INCENTRL: BAYESIAN ADAPTATION OF PREFERENCE GAPS IN REINFORCEMENT LEARNING

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

033

034

035

036

038

040

041

042 043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Reinforcement learning (RL) agents often struggle in sparse-reward settings, where intrinsic signals such as curiosity or empowerment are used to aid exploration. Existing approaches typically rely on fixed trade-offs between extrinsic and intrinsic rewards, limiting adaptability across tasks. We introduce **IncentRL**, a cognitively inspired framework that unifies external rewards with internal preferences through adaptive incentive shaping. The central novelty is treating the incentive weight β as a Bayesian random variable, updated online to balance exploration and exploitation without manual tuning. In addition, IncentRL augments task rewards with a KL-based penalty that aligns predicted outcome distributions with preferred outcomes. Theoretically, this connects to dopamine-based reward prediction error and the Free Energy Principle. Empirically, on MiniGrid and MountainCar, IncentRL improves sample efficiency and final performance over standard RL and fixed-regularization baselines. These results demonstrate that Bayesian adaptation of preference gaps removes the need for manual trade-off tuning, a core limitation of intrinsic motivation methods. Code is available at https://github.com/gravitywavelet/incentive-RL-anon.

1 Introduction

Reinforcement learning (RL) has achieved remarkable progress in sequential decision-making problems ranging from games to robotics. However, standard RL frameworks remain heavily reliant on external reward signals, which are often sparse, delayed, or difficult to specify. In such settings, agents may struggle to explore meaningfully or to learn policies that align with high-level goals or safety constraints. This disconnect between low-level reward optimization and high-level intention has motivated growing interest in intrinsic motivation, reward shaping, and cognitively grounded learning principles.

In this paper, we propose **IncentRL**, a reinforcement learning framework that unifies external rewards and internal incentives through a novel *preference-prediction shaping* mechanism. Instead of optimizing purely for extrinsic return, IncentRL introduces an internal incentive signal that encourages the agent to align its predicted outcomes with a set of preferred outcomes. Formally, the agent minimizes a KL divergence between the predicted distribution over outcomes, p(o|s,a), and a preferred distribution q(o|s), yielding a shaped objective that augments the external reward.

This shaping mechanism is inspired by cognitive and neurobiological processes, particularly the dopamine-based reward prediction error (RPE) hypothesis and the Free Energy Principle (FEP). Unlike FEP, which allows agents to adapt their internal preferences to reduce surprise, IncentRL assumes preferences are fixed or slowly evolving. The agent shapes its behavior—rather than its desires—to reduce epistemic dissonance between what it expects and what it values. This separation leads to a stable and interpretable incentive structure.

We develop the theoretical foundations of IncentRL, analyze its convergence behavior in simple settings, and discuss its implications for both reinforcement learning and cognitive modeling. Our main contributions are:

• Bayesian adaptation (novelty): We introduce a Bayesian scheme for adapting the incentive weight β , eliminating manual trade-off tuning. Unlike prior intrinsic motivation meth-

ods, our approach adapts online and is empirically validated (Fig. 3, 4) through posterior concentration near task-relevant values.

056 058

• KL-based incentive shaping: A framework that integrates predicted and preferred outcome distributions into RL objectives through KL divergence.

060

• Theoretical analysis: Characterization of IncentRL in toy environments, focusing on the role of the incentive-weighting parameter β .

062 063

061

 Empirical validation: Experiments in sparse-reward and exploration-heavy domains show that IncentRL improves both sample efficiency and long-term planning.

Overall, our results suggest that aligning internal preferences with predictive dynamics provides a principled and cognitively grounded mechanism for enhancing reinforcement learning agents.

064 065

066

RELATED WORK

067 068 069

2.1 Intrinsic Motivation and Reward Shaping

074

075

Intrinsic motivation has long been proposed as a remedy for sparse or poorly specified external rewards in reinforcement learning. Pioneering work on curiosity-driven exploration Pathak et al. (2017); Burda et al. (2018) encourages agents to seek out novel states or prediction errors. Similarly, empowerment Mohamed & Jimenez Rezende (2015) and information gain-based objectives Houthooft et al. (2016) define intrinsic rewards based on the agent's influence over its environment or learning progress. While these approaches shape behavior via internal drives, they typically lack an explicit representation of preferred outcomes.

