# Why did I fail? A Causal-based Method to Find Explanations for Robot Failures

Maximilian Diehl and Karinne Ramirez-Amaro

*Abstract*— Robot failures in human-centered environments are inevitable. Therefore, the ability of robots to explain such failures is paramount for interacting with humans to increase trust and transparency. To achieve this skill, the main challenges addressed in this paper are I) acquiring enough data to learn a cause-effect model of the environment and II) generating causal explanations based on that model. We address I) by learning a causal Bayesian network from simulation data. Concerning II), we propose a novel method that enables robots to generate contrastive explanations upon task failures. The explanation is based on setting the failure state in contrast with the closest state that would have allowed for successful execution, which is found through breadth-first search and is based on success predictions from the learned causal model. We assess the sim2real transferability of the causal model on a cube stacking scenario. Based on real-world experiments with two differently embodied robots, we achieve a sim2real accuracy of 70% without any adaptation or retraining. Our method thus allowed real robots to give failure explanations like, 'the upper cube was dropped too high and too far to the right of the lower cube.'

*Index Terms*— Acceptability and Trust, Probabilistic Inference, Learning from Experience

**Fig. 1:** Depicts our method to allow robots to explain their failures. First, a causal model is learned from simulations (steps 1,2). Then, a contrastive explanation is generated by considering the difference between the failure state and closest variable parametrization that would have led to a successful execution (steps 3,4). Finally, the obtained models are transferred and evaluated in two different robots that provide explanations when they commit errors.

## I. INTRODUCTION

One important component in human interactions is the ability to explain one's actions, especially when failures occur [1], [2]. It is argued that robots need this skill if they were to act in human-centered environments on a daily basis [3]. Furthermore, explainability is also shown to increase trust and transparency in robots [1], [2], and the diagnoses capabilities of a robot are crucial for correcting its behavior [4].

There are different types of failures, like task recognition errors (an incorrect action is learned) and task execution errors (the robot drops an object) [5], [6]. In this work, we focus on explaining execution failures. For example, a robot is asked to stack two cubes. Then, the robot will execute this task by picking up a cube and moving its gripper above the goal cube. However, due to sensor and motor inaccuracies, the robot places its gripper slightly shifted to the left, which results in an imperfect cube alignment between the one in hand and the goal cube on the table. Upon opening its gripper, the upper cube lands on the goal but bounces to the table. In such a situation, we expect the robot to reason about what went wrong and generate an explanation based on its previous experience, e.g., 'I failed because the upper cube was dropped too far to the left of the lower cube.'

Maximilian Diehl and Karinne Ramirez-Amaro. Faculty of Electrical Engineering, Chalmers University of Technology, SE-412 96 Gothenburg, Sweden.{diehlm, karinne}@chalmers.se
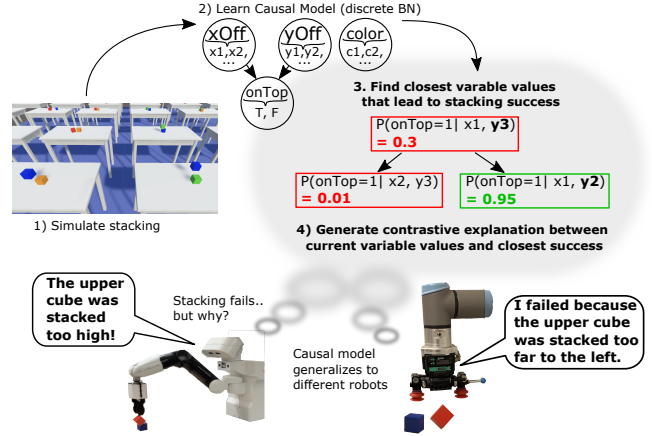
Typically, explanations are based on the concept of causality [7]. Generating explanations is a challenging problem that is addressed based on statistical methods that learn a mapping between possible causes (preconditions) and the action-outcome [4], [8]. However, such statistical models alone are not explanations in itself [1] and require another layer that interprets these models to produce explanations. Another problem is that a considerable amount of data is needed to learn cause-effect relationships. In this case, training such models using a simulated environment will allow a faster and more extensive experience acquisition [4].

In this paper, we propose a method for generating causal explanations of failures based on a causal model that provides robots with a partial understanding of their environment. First, we learn a causal Bayesian network from simulated task executions, tackling the problem of knowledge acquisition. We also show that the obtained model can transfer the acquired knowledge (experience) from simulation to reality and generalize to several real robots with different embodiments. Second, we propose a new method to generate explanations of execution failures based on the learned causal knowledge. Our method is based on a contrastive explanation comparing the variable parametrization associated with the failed action with its closest parametrization that would have led to a successful execution, which is found through breadth-first search (BFS). Finally, we analyze the benefits

of this method on the task of stacking cubes and show that different robots can use the obtained causal model to explain stacking failures (see Fig. 1).

