

Trends in Cognitive Sciences

Review

Reward Prediction Error and Declarative Memory

Kate Ergo,¹ Esther De Loof,¹ and Tom Verguts ¹*

Learning based on reward prediction error (RPE) was originally proposed in the context of nondeclarative memory. We postulate that RPE may support declarative memory as well. Indeed, recent years have witnessed a number of independent empirical studies reporting effects of RPE on declarative memory. We provide a brief overview of these studies, identify emerging patterns, and discuss open issues such as the role of signed versus unsigned RPEs in declarative learning.

Memory and Reward Prediction Error

A tennis player knows how to perform a perfect serve and also knows the opponent's name. But how are these two types of 'knowing' similar, if at all? It is thought that the human brain houses at least two broad and distinct memory systems [1], each with its own learning algorithms and neural correlates. The first is **nondeclarative** (or habit, or implicit) **memory** (see Glossary). The second is **declarative** (or, in humans, propositional, or explicit) **memory**. The computational principle of **reward prediction error** (RPE)-based learning [2,3] is generally thought to drive nondeclarative learning. By contrast, until recently, RPE was not studied in the context of declarative memory. However, several empirical studies have reported effects of RPE on declarative memory as well, suggesting that some of the same computational principles shape nondeclarative and declarative memory systems. Here, we review these recent studies and discuss the most important open questions concerning RPEs and declarative memory. First, however, we provide a brief overview of the computational models linking RPE with nondeclarative learning.

RPE and Nondeclarative Learning

One of the most influential theories in current cognitive neuroscience is predictive coding [4,5]. According to this account, the brain generates predictions about its own percepts, actions, and cognition in order to learn about, build models of, and navigate the world [6]. A key concept in predictive coding is the **prediction error** (PE). Specifically, in order to generate accurate predictions, the brain needs to set a number of parameters (e.g., encoded in its synaptic connections). PEs allow updating such parameters.

Predictions can be made about several variables, such as tomorrow's weather, the next action I (or somebody else) will perform, our partner's mood, and so forth. One particularly relevant variable to make predictions about is reward; a PE in reward (by definition) is an RPE. The concept of RPEs has been very influential in nondeclarative learning. In particular, RPEs have been implemented in a wide range of computational models. For example, to account for **blocking** in nondeclarative learning, Rescorla and Wagner (RW; [7]; Box 1) developed their now-classic model according to which learning depends on PE. Specifically, synaptic strength increases when a reward is better than expected, but synaptic strength decreases when the reward is worse than expected. Hence, the valence of the RPE matters [**signed** reward prediction error (SRPE)].

Further computational development of RW led to the temporal difference (TD; Box 1) reinforcement learning model [3]. The TD model improved on the RW model because it

Highlights

The concept of reward prediction error (RPE) learning has been extremely influential in the nondeclarative memory literature, but the role of RPEs in declarative memory has only recently begun to be explored.

Two main approaches to measure RPE in declarative memory have been developed. In the reward-prediction approach, reward is iteratively sampled by the experimenter from a statistical distribution and predicted by the subject. In the multiple-repetition approach, subjects see the same memoranda repeatedly and estimate their probability of success.

Behaviorally, the reward-prediction approach has mainly yielded signed RPE signatures, whereas the multiplerepetition paradigm has mainly yielded unsigned RPE signatures. However, there are exceptions, and neural (electroencephalography, fMRI) paradigms have yielded both in single studies.

An RPE perspective provides testable hypotheses on why reward and reward prediction manipulations may improve or impair declarative learning.

¹Department of Experimental Psychology, Ghent University, Henri Dunantlaan 2, B-9000 Ghent, Belgium

*Correspondence: tom.verguts@ugent.be (T. Verguts).





Box 1. Models of Learning

Rescorla-Wagner Model

This model [7] describes learning the value (expected reward) of specific events (say, events A and B). This information is encoded in their associative strength to a 'value' unit, symbolized as w_A and w_B for events A and B, respectively. Specifically, based on whether events A and B occur ($x_A = 1$ and $x_B = 1$, respectively) or not ($x_A = 0$ and $x_B = 0$, respectively), an additive prediction is made about the occurrence of reward ($V = x_A \times w_A + x_B \times w_B$). When reward finally occurs (or not), a RPE is calculated (R - V), where occurrence of reward (denoted R) is typically coded as R = 0 (when there is no reward) or R = 1 (when there is reward). This RPE is then used to change the connection strength between cells encoding A and B on the one hand and reward on the other: $\Delta w_i = \alpha \times x_i \times (R - V)$, with $i \in \{A, B\}$. After repeated application of this learning rule, the weights w_A and w_B allow the model to accurately predict reward, based on the (A, B) input combination.