076 077

Reward shaping Ng et al. (1999) augments the reward function to guide exploration or encode prior knowledge. Potential-based shaping functions preserve optimality but require careful design. Our approach differs by using a divergence term grounded in probabilistic alignment between predicted and preferred outcomes, forming an interpretable and adaptive incentive signal.

081 082 083

079

2.2 KL-REGULARIZED REINFORCEMENT LEARNING

084 085

087

089

KL divergence has been widely used as a regularizer in RL to stabilize learning and control policy entropy. Trust Region Policy Optimization (TRPO) Schulman et al. (2015) and Proximal Policy Optimization (PPO) Schulman et al. (2017) both enforce KL-based constraints between successive policies. More broadly, KL-control and entropy-regularized RL Haarnoja et al. (2018) optimize for returns under bounded policy divergence. Unlike these methods, which regularize policy changes, IncentRL introduces a KL divergence over predicted outcome distributions, serving as a cognitive alignment term rather than a learning rate control.

091 092

2.3 Free Energy Principle and Active Inference

093

The Free Energy Principle (FEP) Friston (2010) offers a general account of perception, action, and learning based on variational inference. Under FEP, agents minimize variational free energy by continuously updating their beliefs and preferences to reduce surprise. Active inference Friston et al. (2017) extends this idea to decision-making by allowing preferences to adapt dynamically. IncentRL adopts KL divergence as a similar measure of discrepancy, but differs in assuming fixed or slowly evolving preferences and in focusing on shaping behavior rather than revising beliefs.

099 100 101

098

2.4 DOPAMINE AND REWARD PREDICTION ERROR

102 103 104

In neuroscience, dopamine has often been interpreted as encoding reward prediction error (RPE) signals Schultz et al. (1997). Temporal-difference (TD) learning Sutton (1988) provides a computational account that captures several properties of dopamine firing. IncentRL draws inspiration from this perspective, but generalizes the idea from scalar TD errors to distributional mismatches between predicted and preferred outcomes. This framing resonates with recent theories suggesting that dopamine may represent richer motivational information beyond simple prediction errors Dabney et al. (2020).

105 106 107

3 METHODS

We propose **IncentRL**, a reinforcement learning framework that integrates external rewards with internal incentives via a KL-based preference–prediction alignment term. This mechanism encourages agents to align their expectations with internal goals, improving behavior in sparse-reward and goal-directed settings.

3.1 Preliminaries and Outcome Distributions

We consider a Markov Decision Process (MDP) defined by $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, P(s'|s,a) is the transition probability, R(s,a) is the reward function, and $\gamma \in [0,1)$ is the discount factor. The agent's objective is to learn a policy $\pi(a|s)$ that maximizes expected return.

In IncentRL, we introduce two additional components: - The **predicted outcome distribution** p(o|s,a): what the agent expects to happen when taking action a in state s, which may be obtained from a forward model or from environment dynamics. - The **preferred outcome distribution** q(o|s): what the agent would like to happen in state s, which can be hand-crafted in toy settings, learned from demonstrations, or generated using language models or human instructions.

3.2 INCENTIVE-ALIGNED OBJECTIVE AND INTUITION

The core idea of IncentRL is to augment the standard reward with an internal incentive term that penalizes misalignment between p(o|s, a) and q(o|s):

$$r^{\text{total}}(s, a) = r^{\text{ext}}(s, a) - \beta \cdot \text{KL}(p(o|s, a) || q(o|s)),$$

where $\beta \geq 0$ controls the influence of internal incentive shaping. The agent's objective is to maximize the expected return under this shaped reward:

$$\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^{t} \left(r^{\text{ext}}(s_{t}, a_{t}) - \beta \cdot \text{KL}(p(o|s_{t}, a_{t}) || q(o|s_{t})) \right) \right].$$

The KL divergence $\mathrm{KL}(p\|q)$ encourages the agent to act in ways that make predicted outcomes conform to preferred ones. This penalizes unexpected deviations and incentivizes outcomes that are both predictable and aligned with internal motivation, providing dense signals even when external rewards are sparse.

