

Learning A Risk-Aware Trajectory Planner From Demonstrations Using Logic Monitor

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Risk awareness is an important factor to consider when deploying poli-
2 cies on robots in the real-world. Defining the right set of risk metrics can be
3 difficult. In this work, we use a logic monitor that keeps track of the environ-
4 mental agents’ behaviors and provides a risk metric that the controlled agent can
5 incorporate during planning. We introduce LogicRiskNet, a learning structure that
6 can be constructed from temporal logic formulas describing rules governing a safe
7 agent’s behaviors. The network’s parameters can be learned from demonstration
8 data. By using temporal logic, the network provides an interpretable archite-
9 cture that can explain what risk metrics are important to the human. We integrate
10 LogicRiskNet in an inverse optimal control (IOC) framework and show that we
11 can learn to generate trajectory plans that accurately mimic the expert’s risk han-
12 dling behaviors solely from demonstration data. We evaluate our method on a
13 real-world driving dataset.

14 **Keywords:** Learning from demonstrations, Temporal Logic, Autonomous Driv-
15 ing

16 1 Introduction

17 Imagine a driver that is trying to [turn across traffic](#) at an unprotected intersection where there is
18 oncoming traffic. [When faced with the situation shown in Figure 1\(a\)](#), the driver needs to decide
19 when it is safe to make the turn. The driver makes this decision by evaluating how risky the nearest
20 oncoming vehicle is in terms of the intended turning maneuver. The risk may be dependent on a
21 number of factors, including how far the vehicle is from the intersection, how fast it is driving and
22 whether it has the intention to slow down. If we wish to incorporate such risk-aware decision making
23 into an autonomous driving system, we are faced with two problems: (1) how do we accurately
24 model the risk factors that humans take into account while driving, and (2) how do we generate
25 plans from these risk factors. With the development of data-driven autonomous driving, we have
26 access to a growing amount of large-scale human driving data. Large-scale data contain not only
27 normal driving behaviors, but also what drivers do in order to manage risk. Taking advantage of
28 these data, our goal in this work is to *develop a risk representation that can expressively describe*
29 *the desired/undesired behaviors of road agents and with parameters learnable from demonstration*
30 *data*. We aim to obtain risk metrics that reflect that of the human demonstrator, and generate risk-
31 aware trajectory plans using these metrics.

32 Recent work on risk and uncertainty aware policy learning aim to incorporate common [risk con-](#)
33 [structs](#) such as conditional value at risk, worst-case risk, etc., into the problem formulation either
34 as an auxiliary loss or constraints [1, 2]. Many then use reinforcement learning (RL) to obtain the
35 policy. The main limitation of applying RL to autonomous driving is its exploration requirement.
36 It can be difficult to explore safely especially when we wish to deal with risky scenarios. Work
37 on uncertainty-aware imitation learning [3] and risk-aware offline RL [4] address the exploration
38 problem, but are limited in the types of risks they can express.

39 It is important to be risk-averse in safety critical applications. It is often more important to be averse
40 to the right set of risks. In this paper, we want to answer the question “if we have observed our
41 nearby road agents for a while, can we develop a risk model that quantifies how risky they are

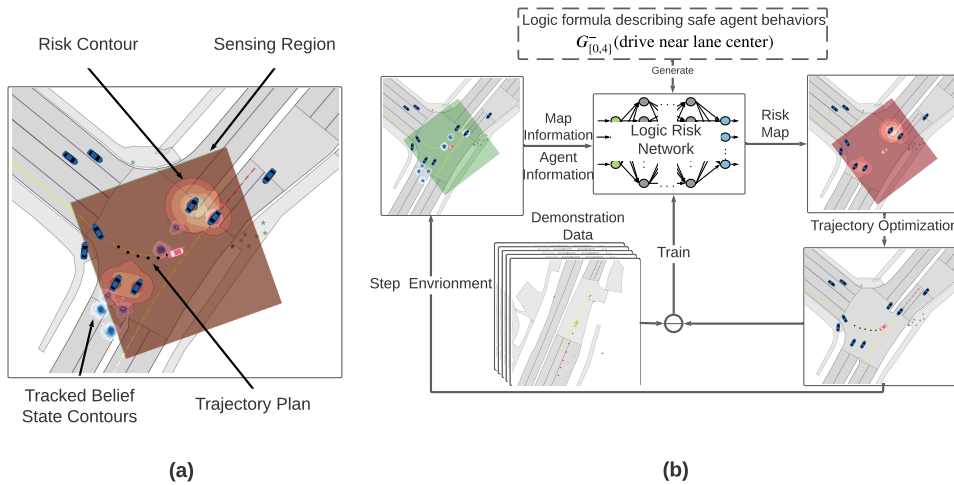


Figure 1: LogicRiskNet-IOC framework. (a) A pictorial depiction of our setup. The red car is the ego vehicle (controlled by our planner). The blue cars move according to the trajectories recorded in the dataset. Their trailing contours represent the tracking distributions. The red contours around the ado vehicles (neighboring vehicles not controlled by us) present their risks (the brighter the contour the higher the risk.) (b) An illustration of our inverse optimal control (IOC) procedure. The outer loop generates trajectory plans according to the risk map. The inner loop updates LogicRiskNet’s parameters using a kind of feature matching.

