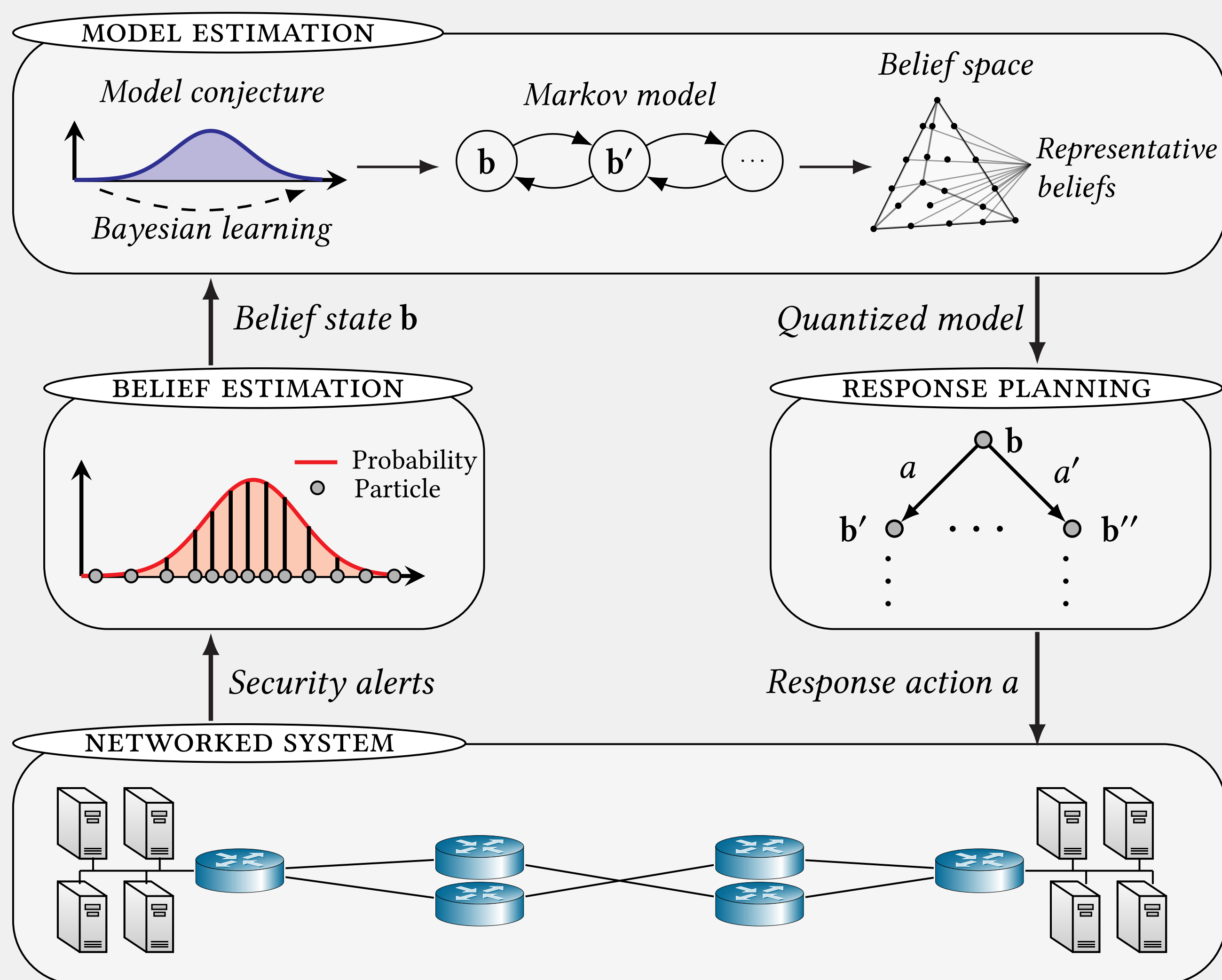


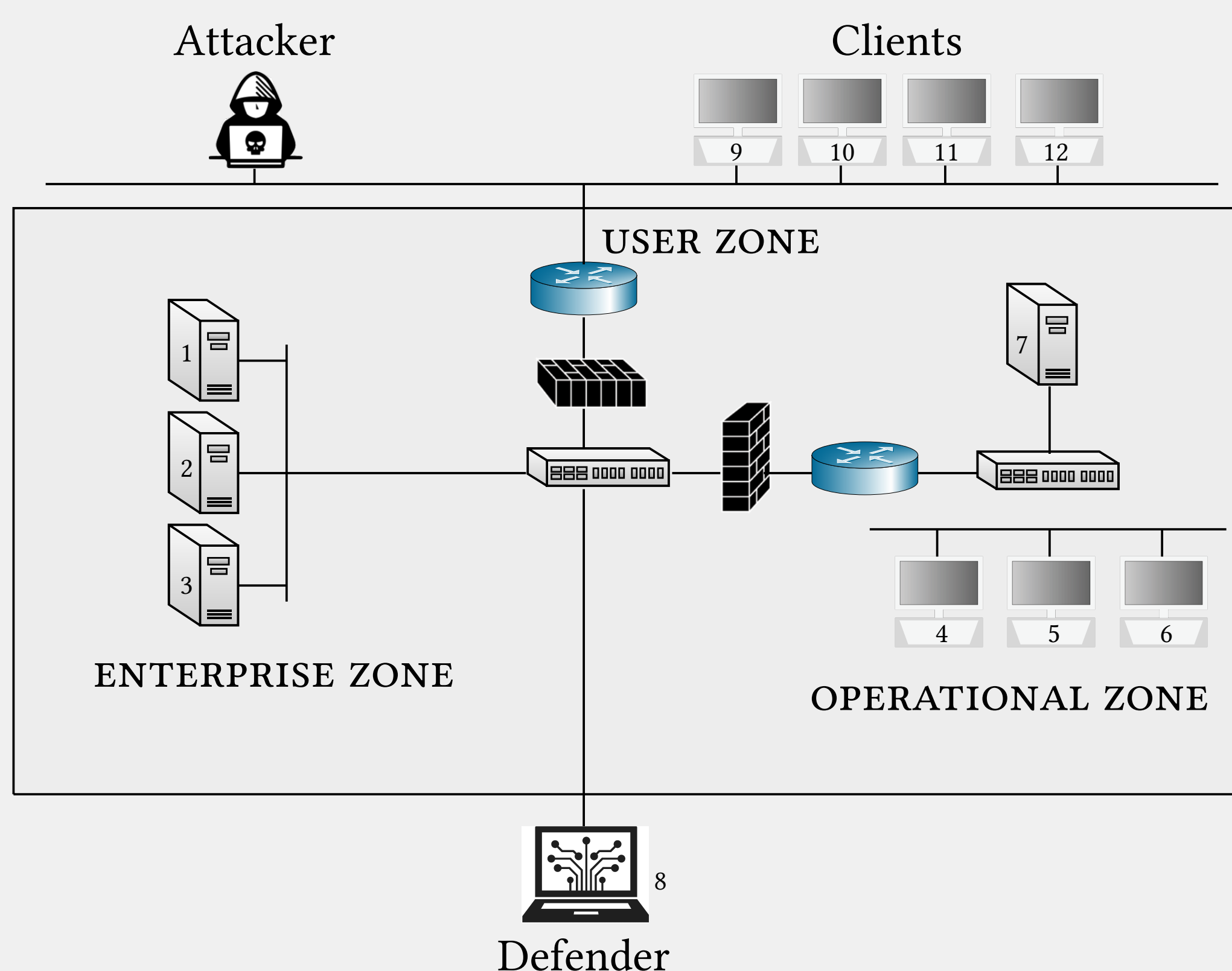
## Contributions

1. We present MOBAL, an online method for incident response planning.
2. We establish bounds on misspecification and quantization errors.
3. We show that MOBAL obtains state-of-the-art performance on CAGE-2.

## Misspecified Online Bayesian Learning (MOBAL)



## POMDP Model of Incident Response



We formulate incident response planning as a POMDP and seek to find a near-optimal response strategy  $\pi$  that maps belief states to response actions.

## Theoretical Results (Informal)

**Proposition 1 (Consistent conjectures).** The model conjecture learned by MOBAL is asymptotically consistent with respect to the information feedback.

**Proposition 2 (Misspecification error bound).** The difference between the conjectured optimal cost function  $\tilde{J}^*$  and the true optimal cost function  $J^*$  is bounded as

$$\|\tilde{J}^* - J^*\|_\infty \leq \frac{\gamma \alpha c_{\max}}{(1 - \gamma)^2},$$

where  $\gamma$  is the discount factor,  $\alpha$  quantifies the difference between the transition probabilities in the conjectured model and the true model, and  $c_{\max}$  is the maximum stage cost.

