

Branching time active inference

Champion, Théophile; Da Costa, Lancelot; Bowman, Howard; Grześ, Marek

DOI:

[10.1016/j.neunet.2022.03.036](https://doi.org/10.1016/j.neunet.2022.03.036)

License:

Creative Commons: Attribution (CC BY)

Document Version

Publisher's PDF, also known as Version of record

Citation for published version (Harvard):

Champion, T, Da Costa, L, Bowman, H & Grześ, M 2022, 'Branching time active inference: the theory and its generality', *Neural Networks*, vol. 151, pp. 295-316. <https://doi.org/10.1016/j.neunet.2022.03.036>

[Link to publication on Research at Birmingham portal](#)

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

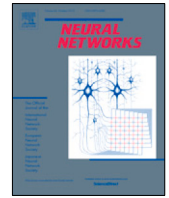
Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.



Branching Time Active Inference: The theory and its generality

Théophile Champion^{a,*}, Lancelot Da Costa^{b,c}, Howard Bowman^{d,a}, Marek Grzes^a

^a University of Kent, School of Computing, Canterbury CT2 7NZ, United Kingdom

^b Imperial College London, Department of Mathematics, London SW7 2AZ, United Kingdom

^c Wellcome Centre for Human Neuroimaging, University College London, London, WC1N 3AR, United Kingdom

^d University of Birmingham, School of Psychology, Birmingham B15 2TT, United Kingdom

ARTICLE INFO

Article history:

Received 23 October 2021

Received in revised form 25 March 2022

Accepted 28 March 2022

Available online 6 April 2022

Keywords:

Active inference

Variational message passing

Tree search

Planning

Free energy principle

ABSTRACT

Over the last 10 to 15 years, active inference has helped to explain various brain mechanisms from habit formation to dopaminergic discharge and even modelling curiosity. However, the current implementations suffer from an exponential (space and time) complexity class when computing the prior over all the possible policies up to the time-horizon. Fountas et al. (2020) used Monte Carlo tree search to address this problem, leading to impressive results in two different tasks. In this paper, we present an alternative framework that aims to unify tree search and active inference by casting planning as a structure learning problem. Two tree search algorithms are then presented. The first propagates the expected free energy forward in time (i.e., towards the leaves), while the second propagates it backward (i.e., towards the root). Then, we demonstrate that forward and backward propagations are related to active inference and sophisticated inference, respectively, thereby clarifying the differences between those two planning strategies.

© 2022 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Active inference is at this point a compelling explanatory approach in cognitive neuroscience, and significant analyses of biologically-realistic implementations in both neural and non-neural communication networks have been conducted. More specifically, active inference extends the free energy principle to generative models with actions (Champion, Grzes, & Bowman, 2021; Da Costa, Parr et al., 2020; Friston, FitzGerald, Rigoli, Schwartenbeck, Doherty, & Pezzulo, 2016) and can be regarded as a form of planning as inference (Botvinick & Toussaint, 2012). This framework has successfully explained a wide range of neuro-cognitive phenomena, such as habit formation (Friston et al., 2016), Bayesian surprise (Itti & Baldi, 2009), curiosity (Schwartenbeck et al., 2018), and dopaminergic discharges (FitzGerald, Dolan, & Friston, 2015). It has also been applied to a variety of tasks, such as animal navigation (Fountas, Sajid, Mediano, & Friston, 2020), robotic control (Pezzato, Corbato, & Wisse, 2020; Sancaktar, van Gerven, & Lanillos, 2020), the mountain car problem (Çatal, Verbelen, Nauta, Boom, & Dhoedt, 2020), the game of DOOM (Cullen, Davey, Friston, & Moran, 2018) and the cart pole problem (Mildridge, 2019). Many of those applications require planning several

steps into the future in order to be solved successfully. However, as explained in more depth in Appendix H, an exhaustive search over all possible sequences of actions will quickly become intractable, i.e., the number of sequences to explore grows exponentially with the time horizon of planning. Fig. 1 illustrates this exponential growth. Exploring only a subset of this exponential number of possible sequences using a tree search therefore becomes a compelling and quite natural alternative.

But what exactly is active inference? Imagine a basketball player at the top of the key (i.e., the area just below the net) ready to take a shot. Intuitively, active inference sees the world as a collection of external states such as the positions of the net, the player and the ball. The player (or agent) is equipped with sensors (such as the eyes) which allow for measurements of the external states. The player is also able to perform actions in the world such as to perform sudden eye movement or simply unfolding his (or her) arms and legs. Furthermore, it is believed that the agent stores an internal representation of the external states, that we shall refer to as the internal states. Importantly, the external and internal states are separated from each other by the Markov blanket (Kirchhoff, Parr, Palacios, Friston, & Kiverstein, 2018), i.e., the sensory information received and actions taken by the agent. In other words, the external states can only modify the internal states indirectly through the observations (also called sensory information) made by the agent, and the internal states can only modify the external states indirectly through the actions taken by the agent.

* Corresponding author.

E-mail addresses: tmac3@kent.ac.uk (T. Champion),

l.da-costa@imperial.ac.uk (L. Da Costa), H.Bowman@kent.ac.uk (H. Bowman), m.grzes@kent.ac.uk (M. Grzes).

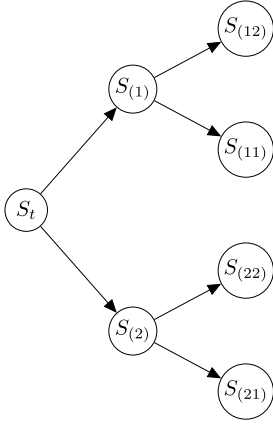


Fig. 1. Illustration of all possible policies up to two time steps in the future when $|U| = 2$. The state at the current time step is denoted by S_t . Additionally, each branch of the tree corresponds to a possible policy, and each node S_j is indexed by a multi-index (e.g. $I = (12)$) representing the sequence of actions that led to this state. This should make it clear that for one time step in the future, there are $|U|$ possible policies, after two time steps there are $|U|^2$ times more policies, and so on until the time-horizon T where there are a total of $|U|^T$ possible policies, i.e., the number of possible policies grows exponentially with the number of time steps for which the agent tries to plan.

More formally, active inference builds on a subfield of Bayesian statistics called variational inference (Fox & Roberts, 2012), in which the true posterior distribution is approximated with a variational distribution. This method provides a way to balance the complexity and accuracy of the posterior distribution. The variational approach is only tractable because some statistical dependencies are ignored during the inference process, i.e., the variational distribution is generally assumed to fully factorize, leading to the well known mean-field approximation:

$$Q(X) = \prod_i Q(X_i), \quad (1)$$

where X is the set of all hidden variables of the model, X_i represents the i th hidden variable, $Q(X)$ is the variational distribution (see below) approximating the posterior $P(X|O)$ where O is the available data, and $Q(X_i)$ is the i th factor of the variational distribution. In 2005, Winn and Bishop (2005) presented a message-based implementation of variational inference, which has naturally been called variational message passing. And more recently, Champion, Grześ et al. (2021) realized an active inference scheme using this variational message passing procedure. By combining the Forney factor graph formalism (Forney, 2001) with the method of Winn and Bishop (2005), it becomes possible to create modular implementations of active inference (Cox, van de Laar, & de Vries, 2019; van de Laar & de Vries, 2019) that allows users to define their own generative models without the burden of deriving update equations.

However, as just stated, there is a major bottleneck to scaling up the active inference framework: the number of action sequences grows exponentially with the time-horizon (see Appendix H for details). In the reinforcement learning literature, this explosion is frequently handled using Monte Carlo tree search (MCTS) (Browne et al., 2012; Schrittwieser et al., 2019; Silver et al., 2016). This approach has been applied to active inference in several papers (Fountas et al., 2020; Maisto, Gregoret, Friston, & Pezzulo, 2021). Fountas et al. (2020) chose to modify the original criterion used during the node selection step in MCTS. This step returns the node that needs to be expanded, and the reinforcement learning community uses the upper confidence

bound for trees (UCT) introduced by Kocsis and Szepesvári (2006) as a selection criterion:

$$UCT_j = \bar{X}_j + 2C_p \sqrt{\frac{2 \ln n}{n_j}}, \quad (2)$$

where n is the number of times the current (parent) node has been explored; n_j stands for the number of times the j th child node has been explored; $C_p > 0$ is the exploration constant and \bar{X}_j is the average reward received by the j th child, i.e., the sum of all rewards received by the current node and its descendants divided by n_j . The child node with the largest UCT_j is selected. In their paper, Fountas et al. (2020) replaced this selection criterion by:

$$U(s, a) = -\tilde{G}(s, a) + C_{\text{explore}} Q(a|s) \frac{1}{1 + N(s, a)} \quad (3)$$

where $U(s, a)$ indicates the utility of selecting action a in state s ; $N(s, a)$ is the number of times that action a was explored in state s ; C_{explore} is an exploration constant equivalent to C_p in the UCT criterion; $Q(a|s)$ is a neural network modelling the posterior distribution over actions, which is trained by minimizing the variational free energy, and $\tilde{G}(s, a)$ is an estimator of the expected free energy (EFE). The EFE is computed from the following equation:

$$\begin{aligned} G(\pi, \tau) = & -\mathbb{E}_{Q(\theta|\pi)} Q(s_\tau | \theta, \pi) Q(o_\tau | s_\tau, \theta, \pi) \left[\ln P(o_\tau | \pi) \right] \\ & + \mathbb{E}_{Q(\theta|\pi)} \left[\mathbb{E}_{Q(o_\tau | \theta, \pi)} H(s_\tau | o_\tau, \pi) - H(s_\tau | \pi) \right] \\ & + \mathbb{E}_{Q(\theta|\pi)} Q(s_\tau | \theta, \pi) H(o_\tau | s_\tau, \theta, \pi) - \mathbb{E}_{Q(s_\tau | \pi)} H(o_\tau | s_\tau, \pi), \end{aligned} \quad (4)$$

where $H(x|y)$ is the entropy of $p(x|y)$. The computation of the EFE is performed by sampling from three distributions whose parameters are predicted by deep neural networks, i.e., the encoder network modelling $Q(s_\tau)$, the decoder network modelling $P(o_\tau | s_\tau)$ and the transition network modelling $P(s_\tau | s_{\tau-1}, a_{\tau-1})$. Note that Eq. (2) was developed by Kocsis and Szepesvári (2006) as a criterion for selecting nodes during planning, such that the selected node minimizes the agent's regret (cf. Appendix G for additional details). Eq. (3) finds its origin in the Predictor Upper Confidence Bound (PUCB) algorithm introduced by Rosin (2010). The idea of the PUCB algorithm is to use contextual information to predict the node to select during planning. Eqs. (2) and (3) both aim to select the node that minimizes the agent's regret, and can therefore be used interchangeably. However, Eq. (3) requires contextual information and a model predicting the node to be selected. Fountas et al. (2020) proposed to use the neural network modelling $Q(a|s)$ as a predictor. This has the advantage of making the predictor very flexible, since neural networks are known to be general function approximators, but neural networks are also expensive to train and lack interpretability.

To avoid the additional complexity brought by the predictor, this paper makes use of (2), which arises from the multi-armed bandit literature (Auer, Cesa-Bianchi, & Fischer, 2002). The idea is to minimize the agent's regret to handle the trade-off between exploration and exploitation at the tree-level in an optimal manner.

A major novelty of our paper is to think about tree search as a dynamical expansion of the generative model, where the past and present is modelled as a partially observable Markov decision process (Sondik, 1971) and the future is modelled by a tree-like generative model. Importantly, our agent treats future states and observations as latent variables over which posterior beliefs are computed, and those beliefs encode the uncertainty of our agent over future states. In contrast, Fountas et al. (2020) are using a maximum a posteriori (MAP) estimate of the future hidden states, while performing MCTS. Lastly, the posterior beliefs held by our agent are computed using variational message passing as

presented in [Champion, Grzeş et al. \(2021\)](#). In comparison, [Fountas et al. \(2020\)](#) perform amortized inference using an encoder network that predicts the mean and variance of the posterior distribution over latent states. Then, (during planning) a MAP estimate is used as input for the neural network modelling the temporal transition. All those neural networks are trained using gradient descent on the variational free energy.

Overall, the key contribution of our paper is to use MCTS to expand or grow the probabilistic graphical model, treat future states and observations as latent variables, and do inference using variational message passing. Indeed, the definition in [Champion, Grzeş et al. \(2021\)](#) of a general message passing procedure for performing active inference makes it possible to construct graphical active inference models in a modular fashion. In turn, this makes it possible to incrementally expand an active inference model as required of our MCTS procedure. It is this message passing procedure that makes our approach possible. To our knowledge, an approach of this kind has never been studied before.

In the following, we first provide the requisite background concerning Forney factor graphs, variational message passing, active inference, and Monte Carlo tree search in Sections 2, 3, 4, and 5, respectively. Next, Section 6 introduces our method that frames planning using a tree as a form of Bayesian model extension. Using terminology from concurrency theory ([Bowman, 2005](#)), we call our new formalism *Branching Time Active Inference* (BTAI). In this domain, models of systems based upon sequences of actions (the format of policies) are described as *linear time*, while models based upon tree and even graph structures are called *branching time* ([Bowman, 2005](#); [van Glabbeek, 1990, 1993](#)). Importantly, BTAI does not consider the generative model and the tree as two different objects, instead, BTAI merges those two objects together into a generative model that can be dynamically expanded. For a detailed analysis of the properties of BTAI, the reader is referred to our companion paper ([Champion, Bowman, & Grzeş, 2021](#)), which provides an empirical demonstration of the benefits of BTAI over standard active inference (Acl) in the context of a graph navigation task. This companion paper also supplies a theoretical comparison of BTAI and standard Acl based upon a complexity class analysis. Briefly, standard Acl has a space complexity class of $O(|\pi| \times T \times |S|)$, where $|\pi| = |U|^T$ is the number of possible policies, $|U|$ is the number of available actions, T is the time horizon of planning, and $|S|$ is the number of values that the hidden state can take. In contrast, the space complexity class of BTAI is $O((K + t) \times |S|)$, where t is the current (i.e. present) time point, and K is the number of expansions of the tree performed during planning. Importantly, even complex applications such as the game of Go can be solved by expanding only a small number of nodes ([Schrittwieser et al., 2019](#); [Silver et al., 2016](#)). Section 6 is followed by Section 7 that explains the connection between our method and the planning strategies used in both active inference and sophisticated inference ([Friston, Da Costa, Hafner, Hesp, & Parr, 2021](#)). Finally, Section 8 concludes this paper and provides ideas for future research.

2. Forney factor graphs

A Forney factor graph ([Forney, 2001](#)) uses three kinds of nodes. The nodes representing hidden and observed variables are depicted by white and grey circles, respectively. And the distribution's factors are represented using white squares, which are linked to variable nodes by arrows or lines. Arrows are used to connect factors to their target variable, while lines link factors to their predictors. [Fig. 2](#) shows an example of a Forney factor graph corresponding to the following generative model:

$$P(O, S) = P_O(O|S)P_S(S). \quad (5)$$

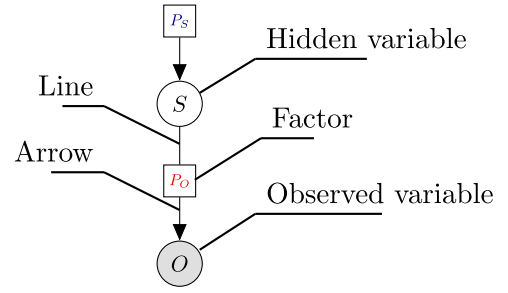


Fig. 2. This figure illustrates the Forney factor graph corresponding to the following generative model: $P(O, S) = P_O(O|S)P_S(S)$. The hidden state is represented by a white circle with the variable's name at the centre, and the observed variable is depicted similarly but with a grey background. The factors of the generative model are represented by squares with a white background and the factor's name at the centre. Finally, arrows connect the factors to their target variable and lines link each factor to its predictor variables.

Generally, factor graphs only describe the model's structure such as the variables and their dependencies, but do not specify the definition of individual factors. For example, the definitions of P_O and P_S are not given by [Fig. 2](#), and additional information is required to remove the ambiguity, e.g., $P_S(S) = \mathcal{N}(S; \mu, \sigma)$ clarifies that P_S is a Gaussian distribution.

3. Variational message passing

We now build on Forney factor graphs and provide an overview of the method of [Winn and Bishop \(2005\)](#). For more details, see [Champion, Grzeş et al. \(2021\)](#), which provided a complete derivation of the equations presented below from Bayes' theorem.

3.1. Winn and Bishop method

Variational message passing as developed by [Winn and Bishop \(2005\)](#) is an approach for inference based upon the mean-field approximation, which assumes that the posterior fully factorizes, i.e.

$$Q(X) = \prod_i Q(X_i), \quad (6)$$

where X is the set of all hidden variables of the model and X_i represents the i th hidden variable. In this section, we focus on the intuition behind the method, starting with the update equation of an arbitrary hidden state x_k :

$$\ln Q_k^*(x_k) = \langle \ln P(x_k | \text{pa}_k) \rangle_{\sim Q_k} + \sum_{c_j \in \text{ch}_k} \langle \ln P(c_j | x_k, \text{cp}_{kj}) \rangle_{\sim Q_k} + C \quad (7)$$

where C is a normalizing constant, and $\langle \cdot \rangle_{\sim Q_k}$ is the expectation over all factors but $Q_k(x_k)$. (7) tells us that the optimal posterior of any hidden states x_k only depends on its Markov blanket, i.e., x_k 's parents pa_k , children ch_k and co-parents cp_{kj} . To make (7) more specific, we assume that each random variable of the model is conjugate to its parents (i.e., the posterior has the same functional form as the prior) and is distributed according to a distribution in the exponential family, i.e.,

$$\ln P(x_k | \text{pa}_k) = \mu_k(\text{pa}_k) \cdot u_k(x_k) + h_k(x_k) + z_k(\text{pa}_k) \quad (8)$$

where $\mu_k(\text{pa}_k)$, $u_k(x_k)$, $h_k(x_k)$ and $z_k(\text{pa}_k)$ are the parameters, the sufficient statistics, the underlying measure and the log partition, respectively. Under those two assumptions, (7) can be re-written as:

$$Q_k^*(x_k) = \exp \left\{ \mu_k^* \cdot u_k(x_k) + h_k(x_k) + \text{Const} \right\} \quad (9)$$

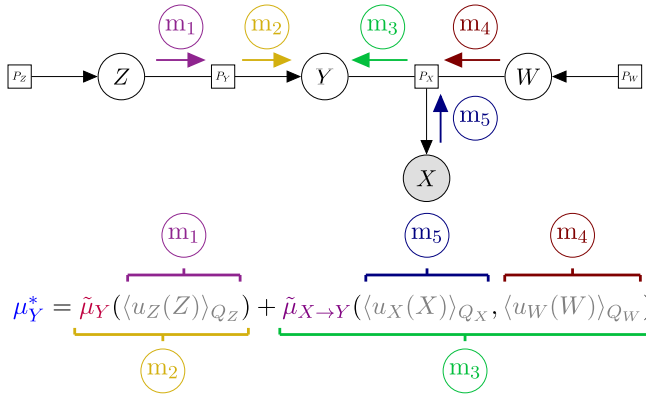


Fig. 3. This figure illustrates the computation of the optimal posterior parameters as a message passing procedure, which requires the transmission of messages from the parent (m_2) and child (m_3) factors. Additionally, the message from the child factor (m_3) requires the computation of messages from the co-parent (m_4) and child (m_5) variables. Also, the message from the parent factor (m_2) requires the computation of a message (m_1) from the parent variable.

with the parameters of the posterior defined as:

$$\mu_k^* = \tilde{\mu}_k(\{\langle u_i(i) \rangle_{Q_i}\}_{i \in \text{pa}_k}) + \sum_{c_j \in \text{ch}_k} \tilde{\mu}_{j \rightarrow k}(\langle u_j(c_j) \rangle_{Q_j}, \{\langle u_l(l) \rangle_{Q_l}\}_{l \in \text{cp}_{kj}}) \quad (10)$$

where $\tilde{\mu}_k$ is a re-parameterization of $\mu_k(\text{pa}_k)$ in terms of the expectation of the sufficient statistics of the parents of x_k , and similarly $\tilde{\mu}_{j \rightarrow k}$ is a re-parameterization of $\mu_{j \rightarrow k}$. Importantly, $u_k(x_k)$ and $h_k(x_k)$ in the optimal posterior (9) are the same as in the prior (8), and only the parameters have changed according to (10).

