



# Simulating Emotions: An Active Inference Model of Emotional State Inference and Emotion Concept Learning

Ryan Smith<sup>1\*</sup>, Thomas Parr<sup>2</sup> and Karl J. Friston<sup>2</sup>

<sup>1</sup> Laureate Institute for Brain Research, Tulsa, OK, United States, <sup>2</sup> Wellcome Centre for Human Neuroimaging, Institute of Neurology, University College London, London, United Kingdom

## OPEN ACCESS

### Edited by:

Alessia Celegghin,  
University of Turin, Italy

### Reviewed by:

Mauro Ursino,  
University of Bologna, Italy  
Carlos Andrés Méndez,  
University of Turin, Italy

### \*Correspondence:

Ryan Smith  
rsmith@laureateinstitute.org

### Specialty section:

This article was submitted to  
Consciousness Research,  
a section of the journal  
Frontiers in Psychology

**Received:** 02 August 2019

**Accepted:** 02 December 2019

**Published:** 19 December 2019

### Citation:

Smith R, Parr T and Friston KJ  
(2019) Simulating Emotions: An Active  
Inference Model of Emotional State  
Inference and Emotion Concept  
Learning. *Front. Psychol.* 10:2844.  
doi: 10.3389/fpsyg.2019.02844

The ability to conceptualize and understand one's own affective states and responses – or “Emotional awareness” (EA) – is reduced in multiple psychiatric populations; it is also positively correlated with a range of adaptive cognitive and emotional traits. While a growing body of work has investigated the neurocognitive basis of EA, the neurocomputational processes underlying this ability have received limited attention. Here, we present a formal Active Inference (AI) model of emotion conceptualization that can simulate the neurocomputational (Bayesian) processes associated with learning about emotion concepts and inferring the emotions one is feeling in a given moment. We validate the model and inherent constructs by showing (i) it can successfully acquire a repertoire of emotion concepts in its “childhood”, as well as (ii) acquire new emotion concepts in synthetic “adulthood,” and (iii) that these learning processes depend on early experiences, environmental stability, and habitual patterns of selective attention. These results offer a proof of principle that cognitive-emotional processes can be modeled formally, and highlight the potential for both theoretical and empirical extensions of this line of research on emotion and emotional disorders.

**Keywords:** emotion concepts, trait emotional awareness, learning, computational neuroscience, active inference

## INTRODUCTION

The ability to conceptualize and understand one's affective responses has become the topic of a growing body of empirical work (McRae et al., 2008; Smith et al., 2015, 2017b,c, 2018c,d,e, 2019a,c; Wright et al., 2017). This body of work has also given rise to theoretical models of its underlying cognitive and neural basis (Wilson-Mendenhall et al., 2011; Lane et al., 2015; Smith and Lane, 2015, 2016; Barrett, 2017; Kleckner et al., 2017; Panksepp et al., 2017; Smith et al., 2018b). Attempts to operationalize this cognitive-emotional ability have led to a range of overlapping constructs, including trait emotional awareness (Lane and Schwartz, 1987), emotion differentiation or granularity (Kashdan and Farmer, 2014; Kashdan et al., 2015), and alexithymia (Bagby et al., 1994a,b).

This work is motivated to a large degree by the clinical relevance of emotion conceptualization abilities. In the literature on the construct of emotional awareness, for example, lower levels of conceptualization ability have been associated with several psychiatric disorders as well as poorer physical health (Levine et al., 1997; Berthoz et al., 2000; Bydłowski et al., 2005; Donges et al., 2005; Lackner, 2005; Subic-Wrana et al., 2005, 2007; Frewen et al., 2008; Baslet et al., 2009; Consoli et al., 2010; Moeller et al., 2014); conversely, higher ability levels have been associated with a range of adaptive emotion-related traits and abilities (Lane et al., 1990, 1996, 2000; Ciarrochi et al., 2003; Barchard and Hakstian, 2004; Bréjard et al., 2012). Multiple evidence-based psychotherapeutic modalities also aim to improve emotion understanding as a central part of psycho-education in psychotherapy (Hayes and Smith, 2005; Barlow et al., 2016).

While there are a number of competing views on the nature of emotions, most (if not all) accept that emotion concepts must be acquired through experience. For example, “basic emotions” theories hold that emotion categories like sadness and fear each have distinct neural circuitry, but do not deny that knowledge about these emotions must be learned (Panksepp and Biven, 2012). Constructivist views instead hold that emotion categories do not have a 1-to-1 relationship to distinct neural circuitry, and that emotion concept acquisition is necessary for emotional experience (Barrett, 2017). While these views focus on understanding the nature of emotions themselves, we have recently proposed a neurocognitive model – termed the “three-process model” (TPM; Smith et al., 2018b, 2019a; Smith, 2019) of emotion episodes – with a primary focus on accounting for individual differences in emotional awareness. This model characterizes a range of emotion-related processes that could contribute to trait differences in both the learning and deployment of emotion concepts in order to understand one’s own affective responses (and in the subsequent use of these concepts to guide adaptive decision-making). The TPM distinguishes the following three broadly defined processes (see **Figure 1**):

1. **Affective response generation:** a process in which somatovisceral and cognitive states are automatically modulated in response to an affective stimulus (whether real, remembered, or imagined) in a context-dependent manner, based on an (often implicit) appraisal of the significance of that stimulus for the survival and goal-achievement of the individual (i.e., predictions about the cognitive, metabolic, and behavioral demands of the situation).
2. **Affective response representation:** a process in which the somatovisceral component of an affective response is subsequently perceived via afferent sensory processing, and then conceptualized as a particular emotion (e.g., sadness, anger, etc.) in consideration of all other available sources of information (e.g., stimulus/context information, current thoughts/beliefs about the situation, etc.).
3. **Conscious access:** a process in which the representations of somatovisceral percepts and emotion concept

representations may or may not enter and be held in working memory – constraining the use of this information in goal-directed decision-making (e.g., verbal reporting, selection of voluntary emotion regulation strategies, etc.).

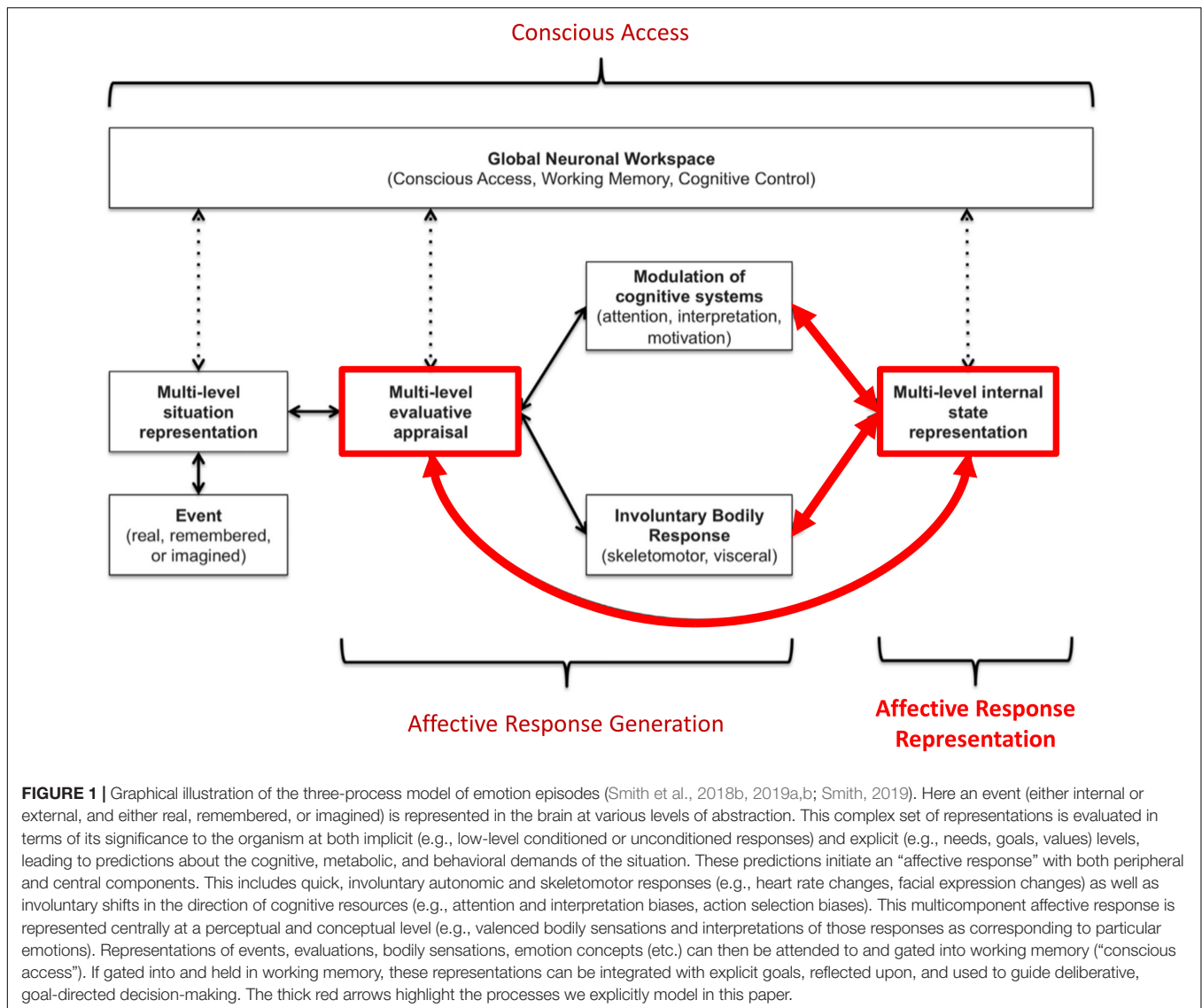
The TPM has also proposed a tentative mapping to the brain in terms of interactions between large-scale neural networks serving domain-general cognitive functions. Some support for this proposal has been found within recent neuroimaging studies (Smith et al., 2017b,c, 2018a,c,d,e). However, the neurocomputational implementation of these processes has not been thoroughly considered. The computational level of description offers the promise of providing more specific and mechanistic insights, which could potentially be exploited to inform and improve pharmacological and psychotherapeutic interventions. While previous theoretical work has applied active inference concepts to emotional phenomena (Joffily and Coricelli, 2013; Seth, 2013; Barrett and Simmons, 2015; Seth and Friston, 2016; Smith et al., 2017d, 2019e; Clark et al., 2018), no formal modeling of emotion concept learning has yet been performed. In this manuscript, we aim to take the first steps in constructing an explicit computational model of the acquisition and deployment of emotion concept knowledge (i.e., affective response representation) as described within the TPM (subsequent work will focus on affective response generation and conscious access processes; see Smith et al., 2019b). Specifically, we present a simple Active Inference model (Friston et al., 2016, 2017a) of emotion conceptualization, formulated as a Markov Decision Process. We then outline some initial insights afforded by simulations using this model.

In what follows, we first provide a brief review of active inference. We will place a special emphasis on deep generative models that afford the capacity to explain multimodal (i.e., interoceptive, proprioceptive, and exteroceptive) sensations that are characteristic of emotional experience. We then introduce a particular model of emotion inference that is sufficiently nuanced to produce synthetic emotional processes but sufficiently simple to be understood from a “first principles” account. We then establish the validity of this model using numerical analyses of emotion concept learning during (synthetic) neurodevelopment. We conclude with a brief discussion of the implications of this work; particularly for future applications.

## AN ACTIVE INFERENCE MODEL OF EMOTION CONCEPTUALIZATION

### A Primer on Active Inference

Active Inference (AI) starts from the assumption that the brain is an inference machine that approximates optimal probabilistic (Bayesian) belief updating across all biopsychological domains (e.g., perception, decision-making, motor control, etc.). AI postulates that the brain embodies an internal model of the world (including the body) that is “generative” in the sense that it is able to simulate the sensory data that it should receive if its model of the world is correct. This simulated (predicted) sensory data



can be compared to actual observations, and deviations between predicted and observed sensations can then be used to update the model. On short timescales (e.g., a single trial in a perceptual decision-making task) this updating corresponds to perception, whereas on longer timescales it corresponds to learning (i.e., updating expectations about what will be observed on subsequent trials). One can see these processes as ensuring the generative model (embodied by the brain) remains an accurate model of the world (Conant and Ashbey, 1970).

Action (be it skeletal motor, visceromotor, or cognitive action) can be cast in similar terms. For example, actions can be chosen to resolve uncertainty about variables within a generative model (i.e., sampling from domains in which the model does not make precise predictions). This can prevent future deviations from predicted outcomes. In addition, the brain must continue to make certain predictions simply in order to survive. For example, if the brain did not in some sense continue to “expect” to observe certain amounts of food, water, shelter, social support, and a

range of other quantities, then it would cease to exist (McKay and Dennett, 2009); as it would not pursue those behaviors leading to the realization of these expectations [c.f. “the optimism bias” (Sharot, 2011)]. Thus, there is a deep sense in which the brain must continually seek out observations that support – or are internally consistent with – its own continued existence. As a result, decision-making can be cast as a process in which the brain infers the sets of actions (policies) that would lead to observations most consistent with its own survival-related expectations (i.e., its “prior preferences”). Mathematically, this can be described as selecting policies that maximize a quantity called “Bayesian model evidence” – that is, the probability that sensory data would be observed under a given model. In other words, because the brain is itself a model of the world, action can be understood as a process by which the brain seeks out evidence for itself – sometimes known as self-evidencing (Hohwy, 2016).

In a real-world setting, directly computing model evidence becomes mathematically intractable. Thus, the brain must

use some approximation. AI proposes that the brain instead computes a statistical quantity called free energy. Unlike model evidence, computing free energy is mathematically tractable. Crucially, this quantity provides a bound on model evidence, such that minimization of free energy is equivalent to maximizing model evidence. By extension, in decision-making an agent can evaluate the *expected* free energy of the alternative policies she could select – that is, the free energy of future trajectories under each policy (i.e., based on predicted future outcomes, given the future states that would be expected under each policy). Therefore, decision-making will be approximately (Bayes) optimal if it operates by inferring (and enacting) the policy that minimizes expected free energy – and thereby maximizes evidence for the brain's internal model. Interestingly, expected free energy can be decomposed into terms reflecting uncertainty and prior preferences, respectively. This decomposition explains why agents that minimize expected free energy will first select exploratory policies that minimize uncertainty in a new environment (often called the “epistemic value” component of expected free energy). Once uncertainty is resolved, the agent then selects policies that exploit that environment to maximize her prior preferences (often called the “pragmatic value” component of expected free energy). The formal mathematical basis for AI has been detailed elsewhere (Friston et al., 2017a), and the reader is referred there for a full mathematical treatment (also see **Figure 2** for some additional detail).