42 based on their past behaviors in the same way that humans would.” We take inspiration from real-
 43 time verification [5] and formal synthesis [6] and introduce a differentiable logic monitor which
 44 we refer to as *LogicRiskNet* as an expressive and learnable risk representation. To summarize our
 45 contributions, we

- 46 • introduce a differentiable logic-based risk metric and show that this metric is able to de-
 47 scribe complex, temporally-extended risk events as logic formulas and learn its parameters
 48 from human driving data;
- 49 • show that we can not only use LogicRiskNet to generate risk-aware trajectories, given its
 50 hierarchical monitoring properties, we can also use it to explain *why* an agent is risky;
- 51 • demonstrate in a real-world driving dataset that our method yields controllable risk-averse
 52 or, alternatively, risk-seeking behaviors.

53 *Note that, while we apply our approach to the domain of autonomous driving, we can equally-well*
 54 *extend our approach to other domains, such as robotics.*

55 2 Background: Past Time Signal Temporal Logic (PtSTL)

56 STL [7] offers a formalism for expressing and reasoning among a rich set of rules. The rules are
 57 defined over predicates (inequalities of real-valued functions). In standard STL, the formulas are
 58 forward-looking, meaning that the formulas are evaluated by looking at future trajectories. In this
 59 work, we need to evaluate formulas using past trajectories (trajectories of road agents we have
 60 tracked over the past), *using this past information to compute a risk of violating these rules in the*
 61 *future*. Therefore, our rules will be encoded as a set of *past time STL* (ptSTL) [8] formulas with the
 62 following syntax

$$\phi := p(x) < \epsilon \mid \neg\phi \mid \phi \wedge \phi \mid \phi \vee \phi \mid F_{[a,b]}^-\phi \mid G_{[a,b]}^-\phi, \quad (1)$$

63 where $a, b \in \mathbb{R}_{\geq 0}$, $a < b$, represent finite time bounds; ϕ is a ptSTL formula; $p(x) < \epsilon$ is a
 64 predicate; $x \in X \subseteq \mathbb{R}^n$ is a state; $p(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$ is the predicate function, $\epsilon \in \mathbb{R}$ is a constant.
 65 \neg (not), \wedge (and) and \vee (or) are Boolean operators. F^- (eventually) and G^- (always) are past time

66 temporal operators. Here $F_{[a,b]}^- \phi$ requires ϕ to be true at least once in the time interval $t \in (-b, -a]$
 67 whereas $G_{[a,b]}^- \phi$ requires ϕ to be true for all of this interval.

68 Let $x^{t_0:t_1} = \{x_{t_0}, \dots, x_{t_1}\}$ denote the state trajectory from time t_0 to t_1 . $x^{t_0:t_1} \models \phi$ denotes that
 69 trajectory $x^{t_0:t_1}$ satisfies ϕ (i.e. ϕ evaluates to true under $x^{t_0:t_1}$). The Boolean semantics of ptSTL
 70 are provided in Appendix A.

71 An example ptSTL formula is $F_{[\tau_0, \tau_1]}^- (|x| < 2) \wedge G_{[\tau_1, \tau_2]}^- (x > 3)$ which means that “ $|x_t| < 2$ is
 72 true for at least one $t \in (-\tau_1, -\tau_0]$ and $x_t > 3$ is true for all $t \in (-\tau_2, -\tau_1]$ ”.

73 Similar to STL, ptSTL is equipped with a *robustness degree* (or robustness for short) that quantifies
 74 the level of satisfaction of a trajectory with respect to a formula. A robustness function, denoted
 75 r , takes in a trajectory and a ptSTL formula and outputs a real-valued number. A robustness greater
 76 than zero signifies that the trajectory satisfies the given formula, i.e. $r(x^{t_0:t_1}, \phi) > 0 \Rightarrow x^{t_0:t_1} \models \phi$.
 77 Negative robustness implies violation of the formula. The formal definition of robustness is also
 78 given in the Appendix A.

79 Note that each predicate relationship $p(x) < \epsilon$ appearing in ϕ has parameters that can be fit to
 80 how humans behave. For example, the parameterized ptSTL formula $\phi = G_{[\tau_0, \tau_1]}^- (\alpha_1 x + \gamma_1 y <$
 81 $\epsilon_1) \wedge F_{[\tau_2, \tau_3]} (\alpha_2 x^2 > \epsilon_2)$ has parameters $\xi = \{(\alpha_1, \gamma_1), \epsilon_1, \alpha_2, \epsilon_2\}$. Hence, we parameterize each
 82 predicate appearing in a formula ϕ , denoting the parameterized formula ϕ^ξ (see, e.g., [9]).