To summarize, our contributions are as follows:

- We present a novel method to generate contrastive causal explanations of action failures based on causal Bayesian networks.
- We demonstrate how causal Bayesian networks can be learned from simulations, exemplified on a cube stacking scenario, and provide extensive real-world experiments that show that the causal model is transferable from simulation to reality without any retraining. Our method can even generalize to various robot platforms with different embodiments. We, thus, show that the simulation-based model serves as an excellent prior experience for the explanations, making them more generally applicable.

## II. RELATED WORK

### A. Causality in Robotics

Despite being acknowledged as an important concept, causality is relatively underexplored in the robotics domain [9], [2]. Some works explore causality to distinguish between task-relevant and -irrelevant variables. For example, CREST [10] uses causal interventions on environment variables to discover which of the variables affect an RL policy. They find that excluding them impacts generalizability and sim-to-real transfer positively. In [11] a set of causal rules is defined to learn to distinguish between unimportant features in physical relations and object affordances. A humanoid iCub robot learns through cumulative experiences that dropping heavy objects into a jar of water will increase the water level, and other variables like color are irrelevant. Brawer et al. present a causal approach to tool affordance learning [9]. One strength of causal models is the ability to learn the causal relevance of variables on a given task. Some works explore Bayesian networks, for example, to learn statistical dependencies between object attributes, grasp actions, and a set of task constraints from simulated data [12]. While the main objective is to use graphical models to generalize task executions, these works don't look into the question of how these models can be utilized for failure explanations. A different paper [13] investigates the problem of learning causal relations between actions in household-related tasks. They discover, for example, that there is a causal connection between opening a drawer and retrieving plates. The learning is based on data that was obtained from human expert demonstrations, which were instructed, for example, to clean the table or wash the dishes in a virtual reality environment. They only retrieve causal links between actions, while we focus on causal relations between different environment variables, like object features and the action outcome.

### B. Learning explainable models of cause-effect relations

In the planning domain, cause-effect relationships are represented through (probabilistic) planning operators [14].

Mitrevksi et al. propose the concept of learning task execution models, which consists of learning symbolic preconditions of a task and a function approximation for the success model [4], based on Gaussian Process models. They noted that a simulated environment could be incorporated for a faster and more extensive experience acquisition, as proposed in [12]. Human virtual demonstrations have been used to construct planning operators to learn cause-effect relationships between actions and observed state-variable changes [14]. However, even though symbolic planning operators are considered human-understandable, they are not explanations in itself, thus requiring an additional layer that interprets the models and generates failure explanations.

Some other works also aim to learn probabilistic action representations experience to generalize the acquired knowledge. For example, learning probabilistic action effects of dropping objects into different containers [8]. Again, the main objective is to find an intelligent way of generalizing the probability predictions for a variety of objects, e.g., bowl vs. bread box, but their method does not include any understanding of why there is a difference in the dropping success probabilities between these different objects. To generalize the obtained models, they consider an ontology to find how closely related these objects are, but they don't consider object variables like the material. A related approach [15] defines object clusters with similar action success probability via ontologies. However, they also note that other information like object affordances or object materials could be used for generalization. Ultimately they also don't look into ways to create explanations.

### C. Contrastive Explanations

Contrastive explanations are deeply rooted in the human way of generation explanations [1]. This also had a significant impact on explanation generation in other fields like Explainable AI Planning (XAIP) [16]. In XAIP, typical questions that a machine should answer are *why a certain plan was generated vs. another one?* or *why the plan contains a particular action $a_1$ and not action $a_2$?* [16], [17]. On the other hand, we are mostly interested in explaining why specific actions failed based on environment variables like object features. Additionally, the type of input information that these methods base their plans on is typically planning operators. Despite containing cause-effect relations, they are different from our causal knowledge representation.

### III. OUR APPROACH TO EXPLAIN FAILURES

Our proposed approach consists of three main steps: A) Identification of the variables used in the analyzed task; B) Learning a Bayesian network which requires to 1) Learn a graphical representation of the variable relations (structure learning) and 2) to learn conditional probability distributions (parameter learning); and C) Our proposed method to explain failures, based on the previously obtained model.

### A. Variable definitions and assumptions

Explaining failures, requires to learn the connections between possible causes and effects of an action. We

describe an action via a set of random variables $\mathbf{X} = \{X_1, X_2, ..., X_n\}$ and require $\mathbf{X}$ to contain a set of treatment variables $C \subset \mathbf{X}$ and outcome (effect) variables $E \subset \mathbf{X}$.