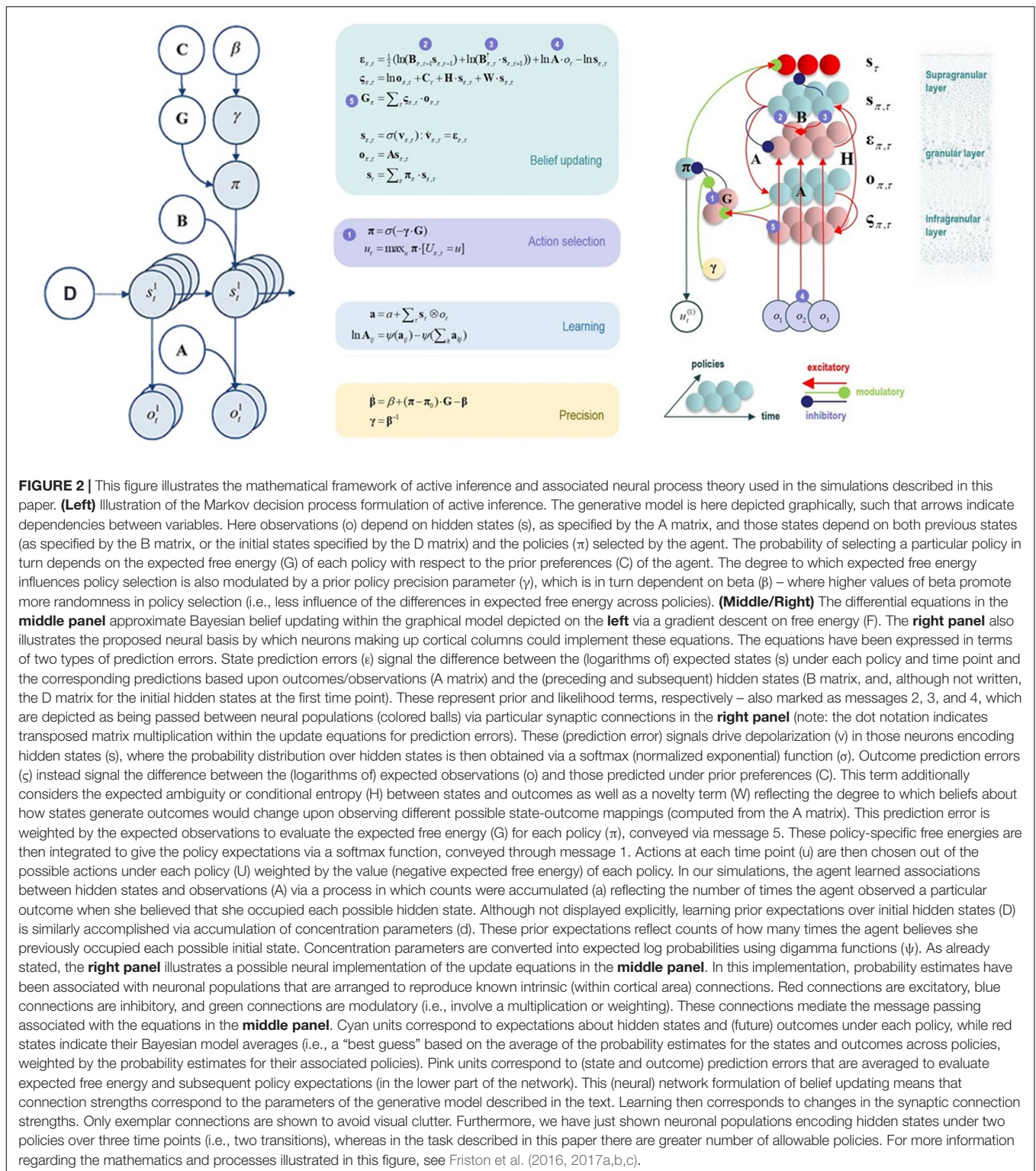
When a generative model is formulated as a partially observable Markov decision process, active inference takes a particular form. Specifically, specifying a generative model in this context requires specifying the allowable policies, hidden states of the world (that the brain cannot directly observe but must infer), and observable outcomes, as well as a number of matrices that define the probabilistic relationships between these quantities (see **Figure 2**). The “A” matrix specifies which outcomes are generated by each combination of hidden states (i.e., a likelihood mapping indicating the probability that a particular set of outcomes would be observed given a particular hidden state). The “B” matrix encodes state transitions, specifying the probability that one hidden state will evolve into another over time. Some of these transitions are controlled by the agent, according to the policy that has been selected. The “D” matrix encodes prior expectations about the initial hidden state of the world. The “E” matrix encodes prior expectations about which policies will be chosen (e.g., frequently repeated habitual behaviors will have higher prior expectation values). Finally, the “C” matrix encodes prior preferences over outcomes. Outcomes and hidden states are generally factorized into multiple outcome *modalities* and hidden state *factors*. This means that the likelihood mapping (the “A” matrix) plays an important role in modeling the interactions among different hidden states at each level of a hierarchical model when generating the outcomes at the level below. One can think of each factor and modality as an independent group of competing states or observations within a given category. For example, one hidden state factor could be “birds,” which includes competing interpretations of sensory input as corresponding to either hawks, parrots, or

pigeons, whereas a separate factor could be “location,” with competing representations of where a bird is in the sky. Similarly, one outcome modality could be size (e.g., is the bird big or small?) whereas another could be color (is the bird black, white, green, etc?).

As shown in the middle and right panels of **Figure 2**, active inference is also equipped with a neural process theory – a proposed manner in which neuronal circuits and their dynamics can invert generative models via a set of linked update equations that minimize prediction errors. In this neuronal implementation, the probability of neuronal firing in specific populations is associated with the expected probability of a state, whereas postsynaptic membrane potentials are associated with the logarithm of this probability. A softmax function acts as an activation function – transforming membrane potentials into firing rates. With this setup, postsynaptic depolarizations (driven by ascending signals) can be understood as prediction errors (free energy gradients) about hidden states – arising from linear mixtures in the firing rates of other neural populations. These prediction errors (postsynaptic currents) in turn drive membrane potential changes (and resulting firing rates). When prediction errors are minimized, postsynaptic influences no longer drive changes in activity (depolarizations and firing rates), corresponding to minimum free energy.

Via similar dynamics, prediction errors about outcomes (i.e., the deviation between preferred outcomes and those predicted under each policy) can also be computed and integrated (i.e., averaged) to evaluate the expected free energy (value) of each policy (i.e., underwriting selection of the policy that best minimizes these prediction errors). Dopamine dynamics also modulate policy selection, by encoding estimates of the expected uncertainty over policies – where greater expected uncertainty promotes less deterministic policy selection. Phasic dopamine responses correspond to updates in expected uncertainty over policies – which occur when there is a prediction error about expected free energy; that is, when there is a difference between the expected free energy of policies before and after a new observation.

Finally, and most centrally for the simulations we report below, learning in this theory corresponds to a form of synaptic plasticity remarkably similar to Hebbian coincidence-based learning mechanisms associated with empirically observed synaptic long-term potentiation and depression (LTP and LTD) processes (Brown et al., 2009). Here one can think of the strength of each synaptic connection as a parameter in one of the matrices described above. For example, the strength of one synapse could encode the amount of evidence a given observation provides for a given hidden state (i.e., an entry in the “A” matrix), whereas another synapse could encode the probability of a state at a later time given a state at an earlier time (i.e., an entry in the “B” matrix). Mathematically, the synaptic strengths correspond to Dirichlet parameters that increase in value in response to new observations. One can think of this process as adding counts to each matrix entry based on coincidences in pre- and post-synaptic activity. For example, if beliefs favor one hidden state, and this co-occurs with a specific observation, then the strength of the value in the “A” matrix encoding



the relationship between that state and that observation will increase. Counts also increase in similar fashion in the “D” matrix encoding prior beliefs about initial states (whenever a given hidden state is inferred at the start of a trial) as well as in the “B” matrix encoding beliefs about transition probabilities

(whenever one specific state is followed by another). For a more detailed discussion, please see the legend for **Figure 2** and associated references.

In what follows, we describe how this type of generative model was specified to perform emotional state inference and emotion

concept learning. We also present simulated neural responses based on the neural process theory described above.

## A Model of Emotion Inference and Concept Learning

In this paper we focus on the second process in the TPM – affective response representation – in which a multifaceted affective response is generated and the ensuing (exteroceptive, proprioceptive, and interoceptive) outcomes are used to infer or represent the current emotional state. The basic idea is to equip the generative model with a space of emotion concepts (i.e., latent or hidden states) that generate the interoceptive, exteroceptive and proprioceptive consequences (at various levels of abstraction) of being in a particular emotional state. Inference under this model then corresponds to inferring that one of several possible emotion concepts is the best explanation for the data at hand (e.g., “my unpleasant feeling of increased heart rate and urge to run away must indicate that I am afraid to give this speech”). Crucially, to endow emotion concept inference with a form of mental action (Metzinger, 2017; Limanowski and Friston, 2018), we also included a state factor corresponding to *selective attention*. Transitions between attentional states were under control of the agent (i.e., “B” matrices were specified for all possible transitions between these states). The “A” matrix mapping emotion concepts to lower-level observations differed in each attentional state, such that precise information about each type of lower-level information was only available in one attentional state (e.g., the agent needed to transition into the “attention to valence” state to gain precise information regarding whether she was feeling pleasant or unpleasant, and so forth; see **Figure 3C**).

The incorporation of selective attention in emotional state inference and learning within our model was motivated by several factors. First, multiple psychotherapeutic modalities improve clients’ understanding of their own emotions in just this way; that is, by having them selectively attend to and record the contexts, bodily sensations, thoughts, action tendencies, and behaviors during emotion-episodes (e.g., Hayes and Smith, 2005; Barlow et al., 2016). Second, low emotional awareness has been linked to biased attention in some clinical contexts (Lane et al., 2018). Third, related personality factors (e.g., biases toward “externally oriented thinking”) are included in leading self-report measures of the related construct of alexithymia (Parker et al., 2003). Finally, emotion learning in childhood appears to involve parent-child interactions in which parents draw attention to (and label) bodily feelings and behaviors during a child’s affective responses [e.g., see work on attunement, social referencing, and related aspects of emotional development (Mumme et al., 1996; Licata et al., 2016; Smith et al., 2018b)] – and the lack of such interactions hinders emotion learning (and mental state learning more generally; Colvert et al., 2008).

In our model, we used relatively high level “outcomes” (i.e., themselves standing in for lower-level representations) to summarize the products of belief updating at lower levels of a hierarchical model. These outcomes were domain-specific, covering interoceptive, proprioceptive and exteroceptive

modalities. A full hierarchical model would consider lower levels, unpacked in terms of sensory modalities; however, the current model, comprising just two levels, is sufficient for our purposes. The bottom portion of **Figure 3A** (in gray) acknowledges the broad form that these lower-level outcomes would be expected to take. The full three-process model would also contain a higher level corresponding to conscious accessibility (for an explicit model and simulations of this higher level, see Smith et al., 2019b). This is indicated by the gray arrows at the top of **Figure 3A**.

Crucially, as mentioned above, attentional focus was treated as a (mental) action that determines the outcome modality or domain to which attention was selectively allocated. Effectively, the agent had to decide which lower-level representations to selectively attend to (i.e., which sequential attention policy to select) in order to figure out what emotional state she was in. Mathematically, this was implemented via interactions in the likelihood mapping – such that being in a particular attentional state selected one and only one precise mapping between the emotional state factor and the outcome information in question (see **Figure 3C**). Formally, this implementation of mental action or attentional focus is exactly the same used to model the exploration of a visual scene using overt eye movements (Mirza et al., 2016). However, on our interpretation, this epistemic foraging was entirely covert; hence mental action (c.f., the premotor theory of attention; (Rizzolatti et al., 1987; Smith and Schenk, 2012; Posner, 2016).

**Figure 3** illustrates the resulting model. The first hidden state factor was a space of (exemplar) emotion concepts (SAD, AFRAID, ANGRY, and HAPPY). The second hidden state factor was attentional focus, and the “B” matrix for this second factor allowed state transitions to be controlled by the agent. The agent could choose to attend to three sources of bodily (interoceptive/proprioceptive) information, corresponding to affective valence (pleasant or unpleasant sensations), autonomic arousal (e.g., high or low heart rate), and motivated proprioceptive action tendencies (approach or avoid). The agent could also attend to two sources of exteroceptive information, including the perceived situation (involving social rejection or a crowded event) and subsequent beliefs about responsibility (attributing agency/blame to self or another). These different sources of information are based on a large literature within emotion research, indicating that they are jointly predictive of self-reported emotions and/or are important factors in affective processing (Russell, 2003; Siemer et al., 2007; Lindquist and Barrett, 2008; Scherer, 2009; Harmon-Jones et al., 2010; Barrett et al., 2011; Barrett, 2017).

Our choice of including valence in particular reflects the fact that our model deals with high levels of hierarchical processing (this choice also enables us to connect more fluently with current literature on emotion concept categories). In this paper, we are using labels like “unpleasant” as pre-emotional constructs. In other words, although affective in nature, we take concepts like “unpleasant” as contributing to elaborated emotional constructs during inference. Technically, valenced states provide evidence for emotional state inference at a higher level (e.g., pleasant sensations provide evidence that one is

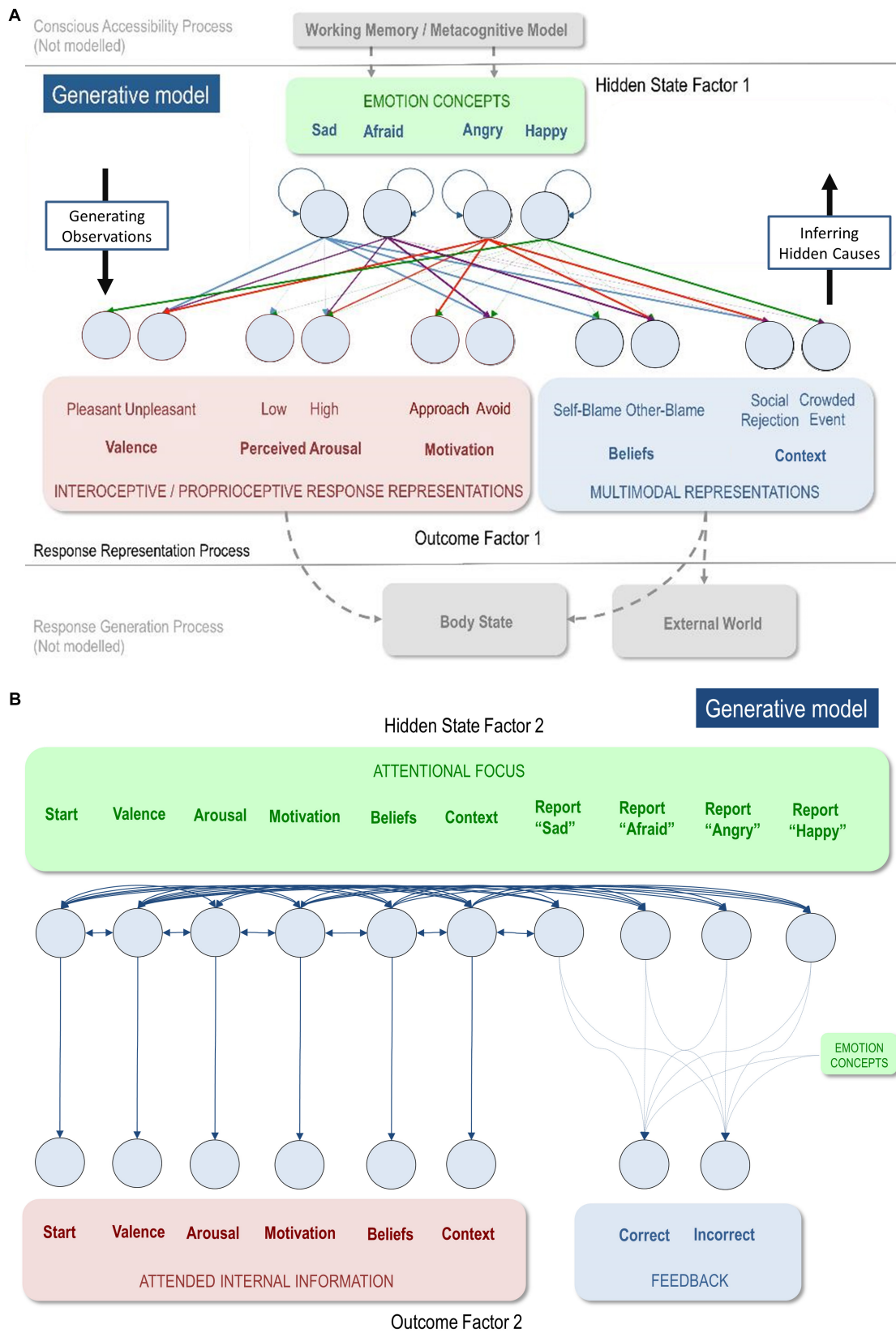
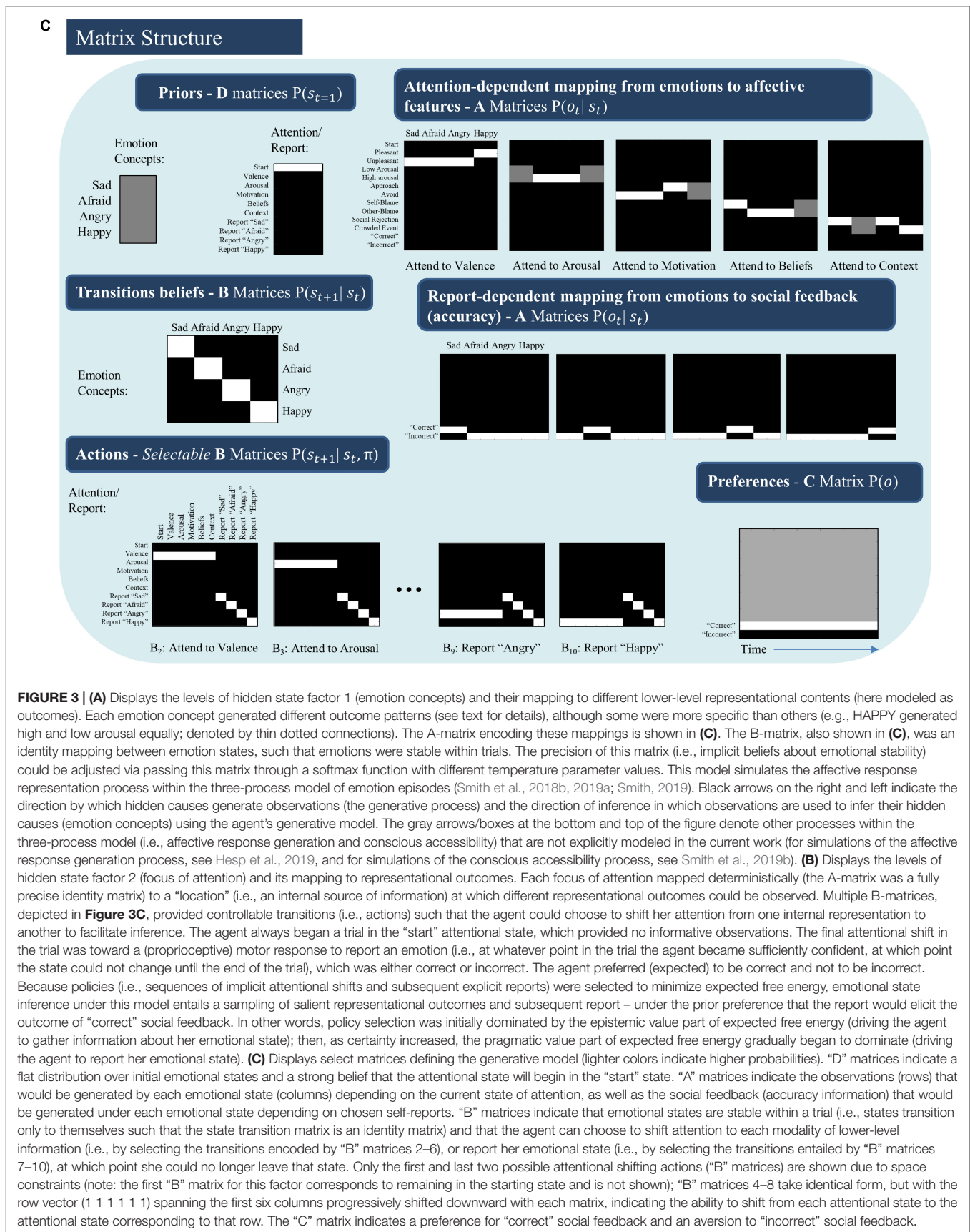