To understand the intuition behind (10), let us suppose that we are given the Forney factor graph illustrated in Fig. 3 and we wish to compute the posterior of Y . Then, the only parent of Y is Z , the only child of Y is X and the only co-parent of Y with respect to X is W . Therefore, applying (10) to our example leads to the equation presented in Fig. 3 whose components can be interpreted as messages. Indeed, each variable (i.e., X , Z and W) sends the expectation of its sufficient statistics (i.e., a message) to the square node in the direction of Y (i.e., either P_X or P_Y). Those messages are then combined using a function (i.e., either $\tilde{\mu}_Y$ or $\tilde{\mu}_{X \rightarrow Y}$) whose output (i.e., another set of messages) are summed to obtain the optimal parameters μ_Y^* . The computation of the optimal parameters (10) can then be understood as a message passing procedure. Also, we provide in Appendix C a concrete instance of the approach presented above.

4. Active inference

This section provides a quick overview of the active inference framework, and Appendix H presents a description of the exponential complexity class that it exhibits. The reader is referred to Appendix F for any notations that might not be explained here. For a more detailed treatment of the active inference framework, we refer the reader to Champion, Grzesiński et al. (2021), Da Costa, Parr et al. (2020), Smith, Friston, and Whyte (2021).

4.1. Generative model

As illustrated in Fig. 4, the classic generative model represents the world as a sequence of hidden states generating observations through the matrix A . The prior over the initial states is defined by the vector D and the transition between time steps is encoded by a 3-tensor B , i.e., one matrix per action. Importantly, the random

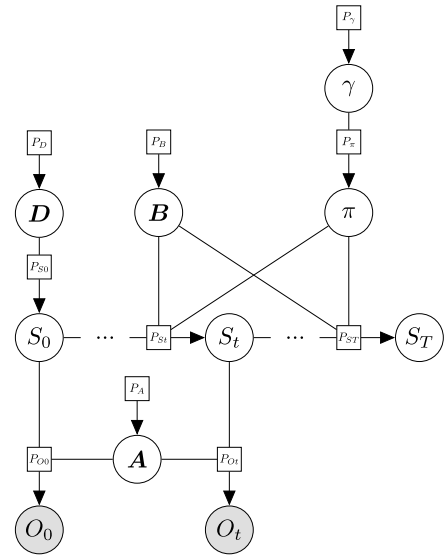


Fig. 4. This figure illustrates the Forney factor graph of the entire generative model presented by Friston et al. (2016). The probability of the initial states is defined by the vector D , and the matrix A defines the probability of the observations given the hidden states. The B matrices define the transition between any successive pair of hidden states. This transition depends on the action performed by the agent, i.e., on the policy π . Furthermore, the prior over the policies has been chosen such that policies minimizing expected free energy are more probable. Finally, the precision parameter γ (which modulates the confidence over which policies to pursue) is distributed according to a gamma distribution.

variable π represents all possible policies up to a given time horizon T and each policy is defined as a sequence of actions, i.e., $\{U_t, \dots, U_{T-1}\}$ where $U_t \in \{1, \dots, |U|\} \forall t \in \{t, \dots, T-1\}$. The prior over the policies is then set such that policies with high probability minimize the EFE, which is defined as follows (Parr & Friston, 2019):

$$\mathbf{G}(\pi) \approx \sum_{\tau=t+1}^T \left[\underbrace{D_{\text{KL}}[Q(O_\tau | \pi) \parallel P(O_\tau)]}_{\text{risk}} + \underbrace{\mathbb{E}_{Q(S_\tau | \pi)}[H[P(O_\tau | S_\tau)]]}_{\text{ambiguity}} \right] \quad (11)$$

where $H[\cdot]$ is the Shannon entropy, the expected outcomes $Q(O_\tau | \pi) \triangleq \sum_{S_\tau} P(O_\tau | S_\tau) Q(S_\tau | \pi)$, \mathbf{G} is a vector containing as many elements as the number of policies, and the i th element of \mathbf{G} represents the cost of the i th policy. The prior preferences over observations $P(O_\tau)$ represent the (categorical) distribution that the agent wants its observations to be sampled from and is traditionally encoded by the vector C . Note that this generalizes the concept of reward from reinforcement learning. Indeed, maximizing reward can be reformulated as sampling observations from a Dirac delta distribution over reward maximizing states (Da Costa, Sajid, Parr, Friston, & Ryan, 2020).

Lastly, the precision parameter γ has been associated to neuromodulators such as dopamine (FitzGerald et al., 2015; Friston, Schwartenbeck, FitzGerald, Moutoussis, Behrens, & Dolan, 2013) and can be understood as modulating the confidence over the information afforded by the expected free energy—e.g., smaller values of γ lead to more stochastic decision-making. Finally, the framework allows A , B and D to be learned by introducing Dirichlet distributions over the columns of these tensors such that the posterior parameters of A , B and D can be reused in a new trial, as parameters of the prior, giving an empirical prior. Finally, the classic generative model is defined as follows:

$$P(O_{0:t}, S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma) = P(\pi|\gamma)P(\gamma)P(\mathbf{A})P(\mathbf{B})P(S_0|\mathbf{D})P(\mathbf{D}) \prod_{\tau=0}^t P(O_\tau|S_\tau, \mathbf{A}) \prod_{\tau=1}^T P(S_\tau|S_{\tau-1}, \pi_{\tau-1}, \mathbf{B}) \quad (12)$$

$$\begin{aligned} P(\pi|\gamma) &= \sigma(-\gamma\mathbf{G}) & P(\gamma) &= \Gamma(1, \boldsymbol{\beta}) \\ P(\mathbf{A}) &= \text{Dir}(\mathbf{a}) & P(\mathbf{B}) &= \text{Dir}(\mathbf{b}) \\ P(S_0|\mathbf{D}) &= \text{Cat}(\mathbf{D}) & P(\mathbf{D}) &= \text{Dir}(\mathbf{d}) \\ P(O_\tau|S_\tau, \mathbf{A}) &= \text{Cat}(\mathbf{A}) & P(S_\tau|S_{\tau-1}, \pi_{\tau-1}, \mathbf{B}) &= \text{Cat}(\mathbf{B}), \end{aligned}$$

where \mathbf{G} is a vector of size $|\pi|$ whose i th element corresponds to the expected free energy of the i th policy, $\sigma(\cdot)$ is the softmax function, $\Gamma(\cdot)$, $\text{Cat}(\cdot)$ and $\text{Dir}(\cdot)$ stand for a gamma, categorical and Dirichlet distribution, respectively, $\pi_{\tau-1} \in \{1, \dots, |U|\}$ is the action prescribed by policy π at time $\tau - 1$, $O_{0:t}$ is the set of (random variables representing) observations between time step 0 and t , and $S_{0:T}$ is the set of (random variables representing) hidden states between time step 0 and T .

4.2. Variational distribution

The most widely used variational distribution (Da Costa, Parr et al., 2020; Friston et al., 2016) is not fully factorized, i.e., the posterior models the influence of the policy on the hidden states, leading to the following factorization:

$$Q(S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma) = Q(\pi)Q(\mathbf{A})Q(\mathbf{B})Q(\mathbf{D})Q(\gamma) \prod_{\tau=0}^T Q(S_\tau|\pi) \quad (13)$$

$$\begin{aligned} Q(S_\tau|\pi) &= \text{Cat}(\hat{\mathbf{D}}_\tau) & Q(\pi) &= \text{Cat}(\hat{\pi}) \\ Q(\gamma) &= \Gamma(1, \hat{\boldsymbol{\beta}}) & Q(\mathbf{D}) &= \text{Dir}(\hat{\mathbf{d}}) \\ Q(\mathbf{A}) &= \text{Dir}(\hat{\mathbf{a}}) & Q(\mathbf{B}) &= \text{Dir}(\hat{\mathbf{b}}) \end{aligned}$$

where all variables with a hat correspond to posterior parameters. Notice that the distributions over \mathbf{A} , \mathbf{B} and \mathbf{D} remain Dirichlet distributions, and the distributions over γ and S_τ remain a gamma and a categorical distribution, respectively. Only the distribution over π changes from a Boltzmann to a categorical distribution but both are discrete distributions.

Remark 1. By definition the generative model $P(O_{0:t}, S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma)$ is a joint probability distribution over both the observed ($O_{0:t}$) and latent ($S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma$) variables. However, the goal of the variational distribution $Q(S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma)$ is to approximate the true posterior $P(S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma|O_{0:t})$, which is a distribution over the latent variables only. Thus, the approximate posterior $Q(S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma)$ is also a distribution over the latent variables only, and does not contain the observed variables.

4.3. Variational free energy

By definition, the variational free energy (VFE) is the Kullback–Leibler divergence between the variational distribution and the generative model, i.e.

$$F = \mathbb{E}_Q[\ln Q(S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma) - \ln P(O_{0:t}, S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma)] \\ = D_{\text{KL}}[Q(x)||P(x|o)] - \ln P(o) \quad (14)$$

$$= \underbrace{D_{\text{KL}}[Q(x)||P(x)]}_{\text{complexity}} - \underbrace{\mathbb{E}_{Q(x)}[\ln P(o|x)]}_{\text{accuracy}} \quad (15)$$

where $x = \{S_{0:T}, \pi, \mathbf{A}, \mathbf{B}, \mathbf{D}, \gamma\}$ refers to the model's hidden variables, and $o = \{O_{0:t}\}$ refers to the sequence of observations made by the agent. (14) shows that minimizing free energy involves moving the variational distribution $Q(x)$ closer to the true posterior $P(x|o)$ in the sense of KL divergence, and that the variational free energy is an upper bound on the negative log evidence. (15) shows the trade-off between complexity and accuracy, where the complexity penalizes the divergence of the posterior $Q(x)$ from the prior $P(x)$ and the accuracy scores how likely the observations are given the generative model and current belief of the hidden states.

To fit the variational distribution as closely as possible to the true posterior, the VFE is minimized w.r.t each variational factor, e.g., $Q(\mathbf{D})$ and $Q(\mathbf{A})$. The minimization process can be solved by iterating the update equations of each factor until convergence of the VFE. More details and intuition about those updates are given by Champion, Grześ et al. (2021).

4.4. Action selection

In active inference, the simplest strategy to select actions is to compute the evidence for all policies under consideration and then choose the most likely action according to these policies. Mathematically, this amounts to a Bayesian model average by executing the action with the highest posterior evidence:

$$u_t^* = \arg \max_u \sum_{m=1}^{|\pi|} [u = \pi_t^m] Q(\pi = m) \quad (16)$$

where $|\pi|$ is the number of policies, π_t^m is the action predicted at the current time step by the m th policy, and $[u = \pi_t^m]$ is an indicator function that equals one if $u = \pi_t^m$ and zero otherwise.

5. Monte Carlo Tree search

By now, the reader should be familiar with the framework of active inference and how variational message passing combined with the Forney factor graph formalism can be used to compute posterior beliefs. We now turn to the last piece of background required to present the method proposed in this paper: Monte Carlo tree search (MCTS), which is based on the multi-armed bandit literature (cf. Appendix G for details).

5.1. A four step process

Monte Carlo tree search has been widely used in the reinforcement learning literature as it enables agents to plan efficiently when the evaluation of every possible action sequence is computationally prohibitive (Browne et al., 2012; Fountas et al., 2020; Schrittwieser et al., 2019; Silver et al., 2016). This algorithm essentially builds a tree in which each node corresponds to a future state and each edge represents the action that led to that state. Initially, the tree is only composed of a root node corresponding to the current state. From here, MCTS is a four step process. First, a node is selected according to a criterion such as the upper confidence bound for trees (UCT):

$$UCT_j = \bar{X}_j + 2C_p \sqrt{\frac{2 \ln n}{n_j}}, \quad (17)$$

where n is the number of times the current (parent) node has been explored, n_j stands for the number of times the j th child node has been explored, $C_p > 0$ is the exploration constant and \bar{X}_j is the average reward received by the j th child. Note, if the rewards are in $[0, 1]$, then $C_p = \frac{1}{\sqrt{2}}$ is known to satisfy the

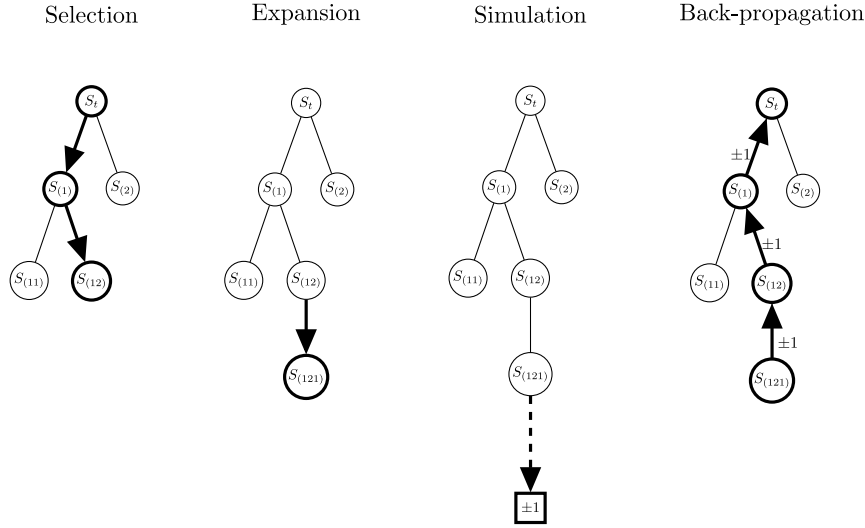


Fig. 5. This figure illustrates the MCTS algorithm as a four step process. First, we start at the node representing the current state S_t and select a node based on the UCT criterion until a leaf node is reached. Second, the tree is expanded to a new node by taking a virtual action from the selected node. Third, the value of this action is estimated by simulating the expected reward following that action. In the simplest version of MCTS, simulations are run until a terminal state is reached, e.g., until the game ends in Go or Chess. Fourth, the expected value is back-propagated to the new node and all of its ancestor nodes. The multi-indices in curly brackets denote action sequences taken from the root node, indicating the current state of the environment.

Hoeffding inequality (Browne et al., 2012) and the UCT criterion reduces to:

$$UCT_j = \bar{X}_j + 2\sqrt{\frac{\ln n}{n_j}}. \quad (18)$$

Importantly, the UCT aims to explore highly rewarding paths (exploitation in first term), while also visiting rarely explored regions (exploration in second term).

As shown in Fig. 5, this criterion is first used at the root level leading to the selection of a node from the root's children. Then, it is used at the level of the root's children, and so on until a leaf node is reached. As explained by Kocsis and Szepesvári (2006), UCT is a direct application of the UCB1 criterion to trees, where at each level, the allocation strategy must pick a node that is expected to lead to the highest reward, and “picking the i th node” can be seen as the i th action of a multi-armed bandit problem. Once a leaf node has been selected, an expansion step is performed by sampling an action from a distribution and adding the node corresponding to this action as a child of the leaf node, i.e., the leaf node is expanded.

The third step consists of performing virtual rollouts into the future to estimate the average future reward obtained from the state corresponding to the newly expanded node. Finally, during the back-propagation step, the average reward obtained from the newly expanded state is used to re-evaluate the average quality of all its ancestors, and the visit counts of all nodes (in the branch explored) are increased. Iterating this four-step process until the time budget has been spent gives a fairly good estimate of the best action to perform next. Fig. 5 summarizes the MCTS procedure. In the next section, we present our approach and show how MCTS can be fused to active inference by performing a dynamical expansion of the generative model.

6. Branching Time Active Inference (BTAI)

In this section, we present a novel active inference agent that frames planning using a tree as a form of Bayesian model extension. Using terminology from concurrency theory (Bowman, 2005), we call our new formalism, *Branching Time Active Inference* (BTAI). In this domain, models of systems based upon sequences

of actions (the format of policies) are described as *linear time*, while models based upon tree and even graph structures are called *branching time* (Bowman, 2005; van Glabbeek, 1990, 1993). Importantly, we do not consider the generative model and the tree as two different objects. Instead, we merge those two objects together into a generative model that can be dynamically expanded.

Fig. 6 illustrates an example of such a model, where for the sake of simplicity, we assume that the matrices \mathbf{A} , \mathbf{B} and \mathbf{D} are given to the agent. Furthermore, the random variable representing the policies has been replaced by random variables representing actions and the precision parameter γ has been removed, which is a common design choice (Fountas et al., 2020). Additionally, we follow Parr and Friston (2019) by viewing future observations as latent random variables. Finally, note that the transition between two consecutive hidden states in the future (S_{Nlast} and S_l where l is a multi-index) will only depend on the matrix $\bar{\mathbf{B}}_l = \bar{\mathbf{B}}(\cdot, \cdot, I_{last})$, i.e., the matrix corresponding to action I_{last} that led to the transition from S_{Nlast} to S_l . The reader is referred to Table 1 for the definition of $\bar{\mathbf{B}}$ and more details about multi-indices can be found in Appendix F.

6.1. Prior, posterior and target distributions

Since the generative model is fairly different from the standard model, we state here its formal definition:

$$P(O_{0:t}, S_{0:t}, U_{0:t-1}, O_{\mathbb{I}_t}, S_{\mathbb{I}_t}, \mathbf{A}, \mathbf{B}, \mathbf{D}, \Theta_{0:t-1}) = \\ P(S_0 | \mathbf{D}) P(\mathbf{A}) P(\mathbf{B}) P(\mathbf{D}) \prod_{\tau=0}^t P(O_\tau | S_\tau, \mathbf{A}) \prod_{\tau=0}^{t-1} P(U_\tau | \Theta_\tau) P(\Theta_\tau) \\ \times \prod_{\tau=1}^t P(S_\tau | S_{\tau-1}, U_{\tau-1}, \mathbf{B}) \prod_{I \in \mathbb{I}_t} P(O_I | S_I) P(S_I | S_{Nlast}) \quad (19)$$

where \mathbb{I}_t is the set of all non-empty multi-indices already expanded by the tree search from the current state S_t , the second product (τ from 0 to $t-1$) models the uncertainty over action, reflecting the focus on actions rather than policies, and S_{Nlast} is the parent state of S_l . Intuitively, the product over all $I \in \mathbb{I}_t$ models the future, while the rest of the above equation models the

Table 1
Update equations notation. Note that Appendix D provides a proof for \mathbf{D} .

