

A Message Passing Realization of Expected Free Energy Minimization

Wouter W. L. Nuijten¹, Mykola Lukashchuk¹, Thijs van de Laar¹, and Bert de Vries^{1,2}

¹ Eindhoven University of Technology, 5612 AP Eindhoven, the Netherlands

² GN Hearing, 5612 AB Eindhoven, The Netherlands

Abstract. We present a message passing approach to Expected Free Energy (EFE) minimization on factor graphs, based on the theory introduced in [37]. By reformulating EFE minimization as Variational Free Energy minimization with epistemic priors, we transform a combinatorial search problem into a tractable inference problem solvable through standard variational techniques. Applying our message passing method to factorized state-space models enables efficient policy inference. We evaluate our method on environments with epistemic uncertainty: a stochastic gridworld and a partially observable Minigrid task. Agents using our approach consistently outperform conventional KL-control agents on these tasks, showing more robust planning and efficient exploration under uncertainty. In the stochastic gridworld environment, EFE-minimizing agents avoid risky paths, while in the partially observable minigrid setting, they conduct more systematic information-seeking. This approach bridges active inference theory with practical implementations, providing empirical evidence for the efficiency of epistemic priors in artificial agents.

Keywords: Active Inference · Epistemic Planning · Expected Free Energy · Factor Graphs · Message Passing

1 Introduction

Expected Free Energy (EFE) minimization, rooted in the Free Energy Principle, provides a framework for modeling intelligent behavior by unifying reward-seeking (pragmatic) and information-seeking (epistemic) drives [16,18]. While control-as-inference approaches have made significant advances in formulating decision-making as probabilistic inference problems [20,1], EFE minimization extends this paradigm by explicitly accounting for epistemic uncertainty [13], though its practical application faces computational challenges for extended planning horizons and high-dimensional state-spaces [30].

Traditional approaches to computing EFE often involve evaluating all possible action sequences, which becomes intractable for non-trivial problems. While various approximations have been developed to address this tractability issue, traditional approaches typically use EFE as a cost function for evaluating policies, rather than as an objective functional for variational optimization of beliefs [29,8,19].

This paper provides empirical validation of the theoretical foundation presented in [37], which reformulates EFE minimization directly as a variational inference problem on factor graphs. By introducing appropriate epistemic priors, we show that minimizing EFE can be achieved through standard Variational Free Energy (VFE) minimization, making it consistent with the Free Energy Principle’s core tenet that all processes are fundamentally based on variational free energy minimization.

We implement this approach through an iterative message passing algorithm on factorized state-space models. We evaluate its performance in environments with different uncertainty characteristics: a stochastic gridworld with perilous transitions and a partially observable Minigrid environment requiring active exploration for successful completion. Our results confirm that agents using our inference-based method exhibit the same characteristic advantages over KL-control agents as direct EFE computation, particularly in handling epistemic uncertainty. This validates our approach while providing a computationally efficient framework for planning under uncertainty.

The remainder of this paper is organized as follows:

- Section 2 provides background on necessary materials.
- Section 3 discusses related work in control as inference and active inference.
- Section 4 presents our methodology for reformulating EFE minimization as an inference problem.
- Section 5 describes our evaluation environments and experimental design.

2 Background

2.1 Variational Inference

Variational inference (VI) provides a principled framework for approximating complex posterior distributions in Bayesian models [24,6,7,38]. The central challenge in Bayesian inference is computing the posterior distribution $p(\mathbf{x}|\mathbf{y})$ of hidden state sequence \mathbf{x} given an observed data sequence \mathbf{y} , which requires evaluating the model evidence $p(\mathbf{y})$ [11,21]. This normalization constant is typically intractable for complex models.

VI reformulates inference as an optimization problem by approximating the Bayesian posterior with a simpler, tractable distribution $q(\mathbf{x})$ from a family of distributions \mathcal{Q} [7]. The functional we will minimize is the Variational Free Energy (VFE). The VFE is defined as $F[q] = D_{KL}(q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})) - \log p(\mathbf{y})$, making it clear that minimizing the VFE is equivalent to minimizing the KL divergence since $\log p(\mathbf{y})$ is constant with respect to q . The VFE also provides a tractable upper bound on the negative log evidence, with $F[q] \geq -\log p(\mathbf{y})$ [23].

2.2 Factor Graphs

Factor graphs are a specific type of probabilistic graphical model that explicitly represents the factorization structure of the model, where factors represent

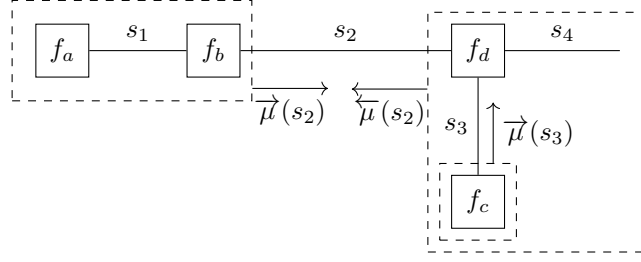


Fig. 1: A Forney-style factor graph representation of the factorization in (2).

(conditional) probability distributions. In our work, we employ Forney-style factor graphs (FFGs) [15], which offer a specific representation approach with notation following [26].

An FFG represents a factorized function $f(\mathbf{s})$ as

$$f(\mathbf{s}) = \prod_{a \in \mathcal{V}} f_a(\mathbf{s}_a), \quad (1)$$

where \mathbf{s} encompasses all variables in the model, and $\mathbf{s}_a \subseteq \mathbf{s}$ represents the subset of variables that participate in factor f_a .

In the FFG representation, nodes ($a \in \mathcal{V}$) correspond to factors in the model, while edges ($\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$) represent variables. An edge connects to a node precisely when the variable appears as an argument in the corresponding factor. We denote the set of edges connected to node $a \in \mathcal{V}$ as $\mathcal{E}(a)$, and the nodes connected to edge $i \in \mathcal{E}$ as $\mathcal{V}(i)$.

To illustrate, the FFG representation of the factorized function

$$f(s_1, s_2, s_3, s_4) = f_a(s_1)f_b(s_1, s_2)f_c(s_3)f_d(s_2, s_3, s_4) \quad (2)$$

is shown in Figure 1.

A common approach to realizing efficient variational inference on factor graphs involves the Bethe assumption, which posits that the posterior distribution factorizes as a product of local marginals associated with the nodes and edges of the graph. This structural assumption on the posterior distributions enables the formulation of message passing algorithms that seek out stationary points of the Bethe free energy [39,40,14].

To illustrate the computational benefits of message passing, consider the generative model in (2), and assume we are interested in computing $p(s_2)$. This marginal distribution can be obtained by summing out all other variables from the joint

$$p(s_2) = \sum_{s_1} \sum_{s_3} \sum_{s_4} f(s_1, s_2, s_3, s_4), \quad (3)$$

which, when each s_i can take 10 values, contains about a thousand terms. However, taking into account the factorization of the generative model and the distributive

law of the product, (3) can be rewritten as

$$p(s_2) = \underbrace{\left(\sum_{s_1} f_a(s_1) f_b(s_1, s_2) \right)}_{\vec{\mu}(s_2)} \cdot \underbrace{\left(\overbrace{\left(\sum_{s_3} f_c(s_3) \right)}^{\vec{\mu}(s_3)} \sum_{s_4} f_d(s_2, s_3, s_4) \right)}_{\overleftarrow{\mu}(s_2)}. \quad (4)$$

The computation in (4) requires only a few hundred summations and is preferable from a computational standpoint. In larger models, the number of computations scale linearly with the number of factor nodes, instead of exponentially. The intermediate results $\vec{\mu}(s_i)$ and $\overleftarrow{\mu}(s_i)$ afford an interpretation as local message in the FFG representation of the model, see Figure 1. For comprehensive treatments of factor graphs and associated (variational) message passing algorithms, we refer readers to [26,27,39,14,40].

3 Related Work

Autonomous decision-making under uncertainty remains a central challenge in control theory and artificial intelligence. This section reviews key developments that contextualize our contribution.

3.1 Control as Inference

The pursuit of efficient and high-performing autonomous systems has driven significant research in control theory. Optimal control [3,4,32] provides a mathematical framework for determining the control inputs that minimize a predefined cost function for a given system. Building upon these foundations, Model Predictive Control (MPC) algorithms address the challenges of real-time control by incorporating a feedback loop and a receding horizon strategy [5,33,34,12]. This approach allows for online adaptation to disturbances and constraints.

A significant paradigm shift in recent years involves viewing control as an inference problem. This perspective allows the application of powerful probabilistic tools to address control challenges, particularly in complex and uncertain environments. Under deterministic dynamics, the sequential decision-making process in closed-loop receding horizon MPC can be elegantly mapped to inference on a factor graph [25,28].