FIGURE 3 | Continued



**FIGURE 3 | (A)** Displays the levels of hidden state factor 1 (emotion concepts) and their mapping to different lower-level representational contents (here modeled as outcomes). Each emotion concept generated different outcome patterns (see text for details), although some were more specific than others (e.g., HAPPY generated high and low arousal equally; denoted by thin dotted connections). The A-matrix encoding these mappings is shown in (C). The B-matrix, also shown in (C), was an identity mapping between emotion states, such that emotions were stable within trials. The precision of this matrix (i.e., implicit beliefs about emotional stability) could be adjusted via passing this matrix through a softmax function with different temperature parameter values. This model simulates the affective response representation process within the three-process model of emotion episodes (Smith et al., 2018b, 2019a; Smith, 2019). Black arrows on the right and left indicate the direction by which hidden causes generate observations (the generative process) and the direction of inference in which observations are used to infer their hidden causes (emotion concepts) using the agent's generative model. The gray arrows/boxes at the bottom and top of the figure denote other processes within the three-process model (i.e., affective response generation and conscious accessibility) that are not explicitly modeled in the current work (for simulations of the affective response generation process, see Hesp et al., 2019, and for simulations of the conscious accessibility process, see Smith et al., 2019b). **(B)** Displays the levels of hidden state factor 2 (focus of attention) and its mapping to representational outcomes. Each focus of attention mapped deterministically (the A-matrix was a fully precise identity matrix) to a "location" (i.e., an internal source of information) at which different representational outcomes could be observed. Multiple B-matrices, depicted in **Figure 3C**, provided controllable transitions (i.e., actions) such that the agent could choose to shift her attention from one internal representation to another to facilitate inference. The agent always began a trial in the "start" attentional state, which provided no informative observations. The final attentional shift in the trial was toward a (proprioceptive) motor response to report an emotion (i.e., at whatever point in the trial the agent became sufficiently confident, at which point the state could not change until the end of the trial), which was either correct or incorrect. The agent preferred (expected) to be correct and not to be incorrect. Because policies (i.e., sequences of implicit attentional shifts and subsequent explicit reports) were selected to minimize expected free energy, emotional state inference under this model entails a sampling of salient representational outcomes and subsequent report – under the prior preference that the report would elicit the outcome of "correct" social feedback. In other words, policy selection was initially dominated by the epistemic value part of expected free energy (driving the agent to gather information about her emotional state); then, as certainty increased, the pragmatic value part of expected free energy gradually began to dominate (driving the agent to report her emotional state). **(C)** Displays select matrices defining the generative model (lighter colors indicate higher probabilities). "D" matrices indicate a flat distribution over initial emotional states and a strong belief that the attentional state will begin in the "start" state. "A" matrices indicate the observations (rows) that would be generated by each emotional state (columns) depending on the current state of attention, as well as the social feedback (accuracy information) that would be generated under each emotional state depending on chosen self-reports. "B" matrices indicate that emotional states are stable within a trial (i.e., states transition only to themselves such that the state transition matrix is an identity matrix) and that the agent can choose to shift attention to each modality of lower-level information (i.e., by selecting the transitions encoded by "B" matrices 2–6), or report her emotional state (i.e., by selecting the transitions entailed by "B" matrices 7–10), at which point she could no longer leave that state. Only the first and last two possible attentional shifting actions ("B" matrices) are shown due to space constraints (note: the first "B" matrix for this factor corresponds to remaining in the starting state and is not shown); "B" matrices 4–8 take identical form, but with the row vector (1 1 1 1 1) spanning the first six columns progressively shifted downward with each matrix, indicating the ability to shift from each attentional state to the attentional state corresponding to that row. The "C" matrix indicates a preference for "correct" social feedback and an aversion to "incorrect" social feedback.



feeling a positive emotion like excitement, joy, or contentment, whereas unpleasant sensations provide evidence that one may be feeling a negative emotion such as sadness, fear, or anger). Based on previous work (Joffily and Coricelli, 2013; Clark et al., 2018; Hesp et al., 2019), we might expect valence to correspond to changes in the precision/confidence associated with lower-level visceromotor and skeletomotor policy selection, or to related internal estimates that can act as indicators of success in uncertainty resolution; see Joffily and Coricelli (2013), de Berker et al. (2016), Peters et al. (2017), Clark et al. (2018). Put another way, feeling good may correspond to high confidence in one's model of how to act, whereas feeling bad may reflect the opposite. Explicitly modeling these lower-level dynamics in a deep temporal model will be the focus of future work.

In the simulations we report here, there were 6 time points in each epoch or trial of emotion inference. At the first time point, the agent always began in an uninformative initial state of attentional focus (the “start” state). The agent's task was to choose what to attend to, and in which order, to infer her most likely emotional state. When she became sufficiently confident, she could choose to respond (i.e., reporting that she felt sad, afraid, angry, or happy). In these simulations the agent selected “shallow” one-step policies, such that she could choose what to attend to next – to gain the most information. Given the number of time points, the agent could choose to attend to up to four of the five possible sources of lower-level information before reporting her beliefs about her emotional state. The “A” matrix mapping attentional focus to attended outcomes was an identity matrix, such that the agent always knew which lower-level information she was currently attending to. This may be thought of as analogous to the proprioceptive feedback consequent on a motor action.

The “B” matrix for hidden emotional states was also an identity matrix, reflecting the belief that emotional states are stable within a trial (i.e., if you start out feeling sad, then you will remain sad throughout the trial). This sort of probability transition matrix in the generative model allows evidence to be accumulated for one state or another over time; here, the emotion concept that provides the best explanation for actively attended evidence in the outcome modalities. The “A” matrix – mapping emotion concepts to outcomes – was constructed such that certain outcome combinations were more consistent with certain emotional states than others: SAD was probabilistically associated with unpleasant valence, either low or high arousal (e.g., lying in bed lethargically vs. intensely crying), avoidance motivation, social rejection, and self-attribution (i.e., self-blame). AFRAID generated unpleasant valence, high arousal, avoidance, other-blame (c.f., fear often being associated with its perceived external cause), and either social rejection or a crowded event (e.g., fear of a life without friends vs. panic in crowded spaces). ANGRY generated unpleasant valence, high arousal, approach, social rejection, and other-blame outcomes. Finally, HAPPY generated pleasant valence, either low or high arousal, either approach or avoidance (e.g., feeling excited to wake up and go to work vs. feeling content in bed and not wanting to go to work), and a crowded event (e.g., having fun at a concert). Because HAPPY does not have strong conceptual links to blame,

we defined a flat mapping between HAPPY and blame, such that either type of blame provided no evidence for or against being happy. Although this mapping from emotional states to outcomes has some face validity, it should not be taken too seriously. It was chosen primarily to capture the ambiguous and overlapping correlates of emotion concepts, and to highlight why adaptive emotional state inference and emotion concept learning can represent difficult problems.

If the “A” matrix encoding state-outcome relationships was completely precise (i.e., if the contingencies above were deterministic as opposed to probabilistic), sufficient information could be gathered through (at most) three attentional shifts; but this becomes more difficult when probabilistic mappings are imprecise (i.e., as they more plausibly are in the real world). **Figure 4** illustrates this by showing how the synthetic subject's confidence about her state decreases as the precision of the mapping between emotional states and outcomes decreases (we measured confidence here in terms of the accuracy of responding in relation to the same setup with infinite precision). Changes in precision were implemented via a temperature parameter of a softmax function applied to a fully precise version of likelihood mappings between emotion concepts and the 5 types of lower-level information that the agent could attend to (where a higher value indicates higher precision). For a more technical account of this type of manipulation, please see Parr and Friston (2017b).

**Figure 4** additionally demonstrates how reporting confidence decreases with decreasing precision of the “B” matrix encoding emotional state transitions, where low precision corresponds to the belief that emotional states are unstable over time. Interestingly, these results suggest that expectations about emotional instability would reduce the ability to understand or infer one's own emotions. From a Bayesian perspective, this result is very sensible: if we are unable to use past beliefs to contextualize the present, it is much harder to accumulate evidence in favor of one hypothesis about emotional state relative to another. Under moderate levels of precision, our numerical analysis demonstrates that the model can conceptualize the multimodal affective responses it perceives with high accuracy.

**Figure 5** illustrates a range of simulation results from an example trial under moderately high levels of “A” and “B” matrix precision (temperature parameter = 2 for each). The upper left plot shows the sequence of (inferred) attentional shifts (note: darker colors indicate higher probability beliefs of the agent, and cyan dots indicate the true states). In this trial, the agent chose a policy in which it attended to valence (observing “unpleasant”), then beliefs (observing “other-blame”), then action (observing “approach”), at which time she became sufficiently confident and chose to report that she was angry. The lower left plot displays the agent's *posterior* beliefs at the end of the trial about her emotional state at each timepoint in the trial, in this case inferring that she had been (and still was) angry. Note that this reflects retrospective inference, and not the agent's beliefs at each timepoint. The lower right and upper right plots display simulated neural responses (based on the neural process theory that accompanies this form of active inference; Friston et al., 2017a), in terms of single-neuron firing rates (raster plots) and local field potentials, respectively. The simulated firing rates in the lower right plot illustrate that

the agent's confidence that she was angry increased gradually with each new observation.

The simulations presented in **Figures 4, 5** make some cardinal points. First, it is fairly straightforward to simulate emotion processing in terms of emotional state inference. This rests upon a particular sort of generative model that can generate outcomes in multiple modalities. The recognition of an emotional state corresponds to the inversion of such models – and therefore necessarily entails multimodal integration. In other words, successfully disambiguating the most likely emotional state here requires consideration of the specific multimodal patterns of experience (i.e., incorporating interoceptive, exteroceptive, and proprioceptive sensations) that would be expected under each emotional state. We have also seen that this form of belief updating – or evidence accumulation – depends sensitively on what sort of evidence is actively attended. This equips the model of emotion concept representation with a form of mental action, which speaks to a tight link between emotion processing and attention to various sources of evidence from within the body – and beyond. Choices to shift attention vs. to self-report are, respectively, driven by the epistemic and pragmatic value of each allowable policy, such that pragmatic value gradually comes to drive the selection of reporting policies as the expected information gain of further attentional

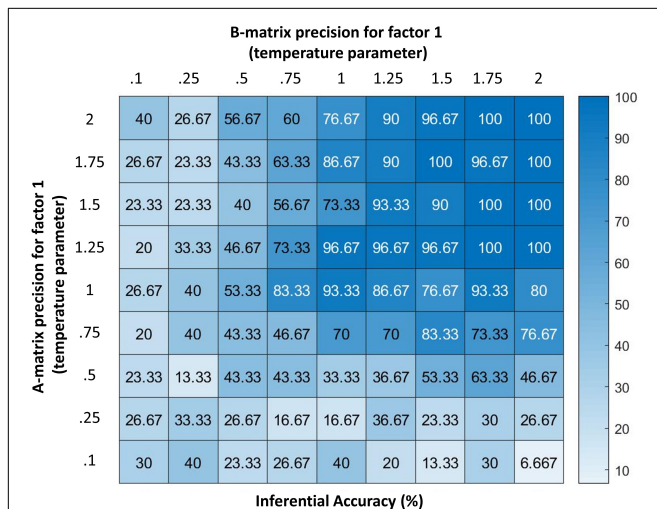
shifts decreases. The physiological plausibility of this emotion inference process has been briefly considered in terms of simulated responses. In the next section, we turn to a more specific construct validation, using empirical phenomenology from neurodevelopmental studies of emotion.

## SIMULATING THE INFLUENCE OF EARLY EXPERIENCE ON EMOTIONAL STATE INFERENCE AND EMOTION CONCEPT LEARNING

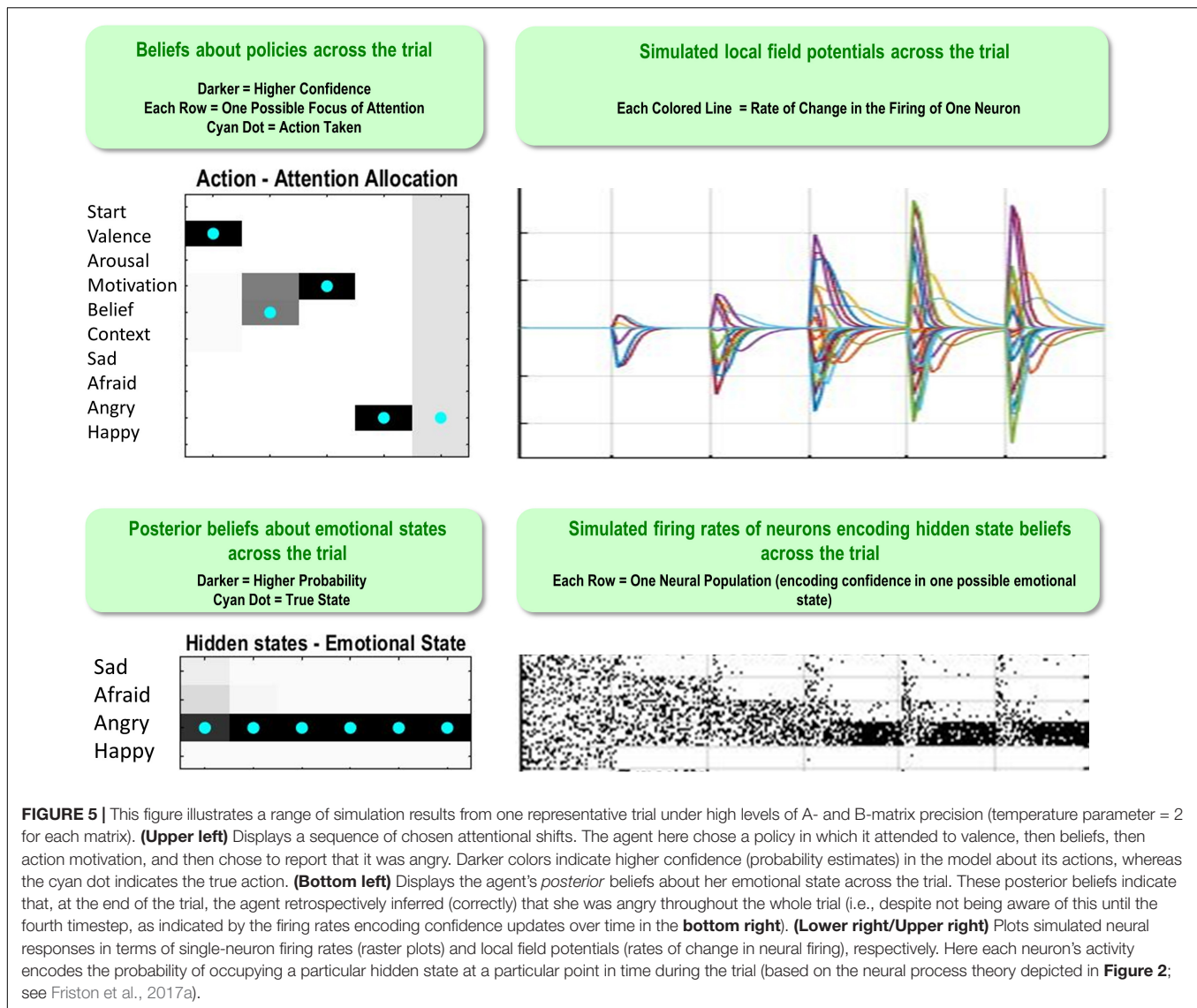
Having confirmed that our model could successfully infer emotional states – if equipped with emotion concepts – we are now in a position to examine emotion concept learning. Specifically, we investigated the conditions under which emotion concepts could be acquired successfully and the conditions under which this type of emotion learning and inference fails.

