

Causal versus Marginal Shapley Values for Robotic Lever Manipulation Controlled using Deep Reinforcement Learning

Sindre Benjamin Remman¹, Inga Strümke² and Anastasios M. Lekkas³

Abstract—We investigate the effect of including domain knowledge about a robotic system’s causal relations when generating explanations. To this end, we compare two methods from explainable artificial intelligence, the popular KernelSHAP and the recent causal SHAP, on a deep neural network trained using deep reinforcement learning on the task of controlling a lever using a robotic manipulator. A primary disadvantage of KernelSHAP is that its explanations represent only the features’ direct effects on a model’s output, not considering the indirect effects a feature can have on the output by affecting other features. Causal SHAP uses a partial causal ordering to alter KernelSHAP’s sampling procedure to incorporate these indirect effects. This partial causal ordering defines the causal relations between the features, and we specify this using domain knowledge about the lever control task. We show that enabling an explanation method to account for indirect effects and incorporating some domain knowledge can lead to explanations that better agree with human intuition. This is especially favorable for a real-world robotics task, where there is considerable causality at play, and in addition, the required domain knowledge is often handily available.

Index Terms—Deep reinforcement learning, robotics, explainable artificial intelligence, Shapley additive explanations, causal SHAP

I. INTRODUCTION

Data-driven control methods have become widespread over the last years due to their ability to capture changes in system dynamics and adapt accordingly. The control system can use such methods to successfully adapt to situations that the engineer cannot envision beforehand. Reinforcement learning-based methods have shown great promise in terms of adaptability in robotics applications. However, this adaptability comes at a cost, as reinforcement learning (RL) methods are often paired with a function approximator such as a deep neural network (DNN), which are in general not interpretable for humans. The combination of RL with DNNs is called deep reinforcement learning (DRL) and has been prominent in the last decade because of its high performance on several difficult tasks [1], [2]. DRL has also had success within robotic manipulation [3]–[5]. However, the non-interpretable nature of DNNs implies that using DRL

to control a real cyber-physical system during safety-critical operation is not prudent.

The problem with interpretability permeates the current machine learning (ML) state-of-the-art. Based on this, scientists are currently researching how to explain the decisions of ML agents, or even how to make the agents explain themselves. The field addressing these issues is called explainable artificial intelligence (XAI), from which a steadily increasing number of methods are being developed. Among the first and most widely used XAI methods is Local Interpretable Model-agnostic Explanations (LIME), presented in [6]. This method locally approximates the uninterpretable model using an interpretable model, for instance, a linear model. Linear LIME in an instance of an *additive feature attribution methods*, whose defining property is having an explanation model that is a linear function of binary variables. This was formalized by [7], who also realized that several explanation methods share this property, thus unifying several explanation methods and introducing SHapley Additive exPlanations (SHAP). The SHAP framework produces feature attributions satisfying the axioms of the Shapley decomposition, a solution concept from cooperative game theory [8]. Adapting linear LIME to satisfy the Shapley axioms results in the feature attribution method KernelSHAP [7].

In the case of SHAP, the binary variables indicate presence or absence of a model feature, as the Shapley value is calculated by considering all possible arrangements of contributors to an outcome. SHAP implementations, therefore, rely on calculating an ML model’s expected outcome in the absence of model features. In KernelSHAP [7], this sampling is done using a marginal distribution of the excluded features, which amounts to assuming independence between the model features. As argued by [9], explanations generated using the marginal distribution can only represent the direct effects of features on the model, not the indirect effects [10]. Taking the causal structure in the data into account, [9] present a modification to the SHAP package, named *causal SHAP*.

The contributions of this paper are the following:

- We employ causal SHAP for explaining a system that involves robotic manipulation. To achieve this, we analyze the system to uncover its causal structure.
- We compare the explanations generated by KernelSHAP and causal SHAP. In doing so, we investigate the effect of taking indirect feature effects into account, thereby obtaining feature attributions based on a more complete physical description of the system at hand.

This paper is organized as follows: in Section II, we present the necessary theory behind DRL, SHAP and causal SHAP;

¹Sindre Benjamin Remman is with the Department of Engineering Cybernetics, Norwegian University of Science and Technology (NTNU), Trondheim, Norway sindre.b.remman@ntnu.no

²Inga Strümke is with the Department of Engineering Cybernetics, Norwegian University of Science and Technology (NTNU), Trondheim, Norway inga.strumke@ntnu.no

³Anastasios M. Lekkas is with the Department of Engineering Cybernetics, Centre for Autonomous Marine Operations and Systems (AMOS), Norwegian University of Science and Technology (NTNU), Trondheim, Norway anastasios.lekkas@ntnu.no