83 3 Problem Formulation and Approach

84 Let $x = (p, \psi, v)$ denote the state of a vehicle that includes position (x-y coordinates), heading
 85 and velocity (which we later use in a unicycle model). We call the vehicle that is controlled by our
 86 planner the *ego agent*, and all other traffic participants not controlled by us *ado agents*. Denote Ω_T
 87 as the set of trajectories with horizon T (i.e. $\omega_T = x^{t:t+T} \in \Omega_T$). Suppose that the ego vehicle
 88 tracks its nearest N ado agents as well as itself via a tracking algorithm (such as a Kalman filter).
 89 This algorithm returns a set of belief states at time step t , $h_t = [x_e^t, \Sigma_e^t, x_{a_i}^t, \Sigma_{a_i}^t]$, $i \in [0, N)$,
 90 which are concatenated into a history of length H as $h_H^t = [h_{t-H}^t, \dots, h_t^t]$. Here, x_{e, a_i}^t and Σ_{e, a_i}^t
 91 are the mean values and covariances for the ego and i -th ado vehicle states. \mathcal{H}_H is the set of
 92 such histories (i.e. $h_H^t \in \mathcal{H}_H$). Define a dataset $\mathcal{D} = [\Omega_T, U_T, \mathcal{H}_H]$ where \mathcal{H}_H is obtained from
 93 sensors, Ω_T is the ego vehicle’s ground truth future trajectory, $u_T = [u^t, \dots, u^{t+T-1}] \in U_T$ is
 94 the corresponding sequence of controls (speed and steering). Let feature $f : \Omega_T \times U_T \rightarrow \mathbb{R}$ be
 95 a function that maps a trajectory of states and controls to a real value (e.g. average distance to the
 96 center lane). $f(\omega_T) = (f^1(\omega_T), \dots, f^n(\omega_T)) \in \mathbb{R}^n$ is a vector of such features. Our problem is
 97 formally defined as follows

98 **Problem 1.** Given (1) a parametric ptSTL formula ϕ^ξ that defines the behavior of a safe ado agent,
 99 (2) a feature f^{ϕ^ξ} representing ado agent risks, (3) a set of user defined features representing behav-
 100 iors of the ego agent \mathbf{f} , (4) a demonstrations dataset \mathcal{D} that contains a human driver’s decisions
 101 and environment information, find a trajectory distribution $\mathcal{P}^\theta : \Omega_T \rightarrow [0, 1]$ parameterized by θ (θ
 102 contains the ptSTL parameters ξ as well as feature combination weights) that best matches human
 103 behavior; i.e., minimizes the following objective functions

$$J_1 = \mathbb{E}_{\mathcal{P}^\theta} [\mathbf{f}(\tilde{\omega}_T)] - \mathbf{f}(\omega_T) \text{ and } J_2 = \mathbb{E}_{\mathcal{P}^\theta} [f^{\phi^\xi}(\tilde{\omega}_T)] - f^{\phi^\xi}(\omega_T) \quad (2)$$

104 where ω_T is a trajectory from the demonstration dataset. $\tilde{\omega}_T$ is a trajectory sampled from \mathcal{P}^θ .

105 Problem 1 defines a feature matching problem similar to many inverse reinforcement learning
 106 (IRL) [10] and IOC [11] formulations. We design a set of features \mathbf{f} to capture driving preferences
 107 in non-risky situations, but also include a risk-based feature f^{ϕ^ξ} capturing the risk management
 108 preferences seen in demonstrations. From the demonstration data we then learn the risk metric pa-
 109 rameters ξ and the weights that combine the feature values. This combination gives us a cost model
 110 that, when solved, yields driving behaviors similar to the demonstrator (in the context of our defined
 111 features and risk models). While existing IOC approaches are able to imitate standard driving be-
 112 haviors such that the generated trajectories from the cost model are exponentially more preferred by
 113 the agent (under maximum-entropy formulations, e.g. [11]). In our case, the additional risk feature

114 provides the capacity to generalize our model better under risky situations without hindering the
 115 ability to replicate the driving style of human demonstrators in normal scenarios.

116 4 Logic Monitor Guided Risk-Aware LfD

117 In this section, we will describe in detail our risk-aware LfD framework. We will first introduce the
 118 *LogicRiskNet* - a differentiable syntax tree generated from a parametric ptSTL formula. The network
 119 takes as input the tracking histories \mathbf{h}_T and outputs a risk assignment for each tracked vehicle based
 120 on their past behaviors. This **risk assignment** can be used to monitor which vehicles are risky in a
 121 given scene. We then show how LogicRiskNet can be incorporated into an IOC framework that is
 122 able to learn this syntax tree from demonstrations while simultaneously learn a cost function that
 123 specifies a risk-aware planning problem, from which a trajectory can be solved efficiently.

124 Our method is depicted in Figure 1(b) where there are 2 nested loops. The outer loop starts with the
 125 map and agent information which are passed through the LogicRiskNet to construct a risk map. This
 126 risk map constitutes the objective function for a trajectory planner. The generated trajectory serves
 127 to control the ego agent in a receding horizon fashion. In the inner loop, the same set optimized
 128 trajectories are also compared against human demonstrated trajectories in the same situations, and
 129 their differences serve as training signals to the LogicRiskNet.

130 **LogicRiskNet - a differentiable parametric ptSTL risk monitor.** We describe how we may con-
 131 struct a ptSTL risk monitor based on a learned stochastic risk measure that can be applied to a
 132 probabilistic description of human behaviors. We extend the definition of robustness degree (intro-
 133 duced in Section 2) to encompass belief states. In order to calculate the risk of a trajectory of belief
 134 states w.r.t a ptSTL formula ϕ , we need to modify the definition of the robustness. Recall that the
 135 basic element of ϕ is a predicate of the form $p(x) < c$, define $\alpha(x) = c - p(x)$ as the predicate (this
 136 is same as the robustness degree for the predicate). Since x is stochastic, rather than deterministic,
 137 we can evaluate ptSTL formulas instead using a *risk measure* $\rho : \mathcal{X} \rightarrow \mathbb{R}$, representing expectation,
 138 mean-variance, value-at-risk, etc.¹ In this work, we assume the expectation risk measure, but can
 139 equally well replace this with others. We can then apply this risk measure to the robustness definition
 140 yielding $\alpha^\rho(x) = \rho(\alpha(x))$, $x \in \mathcal{X}$ (This treatment is similar to [13]).