### Can Emotion Concepts Be Learned in Childhood?

The first question we asked was whether our model could learn about emotions, if it started out with no prior beliefs about how emotions structure its experience. To answer this question, we first ran the model's "A" matrix (mapping emotion concepts to attended outcome information) through a softmax function with a temperature parameter of 0, creating a fully imprecise likelihood mapping. This means that each hidden emotional state predicted all outcomes equally (effectively, none of the hidden states within the emotion factor had any conceptual content). Then we generated 200 sets of observations (i.e., 50 for each emotion concept, evenly interleaved) based on the probabilistic state-outcome mappings encoded in the model described above (i.e., the "generative process"). That is, 50 interleaved learning trials for each emotion were generated by probabilistically sampling from a moderately precise version of the "A" matrix distribution depicted in **Figure 3C** (i.e., temperature parameter = 2). This resulted in 50 sets of observations consistent with the probabilistic mappings for each emotion (e.g., this entailed that roughly 50% of HAPPY trials involved observations of low vs. high arousal, whereas only roughly 1% of HAPPY trials involved the observation of social rejection, etc.). After the 200 learning trials, we then examined the changes in the model's reporting accuracy over time. This meant that the agent, who began with no emotion knowledge (i.e., a fully uninformative "A" matrix), observed patterns of observations consistent with each emotion (as specified above) at 50 timepoints spread out across the 200 trials and needed to learn these associations (i.e., learn the appropriate "A" matrix mapping). This analysis was repeated at several levels of outcome ("A" matrix) and transition ("B" matrix) precision in the generative process – to explore how changes in the predictability or consistency of observed outcome patterns affected the model's ability to learn. In this model, learning was implemented through updating (concentration) parameters for the model's "A" matrix after each trial. The model could also learn prior expectations



**FIGURE 4** | Displays the accuracy of the model (percentage of correct inferences over 30 trials) under different levels of precision for two parameters (denoted by temperature values for a softmax function controlling the specificity of the A and B matrices for hidden state factor 1; higher values indicate higher precision). As can be seen, the model performs with high accuracy at moderate levels of precision. However, its ability to infer its own emotions becomes very poor if the precision of either matrix becomes highly imprecise. Accuracy here is defined in relation to the response obtained from an agent with infinite precision – and can be taken as a behavioral measure of the quality of belief updating about emotional states. These results illustrate how emotion concepts could be successfully inferred despite variability in lower-level observations (e.g., contexts, arousal levels), as would be expected under constructivist theories of emotion (Barrett, 2017); however, they also demonstrate limits in variability, beyond which self-focused emotion recognition would begin to fail.



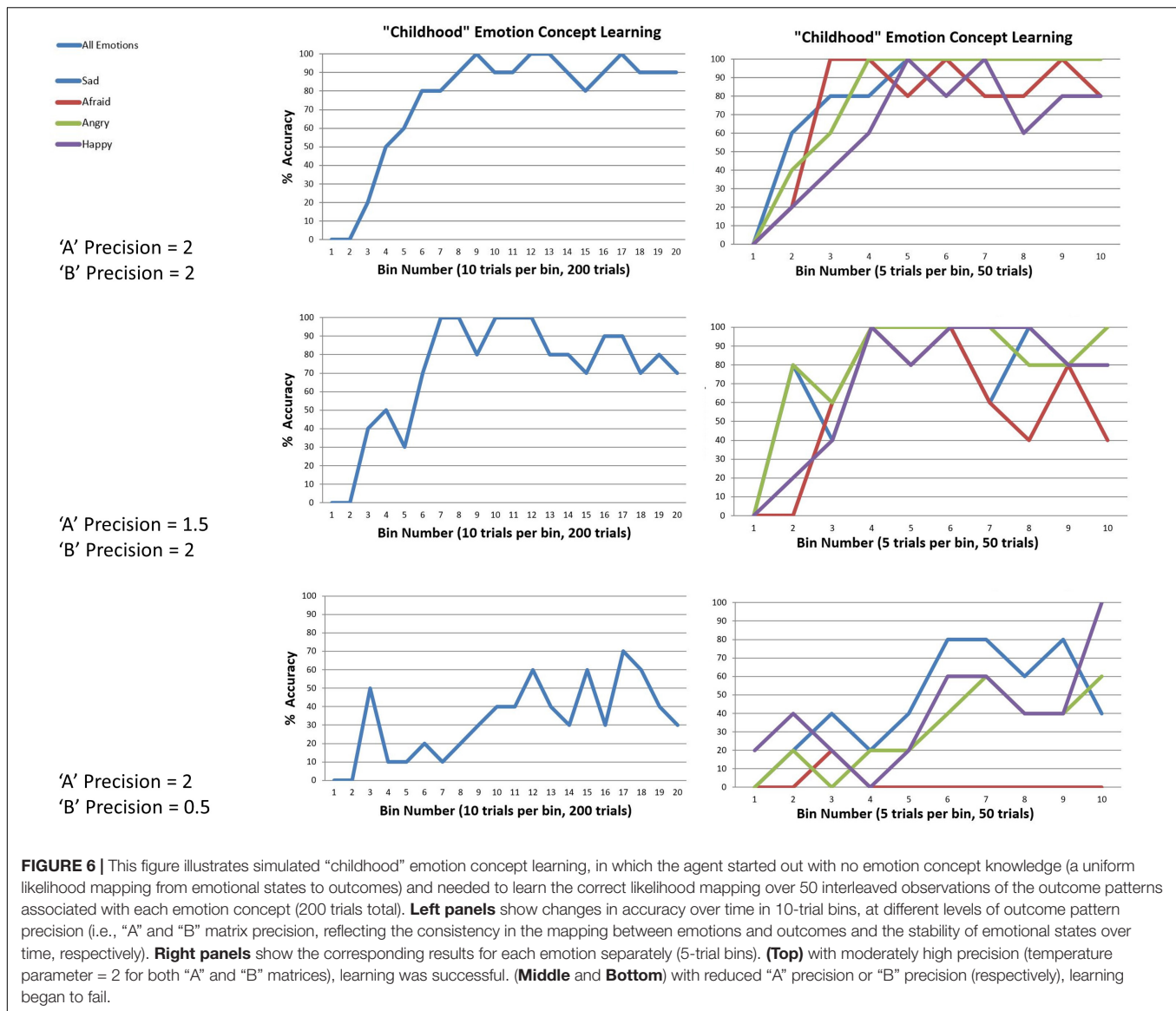
for being in different emotional states, based on updating concentration parameters for its “D” matrix after each trial (i.e., the emotional state it started in). For details of these free energy minimizing learning processes, please see Friston et al. (2016).

We observed that the model could successfully reach 100% accuracy (with minor fluctuation) when both outcome and transition precisions in the generative process were moderately high (i.e., when the temperature parameters for the “A” and “B” matrix of the generative process were 2). The top panel in **Figure 6** illustrates this by plotting the percentage accuracy across all emotions during learning over 200 trials (in bins of 10 trials), and for each emotion (in bins of 5). As can be seen, the model steadily approaches 100% accuracy across trials. The middle and lower panels of **Figure 6** illustrate the analogous results when outcome precision and transition precision were lowered, respectively. The precision values chosen for these illustrations (“A” precision = 1.5, “B” precision = 0.5) represent observed “tipping points” at which learning began to fail (i.e.,

at progressively lower precision values learning performance steadily approached 0% accuracy). As can be seen, lower precision in either the stability of emotions over time or the consistency between observations and emotional states confounded learning. Overall, these findings provide a proof of principle that this sort of model can learn emotion concepts, if provided with a representative and fairly consistent sample of experiences in its “childhood.”

## Can a New Emotion Concept Be Learned in Adulthood?

We then asked whether a new emotion could be learned later, after others had already been acquired (e.g., as in adulthood). To answer this question, we again initialized the model with a fully imprecise “A” matrix (temperature parameter = 0) and set the precision of the “A” and “B” matrices of the generative process to the levels at which “childhood” learning was successful (i.e.,

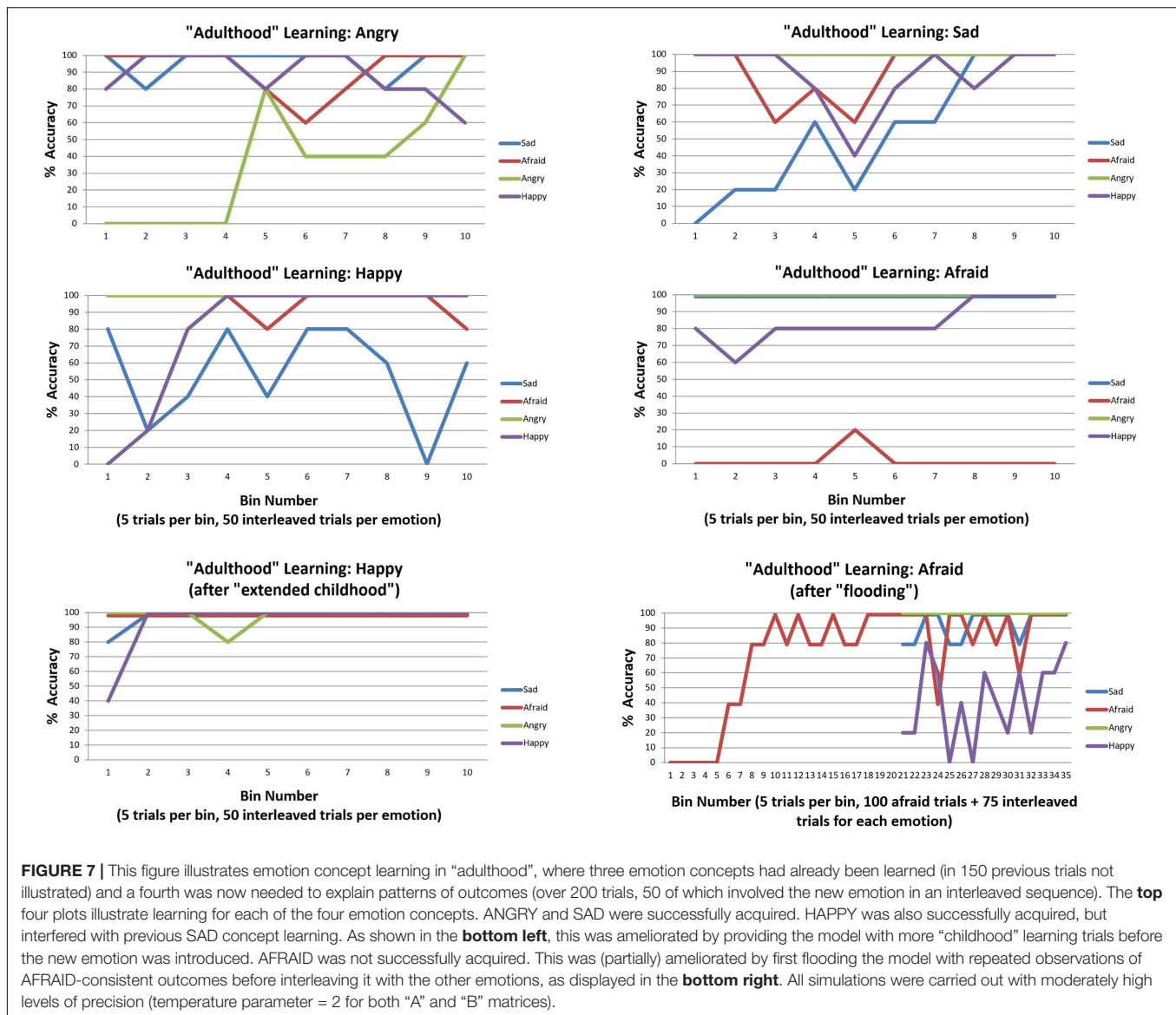


temperature parameter = 2 for each). We then exposed the model to 150 observations that only contained the outcome patterns associated three of the four emotions (50 for each emotion, evenly interleaved). We again allowed the model to accumulate experience in the form of concentration parameters for its "D" matrix – allowing it to learn strong expectations for the emotional states it repeatedly inferred it was in. After these initial 150 trials, we then exposed the model to 200 further trials – using the outcome patterns under all four emotions (50 for each emotion, evenly interleaved). We then asked whether the emotion that was not initially necessary to explain outcomes could be acquired later, when circumstances change.

We first observed that, irrespective of which three emotions were initially presented, accuracy was high by the end of the initial 150 trials (i.e., between 80–100% accuracy for each of the three emotion concepts learned). The upper and middle panels of **Figure 7** illustrate the accuracy over the subsequent 200 trials

as the new emotion was learned. As can be seen in the upper left and right sections of **Figure 7**, ANGRY and SAD were both successfully learned. Interestingly, performance for the other emotions appeared to temporarily drop and then increase again as the new emotion concept was acquired (a type of temporary retroactive interference).

The middle left section of **Figure 7** demonstrates that HAPPY could also be successfully learned; however, it appeared to interfere with prior learning for SAD. Upon further inspection, it appeared that SAD may not have been fully acquired in the first 150 trials (only reaching 80% accuracy near the end). We therefore chose to examine whether an "extended" or "emotionally enriched" childhood might prevent this interference, by increasing the initial learning trial number from 150 to 225 (75 interleaved exposures to SAD, ANGRY, and AFRAID outcome patterns). As can be seen in the lower left panel of **Figure 7**, HAPPY was quickly acquired in the

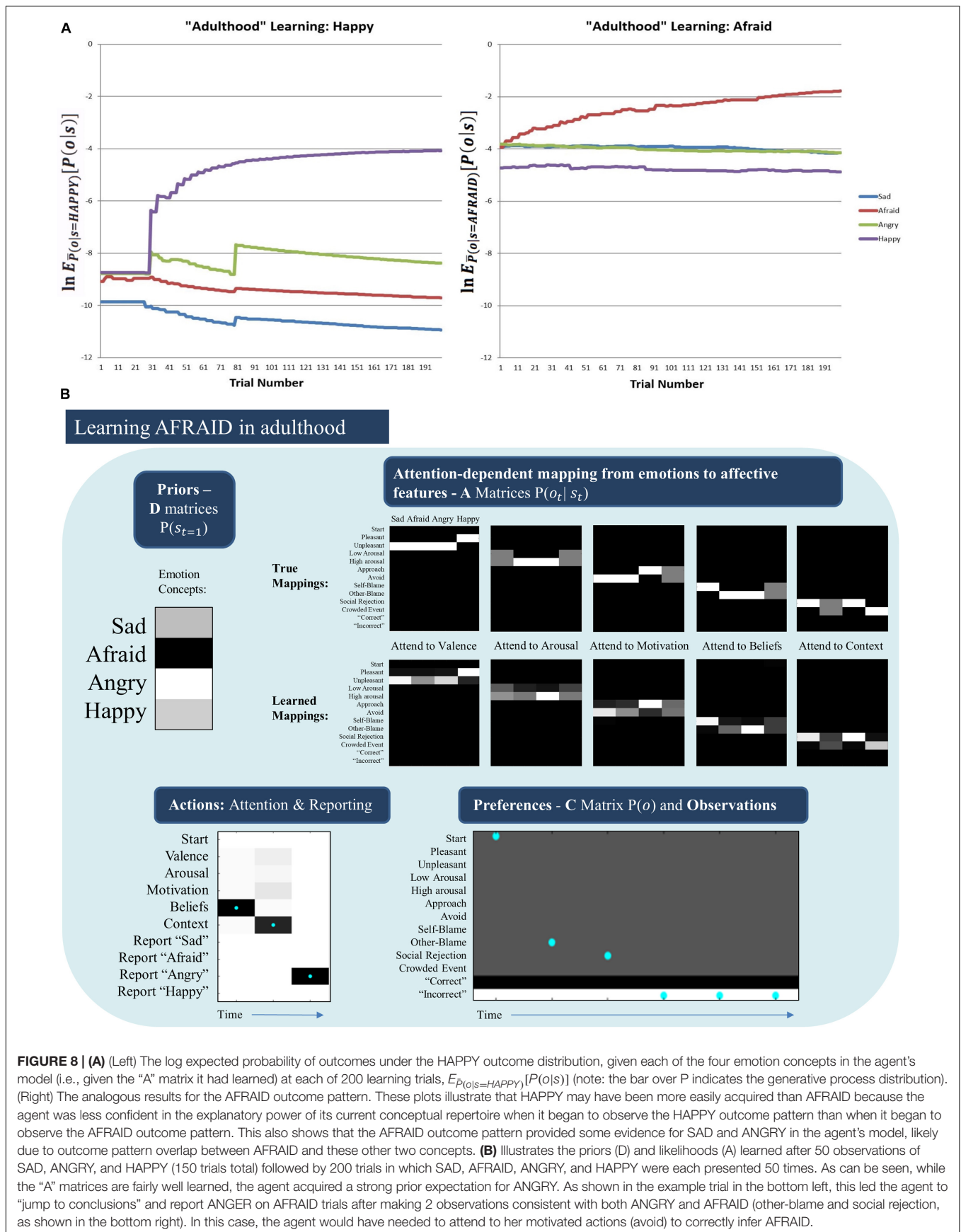


subsequent 200 trials under these conditions, without interfering with previous learning.