Data samples for learning the causal model can in principle be collected in simulation or the real world. However, it is important to sample values for possible causes $C$ randomly. Randomized controlled trials are referred to as the gold standard for causal inference [3] and allow us to exclude the possibility of unmeasured confounders. Consequently, all detected relations between the variables $\mathbf{X}$ are indeed causal and not merely correlations. Besides the apparent advantage of generating truly causal explanations and avoiding the danger of possible confounders, causal models can also answer interventional questions. In contrast, non-causal models can only answer observational queries. The experiment must satisfy the sampled variable values before executing the action for data collection. $E$ is measured at the end of the experiment. Each action execution generates a datasample which results in a particular variable parametrization, which we denote as $d = \{X_1 = x_{1_d}, X_2 = x_{2_d}, ..., X_n = x_{n_d}\}$. The success of a sample $d$ can be determined by checking if the subset of outcome variable parametrization $d_E \subset d$ satisfies a particular set of goal conditions $d_E = x_{\text{success}} \in X_{\text{goal}}$, where $X_{\text{goal}}$ is a set that defines all possible succesfull goal parametrizations. Note, that it is out of scope of this paper, to discuss methods that learn $X_{\text{goal}}$, but rather assume $X_{\text{goal}}$ to be provided a priori. In other words, we assume that the robot knows how an un-succesfull task execution is defined in terms of its outcome variables $E$ and is thus able to detect it by comparing the action execution outcome with $X_{\text{goal}}$. Note, however, that the robot has no a-priori knowledge about which variables in $\mathbf{X} = X_1, X_2, ..., X_n$ are in $C$ or $E$, nor how they are related. This knowledge is generated by learning the Bayesian network.

To efficiently learn a Bayesian network, some assumptions are needed to handle continuous data [18], mainly because many structure learning algorithms do not accept continuous variables as parents of discrete/categorical variables [19]. In our case, this means that the some effect variables from $E$ could not have continuous parent variables out of $C$, which would likely result in an incorrect Bayesian network structure. As a preprocessing step, we therefore discretize all continuous random variables out of $\mathbf{X}$ into equally sized sets of intervals.

### B. Our proposed pipeline to learn the causal model

Formally, Bayesian networks are defined via a graphical structure $\mathcal{G} = (\mathbf{V}, A)$, which is a *directed acyclic graph* (DAG), where $\mathbf{V} = \{X_1, X_2, ..., X_n\}$ represents the set of nodes and $A$ is the set of arcs [19]. Each node $X_i \subseteq \mathbf{X}$ represents a random variable. Based on the dependency structure of the DAG and the *Markov property*, the *joint probability distribution* of a Bayesian network can be factorized into a set of *local probability distributions*, where each random

variable $X_i$ only depends on its direct parents $\Pi_{X_i}$:

$$P(X_1, X_2, ..., X_n) = \prod_{i=1}^{n} P(X_i | \Pi_{X_i}) \qquad (1)$$

Learning a Bayesian network from data consists of two steps:

*1) Structure Learning:* The purpose of this step is to learn the graphical representation of the network $\mathcal{G} = (\mathbf{V}, A)$ and can be achieved by a variety of different algorithms. For the remainder of this paper, we choose the Grow-Shrink [20] algorithm (gs) to learn $\mathcal{G}$. gs falls in the category of *constraint-based-algorithms*, which use statistical tests to learn conditional independence relations (also called "constraints") from the data [21]. Choosing the 'best' algorithm is data-dependent and therefore not further discussed in this paper. However, the interested reader is referred to [21] for a more in-depth discussion on the different learning algorithms. In the following, we assume that the outcome of the structure learning step is indeed a correct Bayesian network graph $\mathcal{G}$. If manual evaluation results in a rebuke of the structure, it might be required to collect more data samples or tune the number of discretization steps.

*2) Parameter Learning:* The purpose of this step is to fit functions that reflect the *local probability distributions*, of the factorization in formula (1). We utilize the maximum likelihood estimator for conditional probabilities (*mle*) to generate a conditional probability table.

### C. Our proposed method to explain failures

Our proposed method to generate contrastive failure explanations uses the obtained causal Bayesian network to compute success predictions and is summarized in algorithm 1.

---

**Algorithm 1** Failure Explanation

**Input:** failure variable parameterization $x_{\text{failure}}$, graphical model $\mathcal{G}$, structural equations $P(X_i | \Pi_{X_i})$, discretization intervals of all model variables $X_{\text{int}}$, success threshold $\epsilon$

**Output:** solution variable parameterization $x_{\text{solution}_{\text{int}}}$, solution success probability prediction $p_{\text{solution}}$

1: $x_{\text{current}_{\text{int}}} \leftarrow$ GETINTERVALFROMVALUES$(x_{\text{failure}}, X_{\text{int}})$
2: $P \leftarrow$ GENERATETRANSITIONMATRIX$(X_{\text{int}})$
3: $q \leftarrow [x_{\text{current}_{\text{int}}}]$
4: $v \leftarrow []$
5: **while** $q \neq \emptyset$ **do**
6:      $node \leftarrow$ POP$(q)$
7:      $v \leftarrow$ APPEND$(v, node)$
8:      **for each** transition $t \in P(node)$ **do**
9:          $child \leftarrow$ CHILD$(P, node)$
10:          **if** $child \notin q, v$ **then**
11:              $p_{\text{solution}} = P(E = x_{\text{success}} | \Pi_E = child)$
12:              **if** $p_{\text{solution}} > \epsilon$ **then**
13:                  $x_{\text{solution}_{\text{int}}} \leftarrow child$
14:                  RETURN$(p_{\text{solution}}, x_{\text{solution}_{\text{int}}})$
15:          $q \leftarrow$ APPEND$(q, x_{\text{current}_{\text{int}}})$

