

# Information utility in the human brain

Ifat Levy<sup>a,b,1</sup>

A funding agency just made a decision about your grant application. The official letter should arrive tomorrow, announcing the outcome, but you can log in to the agency's website and find out today. Would you? Humans and other animals are intrinsically motivated to acquire information, even when this information is of no instrumental value. This makes sense—after all, “knowledge is power” and, outside of the laboratory, information can guide action, to maximize gains and avoid harm. However, in some cases, people prefer to avoid information (1), particularly when this information is likely to validate a negative prediction. Recent theories propose that information carries its own utility—positive utility if the knowledge is likely to confirm a desirable belief (e.g., that your grant will be funded) and negative utility if it is more likely to support an undesirable belief (that your grant will not be funded). In PNAS, Charpentier et al. (2) use a novel experimental paradigm to provide behavioral support for these theories as well as evidence for neural encoding of the utility of information which is consistent with the theories.

In an fMRI experiment (Fig. 1), participants were presented with monetary lotteries. On some trials, the lottery offered a chance for winning \$1 (or winning nothing); on other trials, the lottery presented a chance for losing \$1 (or losing nothing). The probability for winning or losing (between 0.1 and 0.9) was conveyed in the form of a pie chart. Participants could not affect the outcome in any way, and this was made clear to them. They could, however, express their preference for finding out the outcome on each trial. Participants did this by choosing one of two options, offering different probabilities for revealing the outcome. After making the choice, participants were notified whether they would receive the information (knowledge cue) or not (ignorance cue). Knowledge cues were then followed by lottery outcomes (win, lose, or zero); ignorance cues were followed by a null symbol. Participants knew that at the end of the experiment they would receive the accumulated amount of lottery outcomes in all trials, regardless of whether

these outcomes were revealed. Following this task, participants were presented again with all of the lotteries and explicitly asked to indicate how much they would like to know the outcome.

Compatible with the notion of valence-dependent information utility, participants showed greater preference for information on gain trials, compared with loss trials, and also indicated a greater desire to know the outcome of gain trials. Moreover, preference for information increased for increasing gain likelihoods but decreased for increasing likelihoods for loss. That information about rewards had utility in and of itself was also confirmed in an additional behavioral experiment. In that experiment, participants made a series of investments in a simulated stock market and could bid for a chance to know—or not know—the value of their own portfolio. Participants were willing to pay for information, even though this payment had no effect on the outcome of their investments, and they were willing to pay more when they expected this information to be more positive. To economists, there is no better way to show that information indeed does have its own value.

How is the utility of information encoded in the brain? Theories of reinforcement learning typically posit that learning is driven by reward prediction errors (RPEs), or discrepancies between obtained and expected rewards (3). Charpentier et al. (2) hypothesize that the same brain regions that encode RPEs would also encode information prediction errors (IPEs)—the difference between the actual opportunity to gain knowledge and the expectation of such opportunity (Fig. 1). A crucial feature of their experimental design is that the probability for receiving information is completely independent from the outcome probability of the lottery (Fig. 1). This allowed the researchers to identify the neural encoding of IPE, separately from RPE. Note also that IPEs can be computed as soon as the knowledge (or ignorance) cue is presented. Importantly, at that point of the trial, participants know whether the lottery outcome will be revealed, but the outcome itself is still unknown, and

<sup>a</sup>Department of Comparative Medicine, Yale School of Medicine, New Haven, CT 06520; and <sup>b</sup>Department of Neuroscience, Yale School of Medicine, New Haven, CT 06520

Author contributions: I.L. wrote the paper.

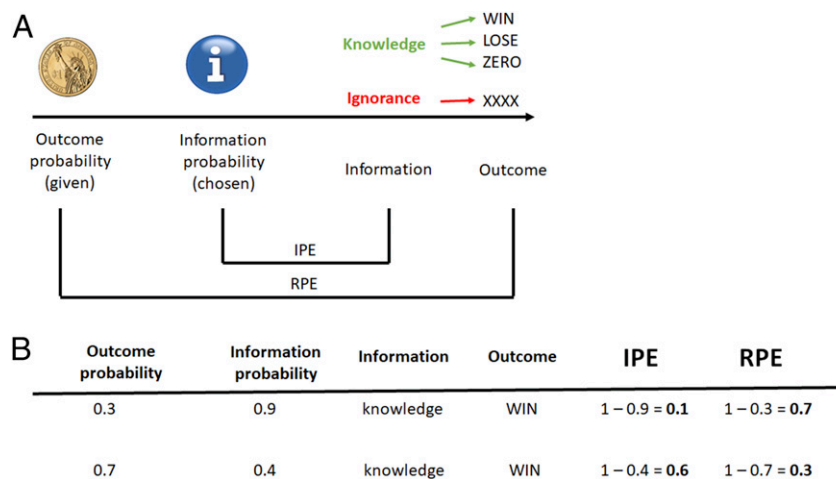
The author declares no conflict of interest.