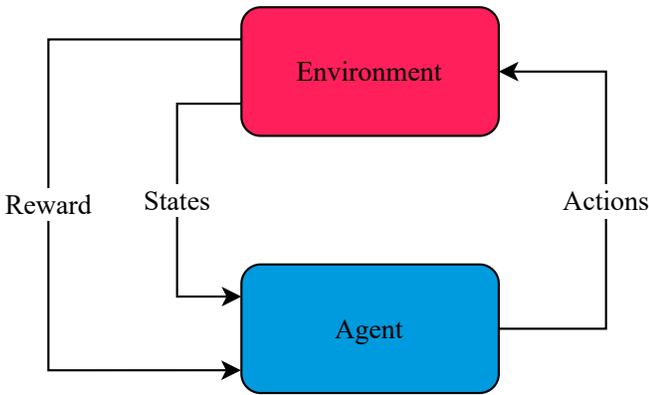


Fig. 1: The reinforcement learning loop.

in Section III, we describe the task to be solved using DRL, the experimental design, and how we use SHAP and causal SHAP to explain the decision-making agent; in Section IV, we present and discuss results obtained; and finally, in Section V, we draw our conclusions.

II. PRELIMINARIES

This section gives an overview of the theory and terminology necessary to understand the remainder of this paper. Firstly, we give an overview of the fundamentals of DRL. Secondly, the theory behind SHAP is explained. Lastly, we look at how SHAP is modified to create causal SHAP values.

A. Deep Reinforcement Learning

In RL, we divide the system into two parts: the *agent* and the *environment*. The interactions between the two are illustrated in Figure 1. The agent receives a state from the environment, performs an action based on this state, and receives a new state together with a reward from the environment. This cycle then repeats for the whole operation. The goal of RL is to find a *policy* that maps states to actions. The so-called *optimal policy* does this in an optimal way, in the sense that it maximizes the long-term expected reward, defined by the discounted infinite horizon model:

$$E\left[\sum_{t=0}^{\infty} \gamma^t r_t\right], \quad (1)$$

where $\gamma \in [0, 1]$ is the discount factor, and r_t is the reward received at time t [11, pp.13-15].

As previously stated, DRL refers to RL where the function approximator is a DNN. Here, we use a DRL agent trained using the Deep Deterministic Policy Gradient (DDPG) [12] algorithm. This is an actor-critic algorithm, which means that it trains two neural networks, the *actor-network* and the *critic-network*. The actor-network functions as the policy, which means that it maps states to actions, and the critic-network is used to guide the training of the actor-network. From a control engineering point of view, the policy is then akin to a controller. DDPG trains a deterministic policy, which means that a specific policy will always give the same output for the same input.

B. Shapley Additive Explanations

The Shapley decomposition, introduced by Lloyd Shapley in 1953 [8], has in recent years been applied extensively in the XAI literature. It is a solution concept from cooperative game theory and distributes a game’s outcome among the participants while uniquely preserving efficiency, monotonicity, and equal treatment, see, for instance, Theorem 2 in [13]. The Shapley value of participant i is calculated as a weighted mean over all subsets $\mathcal{S} \subseteq \mathcal{N}$ of the game’s N participants, not containing participant i :

$$\phi_i = \sum_{\mathcal{S} \subseteq \mathcal{N}} \frac{|\mathcal{S}|!(N - |\mathcal{S}| - 1)!}{N!} (v(\mathcal{S} \cup \{i\}) - v(\mathcal{S})). \quad (2)$$

Here, $v(\mathcal{S})$ is the *characteristic function*, which maps any set of participants in the game to a single real number $2^{\mathcal{N}} \rightarrow \mathbb{R}$, and thus fully characterizes the game.

In the context of XAI, the game can, for instance, be represented by a ML model performing a prediction task, the result of the game by the model prediction, and the participants of the game by the model’s input features. We can then obtain a feature attribution from the Shapley decomposition, quantifying how each of the input features affects the model prediction. The prediction of a ML model f trained on a set of data with features \mathbf{x} , made on the specific input features \mathbf{x}^* , can be decomposed as follows

$$f(\mathbf{x}^*) = \phi_0 + \sum_{i=1}^N \phi_i^*, \quad (3)$$

with ϕ_0 the expected value of the model output across the data set, $E[f(\mathbf{x})]$, and ϕ_i^* the Shapley value for the specific prediction on $\mathbf{x} = \mathbf{x}^*$.

There are two main challenges associated with calculating Shapley values for feature attribution: First, the calculation is very computationally expensive. A model using N features would need to be evaluated 2^N times, once for each feature’s inclusion or exclusion, as is readily seen from Equation (2). Second, it is in general not possible to evaluate a fitted ML model with sets of features missing. To circumvent these challenges, implementations such as the widely used SHAP, introduced by Lundberg and Lee [7], rely on approximations. The SHAP calculation uses as characteristic function an estimate of the expected model prediction, conditional upon the values of the included features, $\mathbf{x}_{\mathcal{S}} = \mathbf{x}_{\mathcal{S}}^*$, namely

$$v(\mathcal{S}) = E[f(\mathbf{x}) | \mathbf{x}_{\mathcal{S}} = \mathbf{x}_{\mathcal{S}}^*], \quad (4)$$

in the notation of [14]. As such, SHAP values attribute the change in the expected model prediction to each model feature by estimating how much each feature contributes to driving the model prediction away from its mean prediction across a data set. KernelSHAP, introduced in [7], estimates Equation (4) with absent features using a marginal distribution, which amounts to the assumption of independence between the included and excluded features.