Temporal Difference Model

The Rescorla-Wagner model can only learn from external feedback (R - V). This is computationally inefficient because reward may not be delivered at each time point when relevant information is provided to the organism. In temporal difference learning [3], learning can also occur if the prediction of reward changes between two time points *t* and *t* + 1. Formally, the learning rule becomes (now with explicit time index *t*): $\Delta w_i(t) = \alpha \times x(t) \times [R(t + 1) + \gamma V(t + 1) - V(t)]$, with $i \in \{A, B\}$. If $\gamma = 0$, the rule reduces to the Rescorla-Wagner rule. In case $\gamma > 0$, learning can also proceed at times *t* when no actual reward was delivered, rendering the algorithm more powerful than the Rescorla-Wagner rule.

Pearce-Hall Model

According to this model [12], learning only occurs when a reward is surprising. Specifically, it uses the absolute value of an RPE ('different from expected' signal), consistent with an unsigned RPE approach. Formally, (one variant of) the learning rule can be written as: $\Delta w_i(t) = x_i(t) \times R(t) \times R(t) - V(t)$.

allows learning also when the reward is not immediately present. However, the main success of the RPE concept as implemented in TD was probably because of its close match to neurophysiological data. In particular, dopaminergic neurons in the ventral tegmental area (VTA) implement a TD-like RPE signature of reward processing [8,9]. In recent years, the role of TD-based RPEs in nondeclarative learning has become well established in psychology, neuroscience, and artificial intelligence. For example, deep reinforcement learning models use TD-based RPEs to solve tasks (e.g., playing Atari games) that were long considered beyond the capacity of artificial agents [10,11].

In contrast to the RW and TD models that are SRPE based, Pearce and Hall proposed that learning occurs whenever reward is surprising [either better or worse; that is, different from expected, consistent with an **unsigned** reward prediction error (URPE); Box 1] [12]. It is noteworthy that normative, Bayesian models of learning exhibit features of both. For example, the Kalman filter [13] updates its estimates on the basis of SRPEs, but its learning rate [i.e., the extent to which parameters (such as synaptic weights) are updated] is driven by uncertainty, which can be estimated via URPEs [14–16]. Empirical signatures of both SRPE and URPE have been observed in the brain [17]. In summary, the concept of PE, and specifically of RPE, has turned out to be fruitful for understanding nondeclarative learning at neurophysiological, behavioral, and computational levels.

RPE in Declarative Learning

Although the role of RPEs in nondeclarative learning has been studied extensively and formalized in a number of computational models, their role in declarative learning has only recently become a topic of interest. Two main approaches exist for elucidating the RPE effect on declarative learning (for an overview, see Table 1). First, in the reward-prediction approach (Box 2), a statistical distribution determines the probability of reward. The participant knows or estimates this reward distribution. Thus, the participant can make a prediction about reward, and, based on the prediction, an RPE can be generated. Studies using reward prediction can be approximately ordered on the basis of the difficulty of this prediction, and we will discuss

Glossary

Blocking: in a blocking experiment, an event A is consistently followed by an unconditional stimulus (US). Subsequently, an additional event B is added to A (again followed by the US). After conditioning, it is observed that the animal in the experiment has not learned the association between B and the US. In the Rescorla-Wagner rule, the interpretation is that event A blocks this B–US association. Specifically, due to event A, the appearance of the US is no longer surprising.

Declarative memory: memory for facts and events ('knowing what') that can (at least in humans) be (consciously) declared. It is typically considered to consist of episodic memory (memory for single episodes) and semantic memory (memory for information aggregated across several episodes). The process of acquisition of declarative memory is called 'declarative learning.' Encoding declarative memories can happen rapidly, typically after only a single exposure, and relies heavily on the hippocampus [66].

Nondeclarative memory:

nondeclarative learning is an umbrella term for the acquisition of different types of knowledge, including procedural memory ('knowing how'). This involves acquiring a motor or cognitive skill (procedure) by means of repeated practice (e.g., learning to play tennis).

Prediction error (PE): difference

between the actual value of some variable and predicted value of that variable (i.e., actual value minus predicted value).

Reward prediction error (RPE):

prediction error where the relevant variable is reward (i.e., actual reward – predicted reward). See also **Prediction error**.

Signed: in mathematics, signed means that the sign of a number is taken into consideration (e.g., -3, +3). In the context of SRPEs, it indicates that we take the valence (positive versus negative RPEs) into account.

Theta phase synchronization:

synchronization of two brain areas in the theta frequency (4–8 Hz). Such synchronization can be achieved by making the theta phase of the two areas identical so that theta waves in both areas 'go up and down' together.



them in that order (easy to difficult). In one of the first studies to use this approach, each of three cues was linked to a different reward value [18]. A medium reward led to improved recognition when it was better than predicted (i.e., when it was preceded by a cue indicating that low or medium reward was equally likely to follow) relative to when it was worse than predicted (i.e., preceded by a cue indicating that high or medium reward was equally likely to follow), consistent with an SRPE effect. However, later work could not replicate the SRPE effect in this specific experimental paradigm [18,19].

Unsigned: unsigned means that the sign is not considered (i.e., absolute value is taken; for example, -3 and +3 both have an unsigned value of 3). See also **Signed**.