141 Given a parametric ptSTL formula ϕ^ξ , a risk measure $\rho(\cdot)$ and the tracking history \mathbf{h}_H , we can
 142 construct the *robustness risk* for ϕ^ξ as

$$R^\phi(\mathbf{h}_H | \xi) = r^\rho(\mathbf{h}_H, \phi^\xi) \quad (3)$$

143 where $r^\rho(\mathbf{h}_H, \phi^\xi)$ is the robustness applied to the risk-based predicate α^ρ . Once the risk predicate
 144 values are calculated at each time step in the trajectory, the robustness calculations for Boolean
 145 and temporal operators can be carried out. In order to tune the parameters ξ using demonstration
 146 data, we follow the approach in [14]. Specifically, if we replace the max and min functions in the
 147 robustness definition with softmax approximations, we can construct a computation graph from R^ϕ
 148 with the predicate parameters at the leaves and the robustness risk at the root. This makes it possible
 149 to learn ξ using backpropagation. The intuition behind Equation 3 is that given a ptSTL formula ϕ
 150 parameterized by ξ , agents' past observations \mathbf{h}_H , and a risk measure ρ , the robustness risk R^ϕ can
 151 be calculated to represent agent risk levels in terms of satisfaction of ϕ . Refer to Appendix F for an
 152 example construction of the LogicRiskNet.

153 **LfD using LogicRiskNet.** With LogicRiskNet R^ϕ , we are able to characterize risky ado vehicles as
 154 using the following risk feature:

$$f^\phi(\mathbf{h}_H, \omega_T | \xi) = \sum_{i=0}^{N-1} \sum_{t=0}^{T-1} -R^\phi(\mathbf{h}_H^i | \xi) \frac{1}{\|\mathbf{p}_e^t - \mathbf{p}_{a_i}^t\|^2} \quad (4)$$

155 where \mathbf{p}_e^t and $\mathbf{p}_{a_i}^t$, $t \in [0, T)$ are the ego and ado vehicles' positions at time t . \mathbf{h}_H^i is the tracking
 156 history of ado vehicle i along with the history of the ego agent. Essentially, (4) describes a collision
 157 feature scaled by the risk R^ϕ , where the inverse-squared distance may be viewed as a generalization

¹We further require that the risk measure is *coherent*; see [12].

Algorithm 1 Trajectory Optimization

```
1: Inputs: initial controls  $\tilde{\mathbf{u}}_T^{init}$ ; initial ego vehicle state  $\mathbf{x}_e^0$ ; tracking history  $\mathbf{h}_H$ ; lane information
    $\mathcal{T}$ ; control effort weight  $\beta^u$ ; lane tracking weight  $\beta^N$ ; ptSTL formula parameters  $\xi$ ;
2:  $\tilde{\mathbf{u}}_T \leftarrow \tilde{\mathbf{u}}_T^{init}$ 
3: for  $i=1 \dots N$  do
4:    $\tilde{\omega}_e \leftarrow \text{ConstructTrajectory}(\tilde{\mathbf{u}}_T, \mathbf{x}_e^0)$ 
5:    $\tilde{f}^\phi \leftarrow \tilde{f}^\phi(\mathbf{h}_H, \tilde{\omega}_e | \xi)$ ,  $\tilde{f}^u \leftarrow \tilde{f}^u(\tilde{\mathbf{u}}_T)$ ,  $\tilde{f}^N \leftarrow \tilde{f}^N(\tilde{\omega}_e, \mathcal{N}_T)$   $\triangleright$  construct features
6:    $\tilde{L} = \tilde{f}^\phi + \beta^u \cdot \tilde{f}^u + \beta^N \cdot \tilde{f}^N$   $\triangleright$  construct objective
7:    $\tilde{\mathbf{u}}_T \leftarrow \tilde{\mathbf{u}}_T - \tilde{\gamma} \nabla_{\mathbf{u}_T} \tilde{L}$   $\triangleright$  update control sequence
8: end for
9:  $\tilde{\omega}_e \leftarrow \text{ConstructTrajectory}(\tilde{\mathbf{u}}_T, \mathbf{x}_e^0)$ 
10: return  $\tilde{\omega}_e, \tilde{f}^\phi, \tilde{f}^u, \tilde{f}^N$ 
```

158 of potential fields [15]. A safe trajectory should yield a low f^ϕ value. A positive R^ϕ signifies a
159 rule abiding ado agent whereas a negative R^ϕ signifies a rule-violating agent. Intuitively, we use
160 the past observations to determine how risky an ado agent is, and by solving for trajectories that
161 minimize f^ϕ we can control the ego agent to stay away from risky ado agents. Our potential-field
162 formulation of risk acts to approximate the behavior of drivers in reacting to potential collisions but
163 do not over-react if collision is not imminent.

