

Published in final edited form as:

*Neural Netw.* 2006 June ; 19(5): 535–546. doi:10.1016/j.neunet.2005.11.002.

## Encoding uncertainty in the hippocampus

L.M Harrison\*, A Duggins†, and Friston K.J\*

\*Wellcome Department of Imaging Neuroscience, Institute of Neurology, 12 Queen Square, London WC1N 3BG, UK

†Department of Neurology, Westmead Hospital, Sydney, Australia

### Abstract

The medial temporal lobe may play a critical role in binding successive events into memory while encoding contextual information in implicit and explicit memory tasks. Information theory provides a quantitative basis to model contextual information engendered by conditional dependence between, or conditional uncertainty about, consecutive events in a sequence. We show that information theoretic indices characterizing contextual dependence within a sequential reaction time task (SRTT) predict regional responses, measured by fMRI, in areas associated with sequence learning and navigation. Specifically, activity of a distributed paralimbic system, centered on the left hippocampus, correlated selectively with predictability as measured with mutual information. This is clear evidence that the brain is sensitive to the probabilistic context in which events are encountered. This is potentially important for theories about how the brain represents uncertainty and makes perceptual inferences, particularly those based on predictive coding and hierarchical Bayes.

### Keywords

Sequential reaction time task (SRTT); 1<sup>st</sup> order Markov sequence; implicit learning; medial temporal lobe; information theory; predictability; Bayesian observer

### Introduction

Causal structure within the physical world induces regularities in the timing and order of events. Such regularities enable an organism to predict outcomes given current information and thereby learn from, and adapt to, the changing world within which it has to survive. A memory system that supports this form of learning is therefore useful. It has been suggested that the medial temporal lobe (MTL) plays a crucial role in generating flexible representations of novel contextual relationships among distinct stimulus features (Chun and Phelps, 1999; Poldrack and Rodriguez, 2003; Rose et al., 2002; Schendan et al., 2003). According to this relational account of memory, the MTL is engaged in associative processes that bind multiple aspects of stimulus events into a memory (Cohen and Eichenbaum, 1993; Wallenstein et al., 1998) whether the content of what has been learned is available to awareness or not.

The aim of this study was to establish a quantitative relationship between neurophysiological responses evoked in the hippocampus, during the presentation of stimulus sequences, and the predictability of those sequences as measured using information

theory. Samples were taken from a discrete conditional probability distribution to generate a 1<sup>st</sup> order Markov sequence (Cox and Miller, 1965) of varying predictability. Information theory measures of conditional uncertainty were then used to model behavioral and functional imaging data acquired during a sequential reaction time task (SRTT) using these sequences. This task is typically used in cognitive psychology to dissociate learning from awareness (Willingham, 1980). We hoped to show that conditional uncertainty and predictability are encoded within the MTL and connected structures.

There are two levels at which hippocampal and related paralimbic structures could be involved in representing the probabilistic structure of sequences. There is considerable evidence that the hippocampus is sensitive to novel events that are, by definition, unpredictable. Here the predictability pertains to the probability of a particular stimulus or event. However, there is a probabilistic context in which events occur that could also be usefully encoded by systems like the hippocampus. This level of representation is the predictability of, or uncertainty about, events before they occur. This uncertainty is not stimulus bound, but reflects the temporal regularity of successive events in a given experimental or environmental context.

The probabilistic context is potentially important from the point of view of perception and representational learning. Theoretical accounts of perceptual inference, based on generative models and predictive coding, emphasize the conjoint influence of bottom-up evidence from sensory inputs and top-down effects that mediate prior expectations. To attain the optimum balance, the relative uncertainty associated with the bottom-up and top-down information must be known, or estimated. This uncertainty clearly changes with the predictability associated with the sensory context. We hypothesized that the neurophysiologic correlate of predictability would be observed with functional neuroimaging, possibly in the hippocampus that has a special role in sequential processing. There is current interest in the neuronal mechanisms that might encode predictability or uncertainty that make these physiological correlates particularly interesting (see for example (Yu and Dayan, 2002)). The theoretical analysis presented in Yu and Dayan is relevant because it implicates cholinergic neurotransmission, that has a key role in regulating hippocampal dynamics (Hasselmo, 1999). Our focus was on encoding uncertainty in stimulus-stimulus relationships. However, it is interesting to note that in stimulus-response learning that the dopamine system, which targets dorsal and ventral striatal, orbital and frontal regions, may encode the discrepancy between predicted and actual reward (prediction error) and uncertainty (Aron et al., 2004; Fiorillo et al., 2003).