A second implementation of the reward-prediction approach is the recent variable-choice paradigm (Figure 1A, Key Figure, and Box 2) [20]. Here, participants learn Dutch–Swahili word associations under different RPE value conditions. (See [20] for an overview of all RPEs in this design.) Predicting the reward probability is again quite easy; participants can deduce it from the number of eligible options. Behaviorally, memory performance showed an SRPE effect in declarative learning: Recognition accuracy and certainty increased linearly with larger and more positive RPEs (Figure 1B). These results were replicated with image–word associations [20] and face– word associations [21].

Table 1. Nonexhaustive Overview of Studies on Reward Prediction Error in Declarative Memory

Approach	Task and stimuli	SRPE/URPE	Effect on memory	Refs
Reward prediction	Each of three cues (colored squares) is followed by one of two potential reward values (medium-low, medium-high, and low-high), so a medium reward can be better or worse than expected. After reward feedback, a novel (indoor or outdoor) scene is presented. Scene recognition is probed after a 1-day delay.	SRPE	Positive	[68]
Reward prediction	On each trial, participants see one Dutch word together with four (trial-novel) Swahili words and choose a translation from either one, two, or four of these Swahili words. Manipulating the number of eligible options (one, two, or four) and whether a trial is rewarded allowed manipulation of RPEs. For example, in the case of a four-option rewarded trial, participants experience an RPE of $1 - \frac{1}{4} = 0.75$; in case of a two-option nonrewarded trial, participants experience an RPE of $0 - \frac{1}{2} = -0.50$.	SRPE	Positive	[20] (see also Figure 1A,B)
Reward prediction	A cue is presented with two targets linked to different reward values. Subjects must (learn to predict and) choose the high-value target. Trial-novel images are shown during subsequent reward feedback. Image memory is probed afterwards via old/new judgments.	SRPE	Positive	[22]
Reward prediction	Participants track the reward associated with different indoor and outdoor scenes. On each trial, participants predict the reward (for a particular scene) and subsequently receive feedback about their estimate. From this difference (feedback – predicted reward), an RPE can be calculated. Scene memory is probed after this initial task via old/new judgments.	URPE	Positive	[24] (see also Figure 1C,D)
Reward prediction	On each trial, participants see a value and a stimulus (animate or inanimate) for that trial and decide to play or pass on that trial (Figure 1E). After each choice, the image is shown with reward feedback. Afterwards, recognition memory for the images is probed via old/new judgments.	SRPE	Positive	[40] (see also Figure 1E,F)
Reward prediction	Participants track the drifting reward probability of colored squares, which are overlaid with incidental trial-unique images and followed by feedback. Recognition memory for the images is probed via old/new judgments after a 1-day delay.	SRPE	Negative	[23]
Multiple repetition	Participants are presented with questions for which they have to generate an answer and rate their confidence, followed by a surprise retest.	URPE	Positive	[31]
Multiple repetition	Participants are presented with general information questions. In a first test phase, participants provide answers and rate their confidence. In the subsequent phase, subjects receive feedback about their answers. Finally, participants are retested on a subset of questions in a second test phase.	URPE	Positive	[67]
Multiple repetition	Participants study a text and are tested after 2 days, at which time they also provide confidence ratings for their answers. On a small fraction of trials, participants receive false feedback (i.e., trials that were answered correctly but labelled as false) and receive novel feedback (i.e., a novel 'correct' answer) on those trials. A second (incidental) test is given after 7 days.	URPE	Positive	[32]

Abbreviations: SRPE, signed reward prediction error; URPE, unsigned reward prediction error.



Box 2. How to Generate and Measure RPEs: Experimental Approaches

Reward-Prediction Approach

Here, participants must both learn declarative information (e.g., word pairs) and simultaneously estimate a (potentially nonstationary) reward distribution throughout the task [23,24,40]. In some cases, the correct RPE can be easily derived analytically; in other cases, RPE can only be calculated after fitting a reinforcement learning model and deriving the RPEs from the model estimates [23,40]. One example of a reward-prediction approach is the variable-choice paradigm. In the variable-choice paradigm [20,33] (see Figure 1A in main text), participants learn stimulus pairs, such as Dutch–Swahili word pairs or image–Swahili stimulus pairs [20]. In the former example, on each trial, a Dutch word is shown together with four Swahili words. Critically, the number of eligible options is manipulated. In the one-option, two-option, and four-option conditions, one, two, or four Swahili words are eligible (framed), respectively, and the probability of choosing the correct translation is thus 100%, 50%, or 25%, respectively. Feedback is given on every trial. Signed and unsigned trial-by-trial RPEs are calculated on the basis of the difference between actual and predicted reward (see Glossary). Memory is probed in a subsequent recognition test.

Multiple-Repetition Approach

Here, general information questions are repeatedly presented, and an RPE is estimated on the basis of previous presentations of each question. For example, in [32], participants first studied a text and subsequently received (multiple-choice) questions about the text. After each question, they rated confidence and received feedback. The trial-by-trial PE was calculated using the confidence rating and feedback. Hypercorrection effect studies also typically use a multiple-repetition paradigm [29,67].