164 To imitate risk-free driving behavior, we additionally introduce a control effort feature and a nominal
165 trajectory feature for the ego vehicle as follows

$$f^u(\mathbf{u}_T) = \sum_{t=0}^{T-1} \|\mathbf{u}^t\|^2 \text{ where } \mathbf{u} = (v, \psi), \quad f^N(\omega_T, \mathcal{N}_T) = \sum_{t=0}^{T-1} \|\mathbf{p}_e^t - \mathbf{p}_{\mathcal{N}}^t\|^2 \quad (5)$$

166 where $\mathbf{p}_{\mathcal{N}}^t \in \mathcal{N}_T$ are points on the nominal trajectory. Nominal trajectories can be the human ego
167 agent’s trajectories (during training) or trajectories on the target lane center (during deployment).
168 These features encourage the ego agent to follow a trajectory (and speed) profile with smooth control
169 efforts. Given features f^ϕ, f^u, f^N , we solve for a trajectory plan ω_T by minimizing $f^\phi + \beta^u \cdot f^u +$
170 $\beta^N \cdot f^N$ using gradient descent. Algorithm 1 describes this trajectory optimization process.

171 In Algorithm 1, line 4 describes the process of using a unicycle model to construct a trajectory from a
172 control sequence. This procedure allows us to generate kinematically feasible trajectories compared
173 to direct trajectory optimization. In calculating the risk feature f^ϕ , we use a constant velocity and
174 heading model estimate of ado agents’ future trajectories $\mathbf{p}_{a_i}^t, t \in [0, T - 1)$. During training, the
175 nominal trajectory \mathcal{N}_T is the ego vehicle’s ground truth future trajectory from the dataset. During
176 deployment, this nominal trajectory is taken from points on the target center lane which can be
177 obtained from route planning.

178 Algorithm 2 describes how the parameters from the LogicRiskNet along with the objective function
179 weights are learned from demonstrations [using a maximum entropy model of uncertainty](#).

180 All entities with superscript d are taken from the demonstrations. Line 4 calculates the sum of
181 features obtained from trajectories from the human driver averaged over the dataset. Lines 6-10
182 calculates the same but with trajectories obtained from Algorithm 1. Line 12 updates the parameters
183 such that $A_{f_D^\phi}$ and $A_{\tilde{f}^\phi}$ can match.

184 5 Experimental Results

185 We train and evaluate our method on the NuScenes dataset. NuScenes [16] is a dataset for au-
186 tonomous driving based in Boston and Singapore. It contains 850 scenes each 20s long, containing
187 23 object classes and HD semantic maps with 11 annotated layers. We chose this particular dataset
188 for the rich semantics it provides which is well suited for rule definitions. We will use 650 scenes
189 for training and 200 scenes for validation. Details in experiment setup, implementation and hyper-
190 parameters used are provided in the Appendix C.

191 **Rules used.** We use the following two rules to describe behaviors of a safe ado agent. Our goal is
 192 to generate ego plans which avoid ado agents deemed risky in the sense that they violate these rules.

$$\phi_1 = G_{[0,H]}^- (\text{driveNearLane} \wedge \text{keepSafeCarDistFromEgo}) \quad (6)$$

$$\phi_2 = G_{[0,H]}^- (\text{egoInIntersection} \rightarrow (\text{farFromIntersection} \vee \text{driveSlowly})) \quad (7)$$

193 where H is the tracking horizon (we use $H = 4$). ϕ_1 describes that a safe ado vehicle should
 194 “*always* drive near the center lane *and* keep a safe distance from the ego vehicle.” ϕ_2 expresses that
 195 at intersections “*always* if the ego vehicle is in the intersection *implies* that a safe ado vehicle should
 196 either be far from the intersection *or* drive slowly.” The final rule $\phi = \phi_1 \wedge \phi_2$ takes the form of a
 197 conjunction of a set of sub-rules. Predicate and parameter definitions can be found in Appendix B.

198 **Methods of evaluation.** We evaluate our method and comparison cases in terms of optimality (time
 199 to reach goal position) and safety (minimum distance to nearby vehicles). Within the dataset, we
 200 will set the human ego vehicle’s start and end positions as the initial and goal positions. Optimality
 201 is measured as the time to travel from the initial position to within a distance to the goal position
 202 averaged over the validation set. Safety is measured as the minimum distance to nearby ado vehicles
 203 in a scene averaged over the validation set. During evaluation, we control the ego vehicle with our
 204 learned planner, the ado vehicles move according to the trajectories recorded in the dataset and time
 205 is synchronized. The advantage is that the ego vehicle can navigate in an environment with realistic
 206 human ado vehicles. The downside to this evaluation approach is that the ado vehicles will not be
 207 able to react to the ego’s actions. In future work, we will complement this evaluation method in
 208 simulated environments with reactive ado vehicles.

209 We use four methods for comparison. *LogicRiskNet-IOC* refers to the proposed method; *Human*
 210 refers the human driver in the dataset; *BC* refers to a behavior cloning agent; and *TrajOpt* refers to
 211 a trajectory optimization agent. Implementation details of *BC* and *TrajOpt* agents can be found in
 212 Appendix C.