Somewhat surprisingly, the model was unable to acquire the AFRAID concept in its "adulthood" (Figure 7, middle right). To better understand this, for each trial we computed the expected evidence for each state, under the distribution of outcomes expected under the generative process (using bar notation to distinguish the process from the model) given a particular state, ( $E_{\bar{P}(o|s)}[P(o|s)]$ ). This was based on the reasoning that, if we treat the different emotional states as alternative models to explain the data, then the likelihood of data given states is equivalent to the evidence for a given state. Figure 8A plots the log transform of this expected evidence for each established emotion concept expected under the distribution of outcomes that would be generated if the "real" emotional state were AFRAID (right panel) and contrasts it with the analogous plot for HAPPY (left panel), which was a more easily acquired concept [we took the

logarithm of this expected evidence (or likelihood) to emphasize the lower evidence values; note that higher (less negative) values correspond to greater evidence in these plots].

As can be seen, in the case of HAPPY the three previously acquired emotion concepts had a relatively low ability to account for all observations, and so HAPPY was a useful construct in providing more accurate explanations of observed outcomes. In contrast, when learning AFRAID the model was already confident in its ability to explain its observations (i.e., the other concepts already had much higher evidence than in the case of HAPPY), and the "AFRAID" outcome pattern also provided moderate evidence for ANGRY and SAD (i.e., the outcome patterns between AFRAID and these other emotion concepts had considerable overlap). In Figure 8B, we illustrate the "A" and "D" matrix values the agent had learned after the total 350 trials (i.e., 150 + 200, as described above), and an exemplar trial in which the agent mistook fear for anger (which was



**FIGURE 8 | (A)** (Left) The log expected probability of outcomes under the HAPPY outcome distribution, given each of the four emotion concepts in the agent's model (i.e., given the "A" matrix it had learned) at each of 200 learning trials,  $E_{P(o|s=HAPPY)} [P(o|s)]$  (note: the bar over P indicates the generative process distribution). (Right) The analogous results for the AFRAID outcome pattern. These plots illustrate that HAPPY may have been more easily acquired than AFRAID because the agent was less confident in the explanatory power of its current conceptual repertoire when it began to observe the HAPPY outcome pattern than when it began to observe the AFRAID outcome pattern. This also shows that the AFRAID outcome pattern provided some evidence for SAD and ANGRY in the agent's model, likely due to outcome pattern overlap between AFRAID and these other two concepts. **(B)** Illustrates the priors (D) and likelihoods (A) learned after 50 observations of SAD, ANGRY, and HAPPY (150 trials total) followed by 200 trials in which SAD, AFRAID, ANGRY, and HAPPY were each presented 50 times. As can be seen, while the "A" matrices are fairly well learned, the agent acquired a strong prior expectation for ANGRY. As shown in the example trial in the bottom left, this led the agent to "jump to conclusions" and report ANGER on AFRAID trials after making 2 observations consistent with both ANGRY and AFRAID (other-blame and social rejection, as shown in the bottom right). In this case, the agent would have needed to attend to her motivated actions (avoid) to correctly infer AFRAID.

the most common confusion). As can be seen, while the “A” matrix mappings were learned fairly well, the agent had learned a strong prior expectation for ANGRY in comparison to its expectation for AFRAID. In the example trial the agent first attends to beliefs (observing other-blame) and then to the context (observing social rejection). These observations are consistent with both ANGRY and AFRAID; however, social rejection is more uniquely associated with ANGRY (i.e., AFRAID is also associated with crowded events, while ANGER is not). Combined with the higher prior expectation for ANGRY, the agent “jumps to conclusions” and becomes sufficiently confident to report ANGER (at which point she receives “incorrect” social feedback). Here, correct inference of AFRAID (i.e., disambiguating AFRAID from ANGRY) would have required that the agent also attend to her action tendencies (where she would have observed avoidance motivation) before deciding which emotion to report.

In this context, greater evidence for an unexplained outcome pattern would be required to “convince” the agent that her currently acquired concepts were not sufficient and that further information gathering (i.e., a greater number of attentional shifts) was necessary before becoming sufficiently confident to report her emotions. Based on this insight, we examined ways in which the model could be given stronger evidence that its current conceptual repertoire was insufficient to account for its observations. We first observed that we could improve model performance by “flooding” the model with an extended pattern of only AFRAID-consistent outcomes (i.e., 100 trials in a row), prior to reintroducing the other emotions in an interleaved fashion. As can be seen in **Figure 7** (bottom right), this led to successful acquisition of AFRAID. However, it temporarily interfered with previous learning of the HAPPY concept. We also observed that by instead increasing the number of AFRAID learning trials from 200 to 600, the model eventually increased its accuracy to between 40 and 80% across the last 10 bins (last 50 trials) – indicating that learning could occur, but at a much slower rate.

Overall, these results confirmed that a new emotion concept could be learned in synthetic “adulthood,” as may occur, for example, in psycho-educational interventions during psychotherapy. However, these results also demonstrate that this type of learning can be more difficult. These results therefore suggest a kind of “sensitive period” early in life where emotion concepts may be more easily acquired.

## Can Maladaptive Early Experiences Bias Emotion Conceptualization?

The final question we asked was whether unfortunate early experiences could hinder our agent’s ability to adaptively infer and/or learn about emotions. Based on the three-process model (Smith et al., 2018b), it has previously been suggested that at least two mechanisms could bring this about:

1. Impoverished early experiences (i.e., not being exposed to the different patterns of observations that would facilitate emotion concept learning).
2. Having early experiences that reinforce maladaptive cognitive habits (e.g., selective attention biases), which can hinder adaptive inference (if the concepts have been

acquired) and learning (if the concepts have not yet been acquired).

We chose to examine both of these possibilities below.

### Non-representative Early Emotional Experiences

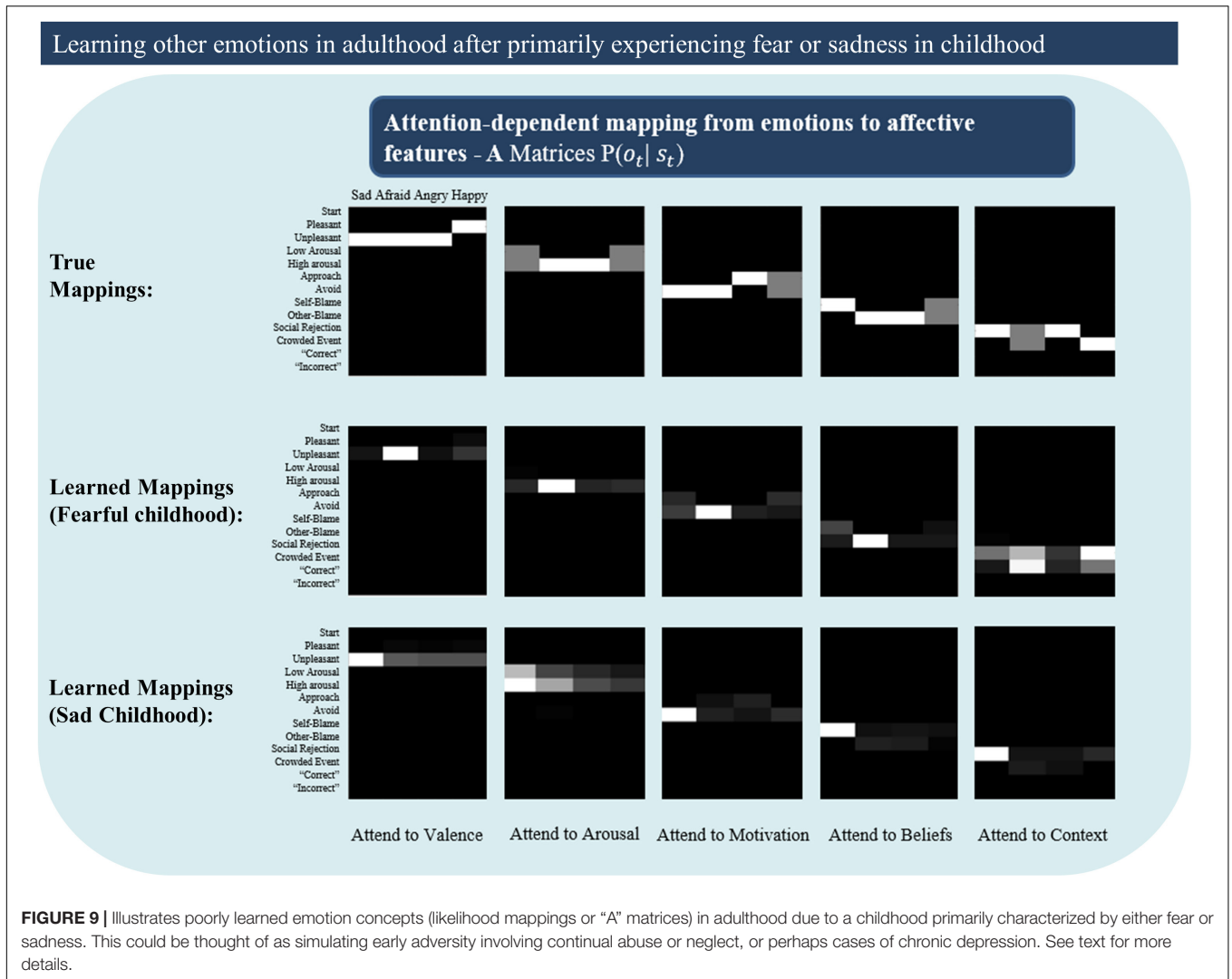
To examine the first mechanism (involving the maladaptive influence of “unrepresentative” early experiences), we used the same learning procedure and parameters described in the previous sections. In this case, however, we exposed the agent to 200 outcomes generated by a generative process where one emotion was experienced 50 times more often than others. Specifically, we examined the cases of a childhood filled with either chronic fear/threat or chronic sadness, as a potential means of simulating the effects of continual childhood abuse or neglect (the sadness simulations might also be relevant to chronic depression over several years). We then examined the model’s ability to learn to infer new emotions in a subsequent 200 trials.

In general, we observed that primarily experiencing fear or sadness during childhood (which could also be thought of as undifferentiated in the sense that they could not be contrasted with other emotions) led the agent to have notable difficulties in learning new emotions later in life. These results were variable upon repeated simulations with different emotions (e.g., verbal reporting continually fluctuated between high and low levels of accuracy for some emotions, while accuracy remained near 0% for others, while yet others were well acquired). For example, in one representative simulation, in which the agent primarily experienced fear during childhood, reporting accuracy continually varied for HAPPY (45% accuracy in the final 20 trials), remained at 0% for ANGRY, remained at 100% for AFRAID, and was stable at 100% for SAD). Whereas primarily experiencing sadness in childhood during a representative simulation led to 0% accuracy for HAPPY, continually varying accuracy for ANGRY and AFRAID (55% accuracy in the final 20 trials for each), and stable high accuracy for SAD (95% accuracy in final 20 trials). Similar patterns of (highly variable) results were observed when performing the same simulations with the other two emotions.

Unlike the results shown in **Figure 8B** – in which likelihood mappings were fairly well acquired (and precise prior expectations for specific emotions hindered correct inference) – poor performance was here explained primarily by poorly acquired likelihood mappings (i.e., the content of the other emotions concepts was often not learned). **Figure 9** illustrates this by presenting the “A” matrices learned by the agent after childhoods dominated by either fear or sadness. As can be seen there, the likelihood mappings do not strongly resemble the true mappings within the generative process. These results in general support the notion that having unrepresentative or insufficiently diverse early emotional experiences could hinder later learning.

### Maladaptive Attention Biases

To examine the second mechanism proposed by the three-process model (involving maladaptive patterns in habitual attention allocation), we equipped the model’s “E” matrix with high prior expectations over specific policies, which meant that it was 50 times more likely to attend to some information and



not to other information. This included: (i) an “external attention bias,” where the agent had a strong habit of focusing on external stimuli (context) and its beliefs about self- and other-blame; (ii) an “internal attention bias,” where the agent had a strong habit of only attending to valence and arousal; and (iii) a “somatic attention bias,” where the agent had a strong habit to attend only to its arousal level and the approach vs. avoid modality.

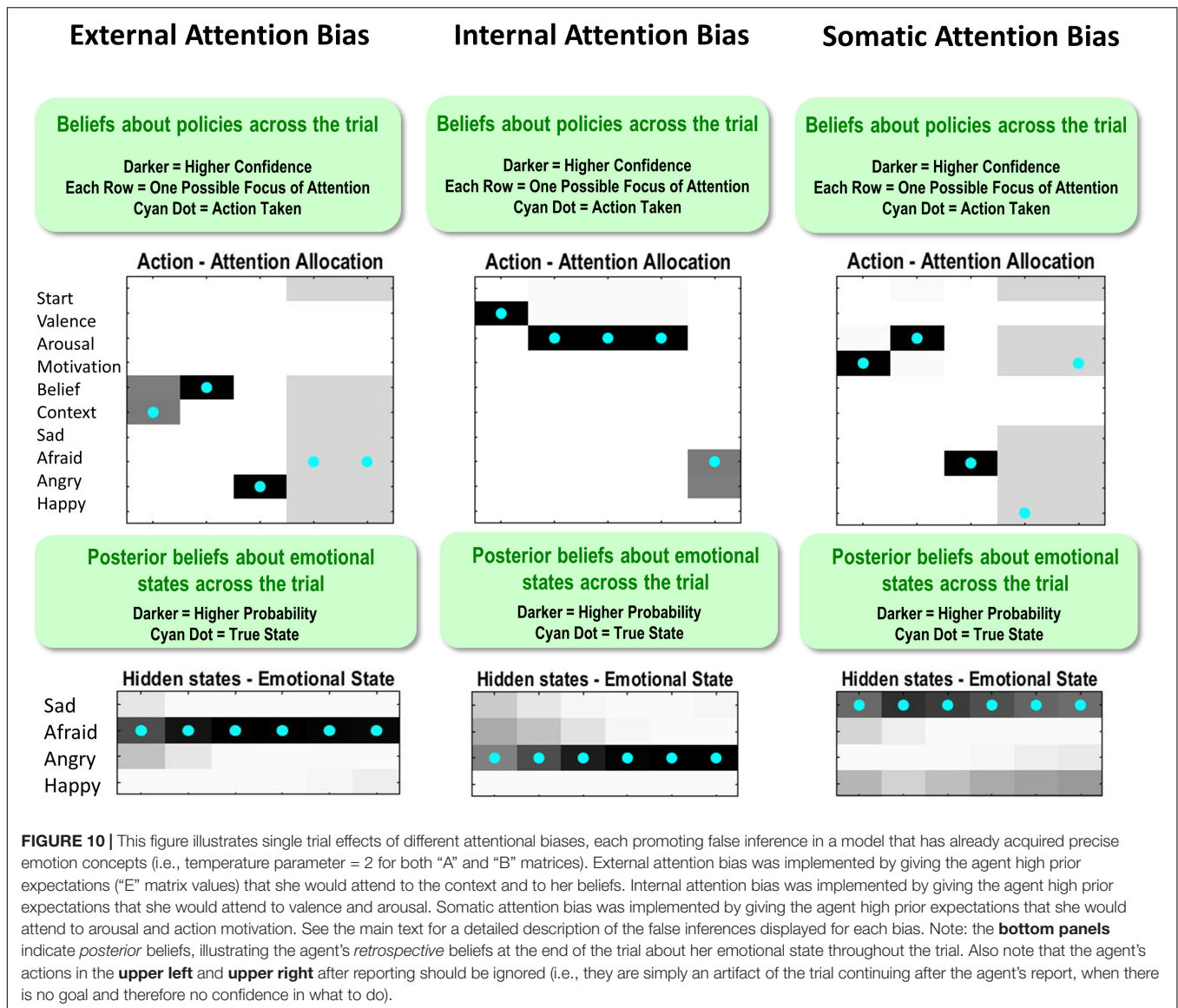
**Figure 10** shows how these different attentional biases promote false inference. On the left, the true state is AFRAID, and the externally focused agent first attends to the stimulus/context (social rejection) and then to her beliefs (other-blame); however, without paying attention to her motivated action (avoid), she falsely reports feeling ANGRY instead of AFRAID (note that, following feedback, there is a retrospective inference that afraid was more probable; similar retrospective inferences after feedback are also shown in the other two examples in **Figure 10**). In the middle, the true state is ANGRY, and the internally focused agent first attends to valence (unpleasant) and then to arousal (high); however, without paying attention to her action tendency (approach), she falsely reports feeling AFRAID. On the right,

the true state is SAD, and the somatically focused agent attends to her motivated action (avoid) and to arousal (high); however, without attending to beliefs (self-blame) she falsely reports feeling AFRAID instead of SAD.

