 460 comprehensive review. *Artificial Intelligence Review*, pp. 1–66, 2022.
- 461
- 462 Amr S. Mohamed and Deepa Kundur. On the use of reinforcement learning for attacking and de-
 463 fending load frequency control. *IEEE Transactions on Smart Grid*, 15(3):3262–3277, 2024. doi:
 464 10.1109/TSG.2023.3343100.
- 465
- 466 Athira M Mohan, Nader Meskin, and Hasan Mehrjerdi. A comprehensive review of the cyber-
 467 attacks and cyber-security on load frequency control of power systems. *Energies*, 13(15):3860,
 468 2020.
- 469
- 470 Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, Wojciech Samek, and Klaus-Robert
 471 Müller. Layer-wise relevance propagation: an overview. *Explainable AI: interpreting, explaining*
 472 *and visualizing deep learning*, pp. 193–209, 2019.
- 473
- 474 Cataldo Musto, Marco de Gemmis, Pasquale Lops, and Giovanni Semeraro. Generating post hoc
 475 review-based natural language justifications for recommender systems. *User Modeling and User-*
 476 *Adapted Interaction*, 31(3):629–673, 2021.
- 477
- 478 Christian Payer, Darko Štern, Horst Bischof, and Martin Urschler. Integrating spatial configuration
 479 into heatmap regression based cnns for landmark localization. *Medical image analysis*, 54:207–
 480 219, 2019.
- 481
- 482 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ” why should i trust you?” explaining the
 483 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference*
 484 *on knowledge discovery and data mining*, pp. 1135–1144, 2016.
- 485
- 486 Patrik Sabol, Peter Sinčák, Pitoyo Hartono, Pavel Kočan, Zuzana Benetinová, Alžbeta Blichárová,
 487 L’udmila Verbóová, Erika Štammová, Antónia Sabolová-Fabianová, and Anna Jašková. Explain-
 488 able classifier for improving the accountability in decision-making for colorectal cancer diagnosis
 489 from histopathological images. *Journal of biomedical informatics*, 109:103523, 2020.
- 490
- 491 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,
 492 and Dhruv Batra. Grad-cam: visual explanations from deep networks via gradient-based local-
 493 ization. *International journal of computer vision*, 128:336–359, 2020.

486 Lei Shi, Zhiyang Teng, Le Wang, Yue Zhang, and Alexander Binder. Deepclue: visual interpretation
487 of text-based deep stock prediction. *IEEE Transactions on Knowledge and Data Engineering*, 31
488 (6):1094–1108, 2018.

489 Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through
490 propagating activation differences. In *International conference on machine learning*, pp. 3145–
491 3153. PMIR, 2017.

492 Zaib Ullah, Fadi Al-Turjman, Leonardo Mostarda, and Roberto Gagliardi. Applications of artificial
493 intelligence and machine learning in smart cities. *Computer Communications*, 154:313–323,
494 2020.

495 Junliang Wang, Chuqiao Xu, Jie Zhang, and Ray Zhong. Big data analytics for intelligent manufac-
496 turing systems: A review. *Journal of Manufacturing Systems*, 62:738–752, 2022.

497 Ling Wang, Zixiao Pan, and Jingjing Wang. A review of reinforcement learning based intelligent
498 optimization for manufacturing scheduling. *Complex System Modeling and Simulation*, 1(4):
499 257–270, 2021.

500 Lindsay Wells and Tomasz Bednarz. Explainable ai and reinforcement learning—a systematic re-
501 view of current approaches and trends. *Frontiers in artificial intelligence*, 4:550030, 2021.

502 Zidong Zhang, Dongxia Zhang, and Robert C Qiu. Deep reinforcement learning for power system
503 applications: An overview. *CSEE Journal of Power and Energy Systems*, 6(1):213–225, 2019.

511 A MICROGRID MODEL

512 The microgrid model is as follows:

$$513 \dot{\mathbf{x}} = \begin{bmatrix} 0 & 0 & 0 & -(kB) & 0 & 0 \\ 514 1/\tau_G & -1/\tau_G & 0 & -d/(\tau_G) & 0 & 0 \\ 515 0 & 1/\tau_T & -1/\tau_T & 0 & 0 & 0 \\ 516 0 & 0 & 1/M & -D/M & 0 & 0 \\ 517 0 & 0 & 0 & 1/\tau_\omega & -1/\tau_\omega & 0 \\ 518 0 & 0 & 1/(M\tau_\nu) & -D/(M\tau_\nu) & 0 & -1/\tau_\nu \end{bmatrix} \mathbf{x} \\ 519 + \begin{bmatrix} 0 & 0 \\ 520 0 & -k \\ 521 0 & 0 \\ 522 -1/M & 0 \\ 523 0 & 0 \\ 524 -1/(M\tau_\nu) & 0 \end{bmatrix} \mathbf{u} + \begin{bmatrix} k \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \mathbf{A}$$

525 Table 1: Model Data

Parameter	Symbol	Value
AGC gain	k	3
Droop gain	d	40
Governor time-constant	τ_G	0.08
Turbine time-constant	τ_T	0.45
Generator inertia	M	6
Damping constant	D	0.03

526 Frequency sensors time-constants $\tau_\omega = \tau_\nu = 0.1$
527 Control center frequency measurement gain $B = 1$

B RL AGENT ARCHITECTURE

Table 2: DDPG neural network architectures and hyperparameters

Actor network			
Layer	# of units	Hyperparameters	
Input	$2 (s_1, s_2)$	$M = 128$	
Normalization	2	$\alpha_\theta = 10^{-4}, \alpha_\phi = 10^{-3}$	
Fully-connected	100	$\gamma = 0.99$	
ReLU		$\tau = 10^{-3}$	
Fully-connected	50	$N \sim \mathcal{N}(0, 0.3)$	
ReLU			
Tanh (or Sigmoid)			
Scaling	1		
Output	1 (A)		
Critic network			
Layer	# of units	Layer	# of units
Input	$2 (s_1, s_2)$	Input	1 (A)
Normalization	2	Normalization	1
Fully-connected	100	Fully-connected	50
ReLU			
Fully-connected	50		
Addition	50	✓	
ReLU			
Fully-connected	1		
Output	1 $Q(s_1, s_2, A)$		