Various changes to the sampling procedure used for estimating the expected model prediction have since been

suggested, and particularly relevant in this context are [9], [14]–[16].

C. Causal SHAP

In order to calculate the expected model prediction, Heskes et al. suggest [9] adapting the sampling procedure in the SHAP calculation, conditioning the absent features upon the values of the included features by *intervention*. First suggested by Aas et al. [14], the SHAP calculation can use the conditional distribution of the excluded features, instead of the marginal. Furthermore, interventional probabilities can be inferred from conditional probabilities using Pearl’s *do*-calculus [17], [18]. Combining conditional SHAP with the *do*-calculus thus allows us to use the interventional distribution in the SHAP calculation. Then, Equation (4) becomes

$$v(\mathcal{S}) = E[f(\mathbf{x}) | do(\mathbf{x}_{\mathcal{S}} = \mathbf{x}_{\mathcal{S}}^*)], \quad (5)$$

which we calculate by integrating over the absent features $\bar{\mathcal{S}}$, as detailed in [9]. The main advantage of the resulting so-called *causal* SHAP values is that both direct as well as *indirect* effects of the model features are taken into account. The direct effects represent the change in the model’s prediction due to a change in a feature without changing the absent features. The indirect effects, on the other hand, represent the change caused in the absent features by the intervention upon a feature, see equations (5) in [9]. The inclusion of these indirect effects constitutes the main difference between causal SHAP values and the marginal SHAP values introduced earlier, as the latter by construction only represent indirect effects.

The causal SHAP implementation used in this paper¹ was created by the first author, by adapting the R implementation by [9] to Python. In the implementation, we specify the causal structure of the data via a (partial) causal ordering in the form of a nested list in which each list is defined as causally dependent on the elements in the preceding list(s). In addition, we use a second, separate list to specify whether the dependencies within each nested list result from a confounding factor or mutual interactions between the components in the nested list.

III. METHODOLOGY

In this section, we describe the task of the DRL agent, its training, our dataset creation, and finally, how we analyze it using the two different SHAP implementations.

A. Lever manipulation task

The robotic manipulator that is used in this paper is the OpenMANIPULATOR-X² by Robotis, which can be seen in Figure 2. This manipulator has five degrees of freedom, four for the joints and one for the gripper. We do not use the first joint during this lever manipulation task, which corresponds to a rotation about the manipulator’s base. This



Fig. 2: The OpenMANIPULATOR-X.

is both because this makes the training more efficient, but also because the angle of this joint is trivial to solve for using

$$\theta_1 = \arctan2(y_{lever}, x_{lever}),$$

where y_{lever} and x_{lever} is the y - and x - coordinates expressed in the inertial frame of the manipulator.

The task involves moving the lever from a randomly selected start angle to a randomly selected target angle. These target and start angles are selected uniformly according to

$$\theta_{start}, \theta_{target} \in \mathbb{R} : \theta_{start}, \theta_{target} \in [-1.0 \text{ rad}, 1.0 \text{ rad}],$$

and $|\theta_{start} - \theta_{target}| > 0.4 \text{ rad}$.

B. States and Actions

In this task, the dimension of the state-space is eight and consists of the joint angles of the manipulator, the distance between the two fingers of the gripper, the horizontal and vertical distance from the end-effector to the lever, and the current and desired angles of the lever. See Figure 3 below for a visualization of the system’s states.

The action vector is of dimension four, where the first three entries correspond to the desired relative movement of the shoulder, elbow and wrist joints respectively. The fourth entry in the action-vector indicates whether the gripper should open or close:

$$a_4 \geq 0, \rightarrow \text{Gripper should open}$$

$$a_4 < 0, \rightarrow \text{Gripper should close.}$$

C. Training procedure

The agent was trained using the DDPG algorithm together with the technique Hindsight Experience Replay (HER) [19]. HER enables the usage of sparse rewards, and we, therefore, give sparse rewards according to

$$r = \begin{cases} -1, & \text{if } |\theta_{lever} - \theta_{target}| \geq 0.025 \text{ rad} \\ 0, & \text{if } |\theta_{lever} - \theta_{target}| < 0.025 \text{ rad} \end{cases},$$

where 0.025 rad is the chosen precision for the lever manipulation task. The agent is trained in simulations, and we use two different simulators for this. The first simulator is PyBullet, which is a fast simulator, but for our purposes, not close enough to the real-world environment. The agent trained in PyBullet is then transfer learned in Gazebo, which