When dealing with stochastic dynamics or the need for state estimation under uncertainty, stochastic optimal control methods can be reformulated using variational inference [22,20]. Here, the intractable posterior distribution over states and/or controls is approximated by a tractable variational distribution.

Active inference [13,10] addresses control under uncertainty by proposing that information gained about the system is also a form of reward. The framework suggests that variational inference naturally balances exploration and exploitation by optimizing the Expected Free Energy [18], which elegantly combines the drive to minimize uncertainty about the environment (information gain) with the need

to achieve desired outcomes. However, a current limitation of active inference lies in the computational cost associated with computing the Expected Free Energy [18], which has spurred recent research into efficient algorithms [29,17,30,8].

Recently, [37] proposed an alternative approach to Expected Free Energy minimization by framing EFE minimization as a regular variational free energy minimization task. This approach is promising for scalable implementation of EFE-minimizing planning algorithms, but offers a theoretical account, without considering practical implementation or empirical validation. In the next section, we will propose a message passing realization of this approach.

4 Methodology

For the main contribution of this paper, we will elaborate on Theorem 1 from [37]. For convenience, we will repeat the theorem here, albeit without the inclusion of model parameters θ :

Theorem 1 (Expected Free Energy Theorem). *Consider an agent with generative model $p(\mathbf{y}, \mathbf{x}, \mathbf{u})$, and prior beliefs $\hat{p}(\mathbf{x})$ about future desired states.*

Consider the Variational Free Energy functional

$$F[q] \triangleq E_{q(\mathbf{y}, \mathbf{x}, \mathbf{u})} \left[\log \frac{\overbrace{q(\mathbf{y}, \mathbf{x}, \mathbf{u})}^{\text{posterior}}}{\underbrace{p(\mathbf{y}, \mathbf{x}, \mathbf{u})}_{\text{generative model}} \underbrace{\hat{p}(\mathbf{x})}_{\text{preference prior}} \underbrace{\tilde{p}(\mathbf{u})\tilde{p}(\mathbf{x})}_{\text{epistemic priors}}} \right], \quad (5)$$

where the generative model in the denominator is augmented by both a preference prior $\hat{p}(\cdot)$ and epistemic priors $\tilde{p}(\cdot)$.

If the epistemic priors are chosen as

$$\tilde{p}(\mathbf{u}) \propto \exp(H[q(\mathbf{x}|\mathbf{u})]) \quad (6a)$$

$$\tilde{p}(\mathbf{x}) \propto \exp(-H[q(\mathbf{y}|\mathbf{x})]) \quad (6b)$$

then $F[q]$ decomposes as

$$F[q] = \underbrace{E_{q(\mathbf{u})}[G(\mathbf{u})]}_{\text{expected policy costs}} + \underbrace{E_{q(\mathbf{y}, \mathbf{x}, \mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}, \mathbf{u})}{p(\mathbf{y}, \mathbf{x}, \mathbf{u})} \right]}_{\text{complexity}} + \text{constant}, \quad (7)$$

where

$$G(\mathbf{u}) = \mathbb{E}_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \left(\frac{q(\mathbf{x}|\mathbf{u})}{\hat{p}(\mathbf{x})} \cdot \frac{1}{q(\mathbf{y}|\mathbf{x})} \right) \right] \quad (8)$$

is the expected free energy as defined in [13]. In (6),

$$H[q(y|x)] = - \int q(y|x) \log q(y|x) dy \quad (9)$$

is the entropy functional.

Proof. The proof of (7) is given in [37, Appendix A]. \square

While (7) shows that minimization of the $F[q]$ leads to minimization of (expected) $G(u)$, the proof of (7) is declarative and does not provide an explicit algorithm for minimizing $F[q]$.

In the following sections, we will describe a message passing algorithm on factor graphs that can be used as a practical approach to search for stationary points of the free energy functional.

4.1 Factorized models and factorized posteriors

Theorem 1 is a general result, however, in practice, we are often interested in factorized state-space models of the form

$$p(\mathbf{y}, \mathbf{x}, \mathbf{u}) = p(x_0) \prod_{t=1}^T p(y_t|x_t)p(x_t|x_{t-1}, u_t)p(u_t) \quad (10)$$

We can make an additional assumption that the posterior distribution factorizes in the same way as the generative model:

$$q(\mathbf{y}, \mathbf{x}, \mathbf{u}) = q(x_0) \prod_{t=1}^T q(y_t|x_t)q(x_t|x_{t-1}, u_t)q(u_t). \quad (11)$$

Note that this is consistent with making the Bethe assumption, which says that the variational posterior distribution can be decomposed into local contributions:

$$q(\mathbf{s}) = \prod_{a \in \mathcal{V}} q_a(\mathbf{s}_a) \prod_{i \in \mathcal{E}} q_i(s_i)^{-1} \quad (12)$$

for $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ the underlying FFG. Under this assumption, we can derive a corollary to Theorem 1 that provides more specific expressions for the epistemic priors.

Corollary 1. *Consider an agent with Variational Free Energy functional as in (5), comprising a generative model (10), a posterior distribution factorized as in (11), and a preference prior $\hat{p}(\mathbf{x}) = \prod_{t=1}^T \hat{p}(x_t)$. If the epistemic priors are chosen as*

$$\tilde{p}(u_t) \propto \exp(H[q(x_t, x_{t-1}|u_t)] - H[q(x_{t-1}|u_t)]) \quad (13a)$$

$$\tilde{p}(x_t) \propto \exp(-H[q(y_t|x_t)]) \quad (13b)$$

then the Variational Free Energy functional (5) decomposes as

$$F[q] = E_{q(\mathbf{u})}[G(\mathbf{u})] + E_{q(\mathbf{y}, \mathbf{x}, \mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}, \mathbf{u})}{p(\mathbf{y}, \mathbf{x}, \mathbf{u})} \right] + \text{constant}. \quad (14)$$

While this corollary is a special case and a direct application of Theorem 1, an elaboration of the proof is given in Appendix A. This corollary states that the preference and epistemic priors can be reduced to local contributions. We will implement the preference and epistemic priors as factor nodes that act as prior distributions during the inference procedure. A timeslice of the augmented factor graph is shown in Figure 2.

The benefit of this approach is that inference on factor graphs is well-understood and can be implemented efficiently using reactive message passing [2]. Effectively, this means that the computational complexity of Expected Free Energy minimization is the same as the computational complexity of variational inference on a factor graph.

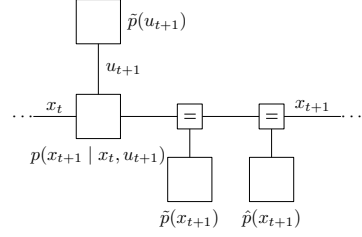


Fig. 2: Slice of the factor graph representation of the augmented generative model. The original generative model (10) is augmented with epistemic priors $\tilde{p}(u_{t+1})$ and $\tilde{p}(x_{t+1})$, and preference priors $\hat{p}(x_{t+1})$ for future timesteps.

4.2 Inferring a policy posterior

Corollary 1 introduces a circular dependency in the model definition: to define the VFE functional with epistemic priors (13), we need access to the variational posterior distribution, but the variational posterior can only be obtained by minimizing the VFE functional given the generative model.

This circular dependency can be resolved through an iterative variational inference procedure implemented as message passing on the factor graph. We first initialize the variational posterior and then iteratively update both the posterior beliefs and epistemic priors until convergence.

On a factor graph, we can implement variational inference using message passing algorithms that iteratively updates posterior distributions [31].

Each message passing iteration τ refines both the posteriors and priors simultaneously. To that extent, let $q_\tau(\cdot)$ be the variational posterior distribution at iteration τ , we then define the epistemic priors as

$$\begin{aligned}\tilde{p}_\tau(u_t) &= \sigma(H[q_{\tau-1}(x_t, x_{t-1}|u_t)] - H[q_{\tau-1}(x_{t-1}|u_t)]) \\ \tilde{p}_\tau(x_t) &= \sigma(-H[q_{\tau-1}(y_t|x_t)])\end{aligned}\quad (15)$$

Here, σ is the softmax function, which guarantees proportionality as in Equations 13a and 13b. A formal description of the algorithm is given in Algorithm 1. While this approach solves the initialization problem, there are some subtleties that need to be addressed. Specifically, although the subtraction of entropies in line Equation 21a results in a constant when using the same variational distribution q for both the epistemic prior \tilde{p} and the optimization, this property no longer holds when we use different distributions - namely, when we use $q_{\tau-1}$ to define

\tilde{p}_τ but optimize with respect to q_τ . While this is not a problem if the inference procedure converges, this convergence is not guaranteed.

