

## 1 Appendix

### 2 A Theoretical Background of Importance Sampling

3 Importance sampling is a variance-reduction technique for Monte Carlo sampling, often used to  
4 obtain more reliable estimations from fewer samples. Given a random variable  $x \sim p(x)$  and a  
5 function  $f(x)$ . Suppose we want to estimate the expected value of  $f(x)$  with sampling:

$$\mu = \mathbb{E}[f(x)] = \int f(x)p(x)dx \approx \frac{1}{n} \sum_i^n f(x_i), \quad (1)$$

6 where  $\{x_i\}_{i=1,\dots,n}$  are sampled from  $p(x)$ , and  $\frac{1}{n} \sum_i^n f(x_i)$  is the estimation using the  $n$  samples.

7 Given an *importance distribution*  $q(x)$ , importance sampling is performed as:

$$\mu = \int f(x)p(x) \frac{q(x)}{q(x)} dx \approx \frac{1}{n} \sum_i^n f(x_i) \frac{p(x_i)}{q(x_i)}, \quad (2)$$

8 where  $\{x_i\}_{i=1,\dots,n}$  are now sampled from the importance distribution  $q(x)$ , and  $\frac{p(x_i)}{q(x_i)}$ , often referred  
9 to as the *importance weights*, are used to correct the point-based estimation. If we define  $\hat{\mu}_q =$   
10  $\sum_i^n f(x_i) \frac{p(x_i)}{q(x_i)}$ ,  $\hat{\mu}_q$  is an unbiased estimator of  $\mu$ .

11 It has also been shown that the variance of  $\hat{\mu}_q$  is,

$$Var[\hat{\mu}_q] = \frac{1}{n} \int \frac{(f(x)p(x) - \mu q(x))^2}{q(x)} dx. \quad (3)$$

12 Therefore, the optimal importance distribution  $q^*$  that offers the lowest variance is:

$$q^*(x) = \frac{|f(x)|p(x)}{\mu} \quad (4)$$

13 The above equation suggests increasing the sampling probability of  $x$  with high absolute function  
14 values  $|f(x)|$ . This theoretical indication supports our core idea of learning to emphasize *critical*  
15 events, which lead to hazards and thus large negative values.

### 16 B Neural Network Details

17 Figure 1 demonstrates network architectures of the attention generator and the critic. The attention  
18 generator network first concatenates numbers in the current belief  $b$  and the observation  $z$  into a  
19 single vector and feeds it to a feature extractor. The feature extractor consists of 10 fully-connected  
20 layers with ReLU activation. The extracted features are input to a Gated Recurrent Unit (GRU) cell  
21 [1], which keeps track of history inputs in its latent memory. Based on the latent memory, the network  
22 uses two fully-connected layers and a soft-max layer to output the importance distribution  $q$ . The  
23 critic network has similar architecture. The belief  $b$ , observation  $z$ , and the generated importance  
24 distribution  $q$  are first concatenated into a single vector. Then, the vector is fed to a feature extractor  
25 containing 6 fully-connected layers with ReLU activation, then input to a GRU cell for tracking  
26 memory. Based on the latent memory, the critic network uses another 3 fully-connected layers to  
27 finally output a single decimal number, representing the estimated planner value  $v$ . See detailed  
28 representations of  $b$ ,  $z$ , and  $q$  in the following section.

### 29 C The POMDP Model for Urban Driving

30 Our POMDP model for driving in an ill-regulated dense urban traffic is defined as follows:

- 31 • **State Modeling:** A world state  $s$  encodes:

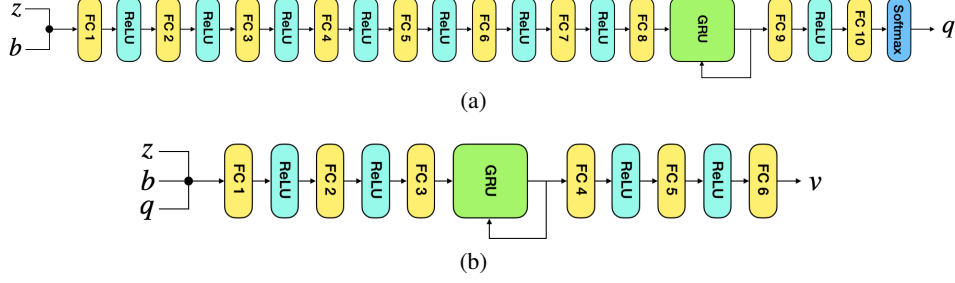


Figure 1: Network architectures of (a) the attention generator and (b) the critic function.