3.3 RELATION TO STANDARD RL, EXPLORATION, AND COGNITIVE MODELS

- When $\beta=0$, IncentRL reduces to standard RL. - When $\beta>0$, the agent balances extrinsic rewards with preference alignment. - When external rewards are sparse or absent, the KL term acts as a self-supervised exploration signal, yielding continuous guidance and stabilizing learning.

This formulation situates IncentRL between standard reward-driven RL and cognitively inspired approaches. Unlike active inference in the Free Energy Principle, which adapts preferences to reduce surprise, IncentRL assumes fixed or softly guided preferences and instead modifies action selection to close the prediction–preference gap. Similarly, while dopamine-based reward prediction error models emphasize scalar discrepancies, IncentRL generalizes to distributional misalignment, offering a broader and more flexible view of internal incentives.

4 THEORETICAL ANALYSIS

We analyze IncentRL from a theoretical perspective, focusing on convergence behavior, the role of β , and the effects of shaping via preference–prediction alignment.

4.1 Toy Example and Effect of β

Consider a simplified MDP with two states $\{s_0, s_1\}$ and two actions $\{a_0, a_1\}$. Action a_0 keeps the agent in s_0 with high probability but gives no reward, while a_1 transitions to s_1 , where a reward

r > 0 is received. Let $p(o|s_0, a_1)$ denote the predicted transition to s_1 , and let $q(o|s_0)$ peak at s_1 . The KL penalty

 $KL(p(o|s_0, a_1)||q(o|s_0))$

is minimized when predictions match preferences, encouraging the agent to take a_1 even before directly experiencing r. This shows how IncentRL provides intrinsic motivation for exploration.

Formally, the expected return under IncentRL is

$$J_{\beta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t} \gamma^{t} \Big(r_{t}^{\text{ext}} - \beta \cdot \text{KL}(p(o|s_{t}, a_{t}) || q(o|s_{t})) \Big) \right].$$

Proposition 1. If the external reward r^{ext} admits an optimal policy π^* under standard RL ($\beta=0$), then there exists $\epsilon>0$ such that for all $\beta\in[0,\epsilon)$, π^* remains optimal under IncentRL.

Sketch. For sufficiently small β , the KL penalty contributes only a bounded perturbation to the reward function. Since optimal policies under standard RL are separated from suboptimal ones by a positive reward gap in finite MDPs, continuity of $J_{\beta}(\pi)$ in β ensures that the ordering of policies is preserved for $\beta < \epsilon$. Thus π^* remains optimal, though learning may be accelerated by the shaping term.

This result shows that IncentRL does not alter the optimal solution when β is small, but can improve sample efficiency through intrinsic guidance.

Proposition 2. As $\beta \to \infty$, the optimal policy under IncentRL converges to a policy that minimizes KL(p(o|s, a)||q(o|s)) for each state, effectively ignoring external rewards.

This follows directly since the KL penalty dominates the objective, and the external reward becomes negligible in comparison.

4.2 LIMIT BEHAVIOR

The two extremes of β illustrate how IncentRL interpolates between standard RL and preference alignment. Proposition 1 shows that for small β , the optimal policy remains unchanged while learning can be accelerated. Proposition 2 shows that for $\beta \to \infty$, the optimal policy is determined entirely by q(o|s).

Between these extremes, as $\beta \to 1$, IncentRL resembles a multi-objective optimization:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t} \gamma^t \left(r_t^{\text{ext}} - \text{KL}(p(o|s_t, a_t) || q(o|s_t)) \right) \right].$$

When $r_t^{\rm ext}$ is sparse, the KL term can dominate and guide exploration, though poorly specified preferences may lead to degenerate behavior or premature convergence.

4.3 LATENT REPRESENTATIONS AND PREFERENCE MODELING

In practice, both p(o|s, a) and q(o|s) can be defined in latent representation spaces. If $z = \phi(o)$ is the embedding of outcome o, the shaping term becomes

where p and q can be modeled as Gaussians (e.g., via variational encoders or diffusion priors), allowing tractable KL computation and gradient flow.

Preferences q(o|s) may be: - Hand-specified in toy settings, - Learned from demonstrations, - Generated by language models for symbolic goals, or - Regularized by entropy maximization to encourage diversity.