Algorithm 1 EFE minimization as VFE minimization

```

1: Input: Factorized generative model  $p(\mathbf{y}, \mathbf{x}, \mathbf{u})$ , preference prior  $\hat{p}(\mathbf{x})$ , number
   of VI iterations  $\tau_{max}$ 
2: Output: Policy posterior  $q_{\tau_{max}}(\mathbf{u})$ 
3:  $q_0(\mathbf{y}, \mathbf{x}, \mathbf{u}) \leftarrow$  Uninformative distribution
4: for  $\tau \leftarrow 1$  to  $\tau_{max}$  do ▷ Iterations of variational inference algorithm
5:   for each time step  $t$  do
6:      $\tilde{p}_\tau(u_t) \leftarrow \sigma(H[q_{\tau-1}(x_t, x_{t-1}|u_t)] - H[q_{\tau-1}(x_{t-1}|u_t)])$ 
7:      $\tilde{p}_\tau(x_t) \leftarrow \sigma(-H[q_{\tau-1}(y_t|x_t)])$ 
8:   end for
9:    $q_\tau(\mathbf{y}, \mathbf{x}, \mathbf{u}) \leftarrow \text{infer}(p(\mathbf{y}, \mathbf{x}, \mathbf{u}))$  ▷ Message passing (4)
10: end for
11: return  $q_{\tau_{max}}(\mathbf{u})$ 

```

5 Evaluation

This section evaluates our EFE-minimizing policy inference method. In this section, we will evaluate the performance of the proposed method. The addition of preference priors is consistent with the literature on KL control [35,36], which means the main point of interest is the influence of the epistemic priors on the policy posterior. To this extent, we will execute the experiments both with and without the epistemic priors, which will correspond to a KL-control and an EFE-minimizing policy, respectively. KL-control is known to be prone to optimistic planning in the face of stochasticity and uncertainty [28,25], so we will explore partially observable Markov decision processes (POMDPs) with stochastic dynamics and observation noise.

For our experimental evaluation, we consider scenarios where the environment dynamics are completely known to the agent, though they may be stochastic or contain inherent uncertainty. This known-dynamics assumption allows us to isolate and evaluate the specific effects of epistemic priors on decision-making, without conflating them with model learning.

5.1 Experimental design

We designed a stochastic grid environment that specifically challenges agents with uncertainty in dynamics and observations. Additionally, we evaluate our method on the Minigrid door-key environment [9], which tests how agents handle partial observability. Both environments highlight the differences between KL-control and EFE-minimizing policies in the presence of epistemic uncertainty.

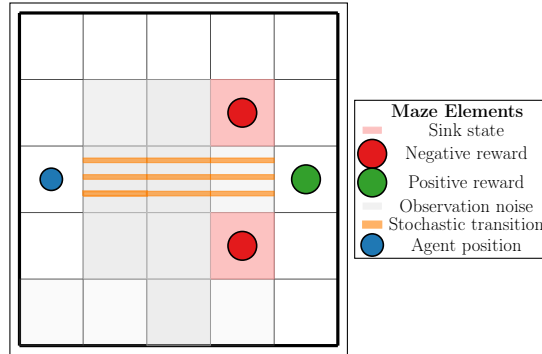


Fig. 3: The stochastic grid environment. The agent should traverse the grid with both stochastic transitions and observation noise. Cells with stochastic transitions appear on the shortest path, creating a risk-reward tradeoff. Opacity for observation noise is used to indicate the uncertainty in the environment.

Stochastic Grid Environment For our first experiment, we focus on a stochastic grid environment. In this environment, the agent has to traverse the grid from one end to the other, with hazards and stochastic transitions. The key challenge is that on the shortest path from the start to the goal, there are cells in which the transition matrix is stochastic, with the risk that the agent will end up in a sink state. The stochasticity presents a direct test of how agents handle uncertainty in dynamics: the KL-control agent is expected to plan optimistically through these uncertain transitions, while the EFE-minimizing agent should recognize the epistemic risk and avoid these cells. This environment also features observation noise, adding another layer of uncertainty that forces the agent to maintain beliefs over possible states rather than having full observability.

A longer but safer path exists that avoids all stochastic transitions. The optimal policy for a risk-aware agent would be to take this safer path, despite it requiring more steps. A visualization of the environment is shown in Figure 3.

The agent receives a reward of 1 for reaching the goal. When ending up in a sink state, the agent receives a penalty of -1 . The full specification of the generative model can be found in Appendix B.

Minigrid Door-Key Environment The second environment we consider is a Minigrid environment, specifically a 4x4 door-key environment. This environment tests a different aspect of epistemic uncertainty, namely, partial observability. The agent has a limited field of view, which means that the agent must actively explore to reduce uncertainty about the environment state.

The task requires the agent to locate and pick up a key, find and open a door, and finally reach the goal square. This multi-step process creates a natural exploration challenge that tests how agents handle partial observability. The agent location, key location, and door location are randomized in each episode,

which means that the agent has epistemic uncertainty about the environment state.

The EFE-minimizing agent should show more directed exploration behavior, actively seeking to reduce uncertainty about the key and door locations. In contrast, the KL-control agent (without epistemic priors) might exhibit less efficient exploration patterns, as it lacks the intrinsic drive to resolve uncertainty.

The Minigrid environment adds another layer of complexity to the task, as the field of view means that the observations are relative to the agent, while the goals are formulated in an external frame of reference. This means that the observation space of the agent is much larger than the state space. The observation space is of size $\approx 5^{49}$, which makes algorithms like Sophisticated Inference [17] intractable. Furthermore, the planning horizon of 22 timesteps makes standard Expected Free Energy computation as policy evaluation intractable. The computational complexity of the door-key environment is where the benefits of the proposed method are most evident.

A visualization of the initial state of the Minigrid environment is shown in Figure 4. The agent receives a reward when reaching the goal, proportional to the number of steps taken. The full specification of the used generative model can be found in Appendix C. The source code and implementation details for all experiments presented in this paper are publicly available in our online repository³.

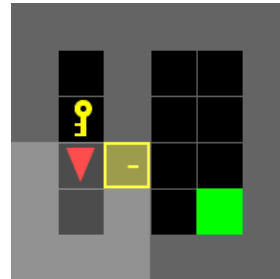


Fig. 4: An initial state of the Minigrid environment. The agent has a limited field of view, indicated by the highlighted cells.

5.2 Results

Stochastic Grid Environment We evaluated the performance of both agents across 100 episodes, Table 1, left, summarizes the quantitative results.

This table suggests distinctly different navigational patterns between both agents. The EFE-minimizing agent consistently chooses the longer but safer path around the stochastic transition cells, demonstrating risk-averse behavior that aligns with theoretical predictions. In contrast, the KL-control agent attempts the shorter path through cells with stochastic transitions, exhibiting the optimistic planning tendency typical of approaches that wrongly account for the system’s aleatoric uncertainty. A more detailed visualization of the trajectories for both agents, as well as an empirical convergence analysis of our algorithm, is provided in Appendix D.

Minigrid Door-Key Environment We evaluated both agents across 200 experimental episodes with a planning horizon of 25 steps. Table 1, right, presents

³ <https://github.com/biaslab/EFEasVFE>

Stochastic Grid			Minigrid Door-Key		
Metric	KL	EFE (ours)	Metric	KL	EFE (ours)
Success Rate	21%	100%	Success Rate	76.5%	88.0%
Avg. Reward	0.22 ± 0.77	1.00 ± 0	Avg. Reward	0.74 ± 0.41	0.85 ± 0.32
-	-	-	Avg. Time to Key Visibility	3.36 ± 6.29	1.34 ± 1.08

Table 1: Performance comparison across environments (100 episodes for Stochastic Grid, 200 episodes for Minigrid).

the quantitative comparison between the EFE-minimizing and KL-control agents in the Minigrid door-key environment.

The EFE-minimizing agent demonstrates more effective exploration patterns, particularly in scenarios requiring active information seeking. This is especially evident in the reduced time needed to locate the key, confirming that epistemic priors enable more directed information-seeking in partially observable environments.

A more detailed visualization of the trajectories for both agents and an empirical convergence analysis of our algorithm is provided in Appendix E.

6 Discussion

Our experimental results demonstrate that agents using the proposed message passing approach for EFE minimization exhibit the characteristic behaviors of active inference: risk-averse path selection in stochastic environments and information-seeking exploration in partially observable settings. These behaviors emerge naturally from the inclusion of epistemic priors in the variational free energy objective, without requiring explicit computation of expected free energy.

The reformulation of EFE minimization as a variational inference problem provides several advantages: it maintains theoretical consistency with the Free Energy Principle’s core tenet; transforms a combinatorial search problem into a tractable inference procedure using message passing on factor graphs; and eliminates the need for ad hoc policy pruning, replacing it with principled reactive processing where the agent minimizes VFE at each point in time. This approach is particularly valuable in complex environments where traditional EFE computation becomes intractable, as demonstrated in our Minigrid experiments.