Importantly, these false reports occur in an agent that has already acquired very precise emotion concepts. Thus, this does not represent a failure to learn about emotions, but simply the effect of having learned poor habits for mental action.

The results of these examples were confirmed in simulations of 40 interleaved emotion trials (10 per emotion) in an agent who had already acquired precise emotion concepts (temperature parameter = 2 for both “A” and “B” matrices; no learning). In these simulations, we observed that the externally focused agent had 100% accuracy for SAD, 10% accuracy for AFRAID, 100% accuracy for ANGRY, and 60% accuracy for HAPPY. The internally focused agent had 100% accuracy for SAD, 50% accuracy for AFRAID, 60% accuracy for ANGRY, and 100% accuracy for HAPPY. The somatically focused agent had 10% accuracy for SAD, 100% accuracy for AFRAID, 100% accuracy for ANGRY, and 0% accuracy for HAPPY. Thus,

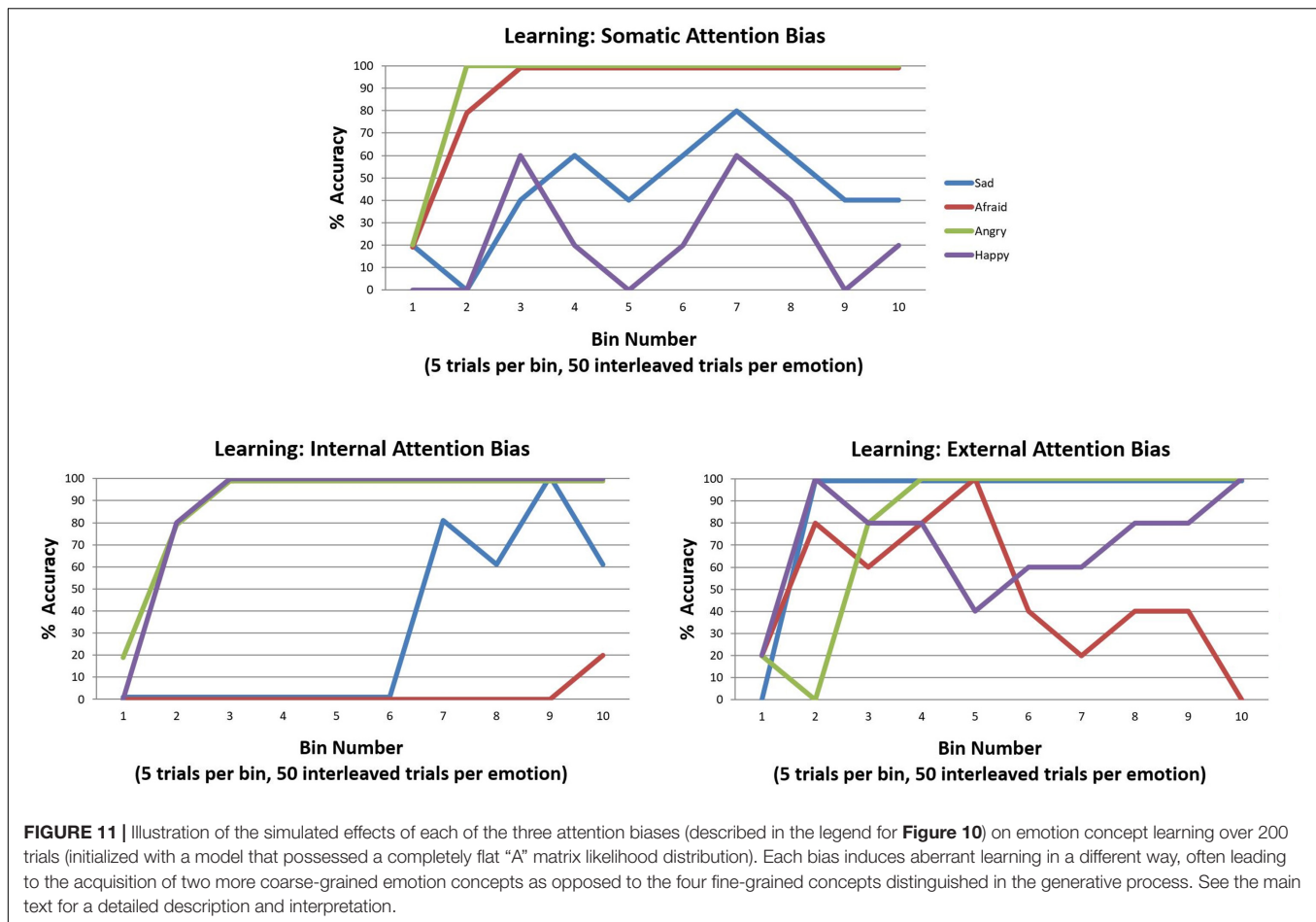




without adaptive attentional habits, the agent was prone to misrepresent her emotions.

In our final simulations, we examined how learning these kinds of attentional biases in childhood could hinder emotion concept learning. To do so, we used the same learning procedure described in the previous section “Can Emotion Concepts Be Learned in Childhood?”. However, in this case we simply equipped the model with the three different attentional biases (“E” matrix prior distributions over policies) and assessed its ability to learn emotion concepts over the 200 trials. The results of these simulations are provided in **Figure 11**. A somatic attention bias primarily allowed the agent to learn two emotion concepts, which corresponded to ANGER and FEAR. However, it is worth highlighting that the low accuracy for the other emotions means that their respective patterns were subsumed under the first two. Thus, it is more accurate to say that the

agent learned two affective concepts, which largely predicted only approach vs. avoidance. In contrast, the internally biased agent easily acquired the distinction between pleasant (HAPPY) and unpleasant emotions (all lumped into ANGER) and began to learn a third concept that distinguished between low vs. high arousal (i.e., SAD vs. ANGRY). However, it did not conceptualize the distinction between approach and avoidance (i.e., ANGRY vs. AFRAID). Lastly, the externally focused agent was somewhat labile in concept acquisition; by the end, she could not predict approach vs. avoidance, but she possessed externally focused concepts with content along the lines of “the state I am in when I’m socially rejected and think it’s my fault vs. someone else’s fault” (SAD vs. ANGRY) and “the state I am in when at a crowded event” (HAPPY). These results provide strong support for the potential role of attentional biases in subverting emotional awareness.



## DISCUSSION

The active inference formulation of emotional processing we have presented represents a first step toward the goal of building quantitative computational models of the ability to learn, recognize, and understand (be “aware” of) one’s own emotions. Although this is clearly a toy model, it does appear to offer some insights, conceptual advances, and possible predictions.

First, in simulating differences in the precision (specificity) of emotion concepts, some intuitive but interesting phenomena emerged. As would be expected, differences in the specificity of the content of emotion concepts – here captured by the precision of the likelihood mapping from states to outcomes (i.e., the precision of what pattern of outcomes each emotion concept predicted) – led to differences in inferential accuracy. This suggests that, as would be expected, those with more precise emotion concepts would show greater understanding of their own affective responses. Perhaps less intuitively, beliefs about the stability of emotion concepts – here captured by the precision of expected state transitions – also influenced inferential accuracy. This predicts that a belief that emotional states are more stable (less labile) over time would also facilitate one’s ability to correctly infer what they are feeling. This appears consistent with

the low levels of emotional awareness or granularity observed in borderline personality disorder, which is characterized by emotional instability (Levine et al., 1997; Suvak et al., 2011).

Next, in simulating emotion concept learning, a few interesting insights emerged. Our simulations first confirmed that emotion concepts could successfully be learned, even when their content was cast (as done here) as complex, probabilistic, and highly overlapping response patterns across interoceptive, proprioceptive, exteroceptive, and cognitive domains. This was true when all emotion concepts needed to be learned simultaneously (as in childhood; see Widen and Russell, 2008; Hietanen et al., 2016), and was also true when a single new emotion concept was learned after others had already been acquired (as in adulthood during psycho-educational therapeutic interventions; e.g., see Hayes and Smith, 2005; Barlow et al., 2016; Burger et al., 2016; Lumley et al., 2017).

These results depended on whether the observed outcomes during learning were sufficiently precise and consistent. One finding worth highlighting was that emotion concept learning was hindered when the precision of transitions among emotional states was too low. This result may be relevant to previous empirical results in populations known to show reduced understanding of emotional states, such as those with autism

(Silani et al., 2008; Erbas et al., 2013) and those who grow up in socially impoverished or otherwise adverse (unpredictable) environments (Colvert et al., 2008; Lane et al., 2018). In autism, it has been suggested that overly imprecise beliefs about state transitions may hinder mental state learning, because such states require tracking abstract behavioral patterns over long timescales (Lawson et al., 2014, 2017; Haker et al., 2016). Children who grow up in impoverished environments may not have the opportunity to interact with others to observe stable patterns in other's affective responses; or receive feedback about their own (Pears and Fisher, 2005; Lane et al., 2018; Smith et al., 2018b). Our results successfully reproduce these phenomena – which represent important examples of mental state learning that may depend on consistently observed outcome patterns that are relatively stable over time.

As emotion concepts are known to differ in different cultures (Russell, 1991), our model and results may also relate to the learning mechanisms allowing for this type of culture-specific emotion categorization learning. Specifically, the “correct” and “incorrect” social feedback in our model could be understood as linguistic feedback from others in one's culture (e.g., a parent labeling emotional reactions for a child using culture-specific categories). If this feedback is sufficiently precise, then emotion concept learning could proceed effectively – even if the probabilistic mapping from emotion categories to other perceptual outcomes is fairly imprecise (i.e., which appears to be the case empirically; Barrett, 2006, 2017).

Another insight worth highlighting was that learning was more difficult when the agent had already acquired previous concepts but entertained the possible existence of new emotions she had not already learned. An interesting observation was that, in some cases (e.g., learning SAD), new learning temporarily interfered with old learning before being fully integrated into the agent's cognitive repertoire (an effect – termed retroactive interference – that has been extensively studied empirically within learning and memory research; Martínez et al., 2014; Darby and Sloutsky, 2015). In the case of learning HAPPY, we found that more extensive learning in the model's “childhood” was necessary to prevent this type of interference with respect to previous acquisition of the concept SAD. A second interesting observation was that AFRAID was very difficult to learn after the other concepts were fully acquired. This appeared to be because the agent had already learned to be highly confident about the explanatory power of her current conceptual repertoire (combined with the fact that AFRAID had considerable outcome overlap with SAD and ANGRY). It was necessary to provide the model with a persistent “flooding” of observations consistent with the new emotion to reduce its confidence sufficiently to acquire a new concept. It is not clear that this type of flooding is realistic, but perhaps resembles the extended periods of fear that occur during exposure-based behavioral therapies (Cooper et al., 2017). We should also emphasize that, due to the oversimplified nature of the mappings between emotion concepts and affective response features in our simulations, the difficulties observed in learning these specific emotions should not be taken too seriously. However, these overall results do predict that (and

illustrate why) emotion concept learning in general should be more difficult in adulthood, and that emotion learning may have a kind of “sensitive period” in childhood (as supported by previous empirical findings; e.g., Pears and Fisher, 2005; Colvert et al., 2008).

The manner in which concept learning was implemented in these simulations may also have more general implications (for considerations of this approach to concept learning more generally, see Smith et al., 2019d). Typically, in Active Inference simulations the state space structure of a model is specified in advance (e.g., Schwartenbeck et al., 2015; Mirza et al., 2016; Parr and Friston, 2017a). Our model was instead equipped with “blank” hidden states devoid of content (i.e., these states started out predicting all outcomes with equal probability in the simulated learning). Over multiple exposures to the observed outcomes, these blank hidden states came to acquire conceptual content that captured distinct statistical patterns in the lower-level affective response components of the model. In some current neural process theories (Bastos et al., 2012; Friston et al., 2017a, 2018; Parr and Friston, 2018), distinct cortical columns are suggested to represent distinct hidden states. Under such theories, our learning model would suggest that the brain might contain “reserve” cortical columns available to capture new patterns of lower-level covariance if/when they begin to be observed in interaction with the world. To our knowledge, no direct evidence of such “reserve neurons” has been observed, although the generation of new neurons (with new synaptic connections) is known to occur in the hippocampus (Chancey et al., 2013). There is also the well-known phenomenon of “silent synapses” in the brain, which can persist into adulthood and become activated when new learning becomes necessary (e.g., see Kerchner and Nicoll, 2008; Chancey et al., 2013; Funahashi et al., 2013). Another interesting consideration is that, during sleep, it appears that many (but not all) synaptic strength increases acquired in the previous day are attenuated (Tononi and Cirelli, 2014). This has been suggested to correspond to a process of Bayesian model reduction, in which redundant model parameters are identified and removed to prevent model over-fitting and promote selection of the most parsimonious model that can successfully account for previous observations (Hobson and Friston, 2012). This also suggests that increases in “reserve” representational resources available for state space expansion (as in concept learning) could perhaps occur after sleep. In short, the acquisition of new concepts, emotion-related or otherwise, speaks to important issues in structure learning. The approach used here offers one solution to the question of how to expand a model, which could complement work on strategies for reducing a model (Friston et al., 2017b).

Although the neural process theory associated with active inference is cast at the level of canonical microcircuits and message passing, and therefore does not make *a priori* predictions about the brain regions that implement the emotion-related processes in our model, it nonetheless can afford empirical testing of macro-anatomical correlates. That is, this process theory can be used to generate predicted neural response time courses during emotional state inference and emotion concept

learning in our model, and the macro-anatomical correlates of these time courses can then be established using neuroimaging methods. At present, the three-process model (Smith et al., 2018b), and supporting evidence (McRae et al., 2008; Smith et al., 2015, 2017b,c, 2018c,d,e, 2019a,b), has identified a number of large-scale networks that plausibly implement the processes we have simulated – and could therefore provide *a priori* hypotheses for future studies along these lines (Yeo et al., 2011; Barrett and Satpute, 2013). For example, “limbic network” regions (including orbitofrontal cortex and amygdala, among others), “salience network” regions (including the anterior insula and dorsal anterior cingulate, among others) and somatomotor/posterior insula regions all appear to be involved in generating affective responses and representing either visceral, somatic, or proprioceptive states at a perceptual level. Regions of the paralimbic cortex (e.g., “default mode network,” with major hubs in the medial prefrontal cortex and posterior cingulate) are in turn most strongly implicated in conceptual inference (Binder et al., 2009) – such as the emotion concept representation processes simulated here. Thus, activity in a number of distinct brain regions/networks would be expected to show associations with distinct belief updating processes in our model.

A final insight offered by our model pertains to the possibility that maladaptive emotional state inference could be due to early experience. We demonstrated this in two ways. First, we simulated exposure to a large number of single emotion-provoking situations in childhood, promoting precise and highly engrained prior expectations for being in a single emotional state, as well as preferential learning of the respective outcome patterns for that state over others. We found that different kinds of “unrepresentative” (single-emotion) outcome patterns in early experience (e.g., chronic fear or sadness in childhood abuse/neglect or severe chronic depression) prevented learning other emotion concepts in somewhat inconsistent ways in repeated simulations. Overall, however, these results supported the idea that later emotional state inferences and emotion concept learning could be compromised by this type of maladaptive early experience. This could potentially relate to cognitive bias learning, such as the negative interpretation biases characteristic of mood and anxiety disorders (which have been interpreted within computational frameworks; Mogg and Bradley, 2005; Smith et al., 2017a).

Second, we examined the possibility that maladaptive cognitive habits could hinder emotional awareness. Here, we demonstrated that such habits can promote false emotional state inference and can hinder emotion concept learning. Specifically, we found that different types of external, internal, and somatic biases led to the acquisition of coarser-grained emotion concepts that failed to distinguish between various elements of affective responses. Aside from its relevance to cognitive biases more generally, these results could also explain certain empirical phenomena in emotional awareness research, such as the finding that males tend to score lower on emotional awareness measures than females (Wright et al., 2017). Specifically, while a genetic contribution to such findings is possible, it is also known that many cultures reinforce emotion avoidance in boys more than

in girls in childhood (Fivush et al., 2000; Diener and Lucas, 2004; Chaplin et al., 2005), and can promote beliefs that paying attention to emotions is a sign of weakness or that emotional information simply carries little practical value. This type of learning could plausibly reinforce biased patterns of attention similar to those simulated here. Thus, our simulations suggest an interesting, testable mechanism by which such (potentially socialization-based) differences may arise.