In another instantiation of the reward-prediction approach, participants actively track and estimate the reward probability distribution. Here, on each trial, they experience an RPE relative to that (estimated) distribution (Figure 1C,D, and Box 2) [22–24]. On the basis of this feedback, participants can update their estimate for subsequent trial estimates. For example, in one study [22], participants estimated the (fixed) probability of reward attached to specific stimuli. At reward feedback, a trial-novel image was presented. Subsequent memory performance for these trialnovel images displayed an SRPE effect, which was more pronounced in adolescents than in adults. In another study [24], participants tracked the reward associated with different indoor and outdoor scenes. Here, a clear URPE effect was observed: scenes associated with a higher URPE during the initial task (i.e., with more surprising rewards in either positive or negative direction) were afterwards better remembered (Figure 1C,D).

A more challenging study [25] used a reward-prediction paradigm to disentangle effects of SRPE, surprise (which corresponds to URPE), and uncertainty. Unlike in the other paradigms just discussed, reward probability was not fixed but instead jumped to a different level at unpredictable time points in the experiment. Only SRPE had an effect on subsequent memory (Figure 1F; see also [26]). Finally, in the study by Wimmer *et al.* [23], the reward probability would fluctuate slowly but unpredictably on each trial, making the reward-prediction task very challenging. In this experiment, unlike the other discussed paradigms, a negative effect of (S) RPE was observed. Specifically, trials (and participants) with stronger and more positive RPEs were associated with impaired declarative learning.

As a second approach, in a multiple-repetition paradigm (Box 2), a set of general information questions are repeated a number of times. Trial-specific confidence ratings ('How certain are you that you answered correctly?') and feedback are used to compute trial-specific PEs. Given that being correct is rewarding [25], these PEs can be considered RPEs. The researchers use these RPEs to predict accuracy on subsequent presentations of the same general information questions. Here, a URPE effect is typically observed. In particular, the hypercorrection effect obtained in this multiple-repetition paradigm entails that errors made with high confidence are beneficial for memory [27–31]. High-confidence errors occur on those trials during which positive feedback was expected but not obtained; thus, this effect is consistent with a URPE effect. Another experiment [32] also showed a hypercorrection effect, which the authors interpreted



Key Figure

Reward Prediction Error in Declarative Memory



(See figure legend at the bottom of the next page.)



as a URPE. Additionally, in this second experiment, participants received false feedback on a small fraction of trials (i.e., trials that were answered correctly but labelled as false) and received novel feedback (i.e., a novel 'correct' answer) on those trials. In those false-feedback trials, a URPE effect was also observed: On trials that were answered with high certainty but that were not rewarded (high URPE), the novel feedback was subsequently recalled more confidently.

Overviewing and categorizing these paradigms, we note that a main difference between the reward-prediction and multiple-repetition approaches is the origin of the RPE: an independent reward generation mechanism in the former and the participant's own confidence in his or her memory in the latter. Another difference is that, in the reward-prediction approach, RPEs are usually computed or estimated, whereas RPEs are deduced from confidence measures in the multiple-repetition approach. There are some exceptions to the latter rule. For example, Rouhani *et al.* [24] implemented a reward-prediction paradigm where confidence is used to calculate an RPE. Finally, in the reward-prediction paradigm, memoranda are usually trial-unique, whereas (by definition) they are not in the multiple-repetition approach. These are just a few of the relevant dimensions; we discuss some other potentially relevant dimensions in the next section.

Open Issues

Despite growing evidence that RPEs that drive declarative memory as well as nondeclarative memory, many uncertainties remain. We discuss a few of them in the next subsections.

RPE: Signed or Unsigned?

Studies with a multiple-repetition paradigm typically observed URPE (i.e., surprise) effects. By contrast, the reward-prediction paradigm has tended to yield SRPE effects, although URPE effects have occasionally been documented as well (Figure 1D) [24]. Why do different designs tend to generate SRPE versus URPE effects on declarative learning? One potentially relevant factor is the range of the RPEs probed. In particular, studies that found a behavioral SRPE effect (i.e., most reward-prediction paradigms) might simply not have investigated the full range of RPEs. For example, in the variable-choice paradigm [20,33], the maximal SRPE equaled 0.75, whereas the minimal SRPE was –0.50 (hence, lower in absolute value). It is possible that a larger range of negative RPEs might lead to a URPE effect. This could be tested by including a few nonrewarded one-option (high-certainty) trials. These highly infrequent events would be accompanied by large negative RPEs.