While our implementation shows promising results, the convergence properties of our iterative approach to handling self-referential epistemic priors require further theoretical investigation. Future research should investigate the inclusion of additional parameters in the generative model, particularly those related to environment dynamics. A natural extension of our work would be to incorporate parameter learning within the epistemic priors. This would allow agents to infer policies that facilitate sample-efficient learning of model parameters. This concept has already been introduced in [37]. However, the exact functional form of the empirical prior has not yet been derived.

7 Conclusion

In this paper, we presented a message passing implementation of Expected Free Energy minimization on factor graphs. Our approach reframes EFE minimization as a variational inference problem, allowing us to use standard message passing algorithms for efficient policy inference. The key insight is that by introducing appropriate epistemic priors, we can transform the expected free energy objective into a modified variational free energy objective that can be optimized through standard inference techniques.

Our experimental results in both stochastic and partially observable environments demonstrate that this approach reproduces the characteristic behaviors of active inference: risk aversion in environments with hazardous stochasticity and information seeking in partially observable environments. The message passing implementation shows significant advantages in computational efficiency compared to traditional methods for computing expected free energy, particularly in complex environments with high-dimensional observation spaces and long planning horizons.

By reformulating EFE minimization as variational inference, our work contributes to unifying the theoretical frameworks of the Free Energy Principle and active inference with practical implementations for decision-making under uncertainty. This bridges the gap between theoretical accounts of intelligent behavior and efficient algorithms for artificial agents, offering a principled approach to balancing pragmatic and epistemic objectives in complex and uncertain environments.

Acknowledgements

This publication is part of the project "ROBUST: Trustworthy AI-based Systems for Sustainable Growth" with project number KICH3.LTP.20.006, which is (partly) financed by the Dutch Research Council (NWO), GN Hearing, and the Dutch Ministry of Economic Affairs and Climate Policy (EZK) under the program LTP KIC 2020-2023.

References

1. Attias, H.: Planning by probabilistic inference. In: International workshop on artificial intelligence and statistics. pp. 9–16. PMLR (2003), <https://proceedings.mlr.press/r4/attias03a.html>
2. Bagaev, D., de Vries, B.: Reactive Message Passing for Scalable Bayesian Inference (Dec 2021). <https://doi.org/10.48550/arXiv.2112.13251>, <http://arxiv.org/abs/2112.13251>, arXiv:2112.13251 [cs]
3. Bellman, R.: The theory of dynamic programming. Bulletin of the American Mathematical Society **60**(6), 503–515 (1954). <https://doi.org/10.1090/S0002-9904-1954-09848-8>, <https://www.ams.org/bull/1954-60-06/S0002-9904-1954-09848-8/>

4. Bellman, R.: Dynamic Programming. *Science* **153**(3731), 34–37 (1966), <https://www.jstor.org/stable/1719695>, publisher: American Association for the Advancement of Science
5. Bertsekas, D.: Dynamic programming and optimal control: Volume I, vol. 4. Athena scientific (2012), <https://books.google.com/books?hl=en&lr=&id=qVBEEAAQBAJ&oi=fnd&pg=PR1&dq=Dynamic+Programming+and+Optimal+Control&ots=x0bAav005n&sig=s3UxthkdznR2UpqCUUsQ7zKgLc>
6. Bishop, C.M., Nasrabadi, N.M.: Pattern recognition and machine learning, vol. 4. Springer (2006), <https://link.springer.com/book/9780387310732>
7. Blei, D.M., Kucukelbir, A., McAuliffe, J.D.: Variational Inference: A Review for Statisticians. *Journal of the American Statistical Association* **112**(518), 859–877 (Apr 2017). <https://doi.org/10.1080/01621459.2017.1285773>, <https://doi.org/10.1080/01621459.2017.1285773>, publisher: Taylor & Francis _eprint: <https://doi.org/10.1080/01621459.2017.1285773>
8. Champion, T., Da Costa, L., Bowman, H., Grześ, M.: Branching Time Active Inference: The theory and its generality. *Neural Networks* **151**, 295–316 (Jul 2022). <https://doi.org/10.1016/j.neunet.2022.03.036>, <https://www.sciencedirect.com/science/article/pii/S0893608022001149>
9. Chevalier-Boisvert, M., Dai, B., Towers, M., Perez-Vicente, R., Willems, L., Lahlou, S., Pal, S., Castro, P.S., Terry, J.: Minigrid & miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks. *Advances in Neural Information Processing Systems* **36**, 73383–73394 (2023), https://proceedings.neurips.cc/paper_files/paper/2023/hash/e8916198466e8ef218a2185a491b49fa-Abstract-Datasets_and_Benchmarks.html
10. Costa, L.D., Tenka, S., Zhao, D., Sajid, N.: Active Inference as a Model of Agency (Jan 2024). <https://doi.org/10.48550/arXiv.2401.12917>, <http://arxiv.org/abs/2401.12917>, arXiv:2401.12917 [cs]
11. Cox, R.T.: Probability, frequency and reasonable expectation. *American journal of physics* **14**(1), 1–13 (1946), <http://www.cs.toronto.edu/~ilya/cox1946.pdf>, publisher: American Association of Physics Teachers
12. Cutler, R.R., Ramaker, B.L.: Dynamic Matrix Control-A Computer Control Algorithm. *Proc. Joint Automatic Control Conference*, 1979 (1979), <https://cir.nii.ac.jp/crid/1570291225777284224>
13. Da Costa, L., Parr, T., Sajid, N., Veselic, S., Neacsu, V., Friston, K.: Active inference on discrete state-spaces: A synthesis. *Journal of Mathematical Psychology* **99**, 102447 (Dec 2020). <https://doi.org/10.1016/j.jmp.2020.102447>, <https://www.sciencedirect.com/science/article/pii/S0022249620300857>
14. Dauwels, J.: On Variational Message Passing on Factor Graphs. In: *IEEE International Symposium on Information Theory*. pp. 2546–2550. Nice, France (Jun 2007). <https://doi.org/10.1109/ISIT.2007.4557602>
15. Forney, G.D.: Codes on graphs: Normal realizations. *IEEE Transactions on Information Theory* **47**(2), 520–548 (2001), <https://ieeexplore.ieee.org/abstract/document/910573/>, publisher: IEEE
16. Friston, K.: The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience* **11**(2), 127–138 (Feb 2010). <https://doi.org/10.1038/nrn2787>, <https://www.nature.com/articles/nrn2787>, number: 2 Publisher: Nature Publishing Group
17. Friston, K., Costa, L.D., Hafner, D., Hesp, C., Parr, T.: Sophisticated Inference (Jun 2020). <https://doi.org/10.48550/arXiv.2006.04120>, <http://arxiv.org/abs/2006.04120>, arXiv:2006.04120 [q-bio]

