# CausalCity: Complex Simulations with Agency for Causal Discovery and Reasoning

Anonymous Author(s) Affiliation Address email

## Abstract

1 The ability to perform causal and counterfactual reasoning are central properties 2 of human intelligence. Decision-making systems that can perform these types of 3 reasoning have the potential to be more generalizable and interpretable. Simulations have helped advance the state-of-the-art in this domain, by providing the ability to 4 systematically vary parameters (e.g., confounders) and generate examples of the 5 outcomes in the case of counterfactual scenarios. However, simulating complex 6 temporal causal events in multi-agent scenarios, such as those that exist in driving 7 and vehicle navigation, is challenging. To help address this, we present a high-8 fidelity simulation environment that is designed for developing algorithms for 9 causal discovery and counterfactual reasoning in the safety-critical context. A core 10 component of our work is to introduce *agency*, such that it is simple to define and 11 12 create complex scenarios using high-level definitions. The vehicles then operate with agency to complete these objectives, meaning low-level behaviors need only 13 be controlled if necessary. We perform experiments with three state-of-the-art 14 methods to create baselines and highlight the affordances of this environment. 15 Finally, we highlight challenges and opportunities for future work. 16

## 17 **1** Introduction

18

Modern machine learning algorithms perform well on clearly defined pattern recognition tasks but 19 still fall short *generalizing* in the ways that human intelligence can [1, 2]. This leads to unsatisfactory 20 results on tasks that require extrapolation from training examples, e.g., out-of-domain recognition [3] 21 22 and open set recognition [4]. Causal reasoning sets human intelligence apart from pattern matching [5] and enables us to answer counterfactual questions such as "what would have happened if..." Reasoning 23 24 such as this is not only important in helping create learning algorithms robust to generalization but is also attractive in applications that require transparent and/or explainable decision making (e.g., safety 25 critical scenarios including medical decision making [6] and autonomous driving [7]). 26

Discovering latent causal mechanisms and handling confounders are the key tasks in causal reasoning. 27 Confounders refer to factors that impact both the intervention and the outcomes [8]. These factors 28 can be "measureable" in some cases and hidden in others. If confounders are hidden then it is 29 difficult to control for them. Ideally, we would have the ability to systematically examine the impact 30 of many different types/classes of confounders both "hidden" and measurable whilst developing 31 causal interference algorithms. Furthermore, we would like the ability to do so in contexts that 32 mirror or match our real-world applications. Recent approaches to causal reasoning involve capturing 33 causal structure and disentangling the underlying factors via an inference algorithm (e.g., neural 34

Project page: https://causalcity.github.io/



Figure 1: We present a high-fidelity simulation environment designed for experiments on causal reasoning in the safety-critical context of driving. The vehicles have agency to "decide" their lowlevel behaviors, which enables scenarios to be designed with simple high-level configurations. Many complex simulated scenarios can be executed with complex causal relationships. Our environment then supports the logging of rich multimodal signals during simulation for forming datasets.

model) [8, 9, 10, 11] and combining this with a graphical representations (e.g., directed acyclic 35 graphs - DAG) to capture the underlying dynamics. These have been employed successfully to make 36 long-term future predictions based on short observations [10]. 37

Causal inference could make a significant impact in safety critical scenarios such as autonomous 38 driving [12, 13] where trajectory prediction is an important component. Researchers have used video 39 and synthetic datasets [14, 7, 15, 16] for analyzing causality in traffic accidents and understanding 40 driving scenes and behaviors. However, some of these datasets are "static" (i.e., comprised of 41 videos that cannot be changed) which means that certain counterfactual scenarios are not present 42 and the distribution of events may be quite uneven/sparse making it difficult to learn relationships. 43 Other datasets have limited diversity in terms of the types of events, e.g., focusing on crashes 44 specifically [15]. Our goal in this work is not to propose a new trajectory prediction algorithm but 45 rather to illustrate how CausalCity can be used and trajectory prediction is a good task to do so. 46

47 48



61 62 63

Figure 2: Simulation is a powerful tool to 64 study causal reasoning. Here we show exam-65 ples of environments used for causal reasoning. 66 A) V-CDN [10], B) CLEVRER [9], C) Causal-67 World [17]. In contract with prior work, D) Causal-68 City (Ours) combines a high-fidelity visual envi-69 ronment with the ability to define and generate 70 complex causal scenarios. 71 72

Simulation has proven helpful as a way of investigating problems involving causal and counterfactual reasoning. The parameters of synthetic environments can be systematically controlled, thereby enabling causal relationships to be established and confounders to be introduced [10, 9, 17]. However, some of this prior work has approached this via a relatively simplistic set of entities and environments (e.g., balls moving in 2D connected via rods and springs [10] or 3D objects moving on a surface and colliding [9] - see Fig. 2) with only a few variables. Other prior work has had a limited number of degrees of freedom [17]. This leaves little room to explore, and control for, different causal relationships among entities. We posit that enabling the *agency* on each entity is crucial to creating simulation environments that reflect the nature and complexity of these types of temporal real-world reasoning tasks. This includes scenarios where each entity makes decisions on its own while interacting with each other, e.g., pedestrians in a crowded street and cars on a busy road. What agency provides is the ability to define scenarios at a higher level, rather than specifying every single low-level action.

To this end, we introduce and publicly release a high-fidelity simulation environment, summarized 73 in Fig. 1, with AI agent controls to create scenarios for causal and counterfactual reasoning. This 74 environment reflects the real-world, safety critical scenario of driving. We want a simulation 75

re environment that enables controllable scenario generation that can be used for temporal and causal

77 reasoning. This environment allows us to create complex scenarios including different types of 78 confounders with relatively little effort.