---

In (L-2 Alg. 1)), a matrix is generated which defines transitions for every single-variable change for all possible

variable parametrizations. For example, if we had two variables $X_1, X_2$ with two intervals $x', x''$. Then, the possible valid transitions for $node = (X_1 = x', X_2 = x')$ would be $child_1 = (X_1 = x', X_2 = x'')$ or $child_2 = (X_1 = x'', X_2 = x')$.

Lines 5-15 (Alg. 1) describe the adapted BFS procedure, which searches for the closest variable parametrization that fulfills the goal criteria of $P(E = x_{success}|\Pi_E = child) > \epsilon$, where $\epsilon$ is the success threshold, which can be heuristically set.

The explanation generation requires further elaboration. The concept of our proposed method is to generate contrastive explanations that compare the current variable parametrization associated with the execution failure $x_{current_{int}}$ with the closest parametrization that would have allowed for a succesfull task execution $x_{solution_{int}}$. Consider Figure 2 for a visualization of the explanation generation, exemplified on two variables $X$ and $Y$, which are both causally influencing the variable $X_{out}$. Furthermore, it is known that $x_{out} = 1 \in X_{goal}$. The resulting explanation would be that the task failed because $X = x_1$ instead of $X = x_2$ and $Y = y_4$ instead of $Y = y_3$.



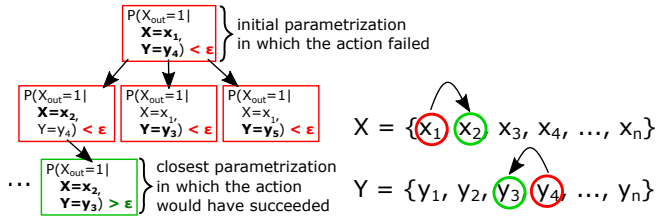**Fig. 2:** Exemplifies how contrastive explanations are generated from the BFS search tree.

## IV. EXPERIMENTS

We evaluate our method based on the stacking cubes scenario. The environment contains two cubes: *CubeUp* and *CubeDown* (see Fig. 3.a). The goal of the stacking action is to place CubeUp *on top* of CubeDown. We define six variables as follows: $\mathbf{X} = \{\texttt{xOff}, \texttt{yOff}, \texttt{dropOff}, \texttt{colorDown}, \texttt{colorUp}, \texttt{onTop}\}$. Both cubes have an edge length of 5cm.
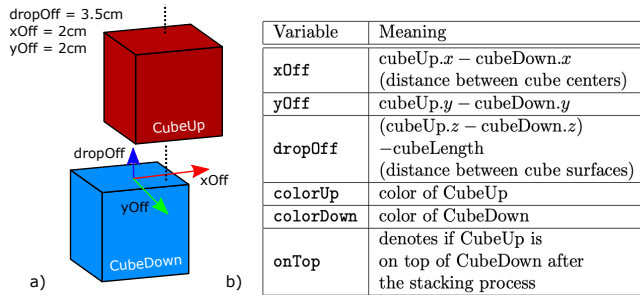


**Fig. 3:** a) visualizes the used variables $\mathbf{X}$ and b) describes their meaning.

### A. Simulation setup

We run the simulations in Unity3d which bases its physics behavior on the Nvidia PhysX engine. For training the Bayesian network we generate 20,000 samples, on 500 parallel table environments (see Fig. 1). We randomly sample values for $\texttt{xOff}$, $\texttt{yOff} \sim \mathcal{U}_{[-0.03,0.03]}$ (in meter), $\texttt{dropOff} \sim \mathcal{U}_{[0.004,0.1]}$ (in meter), $\texttt{colorUp}$, $\texttt{colorDown} = \{\text{Red}, \text{Blue}, \text{Green}, \text{Orange}\}$. $\texttt{onTop} = \{\text{True}, \text{False}\}$ is not sampled but automatically determined after the stacking process.