18. Friston, K., Rigoli, F., Ognibene, D., Mathys, C., Fitzgerald, T., Pezzulo, G.: Active inference and epistemic value. *Cognitive Neuroscience* **6**(4), 187–214 (Oct 2015). <https://doi.org/10.1080/17588928.2015.1020053>, <http://www.tandfonline.com/doi/full/10.1080/17588928.2015.1020053>
19. Friston, K.J., Salvatori, T., Isomura, T., Tschantz, A., Kiefer, A., Verbelen, T., Koudahl, M., Paul, A., Parr, T., Razi, A., Kagan, B.J., Buckley, C.L., Ramstead, M.J.D.: Active Inference and Intentional Behavior. *Neural Computation* **37**(4), 666–700 (Mar 2025). https://doi.org/10.1162/neco_a_01738, https://doi.org/10.1162/neco_a_01738
20. Ito, K., Kashima, K.: Kullback–Leibler control for discrete-time nonlinear systems on continuous spaces. *SICE Journal of Control, Measurement, and System Integration* **15**(2), 119–129 (Jun 2022). <https://doi.org/10.1080/18824889.2022.2095827>, <https://doi.org/10.1080/18824889.2022.2095827>, publisher: Taylor & Francis _eprint: <https://doi.org/10.1080/18824889.2022.2095827>
21. Jaynes, E.T.: *Probability Theory: The Logic of Science*. Cambridge University Press, 1 edn. (Apr 2003). <https://doi.org/10.1017/CB09780511790423>, <https://www.cambridge.org/core/product/identifier/9780511790423/type/book>
22. Kappen, B., Gomez, V., Oppen, M.: Optimal control as a graphical model inference problem. *Machine Learning* **87**(2), 159–182 (May 2012). <https://doi.org/10.1007/s10994-012-5278-7>, <http://arxiv.org/abs/0901.0633>, arXiv:0901.0633 [math]
23. Kingma, D.P., Welling, M.: Auto-Encoding Variational Bayes (Dec 2022). <https://doi.org/10.48550/arXiv.1312.6114>, <http://arxiv.org/abs/1312.6114>, arXiv:1312.6114 [cs, stat]
24. Koller, D., Friedman, N.: *Probabilistic Graphical Models: Principles and Techniques*. MIT Press (Jul 2009), google-Books-ID: 7dzpHCHzNQ4C
25. Levine, S.: Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review (May 2018). <https://doi.org/10.48550/arXiv.1805.00909>, <http://arxiv.org/abs/1805.00909>, arXiv:1805.00909 [cs]
26. Loeliger, H.A.: An introduction to factor graphs. *IEEE Signal Processing Magazine* **21**(1), 28–41 (Jan 2004). <https://doi.org/10.1109/MSP.2004.1267047>, conference Name: IEEE Signal Processing Magazine
27. Loeliger, H.A., Dauwels, J., Hu, J., Korl, S., Ping, L., Kschischang, F.R.: The Factor Graph Approach to Model-Based Signal Processing. *Proceedings of the IEEE* **95**(6), 1295–1322 (Jun 2007). <https://doi.org/10.1109/JPROC.2007.896497>
28. Lázaro-Gredilla, M., Ku, L.Y., Murphy, K.P., George, D.: What type of inference is planning? (Nov 2024). <https://doi.org/10.48550/arXiv.2406.17863>, <http://arxiv.org/abs/2406.17863>, arXiv:2406.17863
29. Paul, A., Sajid, N., Costa, L.D., Razi, A.: On efficient computation in active inference. *Expert Systems with Applications* **253**, 124315 (Nov 2024). <https://doi.org/10.1016/j.eswa.2024.124315>, <http://arxiv.org/abs/2307.00504>, arXiv:2307.00504 [cs]
30. Paul, A., Sajid, N., Gopalkrishnan, M., Razi, A.: Active Inference for Stochastic Control. In: Kamp, M., Koprinska, I., Bibal, A., Bouadi, T., Frénay, B., Galárraga, L., Oramas, J., Adilova, L., Krishnamurthy, Y., Kang, B., Largeron, C., Lijffijt, J., Viard, T., Welke, P., Ruocco, M., Aune, E., Gallicchio, C., Schiele, G., Pernkopf, F., Blott, M., Fröning, H., Schindler, G., Guidotti, R., Monreale, A., Rinzivillo, S., Biecek, P., Ntoutsis, E., Pechenizkiy, M., Rosenhahn, B., Buckley, C., Cialfi, D., Lanillos, P., Ramstead, M., Verbelen, T., Ferreira, P.M., Andresini, G., Malerba, D., Medeiros, I., Fournier-Viger, P., Nawaz, M.S., Ventura, S., Sun, M., Zhou, M.,

- Bitetta, V., Bordino, I., Ferretti, A., Gullo, F., Ponti, G., Severini, L., Ribeiro, R., Gama, J., Gavalda, R., Cooper, L., Ghazaleh, N., Richiardi, J., Roqueiro, D., Saldana Miranda, D., Sechidis, K., Graça, G. (eds.) *Machine Learning and Principles and Practice of Knowledge Discovery in Databases*. pp. 669–680. Springer International Publishing, Cham (2021)
31. Pearl, J.: Reverend Bayes on Inference Engines: A Distributed Hierarchical Approach. In: *AAAI-82 Proceedings*. pp. 133–136. AAAI Press, Carnegie Mellon University, Pittsburgh PA (1982), https://books.google.com/books?hl=nl&lr=&id=e59kEAAAQBAJ&oi=fnd&pg=PA129&ots=qrs53bNhtS&sig=au0q_v4YW1fTTgN0vsmYdDN6RPY, publisher: Morgan & Claypool
 32. Pontryagin, L.S.: *Mathematical Theory of Optimal Processes*. Routledge, London (May 2018). <https://doi.org/10.1201/9780203749319>
 33. RICHALET, J.: Algorithmic control of industrial processes. *Proc. of the 4th IFAC Sympo. on Identification and System Parameter Estimation* pp. 1119–1167 (1976), <https://cir.nii.ac.jp/crid/1570854174674016512>
 34. Richalet, J., Rault, A., Testud, J.L., Papon, J.: Model predictive heuristic control: Applications to industrial processes. *Automatica* **14**(5), 413–428 (Sep 1978). [https://doi.org/10.1016/0005-1098\(78\)90001-8](https://doi.org/10.1016/0005-1098(78)90001-8), <https://www.sciencedirect.com/science/article/pii/0005109878900018>
 35. Todorov, E.: Linearly-solvable Markov decision problems. In: *Advances in Neural Information Processing Systems*. vol. 19. MIT Press (2006), https://proceedings.neurips.cc/paper_files/paper/2006/hash/d806ca13ca3449af72a1ea5aedbed26a-Abstract.html
 36. Todorov, E.: General duality between optimal control and estimation. In: *2008 47th IEEE Conference on Decision and Control*. pp. 4286–4292 (Dec 2008). <https://doi.org/10.1109/CDC.2008.4739438>, <https://ieeexplore.ieee.org/abstract/document/4739438>, ISSN: 0191-2216
 37. Vries, B.d., Nuijten, W., Laar, T.v.d., Kouw, W., Adamiat, S., Nisslbeck, T., Lukashchuk, M., Nguyen, H.M.H., Araya, M.H., Tresor, R., Jenneskens, T., Nikoloska, I., Subramanian, R.G., Erp, B.v., Bagaev, D., Podusenko, A.: Expected Free Energy-based Planning as Variational Inference (Apr 2025). <https://doi.org/10.48550/arXiv.2504.14898>, <http://arxiv.org/abs/2504.14898>, arXiv:2504.14898 [stat]
 38. Winn, J., Bishop, C.: Variational Message Passing. *Journal of Machine Learning Research* **6**, 661–694 (Apr 2005)
 39. Yedidia, J.S., Freeman, W., Weiss, Y.: Constructing free-energy approximations and generalized belief propagation algorithms. *IEEE Transactions on Information Theory* **51**(7), 2282–2312 (Jul 2005). <https://doi.org/10.1109/TIT.2005.850085>
 40. Şenöz, İ.: *Message Passing Algorithms for Hierarchical Dynamical Models*. Phd Thesis 1 (Research TU/e / Graduation TU/e), Eindhoven University of Technology, Eindhoven (Jun 2022), ISBN: 9789038655321

A Proof of Corollary 1

Proof. Proof of Corollary 1. This proof is an adjusted proof of the proof of Theorem 1, which is given in [37].

$$F[q] = E_{q(\mathbf{y}, \mathbf{x}, \mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}, \mathbf{u})}{p(\mathbf{y}, \mathbf{x}, \mathbf{u}) \hat{p}(\mathbf{x}) \tilde{p}(\mathbf{u}) \tilde{p}(\mathbf{x})} \right] \quad (16a)$$

$$= E_{q(\mathbf{u})} \left[\log \frac{q(\mathbf{u})}{p(\mathbf{u})} + \underbrace{E_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}|\mathbf{u})}{p(\mathbf{y}, \mathbf{x}|\mathbf{u}) \hat{p}(\mathbf{x}) \tilde{p}(\mathbf{u}) \tilde{p}(\mathbf{x})} \right]}_{C(\mathbf{u})} \right] \quad (16b)$$

$$= E_{q(\mathbf{u})} \left[\log \frac{q(\mathbf{u})}{p(\mathbf{u})} + \underbrace{G(\mathbf{u}) + E_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}|\mathbf{u})}{p(\mathbf{y}, \mathbf{x}|\mathbf{u})} \right] + \text{constant}}_{C(\mathbf{u}) \text{ if (13) holds}} \right] \quad (16c)$$

$$= E_{q(\mathbf{u})} [G(\mathbf{u})] + E_{q(\mathbf{y}, \mathbf{x}, \mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}, \mathbf{u})}{p(\mathbf{y}, \mathbf{x}, \mathbf{u})} \right] + \text{constant} \quad \text{if (13) holds} \quad (16d)$$

□

In the above derivation, we still need to prove the transition for $C(\mathbf{u})$ from (16b) to (16c), which we address next.