To illustrate this, our simulation engine allows the introduction of any number of vehicles, each 79 of which is controlled at a high-level and has basic AI agency to maneuver avoiding collisions, 80 navigating corners, stopping at traffic lights, etc. The high-level controls for each vehicle allow us to 81 define each agent's behavior in an abstract form controlling their sequence of actions (e.g., turn left at 82 the next intersection, following that merge into the left lane etc.), their speed changes in different 83 legs of the journey, their stopping distance behind other vehicles etc. Furthermore, our simulation 84 can be used to introduce confounders to the environment such as the time of day and the weather 85 conditions, which can be set both changing the visual appearance of the scene but also enabling 86 causal relationships to be introduced (for example between vehicle speed or stopping distance and the 87 amount of water on the roads). Again, agency helps vehicles to change their behavior dynamically 88 from the confounders. Also, traffic lights can be controlled at a low (the timing of each individual 89 light) and high (transition timings for all the lights) levels. All these present opportunities for future 90 work on causal reasoning. In this work we perform a set of experiments on causal discovery; however, 91 we give other examples of how the simulation might be used in Section 6 and on our project page. 92

To summarize, our contributions include: 1) We present a high-fidelity simulation driving environ-93 ment with vehicles (agents) that is designed for developing and testing approaches for causal and 94 counterfactual reasoning. 2) We test benchmark causal inference algorithms on trajectory prediction 95 and causal discovery tasks. 3) We demonstrate how this environment can be used to systematically 96 synthesize data to introduce complex confounders and illustrate how these impact the performance of 97 our baselines. 4) Our environment, a snapshot of the dataset used for analyses, and code are released 98 with this paper (see GITHUB link on the first page). Our simulation allows for the generation of 99 large, multimodal, complex causal datasets within the domain of vehicle navigation and we hope that 100 it will enable researchers to tackle new research problems. 101

# **102 2 Related Work**

**Causal Reasoning.** Schölkopf [2] argues that causality and the "modeling and reasoning about interventions" can help advance machine learning as a whole and contribute to addressing some of the most challenging problems, including domain-transfer, extrapolation and other forms of generalization beyond what is explicitly observed in training datasets. These are some of the reasons that causal reasoning has received growing attention in the machine learning community.

Variational autoencoders (VAE) have been used to capture the causal structure in interactions [8, 18] 108 due to their ability to model uncertainty in data. The encoder can be used to estimate an unknown 109 latent space in order to summarize the causal effects [8] and summarize or disentangle representations 110 of objects or events of interest from confounders [19]. Neural Relational Inference (NRI) [18] train 111 an unsupervised VAE, where the latent representation captures the underlying interaction graph. The 112 approach then learns to simultaneously capture the dynamics and infer interactions. This NRI work 113 114 and others leverage a graph neural network for reconstruction [20, 21]. Bhattacharya et al. [21] cast 115 causal discovery as a continuous optimization problem with differentiable constraints to find the 116 best fitting acyclic directed mixed graph. V-CDN [10] discovers an underlying causal graph without explicit intervention in the scene and identifies interactions between entities from a short sequence of 117 images and make long-term future predictions. 118

Simulation for Causal Reasoning. Computer graphics-based simulations have allowed researchers 119 to explore the causality in video. PhysNet [22] learns physics by using a 3D game engine to create 120 121 small towers of wooden blocks with randomized stability. Happens [23] focuses on understanding the movements of objects as a result of applying external forces to them. A large-scale dataset of forces 122 in scenes is built by reconstructing all images in SUN RGB-D dataset [24] in a physics simulator to 123 estimate the physical movements of objects caused by external forces applied to them. Billiards [25] 124 learns to play a simulated billiards game, which requires planning and executing goal-specific actions 125 in varied and unseen environments. 126

Johnson et al. [26] introduced the CLEVR as a simulation engine for visual question answering (VQA).
While this featured static images of objects with different shapes and colors, it was followed by the
CLEVRER [9], which allows for generating objects that move and collide with one another in a 3D

environment. Specifically, this enables counterfactual reasoning to be conducted. Causal World [17] is

another recent example of simulation for causal reasoning based on a robotic manipulation benchmark.
 This is a 3D environment that exposes high-level variables in the causal generative model, such as

133 properties of blocks, goals, robot links and others like gravity.

**Driving simulation.** Our proposed simulation environment is designed to be applicable to studying 134 causal reasoning generally; however, the specific environment we chose is that of driving. The 135 driving scenario has been used to generate synthetic data, e.g., GTACrash [15] and VIENA [16]. The 136 former involves generating a dataset for detecting car accident, whereas the latter involves generating 137 specific actions for predicting driver maneuvers, pedestrian intentions, front car intentions, traffic rule 138 violations, and accidents scenarios. CARNOVEL [27] is a driving benchmark specifically targeting 139 out-of-distribution generalization using adaptive robust imitative planning (AdaRIP) and DESIRE [28] 140 involves reasoning about scenes, context and past trajectories to predict future trajectories or locations. 141 R2PR [29] and DATF [30] also approach future trajectory forecasting a task that is relevant in the 142 context of causal reasoning in driving simulation. However, these examples are not specifically 143 designed around the idea of causal discovery containing no counterfactual reasoning or control for 144 confounders yet in their simulations or datasets. Also there are simulators [31, 32] that support 145 development of autonomous cars but do not focus on causal reasoning. 146