In this work we were interested in region-specific responses to changes in probabilistic context, as reflected by the conditional uncertainty about sequential events. A sequence that has a simple structure is one where the current event ( $E_t$ ) is conditionally dependent on the previous ( $E_{t-1}$ ). The probability of transition between consecutive events is given by a conditional probability  $p(E_t | E_{t-1})$ , also known as a transition matrix (TM). Serial events that conform to this model are 1<sup>st</sup> order Markov sequences. By presenting different 1<sup>st</sup> order Markov sequences to subjects we were able to vary the uncertainty and quantify it using information theory. We calculated four indices for each sequence: the surprise of each stimulus ( $\hat{h}$ ), the entropy of each sequence ( $\hat{H}$ ), the reduction in surprise afforded by the previous stimulus ( $\hat{i}$ ) and mutual information between consecutive stimuli within a sequence ( $\hat{J}$ ). Surprise and its reduction are stimulus-specific, whereas entropy and mutual information are measures of uncertainty that pertain to the context established by each sequence. Critically, the mutual information provides a natural measure of conditional uncertainty (conditioned on previous stimuli).

We used the SRTT to engage the hippocampal system in a relational task and to model the effect of conditional uncertainty on both behavioral and fMRI responses. Clearly, from the point of view of the subject, the conditional uncertainty had to be learned for each new sequence. As each sequence progressed, conditional uncertainty about the next stimulus falls as the probabilistic structure is disclosed. We modeled this assuming that the subject was an ideal Bayesian observer, who started with flat priors at the beginning of each sequence. In addition to this within-sequence, learning-related change in uncertainty we introduced between-sequence differences by using different probability transition matrices. This increased the statistical efficiency of our experimental design

In brief, we demonstrated a dependence of reaction times on the information theoretic measures above and, critically, showed that BOLD activity increased with mutual information in the left hippocampus, bilateral parieto-occipital sulcus, left retrosplenial cortex and right anterior cingulate. Measuring the correlates of conditional uncertainty in this way represents a quantitative approach to the brain's response to hidden structure within sequences and the encoding of uncertainty<sup>1</sup>. To assess the frequency with which subjects become explicitly aware of the contingencies, we performed an auxiliary behavioral study (without scanning), involving twelve different subjects using the identical paradigm.

## Methods

### Experimental design

The design comprised 12 blocks, each containing a sequence of 40 trials. A trial involved presenting one of four possible colored shapes (displayed at the bottom of the screen; stimulus duration 500ms; stimulus onset asynchrony: 2.2s). Subjects were required to respond by identifying the target and their reaction times were recorded. In all trials two colors and shapes were combined to form four possible events. An example of a trial is shown at the top of Figure 1a. At the beginning of a block subjects were cued for 5 seconds with the four objects in a row at the bottom of the screen, which remained there throughout the block. Following the initial 5 seconds a series of 40 trials were sampled from a transition matrix,  $p(E_t | E_{t-1})$  and presented to subjects as a SRTT. See Figure 1b for an example of a transition matrix, where the top right figure shows its gray-scale representation.

Dependence between consecutive trials is encoded in the transition matrix, which remained constant within a block and varied over blocks. Subjects were asked to respond to each trial by pressing a key to indicate the position of the target in the display at the bottom of the screen as rapidly as possible, but not at the expense of accuracy. A schematic of a block is shown in Figure 1a. No indication as to an underlying pattern within the sequence was given. Thirteen subjects were scanned whilst performing the task and debriefed afterwards to assess their awareness of patterns within the sequence.

### Sample-based estimation of uncertainty

Given the sampling nature of experience an observer can only infer probabilistic structure from events. We used the notion of an 'ideal' observer to estimate the conditional probability,  $p(E_t | E_{t-1})$  using a simple Bayesian update scheme. We assumed that, at the beginning of each block, the observer started with a prior that all current events are equally likely and consecutive events are independent. This is not a trivial assumption and is the topic of current research into intelligent priors given small data sets (see discussion).

<sup>1</sup>Schendan *et al* used second order sequences, *i.e.* where contingencies exist among more than 2 successive trials (see discussion).