Lemma 1 (Proof of equivalence $C(\mathbf{u})$ in (16b) and (16c)).

$$C(\mathbf{u}) = \mathbb{E}_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \frac{\overbrace{q(\mathbf{y}, \mathbf{x}|\mathbf{u})}^{\text{posterior}}}{\underbrace{p(\mathbf{y}, \mathbf{x}|\mathbf{u})}_{\text{predictive}} \underbrace{\hat{p}(\mathbf{x})}_{\text{utility}} \underbrace{\tilde{p}(\mathbf{u}) \tilde{p}(\mathbf{x})}_{\text{epistemic priors}}} \right] \quad (17a)$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \left(\underbrace{\frac{q(\mathbf{x}|\mathbf{u})}{\hat{p}(\mathbf{x})}}_{\text{risk}} \cdot \underbrace{\frac{1}{q(\mathbf{y}|\mathbf{x})}}_{\text{ambiguity}} \right) \right]}_{G(\mathbf{u}) = \text{Expected Free Energy}} + \quad (17b)$$

$$+ \mathbb{E}_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \left(\underbrace{\frac{\hat{p}(\mathbf{x}) q(\mathbf{y}|\mathbf{x})}{q(\mathbf{x}|\mathbf{u})}}_{\text{inverse factors from } G(\mathbf{u})} \cdot \underbrace{\frac{q(\mathbf{y}, \mathbf{x}|\mathbf{u})}{p(\mathbf{y}, \mathbf{x}|\mathbf{u}) \hat{p}(\mathbf{x}) \tilde{p}(\mathbf{u}) \tilde{p}(\mathbf{x})}}_{\text{factors from (17a)}} \right) \right]$$

$$= G(\mathbf{u}) + \underbrace{\mathbb{E}_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}|\mathbf{u})}{p(\mathbf{y}, \mathbf{x}|\mathbf{u})} \right]}_{=B(\mathbf{u})} + \underbrace{\mathbb{E}_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \frac{q(\mathbf{y}|\mathbf{x})}{q(\mathbf{x}|\mathbf{u}) \tilde{p}(\mathbf{u}) \tilde{p}(\mathbf{x})} \right]}_{\text{choose epistemic priors to let this vanish}} \quad (17c)$$

$$= G(\mathbf{u}) + B(\mathbf{u}) + \mathbb{E}_{q(\mathbf{x}|\mathbf{u})} \left[\log \frac{1}{q(\mathbf{x}|\mathbf{u}) \tilde{p}(\mathbf{u})} \right] + \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{q(\mathbf{y}|\mathbf{x})}{\tilde{p}(\mathbf{x})} \right]$$

Now here we can replace the general $q(\mathbf{y}|\mathbf{x})$ and $q(\mathbf{x}|\mathbf{u})$ with the factorised $\prod_t q(y_t|x_t)$ and $\prod_t q(x_t|x_{t-1}, u_t)$.

$$\begin{aligned} C(\mathbf{u}) &= G(\mathbf{u}) + B(\mathbf{u}) + \mathbb{E}_{q(\mathbf{x}|\mathbf{u})} \left[\log \frac{1}{\prod_t q(x_t|x_{t-1}, u_t) \tilde{p}(u_t)} \right] + \\ &\quad + \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log q(y_t|x_t) - \log \tilde{p}(x_t) \right] \end{aligned} \quad (18a)$$

$$\begin{aligned} &= G(\mathbf{u}) + B(\mathbf{u}) \\ &\quad + \sum_{\mathbf{x}} \mathbb{E}_{q(x_t, x_{t-1}|u_t)} \left[-\log q(x_t|x_{t-1}, u_t) - \log \tilde{p}(u_t) \right] + \\ &\quad + \sum_{\mathbf{y}} \mathbb{E}_{q(y_t|x_t)} \left[\log q(y_t|x_t) - \log \tilde{p}(x_t) \right]. \end{aligned} \quad (18b)$$

Now we can recognize the following:

$$\mathbb{E}_{q(x_t, x_{t-1}|u_t)} \left[-\log q(x_t|x_{t-1}, u_t) \right] \quad (19a)$$

$$= \mathbb{E}_{q(x_t, x_{t-1}|u_t)} \left[-\left(\log q(x_t, x_{t-1}|u_t) - \log q(x_{t-1}|u_t) \right) \right] \quad (19b)$$

$$= H[q(x_t, x_{t-1}|u_t)] - H[q(x_{t-1}|u_t)], \quad (19c)$$

and

$$\mathbb{E}_{q(y_t|x_t)} [\log q(y_t|x_t)] = -H[q(y_t|x_t)]. \quad (20)$$

Which, when substituted into (18b), together with the definitions of $\tilde{p}(u_t)$ and $\tilde{p}(x_t)$, yields

$$\begin{aligned} &= G(\mathbf{u}) + B(\mathbf{u}) \\ &\quad + \sum_{\mathbf{x}} \underbrace{H[q(x_t, x_{t-1}|u_t)] - H[q(x_{t-1}|u_t)] - \log \tilde{p}(u_t)}_{=c_x \text{ if } \tilde{p}(u_t) \propto \exp(H[q(x_t, x_{t-1}|u_t)] - H[q(x_{t-1}|u_t)])} + \\ &\quad + \sum_{\mathbf{y}} \underbrace{H[q(y_t|x_t)] - \log \tilde{p}(x_t)}_{=c_y \text{ if } \tilde{p}(x_t) \propto \exp(-H[q(y_t|x_t)])} \end{aligned} \quad (21a)$$

$$= G(\mathbf{u}) + \mathbb{E}_{q(\mathbf{y}, \mathbf{x}|\mathbf{u})} \left[\log \frac{q(\mathbf{y}, \mathbf{x}|\mathbf{u})}{p(\mathbf{y}, \mathbf{x}|\mathbf{u})} \right] + c_x + c_y, \quad \text{if (13) holds.} \quad (21b)$$

B Generative Model for the Gridworld Environment

The generative model for the stochastic grid environment is defined as follows:

$$x_0 \sim p(x_0) \quad (22a)$$

$$x_t \sim \text{Cat}(x_t | x_{t-1}, u_t, B) \quad (22b)$$

$$y_t \sim \text{Cat}(y_t | x_t, A) \quad (22c)$$

$$x_T \sim \hat{p}(x_T). \quad (22d)$$

Here, s_t represents the agent's state at time t , y_t is the observation, and u_t is the action. The transition dynamics are governed by B , and A represents the observation model. The agent starts with a prior belief $p(s_0)$ and aims to reach the goal state by the end of the planning horizon T .

In the case of the KL-control agent, the prior on the control is given by

$$u_t \sim \text{Cat}(u_t | \mathbf{1}/4) \quad \text{for } t = 1, \dots, T. \quad (23)$$

The EFE-minimizing agent uses empirical priors on the control and the states, given by

$$u_t \sim \text{Cat}(u_t | \sigma(H[q(x_t, x_{t-1} | u_t)] - H[q(x_{t-1} | u_t)])) \quad \text{for } t = 1, \dots, T \quad (24a)$$

$$x_t \sim \text{Cat}(x_t | \sigma(-H[q(y_t | x_t)])) \quad \text{for } t = 1, \dots, T. \quad (24b)$$

C Generative Model for the Minigrid Environment

The generative model for the Minigrid environment uses a factorized state and observation space, which makes the model computationally tractable. Here, the location of the agent is denoted by l , the orientation by o , the key-door state by s , the door location by d , and the key location by k . The key-door state is a categorical variable with three possible values: $\{0, 1, 2\}$, where 0 indicates that the key is not picked up yet, 1 indicates that the key is picked up but the door is not opened yet, and 2 indicates that the key is picked up and the door is opened. For the observations, $y_{t,(x,y)}$ is the observation at time t for cell (x, y) of the field of view. The generative model for the Minigrid environment is defined as follows:

$$l_0 \sim p(l_0) \quad (25a)$$

$$o_0 \sim p(o_0) \quad (25b)$$

$$s_0 \sim p(s_0) \quad (25c)$$

$$d \sim p(d) \quad (25d)$$

$$k \sim p(k) \quad (25e)$$

$$l_t \sim \text{Cat}(l_t | l_{t-1}, o_{t-1}, k, d, s_{t-1}, u_t, B^l) \quad (25f)$$

$$o_t \sim \text{Cat}(o_t | o_{t-1}, B^o, u_t) \quad (25g)$$

$$s_t \sim \text{Cat}(s_t | s_{t-1}, l_{t-1}, o_{t-1}, k, d, u_{t-1}, B^s) \quad (25h)$$

$$y_{t,(x,y)} \sim \text{Cat}(y_{t,(x,y)} | l_t, o_t, k, d, s_t, A_{(x,y)}) \quad \forall (x, y) \in \{1, \dots, 7\}^2 \quad (25i)$$

with terminal state goal priors

$$l_T \sim \hat{p}(l_T) \quad (26a)$$

$$s_T \sim \text{Cat}(s_T | [0, 0, 1]) = \hat{p}(s_T), \quad (26b)$$

where B^l is the location transition tensor, B^o is the orientation transition tensor, B^s is the key-door state transition tensor, and $A_{(x,y)}$ are the observation tensors for each cell in the field of view.