Causal reasoning for autonomous driving has been attended to understand the reason of the maneuvers 147 of other vehicles and pedestrians for escaping accidents [14, 7, 15, 16]. This is a popular domain 148 for causal analysis. Drogon is a causal reasoning framework for future trajectory forecasting [12]. 149 The authors use LiDAR data, and design a conditional prediction model to forecast goal-oriented 150 trajectories. Finally, causal reasoning helps to reason about the behavior of vehicles as future locations 151 conditioned on the intention. Ramanishka et al. [14] present a dataset of 104 hours of real human 152 driving for learning driver behavior and causal reasoning. Another benchmark for analyzing causality 153 in traffic accident videos was presented by You et al. [7]. In this work they decompose an accident 154 into a pair of events and analyze the cause and effect. 155

Trajectory Prediction. In temporal reasoning research, trajectory prediction is a common task [9, 10], 156 partly due to the practical utility in numerous applications [33, 34, 35]. In our evaluation, we use 157 trajectory prediction as a key metric for performance and therefore it is helpful to briefly introduce 158 work on this topic. Most of these algorithms have been developed for scenarios with a single type 159 of agents. One such task is predicting pedestrians' future movements [36, 37, 38, 39], which is 160 important for autonomous vehicle and robotics design. Social behaviors have been widely exploited 161 in predicting pedestrians movements; while relevant for pedestrians, they are much less relevant 162 for vehicles. Thus the focus on vehicle trajectory prediction has been on modeling the motion of 163 individual agents (their past trajectory) and the surrounding environment [28, 40, 41]. A notable 164 exception is estimating lane changes on highways [42, 43]; previous efforts have tackled predicting 165 vehicle trajectories in urban scenarios [28, 40, 44]. 166

Compared to single-agent scenarios, multi-agent modeling and prediction is a challenging task for 167 control applications because agents interact with each other. Modeling dependencies between agents 168 is especially critical in scenarios such as modelling vehicles at intersections. Previous approaches 169 have focused on relatively sparse scenarios with only a few heterogeneous interactions. In such cases, 170 the interaction between agents can be modelled using social forces, velocity obstacles [45], or linear 171 trajectory avoidance [46]. When considering more complex interactions, learning-based approaches 172 have been applied between multiple pedestrians [34, 47, 48, 36, 49], vehicles [50, 51, 43, 28, 52, 53], 173 and athletes [54, 51]. These approaches attempt to generalize from previously observed interactions 174 to multi-agent behavior in new situations. To perform prediction without supervision, Ehrhardt et al. 175 [55] learn intuitive physics from visual observations and Kipf et al. [56] adopt contrastive learning to 176 177 perform self-supervision on structured world models.

# **178 3** Simulation Engine

Our goal is to create a high visual-fidelity simulation environment that can be used to systematically implement complex causal relationships in realistic scenarios. For this we focus on city driving scenarios and use a set of downtown city blocks and roads with multiple four way intersections and traffic lights. Fig. 3 shows examples of the visual appearance of the simulation environment, with first person views from close to street level.

Our simulation environment – dubbed *CausalCity* – is built upon AirSim [31], which acts as a 184 plugin for Unreal Engine and allows for obtaining training data using realistic graphics and physics 185 simulations. Our environment contains a city block with multiple four-way intersections and traffic 186 lights with cars navigating through it. The environment is controlled in two primary ways. First, there 187 is a JSON configuration file that defines a set of scenarios. Each scenario lists the vehicles that should 188 be present, their start locations, and the high-level actions that each vehicle should take. Secondly, a 189 python API allows scenarios to be triggered to start, parameters changed in real-time, and enables 190 convenient logging of data as the scenarios progress. While scenarios with the same high-level 191 definitions can be played out identically for reproducibility, it is easy to add variability by altering the 192 number, starting points, actions or velocities of the vehicles or changing other configurations such as 193 the timing of traffic lights, time of the day, and weather conditions. 194

195

197



Figure 3: CausalCity simulation environment.

Environment Features. Our city block includes the typical elements that might be observed in such an environment (e.g., buildings, trees, lamp posts, road works, etc.) (see Fig. 3); along with well defined lanes and traffic lights that can be explicitly targeted. The objects in the environment can be easily added/removed either prior to scene simulation, or dynamically through a Python API to create different configurations of static elements.

Vehicles. We introduce vehicles to the scene in a systematic manner through an AI traffic module, which handles the low-level navigation controls.<sup>1</sup> These vehicles traverse the scene along splines (routes) that are selected via the config-

urable scenario file. The environment contains predefined splines running through each lane and 210 intersection according to general traffic rules (based on right-hand drive). In the configuration file 211 each vehicle is given a starting (spawn) point, identified by a spline ID, and a list of high-level actions 212 to execute post-spawn. Merging actions (mergeL/mergeR) happen along lane splines and turning 213 actions (left/right) happen at intersections. The vehicle states such as positions/velocities can 214 be queried and obtained dynamically during scenario run-time for logging purposes. If desired, 215 information regarding any collisions observed during the scenario can also be recorded. 216

Traffic Lights. The environment also contains traffic lights at every intersection, and the vehicles 217 respond to these traffic signals. While keeping traffic flowing in a realistic manner, this also introduces 218 causal connections at the intersections. The sequence and timing of these lights can be controlled 219 during the scenario run-time. For simplicity, with our environment, we provide scripts to show how 220 to configure the timing of lights in a sequence and how to run these asynchronously. Vehicles can 221 be "forced" to continue driving at a red light, which can be used to simulate dangerous events and 222 increase the likelihood of collisions. 223