Notation	Meaning
$\langle f(X) \rangle_{P_X} \triangleq \mathbb{E}_{P_X} [f(X)]$	The expectation of $f(X)$ over P_X
$\psi(\cdot)$	The digamma function
$\hat{\Theta}_\tau(i) = \langle \ln \Theta_\tau(i) \rangle_{Q_{\Theta_\tau}} = \psi(\hat{\theta}_\tau(i)) - \psi(\sum_k \hat{\theta}_\tau(k))$	The expected logarithm of Θ_τ
$\hat{\mathbf{D}}(i) = \langle \ln \mathbf{D}(i) \rangle_{Q_D} = \psi(\hat{\mathbf{d}}(i)) - \psi(\sum_k \hat{\mathbf{d}}(k))$	The expected logarithm of \mathbf{D}
$\hat{\mathbf{A}}(i, j) = \langle \ln \mathbf{A}(i, j) \rangle_{Q_A} = \psi(\hat{\mathbf{a}}(i, j)) - \psi(\sum_k \hat{\mathbf{a}}(k, j))$	The expected logarithm of \mathbf{A}
$\hat{\mathbf{B}}(i, j, u) = \langle \ln \mathbf{B}(i, j, u) \rangle_{Q_B} = \psi(\hat{\mathbf{b}}(i, j, u)) - \psi(\sum_k \hat{\mathbf{b}}(k, j, u))$	The expected logarithm of \mathbf{B}
$\bar{\mathbf{A}}(i, j) = \langle \mathbf{A}(i, j) \rangle_{Q_A} = \frac{\hat{\mathbf{a}}(i, j)}{\sum_k \hat{\mathbf{a}}(k, j)}$	The expectation of \mathbf{A}
$\bar{\mathbf{B}}(i, j, u) = \langle \mathbf{B}(i, j, u) \rangle_{Q_B} = \frac{\hat{\mathbf{b}}(i, j, u)}{\sum_k \hat{\mathbf{b}}(k, j, u)}$	The expectation of \mathbf{B}

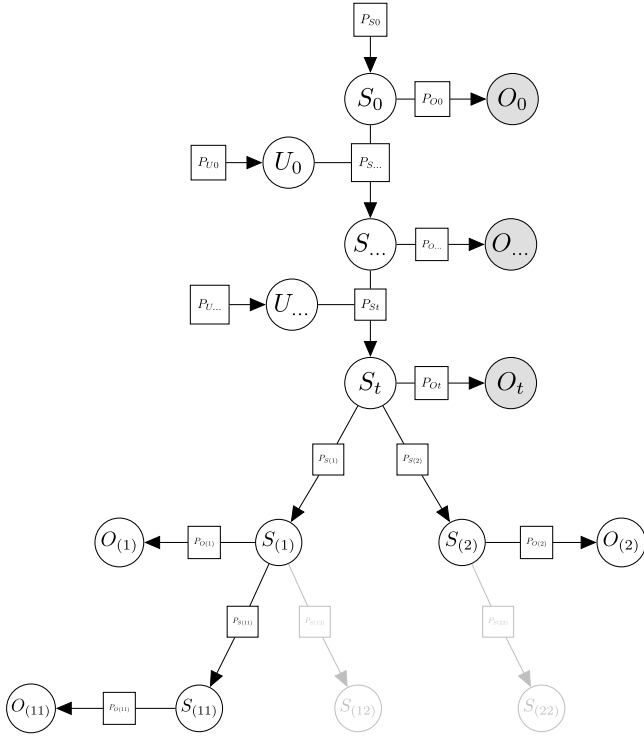


Fig. 6. This figure illustrates the new expandable generative model allowing planning under active inference. The future is now a tree like generative model whose branches correspond to the policies considered by the agent. As we will see, these branches can be dynamically expanded during planning. Here, the nodes in light grey represent possible expansions of the current generative model. For the sake of clarity, the random tensors \mathbf{A} , \mathbf{B} , Θ_τ and \mathbf{D} are not illustrated, i.e., Dirichlet priors over those random tensors are not shown.

past and present. Additionally, we need to define the individual factors:

$$\begin{aligned}
 P(S_0|\mathbf{D}) &= \text{Cat}(\mathbf{D}) & P(U_\tau|\Theta_\tau) &= \text{Cat}(\Theta_\tau) \\
 P(O_\tau|S_\tau, \mathbf{A}) &= \text{Cat}(\mathbf{A}) & P(O_I|S_I) &= \text{Cat}(\bar{\mathbf{A}}) \\
 P(S_\tau|S_{\tau-1}, U_{\tau-1}, \mathbf{B}) &= \text{Cat}(\mathbf{B}) & P(S_I|S_{I_{\text{last}}}) &= \text{Cat}(\bar{\mathbf{B}}_I) \\
 P(\mathbf{D}) &= \text{Dir}(\mathbf{d}) & P(\Theta_\tau) &= \text{Dir}(\theta_\tau) \\
 P(\mathbf{A}) &= \text{Dir}(\mathbf{a}) & P(\mathbf{B}) &= \text{Dir}(\mathbf{b})
 \end{aligned}$$

where $\bar{\mathbf{A}}$ and $\bar{\mathbf{B}}$ are defined in Table 1, $\bar{\mathbf{B}}_I = \bar{\mathbf{B}}(\cdot, \cdot, I_{\text{last}})$ is the matrix corresponding to I_{last} and I_{last} is the last index of the multi-index I , i.e., the last action that led to S_I . Importantly, $\bar{\mathbf{A}}$ and $\bar{\mathbf{B}}$ should not be confused with $\hat{\mathbf{A}}$ and $\hat{\mathbf{B}}$, $\bar{\mathbf{A}}$ is the expectation of \mathbf{A} w.r.t. $Q(\mathbf{A})$, while $\hat{\mathbf{A}}$ is the expectation of the logarithm of \mathbf{A} w.r.t. $Q(\mathbf{A})$.

We now turn to the definition of the variational posterior. Under the mean-field approximation:

$$Q(S_{0:t}, U_{0:t-1}, O_{I_t}, S_{I_t}, \mathbf{A}, \mathbf{B}, \mathbf{D}, \Theta_{0:t-1}) = Q(\mathbf{A})Q(\mathbf{B})Q(\mathbf{D}) \prod_{\tau=0}^{t-1} Q(U_\tau)Q(\Theta_\tau) \prod_{\tau=0}^t Q(S_\tau) \prod_{I \in I_t} Q(O_I)Q(S_I) \quad (20)$$

where the individual factors are defined as:

$$\begin{aligned}
 Q(S_\tau) &= \text{Cat}(\hat{\mathbf{D}}_\tau) & Q(U_\tau) &= \text{Cat}(\hat{\Theta}_\tau) \\
 Q(O_I) &= \text{Cat}(\hat{\mathbf{E}}_I) & Q(S_I) &= \text{Cat}(\hat{\mathbf{D}}_I) \\
 Q(\mathbf{D}) &= \text{Dir}(\hat{\mathbf{d}}) & Q(\Theta_\tau) &= \text{Dir}(\hat{\theta}_\tau) \\
 Q(\mathbf{A}) &= \text{Dir}(\hat{\mathbf{a}}) & Q(\mathbf{B}) &= \text{Dir}(\hat{\mathbf{b}}),
 \end{aligned}$$

where $\hat{\mathbf{D}}_\tau$, $\hat{\Theta}_\tau$, $\hat{\mathbf{E}}_I$, $\hat{\mathbf{D}}_I$, $\hat{\mathbf{a}}$, $\hat{\theta}_\tau$, $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ are the parameters of the factors $Q(S_\tau)$, $Q(U_\tau)$, $Q(O_I)$, $Q(S_I)$, $Q(\mathbf{D})$, $Q(\Theta_\tau)$, $Q(\mathbf{A})$ and $Q(\mathbf{B})$, respectively. Importantly, O_I appears in the variational distribution because observations in the future are treated as hidden variables.

Finally, we follow Millidge, Tschantz, and Buckley (2021) in assuming that the agent aims to minimize the KL divergence between the approximate posterior depicting the state of the environment and a target (desired) distribution. Therefore, our framework allows for the specification of prior preferences over both future hidden states and future observations:

$$V(O_{I_t}, S_{I_t}) = \prod_{I \in I_t} V(O_I)V(S_I) \quad (21)$$

where the individual factors are defined as:

$$V(O_I) = \text{Cat}(\mathbf{C}_O), \quad V(S_I) = \text{Cat}(\mathbf{C}_S). \quad (22)$$

Importantly, by specifying the value of future observations and states, \mathbf{C}_O and \mathbf{C}_S play a similar role to the vector \mathbf{C} in active inference, i.e., they specify which observations and hidden states are rewarding.

To sum up, this framework is defined using three distributions: the prior defines the agent's beliefs before sampling any observation; the posterior is an updated version of the prior which takes into account past observations made by the agent; finally, the target distribution encodes the agent's prior preferences in terms of future observations and hidden states.

6.2. Bayesian belief updates

In this section, we focus on the set of update equations used to perform approximate Bayesian inference. These update equations rely on variational message passing as presented in Section 3, see Champion, Grzesiński et al. (2021) as well as Winn and Bishop (2005) for details. A key strength of the message passing approach is the capacity to derive and implement these updates within an automatic and modular toolbox (Cox et al., 2019; van de Laar & de Vries, 2019), which in a way similar to automatic

differentiation alleviates the final user from the burden of manually deriving complex update equations for each new generative model. To simplify our notation, we use two operators \otimes and \odot that we call generalized outer and inner product, respectively. The generalized outer product creates an N dimensional tensor from N vectors, while the generalized inner product performs a weighted average over one dimension of an N dimensional array, cf [Appendix A](#) for details. Using these notations, the first set of update equations are given by:

$$Q^*(\mathbf{D}) = \text{Dir}(\hat{\mathbf{d}}) \quad \text{where} \quad \hat{\mathbf{d}} = \mathbf{d} + \hat{\mathbf{D}}_0 \quad (23)$$

$$Q^*(\mathbf{A}) = \text{Dir}(\hat{\mathbf{a}}) \quad \text{where} \quad \hat{\mathbf{a}} = \mathbf{a} + \sum_{\tau=0}^t \otimes [\hat{\mathbf{D}}_{\tau}, \mathbf{o}_{\tau}] \quad (24)$$

$$Q^*(\mathbf{B}) = \text{Dir}(\hat{\mathbf{b}}) \quad \text{where} \quad \hat{\mathbf{b}} = \mathbf{b} + \sum_{\tau=1}^t \otimes [\hat{\mathbf{D}}_{\tau-1}, \hat{\Theta}_{\tau-1}, \hat{\mathbf{D}}_{\tau}] \quad (25)$$

$$Q^*(\Theta_{\tau}) = \text{Dir}(\hat{\theta}_{\tau}) \quad \text{where} \quad \hat{\theta}_{\tau} = \theta_{\tau} + \hat{\Theta}_{\tau} \quad (26)$$

where \mathbf{o}_{τ} is the observation made at time τ . Furthermore, this first set of equations count (probabilistically) the number of times, an initial hidden state has been observed, an action has been performed, a state has generated a particular observation or an action has led to the transition between two consecutive hidden states. For example, the posterior parameters $\hat{\mathbf{a}}$ are computed by adding $\sum_{\tau=0}^t \otimes [\hat{\mathbf{D}}_{\tau}, \mathbf{o}_{\tau}]$ (i.e., the number of times a state-observation pair has been observed during this trial) to the prior parameters \mathbf{a} (i.e., the number of times this same pair has been observed during previous trials). The equations for belief updates are given by:

$$Q^*(O_I) = \sigma(\hat{\mathbf{A}} \odot \hat{\mathbf{D}}_I) \quad (27)$$

$$Q^*(S_I) = \sigma(\hat{\mathbf{A}} \odot \hat{\mathbf{E}}_I + \hat{\mathbf{B}}_I \odot \hat{\mathbf{D}}_{N_{\text{last}}} + \sum_{J \in \text{ch}_I} \hat{\mathbf{B}}_J \odot \hat{\mathbf{D}}_J) \quad (28)$$

$$Q^*(U_{\tau}) = \sigma(\hat{\Theta} + \hat{\mathbf{B}} \odot [\hat{\mathbf{D}}_{\tau}, \hat{\mathbf{D}}_{\tau+1}]) \quad (29)$$

$$Q^*(S_{\tau}) = \sigma \left([\tau = 0] \hat{\mathbf{D}}_{\tau} + [\tau \neq 0] \hat{\mathbf{B}} \odot [\hat{\mathbf{D}}_{\tau-1}, \hat{\Theta}_{\tau-1}] + \hat{\mathbf{A}} \odot \mathbf{o}_{\tau} + [\tau = t] \sum_{J \in \text{ch}_t} \hat{\mathbf{B}}_J \odot \hat{\mathbf{D}}_J + [\tau \neq t] \hat{\mathbf{B}} \odot [\hat{\mathbf{D}}_{\tau+1}, \hat{\Theta}_{\tau}] \right) \quad (30)$$

where $\sigma(\cdot)$ is the softmax function, ch_I are the children (states) of the current states S_I , ch_I are the children (states) of the states S_I , [predicate] is an indicator function returning one if the predicate is true and zero otherwise, and the definition of $\hat{\mathbf{A}}$, $\hat{\mathbf{B}}$, $\hat{\mathbf{D}}$ and $\hat{\Theta}_{\tau}$ is given in [Table 1](#). Note that thanks to the operators \otimes and \odot , the perception (i.e., state-estimation) equations can be intuitively understood as a sum of messages, where each message from a factor to a variable is the average over all dimensions except the dimension of the variable, e.g., the message $(\hat{\mathbf{A}} \odot \mathbf{o}_0)$ from P_{00} to S_0 is the vector obtained by weighting the rows of $\hat{\mathbf{A}}$ by the elements of \mathbf{o}_0 . Importantly, the above update equations are almost identical to the ones used in standard active inference, and thus can be implemented efficiently. Indeed, most of the computation required is about addition of matrices $O(n^2)$ and multiplication of matrices $O(n^3)$, or their higher dimensional counterparts.

6.3. Planning as structure learning

In this section, we frame planning as a form of structure learning where the structure of the generative model is modified dynamically. This method is greatly inspired by the Monte Carlo tree search literature, cf., [Section 5](#) for details.

6.3.1. Selection of the node to be expanded

The first step of planning is to select a node to be expanded. The selection process starts at the root node, if the root node still has unexplored children, then one of them is selected. Otherwise, the child node maximizing the *UCT* criterion, where the average reward is replaced by minus the average EFE, is selected, i.e., the selected node maximizes:

$$UCT_J = \underbrace{-\bar{g}_J}_{\text{exploitation}} + \underbrace{C_p \sqrt{\frac{\ln n}{n_J}}}_{\text{exploration}}, \quad (31)$$

where J is a multi-index, n is the number of times the root node has been visited, n_J is the number of times the child corresponding to the multi-index J was selected, and \bar{g}_J is the average cost received when selecting the child S_J . The *UCT* criterion can be understood as a trade-off between exploitation and exploration at the tree level, which is different to the exploitation and exploration dilemma at the model level. This dilemma is handled by the EFE. Also, the notion of cost in the above equation can be defined in many ways and will be the subject of [Section 6.3.3](#). For our purposes, the cost will be equal, or similar, to the expected free energy, which means that the expected free energy drives structure learning. When a root's child is selected, it becomes the new root in the above procedure, which is iterated until a leaf node is reached.

6.3.2. Dynamical expansion of the generative model

Let $S_{N_{\text{last}}}$ denote the leaf node selected for expansion. When $S_{N_{\text{last}}}$ has been selected, the structure of the generative model needs to be modified by expanding all possible actions from that node. For each action, we expand the generative model by adding a future hidden state whose prior distribution is given by

$$P(S_I | S_{N_{\text{last}}}) = \text{Cat}(\bar{\mathbf{B}}_I), \quad (32)$$

where $\bar{\mathbf{B}}_I$ is the matrix corresponding to the last action that led to S_I . Finally, we expand the (future) observation associated with the new hidden state S_I , whose distribution is:

$$P(O_I | S_I) = \text{Cat}(\bar{\mathbf{A}}). \quad (33)$$

To sum up, the expansion step is adding two random variables (S_I and O_I) to the generative model, i.e. the generative model becomes bigger, and I is added to the set of all non-empty multi-indices already expanded by the tree search (\mathbb{I}_I). The prior distributions over those newly added random variables (i.e. S_I and O_I) are defined using the matrices $\bar{\mathbf{B}}_I$ and $\bar{\mathbf{A}}$, which effectively predict the future states and observations. After the expansion step, the posterior distribution over S_I and O_I needs to be computed. At least two kinds of inference strategies can be used. The first—global inference—performs variational message passing over the entire generative model, while the second—local inference—only iterates the update equations of the newly expanded nodes, i.e., S_I and O_I , until convergence to the variational free energy minimum.

6.3.3. Cost evaluation of the expanded nodes

After expanding the model structure, we need to compute the cost of the newly expanded node S_I . As explained in [Section 6.3.1](#), the cost of S_I will influence the probability of expanding S_I during future planning iterations. In active inference, the classic

objective of planning is the expected free energy as defined in Section 4.1, i.e.,

$$g_l^{\text{classic}} \triangleq D_{\text{KL}}[Q(O_l) \parallel V(O_l)] + \mathbb{E}_{Q(S_l)}[H[P(O_l|S_l)]] \quad (34)$$

where g_l^{classic} trades off risk (first summand) and ambiguity (second summand). Alternatively, one could follow Section 5 of [Milledge et al. \(2021\)](#) and define the cost of S_l using the free energy of the expected future:

$$g_l^{\text{feef}} = D_{\text{KL}}[Q(O_l, S_l) \parallel V(O_l, S_l)] \quad (35)$$

where $V(O_l, S_l)$ is the target distribution over states and observations. The target distribution $V(O_l, S_l)$ generalizes the \mathbf{C} matrix in Friston's model by specifying prior preferences over both future observations and future states. Also, this formulation of the cost speaks to the notion of KL divergence minimization proposed by [Hafner et al. \(2020\)](#).

Furthermore, due to the mean-field approximation of the posterior [20](#) and the factorized form of the target distribution [\(21\)](#), the expression of the cost simplifies to

$$g_l^{\text{pcost}} \triangleq D_{\text{KL}}[Q(S_l) \parallel V(S_l)] + D_{\text{KL}}[Q(O_l) \parallel V(O_l)]. \quad (36)$$

Intuitively, the pure cost (g_l^{pcost}) measures how different the predicted future hidden states and observations are from desired states and observations. In future research, it might be interesting to compare the performance and behaviour of g_l^{pcost} , g_l^{feef} and g_l^{classic} empirically and theoretically.

Note that in MCTS, the evaluation of a node's quality is done by performing virtual roll-outs, while g_l^{pcost} , g_l^{feef} and g_l^{classic} are not. If we let g_l be any of those criteria, then we can improve our estimate of the cost, by computing g_l^{average} , i.e., the average cost over N roll-outs of size K . Algorithm 1 presents the pseudo code used to estimate g_l^{average} .