In the case of the KL-control agent, the prior on the control is given by

$$u_t \sim \text{Cat}(u_t | \mathbf{1}/5) \quad \text{for } t = 1, \dots, T. \quad (27)$$

The EFE-minimizing agent uses empirical priors on the control and the states, given by

$$u_t \sim \text{Cat}(u_t | \sigma(\quad (28a)$$

$$H[q(l_t, l_{t-1}, o_{t-1}, k, d, s_{t-1} | u_t)] - H[q(l_{t-1}, o_{t-1}, k, d, s_{t-1} | u_t)] + \quad (28b)$$

$$H[q(s_t, s_{t-1}, l_t, o_{t-1}, k, d | u_t)] - H[q(s_{t-1}, l_t, o_{t-1}, k, d | u_t)] + \quad (28c)$$

$$H[q(o_t, o_{t-1} | u_t)] - H[q(o_{t-1} | u_t)]) \quad \text{for } t = 1, \dots, T \quad (28d)$$

$$l_t \sim \text{Cat} \left(l_t | \sigma \left(\sum_{(x,y)} -H[q(y_{t,(x,y)}, o_t, s_t, k, d, | l_t)] + H[o_t, k, d, s_t | l_t] \right) \right) \quad (28e)$$

$$o_t \sim \text{Cat} \left(o_t | \sigma \left(\sum_{(x,y)} -H[q(y_{t,(x,y)}, l_t, s_t, k, d, | o_t)] + H[l_t, s_t, k, d | o_t] \right) \right) \quad (28f)$$

$$s_t \sim \text{Cat} \left(s_t | \sigma \left(\sum_{(x,y)} -H[q(y_{t,(x,y)}, l_t, o_t, k, d, | s_t)] + H[l_t, o_t, k, d | s_t] \right) \right) \quad (28g)$$

$$k \sim \text{Cat} \left(k | \sigma \left(\sum_t \sum_{(x,y)} -H[q(y_{t,(x,y)}, l_t, o_t, s_t, d, | k)] + H[l_t, o_t, s_t, d | k] \right) \right) \quad (28h)$$

$$d \sim \text{Cat} \left(d | \sigma \left(\sum_t \sum_{(x,y)} -H[q(y_{t,(x,y)}, l_t, o_t, s_t, k, | d)] + H[l_t, o_t, s_t, k | d] \right) \right). \quad (28i)$$

D Additional Results for the Stochastic Grid Environment Experiments

In this section, we will provide further analysis of the results presented in section 5.2. We will go into more detail on the convergence of the Bethe Free Energy over different iterations of the variational inference procedure, and we will elaborate on a trajectory in a specific episode.

D.1 Convergence Analysis

In Figure 5, we plot the Bethe Free Energy over the iterations of the message passing procedure, along the state of the environment at which the inference procedure is being called.

As we can see, even though we have not provided a proof of convergence, in this specific example, the Bethe Free Energy converges to a constant value, indicating that our inference procedure has converged.

Note that the Bethe Free Energy is an approximation of the true Variational Free Energy, and can therefore not be used for model comparison [39]. Although RxInfer minimizes the Bethe Free Energy, this explains the upwards trend in the Bethe Free Energy curve, and we can only use the Bethe Free Energy as a sanity check to check convergence of the inference procedure.

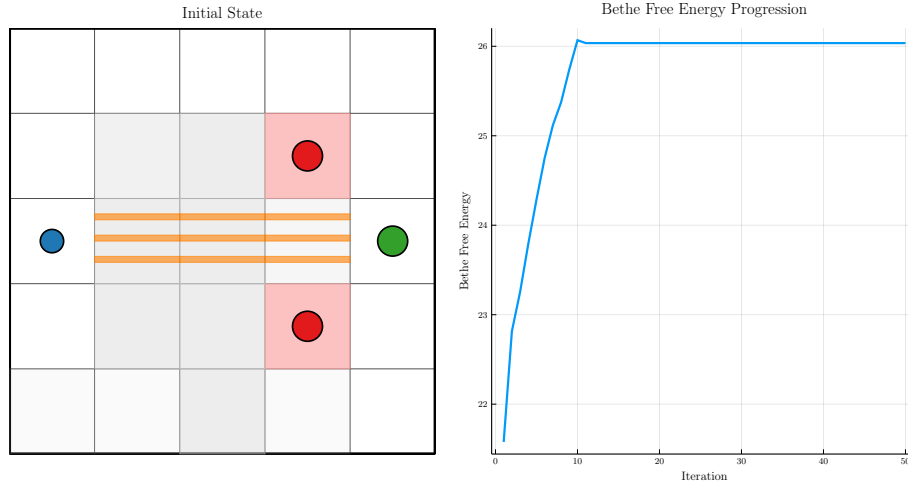


Fig. 5: Visualization of the inference results for the stochastic grid environment. On the left, the initial state of the environment is shown. On the right we show the Bethe Free Energy curve over the iterations of message passing. Convergence to a constant value indicates convergence of the inference procedure.

D.2 Trajectory

Figure 6 provides a frame-by-frame comparison of the trajectories taken by the EFE-minimizing agent (left) and the KL-control agent (right) in the stochastic grid environment. This visualization clearly demonstrates the differences in planning strategies between the two approaches.

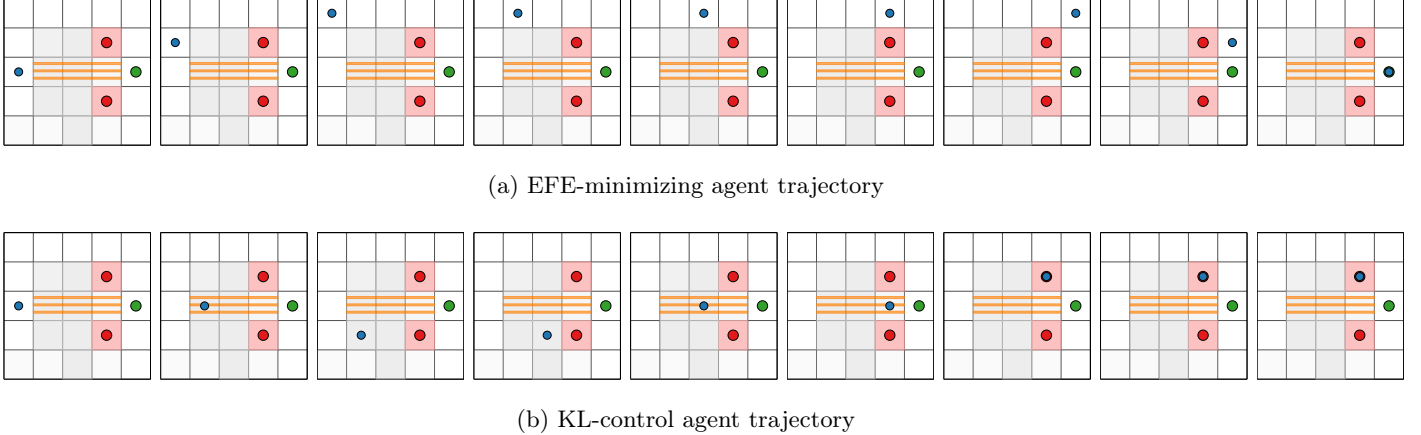


Fig. 6: Comparison of agent trajectories in a stochastic maze environment. Top: EFE-minimizing agent with epistemic priors. Bottom: KL-control agent without epistemic priors.

The EFE-minimizing agent immediately chooses the longer but safer path, moving upward and around the cells with stochastic transitions. This risk-averse behavior is a direct result of the epistemic priors that penalize uncertainty in transitions. By frame $t = 8$, the agent has successfully reached the goal state without encountering any hazardous transitions.

In contrast, the KL-control agent attempts to optimize for the shortest path, moving directly through cells with stochastic transitions. This optimistic planning is characteristic of approaches that don't account for aleatoric uncertainty. While this strategy would be optimal in a deterministic environment, it leads to potential failures in this stochastic setting because the agent cannot manipulate its own luck.

The difference in trajectories directly translates to the performance gap observed across the 100 trial episodes. The EFE-minimizing agent's perfect success rate (100%) compared to the KL-control agent's lower performance (21%) confirms the theoretical prediction that incorporating epistemic uncertainty leads to more robust planning in stochastic environments.

E Additional Results for the Minigrid Environment Experiments

E.1 Convergence Analysis

In Figure 8, we perform inference on the initial state of the Minigrid environment shown in Figure 7. The figure displays the Bethe Free Energy progression during the inference process, along with the agent’s final beliefs about its current location, orientation, and the state of the key and door after the last iteration. We observe that the BFE stabilizes to a constant value, indicating that our inference procedure successfully converges.

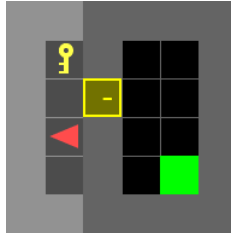


Fig. 7: Initial state of the Minigrid environment.

E.2 Trajectory

Figures 9 and 10 provide a frame-by-frame comparison of the trajectories taken by the EFE-minimizing agent and the KL-control agent in the Minigrid environment. This visualization clearly demonstrates the differences in planning strategies between the two approaches, and highlights the shortcomings of the KL-control approach.