**Environment.** Environmental factors can be modified to create variations in the visual appearance of 224 the scene. These include introducing and controlling the strength of weather effects: rain, fog and 225 snow; changing the time of day, and varying the wetness of the road. This allows for new parameters 226 to be introduced as *confounders* (for example, road wetness during rain can lead to unpredictable 227 steering behavior), and increase the variability in the observed scene in both car behavior as well as 228 visual appearance of the scene, which makes perception tasks more challenging. 229

Views/Cameras. Our environment enables cameras to be placed at any location and moved during 230 a scenario to obtain image data. This means that first person, third person, and bird's eye view 231 perspectives are possible. For simplicity, in our first baselines presented in this paper we use a bird's 232 233 eye view to visualize the scenarios. It is also possible to equip each vehicle with a camera of its own to obtain first person perspectives from the vehicles – presenting opportunities for future work. 234

The environment also allows for recording various modalities of data from multiple cameras for 235 logging/visualization, such as RGB, depth and segmentation maps. In the current version, we record 236 RGB images of the bird's eye view, as well as ground truth instance segmentation maps (generated 237

207

<sup>&</sup>lt;sup>1</sup>https://www.unrealengine.com/marketplace/en-US/product/arch-vis-ai-traffic-system



Figure 4: **Datasets.** We created two datasets, a toy dataset and our CausalCity dataset, both with cars that are connected via A) causal "leader-follower" relationships where one car follows another (A i-iii) and non-causal or random relationships (A iv-vi). B) Shows a heatmap of the paths of the vehicles in the toy dataset and CausalCity dataset. The toy dataset uses a similar city grid structure but has simplified behaviors (constant velocities, straight trajectories, etc.). The CausalCity dataset contains more realistic behaviors that introduce challenging confounders. Notice how there are longer dwell times in different lanes due to traffic patterns etc.

by AirSim), where unique masks corresponding to each car in the scene are drawn for ease of use.
 For simplicity, our segmentation maps contain only the masks corresponding to the cars, and other
 scene objects are ignored.

240 seene objects are ignored.

As described in the following section, we use this framework to generate multiple scenarios, each scenario driven by the corresponding configuration setting, an example of which can be seen in Listing 1. Each scenario involves the vehicles moving through set routes, while the traffic lights and other scene variables can vary as set by the user. Data such as vehicle positions, images etc. are logged as the scenario evolves.

Logging. Our environment allows for rich logging of events. In the current version, for each frame we log the positions of the vehicles  $(x, y, z, \sigma_x, \sigma_y, \sigma_z)$  and the state of the traffic lights (current color and duration since last change). But the positions of other objects, weather events, time of days can all also be recorded as necessary.

# 250 **4 Dataset**

To illustrate the potential of our simulation environment, we generate data and evaluate state-of-the-251 art causal reasoning approaches on it. Previous work has focused on causal discovery in relatively 252 controlled settings (e.g., balls moving in a 2D plane [10] or 3D objects colliding [9]). As we branch 253 out to more realistic and practical scenarios (e.g., autonomous driving) we quickly encounter a 254 number of additional complexities. One way to think about these complexities is in the form of 255 confounders. For example, driving would be extremely difficult and unsafe without traffic signals. If 256 257 we were to attempt to determine a causal relationship between the trajectory of two vehicles (e.g., is 258 a car following an ambulance), the effect of traffic signals on their behavior could be considered a 259 confounder.

Rather than leap directly to a context with multiple confounders, we created two versions of our data:
1) a toy dataset with causal relationships but without agency and no confounders, 2) a complex dataset
created using our high fidelity simulation environment with agency (and therefore the confounders
associated with it). In both cases we generated scenarios with a fixed number of vehicles (4,8,12)
driving in the environment and a fixed number of causal relationships.

To introduce causal relationships between the vehicles we create "links" (or edges in the causal 265 graph). The edges are defined as a "leader-follower" relationship in which two cars are given the 266 same set of actions but one starts ahead of the other. In each scenario we create pairs (e.g., three pairs 267 of six cars = three edges) of "leader-followers"; in the causal reasoning language, the leader vehicles 268 are the *interventions* and the follower vehicles are the *outcomes*; the former causes the latter to move 269 in certain ways. The remaining six vehicles are not causally connected to any other vehicle. This is a 270 sparse graph (only three edges) but that is reasonable as causal relationships in the real life are often 271 sparse. See Fig. 4A for example trajectories. This is just one example of the possible application of 272

our simulation, see Section 6 and the project page for more examples. For simplicity in both datasets
we position the camera from a bird's eye view perspective above the environment looking down. We
record 150 RGB and segmentation frames for each scenario and log the position of each vehicle (6
degrees-of-freedom) at the same rate. Our dataset is available on our project page.

## 277 4.1 Toy Dataset

Our toy dataset uses a road layout that mimics the city block in the simulation environment. The cars 278 do not have agency and thus move at a fixed velocities (2 pixels per frame - similar to the average 279 speed in the CausalCity dataset), and there are no confounders such as traffic lights that influence the 280 velocity of the vehicles. The vehicles do not collide with one another so their paths are uninterrupted. 281 The leader vehicles start, and remain (since both have the same velocity), exactly 30 pixels in front 282 of the follower vehicles in each pair. See our supplementary material for more examples of the 283 trajectories of the cars in our toy dataset. We create a dataset with 4000/500/500 scenarios for the 284 train/validation/test splits. See a heatmap in Fig. 4B that shows the average dwell time across the 285 dataset - notice how uniform it is and contrast that with the heatmap for the CausalCity dataset. 286

#### 287 4.2 CausalCity Dataset

In this dataset the cars have agency, controlled by our simulation engine, and thus have more realistic behaviors than in the toy dataset. Each vehicle has a set of five actions to complete but can drive without manually specified routes. The cars can accelerate when there is space ahead of them and reduce speed when approaching a slower moving vehicle or traffic signal. Their internal controls cause a vehicle to brake when it is approaching another vehicle. However, in some cases, if traveling fast, there may be collisions which can impact the trajectory of the cars.