Algorithm 1: Estimation of g_l^{average}

Input: N the number of virtual roll-outs,
 K the maximal length of each roll-out.
 $g_l^{\text{average}} \leftarrow 0$ // Initialize roll-out estimate to zero
repeat N **times**
 $g_l^{\text{rollout}} \leftarrow g_l$ // Initial cost equals cost of S_l
 for $i \leftarrow 1$ **to** K **do**
 sample a random action U_i uniformly from the set of unexplored actions
 perform the expansion of the current node using U_i (Section 6.3.2)
 perform inference on the newly expanded nodes (Section 6.2)
 $g_l^{\text{rollout}} \leftarrow g_l^{\text{rollout}} + g_J$ where J corresponds to last expanded node
 end
 $g_l^{\text{average}} \leftarrow g_l^{\text{average}} + g_l^{\text{rollout}}$
end
 $g_l^{\text{average}} \leftarrow g_l^{\text{average}} / N$

6.3.4. Propagation of the node cost

In this section, we let $\mathbf{G}_L^{\text{aggr}}$ be a variable that contains the total cost of the node S_L , where L could be any multi-index. According to the previous section, we let g_L be any of the following evaluation criteria g_L^{pcost} , g_L^{feef} and g_L^{classic} . Initially, $\mathbf{G}_L^{\text{aggr}}$ equals g_L . Also, we let S_K be the node that was selected for expansion, and let S_l be an arbitrary hidden state expanded from S_K . The cost of the newly expanded node(s) can be propagated either forward or backward. The forward propagation (towards the leaves) leads to the following equation:

$$\mathbf{G}_l^{\text{aggr}} \leftarrow g_l + \mathbf{G}_K^{\text{aggr}}, \quad (37)$$

where here $\mathbf{G}_K^{\text{aggr}}$ is the aggregated cost of the parent of S_l . Importantly, the symbol \leftarrow refers to a programming-like assignment (i.e., an incremental update) performed each time the tree is expanded. The backward propagation (towards the root) leads to:

$$\mathbf{G}_J^{\text{aggr}} \leftarrow \mathbf{G}_J^{\text{aggr}} + g_l \quad \forall J \in \mathbb{A}_l \quad (38)$$

where \mathbb{A}_l corresponds to all ancestors of the newly expanded node S_l . We will see in Section 7 that these strategies respectively relate to active inference and sophisticated inference ([Friston et al., 2021](#)). Finally, since the agent is free to choose any action, we can back-propagate the (locally) minimum cost, i.e.,

$$\mathbf{G}_J^{\text{aggr}} \leftarrow \mathbf{G}_J^{\text{aggr}} + \min_{a \in \{1, \dots, |U|\}} g_{K::a} \quad \forall J \in \mathbb{A}_l, \quad (39)$$

where $K :: a$ is a multi-index obtained from K by adding the action a to the sequence of actions described by K . In all cases, the propagation step updates the counter n_J associated with each ancestor S_J of the newly expanded hidden state S_l ; this counts the number of times the node S_J has been explored (exactly as in MCTS). This counter will be used for action selection, as well as for the computation of the average cost of $S_J - \bar{g}_J$ —that was left undefined by Section 6.3.1. Formally, \bar{g}_J is given by:

$$\bar{g}_J = \frac{1}{n_J} \mathbf{G}_J^{\text{aggr}}. \quad (40)$$

Remark 2. The forward propagation of the cost presented above will only be used for theoretical purpose in Section 7. Practical implementation of BTAI should use the backward schemes.

6.4. Action selection

The planning procedure presented in the previous section ends after a pre-specified amount of time has elapsed or when a sufficiently good policy has been found. When the planning is over, the agent needs to choose an action to act in its environment. In a companion paper ([Champion, Bowman et al., 2021](#)) that presents empirical results of BTAI, the actions are sampled from $\sigma(-\gamma \frac{g}{N})$, where $\sigma(\cdot)$ is a softmax function, γ is a precision parameter, g is a vector whose elements correspond to the cost of the root's children and N is a vector whose elements correspond to the number of visits of the root's children. Importantly, actions with low average cost are more likely to be selected than actions with high average cost.

Alternative approaches to action selection ([Browne et al., 2012](#)) could be studied. For example, one could imagine sampling actions from a categorical distribution with parameter $\sigma(N)$, where N is a vector containing the n_J of all children of the root node. Or, we could select the action corresponding to the root's child with the highest number of explorations n_J . The fact that it has been visited more often means that it has a lower cost overall. If there were a tie between several actions, the action with the lowest cost would be selected. The study of these strategies is left to future research.

6.5. Action-perception cycle with tree search

In active inference, the action-perception cycle realizes an active inference agent in an infinite loop ([van de Laar & de Vries, 2019](#)). Each loop iteration begins with the agent sampling an observation from the environment. The observation is used to perform inference about the states and contingencies of the world, e.g., an impression on the retina might be used to reconstruct a three dimensional scene with a representation of the objects that it contains. Then, planning is performed by inferring the consequences of alternative action sequences. Importantly,

only a subset of all possible action sequences are evaluated, due to the dynamical expansion of the generative model. Finally, the agent selects an action to perform in the environment by sampling a softmax function of minus the average cost weighted by the precision parameter γ , i.e., $\sigma(-\gamma \frac{G}{N})$. Therefore, actions with low average cost are more likely to be selected than actions with high average cost. We summarize our method using pseudo-code in Algorithm 2.

Algorithm 2: Action-perception cycle with tree search

```

while end of trial not reached do
  sample an observation from the environment
  perform inference using the observation (Section 6.2)
  while maximum planning iteration not reached do
    select a node to be expanded (Section 6.3.1)
    perform the expansion of the node (Section 6.3.2)
    perform inference on the newly expanded nodes
    (Section 6.2)
    evaluate the cost of the newly expanded nodes
    (Section 6.3.3)
    propagate the cost of the nodes through the tree,
    either forward or backward (Section 6.3.4)
  end
  select an action to be performed (Section 6.4)
  execute the action in the environment
end

```

7. Connection between BTAL, active inference and sophisticated inference

In this section, we explore the relationship between BTAL, active inference (Acl) and sophisticated inference (SI). We show that BTAL is a class of algorithms that generalizes Acl and is related to SI. To do so, we focus on the “cost” of a policy for each method. In addition, we need to introduce the notion of localized and aggregated cost. The localized cost of a node S_i , denoted G_i^{local} , is the cost of S_i in and of itself, i.e., without any consideration of the cost of past or future states. The aggregated cost of a node S_i , denoted G_i^{aggre} , is the cost of S_i when taking into account either the cost of future states that can be reached from S_i (which is the case in SI) or the cost of the past states that an agent has to go through in order to reach S_i (which is the case in Acl).

7.1. Active inference

The full framework of active inference was described in Section 4. This section focuses on expressing the expected free energy in a recursive form that highlights the relationship between BTAL and Acl. We start by defining the notion of localized and aggregated EFE with Definitions 3 and 4, respectively. Then, we show that in active inference (under some assumptions described below), the aggregated EFE of a policy of size N is given by the aggregated EFE of a policy of size $N - 1$ plus the localized EFE received at time $t + N$.

In active inference, a policy is a sequence of actions $\pi = (U_t, U_{t+1}, \dots, U_{T-1})$, where T is the time horizon of planning, and for convenience, π_N denotes a policy of size N , obtained by selecting the first N actions of the policy π , i.e., $\pi_N = (U_t, U_{t+1}, \dots, U_{t+N-1})$ with $N \leq T - t$. Recall from Section 4, that (in active inference) the expected free energy of a policy is:

$$G(\pi) = \sum_{\tau=t+1}^T G(\pi, \tau) \quad (41)$$

where: $G(\pi, \tau) = D_{KL}[Q(O_\tau|\pi) \parallel P(O_\tau)] + \mathbb{E}_{Q(S_\tau|\pi)}[H[P(O_\tau|S_\tau)]]$. If instead of letting τ range from $t + 1$ to T , we let N range from

1 to $T - t$, then Eq. (41) can be re-written as:

$$G(\pi) = \sum_{N=1}^{T-t} G(\pi, t + N) \quad (42)$$

where:

$$G(\pi, t + N) = D_{KL}[Q(O_{t+N}|\pi) \parallel P(O_{t+N})] + \mathbb{E}_{Q(S_{t+N}|\pi)}[H[P(O_{t+N}|S_{t+N})]]. \quad (43)$$

Additionally, under the assumption that the probability of observations and states are independent of future actions, i.e., that $\forall j \in \mathbb{N}_{>0}, Q(O_{t+i}|\pi_i) \approx Q(O_{t+i}|\pi_{i+j})$ and $\forall j \in \mathbb{N}_{>0}, Q(S_{t+i}|\pi_i) \approx Q(S_{t+i}|\pi_{i+j})$, π can be replaced by π_N in the RHS of the above equation, leading to:

$$G(\pi) = \sum_{N=1}^{T-t} G(\pi_N, t + N) \quad (44)$$

where:

$$G(\pi_N, t + N) = D_{KL}[Q(O_{t+N}|\pi_N) \parallel P(O_{t+N})] + \mathbb{E}_{Q(S_{t+N}|\pi_N)}[H[P(O_{t+N}|S_{t+N})]]. \quad (45)$$

Importantly, the elements of the above summation constitute the localized cost presented in Definition 3.

Definition 3. We define the **localized cost** received at time $t + N$ after selecting policy π_N as:

$$G_{\pi_N}^{local} = G(\pi_N, t + N) = D_{KL}[Q(O_{t+N}|\pi_N) \parallel P(O_{t+N})] + \mathbb{E}_{Q(S_{t+N}|\pi_N)}[H[P(O_{t+N}|S_{t+N})]]. \quad (46)$$

Importantly, the localized cost quantifies the amount of risk and ambiguity received by the agent at time step $t + N$, assuming that it will follow the policy π_N . We now turn to the notion of aggregated cost of a policy of size N . Definition 4 states that the aggregated cost of a policy is defined recursively. Indeed, by definition, a policy of size zero has an aggregated cost of zero, and then, the aggregated cost of a policy π_N (of size N) is equal to the aggregated cost of π_{N-1} (of size $N - 1$) plus the localized cost received at time $t + N$.

Definition 4. We define the **aggregated cost** of a policy π_N of size N as:

$$G_{\pi_N}^{aggre} = \begin{cases} 0 & \text{if } N = 0 \\ G_{\pi_{N-1}}^{aggre} + G_{\pi_N}^{local} & \text{otherwise} \end{cases} \quad (47)$$

Equipped with Definitions 3 and 4, we are now ready to state and prove Theorem 5 using the two Lemmas of Appendix E.

Theorem 5. Under the assumption that the probability of observations and states are independent of future actions, i.e., $\forall j \in \mathbb{N}_{>0}, Q(O_{t+i}|\pi_i) \approx Q(O_{t+i}|\pi_{i+j})$ and $\forall j \in \mathbb{N}_{>0}, Q(S_{t+i}|\pi_i) \approx Q(S_{t+i}|\pi_{i+j})$, the expected free energy can be written as:

$$G(\pi_N) \approx G_{\pi_N}^{aggre} = G_{\pi_{N-1}}^{aggre} + G_{\pi_N}^{local} \quad (48)$$

Proof. This proof is based on two lemmas demonstrated in Appendix E. Note that in active inference the expected free energy is defined as:

$$G(\pi) = \sum_{\tau=t+1}^T G(\pi, \tau). \quad (49)$$

Let N denote the size of the policy π , i.e. $N = T - t$. Note that because π is of size N , then by definition $\pi = \pi_N$, and the above

equation can be re-written as:

$$\mathbf{G}(\pi) = \mathbf{G}(\pi_N) = \sum_{\tau=t+1}^{t+N} \mathbf{G}(\pi_N, \tau). \quad (50)$$

Expanding the summation and using [Definition 3](#):

$$\mathbf{G}(\pi_N) = \sum_{\tau=t+1}^{t+N-1} \mathbf{G}(\pi_N, \tau) + \mathbf{G}(\pi_N, t+N) \quad (51)$$

$$= \sum_{\tau=t+1}^{t+N-1} \mathbf{G}(\pi_N, \tau) + \mathbf{G}_{\pi_N}^{local}. \quad (52)$$

If, instead of letting τ range from $t+1$ to $t+N-1$, we let i range from 1 to $N-1$, then the above equation can be re-written as:

$$\mathbf{G}(\pi_N) = \sum_{i=1}^{N-1} \mathbf{G}(\pi_N, t+i) + \mathbf{G}_{\pi_N}^{local}. \quad (53)$$

Note that $\forall i \in \{1, \dots, N-1\}$, $N > i$, and thus there exists a $k_i \in \mathbb{N}_{>0}$ such that $\forall i \in \{1, \dots, N-1\}$, $N = i + k_i$. Therefore, we replace N by $i + k_i$ in the above summation:

$$\mathbf{G}(\pi_N) = \sum_{i=1}^{N-1} \mathbf{G}(\pi_{i+k_i}, t+i) + \mathbf{G}_{\pi_N}^{local}. \quad (54)$$

[Lemma 11](#) tells us that under the assumption that the probability of observations and states are independent of future actions, $\forall k_i \in \mathbb{N}_{>0}$, $\mathbf{G}(\pi_{i+k_i}, t+i) \approx \mathbf{G}(\pi_i, t+i)$, which allows us to remove the k_i to get:

$$\mathbf{G}(\pi_N) \approx \sum_{i=1}^{N-1} \mathbf{G}(\pi_i, t+i) + \mathbf{G}_{\pi_N}^{local}. \quad (55)$$

Finally, [Lemma 12](#) states that $\sum_{i=1}^{N-1} \mathbf{G}(\pi_i, t+i) = \mathbf{G}_{\pi_{N-1}}^{aggre}$, and thus:

$$\mathbf{G}(\pi_N) \approx \mathbf{G}_{\pi_{N-1}}^{aggre} + \mathbf{G}_{\pi_N}^{local} \triangleq \mathbf{G}_{\pi_N}^{aggre}. \quad (56)$$

The above equation will be used in [Section 7.4](#) to show that BTAI generalizes active inference. ■

7.2. Sophisticated inference

Sophisticated inference ([Friston et al., 2021](#)) is a new type of active inference that defines the EFE recursively from the time horizon backward. Intuitively, the agent does not simply ask “what would happen if I did that”, but instead wonders “what would I believe about what would happen if I did that”. In other words, the agent is exhibiting a form of sophistication, which refers to the fact of having beliefs about one’s own or another’s beliefs. [Friston et al. \(2021\)](#) also replaced variational message passing by an alternative inference scheme called Bayesian Filtering ([Fox, Hightower, Liao, Schulz, & Borriello, 2003](#)). While the change of inference method is of little relevance to us here, the recursive definition of the EFE is at the core of this section. As explained in [Section 4.3](#) of [Da Costa, Sajid et al. \(2020\)](#), the (recursive) EFE of a Markov decision process is given by:

$$G(U_{T-1}, S_{T-1}) = D_{KL} [Q(S_T|U_{T-1}, S_{T-1})||V(S_T)] \quad (57)$$

$$G(U_\tau, S_\tau) = D_{KL} [Q(S_{\tau+1}|U_\tau, S_\tau)||V(S_{\tau+1})] + \mathbb{E}_{Q(U_{\tau+1}, S_{\tau+1}|U_\tau, S_\tau)}[G(U_{\tau+1}, S_{\tau+1})] \quad (58)$$

where U_τ and S_τ are the action and state at time τ , and $V(S_\tau)$ is the target (i.e., desired) distribution over states at time τ . Using our terminology of localized and aggregated cost, this can be

rewritten as:

$$\underbrace{G(U_{T-1}, S_{T-1})}_{\mathbf{G}_{\pi_{T-1}, S_{T-1}}^{aggre}} = \underbrace{D_{KL} [Q(S_T|U_{T-1}, S_{T-1})||V(S_T)]}_{\mathbf{G}_{\pi_{T-1}, S_{T-1}}^{local}} \quad (59)$$

$$\underbrace{G(U_\tau, S_\tau)}_{\mathbf{G}_{\pi_\tau, S_\tau}^{aggre}} = \underbrace{D_{KL} [Q(S_{\tau+1}|U_\tau, S_\tau)||V(S_{\tau+1})]}_{\mathbf{G}_{\pi_\tau, S_\tau}^{local}} + \mathbb{E}_{Q(U_{\tau+1}, S_{\tau+1}|U_\tau, S_\tau)}[\underbrace{G(U_{\tau+1}, S_{\tau+1})}_{\mathbf{G}_{\pi_{\tau+1}, S_{\tau+1}}^{aggre}}] \quad (60)$$

Put simply, the aggregated cost of taking action U_τ in state S_τ can be computed by summing the localized cost at time step τ and the expected aggregated cost at time step $\tau+1$, i.e.,

$$\mathbf{G}_{\pi_\tau, S_\tau}^{aggre}(U_\tau, S_\tau) = \mathbf{G}_{\pi_\tau, S_\tau}^{local}(U_\tau, S_\tau) + \mathbb{E}_{Q(U_{\tau+1}, S_{\tau+1}|U_\tau, S_\tau)}[\mathbf{G}_{\pi_{\tau+1}, S_{\tau+1}}^{aggre}(U_{\tau+1}, S_{\tau+1})] \quad (61)$$

Note that for $\tau = T-1$ the second term vanishes because future states beyond the temporal horizon are ignored, and thus $\mathbf{G}_{\pi_{T-1}, S_{T-1}}^{aggre}(U_{T-1}, S_{T-1}) = \mathbf{G}_{\pi_{T-1}, S_{T-1}}^{local}(U_{T-1}, S_{T-1})$. Also, the above equation will be useful in [Section 7.5](#) to show that BTAI is related to sophisticated inference.

Remark 6. The recursive aspect of [Eq. \(58\)](#) is deeply related to dynamic programming and the interested reader is referred to [Da Costa, Sajid et al. \(2020\)](#) for details about this relationship.

7.3. Branching Time Active Inference (BTAI)

In BTAI, the (localized) cost of the hidden state S_j is defined as $\mathbf{G}_I^{local} = g_I$, where g_I can be equal to $g_I^{classic}$, g_I^{feef} or g_I^{pcost} , and there are two ways of computing the aggregated cost of S_j . We can either propagate the localized cost towards the leaves (forward):

$$g_I \leftarrow g_I + g_{Nlast} \quad (62)$$

where the g_{Nlast} is the cost of the parent of S_j . Alternatively, we can back-propagate the cost towards the root

$$g_J \leftarrow g_J + g_I \quad \forall J \in \mathbb{A}, \quad (63)$$

where \mathbb{A} corresponds to all ancestors of the newly expanded node S_j .

7.4. BTAI as a generalization of active inference

To understand the relationship between BTAI and active inference, we need to focus on the forward propagation of the cost where the cost is given by $g_I^{classic}$. Recall that the update for forward propagation is given by:

$$g_I^{classic} \leftarrow g_I^{classic} + g_{Nlast}^{classic}, \quad (64)$$

where the g_{Nlast} is the cost of the parent of S_I , i.e., the parent of the newly expanded node. This equation tells us that the aggregated cost of S_I is equal to the localized cost of S_I plus the aggregated cost of S_{Nlast} , i.e.,

$$\underbrace{g_I^{classic}}_{\mathbf{G}_I^{aggre}} \leftarrow \underbrace{g_I^{classic}}_{\mathbf{G}_I^{local}} + \underbrace{g_{Nlast}^{classic}}_{\mathbf{G}_{Nlast}^{aggre}} \Leftrightarrow \mathbf{G}_I^{aggre} = \mathbf{G}_I^{local} + \mathbf{G}_{Nlast}^{aggre}, \quad (65)$$

but, then, we also recall [\(48\)](#), i.e.,

$$\mathbf{G}_{\pi_N}^{aggre} \approx \mathbf{G}_{\pi_{N-1}}^{aggre} + \mathbf{G}_{\pi_N}^{local} \Leftrightarrow \mathbf{G}_{\pi_N}^{aggre} \approx \mathbf{G}_{\pi_N}^{local} + \mathbf{G}_{\pi_{N-1}}^{aggre}. \quad (66)$$

The only difference between [Eqs. \(48\)](#) and [\(65\)](#) is notational. Indeed, in BTAI ([Eq. \(48\)](#)) a policy is represented by a multi-index denoting the sequence of actions selected, e.g., $I = (1, 2)$ corresponds to a policy of size two consisting of action one

followed by action two. In contrast, in active inference, a policy is a sequence of actions, e.g., $\pi_2 = (1, 2)$ corresponds to the same policy as the one described by I .