The EFE-minimizing agent is able to solve the task at hand, while the KL-control agent stays in the corner of the grid facing the wall. As shown in Figure 9, the EFE-minimizing agent reaches the goal state, while the KL-control agent does not.

The difference in trajectories directly translates to the performance gap observed across the test episodes. The EFE-minimizing agent’s superior performance confirms our theoretical prediction that incorporating epistemic uncertainty leads to more efficient planning in partially observable environments like Minigrid.

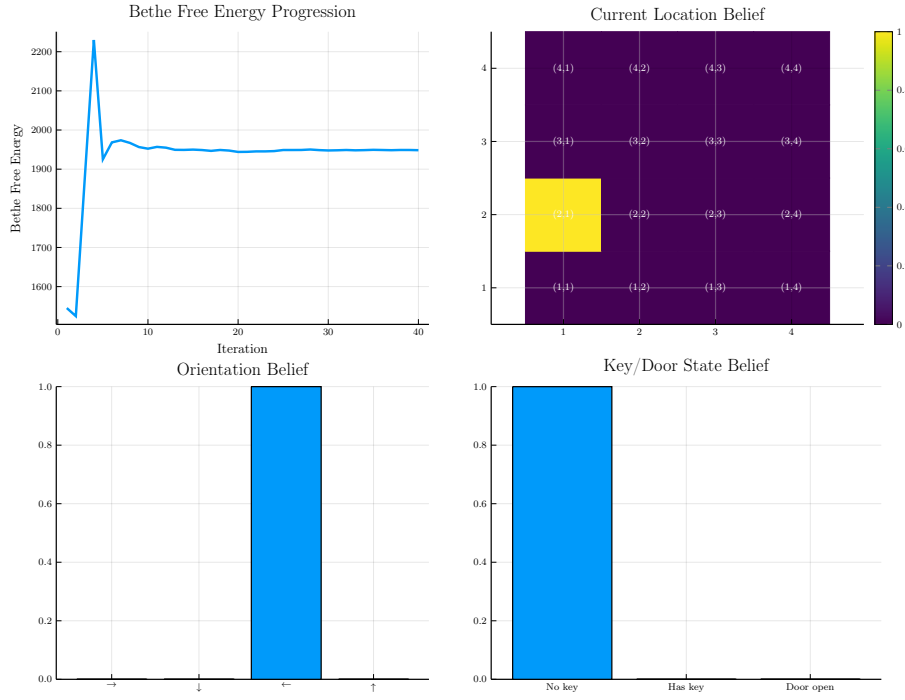


Fig. 8: Visualization of the inference results for the Minigrid environment. Top left: Bethe Free Energy curve over the iterations of message passing. Top right: Agent's belief of its current location after the last iteration. Bottom left: Agent's belief of its current orientation after the last iteration. Bottom right: Agent's belief of the state of the key and door after the last iteration.

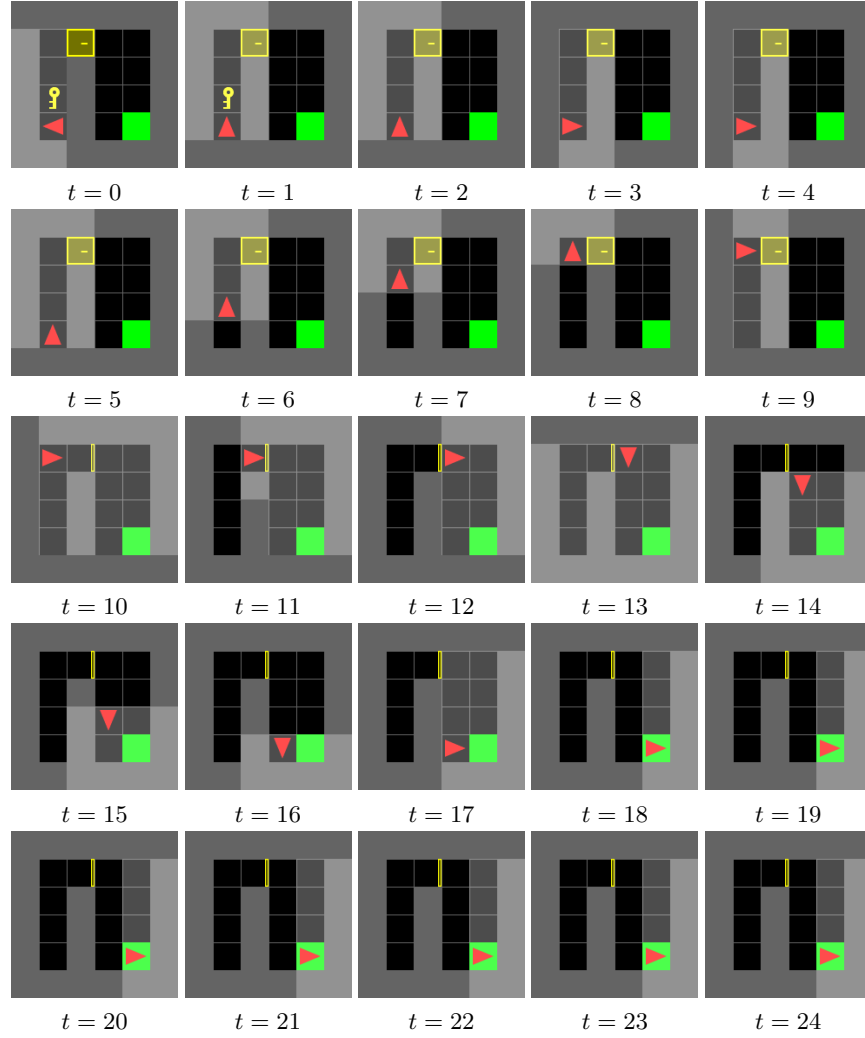


Fig. 9: Visualization of the agent's trajectory in the Minigrid environment using EFE-based control. The 5×5 grid shows the sequential frames of the agent's movement.

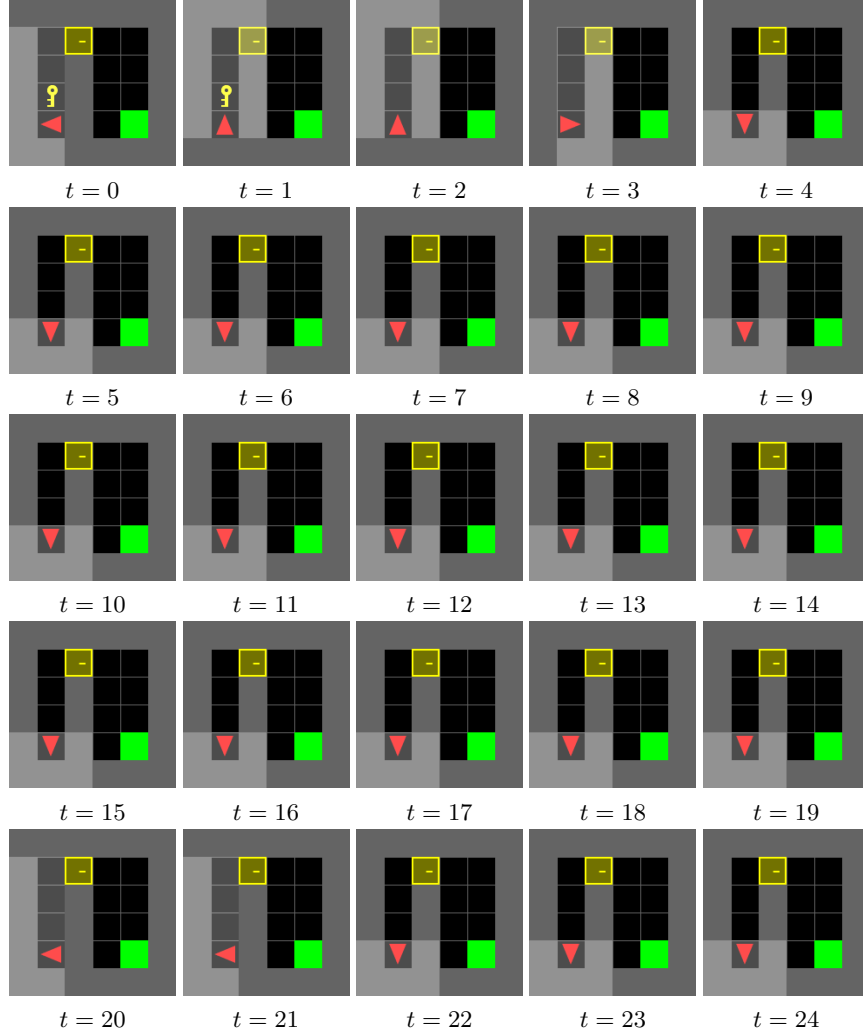


Fig. 10: Visualization of the agent's trajectory in the Minigrid environment using KL control. The 5×5 grid shows the sequential frames of the agent's movement throughout the episode.