However, this is unlikely to be the full story, because both RPE signatures have been observed even within a single study. In an electroencephalographic (EEG) study with the variable-choice paradigm [33], a URPE pattern was observed during reward feedback in the theta (4–8 Hz) frequency band, consistent with literature implicating theta in URPE processing [34]. By contrast, SRPE signatures were found in the high-beta (20–30 Hz) and high-alpha (10–15 Hz) frequency ranges, consistent with a functional role of both beta and alpha power in reward feedback processing [35,36]. Furthermore, in an fMRI study using a multiple-repetition paradigm [32], SRPE-consistent activation was found in several areas (including striatum), but URPE signatures were found in others (including insula). Together, these findings suggest that both SRPE and URPE are important for declarative learning and that we need an account identifying the functional

Figure 1. Reward-prediction approach applied in three paradigms and typical findings. (A) Variable-choice paradigm [20]. (B) Variable-choice paradigm behavioral results [20] show a signed reward prediction error (SRPE) signature for recognition in both the immediate and delayed test groups. Recognition of word pairs increased linearly with larger and more positive reward prediction errors (RPEs). (C) Paradigm reproduced from [24]. (D) Rouhani *et al.* [24] found a unsigned RPE signature, with memory improving for both large negative and large positive RPEs. (E) Paradigm reproduced from [40]. (F) Jang *et al.* [40] found an SRPE signature: memory score increased with increasing RPE.



role of each in time, (neural) space, and frequency band. The Bayesian learning model mentioned in the introduction, which naturally incorporates both, may be a useful starting point in this respect. Specifically, as this model suggests, it may be that URPE drives learning rate, SRPE drives update, and their combination (learning rate × update) determines a PE that drives declarative learning.

Timing Issues of RPEs

In most paradigms, a novel declarative memorandum is presented on each trial, followed by an RPE, followed by declarative feedback about what the correct answer should have been (see Figure 1A, word pair encoding, for an example). Here, RPE can have either a retrograde effect (if it interacts with the originally presented memoranda), or instead an anterograde effect (if it interacts with the declarative feedback). Concerning the anterograde effect, in studies using the variable-choice paradigm, the declarative feedback appeared either simultaneously with the RPE (delay of 0 ms [20]) or with a delay of 3000 ms [33]. The fact that we find very similar results in the two cases suggests that the timing of the RPE-feedback interval may not be crucial, at least within the first few hundreds of milliseconds. An interesting parallel can be drawn here with the test-potentiated learning effect from the declarative memory literature. Here, taking a test potentiates the learning of (old or novel) material that is subsequently presented [37,38]. Also, for a retrograde effect (of RPE on originally presented memorandum), an interesting analogy can be made with earlier literature. In particular, Braun et al. [39] found a retrograde effect of reward on declarative memory, with objects that were (temporarily) closer to (subsequent) reward being better remembered afterwards. In the reward-prediction approach, it remains to be shown which of these two (anterograde or retrograde effect of RPE) is crucial for driving the RPEbased declarative memory improvement.

An RPE can also appear at cue rather than at feedback. Only a single paper thus far has investigated both cue- and feedback-locked RPE effects [40]. Those authors observed cue- but not feedback-locked RPE effects; however, in their experiment, there was both a cue- and a feedback-locked RPE on each trial. It is very well possible that an initial RPE suppresses a second RPE occurring (e.g., a few hundred milliseconds later) in that same trial. We conclude that RPE timing issues need to be studied more systematically. In particular, if this research is to have practical application in education, such studies will be imperative.

RPE: Why and How?

In nondeclarative learning, a normative argument for why RPE is useful is well established: Calculating RPE is necessary for online (i.e., while interacting with the world) reward maximization [3]; this idea is inherent in the RW, TD, and Pearce-Hall models (Box 1). Does this argument apply to declarative memory as well? An intuitive argument is that it makes sense to only remember stimuli (or more generally, episodes) that are associated with a reward level that is sufficiently different from what is already expected. Indeed, if a stimulus from some category is accompanied by reward each time it is encountered, it makes little sense to explicitly remember each novel stimulus instance as a separate event once it has already been learned.

Another issue is how RPE improves memory. One potential mechanism is via phase-locking to neural oscillations in specific frequency bands. In particular, neural **theta phase synchronization** may provide one (but not an exclusive) solution: Brain areas in theta phase synchrony are thought to communicate and learn more efficiently [41], thus facilitating memory integration [42]. Indeed, episodic memory is enhanced when multimodal (audiovisual) stimuli are synchronously presented in theta phase, with stronger theta phase synchronization predicting better memory performance [43,44]. Dopaminergic midbrain neurons have also been found to phase-lock to (cortical) theta



during encoding, with stronger phase-locking during subsequently remembered (versus forgotten) memoranda [45]. Thus, it is possible that RPEs (via neuromodulatory signaling) increase theta synchrony, which subsequently allows the relevant brain areas to 'glue' the episode together more efficiently [46]. The EEG variable-choice paradigm study mentioned above [33] provides preliminary evidence for this view. Further, computational models that consider RPE-theta interactions to drive learning have started to appear [47].