¹https://github.com/sbremman/causal_shap_python

²https://emanual.robotis.com/docs/en/platform/openmanipulator_x/overview/

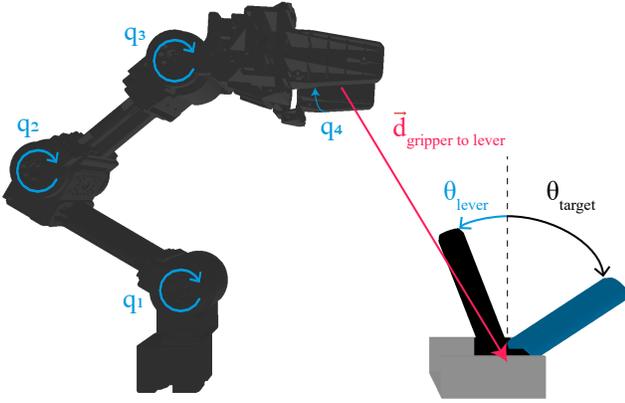


Fig. 3: Visualization of the system’s states.

is slower, but more similar to the real-world environment. After transfer learning in Gazebo, we can deploy the agent in the real-world environment. This training procedure is described in more detail in [20], where the main difference here is that we have reduced the number of input features, which was done primarily to make the feature attributions simpler to interpret.

D. Dataset and explanations

In order to sample the excluded features, all implementations of SHAP rely on background datasets. Furthermore, we wish to choose interesting decisions by our agent to explain. To this end, we collect a dataset by letting the fully trained DRL agent operate in the real-world environment by running 15 *test episodes* with randomly selected target and start lever angles. From these test episodes, we identify a collection of interesting events and compare the explanations generated by the two different SHAP implementations on these. We remove the episodes where these events occur and use the resulting dataset as our background data set. We chose one event from episode 1 and two events from episode 3, meaning that our background dataset consists of episodes 2 and 4–15.

To display the results, we create a plot similar to the force plot available in the SHAP package³. The force plot illustrates the “force” of each feature on the prediction, showing how the features force the prediction away from the mean prediction and towards the model’s prediction. In our adaptation, we show one force plot for each of the agent’s actions on the same figure. This is primarily done to compactly convey the information and compare the different actions for the same decision.

As stated in Section II-C, the causal SHAP implementation requires the causal structure of the data to be specified by a causal ordering. The causal ordering we use is: $[\theta_{target}], [q_1, q_2, q_3], [q_4, d_x, d_z], [\theta_{lever}]$. In addition, we assume that none of the features are influenced by a confounding factor but instead have only mutual interactions. A visualization of our causal ordering is shown in Figure 4, where blue arrows indicate each feature’s direct effect on the target, purple arrows indirect effects via other features, and the red

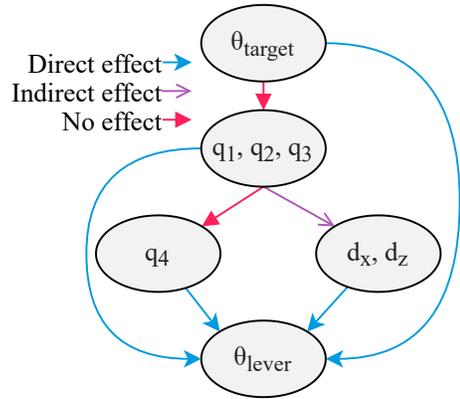


Fig. 4: Visualization of the chosen causal ordering.

arrows indirect effects that we know are not present. This is described below.

In our setup, θ_{target} reveals a weakness in the causal SHAP implementation: This feature does not have any causal connection with any of the other features, but because of the way causal SHAP is implemented, it still has to be defined in the causal ordering. We choose to list this feature first in the causal ordering to avoid the indirect effects of all the other features flowing through this independent feature (which would happen if we placed it after other features). We assume that although this feature is put first in the causal ordering, the causal SHAP algorithm will uncover that this feature only has a direct effect on the prediction, and therefore assign it an indirect causal connection strength close to 0 to the following features in the causal ordering.

Another weakness is that, according to our causal ordering, q_1, q_2, q_3 (the joint features), have a causal effect on q_4 (the gripper feature). This is not necessarily true, but the current implementation does not allow two features to affect a third feature without affecting each other in the causal ordering. These two issues are highlighted in Figure 4, indicating that θ_{target} does not influence features succeeding it in the causal ordering and that the joint variables do not influence q_1 .

IV. RESULTS AND DISCUSSION

In this section, we discuss and compare the results from using causal SHAP and KernelSHAP to explain the actions of the agent performing the robotic lever manipulation task. To do this, we select three events from the dataset described above. The first of these events is from episode 1, where the manipulator pushes the lever from the start angle to the target angle with the gripper closed. We select a time-step in which the manipulator is actively pushing the manipulator and analyze it. This event is hereafter referred to as the *pushing event*. The second event selected is from episode 3, in which the manipulator is grasping the lever before pulling it. We here analyze the exact time-step in which the manipulator grasps the lever. We refer to this event as the *grasping event*. The third and final event we analyze also belongs to episode 3 and takes place just after the grasping event when the manipulator is being used to pull the lever from the start angle to the target angle. This event is referred