Traffic lights help control the flow of traffic and also impact the velocity of vehicles regardless of causal relationships (thus they are confounders). These factors mean that even if two cars are causally linked, they will not remain a fixed distance away from each other. The "follower" vehicles may catch up with the "leader" if the "leader" is stopped at a red traffic signal, or could fall further behind it if leader makes it past a light but it turns red as the follower approaches it. See Fig. 4A and our supplementary for examples of the trajectories of the cars in our CausalCity dataset.

We observe that introducing agency to a simulation that enables highlevel scenario definitions will inevitably introduce confounders to the environment make the causal relationships more difficult to recover using the baseline algorithms (as we will see in the results). Once again, we create dataset with 4000/500/500 scenarios for the train/validation/test splits.

# 304 5 Experiments

We evaluate three state-of-the-art causal inference algorithms – NRI [18], NS-DR [9] and V-CDN [10]
 – on both the toy and CausalCity datasets. Training was carried out on a single Nvidia P100 GPU.
 Each experiment typically required 10 hours of training and evaluation time.

#### 308 5.1 Models

**NRI** [18]. NRI is a variational autoencoder (VAE) optimized to discover a relational structure 309 while learning the dynamical model of the underlying system. The interaction structure is explicitly 310 modeled using a node-to-node message passing operation similar to Gilmer et al. [57]. Given 311 sequences of locations and velocities, NRI reconstructs the original trajectory based on the predicted 312 interaction graph. As such, the encoder learns to predict a probability distribution of edges between 313 nodes without knowing the underlying interaction graph apriori. We leverage a recurrent neural 314 network as a decoder to predict multiple time steps into the future and a fully-connected network as 315 the encoder. The directed causal graph is inferred through the encoder using 100 frames of "historical" 316 data, and then the decoder is used to predict future trajectory (up to 20 frames in our experiments) 317 conditioned on the causal graph. We reuse predicted trajectories as inputs of the decoder to estimate 318 the further steps. Further specifics of the implementation can be found in [18]. 319

**NS-DR** [9]. CLEVRER is a dataset designed for reasoning about causal relationships between 320 objects and events in a video. To accomplish causal reasoning in their work, Yi et al. [9] used a 321 propagation network (PropNet) [58] to learn object dynamics from videos and predict object motion 322 and collision events. We adapt this dynamics predictor model to our scenario. First, the input to the 323 PropNet is segmentation masks of all cars in all frames of a video. These segmentation masks could 324 be generated by popular semantic segmentation approaches. However, in our current setup, we use 325 326 the segmentation masks provided by the CausalCity simulation engine, which we assume to be the upper-bound in terms of segmentation performance. Next, PropNet builds a directed graph where 327 vertices and edges represent cars and their relationship, respectively. Each vertex encodes information 328 about states and attributes of a car, where the states denote mask patches taken from a history of 329 images, and the attributes denote the color of a car. We adapt our approach by assigning one unique 330 color to each car in the scene. Finally, since we do not have collision event like in the CLEVRER 331 dataset, our edge relationships do not contain collision state. However, including collision events 332 would be an interesting direction for future work in the autonomous driving context. Otherwise the 333 implementation matches that of Yi et al. [9] and their associated code base. 334



Figure 5: Trajectory Prediction Error. Mean
 square error in future trajectory prediction for eight
 car scenarios with two causal connection for a) NS DR, b) NRI, b) V-CDN algorithms. Shaded regions
 reflect standard error. Notice the different scales
 on y-axes for the two plots.

V-CDN [10]. Deep graph neural networks are often used to represent underlying properties (e.g., dynamics) of physical interactions. V-CDN infers the structural causal model (SCM) from visual inputs for future prediction without supervision from the ground-truth graph structure. Li et al. [10] show that this can help models perform counterfactual reasoning about unseen scenarios. V-CDN consists of three parts; (1) a perception module (2) an inference module and (3) a dynamics module. Specifically, the perception module takes a sequence of images and finds keypoints, and the inference module takes these keypoints to discover a causal graph that represents the causal relationships, i.e., a physical connection in their scenario. The dynamics module is a graph recurrent network that predicts the future location of keypoints conditioned on the estimated causal graph.

We use the official implementation of Li et al. [10] and adapt their inference and dynamics modules, while "bypassing" the perception module for fair comparisons with other models. We take location and acceleration of the cars as input to compare the results with other baselines. Following Li et al. [10], for a set of cars, we construct a directed causal graph and predict the future movement of cars by conditioning on the current state and the inferred causal graph. In 5, we report the mean squared error of car locations (normalized by image dimensions from 0 to 1).

# 360 6 Discussion

How do the different methods perform? As expected we observe a gradual increase in trajectory prediction error as we attempt to predict the vehicle locations further into the future. Fig. 5 shows the normalized mean squared pixel error for time steps 1 to 20 into the future for each of the baseline methods (this corresponds to approximately 1 to 20 seconds into the future, as we sample at 1 Hz on average). The baselines show differing performance and we observe that for NS-DR the trajectory prediction errors increase more rapidly. This is consistent with previous results that found trajectory prediction to be poorer without an explicit causal discovery step [10].