Whereas dopaminergic RPEs likely support nondeclarative learning via basal ganglia pathways, dopaminergic RPEs may support declarative memory via the hippocampus [48]. Standard theory holds that (dopaminergic) VTA calculates SRPE, but a substantial number of URPE neurons have also been observed in VTA and nearby midbrain areas [49]. Moreover, also noradrenergic locus coeruleus projects to the hippocampus and may thus exert URPE effects [50]. Earlier authors proposed that VTA-hippocampus interactions originate in the hippocampus [51]. We propose that VTA-hippocampus interactions may also originate in VTA and that SRPEs (encoded by VTA, possibly based on input from ventral striatum [52]) and URPEs (encoded in VTA and locus coeruleus) may modulate the hippocampus for episodic memory encoding. Consistently, a number of studies have demonstrated that midbrain VTA activation (triggered by reward or by RPE) is associated with improved episodic learning [21,53,54].

Effect of Test Delay on Declarative Memory

In declarative memory studies, participants are typically subjected to an implicit or explicit memory test, either on the same day or after a considerable delay (ranging from a few hours to a few weeks). If, as suggested above, SRPEs are encoded by dopaminergic neurons, then effects should be stronger with longer delays. Indeed, although early and late long-term memory effects both rely on dopamine, late effects have a stronger dependency on dopamine [48]. Consistently, an effect of reward in declarative learning is typically stronger after a delay [55,56]. However, a systematic comparison of the delay-by-RPE interaction on declarative memory remains to be carried out.

Reconsolidation

When information is retrieved from memory, it enters a plastic, labile state, allowing the information to be changed, strengthened, or weakened, a process called 'reconsolidation' [57,58]. Reconsolidation is most intensively studied in nondeclarative memory [59], but it is observed in declarative memory as well [60]. PE is required for reconsolidation [61], both in nondeclarative memory [62] and in declarative memory [60,63]. Given the important role of RPE in declarative learning, and given that similar principles drive learning and reconsolidation [63], we predict that RPE may modulate reconsolidation, too. The multiple-repetition approach, where declarative memory is probed iteratively, can be considered as a first attempt at investigating RPEs in the context of reconsolidation. This remains, however, to be further investigated.

Concluding Remarks

Learning, RPEs, and declarative memory are sometimes treated as separate topics, each with their own prominent paradigms, findings, and theories. The current perspective suggests instead that they are intimately related. Briefly, learning is modulated by RPEs and leads to (declarative) memory traces in the brain. We discussed a few recent paradigms that have begun to explore such interactions. In the Open Issues section, we highlighted a number of dimensions of those paradigms that, if addressed, could greatly facilitate further development of the research field. Although much remains to be found out (see Outstanding Question), concrete models and predictions are beginning to emerge, with relevance for both natural and artificial intelligence. Concerning the latter, it is of interest that recent deep neural networks integrate nondeclarative

Outstanding Questions

Under which circumstances (e.g., experimental paradigm, population, learning-test delay) is an SRPE versus URPE signature in declarative memory observed?

Most studies observe a beneficial effect of (signed or unsigned) RPEs, but RPEs may also impair subsequent memory. Under which circumstances are RPEs beneficial versus harmful for learning?

Can we identify computational principles from the TD learning model in episodic memory? In particular, the TD model predicts that learning occurs when the reward prediction changes across successive time steps (even in the absence of actual reward).

Can the differentiation predicted by Bayesian models of (nondeclarative) learning (SRPE and URPE drive updates and learning rates, respectively) be observed in declarative learning?

What is the effect of test delay on declarative memory? Does the RPE effect increase or decrease after sleep?

Does RPE support memory via theta phase synchronization?

What is the relationship between RPE and testing effect for improving declarative memory?

Can the concept of RPE advance learning in education contexts?

How do reward-related disorders and ageing influence the RPE effect on declarative memory?

What are commonalities and differences (e.g., in neural structures, behavioral commonalities) between RPEs and more general PEs?

Do RPEs have differential effects when generated by the agent itself or by another agent? when they are generated for the agent or for another one? learning (as in standard neural networks) with declarative memory [64,65]. In such artificial systems, RPEs may determine when (or how strongly) to store a declarative memorandum. We are excited about what the (near) future will bring in this domain, not only because of its conceptual unification but also because of its promise for informing education policy and practice.

Acknowledgments

K.E. and T.V. were supported by project G012816N from the Research Council Flanders (FWO). K.E. is a research fellow at the Research Council Flanders. E.D.L. and T.V. were supported by project BOF17-GOA-004 awarded by the Ghent University Research Council.