7.5. Relationship between BTAI and sophisticated inference

The relationship between BTAI and sophisticated inference is slightly more involved. The backward propagation equation, i.e.,

$$g_J \leftarrow g_J + g_I \quad \forall J \in \mathbb{A}, \quad (67)$$

tells us that when expanding a node S_I , we first need to compute its localized cost g_I and then add g_I to the aggregated cost of its ancestors S_J where $J \in \mathbb{A}$. In other words, we can rewrite the backward propagation equation as the following: the aggregated cost of an arbitrary node S_J will equal the sum of its localized cost g_J and that of its descendants \mathbb{D}_J that have already been evaluated

$$\mathbf{G}_J^{\text{aggre}} = \mathbf{G}_J^{\text{local}} + \sum_{S_K \in \mathbb{D}_J} \mathbf{G}_K^{\text{local}}, \quad (68)$$

where the descendants are the children, children of children, etc. We can further simplify this expression by grouping the summands by children of S_J . This leads us to:

$$\mathbf{G}_J^{\text{aggre}} = \mathbf{G}_J^{\text{local}} + \sum_{S_I \in \text{ch}_J} \left(\mathbf{G}_I^{\text{local}} + \sum_{S_K \in \mathbb{D}_I} \mathbf{G}_K^{\text{local}} \right) \quad (69)$$

$$= \mathbf{G}_J^{\text{local}} + \sum_{S_I \in \text{ch}_J} \mathbf{G}_I^{\text{aggre}}, \quad (70)$$

where ch_J are the children of S_J and \mathbb{D}_I are the descendants of S_I . The above equation has clear similarities to Eq. (61), which is,

$$\mathbf{G}^{\text{aggre}}(U_\tau, S_\tau) = \mathbf{G}^{\text{local}}(U_\tau, S_\tau) + \mathbb{E}_{Q(U_{\tau+1}, S_{\tau+1} | U_\tau, S_\tau)} [\mathbf{G}^{\text{aggre}}(U_{\tau+1}, S_{\tau+1})]. \quad (71)$$

However, the second term of the RHS of (61) is an expectation, while the second term of the RHS of (69) is a summation over the children that have already been expanded. The expectation in (61) is w.r.t.

$$Q(U_{\tau+1}, S_{\tau+1} | U_\tau, S_\tau) = Q(U_{\tau+1} | S_{\tau+1}) Q(S_{\tau+1} | U_\tau, S_\tau), \quad (72)$$

where:

$$Q(U_{\tau+1} | S_{\tau+1}) > 0 \Leftrightarrow U_{\tau+1} \in \arg \min_{U \in \mathbb{U}} \mathbf{G}^{\text{aggre}}(U, S_{\tau+1}), \quad (73)$$

where \mathbb{U} is the set of all possible actions, and $\arg \min$ is defined as:

$$M = \arg \min_{U \in \mathbb{U}} f(U) \Leftrightarrow M = \left\{ U \in \mathbb{U} \mid f(U) = \min_{U' \in \mathbb{U}} f(U') \right\}. \quad (74)$$

This means that an action is assigned positive probability mass if and only if it minimizes the aggregated cost at the next time point $\mathbf{G}^{\text{aggre}}(U_{\tau+1}, S_{\tau+1})$ and a set is required, because multiple actions could have the same minimum cost. Note that if there is a unique minimum, then $Q(U_{\tau+1} | S_{\tau+1})$ will be a one-hot like distribution with a probability of one for the best action.

To conclude, Eqs. (61) and (69) suggest that BTAI and SI share a similar notion of EFE, where the immediate (or localized) EFE is added to the future (or aggregated) EFE. Both BTAI and SI propagate the cost backward, however, in SI the aggregated EFE (i.e., the back-propagated cost) is weighted by the probability of the next action and states, i.e., $Q(U_{\tau+1}, S_{\tau+1} | U_\tau, S_\tau)$. Intuitively, the weighting terms in SI discounts the impact of the back-propagated cost for unlikely states and (locally) sub-optimal actions. Importantly, those weighting terms emerge from

the recursive definition of the EFE that relates to the Bellman equation (Da Costa, Sajid et al., 2020). In contrast, there are no such weights in BTAI because BTAI finds its inspiration in active inference.

8. Conclusion and future works

In this paper, we have presented a new approach where planning is cast as structure learning. Simply put, this approach consists of dynamically expanding branches of the generative model by evaluating alternative futures under different action sequences. The dynamic expansion trades off evaluating promising (with respect to the target distribution) policies with exploring policies whose outcomes are uncertain. We proposed two different tree search methods: the first in which the nodes' cost is propagated forward from the root node to the leaves; the second in which the nodes' cost is propagated backward from the leaves to the root node. Then, in Section 7 we showed that forward propagation of the EFE leads to active inference (Acl) under the assumption that the probability of observations and states are independent of future actions, and that backward propagation relates to sophisticated inference (SI). This clarifies the link between Acl and BTAI, and helps to understand the relationship between Acl and SI.

Importantly, by performing a complexity class analysis, we have shown that while Active Inference suffers from an exponential complexity class, our approach scales nicely (linearly) with the number of tree expansions, cf., Section 3.4.2 of our companion paper (Champion, Bowman et al., 2021). Of course, the total number of possible expansions grows exponentially but as has been empirically shown in the reinforcement learning literature, even complex tasks, such as chess or the game of go, can be performed efficiently with MCTS (Schrittwieser et al., 2019; Silver et al., 2016), i.e., the efficiency of MCTS-like approaches relies on the ability to guide the expansion procedure using either powerful heuristics (like the EFE) or neural networks or both.

We also know that humans engage in counterfactual reasoning (Rafetseder, Schrittwieser, & Perner, 2013), which, in our planning context, could involve the consideration and evaluation of alternative (non-selected) sequences of decisions. It may be that, because of the more exhaustive representation of possible trajectories, the classic active inference can more efficiently engage in counterfactual reasoning. In contrast, branching-time active inference would require these alternatives to be generated “a fresh” for each counterfactual deliberation. In this sense, one might argue that there is a trade-off: branching-time active inference provides considerably more efficient planning to attain current goals, classic active inference provides a more exhaustive assessment of paths not taken.

Now that we have laid out the mathematics of BTAI, many directions of research could be investigated. One could for example obtain an intuitive understanding of the model's parameters through experimental study. At first it might be necessary to restrict oneself to agents without learning, i.e., inference only. This step should help answer questions such as: How does the number of expansions of the tree and the quality of the prior preferences impact the quality of planning? What is the best inference method (i.e., local or global inference) to use during planning?

Then, one could consider learning of the transition and likelihood matrices as well as the vector of initial states. This can be done in at least two ways. The first is to add Dirichlet priors over those matrices/vectors and the second would be to use neural networks as function approximators. The second option will lead to a deep active inference agent (Millidge, 2020; Sancaktar & Lanillos, 2020) equipped with tree search that could be directly

compared to the method of Fountas et al. (2020). Including deep neural networks in the framework will also enable direct comparison with the deep reinforcement learning literature (Haarnoja, Zhou, Abbeel, & Levine, 2018; Lample & Chaplot, 2017; Mnih et al., 2013; Silver et al., 2016; van Hasselt, Guez, & Silver, 2016). These comparisons will enable the impact of epistemic terms to be studied when the agent is composed of deep neural networks.

Another, very important direction for future research would be the creation of a biologically plausible implementation of BTAI. For example, using artificial neural networks to model the various mappings of the framework may provide a neural-based implementation of BTAI that is closer to biology. This would especially be the case, if the back-propagation algorithm frequently used for learning is replaced by the (more biologically plausible) generalized recirculation algorithm (O'Reilly, 1996). Another possible approach would be to use populations of neurons to encode the update equations of the framework, as was proposed by Friston, Parr, and de Vries (2017).

Whatever technique is chosen for learning and inference, implementing MCTS in a biologically plausible way will be challenging. Indeed, MCTS requires a dynamic expansion of the search tree used to explore the space of possible policies. Each time an expansion is performed, the agent needs to store the associated variables such as: the number of visits, the aggregated expected free energy, and the posterior beliefs of the newly expanded node. Given the fast pace at which planning must be performed to be useful, slow mechanisms such as synaptic plasticity and neurogenesis are likely to be unsuitable for the task. A more plausible approach might rely upon a change of neuronal activation, which can occur within a few hundred milliseconds. One such approach uses a binding pool (Bowman & Wyble, 2007) and provides a notion of variable. In this framework, a variable is composed of two parts. First, a *token* that can be intuitively understood as the variable's name, and second, a *type* corresponding to the variable's value. The binding pool is then composed of neurons representing the fact that a variable's name is *bound* (or set) to a specific value. A localist realization of a binding pool could be implemented as a 2D array of neurons of size “number of tokens” \times “number of types”. However, such a representation is quite inefficient and a more compact (i.e. distributed) representation has been developed (Wyble & Bowman, 2006). If variables are complex data structures, such as those required by BTAI, they can be realized in a neural substrate.

Finally, in this paper, we focused on the UCT criterion for node selection because it is a standard choice in the reinforcement learning literature (Schrittwieser et al., 2019; Silver et al., 2016). However, it would be valuable to consider alternative criteria such as Thompson Sampling (Russo, Van Roy, Kazerouni, Osband, & Wen, 2018; Sajid, Ball, Parr, & Friston, 2021; Thompson, 1933) or expected improvement (Bergstra, Bardenet, Bengio, & Kégl, 2011; Brochu, Cora, & de Freitas, 2010). For example, Thompson Sampling has been shown to improve upon the standard UCT criterion when applied to MCTS (Bai, Wu, & Chen, 2013), but requires additional modelling and we leave this for future research.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

We would like to thank the reviewers for their valuable feedback, which greatly improved the quality of the present paper.

LD is supported by the Fonds National de la Recherche, Luxembourg (Project code: 13568875). This publication is based on work partially supported by the EPSRC, United Kingdom Centre for Doctoral Training in Mathematics of Random Systems: Analysis, Modelling and Simulation (EP/S023925/1).

Appendix A. Generalized inner and outer products

Generalized outer products:. Given N vectors V^i , the generalized outer product returns an N dimensional array W , whose element in position (x_1, \dots, x_N) is given by $V_{x_1}^1 \times \dots \times V_{x_N}^N$, where $V_{x_j}^i$ is the x_j -th element of the i th vector. In other words:

$$W = \otimes [V^1, \dots, V^N] \Leftrightarrow W(x_1, \dots, x_N) = V_{x_1}^1 \times \dots \times V_{x_N}^N \quad \forall x_j \in \{1, \dots, |V^j|\} \quad \forall j \in \{1, \dots, N\}, \quad (75)$$

where $|V^j|$ is the number of elements in V^j . Also, note that by definition W is a N -tensor of size $|V^1| \times \dots \times |V^N|$. Fig. 7 illustrates the generalized outer product for $N = 3$.

Generalized inner products:. Given an N -tensor W and $M = N - 1$ vectors V^i , the generalized inner product returns a vector Z obtained by performing a weighted average (with weighting coming from the vectors) over all but one dimension. In other words:

$$Z = W \odot [V^1, \dots, V^M] \Leftrightarrow \forall x_j \in \{1, \dots, |Z|\}, \quad Z(x_j) = \sum_{\substack{x_1 \in \{1, \dots, |V^1|\} \\ \dots \\ x_M \in \{1, \dots, |V^M|\}}} V_{x_1}^1 \times \dots \times W(x_1, \dots, x_j, \dots, x_M) \times \dots \times V_{x_M}^M \quad (76)$$

where $|Z|$ denotes the number of elements in Z , and the large summand is over all x_r for $r \in \{1, \dots, M\} \setminus \{j\}$, i.e., excluding j . Also, note that if $|W|_{|V^i|} \forall i \in \{1, \dots, M\}$ is the number of elements in the dimension corresponding to V^i , then for $W \odot [V^1, \dots, V^M]$ to be properly defined, we must have $|W|_{|V^i|} = |V^i| \forall i \in \{1, \dots, M\}$ where $|V^i|$ is the number of elements in V^i . Fig. 8 illustrates the generalized inner product for $N = 3$.

Naming of the dimensions:. Importantly, we should imagine that each side of W has a name, e.g., if W is a 3×2 matrix, then the i th dimension of W could be named: “the dimension of V_i ”. This enables us to write: $Z^1 = W \odot V^1$ and $Z^2 = W \odot V^2$, where Z^1 is a 1×2 matrix (i.e., a vector with two elements) and Z^2 is a 3×1 matrix (i.e., a vector with three elements). The operator \odot knows (thanks to the dimension name) that $W \odot V^1$ takes the weighted average w.r.t “the dimension of V_1 ”, while $W \odot V^2$ must take the weighted average over “the dimension of V_2 ”.

In the context of active inference, the matrix A has two dimensions that we could call “the observation dimension” (i.e., row-wise) and “the state dimension” (i.e., column-wise). Trivially, $A \odot \mathbf{o}_\tau$ will then correspond to the average of A along the observation dimension and $A \odot \hat{D}_\tau$ will correspond to the average of A along the state dimension.

Appendix B. Generalized inner/outer products and other well-known products

In this section, we explore the relationship between our generalized inner and outer products—presented in Appendix A—and other well known products in the literature.

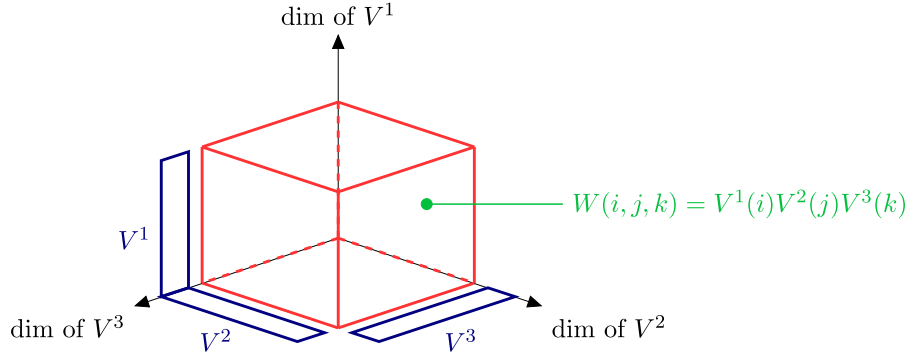


Fig. 7. This figure illustrates the generalized outer product $W = \otimes[V^1, V^2, V^3]$, where W is a cube of values illustrated in red, whose typical element $W(i, j, k)$ is the product of $V^1(i)$, $V^2(j)$ and $V^3(k)$. Also, the vectors $V^i \forall i \in \{1, \dots, 3\}$ are drawn in blue along the dimension of the cube they correspond to. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

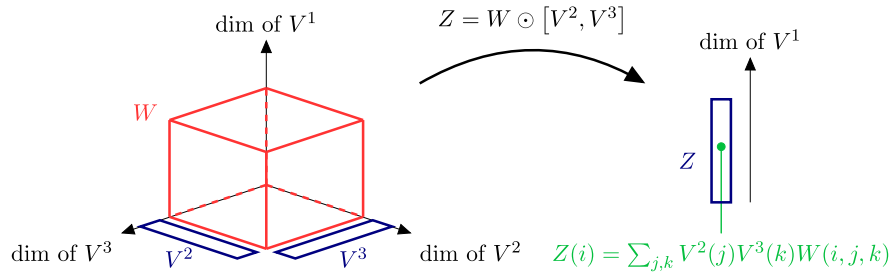


Fig. 8. This figure illustrates the generalized inner product $Z = W \odot [V^2, V^3]$, where W is a cube of values illustrated in red with typical element $W(i, j, k)$. Also, the vectors Z and $V^i \forall i \in \{2, 3\}$ are drawn in blue along the dimension of the cube they correspond to. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Inner product of two vectors

The inner product of two vectors \vec{a} and \vec{b} of the same size is:

$$\vec{a} \cdot \vec{b} = \sum_{i=1}^{|\vec{a}|} a_i b_i, \quad (77)$$

where a_i and b_i are the elements of the vectors \vec{a} and \vec{b} , respectively, and $|\vec{a}|$ is the number of elements in \vec{a} . This product is a special case of our generalized inner product, i.e.,

$$Z = W \odot V \Leftrightarrow Z = \sum_{i=1}^{|W|} W_i V_i, \quad (78)$$

where Z is a scalar (a 0-tensor), W and V are two vectors (two 1-tensors) of the same size, and $|W|$ is the number of elements in W .

Inner product of two matrices (Frobenius inner product)

The inner product of two matrices \mathbf{A} and \mathbf{B} of same sizes is:

$$\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i=1}^{|\mathbf{A}|_1} \sum_{j=1}^{|\mathbf{A}|_2} a_{ij} b_{ij}, \quad (79)$$

where $|\mathbf{A}|_1$ is the number of rows in \mathbf{A} , $|\mathbf{A}|_2$ is the number of columns in \mathbf{A} , and a_{ij} (b_{ij}) are the elements of the matrices \mathbf{A} (respectively \mathbf{B}). This product is not a special case of our generalized inner product.

Inner product of two tensors (Frobenius inner product)

The inner product of two tensors \mathbf{A} and \mathbf{B} of same sizes is:

$$\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i_1=1}^{|\mathbf{A}|_1} \dots \sum_{i_n=1}^{|\mathbf{A}|_n} a(i_1, \dots, i_n) b(i_1, \dots, i_n), \quad (80)$$

where $a(i_1, \dots, i_n)$ and $b(i_1, \dots, i_n)$ are the elements of the tensors \mathbf{A} and \mathbf{B} , respectively, and $|\mathbf{A}|_i$ is the number of elements in the i th dimension of \mathbf{A} . This product is not a special case of our generalized inner product.

Standard matrix multiplication

Let \mathbf{A} be an $n \times m$ matrix and \vec{b} be a vector of size m . The standard matrix multiplication of \mathbf{A} by \vec{b} is given by:

$$\vec{c} = \mathbf{A}\vec{b} \Leftrightarrow \vec{c}_i = \sum_{j=1}^m \mathbf{A}_{ij} \vec{b}_j. \quad (81)$$

This is a special case of our generalized inner product, i.e.,

$$\vec{c} = \mathbf{A}\vec{b} = \mathbf{A} \odot \vec{b}. \quad (82)$$

Additionally, let \mathbf{A} be an $n \times m$ matrix and \vec{a} be a vector of size n , then:

$$\vec{c} = \mathbf{A}^T \vec{a} \Leftrightarrow \vec{c}_i = \sum_{j=1}^n \mathbf{A}_{ji} \vec{a}_j. \quad (83)$$

This is a special case of our generalized inner product, i.e.,

$$\vec{c} = \mathbf{A}^T \vec{a} = \mathbf{A} \odot \vec{a}. \quad (84)$$

Note that because the dimensions are “named” (cf. Appendix A) the operator performs the transposition implicitly.

Outer product of two vectors

Given two vectors \vec{a} and \vec{b} , their outer product—denoted $\vec{a} \otimes \vec{b}$ —is a matrix defined as:

$$\vec{a} \otimes \vec{b} = \begin{bmatrix} a_1 b_1 & \dots & a_1 b_n \\ \vdots & \ddots & \vdots \\ a_m b_1 & \dots & a_m b_n \end{bmatrix}. \quad (85)$$

This outer product is a special case of our generalized outer product, where the operator is applied to only two vectors, i.e.

$$\vec{a} \otimes \vec{b} = \otimes[\vec{a}, \vec{b}]. \quad (86)$$

Outer product of two tensors

Given two tensors \mathbf{U} and \mathbf{V} , the outer product of \mathbf{U} and \mathbf{V} is another tensor \mathbf{W} such that:

$$\begin{aligned} \mathbf{W}(i_1, \dots, i_n, j_1, \dots, j_m) &= \mathbf{V}(i_1, \dots, i_n) \mathbf{U}(j_1, \dots, j_m), \\ \forall \alpha \in \{1, \dots, |\mathbf{V}|_\alpha\} \forall \alpha \in \{1, \dots, n\}, \\ \forall \beta \in \{1, \dots, |\mathbf{U}|_\beta\} \forall \beta \in \{1, \dots, m\}. \end{aligned} \quad (87)$$

Given N vectors $V^i \forall i \in \{1, \dots, N\}$, our outer product is a sequence of outer tensor products, i.e.,

$$\otimes[V^1, \dots, V^N] = \left[\left[V^1 \otimes_{\text{tensor}} V^2 \right] \otimes_{\text{tensor}} \dots \right] \otimes_{\text{tensor}} V^N, \quad (88)$$

where \otimes_{tensor} and \otimes are the tensor and generalized outer products, respectively.