**Does agency impact trajectory prediction and causal discovery?** When we contrast the performance on the toy dataset with performance on the CausalCity dataset we observe that trajectory prediction errors are larger on the CausalCity dataset as we attempt more distant future predictions. The key difference between the two datasets is the lack of *confounders* in the toy dataset (see Fig. 4B where the heatmaps contrast the trajectories and dwell times of the vehicles). In the toy dataset the trajectories have fixed velocities and the vehicles travel on quite predictable - but less realistic - routes. In the CausalCity dataset the vehicles have agency that allows them to follow the rules of the road (e.g., drive in the correct lanes), to avoid collisions with other vehicles (i.e., brake if they are approaching another car), stop at traffic lights to reduce the risk of accidents etc.



Figure 6: Bar chart showing the F1 score for edge prediction. i) NRI and ii) V-CDN results for scenarios with 4, 8, 12 cars with a fixed proportion of edges (50% of cars having a causal connection). iii) NRI and iv) V-CDN results for scenarios with 20%, 50% and 80% of cars with a fixed number of cars (8). Error bars reflect standard error. When considering the resulting trajectories of the vehicles, this adds significant - but much more realistic - confounders. The effect is a more challenging task with some room for improvement. It is clear that current state-of-the-art benchmarks struggle with trajectory prediction to some degree. This may be partially explained by the fact that causal discovery tends to be more difficult too. Fig. 6 shows the F1 scores for edge type discovery on our toy dataset and CausalCity dataset. We performed experiments varying the number of cars (4, 8, 12) and the percentage of causal connections (20%, 50%, 80%). The NS-DR method does not perform causal discovery and therefore we show the results for NRI and V-CDN. We observe that overall these are lower for the CausalCity dataset and in particular for the case with 8 cars.

How does the number of cars and causal relationships impact results? Discovering causal relationships is important as it can help us learn the structure of the world and make better predictions about the future. Fig. 6 shows how causal discovery performance varies with the number of vehicles and proportion of cars that have a causal connection. We observe that causal discovery becomes more difficult as the number of cars increases (holding the proportion of cars that have a causal connection constant). Greater

406 proportions of causal connections (a less sparse causal graph) aid in causal discovery. One aspect of 407 our task that makes causal discovery particularly difficult is how sparse the causal graph is.

What other tasks can CausalCity support? We chose to demonstrate the capabilities of the 408 CausalCity simulation on the task of causal discovery with vehicles in leader-follower style context. 409 However, there are many other tasks in the domain of causal reasoning, discovery and counterfactual 410 reasoning that the simulation could be used for. For example, our simulation enables a "hero" vehicle 411 to be used to create targeted interventions in the scene and such a method could be used to test the 412 ability for algorithms to reason counterfactually (i.e., what would have happened if the hero vehicle 413 did not stop at the traffic signal?). As this is a simulation it is possible to generate a scene with 414 the same initial starting conditions but to strategically intervene with a specific action at a specific 415 416 moment. We have included an example of how to conduct this type of experiment in our repo.

## 417 **7 Broader Impacts**

Causal reasoning presents promising opportunities for machine learning. Specifically, causal discov-418 ery and counterfactual reasoning could help create models that are more explainable. Therefore, tools 419 that help advance this understanding will be valuable to the research community. However, we must 420 421 acknowledge some limitations of our system. Our simulation environment is designed around the task of driving; however, this does not mean that a system trained on these data will be appropriate for 422 real-world applications. The scenarios created in our simulation are complex and do have reasonable 423 visual fidelity, but the are still a long way from simulating realistic behavior of drivers. We are adding 424 pedestrians to the simulation engine but it does not currently feature animals (e.g., birds) which are 425 another commonly occurring element in everyday environments. This environment was designed for 426 experimentation, specifically in the domain of causal discovery; generalization to real-world tasks -427 especially safety critical ones like driving - would require greater testing on real-world data. 428

## 429 **References**

- 430 [1] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [2] Bernhard Schölkopf. Causality for machine learning. arXiv preprint arXiv:1911.10500, 2019.
- [3] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation.
   In *International conference on machine learning*, pages 1180–1189. PMLR, 2015.
- [4] Walter J Scheirer, Anderson de Rezende Rocha, Archana Sapkota, and Terrance E Boult.
   Toward open set recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(7):1757–1772, 2012.
- [5] Elizabeth S Spelke. Core knowledge. American psychologist, 55(11):1233, 2000.
- [6] Jonathan G Richens, Ciarán M Lee, and Saurabh Johri. Improving the accuracy of medical
   diagnosis with causal machine learning. *Nature communications*, 11(1):1–9, 2020.
- [7] Tackgeun You and Bohyung Han. Traffic accident benchmark for causality recognition. 2020.
- [8] Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling.
   Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, pages 6446–6456, 2017.
- [9] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B
   Tenenbaum. Clevrer: Collision events for video representation and reasoning. *arXiv preprint arXiv:1910.01442*, 2019.
- [10] Yunzhu Li, Antonio Torralba, Animashree Anandkumar, Dieter Fox, and Animesh Garg. Causal
   discovery in physical systems from videos. *arXiv preprint arXiv:2007.00631*, 2020.
- [11] Mengyue Yang, Furui Liu, Zhitang Chen, Xinwei Shen, Jianye Hao, and Jun Wang. Causal vae: Disentangled representation learning via neural structural causal models. *arXiv preprint arXiv:2004.08697*, 2020.
- [12] Chiho Choi, Abhishek Patil, and Srikanth Malla. Drogon: A causal reasoning framework for
   future trajectory forecast. *arXiv preprint arXiv:1908.00024*, 2019.
- [13] Xuanpeng Li, Qifan Xue, Jingwen Zhao, and Dong Wang. Causal reasoning in multi-object
   interaction on the traffic scene: Occlusion-aware prediction of visibility fluent. *IEEE Access*,
   8:80527–80535, 2020.
- [14] Vasili Ramanishka, Yi-Ting Chen, Teruhisa Misu, and Kate Saenko. Toward driving scene understanding: A dataset for learning driver behavior and causal reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7699–7707, 2018.
- [15] Hoon Kim, Kangwook Lee, Gyeongjo Hwang, and Changho Suh. Crash to not crash: Learn
   to identify dangerous vehicles using a simulator. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 978–985, 2019.
- [16] Mohammad Sadegh Aliakbarian, Fatemeh Sadat Saleh, Mathieu Salzmann, Basura Fernando,
   Lars Petersson, and Lars Andersson. Viena: A driving anticipation dataset. In *Asian Conference on Computer Vision*, pages 449–466. Springer, 2018.
- [17] Ossama Ahmed, Frederik Träuble, Anirudh Goyal, Alexander Neitz, Manuel Wüthrich, Yoshua
   Bengio, Bernhard Schölkopf, and Stefan Bauer. Causalworld: A robotic manipulation bench mark for causal structure and transfer learning. *arXiv preprint arXiv:2010.04296*, 2020.
- [18] T Kipf, E Fetaya, K-C Wang, M Welling, and R Zemel. Neural relational inference for
   interacting systems. 2018.
- [19] Yuval Atzmon, Felix Kreuk, Uri Shalit, and Gal Chechik. A causal view of compositional
   zero-shot recognition. *arXiv preprint arXiv:2006.14610*, 2020.