213 **Results and discussions.** LogicRiskNet does not only output the final risk measure but can also
 214 output the risk of all intermediate sub-formulas without needing additional computation. This allows
 215 us to explain the reason of its decisions. Figure 2 shows an example of an unprotected right turn
 216 at an intersection and the vehicle that we are monitoring is highlighted by the yellow dash circle in
 217 Figures 2(a) and 2(c). Here we are showing two consecutive steps in the scenes. Looking at the
 218 top plot in Figure 2(b) we can see that this vehicle is labelled as risky (a negative risk value denotes
 219 violation of the safety rule ϕ .) To gain more insights on why this is the case, we can look at the
 220 risk values of sub-formulas. The middle and bottom plot shows that this vehicle is neither driving
 221 slowly nor far away from the intersection. Figure 2(d) shows that after one time-step the risk for the

Algorithm 2 LogicRiskNet Guided LfD

```

1: Inputs: parametric ptSTL formula  $\phi^\xi$ ; parameter learning rate  $\gamma^\phi$ ; demonstration dataset  $\mathcal{D}$ 
2:  $\xi \leftarrow \xi^{init}$ ,  $\beta^u \leftarrow \beta^{u,init}$ ,  $\beta^N \leftarrow \beta^{N,init}$ 
3: for  $i=1 \dots N$  do
4:    $A_{f_D^\phi} = \frac{1}{|\mathcal{D}|} \sum_{d=0}^{|\mathcal{D}|} \left( f^\phi \left( \mathbf{h}_H^d, \boldsymbol{\omega}_T^d \mid \xi \right) + f^u \left( \mathbf{u}_T^d \right) + f^N \left( \boldsymbol{\omega}_T^d, \mathcal{N}_T^d \right) \right)$ 
5:    $S_{\tilde{f}^\phi} = 0$ 
6:   for  $j=1 \dots |\mathcal{D}|$  do
7:      $\tilde{\boldsymbol{\omega}}_e, \tilde{f}^\phi, \tilde{f}^u, \tilde{f}^N = \text{TrajectoryOptimization} \left( \tilde{\mathbf{u}}_T^{init}, \mathbf{x}_e^{d,0}, \mathbf{h}_H^d, \mathcal{N}_T^d, \beta^u, \beta^N, \xi \right)$ 
8:      $S_{\tilde{f}^\phi} += \tilde{f}^\phi + \tilde{f}^u + \tilde{f}^N$ 
9:   end for
10:   $A_{\tilde{f}^\phi} = \frac{1}{|\mathcal{D}|} S_{\tilde{f}^\phi}$ 
11:   $L^\phi = (A_{f_D^\phi} - A_{\tilde{f}^\phi})^2$ 
12:   $(\xi, \beta^u, \beta^N) \leftarrow (\xi, \beta^u, \beta^N) - \gamma^\phi \nabla_{\xi, \beta^u, \beta^N} L^\phi$ 
13: end for
14: return  $\xi, \beta^u, \beta^N$ 

```

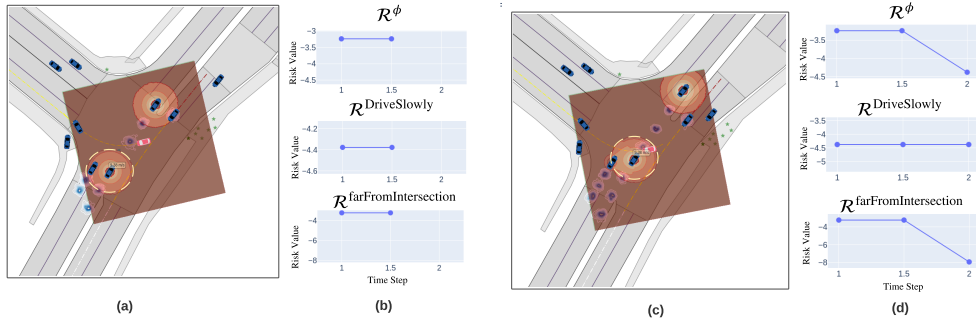


Figure 2: LogicRiskNet monitor results. (a) Monitoring the highlighted ado vehicle at time 1.5 sec. (b) Risk traces for ϕ and selected sub-formulas. (c) Monitoring at time 2 sec. (d) Corresponding risk traces.

Table 1: Performance comparison results

	Time to goal				Safety distance				% Reach goal
	min	mean	max	90th	min	mean	max	90th	
BC	13.21	19.98	20.0	20.0	0.34	4.22	10.86	8.66	37%
TrajOpt	16.92	18.50	20.0	19.76	3.89	6.37	12.43	8.34	87%
RiskLogicNet-IOC	15.31	17.35	20.0	19.0	2.72	3.98	9.01	6.12	93%
Human	19.0	19.75	20.0	20.0	3.03	4.62	9.45	6.84	100%

222 monitored vehicle has increased (risk value becomes more negative), and the cause is not because
 223 the vehicle has sped up but because it has driven closer to the center of the intersection. The structure
 224 of the entire LogicRiskNet and the monitoring traces for all sub-formulas are shown in Appendix D.

225 Table 1 shows comparison results in terms of the time it takes to reach the goal and the minimum
 226 distance to neighboring vehicles during this navigation. An aggressive agent will achieve a low time-
 227 to-goal (drive faster) as well as a small safety distance. The %Reach goal measures the percentage
 228 of scenes in the validation set the corresponding agent is able to reach the goal state. Failure is called
 229 if the agent does not reach the goal in the duration of the scene. All results are averaged over scenes
 230 in the validation set (each scene is ~ 20 seconds). From the table we can observe that *LogicRiskNet-*
 231 *IOC* achieves the most similar behavior to *Human* with the former slightly more aggressive. In
 232 comparison, *TrajOpt* is much more conservative. This is because it treats all agents as the same
 233 collision object that it needs to avoid whereas *LogicRiskNet-IOC* agent will place less emphasis on
 234 low risk agents (i.e. vehicles that are stopped). This gives the latter more free space to plan and also
 235 behaves more human-like (as humans place different amount of risk on vehicles based on their past
 236 behaviors). *BC* agent has a wide span of behaviors but mainly suffers from distribution shift [17]
 237 shown by its low success rate at reaching the goal.