³<https://github.com/slundberg/shap>

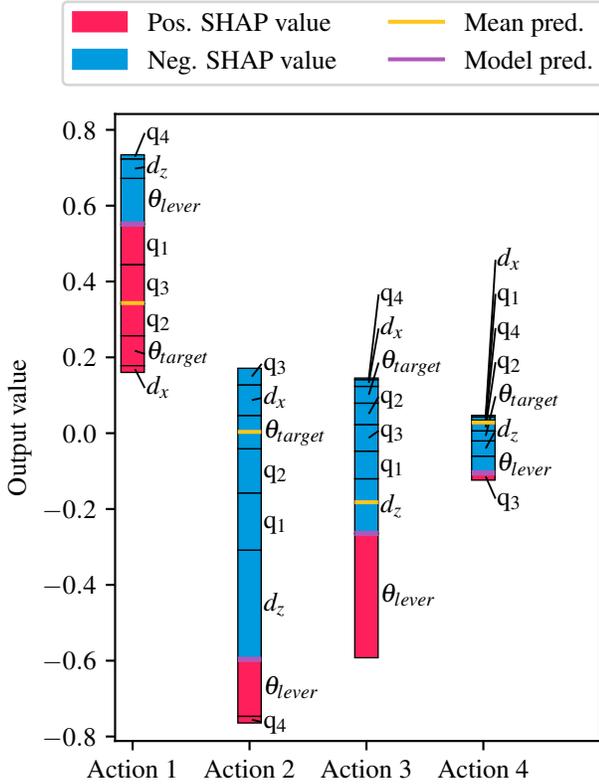


Fig. 5: Episode 1: pushing event, causal SHAP

to as the *pulling event*.

A. Pushing event

For the pushing event, the force plots for causal SHAP are shown in Figure 5 and for KernelSHAP in Figure 6. In both plots, for action 4, we can see that all but one feature have negative SHAP values, which tells us that most aspects of this situation inform the agent to keep the gripper closed. For both methods, and for the actions corresponding to moving the joints (actions 1 to 3), we see that q_4 is the feature with the lowest SHAP value. This means that both methods agree that it is not important whether the gripper is closed for the movement of the joints. Both methods assign similar SHAP values to θ_{lever} and θ_{target} . However, we can see that causal SHAP assigns slightly lower values to θ_{lever} compared to KernelSHAP. This is likely because this feature is at the bottom of the causal ordering. Nevertheless, θ_{lever} is still among the features with the highest SHAP values. This tells us that θ_{lever} is an important predictor variable for how the manipulator should be used to move the lever; in other words, it is important to know where the lever is in order to move it.

In general, we see that causal SHAP gives higher SHAP values to the joint variables, q_1, q_2, q_3 . In contrast, KernelSHAP gives higher values to d_x and d_z , the features that together form a Cartesian vector from the end-effector to the lever's base. The joint variables are higher in the causal ordering, which is why causal SHAP gives more importance to the joint variables. These two sets of features

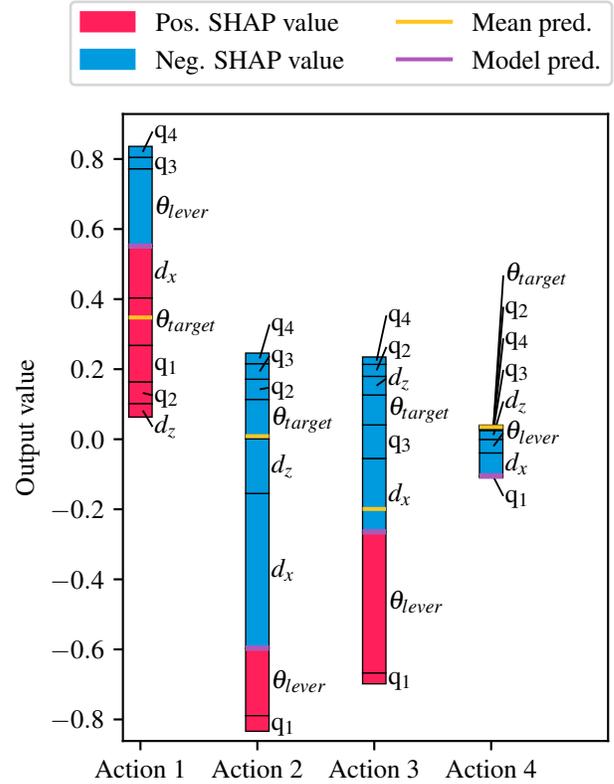


Fig. 6: Episode 1: pushing event, KernelSHAP

contain much of the same information, that is, information about the manipulator's position. However, in addition to this information about the manipulator's position, the joint variables contain information about the manipulator's orientation, while d_x and d_z contain information about where the lever is situated in relation to the manipulator. The exclusive information that these two sets of features contain makes each of them valuable for performing the lever control task. However, because the joint variables are, in fact, what is being controlled by the actions 1 – 3, and it is difficult to control something one does not know where is, it stands to reason that these features should have among the highest SHAP values.