Published under the PNAS license.

See companion article on page E7255.

<sup>1</sup>Email: ifat.levy@yale.edu.

Published online July 6, 2018.



**Fig. 1. Experimental design. (A) Trial structure. (B) Examples of two trials with different outcome and information probabilities.**

RPE cannot be computed; RPEs are in turn calculated once the outcome is revealed (Fig. 1). Thus, in this experiment, IPEs are generated at a different time than RPEs, and their magnitudes are not correlated with those of RPEs.

Based on the behavioral results, it is reasonable to expect valence-dependent encoding of IPE—that is, neural representation that reflects greater utility for information about gains, compared with losses; alternatively, it could also be that any opportunity to gain knowledge is coded by the brain as reward, irrespective of valence. Charpentier et al. (2) test both of these possibilities by calculating both IPE and valence-dependent IPE (IPE multiplied by the expected value of the lottery) for each trial and including both measures in the same model for analyzing the fMRI data. The analysis focused on brain areas that are most associated with reward: the ventral tegmental area and substantia nigra (VTA-SN)—midbrain regions rich in dopaminergic neurons—and the nucleus accumbens (NAc)—a major target of these dopamine neurons. Results show that valence-dependent IPE, but not valence-independent IPE, was encoded in VTA-SN. In addition, NAc tracking of valence-dependent IPE predicted individual preference for information—the more tightly was NAc activity associated with this prediction error the more sensitive was the individual to the expected value of the lottery. Thus, reward-related structures seem to treat information about rewards in a manner very similar to how they treat the rewards themselves. While fMRI cannot tell us whether the neurons encoding IPEs are the same neurons that encode RPEs, findings in monkeys suggest that this is the case (4). As if this was not complex enough (how does the brain differentiate between IPEs and RPEs?), a recent study in humans (5) implicated the NAc in encoding yet another type of prediction error. Participants in that study were required to acquire and update declarative information (regarding the Falklands War). Activity in NAc—as well as in several other areas—encoded the degree to which new factual information violated expectations based on prior knowledge and beliefs. In other words, IPEs that have nothing to do with gains and losses may be encoded in a manner similar to IPEs and RPEs. Whether the same neurons participate in all of these representations and, more broadly, how the brain distinguishes between the different types of information is an open question.

Why does the brain encode the discrepancy between the actual and expected opportunity to gain knowledge, rather than simply encoding the value of the opportunity itself? In the reported study, IPEs (as well as RPEs) were of no use to the participants.

There was nothing to learn, or update, based on these errors; information provided in one trial was no longer useful in the next trial. This was a sensible choice of design, because it allowed a very simple and precise calculation of prediction errors, irrespective of individual differences in the rate of learning from feedback. An intriguing question is, however, how IPEs are used by the brain in a more naturalistic setting, where learning is not only possible but is also desirable for adaptive behavior. Like RPEs, which are used to update the value of cues that predict rewards, IPEs can be used to update the value of potential sources of information. Evidence from monkeys suggests that neurons in posterior parietal cortex encode the expected value of information, independently from the expected value of the action chosen based on that information (6). IPEs may be used to update such representations—whether and how this is done is a matter for future research.

If we stopped here, it would seem that information about probable losses invariably carries negative utility. However, although participants expressed lower desire for information about losses, compared with gains, they still chose the more informative option (which offered higher probability for information) on the majority of loss trials. This suggests that preference for information is driven both by the desirability of the expected outcome (which motivates information seeking for probable gains and information avoidance for probable losses) and by general curiosity, or preference for reducing uncertainty (which motivates information seeking for highly uncertain gains and losses) (7). Indeed, a formal statistical test showed significant effects of both the level of uncertainty and the expected value of the lottery. Because participants requested more information on gain trials, uncertainty about the outcomes of these trials was reduced, compared with uncertainty about outcomes of loss trials. Interestingly, this asymmetry between gains and losses echoes the widely observed differences in uncertainty attitudes between the gain and loss domains. Most individuals are risk-averse when choosing between potential gains but risk-seeking when choosing between potential losses (8, 9). Similarly, ambiguity (missing information about outcome probabilities) affects choices in the realm of gains much more than choices about losses (10). The behavioral results reported by Charpentier et al. (2) suggest a link between information seeking and individual uncertainty attitudes. The greater tolerance for uncertainty in the loss domain may have developed in response to the greater experienced uncertainty around losses. Alternatively, greater tolerance to uncertainty in the loss domain may be one of the underlying causes

- Levy