**Proposition 3 (Approximation error bound).** The difference between the cost function approximation  $\tilde{J}$  obtained through quantization and the conjectured optimal cost function  $\tilde{J}^*$  is bounded as

$$|\tilde{J}(b) - \tilde{J}^*(b)| \leq \frac{\epsilon}{1 - \gamma},$$

where  $\gamma$  is the discount factor and  $\epsilon$  is the maximum variation of  $\tilde{J}^*$  within each belief space partition.

**Proposition 4 (Asymptotic (conjectured) optimality).** The cost function approximation  $\tilde{J}$  obtained through quantization converges to the conjectured optimal cost function  $\tilde{J}^*$  as  $r \rightarrow \infty$ , where  $r$  is the quantization resolution.

**Theorem 1 (Sub-optimality bound of MOBAL).** The sub-optimality of the cost function approximation  $\tilde{J}$  obtained through MOBAL is bounded as

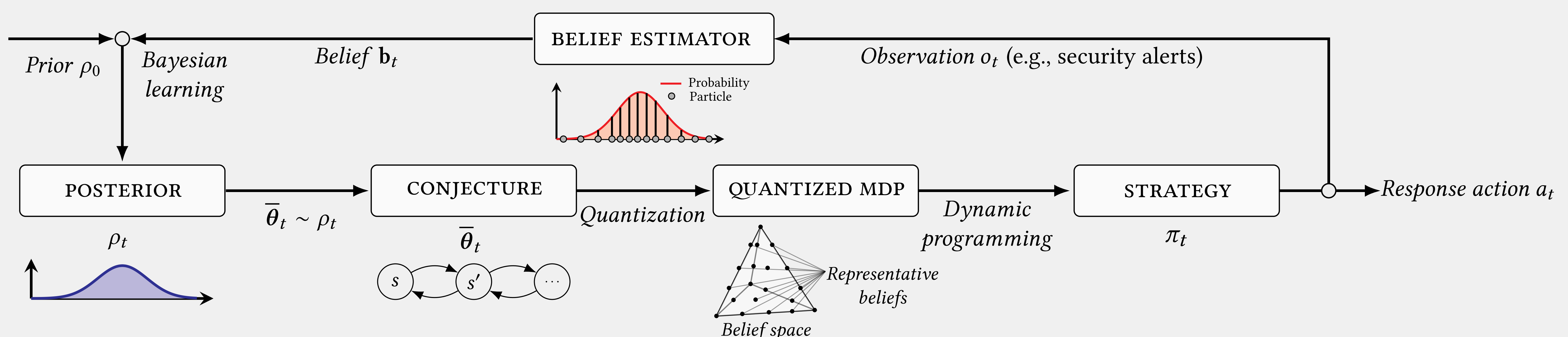
$$\|\tilde{J} - J^*\|_\infty \leq \frac{\epsilon}{1 - \gamma} + \frac{\gamma \alpha c_{\max}}{(1 - \gamma)^2}.$$

## Evaluation Results on the CAGE-2 Benchmark

Method	Offline/Online compute time (min)	Cost ( $\downarrow$ better)
<b>No misspecification</b>		
MOBAL	0/8.50	15.19 $\pm$ 0.82
CARDIFF	300/0.01	<b>13.69</b> $\pm$ 0.53
PPO	1000/0.01	119.02 $\pm$ 58.11
C-POMCP	0/0.50	<b>13.32</b> $\pm$ 0.18
POMCP	0/0.50	29.51 $\pm$ 2.00
<b>Misspecification</b>		
MOBAL	0/8.50	<b>35.91</b> $\pm$ 9.01
CARDIFF	300/0.01	94.28 $\pm$ 33.27
PPO	1000/0.01	124.38 $\pm$ 55.49
C-POMCP	0/0.50	92.71 $\pm$ 27.67
POMCP	0/0.50	91.51 $\pm$ 28.23

(C-POMCP and CARDIFF are state-of-the-art methods.)

## Online Response Planning, Belief Estimation, and Bayesian Learning



**Figure:** MOBAL: an iterative method for online learning of incident response strategies under model misspecification. The figure illustrates a time step during which (i) the posterior distribution over possible system models is updated via Bayesian learning based on feedback from the system; (ii) a conjectured model is sampled from the posterior and quantized into a computationally tractable MDP; and (iii) a response strategy is computed using dynamic programming. **Preprint:** <https://arxiv.org/pdf/2508.14385>.