Kronecker product

Given two matrices \mathbf{A} and \mathbf{B} , the Kronecker product of \mathbf{A} and \mathbf{B} —denoted \otimes_K —is a generalization of the outer product from vectors to matrices defined as:

$$\mathbf{A} \otimes_K \mathbf{B} = \begin{bmatrix} a_{11} \mathbf{B} & \dots & a_{1n} \mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1} \mathbf{B} & \dots & a_{mn} \mathbf{B} \end{bmatrix}, \quad (89)$$

where a_{ij} are the elements of \mathbf{A} . Note that even if the Kronecker product is a generalization of the outer product, it is neither a special case nor a generalization of our generalized outer product.

Hadamard product

Given two matrices \mathbf{A} and \mathbf{B} of the same size, the Hadamard product of \mathbf{A} and \mathbf{B} is an element-wise product defined by:

$$\mathbf{C} = \mathbf{A} \odot \mathbf{B} \Leftrightarrow c_{ij} = a_{ij} b_{ij} \forall i \in \{1, \dots, |\mathbf{A}|_1\} \forall j \in \{1, \dots, |\mathbf{A}|_2\}, \quad (90)$$

where a_{ij} , b_{ij} and c_{ij} are the elements in the i th row and j th column of \mathbf{A} , \mathbf{B} and \mathbf{C} , respectively, $|\mathbf{A}|_1$ is the number of rows in \mathbf{A} , and $|\mathbf{A}|_2$ is the number of columns in \mathbf{A} . This product is unrelated to both our generalized inner and outer products.

Appendix C. Instance of variational message passing

This appendix provides a concrete instance of the method of Winn and Bishop discussed in Section 3. The generative model is as follows:

$$P(S, \mathbf{D}) = P(S|\mathbf{D})P(\mathbf{D}) \quad (91)$$

where:

$$P(S|\mathbf{D}) = \text{Cat}(\mathbf{D}) \quad (92)$$

$$P(\mathbf{D}) = \text{Dir}(\mathbf{d}). \quad (93)$$

Additionally, the variational distribution is given by:

$$Q(\mathbf{D}) = \text{Dir}(\hat{\mathbf{d}}), \quad (94)$$

which means that we assume that S is an observed random variable. Let us start with the definition of the Dirichlet and categorical distributions written in the form of the exponential family:

$$\ln P(\mathbf{D}) = \underbrace{\left[\begin{matrix} \mathbf{d}_1 - 1 \\ \dots \\ \mathbf{d}_{|S|} - 1 \end{matrix} \right]}_{\mu_D(\mathbf{D})} \cdot \underbrace{\left[\begin{matrix} \ln \mathbf{D}_1 \\ \dots \\ \ln \mathbf{D}_{|S|} \end{matrix} \right]}_{u_D(\mathbf{D})} - \underbrace{\ln B(\mathbf{d})}_{z_D(\mathbf{d})} \quad (95)$$

$$\ln P(S|\mathbf{D}) = \underbrace{\left[\begin{matrix} \ln \mathbf{D}_1 \\ \dots \\ \ln \mathbf{D}_{|S|} \end{matrix} \right]}_{\mu_S(\mathbf{D})} \cdot \underbrace{\left[\begin{matrix} [S = 1] \\ \dots \\ [S = |S|] \end{matrix} \right]}_{u_S(S)} \quad (96)$$

where \cdot performs an inner product of the two vectors it is applied to, $B(\mathbf{d})$ is the Beta function and $|S|$ is the number of values a state can take. The first step requires us to re-write Eq. (96) as a function of $u_D(\mathbf{D})$, which is straightforward because $\mu_S(\mathbf{D})$ is just another name for $u_D(\mathbf{D})$. Using the fact that the inner product is commutative:

$$\ln P(S|\mathbf{D}) = \underbrace{\left[\begin{matrix} [S = 1] \\ \dots \\ [S = |S|] \end{matrix} \right]}_{\mu_{S \rightarrow D}(S)} \cdot \underbrace{\left[\begin{matrix} \ln \mathbf{D}_1 \\ \dots \\ \ln \mathbf{D}_{|S|} \end{matrix} \right]}_{u_D(\mathbf{D})}. \quad (97)$$

The second step aims to substitute (95) and (97) into the variational message passing Eq. (7), i.e.

$$\begin{aligned} \ln Q^*(\mathbf{D}) &= \left\langle \underbrace{\left[\begin{matrix} \mathbf{d}_1 - 1 \\ \dots \\ \mathbf{d}_{|S|} - 1 \end{matrix} \right]}_{\mu_D(\mathbf{d})} \cdot \underbrace{\left[\begin{matrix} \ln \mathbf{D}_1 \\ \dots \\ \ln \mathbf{D}_{|S|} \end{matrix} \right]}_{u_D(\mathbf{D})} - \underbrace{\ln B(\mathbf{d})}_{z_D(\mathbf{d})} \right\rangle \\ &\quad + \left\langle \underbrace{\left[\begin{matrix} [S = 1] \\ \dots \\ [S = |S|] \end{matrix} \right]}_{\mu_{S \rightarrow D}(S)} \cdot \underbrace{\left[\begin{matrix} \ln \mathbf{D}_1 \\ \dots \\ \ln \mathbf{D}_{|S|} \end{matrix} \right]}_{u_D(\mathbf{D})} \right\rangle + \text{Const}, \end{aligned} \quad (98)$$

where $\langle \cdot \rangle$ refers to $\langle \cdot \rangle_{\sim Q_D}$. Note that in the above equation, \mathbf{d}_i are fixed parameters, therefore there is no posterior over \mathbf{d} and the first expectation $\langle \cdot \rangle_{\sim Q_D}$ can be removed. The third step rests on taking the exponential of both sides, using the linearity of expectation and factorizing by $u_D(\mathbf{D})$ to obtain:

$$Q^*(\mathbf{D}) = \exp \left\{ \left[\begin{matrix} \mathbf{d}_1 - 1 + \langle [S = 1] \rangle \\ \dots \\ \mathbf{d}_{|S|} - 1 + \langle [S = |S|] \rangle \end{matrix} \right] \cdot u_D(\mathbf{D}) + \text{Const} \right\}, \quad (99)$$

where $z_D(\mathbf{d})$ have been absorbed into the constant term because it does not depend on \mathbf{D} . The fourth step is a re-parameterization done by observing that $\langle [S = i] \rangle$ is the i th element of the expectation of the vector $u_S(S)$, i.e. $\langle u_S(S) \rangle_i = \langle [S = i] \rangle$:

$$Q^*(\mathbf{D}) = \exp \left\{ \underbrace{\left[\begin{matrix} \mathbf{d}_1 - 1 + \langle u_S(S) \rangle_1 \\ \dots \\ \mathbf{d}_{|S|} - 1 + \langle u_S(S) \rangle_{|S|} \end{matrix} \right]}_{\bar{\mu}_D(\dots) + \bar{\mu}_{S \rightarrow D}(\dots)} \cdot u_D(\mathbf{D}) + \text{Const} \right\}. \quad (100)$$

The last step consists of computing the expectation of $\langle u_S(S) \rangle_i$ for all i . This can be achieved by realizing that the probability of an indicator function for an event is the probability of this event, i.e. $\langle u_S(S) \rangle_i = \langle [S = i] \rangle = Q(S = i) = \hat{\mathbf{D}}_i$ where $\hat{\mathbf{D}}$ is a one hot

representation of the observed value for S . Substituting this result in Eq. (100), leads to the final result:

$$Q^*(\mathbf{D}) = \exp \left\{ \left[\begin{array}{c} \mathbf{d}_1 - 1 + \hat{\mathbf{D}}_1 \\ \dots \\ \mathbf{d}_{|S|} - 1 + \hat{\mathbf{D}}_{|S|} \end{array} \right] \cdot u_D(\mathbf{D}) + \text{Const} \right\}. \quad (101)$$

Indeed, the above equation is in fact a Dirichlet distribution in exponential family form, and can be re-written into its usual form to obtain the final update equation:

$$Q^*(\mathbf{D}) = \text{Dir}(\mathbf{d} + \hat{\mathbf{D}}). \quad (102)$$

Appendix D. Expected log of Dirichlet distribution

Definition 7. A probability distribution over \mathbf{x} parameterized by $\boldsymbol{\mu}$ is said to belong to the exponential family if its probability mass function $P(\mathbf{x}|\boldsymbol{\mu})$ can be written as:

$$P(\mathbf{x}|\boldsymbol{\mu}) = h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x}) - A(\boldsymbol{\mu})], \quad (103)$$

where $h(\mathbf{x})$ is the base measure, $\boldsymbol{\mu}$ is the vector of natural parameters, $T(\mathbf{x})$ is the vector of sufficient statistics, and $A(\boldsymbol{\mu})$ is the log partition.

Lemma 8. The log partition is given by:

$$A(\boldsymbol{\mu}) = \ln \int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] d\mathbf{x}. \quad (104)$$

Proof. Starting with the fact that $P(\mathbf{x}|\boldsymbol{\mu})$ integrate to one:

$$\int P(\mathbf{x}|\boldsymbol{\mu}) d\mathbf{x} = \int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x}) - A(\boldsymbol{\mu})] d\mathbf{x} = 1 \quad (105)$$

$$\Leftrightarrow \frac{1}{\exp A(\boldsymbol{\mu})} \int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] d\mathbf{x} = 1 \quad (106)$$

$$\Leftrightarrow A(\boldsymbol{\mu}) = \ln \int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] d\mathbf{x} \quad (107)$$

■

Lemma 9. The gradient of the log partition function is the expectation of the sufficient statistics, i.e.,

$$\frac{\partial A(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = \mathbb{E}_{P(\mathbf{x}|\boldsymbol{\mu})}[T(\mathbf{x})]. \quad (108)$$

Proof. Restarting with the derivative of the result of Lemma 8:

$$\frac{\partial A(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = \frac{\partial}{\partial \boldsymbol{\mu}} \left[\ln \int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] d\mathbf{x} \right], \quad (109)$$

and using the chain rule:

$$\frac{\partial A(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = \frac{1}{\int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] d\mathbf{x}} \frac{\partial}{\partial \boldsymbol{\mu}} \left[\int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] d\mathbf{x} \right]. \quad (110)$$

Note that the denominator of the first term is equal to the exponential of $A(\boldsymbol{\mu})$, and we can swap the derivative and the integral because the limit of integration does not depend on the parameters $\boldsymbol{\mu}$:

$$\frac{\partial A(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = \frac{1}{\exp A(\boldsymbol{\mu})} \int \frac{\partial}{\partial \boldsymbol{\mu}} \left[h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] \right] d\mathbf{x} \quad (111)$$

$$= \frac{1}{\exp A(\boldsymbol{\mu})} \int h(\mathbf{x}) \frac{\partial}{\partial \boldsymbol{\mu}} \left[\exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] \right] d\mathbf{x}. \quad (112)$$

Using the chain rule again:

$$\frac{\partial A(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = \frac{1}{\exp A(\boldsymbol{\mu})} \int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x})] T(\mathbf{x}) d\mathbf{x} \quad (113)$$

$$= \int h(\mathbf{x}) \exp[\boldsymbol{\mu} \cdot T(\mathbf{x}) - A(\boldsymbol{\mu})] T(\mathbf{x}) d\mathbf{x} \quad (114)$$

$$= \int P(\mathbf{x}|\boldsymbol{\mu}) T(\mathbf{x}) d\mathbf{x} \quad (115)$$

$$= \mathbb{E}_{P(\mathbf{x}|\boldsymbol{\mu})}[T(\mathbf{x})]. \quad (116)$$

■

Theorem 10. If \mathbf{D} is distributed according to a Dirichlet distribution $Q(\mathbf{D}) = \text{Dir}(\mathbf{D}; \hat{\mathbf{d}})$, then: $\dot{\mathbf{D}} = \mathbb{E}_{Q(\mathbf{D})}[\ln \mathbf{D}] \Leftrightarrow \dot{\mathbf{D}}_i = \psi(d_i) - \psi(\sum_j d_j)$.

Proof. Let $\boldsymbol{\mu}$ be equal to $\hat{\mathbf{d}} - \bar{\mathbf{1}}$. Taking the exponential of both sides in Eq. (95) and using that $\hat{\mathbf{d}} = \boldsymbol{\mu} + \bar{\mathbf{1}}$, we obtain:

$$Q(\mathbf{D}) = \exp \left\{ \underbrace{\left[\begin{array}{c} \hat{\mathbf{d}}_1 - 1 \\ \dots \\ \hat{\mathbf{d}}_{|S|} - 1 \end{array} \right]}_{\boldsymbol{\mu}} \cdot \underbrace{\left[\begin{array}{c} \ln \mathbf{D}_1 \\ \dots \\ \ln \mathbf{D}_{|S|} \end{array} \right]}_{T(\mathbf{D})} - \underbrace{\ln B(\boldsymbol{\mu} + \bar{\mathbf{1}})}_{A(\boldsymbol{\mu})} \right\}, \quad (117)$$

where $\boldsymbol{\mu}$ is the vector of natural parameters, $T(\mathbf{D})$ is the vector of sufficient statistics, $A(\boldsymbol{\mu})$ is the log partition, and $B(\cdot)$ is the beta function. Using the result of Lemma 9:

$$\begin{aligned} \dot{\mathbf{D}} \triangleq \mathbb{E}_{Q(\mathbf{D})}[T(\mathbf{D})] &= \mathbb{E}_{Q(\mathbf{D})}[\ln \mathbf{D}] = \frac{\partial A(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} \\ &= \frac{\partial}{\partial \boldsymbol{\mu}} [\ln B(\boldsymbol{\mu} + \bar{\mathbf{1}})]. \end{aligned} \quad (118)$$

We now focus on a typical element of $\dot{\mathbf{D}}$:

$$\dot{\mathbf{D}}_i = \frac{\partial}{\partial \mu_i} [\ln B(\boldsymbol{\mu} + \bar{\mathbf{1}})], \quad (119)$$

and use the definition of the beta function:

$$\dot{\mathbf{D}}_i = \frac{\partial}{\partial \mu_i} \left[\ln \frac{\prod_k \Gamma(\mu_k + 1)}{\Gamma(\sum_k \mu_k + 1)} \right] \quad (120)$$

$$= \frac{\partial}{\partial \mu_i} \left[\sum_k \ln \Gamma(\mu_k + 1) - \ln \Gamma(\sum_k \mu_k + 1) \right] \quad (121)$$

$$= \frac{\partial}{\partial \mu_i} [\ln \Gamma(\mu_i + 1)] - \frac{\partial}{\partial \mu_i} [\ln \Gamma(\sum_k \mu_k + 1)], \quad (122)$$

where $\Gamma(\cdot)$ is the gamma function. The last step relies on the chain rule:

$$\begin{aligned} \dot{\mathbf{D}}_i &= \frac{\partial}{\partial (\mu_i + 1)} [\ln \Gamma(\mu_i + 1)] \underbrace{\frac{\partial \mu_i + 1}{\partial \mu_i}}_{=1} \\ &\quad - \frac{\partial}{\partial (\sum_k \mu_k + 1)} [\ln \Gamma(\sum_k \mu_k + 1)] \underbrace{\frac{\partial \sum_k \mu_k + 1}{\partial \mu_i}}_{=1} \end{aligned} \quad (123)$$

$$\begin{aligned} &= \frac{\partial}{\partial (\mu_i + 1)} [\ln \Gamma(\mu_i + 1)] \\ &\quad - \frac{\partial}{\partial (\sum_k \mu_k + 1)} [\ln \Gamma(\sum_k \mu_k + 1)] \end{aligned} \quad (124)$$

$$= \psi(\mu_i + 1) - \psi(\sum_k \mu_k + 1) \quad (125)$$

$$= \psi(\hat{\mathbf{d}}_i) - \psi(\sum_k \hat{\mathbf{d}}_k), \quad (126)$$

where we used that $\hat{\mathbf{d}} = \boldsymbol{\mu} + \bar{\mathbf{1}}$ and the definition of the digamma function, i.e., $\psi(x) = \frac{\partial \ln \Gamma(x)}{\partial x}$.

Appendix E. Relationship between BTAI and active inference (Lemmas)

Lemma 11. Under the assumption that the probability of observations and states are independent of future actions, i.e.,

$$\forall j \in \mathbb{N}_{>0}, Q(O_{t+i}|\pi_i) \approx Q(O_{t+i}|\pi_{i+j}) \text{ and } Q(S_{t+i}|\pi_i) \approx Q(S_{t+i}|\pi_{i+j}), \quad (127)$$

then:

$$\forall j \in \mathbb{N}_{>0}, \mathbf{G}(\pi_{i+j}, t+i) \approx \mathbf{G}(\pi_i, t+i). \quad (128)$$

Proof. The proof is straightforward, we start with the following definition:

$$\mathbf{G}(\pi_{i+j}, t+i) = D_{\text{KL}}[Q(O_{t+i}|\pi_{i+j}) \parallel P(O_{t+i})] + \mathbb{E}_{Q(S_{t+i}|\pi_{i+j})}[\mathbf{H}[P(O_{t+i}|S_{t+i})]]. \quad (129)$$

Then, using the assumption that the probability of observations and states are independent of future actions, i.e., $\forall j \in \mathbb{N}_{>0}, Q(O_{t+i}|\pi_{i+j}) \approx Q(O_{t+i}|\pi_i)$ and $\forall j \in \mathbb{N}_{>0}, Q(S_{t+i}|\pi_{i+j}) \approx Q(S_{t+i}|\pi_i)$, we get:

$$\mathbf{G}(\pi_{i+j}, t+i) \approx D_{\text{KL}}[Q(O_{t+i}|\pi_i) \parallel P(O_{t+i})] + \mathbb{E}_{Q(S_{t+i}|\pi_i)}[\mathbf{H}[P(O_{t+i}|S_{t+i})]] \triangleq \mathbf{G}(\pi_i, t+i). \quad (130)$$

■

Lemma 12. The aggregated cost for an arbitrary N is given by:

$$\mathbf{G}_{\pi_N}^{\text{aggre}} = \sum_{i=1}^N \mathbf{G}(\pi_i, t+i), \quad (131)$$

Proof. The proof is done by induction. The initialization holds for $N = 1$, indeed, $\pi_1 = \{U_t\}$ and by definition:

$$\underbrace{\mathbf{G}_{\pi_1}^{\text{aggre}}}_{\mathbf{G}_{\pi_1}^{\text{aggre}}} = \sum_{\tau=t+1}^{t+1} \mathbf{G}(\pi_1, \tau) = D_{\text{KL}}[Q(O_{t+1}|\pi_1) \parallel P(O_{t+1})] + \underbrace{\mathbb{E}_{Q(S_{t+1}|\pi_1)}[\mathbf{H}[P(O_{t+1}|S_{t+1})]]}_{\mathbf{G}_{\pi_1}^{\text{local}}}, \quad (132)$$

$$\Leftrightarrow \mathbf{G}_{\pi_1}^{\text{aggre}} = \mathbf{G}_{\pi_0}^{\text{aggre}} + \mathbf{G}_{\pi_1}^{\text{local}}, \quad (133)$$

because by definition $\mathbf{G}_{\pi_0}^{\text{aggre}} = 0$. Then, assuming that $\mathbf{G}_{\pi_N}^{\text{aggre}} = \sum_{i=1}^N \mathbf{G}(\pi_i, t+i)$ holds for some N , we show that it holds for $N+1$ as well. By definition:

$$\mathbf{G}_{\pi_{N+1}}^{\text{aggre}} = \mathbf{G}_{\pi_N}^{\text{aggre}} + \mathbf{G}_{\pi_{N+1}}^{\text{local}}, \quad (134)$$

and:

$$\mathbf{G}_{\pi_{N+1}}^{\text{local}} = \mathbf{G}(\pi_{N+1}, t+N+1). \quad (135)$$

Using the inductive hypothesis and the above two equations:

$$\mathbf{G}_{\pi_{N+1}}^{\text{aggre}} = \sum_{i=1}^N \mathbf{G}(\pi_i, t+i) + \mathbf{G}(\pi_{N+1}, t+N+1) \quad (136)$$

$$= \sum_{i=1}^{N+1} \mathbf{G}(\pi_i, t+i). \quad (137)$$

■

Appendix F. Notation

In this appendix, we introduce the notation used throughout this paper. The following sub-sections describe the notation related to sets of numbers, tensors, probability distributions, global variables, multi-indices and random variables, respectively.