### B. Robot Experiments setup

We run and assess our experiments on two different robotic platforms (Fig. 1): the TIAGo service robot with one arm and parallel gripper, and the UR3 with a piCOBOT suction gripper. The real cubes are 3D printed out of PLA (polylactic acid) and weight around 25 grams each. For each robot we run 180 stacking trials. Instead of randomly sampling values for the variables, as we do for training the causal model, we evaluate the real world behavior at 36 different points, where $\texttt{xOff}_e, \texttt{yOff}_e = \{0.0, 0.01, 0.02\}$ (in meter), $\texttt{dropOff}_e = \{0.005, 0.02, 0.035, 0.05\}$ (in meter), $\texttt{colorUp}_e = \{\text{Red}\}$ and $\texttt{colorDown}_e = \{\text{Blue}\}$. For each unique stacking setup instantiation we conduct 5 iterations. After each trial, the cubes are re-adjusted into an always similar pre-stack position by the operator. The stacking outcome ($\texttt{onTop}$ value) was also determined by the operator. Note, that the purpose of the robot experiments is not to modify the causal model that we learned from the simulation but to evaluate the model transferability to the real environment.

## V. RESULTS AND DISCUSSION

### A. Analysis of the obtained causal model

We first present and discuss the learned causal model of the stacking task simulations. The structure learning was performed based on the Grow-Shrink algorithm [20]. 10-fold cross validation reports an average loss of 0.10269 and a standard deviation of 0.00031. Figure 4 displays the resulting DAG, which shows the relations between the variables $X = \{\texttt{xOff}, \texttt{yOff}, \texttt{dropOff}, \texttt{colorDown}, \texttt{colorUp}, \texttt{onTop}\}$. The graph indicates, that there are causal relations from $\texttt{xOff}, \texttt{yOff}$ and $\texttt{dropOff}$ to $\texttt{onTop}$, while the two color variables $\texttt{colorDown}$ and $\texttt{colorUp}$ are independent. In other words, it makes a difference from which position the cube is dropped, but the cube color has no impact on the stacking success.
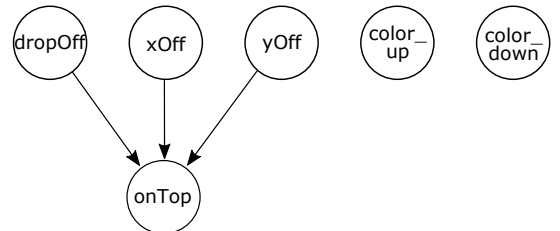


**Fig. 4:** Obtained bayesian network structure.

The following table shows the obtained intervals for the three variables that affects `onTop`:

| dropOff (in m) | xOff (in m) | yOff (in m) |
|---|---|---|
| $z_1 : [.004, 0.018]$ | $x_1 : [-0.03, -0.018]$ | $y_1 : [-0.03, -0.018]$ |
| $z_2 : (.018, 0.032]$ | $x_2 : (-0.018, -0.006]$ | $y_2 : (-0.018, -0.006]$ |
| $z_3 : (.032, 0.045]$ | $x_3 : (-0.006, 0.006]$ | $y_3 : (-0.006, 0.006]$ |
| $z_4 : (.045, 0.059]$ | $x_4 : (0.006, 0.018]$ | $y_4 : (0.006, 0.0179]$ |
| $z_5 : (.059, 0.073]$ | $x_5 : (0.018, 0.03]$ | $y_5 : (0.0179, 0.03]$ |
| $z_6 : (.073, 0.086]$ | | |
| $z_7 : (.086, 0.099]$ | | |

The conditional probabilities $P(\texttt{onTop} = 1|\Pi_{\texttt{onTop}})$ are visualized in Fig. 5. These plots allow us to conclude that stacking success decreases the greater the drop-offset and the more offset in both x- and y-direction. In particular, there is a diminishing chance of stacking success for the values |xOff| > 0.018 or |yOff| > 0.018, no matter the dropOff height. Therefore, these obtained probabilities resemble our intuitive understanding of the physical processes and interactions involved in the cube stacking task. Nevertheless, real-world experiments have a higher complexity due to the many environment uncertainties. We therefore, expect the simulation to be less conservative than reality, as we have higher control over the variables involved in the stacking process. The upper cube, for example, will always hit the lower cube perfectly flat. In reality, small grasp position offsets from the cube's center of mass can lead to slightly rotated cube impact positions, with consequences regarding the stacking success.

### B. Assessment of the causal models in the real robots

To evaluate how well the causal model and the real-world match, we now evaluate the experiment outcome of the robot stacking trials. The results are presented in Fig. 6. The dark points indicate the nine stacking locations (all possible combinations of x- and y-offset values) for each of the four drop-off heights. The plots show contours of the probabilities, meaning the stacking success probabilities are interpolated between the nine measurement points. For the comparison between Fig. 5 and Fig. 6, note that `xOff` and `yOff` only cover positive values in the real-world experiments and, therefore, only cover the first quadrant of the simulation plots. Furthermore, the simulation data was simulated till 3 cm offset while the real-world experiments only cover a maximal offset of 2cm in x- and y-direction. We compute a sim2real accuracy score, which reflects how much the probability predications of the causal model coincide with the real-world measurements via to following formula:

$$
\text{sim2real} = \frac{1}{|\texttt{xOff}_e||\texttt{yOff}_e||\texttt{dropOff}_e|} \sum_{x=1}^{|\texttt{xOff}_e|} \sum_{y=1}^{|\texttt{yOff}_e|} \sum_{d=1}^{|\texttt{dropOff}_e|}
$$
$$
\Big( 1 - \Big( |P_{real}(\texttt{onTop} = 1|\texttt{xOff}_e = x, \texttt{yOff}_e = y,
$$
$$
\texttt{dropOff}_e = d) - P_{sim}(\texttt{onTop} = 1|\texttt{xOff} = x,
$$
$$
\texttt{yOff} = y, \texttt{dropOff} = d)| \Big) \Big)
$$
$$(2)$$