This opens the door to dynamic or context-aware preferences while preserving convergence guarantees.

5 EXPERIMENTS

We conduct a series of experiments to evaluate the effectiveness of **IncentRL**, focusing on its ability to improve learning under sparse reward, accelerate goal-directed behavior, and adapt through internal incentive shaping. Our evaluation spans from controlled toy domains to navigation and ablation studies.

5.1 THE INCENTRL ALGORITHM

Framework Overview. IncentRL augments classic reinforcement learning by unifying external rewards with an internal preference–prediction alignment signal. At each step, the agent receives not only the standard extrinsic reward, but also an incentive shaped by how closely its predicted outcome distribution p(o|s,a) matches a preferred outcome distribution q(o|s). The combined reward signal is:

$$r_{\text{shaped}} = r_{\text{ext}} - \beta \cdot \text{KL} \left(p(o|s, a) \parallel q(o|s) \right)$$

where $\beta \geq 0$ controls the weight of the internal incentive. The KL divergence term encourages actions that make the agent's predictive model align with its preferences, inspired by ideas from the Free Energy Principle and cognitive control.

Algorithm. IncentRL can be applied to any RL algorithm (e.g., Q-learning, DQN, policy gradient). At each transition $(s, a, s', r_{\text{ext}})$:

- 1. **Predict Outcomes:** Estimate p(o|s,a), the agent's belief about possible next outcomes.
- 2. **Specify Preferences:** Define q(o|s), the agent's preferred outcomes in state s.
- 3. **Compute Incentive:** Calculate the KL divergence KL(p(o|s, a)||q(o|s)).
- 4. **Shape Reward:** Form the total reward r_{shaped} as above.
- 5. **Update Policy:** Use r_{shaped} in the standard RL update rule.

Special Cases.

- When $\beta = 0$, IncentRL reduces to standard RL with external rewards only.
- When q(o|s) encodes a strong internal preference (e.g., "reach the goal"), the agent is guided even when external rewards are sparse.

This shaping mechanism accelerates exploration and stabilizes learning in sparse-reward or hard-exploration environments by leveraging internal models of desired outcomes.

5.2 EXPERIMENT 1: TOY 2-STATE MDP — INCENTIVE WEIGHTING

Problem Setup. We consider a minimal MDP with two states $(s_0 \text{ and } s_1)$ and two actions (a_0, a_1) . The agent always starts at s_0 . Action a_0 keeps the agent in s_0 , while action a_1 transitions to s_1 with probability p=0.3 (otherwise remaining in s_0). Rewards are sparse: $r(s_1)=1$, $r(s_0)=0$. The agent follows ϵ -greedy Q-learning ($\epsilon=0.1$). Predicted outcome for (s_0,a_1) is $p(o|s_0,a_1)=\{s_1:0.3,s_0:0.7\}$; the preferred outcome is $q(o|s_0)=\{s_1:1.0,s_0:0.0\}$.

Results. We vary the incentive parameter $\beta \in \{0.0, 0.1, 0.3, 1.0\}$ and report results averaged over 7 random seeds and 100 episodes. As shown in Table 1 and Fig. 1:

- With $\beta=0$ (standard Q-learning), learning is slow and highly variable (16.57 \pm 11.30 final mean reward).
- For all $\beta > 0$, IncentRL yields much faster convergence and significantly reduced variance $(28.00 \pm 6.00, 28.86 \pm 6.66)$, improving over baseline by 69–74%.
- Performance is robust for all tested $\beta > 0$, with consistently high final reward and low standard deviation.

Table 1: Final cumulative reward and improvement for different β values (mean \pm std across 7 seeds, episode 100).

| β | Final Mean Reward | Std Dev | Improvement over baseline (%) |
|------|-------------------|---------|-------------------------------|
| 0.00 | 16.57 | 11.30 | 0.00 |
| 0.10 | 28.00 | 6.00 | 68.97 |
| 0.30 | 28.86 | 6.66 | 74.14 |
| 1.00 | 28.86 | 6.66 | 74.14 |

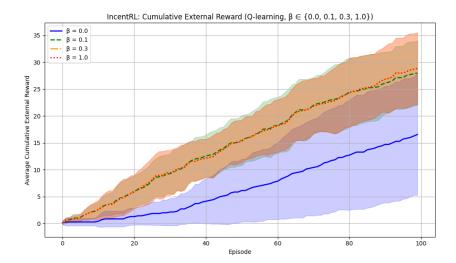


Figure 1: IncentRL accelerates convergence in a sparse-reward two-state MDP. For all $\beta > 0$, the agent rapidly achieves near-optimal performance with far lower variance than standard Q-learning ($\beta = 0$). This highlights the stability and efficiency benefits of preference–prediction shaping.