Sets of numbers

Definition 13. Let $\mathbb{N}_{>0}$ be the set of all strictly positive integers defined as:

$$\mathbb{N}_{>0} = \{x \in \mathbb{N} \mid x > 0\}, \quad (138)$$

where \mathbb{N} is the set of all natural numbers.

Definition 14. Let $\mathbb{R}_{>0}$ be the set of all strictly positive real numbers defined as:

$$\mathbb{R}_{>0} = \{x \in \mathbb{R} \mid x > 0\}, \quad (139)$$

where \mathbb{R} is the set of all real numbers.

Tensors

Definition 15. An n -tensor is an n -dimensional array of values. Each element of an n -tensor is indexed by an n -tuple of non-negative integers, i.e., (x_1, \dots, x_n) where $x_i \in \mathbb{N}_{>0} \forall i \in \{1, 2, \dots, n\}$.

Definition 16. Let T be an n -tensor. The **element** of T indexed by the n -tuple (x_1, \dots, x_n) is a real number denoted by $T(x_1, \dots, x_n)$.

Example 1. Let T be a 2-tensor defined as:

$$T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}. \quad (140)$$

Then $T(1, 1) = 1$, $T(1, 2) = 2$, $T(2, 1) = 3$ and $T(2, 2) = 4$.

Remark 17. A 0-tensor is a scalar, a 1-tensor is a vector, and a 2-tensor is a matrix.

Definition 18. Let T be an n -tensor. The **size** of T is a vector of size n denoted $|T|$ whose i th element corresponds to the size of the i th dimension of T .

Example 2. Let T be a 2-tensor defined as:

$$T = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}. \quad (141)$$

Then $|T|_1 = 2$ and $|T|_2 = 3$.

Definition 19. Let T be an n -tensor. A **1-sub-tensor** of T is a 1-tensor obtained by selecting a 1-dimensional slice of T , i.e.,

$$T(x_1, \dots, x_{i-1}, \cdot, x_{i+1}, \dots, x_n), \quad (142)$$

where \cdot represents the selection of a 1-dimensional slice of T , and the values of all $x_{j \neq i}$ must be set to specific values in $\{1, \dots, |T|_j\}$. Fig. 9 (left) illustrates the notion of a 1-dimensional slice.

Definition 20. Let T be an n -tensor and $m < n$. An **m -sub-tensor** of T (denoted W) is an m -tensor obtained by selecting an m -dimensional slice of T , i.e.,

$$W = T(x_1, \dots, x_{i_1-1}, \cdot, x_{i_1+1}, \dots, x_{i_m-1}, \cdot, x_{i_m+1}, \dots, x_n), \quad (143)$$

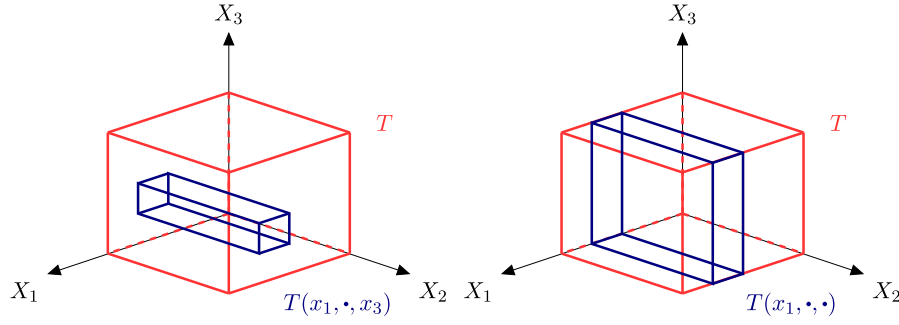


Fig. 9. This figure illustrates the notion of a 1-dimensional slice (on the left) and of a 2-dimensional slice (on the right).

where $i_k \in \{1, \dots, n\} \forall k \in \{1, \dots, m\}$ are indices representing the dimension being selected. Naturally, the k th dimension of W corresponds to the i_k -th dimension of T for $k \in \{1, \dots, m\}$. Importantly, the sextuple of dots in the middle of the expression represents that there will be m symbols “ \cdot ”, i.e., one for each dimension selected. Fig. 9 (right) illustrates the special case of a 2-dimensional slice, i.e., $m = 2$.

Example 3. Let T be a 3-tensor such that:

- $|T|_1 = 2$
- $T(1, x_2, x_3) = 1, \forall (x_2, x_3) \in \{1, \dots, |T|_2\} \times \{1, \dots, |T|_3\}$
- $T(2, x_2, x_3) = 2, \forall (x_2, x_3) \in \{1, \dots, |T|_2\} \times \{1, \dots, |T|_3\}$.

Then $T(1, \cdot, \cdot)$ is a 2-sub-tensor of T full of ones, and $T(2, \cdot, \cdot)$ is a 2-sub-tensor of T full of twos.

Probability distributions

Definition 21. A random n -tensor is an n -tensor over which we have an n -dimensional probability distribution.

Remark 22. A random variable is a random 0-tensor, a random vector is a random 1-tensor, and a random matrix is a random 2-tensor.

Definition 23. An n -tensor T is said to **represent a joint distribution** over a set of n random variables $\{X_1, \dots, X_n\}$ if:

$$P(X_1 = x_1, \dots, X_n = x_n) = T(x_1, \dots, x_n). \quad (144)$$

For conciseness, if T represents $P(X_1, \dots, X_n)$ we let:

$$P(X_1, \dots, X_n) = \text{Cat}(T). \quad (145)$$

Remark 24. If T represents $P(X_1, \dots, X_n)$, then the sum of its elements must equal one.

Remark 25. In contrast, if T is a **random** n -tensor, then:

$$P(X_1, \dots, X_n | T) = \text{Cat}(T), \quad (146)$$

which means that the joint probability over X_1, \dots, X_n is represented by T , and because T is a **random** tensor (taking values in the set of valid n -tensors \mathbb{T}_n , i.e., the set of all n -tensors whose elements sum up to one), we must specify which instance of $T \in \mathbb{T}_n$ should be used to define the joint distribution over X_1, \dots, X_n .

Definition 26. An m -tensor R is said to **represent a conditional distribution** over a set of m random variables $\{X_1, \dots, X_m\}$ if:

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n | X_{n+1} \\ = x_{n+1}, \dots, X_m = x_m) = R(x_1, \dots, x_m). \end{aligned} \quad (147)$$

For conciseness, if R represents $P(X_1, \dots, X_n | X_{n+1}, \dots, X_m)$ we let:

$$P(X_1, \dots, X_n | X_{n+1}, \dots, X_m) = \text{Cat}(R). \quad (148)$$

Remark 27. If R represents $P(X_1, \dots, X_n | X_{n+1}, \dots, X_m)$, then the elements of the n -sub-tensor $R(\cdot, \dots, \cdot, x_{n+1}, \dots, x_m)$ must sum to one $\forall x_i \in \{1, \dots, |R|_i\} \forall i \in \{n+1, \dots, m\}$.

Remark 28. If R is a random m -tensor, then:

$$P(X_1, \dots, X_n | X_{n+1}, \dots, X_m, R) = \text{Cat}(R). \quad (149)$$

Remark 29. Importantly, Definition 26 uses the symbol T to represent $P(X_1, \dots, X_n)$ and Definition 23 uses the symbol R to represent $P(X_1, \dots, X_n | X_{n+1}, \dots, X_m)$. Throughout this document, different symbols will be used for representing joint and conditional distributions.

Definition 30. Let R be a random m -tensor representing $P(X_1, \dots, X_n | X_{n+1}, \dots, X_m, R)$ and $k = m - n$ be the number of variables upon which the variables X_1, \dots, X_n are conditioned. Having a **Dirichlet prior** over R means that:

$$P(R) = \prod_{i_1=1}^{|R|_{n+1}} \dots \prod_{i_k=1}^{|R|_m} \text{Dir}(\mathbf{r}(i_1, \dots, i_k, \cdot)), \quad (150)$$

where \mathbf{r} is an $(m+1)$ -tensor such that the 1-sub-tensor $\mathbf{r}(i_1, \dots, i_k, \cdot)$ contains the parameters of the Dirichlet prior over $P(X_1, \dots, X_n | X_{n+1} = i_1, \dots, X_m = i_k)$. For conciseness, we denote the prior over R as:

$$P(R) = \text{Dir}(\mathbf{r}). \quad (151)$$

Remark 31. Definition 30 implicitly means that if V is a 1-tensor then $\text{Dir}(V)$ represents a Dirichlet distribution. However, if V is an m -tensor (with $m \neq 1$) then $\text{Dir}(V)$ represents a product of Dirichlet distributions.

Remark 32. If R is a random m -tensor representing $P(X_1, \dots, X_n | X_{n+1}, \dots, X_m, R)$, then its prior will be a product of $|X_{n+1}| \times \dots \times |X_m|$ Dirichlet distributions, where $|X_i|$ is the number of values that X_i can take. Additionally, each Dirichlet distribution will have $|X_1| \times \dots \times |X_n|$ parameters stored in the last dimension of \mathbf{r} , where \mathbf{r} is the tensor storing the parameters of the prior over R , i.e. $P(R) = \text{Dir}(\mathbf{r})$.

Example 4. Let \mathbf{A} be a random 2-tensor representing $P(O|S, \mathbf{A})$, then the Dirichlet prior over \mathbf{A} is given by:

$$P(\mathbf{A}) = \text{Dir}(\mathbf{a}) \triangleq \prod_{i=1}^{|A|_2} \text{Dir}(\mathbf{a}(i, \cdot)), \quad (152)$$

where $|A|_2$ is the number of values that S can take (i.e., the number of hidden states).

Example 5. Let \mathbf{B} be a random 3-tensor representing $P(S_{\tau+1}|S_\tau, U_\tau, \mathbf{B})$, then the Dirichlet prior over \mathbf{B} is given by:

$$P(\mathbf{B}) = \text{Dir}(\mathbf{b}) \triangleq \prod_{i_1=1}^{|B|_2} \prod_{i_2=1}^{|B|_3} \text{Dir}(\mathbf{b}(i_1, i_2, \cdot)), \quad (153)$$

where $|B|_2$ is the number of values that S_τ can take (i.e., the number of hidden states) and $|B|_3$ is the number of values that U_τ can take (i.e., the number of actions).

Global labels

Definition 33. The **number of actions** available to the agent is denoted $|U|$.

Definition 34. The **number of states** in the environment is denoted $|S|$.

Definition 35. The **number of observations** that the agent can make is denoted $|O|$.

Definition 36. The **number of policies** that the agent can pick from is denoted $|\pi|$.

Definition 37. The time point representing the **present** is a natural number denoted t .

Definition 38. The **time-horizon** (i.e., the time point after which the agent stops modelling the sequence of hidden states) is a natural number denoted T .

Multi-indices

Definition 39. A **multi-index** is a sequence of indices denoted by: $I = (i_1, \dots, i_n)$, where $i_j \in \mathcal{D} \forall j \in \{1, \dots, n\}$, and in this paper $\mathcal{D} = \{1, \dots, |U|\}$.

Remark 40. Multi-indices are used to index random variables such that S_I is the hidden state obtained after taking the sequence of actions described by I , and O_I is the random variable representing the observation generated by S_I .

Definition 41. The **last index** of a multi-index is denoted I_{last} , i.e., I_{last} is the last element of the sequence I .

Definition 42. The **one-hot representation** of the action corresponding to I_{last} is denoted \bar{I}_{last} .

Definition 43. Given a multi-index I , $I \setminus last$ corresponds to the sequence of actions described by I **without the last element**.

Remark 44. In Section 6, when a hidden state (i.e., S_I) is indexed by I , then $S_{I_{last}}$ will be the parent of S_I .

Definition 45. Given an expandable generative model, \mathbb{I}_t is the set of **all multi-indices already expanded** from the current state S_t .

Remark 46. In Section 6 each time a hidden state (i.e., S_I) is added to the generative model, I is added to the set of all multi-indices already expanded \mathbb{I}_t .

Random variables and parameters of their distributions

Remark 47. Parameters of the posterior distributions are recognizable by the hat notation, e.g., $\hat{\mathbf{a}}$, $\hat{\mathbf{b}}$ and $\hat{\mathbf{d}}$ will be posterior parameters, while \mathbf{a} , \mathbf{b} and \mathbf{d} will be prior parameters.

Remark 48. The expected logarithm of an arbitrary tensor \mathbf{X} representing a conditional or a joint distribution is denoted $\bar{\mathbf{X}}$, e.g. $\bar{\mathbf{A}} = \mathbb{E}_{Q(\mathbf{A})}[\ln \mathbf{A}]$ and $\bar{\mathbf{B}} = \mathbb{E}_{Q(\mathbf{B})}[\ln \mathbf{B}]$.

Definition 49. Let U_τ be a random variable taking values in $\{1, \dots, |U|\}$ indexing all possible actions. The prior distribution over U_τ is a categorical distribution represented by Θ_τ . The posterior distribution over U_τ is a categorical distribution represented by $\hat{\Theta}_\tau$.

Definition 50. Let S_τ be a random variable taking values in $\{1, \dots, |S|\}$ indexing all possible states. The prior and posterior distributions over S_τ are categorical distributions represented by different tensors depending on the generative model being considered. Therefore, those distributions are defined in Sections 4 and 6.

Definition 51. Let O_τ be an observed random variable taking values in $\{1, \dots, |O|\}$ indexing all possible observations. The prior distribution over O_τ is a categorical distribution conditioned on S_τ and represented by the random matrix \mathbf{A} . There is no posterior distribution over O_τ because O_τ is observed, i.e., realized.

Definition 52. Let O_I be a random variable taking values in $\{1, \dots, |O|\}$ indexing all possible observations. The prior distribution over O_I is a categorical distribution conditioned on S_I and represented by the random matrix $\bar{\mathbf{A}}$. Note that O_I refers to an observation in the future and is therefore a hidden variable. The posterior distribution over O_I is a categorical distribution represented by $\bar{\mathbf{E}}_I$.

Definition 53. Let S_I be a random variable taking values in $\{1, \dots, |S|\}$ indexing all possible states. The prior distribution over S_I is a categorical distribution conditioned on $S_{I_{last}}$ and represented by the matrix $\bar{\mathbf{B}}_I \triangleq \bar{\mathbf{B}}(\cdot, \cdot, I_{last})$. The posterior distribution over S_I is a categorical distribution represented by $\hat{\mathbf{D}}_I$.

Definition 54. Let \mathbf{A} be a $|O| \times |S|$ random matrix defining the probability of an observation O_τ given the hidden state S_τ . The prior distribution over \mathbf{A} is a product of Dirichlet distributions whose parameters are stored in a $|S| \times |O|$ matrix \mathbf{a} . The posterior distribution over \mathbf{A} is also a product of Dirichlet distribution, but the parameters are stored in a $|S| \times |O|$ matrix $\hat{\mathbf{a}}$.

Definition 55. Let \mathbf{B} be a $|S| \times |S| \times |U|$ random 3-tensor defining the probability of transiting from S_τ to $S_{\tau+1}$ when taking action U_τ . The prior distribution over \mathbf{B} is a product of Dirichlet distributions whose parameters are stored in a $|S| \times |U| \times |S|$ 3-tensor \mathbf{b} . The posterior distribution over \mathbf{B} is also a product of Dirichlet distributions, but the parameters are stored in a $|S| \times |U| \times |S|$ 3-tensor $\hat{\mathbf{b}}$.

Definition 56. Let \mathbf{D} be a random vector of size $|S|$ defining the probability of the initial state S_0 . The prior distribution over \mathbf{D} is a Dirichlet distribution whose parameters are stored in a vector \mathbf{d} of size $|S|$. The posterior distribution over \mathbf{D} is also a Dirichlet distribution, but the parameters are stored in a vector $\hat{\mathbf{d}}$ of size $|S|$.

Definition 57. Let $\Theta_\tau \forall \tau \in \{0, \dots, t-1\}$ be a random vector of size $|U|$ defining the probability of the action U_τ . The prior distribution over Θ_τ is a Dirichlet distribution whose parameters are stored in a vector θ_τ of size $|U|$. The posterior distribution over Θ_τ is also a Dirichlet distribution, but the parameters are stored in a vector $\hat{\theta}_\tau$ of size $|U|$.

Definition 58. Let γ be a random variable taking values in $\mathbb{R}_{>0}$. The prior distribution over γ is a gamma distribution with shape parameter $\alpha = 1$ and rate parameter $\beta \in \mathbb{R}_{>0}$. The posterior distribution over γ is a gamma distribution with shape parameter $\hat{\alpha} = 1$ and rate parameter $\hat{\beta} \in \mathbb{R}_{>0}$.

Definition 59. Let π be a random variable taking values in $\{1, \dots, |\pi|\}$ indexing all possible policies. The prior distribution over π is a softmax function of the vector \mathbf{G} multiplied by minus the precision γ . Note that \mathbf{G} is a vector of size $|\pi|$ whose i th element is the expected free energy of the i th policy. The posterior distribution over π is a categorical distribution whose parameters are stored in a vector $\hat{\pi}$ of size $|\pi|$.

Appendix G. Multi-armed bandit problem

In the multi-armed bandit problem, the agent is prompted with K actions (one for each bandit's arm). Pulling the i th arm returns a reward sampled from the reward distribution $P_i(X)$ associated to this arm. Let μ_i be the mean of the i th reward distribution and $T_i(n)$ be the number of times the i th bandit has been selected after n plays. To solve the bandit problem, one needs to come up with an *allocation strategy* that selects the action that minimizes the agent's regret defined as:

$$R_n = \mu^* n - \sum_{i=1}^K \mu_i \mathbb{E}[T_i(n)], \quad (154)$$

where μ^* is the average reward of the best action. Note that an upper bound of $\mathbb{E}[T_i(n)]$ is derived by first upper bounding $T_i(n)$, and then using: the Chernoff-Hoeffding bound, the Bernstein inequality and some properties of p -series, cf., proof of Theorem 1 in Auer et al. (2002) for details. So, if we first assume that,

$$UCB1_i = \underbrace{\bar{X}_i}_{\text{exploitation}} + \underbrace{\sqrt{\frac{2 \ln n}{n_i}}}_{\text{exploration}}, \quad (155)$$

where n_i is the number of times the i th action has been selected, and \bar{X}_i is the average reward received after taking the i th action. Then, the main result of Auer et al. (2002) was to show that if an allocation strategy was using the UCB1 criterion to select the next action, the expected regret of this allocation strategy will grow at most logarithmically in the number of plays n , i.e., $\mathcal{O}(\ln n)$. In addition, since it is known that the expected regret of the (best) allocation strategy grows at least logarithmically in n (Lai & Robbins, 1985), we say that the UCB1 criterion resolves the exploration/exploitation trade-off, i.e., the UCB1 criterion ensures that the expected regret grows as slowly as possible.