References

- Squire, L.R. (2004) Memory systems of the brain: a brief history and current perspective. Neurobiol. Learn. Mem. 82, 171–177
- Wang, J.X. et al. (2018) Prefrontal cortex as a metareinforcement learning system. Nat. Neurosci. 21, 860–868
- Sutton, R.S. and Barto, A.G. (2018) Reinforcement Learning: An Introduction, MIT Press
- Rao, R.P. and Ballard, D.H. (1999) Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nat. Neurosci.* 2, 79–87
- Friston, K.J. (2003) Learning and inference in the brain. Neural Netw. 16, 1325–1352
- Den Ouden, H.E.M. et al. (2012) How prediction errors shape perception, attention, and motivation. Front. Psychol. 3, 548
- Rescorla, R.A. and Wagner, A.R. (1972) A theory of Pavlovian conditioning: variations in the effectiveness of reinforcement and nonreinforcement. In *Classical Conditioning II: Current Research and Theory* (Blake, A.H. and Prokasy, W.F., eds), pp. 64–99, Appleton-Century-Croft
- Ljungberg, T. et al. (1992) Responses of monkey dopamine neurons during learning of behavioral reactions. J. Neurophysiol. 67, 145–163
- Eshel, N. *et al.* (2016) Dopamine neurons share common response function for reward prediction error. *Nat. Neurosci.* 19, 479–486
- Mnih, V. et al. (2015) Human-level control through deep reinforcement learning. Nature 518, 529–533
- Silver, D. et al. (2016) Mastering the game of Go with deep neural networks and tree search. Nature 529, 484–489
- Pearce, J.M. and Hall, G. (1980) A model for Pavlovian learning: variations in the effectiveness of conditioned but not of unconditioned stimuli. *Psychol. Rev.* 87, 532–552
- Dayan, P. et al. (2000) Learning and selective attention. Nat. Neurosci. 3 Suppl, 1218–1223
- 14. Behrens, T.E.J. *et al.* (2007) Learning the value of information in an uncertain world. *Nat. Neurosci.* 10, 1214–1221
- Silvetti, M. et al. (2018) Dorsal anterior cingulate-brainstem ensemble as a reinforcement meta-learner. PLoS Comput. Biol. 14, e1006370
- Courville, A.C. et al. (2006) Bayesian theories of conditioning in a changing world. Trends Cogn. Sci. 10, 294–300
- Roesch, M.R. et al. (2012) Surprise! Neural correlates of Pearce Hall and Rescorla – Wagner coexist within the brain. *Eur. J. Neurosci.* 35, 1190–1200
- Mason, A. et al. (2017) Adaptive scaling of reward in episodic memory: a replication study. Q. J. Exp. Psychol. 70, 2306–2318
- Mason, A. et al. (2017) The role of reward and reward uncertainty in episodic memory. J. Mem. Lang. 96, 62–77
- De Loof, E. et al. (2018) Signed reward prediction errors drive declarative learning. PLoS One 13, e0189212
- Buc Calderon, C. et al. (2020) Signed reward prediction errors in the ventral striatum drive episodic memory. *BioRxiv* 2020.01.03.893578
- Davidow, J.Y. et al. (2016) An upside to reward sensitivity: the hippocampus supports enhanced reinforcement learning in adolescence. Neuron 92, 93–99
- Wimmer, G.E. et al. (2014) Episodic memory encoding interferes with reward learning and decreases striatal prediction errors. J. Neurosci. 34, 14901–14912
- Rouhani, N. *et al.* (2018) Dissociable effects of surprising rewards on learning and memory. *J. Exp. Psychol. Learn. Mem. Cogn.* 44, 1430–1443

- Satterthwaite, T.D. et al. (2012) Being right is its own reward: load and performance related ventral striatum activation to correct responses during a working memory task in youth. *Neuroimage* 61, 723–729
- Aberg, K.C. et al. (2017) Trial-by-trial modulation of associative memory formation by reward prediction error and reward anticipation as revealed by a biologically plausible computational model. Front. Hum. Neurosci. 11, 56
- 27. Metcalfe, J. (2017) Learning from errors. Annu. Rev. Psychol. 68, 465–489
- Metcalfe, J. and Finn, B. (2011) People's hypercorrection of high-confidence errors: did they know it all along? *J. Exp. Psychol. Leam. Mem. Cogn.* 37, 437–448
- 29. Butterfield, B. and Metcalfe, J. (2006) The correction of errors committed with high confidence. *Metacognition Learn.* 1, 69–84
- Fazio, L.K. and Marsh, E.J. (2009) Surprising feedback improves later memory. *Psychon. Bull. Rev.* 16, 88–92
- Butterfield, B. and Metcalfe, J. (2001) Errors committed with high confidence are hypercorrected. J. Exp. Psychol. Learn. Mem. Cogn. 27, 1491–1494
- Pine, A. *et al.* (2018) Knowledge acquisition is governed by striatal prediction errors. *Nat. Commun.* 9, 1673
- Ergo, K. et al. (2019) Oscillatory signatures of reward prediction errors in declarative learning. *Neuroimage* 186, 137–145
- Cavanagh, J.F. and Frank, M.J. (2014) Frontal theta as a mechanism for cognitive control. *Trends Cogn. Sci.* 18, 414–421
- Kleberg, F.I. *et al.* (2014) Ongoing theta oscillations predict encoding of subjective memory type. *Neurosci. Res.* 83, 69–80
- HajiHosseini, A. et al. (2012) The role of beta-gamma oscillations in unexpected rewards processing. *Neuroimage* 60, 1678–1685
 Arnold, K.M. and Mcdermott, K.B. (2013) Test-potentiated
- learning: distinguishing between direct and indirect effects of tests. J. Exp. Psychol. Learn. Mem. Cogn. 39, 940–945
 38. Pastötter, B. and Bäuml, K.T. (2014) Retrieval practice enhances
- Pastotter, B. and Baumi, K. I. (2014) Hetneval practice enhances new learning: the forward effect of testing. *Front. Psychol.* 5, 286
 Braun, E.K. *et al.* (2018) Retroactive and graded prioritization of
- memory by reveal. *Nat. Commun.* 9, 4886
- Jang, A.I. *et al.* (2019) Positive reward prediction errors during decision-making strengthen memory encoding. *Nat. Hum. Behav.* 3, 719–732
- Fries, P. (2015) Rhythms for cognition: communication through coherence. *Neuron* 88, 220–235
- Backus, A.R. et al. (2016) Hippocampal-prefrontal theta oscillations support memory integration. Curr. Biol. 26, 450–457
- Wang, D. *et al.* (2018) Single-trial phase entrainment of theta oscillations in sensory regions predicts human associative memory performance. *J. Neurosci.* 38, 6299–6309
- Clouter, A. *et al.* (2017) Theta phase synchronization is the glue that binds human associative memory. *Curr. Biol.* 27, 3143–3148
- Kaminski, J. *et al.* (2018) Novelty-sensitive dopaminergic neurons in the human substantia nigra predict success of declarative memory formation. *Curr. Biol.* 28, 1333–1343
- Berens, S.C. and Horner, A.J. (2017) Theta rhythm: temporal glue for episodic memory. *Curr. Biol.* 27, R1110–R1112
- Verbeke, P. and Verguts, T. (2019) Learning to synchronize: how biological agents can couple neural task modules for dealing with the stability-plasticity dilemma. *PLoS Comput. Biol.* 15, e1006604