The sim2real accuracy amounts to 71% for the TIAGo and 69% for the UR3. We can conclude that the probability model obtained from simulated data can be utilized to make reasonable predictions regarding real-world stacking success and that the model generalizes well to differently embodied robots. We want to emphasize that the causal model was not retrained or adapted. We collected the real-world data purely to evaluate the sim2real accuracy of the model. The largest discrepancy between model and reality can be determined for the higher drop off positions. For the real world measurements the stacking success falls significantly, already for dropOff heights of 2cm or 3.5cm. In the model, there is a significantly drop much later, starting at around 5cm. This can be mainly attributed to the increased complexity of real-world action executions.

It is also interesting to compare similarities regarding probability outcomes between the two differently embodied robots. The correspondence concerning the 36 measured positions amounts to 85%. If we had a lower sim2real accuracy or more significant differences between the two robots, it would be advisable to include robot-specific variables (such as the gripper type and orientation) and adapt the model with real-world data. But even then, the model that we obtain from the simulation can be used as an excellent experience prior, allowing for faster applicability and learning.

### C. Explanability capabilities

Finally, we tested the explanation of failures on several examples (three examples are shown in Tab. I). We set the probability threshold which distinguishes a failure from a success to $\epsilon = 0.8$.

| input | input interval | curr. succ. probability | closest solution interval | exp. succ. probability |
|---|---|---|---|---|
| **Example 1:** | | | | |
| xOff = 0.0<br>yOff = 0.02<br>dropOff = 0.08 | $x_3$<br>$\mathbf{y_5}$<br>$\mathbf{z_6}$ | 0.0 | $x_3$<br>$\mathbf{y_4}$<br>$\mathbf{z_3}$ | 0.85 |
| **Explanation**: The upper cube was dropped too high and too far to the front of the lower cube. | | | | |
| **Example 2:** | | | | |
| xOff = 0.015<br>yOff = 0.0<br>dropOff = 0.05 | $x_4$<br>$y_3$<br>$\mathbf{z_4}$ | 0.58 | $x_4$<br>$y_3$<br>$\mathbf{z_3}$ | 0.91 |
| **Explanation**: The upper cube was dropped too high. | | | | |
| **Example 3:** | | | | |
| xOff = -0.015<br>yOff = 0.015<br>dropOff = 0.02 | $\mathbf{x_1}$<br>$\mathbf{y_1}$<br>$z_2$ | 0.017 | $\mathbf{x_2}$<br>$\mathbf{y_2}$<br>$z_2$ | 1.0 |
| **Explanation**: The upper cube was dropped too far to the left and too far to the back of the lower cube. | | | | |