Appendix H. The exponential complexity class

In this appendix, we precisely pinpoint the exponential complexity class that is addressed in this paper, but first, we introduce a multi-index notation. Multi-indices will help us to refer to hidden states in the future. Naturally enough the indices inside the multi-indices will correspond to the actions the agent will have to perform to reach the hidden state, e.g., $S_{(123)}$ corresponds to a hidden state at time $t+3$ obtained by performing action 1 at time

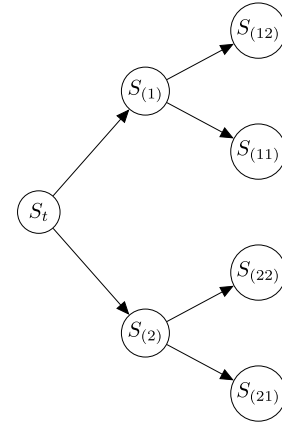


Fig. 10. Illustration of all possible policies up to two time steps in the future when $|U| = 2$. The state at the current time step is denoted by S_t . Additionally, each branch of the tree corresponds to a possible policy, and each node S_l is indexed by a multi-index (e.g. $l = (12)$) representing the sequence of actions that led to this state. This should make it clear that for one time step in the future, there are $|U|$ possible policies, after two time steps there are $|U|$ times more policies, and so on until the time-horizon T where there are a total of $|U|^T$ possible policies, i.e., the number of possible policies grows exponentially with the number of time steps for which the agent tries to plan.

$t, 2$ at time $t+1$ and 3 at time $t+2$. Using this notation, Fig. 10 depicts all the possible policies up to two time steps in the future and the associated hidden states. Importantly, Fig. 10 shows that the number of policies grows exponentially with the number of time steps for which the agent tries to plan. Therefore, the definition of the prior over the policies, i.e., $P(\pi|\gamma) = \sigma(-\gamma\mathbf{G})$, exhibits an exponential space and time complexity class because the agent needs to store and compute the $|U|^T$ parameters of $P(\pi|\gamma)$, where T is the time-horizon.

To show that this exponential explosion is not only a theoretical problem and also appears in practice, we modified the *DEMO_MDP_maze.m* of the SPM¹ package in two ways. First, we allowed the agent to plan N time steps in the future as follows:

```

1 % V = Allowable policies (N moves into the future),
  nu = number of actions
2 V = [];
3 for i1 = 1:nu, i2 = 1:nu, ..., iN = 1:nu
4     V(:,end + 1) = [i1;i2; ... ;iN];
5 end

```

Second, we used the “tic” and “toc” functions provided by Matlab to track the execution time required to execute the function “*spm_maze_search*” where the parameters of the model are assumed to be known already by the agent:

¹ Statistical parametric mapping (SPM) is a software package created by the Wellcome Department of Imaging Neuroscience at University College London, which was initially developed to carry out statistical analyses of functional neuroimaging data. Today, SPM also contains Matlab simulations implementing active inference agents (among other things), cf. <https://www.fil.ion.ucl.ac.uk/spm/>.

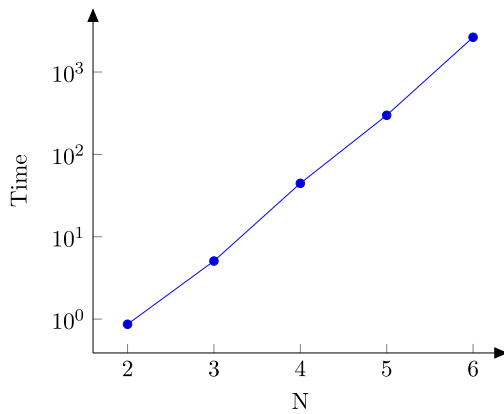


Fig. 11. This figure shows the time required to execute the function “spm_maze_search” when the agent is allowed to plan N time steps in the future (for N from 2 to 6). A logarithmic scale is used on the time axis.

```

1 tic % Start a timer used to evaluate the execution
  time of spm_maze_search
2 MDP = spm_maze_search(mdp,8,START,END,128,1);
3 toc % Display the time elapsed since the last call to tic

```

Fig. 11 presents the results of our simulations for N from 2 to 6. Under a logarithmic scale on the time axis, the experimental results show that the graph is almost a perfect line, which provides empirical evidence for an exponential time explosion. Note that the simulation for $N = 7$ crashed after trying to allocate an array of 9.5 GB (space explosion). In Section 5, we presented an approach proposed to deal with the exponential complexity class that arises during planing, yet is fundamentally similar to active inference. This effectively means that our paper shows how standard active inference can be made more efficient and scale to longer time horizons.

References

- Auer, P., Cesa-Bianchi, N., & Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2), 235–256. <http://dx.doi.org/10.1023/A:1013689704352>, <https://doi.org/10.1023/A:1013689704352>.
- Bai, A., Wu, F., & Chen, X. (2013). Bayesian mixture modelling and inference based Thompson sampling in Monte-Carlo tree search. In *Proceedings of the advances in neural information processing systems (NIPS)* (pp. 1646–1654). Lake Tahoe, United States.
- Bergstra, J., Bardenet, R., Bengio, Y., & Kégl, B. (2011). Algorithms for hyperparameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, & K. Q. Weinberger (Eds.), *Advances in neural information processing systems* (Vol. 24). Curran Associates, Inc., URL <https://proceedings.neurips.cc/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf>.
- Botvinick, M., & Toussaint, M. (2012). Planning as inference. *Trends in Cognitive Sciences*, 16(10), 485–488, <https://doi.org/10.1016/j.tics.2012.08.006>.
- Bowman, H. (2005). *Concurrency theory: Calculi an automata for modelling untimed and timed concurrent systems*. Dordrecht: Springer, URL <https://cds.cern.ch/record/1250124>.
- Bowman, H., & Wyble, B. (2007). The simultaneous type, serial token model of temporal attention and working memory. *Psychological Review*, 114(1), 38–70.
- Brochu, E., Cora, V. M., & de Freitas, N. (2010). A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. Citeseer, ArXiv preprint [arXiv:1012.2599](https://arxiv.org/abs/1012.2599).
- Browne, C. B., Powley, E., Whitehouse, D., Lucas, S. M., Cowling, P. I., Rohlfshagen, P., et al. (2012). A survey of Monte Carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(1), 1–43.
- Çatal, O., Verbelen, T., Nauta, J., Boom, C. D., & Dhoedt, B. (2020). Learning perception and planning with deep active inference. In *2020 IEEE international conference on acoustics, speech and signal processing, ICASSP 2020, Barcelona, Spain, May 4–8, 2020* (pp. 3952–3956). IEEE, <http://dx.doi.org/10.1109/ICASSP40776.2020.9054364>, <https://doi.org/10.1109/ICASSP40776.2020.9054364>.
- Champion, T., Bowman, H., & Grzes, M. (2021). Branching time active inference: empirical study. Available at <https://arxiv.org/abs/2111.11276>.
- Champion, T., Grzes, M., & Bowman, H. (2021). Realizing active inference in variational message passing: The outcome-blind certainty seeker. *Neural Computation*, 1–65. http://dx.doi.org/10.1162/neco_a_01422, https://doi.org/10.1162/neco_a_01422, arXiv:https://direct.mit.edu/neco/article-pdf/doi/10.1162/neco_a_01422/1930278/neco_a_01422.pdf.
- Cox, M., van de Laar, T., & de Vries, B. (2019). A factor graph approach to automated design of Bayesian signal processing algorithms. *International Journal of Approximate Reasoning*, 104, 185–204. <http://dx.doi.org/10.1016/j.ijar.2018.11.002>, <https://doi.org/10.1016/j.ijar.2018.11.002>.
- Cullen, M., Davey, B., Friston, K. J., & Moran, R. J. (2018). Active inference in openai gym: A paradigm for computational investigations into psychiatric illness. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 3(9), 809–818, Computational Methods and Modeling in Psychiatry. <https://doi.org/10.1016/j.bpsc.2018.06.010>. URL <http://www.sciencedirect.com/science/article/pii/S2451902218301617>.
- Da Costa, L., Parr, T., Sajid, N., Veselic, S., Neacsu, V., & Friston, K. (2020). Active inference on discrete state-spaces: A synthesis. *Journal of Mathematical Psychology*, 99, Article 102447, <https://doi.org/10.1016/j.jmp.2020.102447>. URL <https://www.sciencedirect.com/science/article/pii/S0022249620300857>.
- Da Costa, L., Sajid, N., Parr, T., Friston, K., & Smith, R. (2020). The relationship between dynamic programming and active inference: the discrete, finite-horizon case. [arXiv:2009.08111](https://arxiv.org/abs/2009.08111).
- FitzGerald, T. H. B., Dolan, R. J., & Friston, K. (2015). Dopamine, reward learning, and active inference. *Frontiers in Computational Neuroscience*, 9, 136. <http://dx.doi.org/10.3389/fncom.2015.00136>, URL <https://www.frontiersin.org/article/10.3389/fncom.2015.00136>.
- Forney, G. D. (2001). Codes on graphs: normal realizations. *IEEE Transactions on Information Theory*, 47(2), 520–548.
- Fountas, Z., Sajid, N., Mediano, P. A. M., & Friston, K. J. (2020). Deep active inference agents using Monte-Carlo methods. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, & H. Lin (Eds.), *Advances in neural information processing systems 33: Annual conference on neural information processing systems 2020, NeurIPS 2020, December 6–12, 2020, Virtual*. URL <https://proceedings.neurips.cc/paper/2020/hash/865dfbde8a344b44095495f3591f7407-Abstract.html>.
- Fox, V., Hightower, J., Liao, L., Schulz, D., & Borriello, G. (2003). Bayesian filtering for location estimation. *IEEE Pervasive Computing*, 2(3), 24–33. <http://dx.doi.org/10.1109/MPRV.2003.1228524>.
- Fox, C. W., & Roberts, S. J. (2012). A tutorial on variational Bayesian inference. *Artificial Intelligence Review*, 38(2), 85–95. <http://dx.doi.org/10.1007/s10462-011-9236-8>, <https://doi.org/10.1007/s10462-011-9236-8>.
- Friston, K., Da Costa, L., Hafner, D., Hesp, C., & Parr, T. (2021). Sophisticated inference. *Neural Computation*, 33(3), 713–763. http://dx.doi.org/10.1162/neco_a_01351, https://doi.org/10.1162/neco_a_01351. arXiv:https://direct.mit.edu/neco/article-pdf/33/3/713/1889421/neco_a_01351.pdf.
- Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., Doherty, J. O., & Pezzulo, G. (2016). Active inference and learning. *Neuroscience & Biobehavioral Reviews*, 68, 862–879, <https://doi.org/10.1016/j.neubiorev.2016.06.022>.
- Friston, K. J., Parr, T., & de Vries, B. (2017). The graphical brain: Belief propagation and active inference. *Network Neuroscience*, 1(4), 381–414. http://dx.doi.org/10.1162/NETN_a_00018, https://doi.org/10.1162/NETN_a_00018.
- Friston, K., Schwartenbeck, P., Fitzgerald, T., Moutoussis, M., Behrens, T., & Dolan, R. (2013). The anatomy of choice: active inference and agency. *Frontiers in Human Neuroscience*, 7, 598. <http://dx.doi.org/10.3389/fnhum.2013.00598>, URL <https://www.frontiersin.org/article/10.3389/fnhum.2013.00598>.
- Haarnoja, T., Zhou, A., Abbeel, P., & Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. CoRR [arXiv:1801.01290](https://arxiv.org/abs/1801.01290).
- Hafner, D., Ortega, P. A., Ba, J., Parr, T., Friston, K. J., & Heess, N. (2020). Action and perception as divergence minimization. CoRR [arXiv:2009.01791](https://arxiv.org/abs/2009.01791).
- Itti, L., & Baldi, P. (2009). Bayesian surprise attracts human attention. *Vision Research*, 49(10), 1295–1306, Visual attention: Psychophysics, electrophysiology and neuroimaging. <https://doi.org/10.1016/j.visres.2008.09.007> URL <http://www.sciencedirect.com/science/article/pii/S0042698908004380>.
- Kirchhoff, M., Parr, T., Palacios, E., Friston, K., & Kiverstein, J. (2018). The Markov blankets of life: autonomy, active inference and the free energy principle. *Journal of the Royal Society Interface*, 15(138), Article 20170792. <http://dx.doi.org/10.1098/rsif.2017.0792>, arXiv:<https://royalsocietypublishing.org/doi/pdf/10.1098/rsif.2017.0792>. URL <https://royalsocietypublishing.org/doi/abs/10.1098/rsif.2017.0792>.
- Kocsis, L., & Szepesvári, C. (2006). Bandit based Monte-Carlo planning. In J. Fürnkranz, T. Scheffer, & M. Spiliopoulou (Eds.), *Lecture notes in computer science: Vol. 4212, Machine learning: ECML 2006, 17th European conference on machine learning, Berlin, Germany, September 18–22, 2006, proceedings* (pp. 282–293). Springer, http://dx.doi.org/10.1007/11871842_29, https://doi.org/10.1007/11871842_29.

- Lai, T., & Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1), 4–22. [http://dx.doi.org/10.1016/0196-8858\(85\)90002-8](http://dx.doi.org/10.1016/0196-8858(85)90002-8), [https://doi.org/10.1016/0196-8858\(85\)90002-8](https://doi.org/10.1016/0196-8858(85)90002-8).
- Lample, G., & Chaplot, D. S. (2017). Playing FPS games with deep reinforcement learning. In S. P. Singh, & S. Markovitch (Eds.), *Proceedings of the thirty-first AAAI conference on artificial intelligence*, February 4–9, 2017, San Francisco, California, USA (pp. 2140–2146). AAAI Press, URL <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14456>.
- Maisto, D., Gregoratti, F., Friston, K. J., & Pezzulo, G. (2021). Active tree search in large POMDPs. *CoRR* [arXiv:2103.13860](https://arxiv.org/abs/2103.13860).
- Millidge, B. (2019). Combining active inference and hierarchical predictive coding: A tutorial introduction and case study. <https://doi.org/10.31234/osf.io/kf6wc>.
- Millidge, B. (2020). Deep active inference as variational policy gradients. *Journal of Mathematical Psychology*, 96, Article 102348, <https://doi.org/10.1016/j.jmp.2020.102348>. URL <http://www.sciencedirect.com/science/article/pii/S0022249620300298>.
- Millidge, B., Tschantz, A., & Buckley, C. L. (2021). Whence the expected free energy? *Neural Computation*, 33(2), 447–482. http://dx.doi.org/10.1162/neco_a_01354, https://doi.org/10.1162/neco_a_01354. [arXiv:https://direct.mit.edu/neco/article-pdf/33/2/447/1896836/neco_a_01354.pdf](https://direct.mit.edu/neco/article-pdf/33/2/447/1896836/neco_a_01354.pdf).
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., et al. (2013). Playing atari with deep reinforcement learning. *CoRR* [arXiv:1312.5602](https://arxiv.org/abs/1312.5602).
- O'Reilly, R. C. (1996). Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm. *Neural Computation*, 8(5), 895–938. <http://dx.doi.org/10.1162/neco.1996.8.5.895>, <https://doi.org/10.1162/neco.1996.8.5.895>.
- Parr, T., & Friston, K. J. (2019). Generalised free energy and active inference. *Biological Cybernetics*, 113(5), 495–513. <http://dx.doi.org/10.1007/s00422-019-00805-w>, <https://doi.org/10.1007/s00422-019-00805-w>.
- Pezzato, C., Corbato, C. H., & Wisse, M. (2020). Active inference and behavior trees for reactive action planning and execution in robotics. *CoRR* [arXiv:2011.09756](https://arxiv.org/abs/2011.09756).
- Rafetseder, E., Schwitalla, M., & Perner, J. (2013). Counterfactual reasoning: From childhood to adulthood. *Journal of Experimental Child Psychology*, 114(3), 389–404.
- Rosin, C. D. (2010). Multi-armed bandits with episode context. In *International symposium on artificial intelligence and mathematics, ISAIM 2010, Fort Lauderdale, Florida, USA, January 6–8, 2010*. URL <http://gauss.ececs.uc.edu/Workshops/isaim2010/papers/rosin.pdf>.
- Russo, D. J., Van Roy, B., Kazerouni, A., Osband, I., & Wen, Z. (2018). A tutorial on thompson sampling. *Foundation in Trends Machine Learning*, 11(1), 1–96. <http://dx.doi.org/10.1561/22000000070>, <https://doi.org/10.1561/22000000070>.
- Sajid, N., Ball, P. J., Parr, T., & Friston, K. J. (2021). Active inference: Demystified and compared. *Neural Computation*, 33(3), 674–712. http://dx.doi.org/10.1162/neco_a_01357, https://doi.org/10.1162/neco_a_01357. [arXiv:https://direct.mit.edu/neco/article-pdf/33/3/674/1889396/neco_a_01357.pdf](https://direct.mit.edu/neco/article-pdf/33/3/674/1889396/neco_a_01357.pdf).
- Sancaktar, C., van Gerven, M. A. J., & Lanillos, P. (2020). End-to-end pixel-based deep active inference for body perception and action. In *Joint IEEE 10th international conference on development and learning and epigenetic robotics, ICDL-EpiRob 2020, Valparaiso, Chile, October 26–30, 2020* (pp. 1–8). IEEE, <http://dx.doi.org/10.1109/ICDL-EpiRob48136.2020.9278105>, <https://doi.org/10.1109/ICDL-EpiRob48136.2020.9278105>.
- Sancaktar, C., & Lanillos, P. (2020). End-to-end pixel-based deep active inference for body perception and action. [arXiv:2001.05847](https://arxiv.org/abs/2001.05847).
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., et al. (2019). Mastering atari, go, chess and shogi by planning with a learned model. [arXiv:1911.08265](https://arxiv.org/abs/1911.08265).
- Schwartenbeck, P., Passecker, J., Hauser, T. U., FitzGerald, T. H. B., Kronbichler, M., & Friston, K. (2018). Computational mechanisms of curiosity and goal-directed exploration. *BioRxiv*, <http://dx.doi.org/10.1101/411272>, URL <https://www.biorxiv.org/content/early/2018/09/07/411272>.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., et al. (2016). Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587), 484–489. <http://dx.doi.org/10.1038/nature16961>, <https://doi.org/10.1038/nature16961>.
- Smith, R., Friston, K. J., & Whyte, C. J. (2021). A step-by-step tutorial on active inference and its application to empirical data. URL <https://psyarxiv.com/b4jm6/>.
- Sondik, E. J. (1971). *The optimal control of partially observable Markov processes* (Ph.D. thesis), Stanford University, URL <https://ci.nii.ac.jp/naid/20000916958/en/>.
- Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4), 285–294, URL <http://www.jstor.org/stable/2332286>.
- van de Laar, T., & de Vries, B. (2019). Simulating active inference processes by message passing. *Frontiers Robotics and AI*, 2019, <http://dx.doi.org/10.3389/frobt.2019.00020>, <https://doi.org/10.3389/frobt.2019.00020>.
- van Glabbeek, R. J. (1990). The linear time-branching time spectrum (extended abstract). In *Proceedings of the theories of concurrency: Unification and extension* (pp. 278–297). Berlin, Heidelberg: Springer-Verlag.
- van Glabbeek, R. J. (1993). The linear time – Branching time spectrum II. In E. Best (Ed.), *CONCUR'93* (pp. 66–81). Berlin, Heidelberg: Springer Berlin Heidelberg.
- van Hasselt, H., Guez, A., & Silver, D. (2016). Deep reinforcement learning with double Q-learning. In D. Schuurmans, & M. P. Wellman (Eds.), *Proceedings of the thirtieth AAAI conference on artificial intelligence, February 12–17, 2016, Phoenix, Arizona, USA* (pp. 2094–2100). AAAI Press, URL <http://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/12389>.
- Winn, J., & Bishop, C. (2005). Variational message passing. *Journal of Machine Learning Research*, 6, 661–694.
- Wyble, B., & Bowman, H. (2006). A neural network account of binding discrete items into working memory using a distributed pool of flexible resources. *Journal of Vision*, 6(6), 33–33a.