B. Grasping event

For the grasping event, the plots for causal SHAP and KernelSHAP are shown in Figure 7 and Figure 8, respectively. In contrast to the plots for the pushing event, q_4 is not the feature with the lowest SHAP value. This is presumably because the agent intends to pull the lever, and q_4 is important for knowing that the lever still needs to be grasped by the manipulator.

At this point, we see more significant differences between the two methods than we did during the pushing event, with regards to the magnitude of the θ_{lever} feature's SHAP value. Again, KernelSHAP assigns the most importance to this feature. In fact, for causal SHAP, θ_{lever} overall has the lowest SHAP value in the grasping event. This is curious since the position of the lever should be highly important

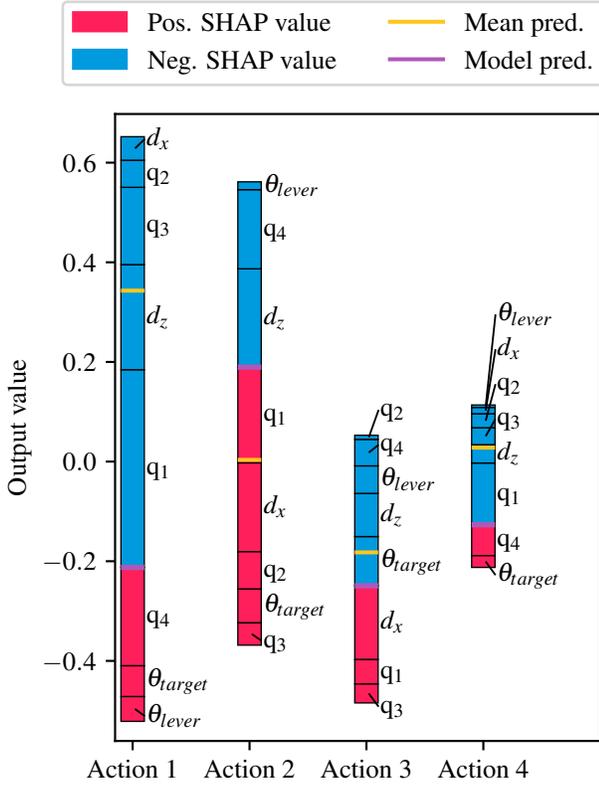


Fig. 7: Episode 3: grasping event, causal SHAP

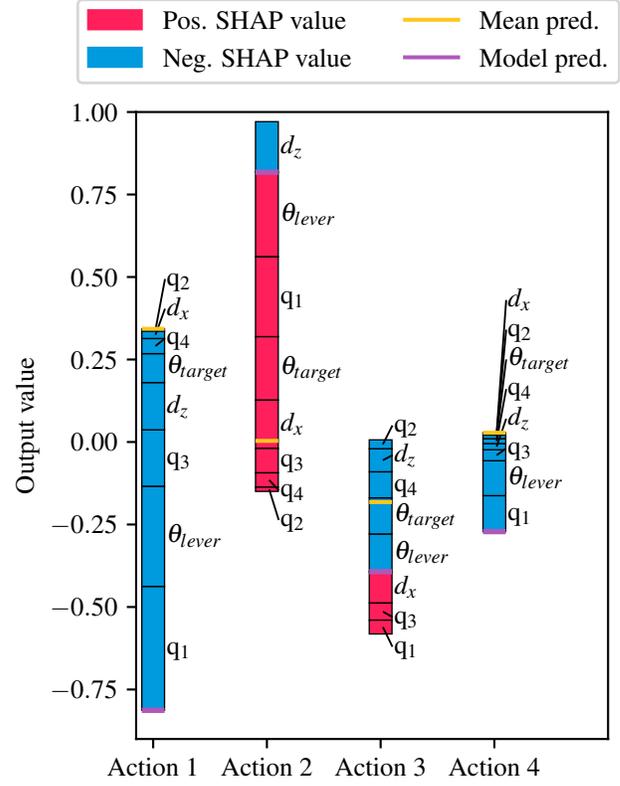


Fig. 9: Episode 3: pulling event, causal SHAP

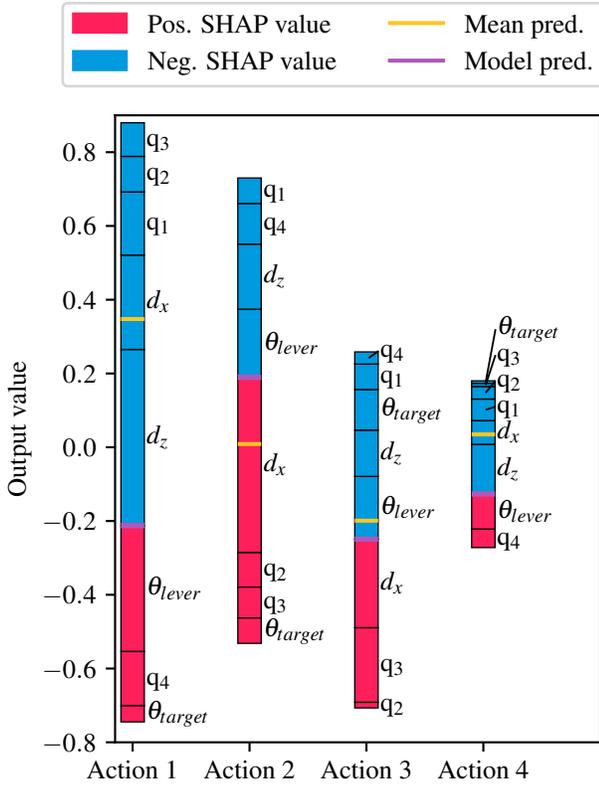


Fig. 8: Episode 3: grasping event, KernelSHAP

for deciding how to grasp the lever. This is likely because causal SHAP assigns some of the contribution from θ_{lever} to features above it in the causal ordering.

Similar to the pushing event, the joint variables are generally more important according to causal SHAP than they are according to KernelSHAP. Conversely, d_x and d_z are generally more important according to KernelSHAP than what they are to causal SHAP. However, there are some exceptions to this: Consider q_2 , which has approximately the same SHAP values for both plots' corresponding actions, and action 3, for which the SHAP value of q_3 has a larger magnitude from KernelSHAP than from causal SHAP.

C. Pulling event

Figure 9 and Figure 10 show the results from causal SHAP and KernelSHAP, respectively. Here, we see even more clearly what we have seen in the two previous events regarding KernelSHAP prioritizing d_x and d_z over the joint variables, and vice versa for causal SHAP. In fact, d_x and d_z account for over half of the total magnitude of all the SHAP values for actions 1, 2 and 4 in Figure 10.