In closing, it is important to note that this model is deliberately simple and is meant only to represent a proof of principle that emotion inference and learning can be modeled within a neurocomputational framework from first principles. We chose a particular pattern of state-outcome mappings to simulate the content of emotion concepts, but this is unlikely to represent a fully accurate depiction of human emotion concepts or the outcomes they predict. Human emotion concepts likely draw on much higher-dimensional patterns of somatic and visceral sensations, behavioral motivations, and cognitive appraisal patterns. There are also “secondary” emotion concepts like jealousy or embarrassment, which may require including more specific context and appraisal observations in a model (e.g., observing a lover with a competing suitor, observing oneself committing actions that break social norms, etc.). Further, human agents (or at least some of them) are likely to have a much richer space of both emotion and non-emotion concepts available for explaining their patterns of internal experience in conjunction with other beliefs and exteroceptive evidence (e.g., a pattern of low arousal, unpleasant valence, and avoidance in many contexts could also be explained by the concept of sickness rather than sadness; see Smith et al., 2019b). A more complete model would take into account many different possible conceptual interpretations of this sort. In addition, our simulations only attempted to capture the second process within the three-process model (i.e., affective response representation; Smith et al., 2018b). Incorporating the other two processes (affective response generation and conscious access) would undoubtedly induce additional dynamics (including explicit brain-body interactions) that could alter or nuance the simulation results we have provided. Modeling these additional processes will be an important goal of future work (see Smith et al., 2019b).

A final more general limitation with this type of modeling is that, in its current form, there are limited means of evaluating how well it represents the true form of emotional state inference and emotion concept learning implemented in the human brain. Here, we have focused on reproducing and validating a minimal model that evinces emotional state inference and learning within the active inference framework. Crucially, this model has – by construction – a construct validity with the three-process model and associated empirical evidence. As noted above, external validation of the model’s ability to capture human brain processes will be an important next step, and can be done, for example, by examining whether the simulated neural responses we have presented are observable within particular brain regions during future neuroimaging studies of attending to – and reporting – one’s own emotions (e.g., Lane et al., 1997; Gusnard et al., 2001; Smith et al., 2014, 2018c,d).

With these limitations in mind, however, this approach to computationally modeling emotion-related processes appears promising with respect to the initial insights it can offer. It can illustrate selective information integration in the service of conceptual inference, it can successfully simulate concept learning and some of its known vulnerabilities, and it can highlight maladaptive interactions between cognitive habits, early experience, and the ability to understand and be aware of one's own emotions later in life, all of which may play important roles in the development of emotional pathology. Finally, it highlights the potential for future empirical work in which tasks could be adapted to the broad structure of such models, which would allow investigation of individual differences in emotion processing as well as its neural basis. In other words, once we have a validated model of these emotion-related processes – at the subjective and neuronal level – we can, in principle, fit the model to observed responses and thereby phenotype subjects in terms of their emotion-related beliefs states (Schwartenbeck and Friston, 2016).

## Software Note

Although the generative model – specified by the various matrices described in this paper – changes from application to application, the belief updates are generic and can be implemented using standard routines (here `spm_MDP_VB_X.m`). These routines are available as Matlab code in the latest version of SPM academic

software<sup>1</sup>. The simulations in this paper can be reproduced (and customized) via running the Matlab code included here in the **Supplementary Material** (`Emotion_learning_model.m`).

## AUTHOR CONTRIBUTIONS

RS took the lead in writing the manuscript and constructing the model. TP assisted with programming the model and running simulations, and also assisted in writing and editing the manuscript. KF provided guidance in constructing the model and also contributed to the writing and editing of the manuscript.

## FUNDING

TP is supported by the Rosetrees Trust (Award Number 173346). KF is a Wellcome Trust Principal Research Fellow (Ref: 088130/Z/09/Z).

## SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fpsyg.2019.02844/full#supplementary-material>

<sup>1</sup><http://www.fil.ion.ucl.ac.uk/spm/>

## REFERENCES

- Bagby, R., Parker, J., and Taylor, G. (1994a). The twenty-item toronto alexithymia scale—I. Item selection and cross-validation of the factor structure. *J. Psychosom. Res.* 38, 23–32. doi: 10.1016/0022-3999(94)90005-1
- Bagby, R., Parker, J., and Taylor, G. (1994b). The twenty-item toronto alexithymia scale—II. Convergent, discriminant, and concurrent validity. *J. Psychosom. Res.* 38, 33–40. doi: 10.1016/0022-3999(94)90006-x
- Barchard, K., and Hakstian, A. (2004). The nature and measurement of emotional intelligence abilities; basic dimensions and their relationships with other cognitive abilities and personality variables. *Educ. Psychol. Measure.* 64, 437–462. doi: 10.1177/0013164403261762
- Barlow, D., Allen, L., and Choate, M. (2016). Toward a Unified Treatment for Emotional Disorders - Republished Article. *Behav. Ther.* 47, 838–853. doi: 10.1016/j.beth.2016.11.005
- Barrett, L. (2006). Are emotions natural kinds? *Perspect. Psychol. Sci.* 1, 28–58. doi: 10.1111/j.1745-6916.2006.00003.x
- Barrett, L. (2017). *How Emotions are Made: The Secret Life of the Brain*. New York, NY: Houghton Mifflin Harcourt.
- Barrett, L., Mesquita, B., and Gendron, M. (2011). Context in emotion perception. *Curr. Direct. Psychol. Sci.* 20, 286–290. doi: 10.1177/09637214111422522
- Barrett, L., and Satpute, A. (2013). Large-scale brain networks in affective and social neuroscience: towards an integrative functional architecture of the brain. *Curr. Opin. Neurobiol.* 23, 361–372. doi: 10.1016/j.conb.2012.12.012
- Barrett, L., and Simmons, W. (2015). Interoceptive predictions in the brain. *Nat. Rev. Neurosci.* 16, 419–429. doi: 10.1038/nrn3950
- Baslet, G., Termini, L., and Herbener, E. (2009). Deficits in emotional functioning in schizophrenia and their relationship with other measures of functioning. *J. Nerv. Mental Dis.* 197, 655–660. doi: 10.1097/NMD.0b013e3181b3b20f
- Bastos, A., Usrey, W., Adams, R., Mangun, G., Fries, P., and Friston, K. (2012). Canonical microcircuits for predictive coding. *Neuron* 76, 695–711. doi: 10.1016/j.neuron.2012.10.038
- Berthoz, S., Ouhayoun, B., and Parage, N. (2000). Etude préliminaire des niveaux de conscience émotionnelle chez des patients déprimés et des contrôles. (Preliminary study of the levels of emotional awareness in depressed patients and controls.). *Ann. Med. Psychol.* 158, 665–672.
- Binder, J., Desai, R., Graves, W., and Conant, L. (2009). Where is the semantic system? a critical review and meta-analysis of 120 functional neuroimaging studies. *Cereb. Cortex* 19, 2767–2796. doi: 10.1093/cercor/bhp055
- Bréjard, V., Bonnet, A., and Pedinielli, J. (2012). The role of temperament and emotional awareness in risk taking in adolescents. *L'Encéphale* 38, 1–9. doi: 10.1016/j.enceph.2011.04.005
- Brown, T. H., Zhao, Y., and Leung, V. (2009). “Hebbian plasticity A2–squire,” in *Encyclopedia of Neuroscience*, ed. R. Larry (Oxford: Academic Press), 1049–1056.
- Burger, A., Lumley, M., Carty, J., Latsch, D., Thakur, E., Hyde-Nolan, M., et al. (2016). The effects of a novel psychological attribution and emotional awareness and expression therapy for chronic musculoskeletal pain: a preliminary, uncontrolled trial. *J. Psychosom. Res.* 81, 1–8. doi: 10.1016/j.jpsychores.2015.12.003
- Bydlowski, S., Corcos, M., Jeammet, P., Paterniti, S., Berthoz, S., Laurier, C., et al. (2005). Emotion-processing deficits in eating disorders. *Int. J. Eat. Disord.* 37, 321–329. doi: 10.1002/eat.20132
- Chancey, J., Adlaf, E., Sapp, M., Pugh, P., Wadiche, J., and Overstreet-Wadiche, L. (2013). GABA depolarization is required for experience-dependent synapse unsilencing in adult-born neurons. *J. Neurosci.* 33, 6614–6622. doi: 10.1523/JNEUROSCI.0781-13.2013
- Chaplin, T., Cole, P., and Zahn-Waxler, C. (2005). Parental socialization of emotion expression: gender differences and relations to child adjustment. *Emotion* 5, 80–88. doi: 10.1037/1528-3542.5.1.80
- Ciarrochi, J., Caputi, P., and Mayer, J. (2003). The distinctiveness and utility of a measure of trait emotional awareness. *Pers. Individ. Differ.* 34, 1477–1490. doi: 10.1016/s0191-8869(02)00129-0
- Clark, J., Watson, S., and Friston, K. (2018). What is mood? A computational perspective. *Psychol. Med.* 48, 2277–2284. doi: 10.1017/S0033291718000430

- Colvert, E., Rutter, M., Kreppner, J., Beckett, C., Castle, J., Groothues, C., et al. (2008). Do theory of mind and executive function deficits underlie the adverse outcomes associated with profound early deprivation?: findings from the english and romanian adoptees study. *J. Abnorm. Child Psychol.* 36, 1057–1068. doi: 10.1007/s10802-008-9232-x
- Conant, C., and Ashbey, W. (1970). Every good regulator of a system must be a model of that system. *Int. J. Syst. Sci.* 1, 89–97. doi: 10.1080/00207727008920220
- Consoli, S., Lemogne, C., Roch, B., Laurent, S., Plouin, P., and Lane, R. (2010). Differences in emotion processing in patients with essential and secondary hypertension. *Am. J. Hyperten.* 23, 515–521. doi: 10.1038/ajh.2010.9
- Cooper, A. A., Clifton, E. G., and Feeny, N. C. (2017). An empirical review of potential mediators and mechanisms of prolonged exposure therapy. *Clin. Psychol. Rev.* 56, 106–121. doi: 10.1016/j.cpr.2017.07.003
- Darby, K., and Sloutsky, V. (2015). The cost of learning: interference effects in memory development. *J. Exp. Psychol. Gen.* 144, 410–431. doi: 10.1037/xge0000051
- de Berker, A., Rutledge, R., Mathys, C., Marshall, L., Cross, G., Dolan, R., et al. (2016). Computations of uncertainty mediate acute stress responses in humans. *Nat. Commun.* 7:10996. doi: 10.1038/ncomms10996
- Diener, M., and Lucas, R. (2004). Adults desires for childrens emotions across 48 countries: \$sociations with individual and national characteristics. *J. Cross Cult. Psychol.* 35, 525–547. doi: 10.1177/0022022104268387
- Donges, U., Kersting, A., Dannlowski, U., Lalee-Mentzel, J., Arolt, V., and Suslow, T. (2005). Reduced awareness of others' emotions in unipolar depressed patients. *J. Nerv. Mental Dis.* 193, 331–337. doi: 10.1097/01.nmd.0000161683.02482.19
- Erbas, Y., Ceulemans, E., Boonen, J., Noens, I., and Kuppens, P. (2013). Emotion differentiation in autism spectrum disorder. *Res. Autism Spectr. Disord.* 7, 1221–1227. doi: 10.1016/j.rasd.2013.07.007
- Fivush, R., Brotman, M., Buckner, J., and Goodman, S. (2000). Gender differences in parent-child emotion narratives. *Sex Roles* 42, 233–253. doi: 10.1023/A:1007091207068
- Frewen, P., Lane, R., Neufeld, R., Densmore, M., Stevens, T., and Lanius, R. (2008). Neural correlates of levels of emotional awareness during trauma script-imagery in posttraumatic stress disorder. *Psychosom. Med.* 70, 27–31. doi: 10.1097/psy.0b013e31815f66d4
- Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., O Doherty, J., and Pezzulo, G. (2016). Active inference and learning. *Neurosci. Biobehav. Rev.* 68, 862–879. doi: 10.1016/j.neubiorev.2016.06.022
- Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., and Pezzulo, G. (2017a). Active inference: a process theory. *Neural Comput.* 29, 1–49. doi: 10.1162/NECO\_a\_00912
- Friston, K., Lin, M., Frith, C., Pezzulo, G., Hobson, J., and Ondobaka, S. (2017b). Active inference. Curiosity and insight. *Neural Comput.* 29, 2633–2683. doi: 10.1162/neco\_a\_00999
- Friston, K., Parr, T., and de Vries, B. (2017c). The graphical brain: belief propagation and active inference. *Netw. Neurosci.* 1, 381–414. doi: 10.1162/NETN\_a\_00018
- Friston, K., Rosch, R., Parr, T., Price, C., and Bowman, H. (2018). Deep temporal models and active inference. *Neurosci. Biobehav. Rev.* 90, 486–501. doi: 10.1016/J.NEUBIOREV.2018.04.004
- Funahashi, R., Maruyama, T., Yoshimura, Y., and Komatsu, Y. (2013). Silent synapses persist into adulthood in layer 2/3 pyramidal neurons of visual cortex in dark-reared mice. *J. Neurophysiol.* 109, 2064–2076. doi: 10.1152/jn.00912.2012
- Gusnard, D., Akbudak, E., Shulman, G., and Raichle, M. (2001). Medial prefrontal cortex and self-referential mental activity: relation to a default mode of brain function. *Proc. Natl. Acad. Sci. U.S.A.* 98, 4259–4264. doi: 10.1073/pnas.071043098
- Haker, H., Schneebeli, M., and Stephan, K. (2016). Can bayesian theories of autism spectrum disorder help improve clinical practice? *Front. Psychiatry* 7:107. doi: 10.3389/fpsy.2016.00107
- Harmon-Jones, E., Gable, P., and Peterson, C. (2010). The role of asymmetric frontal cortical activity in emotion-related phenomena: a review and update. *Biol. Psychol.* 84, 451–462. doi: 10.1016/J.BIOPSYCHO.2009.08.010
- Hayes, S., and Smith, S. (2005). *Get Out of Your Mind and Into Your Life: The New Acceptance and Commitment Therapy*. Oakland, CA: New Harbinger Publications.
- Hesp, C., Smith, R., Allen, M., Friston, K., and Ramstead, M. (2019). Deeply felt affect: the emergence of valence in deep active inference. *PsyArXiv* [Preprint]. doi: 10.31234/osf.io/62pfd
- Hietanen, J., Glerean, E., Hari, R., and Nummenmaa, L. (2016). Bodily maps of emotions across child development. *Develop. Sci.* 19, 1111–1118. doi: 10.1111/desc.12389
- Hobson, J., and Friston, K. (2012). Waking and dreaming consciousness: neurobiological and functional considerations. *Prog. Neurobiol.* 98, 82–98. doi: 10.1016/j.pneurobio.2012.05.003
- Hohwy, J. (2016). The self-evidencing brain. *Noûs* 50, 259–285. doi: 10.1111/nous.12062
- Joffily, M., and Coricelli, G. (2013). Emotional valence and the free-energy principle. *PLoS Comput. Biol.* 9:e1003094. doi: 10.1371/journal.pcbi.1003094
- Kashdan, T., Barrett, L., and McKnight, P. (2015). Unpacking emotion differentiation: transforming unpleasant experience by perceiving distinctions in negativity. *Curr. Direct. Psychol. Sci.* 24, 10–16. doi: 10.1177/0963721414550708
- Kashdan, T., and Farmer, A. (2014). Differentiating emotions across contexts: comparing adults with and without social anxiety disorder using random, social interaction, and daily experience sampling. *Emotion* 14, 629–638. doi: 10.1037/a0035796
- Kerchner, G., and Nicoll, R. (2008). Silent synapses and the emergence of a postsynaptic mechanism for LTP. *Nat. Rev. Neurosci.* 9, 813–825. doi: 10.1038/nrn2501
- Kleckner, I., Zhang, J., Touroutoglou, A., Chanes, L., Xia, C., Simmons, W., et al. (2017). Evidence for a large-scale brain system supporting allostasis and interoception in humans. *Nature Hum. Behav.* 1:0069. doi: 10.1038/s41562-017-0069
- Lackner, J. (2005). "Is IBS a problem of emotion dysregulation? Testing the levels of emotional awareness model," in *Proceeding of the Presented at the Annual Meeting of the American Psychosomatic Society*, McLean, VA.
- Lane, R., Anderson, F., and Smith, R. (2018). Biased competition favoring physical over emotional pain: a possible explanation for the link between early adversity and chronic pain. *Psychosom. Med.* 80, 880–890. doi: 10.1097/PSY.0000000000000640
- Lane, R., Fink, G., Chua, P., and Dolan, R. (1997). Neural activation during selective attention to subjective emotional responses. *Neuroreport* 8, 3969–3972. doi: 10.1097/00001756-199712220-00024
- Lane, R., Quinlan, D., Schwartz, G., Walker, P., and Zeitlin, S. (1990). The levels of emotional awareness scale: a cognitive-developmental measure of emotion. *J. Pers. Assess.* 55, 124–134. doi: 10.1080/00223891.1990.9674052
- Lane, R., and Schwartz, G. (1987). Levels of emotional awareness: a cognitive-developmental theory and its application to psychopathology. *Am. J. Psychiatry* 144, 133–143. doi: 10.1176/ajp.144.2.133
- Lane, R., Sechrest, L., Reidel, R., Weldon, V., Kaszniak, A., and Schwartz, G. (1996). Impaired verbal and nonverbal emotion recognition in alexithymia. *Psychosom. Med.* 58, 203–210. doi: 10.1097/00006842-199605000-00002
- Lane, R., Sechrest, L., Riedel, R., Shapiro, D., and Kaszniak, A. (2000). Pervasive emotion recognition deficit common to alexithymia and the repressive coping style. *Psychosom. Med.* 62, 492–501. doi: 10.1097/00006842-200007000-00007
- Lane, R., Weihs, K., Herring, A., Hishaw, A., and Smith, R. (2015). Affective agnosia: expansion of the alexithymia construct and a new opportunity to integrate and extend Freud's legacy. *Neurosci. Biobehav. Rev.* 55, 594–611. doi: 10.1016/j.neubiorev.2015.06.007
- Lawson, R., Mathys, C., and Rees, G. (2017). Adults with autism overestimate the volatility of the sensory environment. *Nat. Neurosci.* 20, 1293–1299. doi: 10.1038/nn.4615
- Lawson, R., Rees, G., and Friston, K. (2014). An aberrant precision account of autism. *Front. Hum. Neurosci.* 8:302. doi: 10.3389/fnhum.2014.00302
- Levine, D., Marziali, E., and Hood, J. (1997). Emotion processing in borderline personality disorders. *J. Nerv. Mental Dis.* 185, 240–246.
- Licata, M., Kristen, S., and Sodian, B. (2016). Mother-child interaction as a cradle of theory of mind: the role of maternal emotional availability. *Soc. Develop.* 25, 139–156. doi: 10.1111/sode.12131