Trends in Cognitive Sciences



- Lisman, J.E. et al. (2011) A neoHebbian framework for episodic memory: role of dopamine-dependent late LTP. Trends Neurosci. 34, 536–547
- Matsumoto, M. and Hikosaka, O. (2009) Two types of dopamine neuron distinctly convey positive and negative motivational signals. *Nature* 459, 837–841
- Wagatsuma, A. et al. (2018) Locus coeruleus input to hippocampal CA3 drives single-trial learning of a novel context. Proc. Natl. Acad. Sci. U. S. A. 115, E310–E316
- Lisman, J.E. and Grace, A. (2005) The hippocampal-VTA loop: controlling the entry of information into long-term memory. *Neuron* 46, 703–713
- Takahashi, Y.K. et al. (2016) Temporal specificity of reward prediction errors signaled by putative dopamine neurons in rat VTA depends on ventral striatum. *Neuron* 91, 182–193
- Wittmann, B.C. *et al.* (2005) Reward-related FMRI activation of dopaminergic midbrain is associated with enhanced hippocampus-dependent long-term memory formation. *Neuron* 45, 459–467
- Gruber, M.J. *et al.* (2016) Post-learning hippocampal dynamics promote preferential retention of rewarding events. *Neuron* 89, 1110–1120
- Patil, A. et al. (2016) Reward retroactively enhances memory consolidation for related items. Learn. Mem. 24, 65–69
- Miendlarzewska, E.A. *et al.* (2016) Influence of reward motivation on human declarative memory. *Neurosci. Biobehav. Rev.* 61, 156–176
- 57. Alberini, C.M. and Ledoux, J.E. (2013) Memory reconsolidation. *Curr. Biol.* 23, R746–R750

- Fernández, R.S. et al. (2016) The fate of memory: reconsolidation and the case of prediction error. Neurosci. Biobehav. Rev. 68, 423–441
- Nader, K. et al. (2000) Fear memories require protein synthesis in the amygdala for reconsolidation after retrieval. Nature 406, 722–726
- Sinclair, A.H. and Barense, M.D. (2018) Surprise and destabilize: prediction error influences episodic memory reconsolidation. *Learn. Mem.* 25, 369–381
- Exton-McGuinness, M.T.J. et al. (2015) Updating memories the role of prediction errors in memory reconsolidation. *Behav. Brain Res.* 278, 375–384
- Sevenster, D. et al. (2013) Prediction error governs pharmacologically induced amnesia for learned fear. Science 339, 830–833
- Sinclair, A.H. and Barense, M.D. (2019) Prediction error and memory reactivation: how incomplete reminders drive reconsolidation. *Trends Neurosci.* 42, 727–739
- Graves, A. *et al.* (2016) Hybrid computing using a neural network with dynamic external memory. *Nature* 538, 471–476
- Botvinick, M. *et al.* (2019) Reinforcement learning, fast and slow. *Trends Cogn. Sci.* 23, 408–422
- Eichenbaum, H. (2004) Hippocampus: cognitive processes and neural representations that underlie declarative memory. *Neuron* 44, 109–120
- Metcalfe, J. et al. (2012) Neural correlates of people's hypercorrection of their false beliefs. J. Cogn. Neurosci. 24, 1571–1583
- Bunzeck, N. et al. (2010) A common mechanism for adaptive scaling of reward and novelty. *Hum. Brain Mapp.* 1394, 1380–1394