- 473 [20] Amy Fire and Song-Chun Zhu. Inferring hidden statuses and actions in video by causal
   474 reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* 475 Workshops, pages 48–56, 2017.
- 476 [21] Rohit Bhattacharya, Tushar Nagarajan, Daniel Malinsky, and Ilya Shpitser. Differentiable causal
   477 discovery under unmeasured confounding, 2020.
- 478 [22] Adam Lerer, Sam Gross, and Rob Fergus. Learning physical intuition of block towers by
   479 example. *arXiv preprint arXiv:1603.01312*, 2016.
- [23] Roozbeh Mottaghi, Mohammad Rastegari, Abhinav Gupta, and Ali Farhadi. "what happens
   if..." learning to predict the effect of forces in images. In *European conference on computer vision*, pages 269–285. Springer, 2016.
- [24] S. Song, S. P. Lichtenberg, and J. Xiao. Sun rgb-d: A rgb-d scene understanding benchmark
   suite. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages
   567–576, 2015.
- [25] Katerina Fragkiadaki, Pulkit Agrawal, Sergey Levine, and Jitendra Malik. Learning visual
   predictive models of physics for playing billiards. *arXiv preprint arXiv:1511.07404*, 2015.
- Iustin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick,
   and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary
   visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2901–2910, 2017.
- [27] Angelos Filos, Panagiotis Tigkas, Rowan McAllister, Nicholas Rhinehart, Sergey Levine, and
   Yarin Gal. Can autonomous vehicles identify, recover from, and adapt to distribution shifts? In
   *International Conference on Machine Learning*, pages 3145–3153. PMLR, 2020.
- [28] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmo han Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents.
   In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
   336–345, 2017.
- [29] Nicholas Rhinehart, Kris M Kitani, and Paul Vernaza. R2p2: A reparameterized pushfor ward policy for diverse, precise generative path forecasting. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 772–788, 2018.
- [30] Seong Hyeon Park, Gyubok Lee, Jimin Seo, Manoj Bhat, Minseok Kang, Jonathan Francis,
   Ashwin Jadhav, Paul Pu Liang, and Louis-Philippe Morency. Diverse and admissible trajectory
   forecasting through multimodal context understanding. In *European Conference on Computer Vision*, pages 282–298. Springer, 2020.
- [31] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and service robotics*, pages 621–635.
   Springer, 2018.
- [32] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun.
   CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017.
- [33] Alexandre Alahi, Vignesh Ramanathan, and Li Fei-Fei. Socially-aware large-scale crowd fore casting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
   pages 2203–2210, 2014.
- [34] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and
   Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–971, 2016.
- [35] Boris Ivanovic and Marco Pavone. The trajectron: Probabilistic multi-agent trajectory modeling
   with dynamic spatiotemporal graphs. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2375–2384, 2019.