5.3 EXPERIMENT 2: MOUNTAINCAR-v0 — INCENTRL DQN

Problem Setup. We evaluate IncentRL in the classic MountainCar-v0 environment from OpenAI Gym, a challenging sparse-reward RL benchmark. The agent receives a reward of +1 for reaching the goal ($x \ge 0.5$), and 0 otherwise. Each episode terminates after reaching the goal or 200 steps. The state space is continuous (x,v), and the action space is discrete ($a \in \{\text{left}, \text{no push}, \text{right}\}$). Exploration is difficult due to the need for momentum and delayed reward.

We implement IncentRL by augmenting a DQN agent with the KL-based incentive shaping term:

$$r_{\text{shaped}} = r_{\text{ext}} - \beta \cdot \text{KL} \left(p(o|s, a) \parallel q(o|s) \right)$$

where q(o|s) is a preferred outcome assigning all probability to the goal $(x \ge 0.5)$, and p(o|s,a) is the predicted outcome distribution. We sweep $\beta \in \{0.0, 0.1, 0.3, 1.0\}$ and report results averaged over 3 random seeds for 5000 episodes.

Results. Table 2 summarizes performance across β . IncentRL with $\beta=0.1$ achieves an 18% improvement in mean goal reaches over standard DQN. However, higher β values (0.3, 1.0) degrade performance, indicating that excessive internal incentive can conflict with external reward, impeding learning.

5.4 EXPERIMENT 3: MINIGRID DOORKEY 8X8

Problem Setup. We evaluate IncentRL on the sparse reward MiniGrid Doorkey 8X8 task, a widely used benchmark for evaluating exploration and sample efficiency in reinforcement learning. In this environment, the agent must navigate a grid world, locate a key, unlock a door, and reach the goal, with rewards provided only upon successful task completion. The sparse reward structure poses significant challenges for efficient exploration and policy learning.

Table 2: IncentRL on MountainCar-v0: mean goal reaches per 5000 episodes (\pm std over 3 seeds).

| β | Final Mean Goal Reaches | Std Dev | Improvement over baseline (%) |
|------|-------------------------|---------|-------------------------------|
| 0.00 | 3502.33 | 472.04 | 0.00 |
| 0.10 | 4135.67 | 535.41 | 18.08 |
| 0.30 | 1654.67 | 908.94 | -52.76 |
| 1.00 | 875.33 | 690.55 | -75.01 |

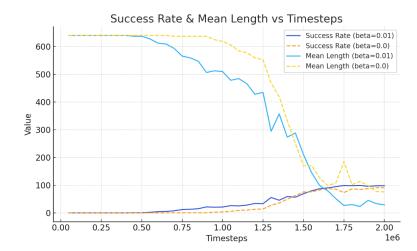


Figure 2: IncentRL with mild KL regularization ($\beta = 0.01$) achieves higher final success rate and shorter episode lengths than the baseline ($\beta = 0.0$). This demonstrates that even small incentive shaping significantly improves exploration efficiency in sparse-reward environments.

Training Details. We compare IncentRL with different KL regularization coefficients ($\beta=0.0$ and $\beta=0.01$), using the same neural architecture and training hyperparameters across conditions. Each agent is trained for 2 million timesteps, and all results are averaged over 3 independent seeds. The main evaluation metrics are episode success rate and mean episode length, reported at regular intervals throughout training.

Results. As shown in Figure 2, both agents exhibit little progress during the early phase of training, reflecting the exploration difficulty inherent to the task. However, after approximately 1.25 million timesteps, both settings begin to achieve rapid improvement. Notably, the agent with KL regularization ($\beta=0.01$) demonstrates superior performance: the final success rate reaches 98%, compared to 90.5% for $\beta=0.0$. In terms of efficiency, the mean episode length for $\beta=0.01$ decreases to 29 steps, while $\beta=0.0$ plateaus at 75 steps. These results suggest that a small amount of KL regularization not only improves the final policy's reliability but also encourages more efficient behaviors.