- 32 – the state of the ego-vehicle,  $s_c = (p_c, v_c, \alpha_c, P_c)$ , where  $p_c$ ,  $v_c$  and  $\alpha_c$  denote its
- 33 position, velocity, and heading direction, and  $P_c$  denotes its reference path.
- 34 – observable states of 20 nearest exo-agents,  $s_{exo} = \{p_i, v_i, \alpha_i\}_{i=1, \dots, 20}$ , where  $p_i$ ,  $v_i$ ,
- 35  $\alpha_i$  are the position, velocity, and heading direction of the  $i$ th exo-agent.
- 36 – hidden states of 20 nearest exo-agents,  $\theta_{exo} = \{\theta_i\}_{i=1, \dots, 20}$ , where  $\theta_i$  is the intention of
- 37 the  $i$ th exo-agent. Suppose an exo-agent has  $M$  potential paths to undertake according
- 38 to the lane network, the value of its intention  $\theta$  will be taken from  $\{0, \dots, M - 1\}$ .

39 A belief  $b$  is thus a discrete probability distribution defined over the hidden states or intentions  
 40 of exo-agents, assuming probabilistic independence between different participants. It is  
 41 represented using  $\sum_{i=1}^{20} M_i$  probability values, where  $M_i$  is the number of intentions for  
 42 the  $i$ th exo-agent. An importance distribution  $q$  is specified in the same way.

- 43 • **Action Modeling:** An action  $a$  of the ego-vehicle is its acceleration discretized to three  
 44 values, *ACC*, *CUR*, and *DEC*, meaning to accelerate, keep the current speed, and decelerate.  
 45 The acceleration and deceleration are  $3m/s^2$  and  $-3m/s^2$ , respectively.
- 46 • **Observation Modeling:** An observation  $z$  from the environment includes all observable  
 47 parts of the state  $s$  and excludes the hidden intentions. Namely,  $z = (s_c, s_{exo})$ . Due to  
 48 perceptual uncertainty, these observations often come with noise. However, in this work,  
 49 we particularly focused on the uncertainty in human behaviors and ignored perceptual  
 50 uncertainty, because the latter often has a secondary influence on decision-making.
- 51 • **Transition Modeling:** Our transition model assumes the ego-vehicle follows its reference  
 52 path using a pure-pursuit steering controller and the input acceleration. Exo-agents are not  
 53 controlled by the algorithm. We assume they take one of their hypothetical intended paths,  
 54 using the GAMMA motion model [2] to avoid collision with surrounding participants. At  
 55 each time step, all agents are simulated forward by a fixed duration of  $1/3s$ . Afterward,  
 56 small Gaussian noises are added to all transitions to model uncertain human control.
- 57 • **Reward Modeling:** The reward function takes into account driving safety, efficiency, and  
 58 smoothness. When the ego-vehicle collides with any exo-agent, it imposes a severe penalty  
 59 of  $-20 \times (v^2 + 0.5)$  depending on the driving speed  $v$ . To encourage driving smoother, we  
 60 also add a small penalty of  $-0.1$  for the actions *ACC* and *DEC* to penalize excessive speed  
 61 changes. Finally, to encourage the vehicle to drive at a speed closer to its maximum speed  
 62  $v_{max}$ , we give it a penalty of  $\frac{v - v_{max}}{v_{max}}$  at every time step.

## 63 References

- 64 [1] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural  
 65 networks on sequence modeling. In *NIPS 2014 Workshop on Deep Learning, December 2014*,  
 66 2014.
- 67 [2] Y. Luo, P. Cai, Y. Lee, and D. Hsu. Gamma: A general agent motion model for autonomous  
 68 driving. *IEEE Robotics and Automation Letters*, 2022.