As discussed in Section III-D, θ_{target} was put on the top of the causal ordering by necessity. However, we can see that this feature has approximately the same SHAP value for all actions in the pulling event for both of the methods. This was also the case for the pushing event and the grasping event described above. This suggests that the causal SHAP algorithm has discovered that this feature only has a direct effect, and therefore has given an indirect causal connection

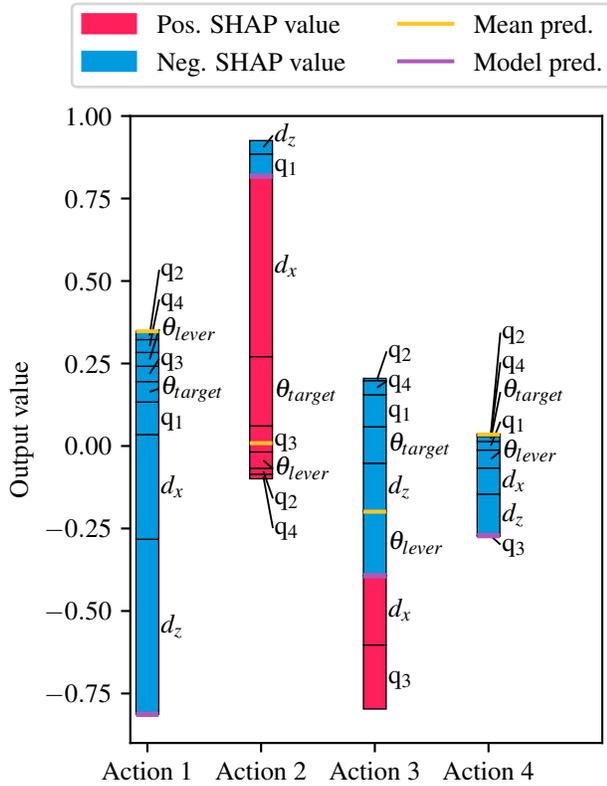


Fig. 10: Episode 3: pulling event, KernelSHAP

strength of approximately 0 to the features succeeding it in the causal graph, in agreement with our expectation.

V. CONCLUSION

We have shown how causal SHAP can be used to explain not only the direct effects but also the indirect effects a feature can have on the decisions of a DRL agent controlling a real-world robotics system. In addition, we have shown how causal SHAP allows for the incorporation of domain knowledge in the explanation generation process. We expect that the causal relations will be an essential part of XAI going forward, especially for explaining models of physical systems.

Because of the complex nature of the explanations shown in this paper, we recognize that these are most useful for data scientists, model developers, and others with experience in data analysis and XAI. When it comes to explanations for end-users that are non-experts, our explanations would likely need to be processed. According to [21], explanations should be contrastive, and explanations based on *counterfactuals* [22] are more intuitive to humans. Work has been done towards unifying feature attribution and counterfactuals in [23], and also towards generating counterfactuals from SHAP, as in [24]. Further work could, therefore, consist of researching whether transforming these feature attribution explanations to counterfactual explanations could make them more convenient for non-experts.

More specific to causal SHAP, further work can, among others, consist in altering the implementation of causal

SHAP to account for fully independent features (such as the θ_{target} in this paper). Another valuable addition to the implementation would be the possibility to specify manually which indirect causal connections should have strength 0 (corresponding to causally independent subsystems).