- Limanowski, J., and Friston, K. (2018). "Seeing the Dark": grounding phenomenal transparency and opacity in precision estimation for active inference. *Front. Psychol.* 9:643. doi: 10.3389/fpsyg.2018.00643
- Lindquist, K., and Barrett, L. (2008). Constructing emotion: the experience of fear as a conceptual act. *Psychol. Sci.* 19, 898–903. doi: 10.1111/j.1467-9280.2008.02174.x
- Lumley, M., Schubiner, H., Lockhart, N., Kidwell, K., Harte, S., Clauw, D., et al. (2017). Emotional awareness and expression therapy, cognitive behavioral therapy, and education for fibromyalgia: a cluster-randomized controlled trial. *Pain* 158, 2354–2363. doi: 10.1097/j.pain.0000000000001036
- Martínez, M., Villar, M. E., Ballarini, F., and Viola, H. (2014). Retroactive interference of object-in-context long-term memory: role of dorsal hippocampus and medial prefrontal cortex. *Hippocampus* 24, 1482–1492. doi: 10.1002/hipo.22328
- McKay, R., and Dennett, D. (2009). The evolution of misbelief. *Behav. Brain Sci.* 32, 493–510. doi: 10.1017/S0140525X09990975
- McRae, K., Reiman, E., Fort, C., Chen, K., and Lane, R. (2008). Association between trait emotional awareness and dorsal anterior cingulate activity during emotion is arousal-dependent. *Neuroimage* 41, 648–655. doi: 10.1016/j.neuroimage.2008.02.030
- Metzinger, T. (2017). "The problem of mental action predictive control without sensory sheets," in *Philosophy and Predictive Processing*, eds T. Metzinger, and W. Wiese. (Hong Kong: MIND Group).
- Mirza, M., Adams, R., Mathys, C., and Friston, K. (2016). Scene construction, visual foraging, and active inference. *Front. Comput. Neurosci.* 10:56. doi: 10.3389/fncom.2016.00056
- Moeller, S., Konova, A., Parvaz, M., Tomasi, D., Lane, R., Fort, C., et al. (2014). Functional, structural, and emotional correlates of impaired insight in cocaine addiction. *JAMA Psychiatry* 71, 61–70. doi: 10.1001/jamapsychiatry.2013.2833
- Mogg, K., and Bradley, B. (2005). Attentional bias in generalized anxiety disorder versus depressive disorder. *Cogn. Ther. Res.* 29, 29–45. doi: 10.1007/s10608-005-1646-y
- Mumme, D., Fernald, A., and Herrera, C. (1996). Infants' responses to facial and vocal emotional signals in a social referencing paradigm. *Child Develop.* 67, 3219–3237. doi: 10.1111/j.1467-8624.1996.tb01910.x
- Panksepp, J., and Biven, L. (2012). *The Archaeology of Mind: Neuroevolutionary Origins of Human Emotions*. New York, NY: W.W. Norton & Company.
- Panksepp, J., Lane, R., Solms, M., and Smith, R. (2017). Reconciling cognitive and affective neuroscience perspectives on the brain basis of emotional experience. *Neurosci. Biobehav. Rev.* 76, 187–215. doi: 10.1016/j.neubiorev.2016.09.010
- Parker, J., Taylor, G., and Bagby, R. (2003). The 20-item Toronto alexithymia scale: III. Reliability and factorial validity in a community population. *J. Psychosom. Res.* 55, 269–275. doi: 10.1016/S0022-3999(02)00578-0
- Parr, T., and Friston, K. (2017a). Working memory, attention, and salience in active inference. *Sci. Rep.* 7:14678. doi: 10.1038/s41598-017-15249-0
- Parr, T., and Friston, K. (2017b). Uncertainty, epistemics and active inference. *J. R. Soc. Interf.* 14:20170376. doi: 10.1098/rsif.2017.0376
- Parr, T., and Friston, K. (2018). The anatomy of inference: generative models and brain structure. *Front. Comput. Neurosci.* 12:90. doi: 10.3389/fncom.2018.00090
- Pears, K., and Fisher, P. (2005). Emotion understanding and theory of mind among maltreated children in foster care: evidence of deficits. *Develop. Psychopathol.* 17, 47–65. doi: 10.1017/S0954579405050030
- Peters, A., McEwen, B. S., and Friston, K. (2017). Uncertainty and stress: why it causes diseases and how it is mastered by the brain. *Prog. Neurobiol.* 156, 164–188. doi: 10.1016/j.pneurobio.2017.05.004
- Posner, M. (2016). Orienting of attention: then and now. *Q. J. Exp. Psychol.* 69, 1864–1875. doi: 10.1080/17470218.2014.937446
- Rizzolatti, G., Riggio, L., Dascola, I., and Umiltà, C. (1987). Reorienting attention across the horizontal and vertical meridians: evidence in favor of a premotor theory of attention. *Neuropsychologia* 25, 31–40. doi: 10.1016/0028-3932(87)90041-8
- Russell, J. (1991). Culture and the categorization of emotions. *Psychol. Bull.* 110, 426–450. doi: 10.1037/0033-2909.110.3.426
- Russell, J. (2003). Core affect and the psychological construction of emotion. *Psychol. Rev.* 110, 145–172. doi: 10.1037/0033-295x.110.1.145
- Scherer, K. (2009). The dynamic architecture of emotion: Evidence for the component process model. *Cogn. Emot.* 23, 1307–1351. doi: 10.1080/02699930902928969
- Schwarzenbeck, P., FitzGerald, T., Mathys, C., Dolan, R., and Friston, K. (2015). The dopaminergic midbrain encodes the expected certainty about desired outcomes. *Cereb. Cortex* 25, 3434–3445. doi: 10.1093/cercor/bh u159
- Schwarzenbeck, P., and Friston, K. (2016). Computational phenotyping in psychiatry: a worked example. *ENeuro* 3: ENEURO.0049-16.2016. doi: 10.1523/ENeuro.0049-16.2016
- Seth, A. (2013). Interoceptive inference, emotion, and the embodied self. *Trends Cogn. Sci.* 17, 565–573. doi: 10.1016/j.tics.2013.09.007
- Seth, A., and Friston, K. (2016). Active interoceptive inference and the emotional brain. *Philos. Trans. R. Soc. Lond. B* 371:20160007. doi: 10.1098/rstb.2016.0007
- Sharot, T. (2011). The optimism bias. *Curr. Biol.* 21, R941–R945. doi: 10.1016/J.CUB.2011.10.030
- Siemer, M., Mauss, I., and Gross, J. (2007). Same situation–different emotions: how appraisals shape our emotions. *Emotion* 7, 592–600. doi: 10.1037/1528-3542.7.3.592
- Silani, G., Bird, G., Brindley, R., Singer, T., Frith, C., and Frith, U. (2008). Levels of emotional awareness and autism: an fMRI study. *Psychology* 3, 97–112. doi: 10.1080/17470910701577020
- Smith, D., and Schenk, T. (2012). The premotor theory of attention: time to move on? *Neuropsychologia* 50, 1104–1114. doi: 10.1016/J.NEUROPSYCHOLOGIA.2012.01.025
- Smith, R. (2019). "The three-process model of implicit and explicit emotion," in *Neuroscience of Enduring Change: Implications for Psychotherapy*, eds R. Lane, and L. Nadel (Oxford: Oxford University Press).
- Smith, R., Akozei, A., Killgore, W., and Lane, R. (2017a). Nested positive feedback loops in the maintenance of major depression: an integration and extension of previous models. *Brain Behav. Immun.* 67, 374–397. doi: 10.1016/j.bbi.2017.09.011
- Smith, R., Akozei, A., Bao, J., Smith, C., Lane, R., and Killgore, W. (2017b). Resting state functional connectivity correlates of emotional awareness. *NeuroImage* 159, 99–106. doi: 10.1016/j.neuroimage.2017.07.044
- Smith, R., Lane, R., Akozei, A., Bao, J., Smith, C., Sanova, A., et al. (2017c). Maintaining the feelings of others in working memory is associated with activation of the left anterior insula and left frontal-parietal control network. *Soc. Cogn. Affect. Neurosci.* 12, 848–860. doi: 10.1093/scan/nsx011
- Smith, R., Thayer, J., Khalsa, S., and Lane, R. (2017d). The hierarchical basis of neurovisceral integration. *Neurosci. Biobehav. Rev.* 75, 274–296. doi: 10.1016/j.neubiorev.2017.02.003
- Smith, R., Bajaj, S., Dailey, N., Akozei, A., Smith, C., Sanova, A., et al. (2018a). Greater cortical thickness within the limbic visceromotor network predicts higher levels of trait emotional awareness. *Conscious. Cogn.* 57, 54–61. doi: 10.1016/j.concog.2017.11.004
- Smith, R., Killgore, W., and Lane, R. (2018b). The structure of emotional experience and its relation to trait emotional awareness: a theoretical review. *Emotion* 18, 670–692. doi: 10.1037/emo0000376
- Smith, R., Lane, R., Akozei, A., Bao, J., Smith, C., Sanova, A., et al. (2018c). The role of medial prefrontal cortex in the working memory maintenance of one's own emotional responses. *Sci. Rep.* 8:3460. doi: 10.1038/s41598-018-21896-8
- Smith, R., Lane, R., Sanova, A., Akozei, A., Smith, C., and Killgore, W. W. D. (2018d). Common and unique neural systems underlying the working memory maintenance of emotional vs. bodily reactions to affective stimuli: the moderating role of trait emotional awareness. *Front. Hum. Neurosci.* 12:370. doi: 10.3389/fnhum.2018.00370
- Smith, R., Sanova, A., Akozei, A., Lane, R., and Killgore, W. (2018e). Higher levels of trait emotional awareness are associated with more efficient global information integration throughout the brain: a graph-theoretic analysis of resting state functional connectivity. *Soc. Cogn. Affect. Neurosci.* 13, 665–675. doi: 10.1093/scan/nsy047
- Smith, R., Braden, B., Chen, K., Ponce, F., Lane, R., and Baxter, L. (2015). The neural basis of attaining conscious awareness of sad mood. *Brain Imaging Behav.* 9, 574–587. doi: 10.1007/s11682-014-9318-8

- Smith, R., Fass, H., and Lane, R. (2014). Role of medial prefrontal cortex in representing one's own subjective emotional responses: a preliminary study. *Conscious. Cogn.* 29, 117–130. doi: 10.1016/j.concog.2014.08.002
- Smith, R., Kaszniak, A., Katsanis, J., Lane, R., and Nielsen, L. (2019a). The importance of identifying underlying process abnormalities in alexithymia: implications of the three-process model and a single case study illustration. *Conscious. Cogn.* 68, 33–46. doi: 10.1016/J.CONCOG.2018.12.004
- Smith, R., Lane, R., Parr, T., and Friston, K. (2019b). Neurocomputational mechanisms underlying emotional awareness: insights afforded by deep active inference and their potential clinical relevance. *Neurosci. Biobehav. Rev.* 107, 473–491. doi: 10.1101/681288
- Smith, R., Quinlan, D., Schwartz, G. E., Sanova, A., Alkozei, A., and Lane, R. D. (2019c). Developmental contributions to emotional awareness. *J. Pers. Assess.* 101, 150–158. doi: 10.1080/00223891.2017.1411917
- Smith, R., Schwartenbeck, P., Parr, T., and Friston, K. (2019d). An active inference approach to modeling structure learning: concept learning as an example case. *BioRxiv* [Preprint]. doi: 10.1101/633677
- Smith, R., Weihs, K., Alkozei, A., Killgore, W., and Lane, R. (2019e). An embodied neurocomputational framework for organically integrating biopsychosocial processes: an application to the role of social support in health and disease. *Psychosom. Med.* 81, 125–145. doi: 10.1097/PSY.0000000000000661
- Smith, R., and Lane, R. (2015). The neural basis of one's own conscious and unconscious emotional states. *Neurosci. Biobehav. Rev.* 57, 1–29. doi: 10.1016/j.neubiorev.2015.08.003
- Smith, R., and Lane, R. (2016). Unconscious emotion: a cognitive neuroscientific perspective. *Neurosci. Biobehav. Rev.* 69, 216–238. doi: 10.1016/j.neubiorev.2016.08.013
- Subic-Wrana, A., Beetz, M., Paulussen, J., Wiltnik, J., and Beutel, M. (2007). "Relations between attachment, childhood trauma, and emotional awareness in psychosomatic inpatients," in *Proceedings of the Presented at the Annual Meeting of the American Psychosomatic Society*. Budapest.
- Subic-Wrana, C., Bruder, S., Thomas, W., Lane, R., and Köhle, K. (2005). Emotional awareness deficits in inpatients of a psychosomatic ward: a comparison of two different measures of alexithymia. *Psychosom. Med.* 67, 483–489. doi: 10.1097/01.psy.0000160461.19239.13
- Suvak, M., Litz, B., Sloan, D., Zanarini, M., Barrett, L., and Hofmann, S. (2011). Emotional granularity and borderline personality disorder. *J. Abnorm. Psychol.* 120, 414–426. doi: 10.1037/a0021808
- Tononi, G., and Cirelli, C. (2014). Sleep and the price of plasticity: from synaptic and cellular homeostasis to memory consolidation and integration. *Neuron* 81, 12–34. doi: 10.1016/J.NEURON.2013.12.025
- Widen, S., and Russell, J. (2008). Children acquire emotion categories gradually. *Cogn. Develop.* 23, 291–312. doi: 10.1016/j.cogdev.2008.01.002
- Wilson-Mendenhall, C., Barrett, L., Simmons, W., and Barsalou, L. (2011). Grounding emotion in situated conceptualization. *Neuropsychologia* 49, 1105–1127. doi: 10.1016/j.neuropsychologia.2010.12.032
- Wright, R., Riedel, R., Sechrest, L., Lane, R., and Smith, R. (2017). Sex differences in emotion recognition ability: the mediating role of trait emotional awareness. *Motiv. Emot.* 42, 149–160. doi: 10.1007/s11031-017-9648-0
- Yeo, B., Krienen, F., Sepulcre, J., Sabuncu, M., Lashkari, D., Hollinshead, M., et al. (2011). The organization of the human cerebral cortex estimated by intrinsic functional connectivity. *J. Neurophysiol.* 106, 1125–1165. doi: 10.1152/jn.00338.2011

**Conflict of Interest:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2019 Smith, Parr and Friston. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.