**TABLE I:** Three examples of failure explanations.

Example 2 is particularly interesting, as it showcases that there are often multiple correct explanations for the error. In this case it would have been possible to achieve a successful stacking by either going from `dropOff` = $z_4$ to `dropOff` = $z_3$ or by changing `xOff` = $z_4$ to `xOff` = $z_3$ (search tree is visualised in Fig. 7). Which solution is found first depends on the variable prioritization within the tree search due to the used BFS algorithm.
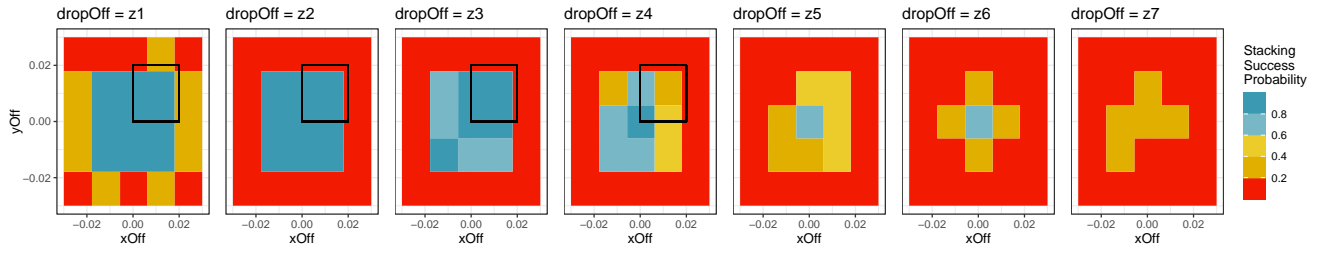
**Fig. 5:** Visualisation of the conditional probability table for $P(\texttt{onTop} = 1|\Pi_{\texttt{onTop}})$. $\texttt{xOff}$, $\texttt{yOff}$ are discretized into 5 intervals and $\texttt{dropOff}$ into 7. Values for $\texttt{xOff}$, $\texttt{yOff}$ are in meter. For easier comparison between the causal model and the real-world experiments (Figure 6), the black rectangles denote the $\texttt{xOff}$ and $\texttt{yOff}$ value range evaluated in the real-world experiments.
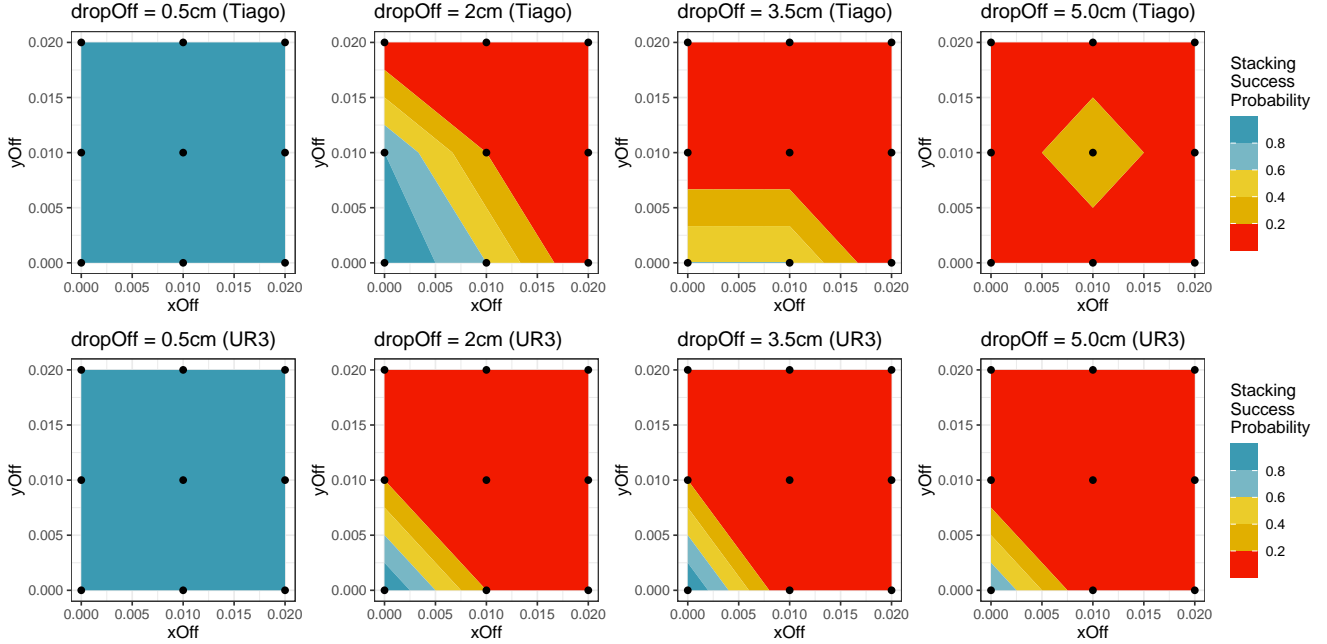


**Fig. 6:** Shows real-world probabilities for $P(\texttt{onTop} = 1|\Pi_{\texttt{onTop}})$, evaluated for two robots (TIAGo with parallel gripper and UR3 with suction gripper) at 36 different points, with $\texttt{xOff}_e, \texttt{yOff}_e = \{0.0, 0.01, 0.02\}$ (in meter), $\texttt{dropOff}_e = \{0.005, 0.02, 0.035, 0.05\}$ (in meter). For each variable parametrization, 5 experiments were conducted. The plots show contours of the probabilities, meaning the stacking success probabilities are interpolated between the nine measurement points for each $\texttt{dropOff}$ value. For easier comparison with the causal model probabilities, the black rectangles in Figure 5 denote the $\texttt{xOff}$ and $\texttt{yOff}$ value range of the real-world experiments.
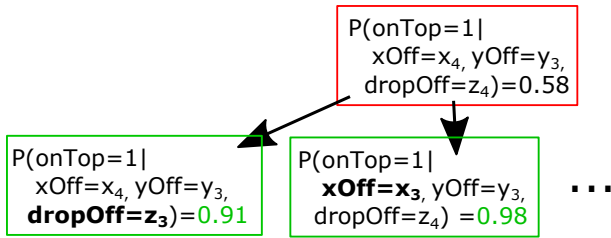


**Fig. 7:** BFS for explaining example 2 from Tab. I

## VI. CONCLUSION

This paper presents and demonstrates our novel approach to provide robots with the ability to explain failures based on a causal model of the environment. First, we learn a causal Bayesian network from simulated task executions. We show that it is possible to transfer this knowledge from simulation to reality and generalize the model to several real robots with a different embodiment with 70% sim2real accuracy. Furthermore, we propose a new method to generate explanations of execution failures based on the causal knowledge. This method is based on a contrastive explanation comparing the action parametrization of the failure with its closest parametrization that would have led to a successful execution, which is found through breadth-first search (BFS).