ACKNOWLEDGMENT

The Research Council of Norway supported this work through the EXAIGON project, project number 304843.

REFERENCES

- [1] D. Silver, T. Hubert, J. Schrittwieser, and D. Hassabis, *AlphaZero Shedding new light on chess, shogi, and Go*, 2018 (accessed 15-Oct-2020). [Online]. Available: <https://deepmind.com/blog/article/alphazero-shedding-new-light-grand-games-chess-shogi-and-go>
- [2] A. P. Badia, B. Piot, S. Kapturovski, P. Sprechmann, A. Vitvitskiy, D. Guo, and C. Blundell, “Agent57: Outperforming the Atari Human Benchmark,” *arXiv:2003.13350 [cs, stat]*, Mar. 2020.
- [3] S. Levine, C. Finn, T. Darrell, and P. Abbeel, “End-to-End Training of Deep Visuomotor Policies,” *arXiv:1504.00702 [cs]*, Apr. 2016.
- [4] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine, “QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation,” *arXiv:1806.10293 [cs, stat]*, Nov. 2018, arXiv: 1806.10293. [Online]. Available: <http://arxiv.org/abs/1806.10293>
- [5] S. Gu, E. Holly, T. Lillicrap, and S. Levine, “Deep Reinforcement Learning for Robotic Manipulation with Asynchronous Off-Policy Updates,” *arXiv:1610.00633 [cs]*, Nov. 2016, arXiv: 1610.00633. [Online]. Available: <http://arxiv.org/abs/1610.00633>
- [6] M. T. Ribeiro, S. Singh, and C. Guestrin, “‘‘Why Should I Trust You?’’: Explaining the Predictions of Any Classifier,” *arXiv:1602.04938 [cs, stat]*, Aug. 2016.
- [7] S. M. Lundberg and S.-I. Lee, “A Unified Approach to Interpreting Model Predictions,” in *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2017, pp. 4765–4774.
- [8] L. S. Shapley, *A VALUE FOR N-PERSON GAMES*. Defense Technical Information Center, 1952.
- [9] T. Heskes, E. Sijben, I. G. Bucur, and T. Claassen, “Causal Shapley Values: Exploiting Causal Knowledge to Explain Individual Predictions of Complex Models,” in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 4778–4789.
- [10] J. Pearl, “Direct and indirect effects,” *arXiv preprint arXiv:1301.2300*, 2013.
- [11] M. Wiering and M. van Otterlo, *Reinforcement Learning: State-of-the-Art*. Springer Publishing Company, Incorporated, 2014.
- [12] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *arXiv:1509.02971 [cs, stat]*, Jul. 2019.
- [13] H. P. Young, “Monotonic solutions of cooperative games,” *International Journal of Game Theory*, vol. 14, pp. 65–72, 1985.
- [14] K. Aas, M. Jullum, and A. Løland, “Explaining individual predictions when features are dependent: More accurate approximations to Shapley values,” *Artificial Intelligence*, vol. 298, p. 103502, 2021.
- [15] C. Frye, C. Rowat, and I. Feige, “Asymmetric shapley values: incorporating causal knowledge into model-agnostic explainability,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1229–1239. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/0d770c496aa3da6d2c3f2bd19e7b9d6b-Paper.pdf>
- [16] D. Janzing, L. Minorics, and P. Bloebaum, “Feature relevance quantification in explainable ai: A causal problem,” in *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, S. Chiappa and R. Calandra, Eds., vol. 108. PMLR, 26–28 Aug 2020, pp. 2907–2916.
- [17] J. Pearl, “Causal diagrams for empirical research,” *Biometrika*, vol. 82, no. 4, pp. 669–688, 1995.
- [18] —, “The do-calculus revisited,” *arXiv preprint arXiv:1210.4852*, 2012.

- [19] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba, "Hindsight Experience Replay," *arXiv:1707.01495 [cs]*, Feb. 2018.
- [20] S. B. Remman and A. M. Lekkas, "Robotic Lever Manipulation using Hindsight Experience Replay and Shapley Additive Explanations," *arXiv:2110.03292 [cs]*, Oct. 2021, arXiv: 2110.03292. [Online]. Available: <http://arxiv.org/abs/2110.03292>
- [21] T. Miller, "Explanation in artificial intelligence: Insights from the social sciences," *Artificial Intelligence*, vol. 267, pp. 1–38, 2019. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0004370218305988>
- [22] D. Hume *et al.*, *An enquiry concerning human understanding: A critical edition*. Oxford University Press, 2000, vol. 3.
- [23] R. Kommiya Mothilal, D. Mahajan, C. Tan, and A. Sharma, "Towards unifying feature attribution and counterfactual explanations: Different means to the same end," *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, May 2021. [Online]. Available: <http://dx.doi.org/10.1145/3461702.3462597>
- [24] S. Rathi, "Generating counterfactual and contrastive explanations using shap," *arXiv preprint arXiv:1906.09293*, 2019.