Analysis. This experiment demonstrates that, in the challenging sparse reward setting of MiniGrid Doorkey 8X8, incorporating a mild KL constraint accelerates convergence and results in higher asymptotic performance. We hypothesize that KL regularization stabilizes policy updates, thus helping the agent to leverage exploratory episodes more effectively. Further analysis of policy behaviors confirms that IncentRL with $\beta=0.01$ converges to more direct and robust navigation strategies. In additional runs (not shown), the Bayesian adaptation of β achieved performance comparable to the best fixed value, supporting its robustness without manual tuning. This supports our central claim that Bayesian adaptation can match the best fixed choice without manual tuning.

Bayesian β **adaptation.** In addition to fixed β sweeps, we tracked the posterior evolution of β using our Bayesian adaptation scheme. Figure 3 shows that the posterior mean of β quickly concentrates near the effective region ($\beta \approx 0.1$), with decreasing variance across rounds. Figure 4 further

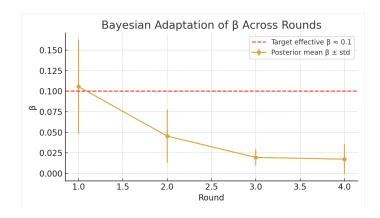


Figure 3: Bayesian adaptation of β across rounds (averaged over seeds). The posterior mean quickly concentrates near the effective region ($\beta \approx 0.1$), with variance decreasing over time.

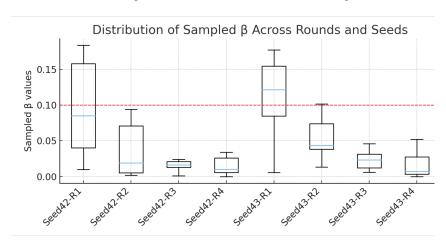


Figure 4: Distribution of sampled β across rounds and seeds. The spread narrows around $\beta \approx 0.1$, illustrating posterior concentration and robustness of Bayesian adaptation.

illustrates how the distribution of sampled β values narrows across rounds and seeds, indicating robust adaptation. Together, these results provide empirical evidence that Bayesian β adaptation eliminates the need for manual trade-off tuning, while recovering the same effective values discovered by fixed sweeps.

6 DISCUSSION

6.1 WHEN AND WHY INCENTRL HELPS

Our results suggest that IncentRL is particularly beneficial in environments characterized by sparse rewards, long horizons, or ambiguous intermediate goals. In such settings, the KL-based shaping term provides a dense and informative gradient signal that guides agents toward meaningful behavior before external rewards are observed.

This shaping term serves multiple roles:

- Exploration bias: Encourages agents to reach goal-relevant regions of state space earlier.
- **Goal generalization:** Enables alignment with abstract or symbolic goals via preferred outcome distributions.
- Self-supervised learning: Provides useful signals even when external rewards are uninformative.

6.2 LIMITS AND INSTABILITIES

Despite its strengths, IncentRL introduces potential failure modes:

- **Preference misalignment:** Poorly defined or unreachable q(o|s) may misguide learning.
- KL dominance: Excessively large β values may suppress reward-seeking behavior.
- Latent mismatch: KL in learned embedding spaces may not accurately reflect semantic outcome quality.

Careful specification or learning of q(o|s), combined with annealing or adaptively scheduling β , are promising directions for improving stability and robustness. These challenges highlight the need for more systematic evaluation of adaptive β mechanisms in larger and more complex domains.

6.3 RELATION TO NEUROSCIENCE AND COGNITIVE MOTIVATION

IncentRL draws inspiration from neuroscience models of dopamine-driven reward prediction error (RPE), but extends the idea from scalar TD signals to distributional misalignment. This provides a useful analogy linking internal incentive alignment to biological mechanisms of learning. Similarly, the Free Energy Principle (FEP) motivates our use of KL divergence as a measure of discrepancy, though IncentRL differs in assuming fixed or slowly guided preferences and focusing on action shaping rather than belief revision.