- [36] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan:
   Socially acceptable trajectories with generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2255–2264, 2018.
- [37] Anirudh Vemula, Katharina Muelling, and Jean Oh. Social attention: Modeling attention in
   human crowds. In 2018 IEEE international Conference on Robotics and Automation (ICRA),
   pages 1–7. IEEE, 2018.
- [38] Matteo Lisotto, Pasquale Coscia, and Lamberto Ballan. Social and scene-aware trajectory
   prediction in crowded spaces. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019.
- [39] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Rezatofighi, and
   Silvio Savarese. Sophie: An attentive gan for predicting paths compliant to social and physical
   constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
   pages 1349–1358, 2019.
- [40] Shashank Srikanth, Junaid Ahmed Ansari, Sarthak Sharma, et al. Infer: Intermediate representations for future prediction. *arXiv preprint arXiv:1903.10641*, 2019.
- [41] Nachiket Deo and Mohan M Trivedi. Convolutional social pooling for vehicle trajectory
   prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* Workshops, pages 1468–1476, 2018.
- [42] Alex Kuefler, Jeremy Morton, Tim Wheeler, and Mykel Kochenderfer. Imitating driver behavior
   with generative adversarial networks. In 2017 IEEE Intelligent Vehicles Symposium (IV), pages
   204–211. IEEE, 2017.
- [43] N. Deo and M. M. Trivedi. Multi-modal trajectory prediction of surrounding vehicles with
   maneuver based lstms. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 1179–1184,
   2018.
- [44] A. Zyner, S. Worrall, and E. Nebot. Naturalistic driver intention and path prediction using
   recurrent neural networks. *IEEE Transactions on Intelligent Transportation Systems*, 21(4):1584–
   1594, 2020.
- [45] Jur van den Berg, Stephen J. Guy, Ming Lin, and Dinesh Manocha. Reciprocal n-body collision avoidance. In Cédric Pradalier, Roland Siegwart, and Gerhard Hirzinger, editors, *Robotics Research*, 2011.
- [46] S. Pellegrini, A. Ess, K. Schindler, and L. van Gool. You'll never walk alone: Modeling social
   behavior for multi-target tracking. In 2009 IEEE 12th International Conference on Computer
   Vision, pages 261–268, 2009.
- [47] F. Bartoli, G. Lisanti, L. Ballan, and A. Del Bimbo. Context-aware trajectory prediction. In
   2018 24th International Conference on Pattern Recognition (ICPR), pages 1941–1946, 2018.
- [48] Tharindu Fernando, Simon Denman, Sridha Sridharan, and Clinton Fookes. Soft+ hardwired
   attention: An 1stm framework for human trajectory prediction and abnormal event detection.
   *Neural networks*, 108:466–478, 2018.
- [49] Wei-Chiu Ma, De-An Huang, Namhoon Lee, and Kris M Kitani. Forecasting interactive dynamics of pedestrians with fictitious play. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 774–782, 2017.
- [50] Rohan Chandra, Uttaran Bhattacharya, Aniket Bera, and Dinesh Manocha. Traphic: Trajectory
   prediction in dense and heterogeneous traffic using weighted interactions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8483–8492, 2019.
- [51] Tianyang Zhao, Yifei Xu, Mathew Monfort, Wongun Choi, Chris Baker, Yibiao Zhao, Yizhou
   Wang, and Ying Nian Wu. Multi-agent tensor fusion for contextual trajectory prediction. In
   *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
   12126–12134, 2019.

- Seong Hyeon Park, ByeongDo Kim, Chang Mook Kang, Chung Choo Chung, and Jun Won Choi.
   Sequence-to-sequence prediction of vehicle trajectory via lstm encoder-decoder architecture. In
   2018 IEEE Intelligent Vehicles Symposium (IV), pages 1672–1678. IEEE, 2018.
- [53] Wenchao Ding, Jing Chen, and Shaojie Shen. Predicting vehicle behaviors over an extended
   horizon using behavior interaction network. In 2019 International Conference on Robotics and
   Automation (ICRA), pages 8634–8640. IEEE, 2019.
- [54] Chen Sun, Per Karlsson, Jiajun Wu, Joshua B Tenenbaum, and Kevin Murphy. Stochastic pre diction of multi-agent interactions from partial observations. *arXiv preprint arXiv:1902.09641*, 2019.
- [55] Sebastien Ehrhardt, Aron Monszpart, Niloy Mitra, and Andrea Vedaldi. Unsupervised intuitive
   physics from visual observations. In *Asian Conference on Computer Vision*, pages 700–716.
   Springer, 2018.
- [56] Thomas Kipf, Elise van der Pol, and Max Welling. Contrastive learning of structured world
   models. *arXiv preprint arXiv:1911.12247*, 2019.
- [57] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
   message passing for quantum chemistry. In *ICML*, 2017.
- [58] Yunzhu Li, Jiajun Wu, Jun-Yan Zhu, Joshua B Tenenbaum, Antonio Torralba, and Russ Tedrake.
   Propagation networks for model-based control under partial observation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 1205–1211. IEEE, 2019.

## 588 Checklist

589	1. For all authors
590 591	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
592 593	<ul><li>(b) Did you describe the limitations of your work? [Yes] We describe limitations in Section 7.</li></ul>
594 595	(c) Did you discuss any potential negative societal impacts of your work? [Yes] We have included a broader impacts statement that covers potential negative impact.
596 597	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
598	2. If you are including theoretical results
599	(a) Did you state the full set of assumptions of all theoretical results? [N/A]
600	(b) Did you include complete proofs of an incoretical results? [IV/A]
601 602 603	<ul> <li>a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] Our project</li> </ul>
604 605	page: https://causalcity.github.io/ contains links to the code, instructions, datasets and simulation used in this work.
606 607	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
608 609	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] See Fig. 5 and 6.
610 611	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
612	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
613	(a) If your work uses existing assets, did you cite the creators? [Yes]
614	(b) Did you mention the license of the assets? [N/A]
615 616	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes] Yes.

(d) Did you discuss whether and how consent was obtained from people whose data you're 617 using/curating? [N/A] 618 (e) Did you discuss whether the data you are using/curating contains personally identifiable 619 information or offensive content? [N/A] 620 5. If you used crowdsourcing or conducted research with human subjects... 621 (a) Did you include the full text of instructions given to participants and screenshots, if 622 applicable? [N/A] 623 (b) Did you describe any potential participant risks, with links to Institutional Review 624 Board (IRB) approvals, if applicable? [N/A] 625 (c) Did you include the estimated hourly wage paid to participants and the total amount 626 spent on participant compensation? [N/A] 627