238 Once the network parameters ξ and the objective parameters β are learned, we obtain a system that
 239 behaves similar to the human demonstrator. In practice, this may not be satisfactory and the user
 240 may wish to further fine tune the system. We conduct a study to see how varying the magnitude
 241 of the risk feature in Equation (4) changes the behavior of the ego agent. In Figure 3, the x-axis
 242 represents a risk coefficient multiplied to \tilde{f}^ϕ on line 6 in Algorithm 1. We run Algorithm 1 across
 243 all scenes in the validation set and record the time-to-goal and safety distance. Figure 3 shows
 244 that with the increasing relative importance placed on the risk term, the ego vehicle becomes more
 245 conservative (driving slower and keeping a further distance to ado vehicles). This finding allows
 246 the user to control the ego vehicle to strike a desired balance between risk-averse and risk-seeking
 247 behavior.

248 6 Related Work

249 *Risk and uncertainty aware policy learning* has been well-studied in the context of reinforcement
 250 learning. There have been several comprehensive surveys on this topic; [1] provide a general survey,
 251 while [2] provide greater focus on autonomous driving (AD). In general, common risk metrics such

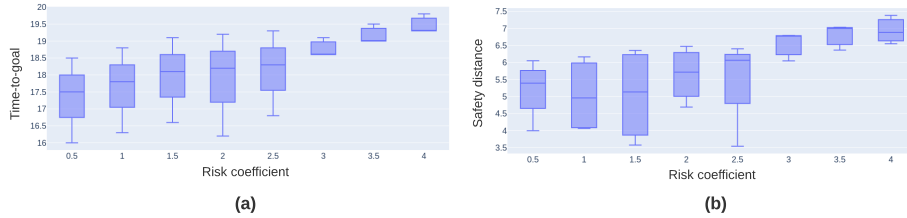


Figure 3: Controlling risk management behaviors by tuning risk feature coefficients. Both (a) Time-to-goal and (b) safety distance increases with the risk coefficient indicating that the ego vehicle becomes more conservative as we tune up this coefficient.

252 as conditional value at risk, mean-variance, worst-case analysis, etc. have been used as either an
 253 auxiliary loss in the objective function or as a set of constraints. In many robotic and driving tasks,
 254 however, RL is not readily applicable given its exploration requirement. Imitation learning and
 255 offline RL helps with addressing this problem. Specifically, in [3], kernelized movement primitives
 256 are used to estimate uncertainties in the demonstrations in finding optimal gains in a controller. The
 257 authors of [4] use offline distributional RL to learn a policy that is averse to conditional value at risk.
 258 Our work, in contrast, aims to learn from demonstrations to learn a human’s notion of risk. It does
 259 not require exploration beyond the demonstrations, yet it is able to generalize well to new scenarios,
 260 and has the capacity to learn rare situations such as rule violations and near collisions.

261 *Temporal logic (TL)-guided policy learning* is an area that we take have taken inspiration from.
 262 In this area, TL is often used to specify the ego agent’s desired high-level behavior and used to
 263 generate rewards. The authors of [18, 19, 20] provide surveys of recent work on the use of TL in
 264 RL. The exploration problem still exists in these methods. To address these challenges, the authors
 265 of [21, 22] learn finite state automata from demonstration and use them to guide planning with the
 266 value iteration network which avoids exploration. It can sometimes be tedious to manually design
 267 a TL formula that yields satisfying behaviors. Work has been done to make components of the
 268 formula learnable from data. In [23], the authors propose learning linear temporal logic (LTL)
 269 formulas from demonstrations. Given the close relationship between TL and automata[24], the
 270 authors of [25] propose a method that learns reward machines (an automata-like reward presentation)
 271 from demonstrations. A shortcoming of these methods is that the LTL that they use operates on
 272 propositions (binary variables with values true or false). Unlike our STL-based approach, these
 273 approaches require discrete state and action spaces and require the demonstrations themselves to
 274 have the same discrete representations.

275 *Safe and risk-aware formal synthesis* is an area that focuses on incorporating uncertainty and risks
 276 in planning with TL constraints. We have taken much inspiration from [13, 12] in their techniques of
 277 incorporating belief states into STL risk definitions. The resulting controls are obtained by solving
 278 a mixed-integer linear program with the STL risk as constraints. This can be difficult to scale with
 279 the complexity of the STL risk formulas. In [26, 27], the authors consider the uncertainties in the
 280 agent trajectories via probabilistic STL. Methods of computing specification-satisfying controllers
 281 under probabilistic STL constraints are also introduced. In these works, the probability of events
 282 happening at each timestep is required, which is difficult to obtain and does not reflect how humans
 283 reason about risks in the world. Our approach, in contrast, is data-driven and our logic risk metric is
 284 compatible with existing model-based/model-free policy learning paradigms.

285 7 Conclusion

286 In this work, we introduced the LogicRiskNet - an expressive and differentiable risk representation
 287 based on temporal logic descriptions of risky agent behaviors. We show that with a demonstration
 288 dataset we are able to learn risk parameters that results in trajectories similar to the demonstrator. We
 289 also show that given the structure of LogicRiskNet we can monitor and explain why an agent is risky
 290 without additional computation. Given our choice of the IOC framework and the interpretability of
 291 our method, users are able to tune the objective function post-training to obtain the desired level of
 292 risk-averse behaviors.