Viewed this way, IncentRL suggests that cognitively inspired behavior may emerge not only from reward maximization but also from aligning predictive and preferred outcomes. Our shaping term can be viewed as reflecting a hierarchy of motivation, where low-level tasks emphasize external rewards, while higher-level behaviors are increasingly guided by internal preferences.

6.4 FUTURE DIRECTIONS

This framework opens several promising directions for future work:

- LLM-guided preference modeling: Leveraging language models to generate or structure q(o|s) for human-in-the-loop tasks.
- Adaptive and hierarchical β: Learning task-dependent schedules and organizing incentives across multiple levels of abstraction.
- Deeper evaluation of Bayesian β: Comparing against meta-gradient or annealing strategies on larger benchmarks.

7 Conclusion

We introduced IncentRL, a reinforcement learning framework that integrates external rewards with internal incentives via KL-based preference–prediction alignment. By treating the incentive weight β as a Bayesian random variable, IncentRL adapts online to task demands. Theoretical analysis connects IncentRL to dopamine-inspired reward prediction error and the Free Energy Principle, while experiments on sparse-reward benchmarks show consistent gains in sample efficiency and stability. Our results highlight adaptive incentive shaping as a principled approach to building agents that are not only reward-maximizing but also intrinsically aligned with cognitive motivation. A key takeaway is that Bayesian adaptation of preference gaps removes the long-standing need for manual trade-off tuning in intrinsic motivation, paving the way toward more autonomous and general-purpose RL agents.

REFERENCES

Yuri Burda, Harri Edwards, Deepak Pathak, Amos Storkey, Trevor Darrell, and Alexei A Efros. Large-scale study of curiosity-driven learning. *arXiv preprint arXiv:1808.04355*, 2018.

Will Dabney, Zeb Kurth-Nelson, Naoshige Uchida, Clara Kwon Starkweather, Demis Hassabis, Rémi Munos, and Matthew Botvinick. A distributional code for value in dopamine-based reinforcement learning. *Nature*, 577(7792):671–675, 2020.

- Karl Friston. The free-energy principle: a unified brain theory? *Nature reviews neuroscience*, 11 (2):127–138, 2010.
- Karl Friston, Thomas FitzGerald, Francesco Rigoli, Philipp Schwartenbeck, and Giovanni Pezzulo. Active inference: a process theory. *Neural computation*, 29(1):1–49, 2017.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. Pmlr, 2018.
- Rein Houthooft, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. Vime: Variational information maximizing exploration. *Advances in neural information processing systems*, 29, 2016.
- Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsically motivated reinforcement learning. *Advances in neural information processing systems*, 28, 2015.
- Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Icml*, volume 99, pp. 278–287. Citeseer, 1999.
- Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *International conference on machine learning*, pp. 2778–2787. PMLR, 2017.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Wolfram Schultz, Peter Dayan, and P Read Montague. A neural substrate of prediction and reward. *Science*, 275(5306):1593–1599, 1997.
- Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3 (1):9–44, 1988.

A ADDITIONAL RESULTS ON BAYESIAN ADAPTION

Summary statistics. Table 3 reports the mean, standard deviation, and range of sampled β values across rounds, averaged over multiple seeds. In early rounds (Round 1), β samples are dispersed over a wide range (mean ≈ 0.11 , std ≈ 0.06). As adaptation proceeds, the distribution narrows (Rounds 2–4), with mean values stabilizing below 0.05 and variance decreasing substantially. This posterior concentration around task-relevant values illustrates that Bayesian adaptation can efficiently eliminate poor β regions while retaining effective ones, complementing the qualitative evidence from Figures 3 and 4.

Table 3: Summary statistics of Bayesian β adaptation across rounds (averaged over seeds).

| Round | Mean β | Std β | Min β | Max β |
|-------|--------------|-------------|-------------|-------------|
| 1 | 0.1057 | 0.0573 | 0.0078 | 0.1805 |
| 2 | 0.0454 | 0.0322 | 0.0074 | 0.0977 |
| 3 | 0.0194 | 0.0098 | 0.0035 | 0.0351 |
| 4 | 0.0173 | 0.0182 | 0.0000 | 0.0587 |