The obtained causal models could be used not just to explain failures but, in general, justify why the robot has acted in a specific way. For example, if asked why the robot has stacked a cube the way it did, it could answer that the cube would have dropped to the table if it was stacked more to the right. This is considered as future work.

## REFERENCES

[1] T. Miller, "Explanation in artificial intelligence: Insights from the social sciences," *Artif. Intell.*, vol. 267, pp. 1–38, 2019.

[2] T. Hellström, "The relevance of causation in robotics: A review, categorization, and analysis," *Paladyn, Journal of Behavioral Robotics*, vol. 12, no. 1, pp. 238–255, 2021. [Online]. Available: https://doi.org/10.1515/pjbr-2021-0017

[3] J. Pearl and D. Mackenzie, *The Book of Why: The New Science of Cause and Effect*, 1st ed. USA: Basic Books, Inc., 2018.

[4] A. Mitrevski, P. G. Plöger, and G. Lakemeyer, "Representation and experience-based learning of explainable models for robot action execution," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 5641–5647.

[5] L. K., S. Y., K. T.-K., and D. Y., "A Syntactic Approach to Robot Imitation Learning Using Probabilistic Activity Grammars," *Robot. Auton. Syst.*, vol. 61, no. 12, pp. 1323–1334, Dec. 2013. [Online]. Available: http://dx.doi.org/10.1016/j.robot.2013.08.003

[6] K. S. and S. S., "Cognitive robots learning failure contexts through real-world experimentation," *Autonomous Robots*, no. 4, pp. 469–485, Dec. 2015.

[7] D. Lewis, "Causal explanation," in *Philosophical Papers Vol. Ii*, D. Lewis, Ed. Oxford University Press, 1986, pp. 214–240.

[8] A. S. Bauer, P. Schmaus, F. Stulp, and D. Leidner, "Probabilistic effect prediction through semantic augmentation and physical simulation," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 9278–9284.

[9] J. Brawer, M. Qin, and B. Scassellati, "A causal approach to tool affordance learning," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 8394–8399.

[10] T. E. Lee, J. A. Zhao, A. S. Sawhney, S. Girdhar, and O. Kroemer, "Causal reasoning in simulation for structure and transfer learning of robot manipulation policies," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 4776–4782.

[11] A. A. Bhat, V. Mohan, G. Sandini, and P. G. Morasso, "Humanoid infers archimedes' principle: understanding physical relations and object affordances through cumulative learning experiences," *Journal of the Royal Society Interface*, vol. 13, 2016.

[12] D. Song, K. Huebner, V. Kyrki, and D. Kragic, "Learning task constraints for robot grasping using graphical models," in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2010, pp. 1579–1585.

[13] C. Uhde, N. Berberich, K. Ramirez-Amaro, and G. Cheng, "The robot as scientist: Using mental simulation to test causal hypotheses extracted from human activities in virtual reality," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 8081–8086.

[14] M. Diehl, C. Paxton, and K. Ramirez-Amaro, "Automated generation of robotic planning domains from observations," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 6732–6738.

[15] A. Mitrevsk, P. G. Plöger, and G. Lakemeyer, "Ontology-assisted generalisation of robot action execution knowledge," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 6763–6770.

[16] T. Chakraborti, S. Sreedharan, and S. Kambhampati, "The emerging landscape of explainable automated planning & decision making," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, ser. IJCAI'20, 2021.

[17] B. Seegebarth, F. Müller, B. Schattenberg, and S. Biundo, "Making hybrid plans more clear to human users — a formal approach for generating sound explanations," in *Proceedings of the Twenty-Second International Conference on International Conference on Automated Planning and Scheduling*, ser. ICAPS'12. AAAI Press, 2012, p. 225–233.

[18] Y.-C. Chen, T. A. Wheeler, and M. J. Kochenderfer, "Learning discrete bayesian networks from continuous data," *J. Artif. Int. Res.*, vol. 59, no. 1, p. 103–132, may 2017.

[19] M. Scutari, "Learning bayesian networks with the bnlearn R package," *Journal of Statistical Software*, vol. 35, no. 3, pp. 1–22, 2010.

[20] D. Margaritis, "Learning bayesian network model structure from data," Tech. Rep., 2003.

[21] R. Nagarajan and M. Scutari, *Bayesian Networks in R with Applications in Systems Biology*. New York: Springer, 2013, iSBN 978-1-4614-6445-7, 978-1-4614-6446-4.