293 **References**

- 294 [1] J. Garcia and F. Fernández. A comprehensive survey on safe reinforcement learning. J. Mach.
295 Learn. Res., 16:1437–1480, 2015.
- 296 [2] Z. yu Zhu and H. Zhao. A survey of deep rl and il for autonomous driving policy learning.
297 ArXiv, abs/2101.01993, 2021.
- 298 [3] J. Silvério, Y. Huang, F. J. Abu-Dakka, L. Rozo, and D. Caldwell. Uncertainty-aware imitation
299 learning using kernelized movement primitives. 2019 IEEE/RSJ International Conference on
300 Intelligent Robots and Systems (IROS), pages 90–97, 2019.
- 301 [4] N. A. Urp’i, S. Curi, and A. Krause. Risk-averse offline reinforcement learning. ArXiv,
302 abs/2102.05371, 2021.
- 303 [5] J. V. Deshmukh, A. Donzé, S. Ghosh, X. Jin, G. Juniwal, and S. Seshia. Robust online moni-
304 toring of signal temporal logic. ArXiv, abs/1506.08234, 2015.
- 305 [6] C. Belta, B. Yordanov, and E. A. Gol. Formal methods for discrete-time dynamical systems.
306 2017.
- 307 [7] A. Donzé and O. Maler. Robust satisfaction of temporal logic over real-valued signals. In
308 FORMATS, 2010.
- 309 [8] E. A. Gol. Efficient online monitoring and formula synthesis with past stl. 2018 5th
310 International Conference on Control, Decision and Information Technologies (CoDIT), pages
311 916–921, 2018.
- 312 [9] E. Asarin, A. Donzé, O. Maler, and D. Nickovic. Parametric identification of temporal proper-
313 ties. In RV, 2011.
- 314 [10] G. Neu and C. Szepesvari. Apprenticeship learning using inverse reinforcement learning and
315 gradient methods. In UAI, 2007.
- 316 [11] M. Kuderer, S. Gulati, and W. Burgard. Learning driving styles for autonomous vehicles from
317 demonstration. 2015 IEEE International Conference on Robotics and Automation (ICRA),
318 pages 2641–2646, 2015.
- 319 [12] L. Lindemann, N. Matni, and G. J. Pappas. Stl robustness risk over discrete-time stochastic
320 processes. ArXiv, abs/2104.01503, 2021.
- 321 [13] S. Safaoui, L. Lindemann, D. Dimarogonas, I. Shames, and T. Summers. Control design for
322 risk-based signal temporal logic specifications. IEEE Control Systems Letters, 4:1000–1005,
323 2020.
- 324 [14] K. Leung, N. Aréchiga, and M. Pavone. Back-propagation through signal tempo-
325 ral logic specifications: Infusing logical structure into gradient-based methods. In
326 Workshop on Algorithmic Foundations of Robotics, 2020.
- 327 [15] A. Pierson, W. Schwarting, S. Karaman, and D. Rus. Navigating congested environments with
328 risk level sets. In 2018 IEEE International Conference on Robotics and Automation (ICRA),
329 pages 5712–5719, 2018. doi:10.1109/ICRA.2018.8460697.
- 330 [16] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan,
331 and O. Beijbom. nuscenes: A multimodal dataset for autonomous driving. arXiv preprint
332 arXiv:1903.11027, 2019.
- 333 [17] F. Codevilla, E. Santana, A. M. López, and A. Gaidon. Exploring the limitations of behav-
334 ior cloning for autonomous driving. 2019 IEEE/CVF International Conference on Computer
335 Vision (ICCV), pages 9328–9337, 2019.
- 336 [18] H.-C. Liao. A survey of reinforcement learning with temporal logic rewards.
- 337 [19] Q. Gao. Deep reinforcement learning with temporal logic specifications. 2018.

- 338 [20] C.-J. Heiker. Temporal logic specifications in reinforcement learning. 2021.
- 339 [21] B. Araki, K. Vodrahalli, T. Leech, C. Vasile, M. Donahue, and D. Rus. Learning to plan with
340 logical automata. In Robotics: Science and Systems, 2019.
- 341 [22] B. Araki, K. Vodrahalli, T. Leech, C. Vasile, M. Donahue, and D. Rus. Deep bayesian non-
342 parametric learning of rules and plans from demonstrations with a learned automaton prior. In
343 AAAI, 2020.
- 344 [23] D. Kasenberg and M. Scheutz. Interpretable apprenticeship learning with temporal logic spec-
345 ifications. 2017 IEEE 56th Annual Conference on Decision and Control (CDC), pages 4914–
346 4921, 2017.
- 347 [24] C. Baier and J. Katoen. Principles of model checking. 2008.
- 348 [25] R. T. Icarte, E. Waldie, T. Q. Klassen, R. Valenzano, M. P. Castro, and S. A. McIlraith. Learning
349 reward machines for partially observable reinforcement learning. In NeurIPS, 2019.
- 350 [26] D. Sadigh and A. Kapoor. Safe control under uncertainty with probabilistic signal temporal
351 logic. In Robotics: Science and Systems, 2016.
- 352 [27] C. Yoo and C. Belta. Control with probabilistic signal temporal logic. ArXiv, abs/1510.08474,
353 2015.