The marginal distribution,  $p(E_t)$ , was estimated from the number of occurrences of event  $i$  up to sample  $t$  (written as  $n_i^t$ , where  $i$  indexes the current event type and  $t$  the trial number). The estimate at sample  $t$  ( $t > 0$  and  $t = 0$  respectively) is given by Eq.1, from which entropy is calculated.

$$f_t(E_t=i) = \frac{(n_i^t+1)}{\sum_i (n_i^t+1)}, \quad \left( f_0(E_0=i) = \frac{1}{4} \right) \quad 1)$$

Figure 2 shows the estimate (dashed line) of one sequence, where entropy is initially maximal and decreases towards the true value (solid line) with increased sampling. Similarly an estimate of the joint probability distribution can be estimated from a count of event pairs up to sample  $t$  (written as  $n_{ij}^t$ , where  $i$  and  $j$  index the current and previous event type) and is given by

$$f_t(E_t=i, E_{t-1}=j) = \frac{(n_{ij}^t+1)}{\sum_{i,j} (n_{ij}^t+1)}, \quad \left( f_0(E_0=i, E_{-1}=j) = \frac{1}{16} \right) \quad 2)$$

Initially, the observer is ignorant of all contingencies ( $f_0$  in parentheses). An estimate of the transition matrix,  $f_t(E_t|E_{t-1})$ , is easily calculated from  $f_t(E_t, E_{t-1})$ , which approaches the true transition matrix with more samples (shown for one block in Figure 3). The estimate of mutual information is therefore zero. As sampling begins, a tally of consecutive event pairs ( $E_t, E_{t-1}$ ) is used to update  $f_t(E_t|E_{t-1})$  with each sample. The final estimate of the true transition matrix is shown in the middle figure after 40 trials, *i.e.*  $f_{41}(E_t|E_{t-1})$ . This can be compared to the true distribution,  $p(E_t|E_{t-1})$ . A plot of mutual information is shown in the lower figure, where the true value (solid line) is constant over the block, while the estimate (dashed line) rises towards the true value.

The relationship between entropy, conditional entropy and mutual information is illustrated in Figure 4 with a Venn diagram (Cover, 1991) and Figure 5 with time series of estimates over 1 block. The reduction in uncertainty in the current trial, afforded by the previous, is apparent on comparing  $H(E_t)$  and  $H(E_t|E_{t-1}) = H(E_t) - I(E_t; E_{t-1})$  in the top graph, Figure 5. The difference is the mutual information and is shown at the bottom of the figure.

Once entropy and mutual information are estimated the trial-by-trial surprise and reduction in surprise can be calculated. The expressions used to estimate the surprise ( $\hat{h}$ ), entropy ( $\hat{H}$ ), surprise reduction ( $\hat{i}$ ) and mutual information ( $\hat{I}$ ), are

$$\begin{aligned} \hat{h}(E_t=i) &= -\log(f_t(E_t=i)) \\ \hat{H}(E_t) &= \langle -\log(f_t(E_t)) \rangle \\ \hat{i}(E_t=i; E_{t-1}=j) &= \log\left(\frac{f_t(E_t=i|E_{t-1}=j)}{f_t(E_t=i)}\right) \\ \hat{I}(E_t; E_{t-1}) &= \left\langle \log\left(\frac{f_t(E_t|E_{t-1})}{f_t(E_t)}\right) \right\rangle_{f_t(E_t, E_{t-1})} \end{aligned} \quad 3)$$

Examples of all four quantities over the duration of an experiment, for one subject, are shown in Figure 6. Initially all events are equally uninformative, but as contingencies are learned, some events become more predictive than others. The same is true of surprise, where initially all events are equally surprising, but eventually some events are more surprising than others. This explains the increased variability in both trial-by-trial surprise and reduction in surprise with increasing sample number.

We were interested in regional responses that correlated with conditional entropy  $H(E_t | E_{t-1})$  *i.e.* conditional uncertainty. However, entropy and conditional entropy are themselves highly correlated, which could confound our interpretation. However, the expression  $H(E_t | E_{t-1}) = H(E_t) - I(E_t; E_{t-1})$ , allows us to decompose conditional entropy into an instantaneous component (entropy), and one that models the temporal relationship (mutual information) among consecutive trials. This decomposition disambiguates the interpretation of cortical correlates. Heuristically, this says that conditional uncertainty about the next stimulus, given the current stimulus, can be partitioned into two components. The first is simply the uncertainty about the next stimulus irrespective of the preceding stimulus. The second component is the reduction in this uncertainty afforded by its precedent. This second component is specifically related to the probabilistic structure of the temporal contingencies and was the focus of our analysis.

### Assessing explicit learning

In the behavioral experiment, after each block subjects were given a free generation task used by Honda *et al* (Honda et al., 1998) to assess awareness of a deterministic sequence. Subjects were asked: ‘Did you notice anything about the task?’ If they answered yes, they were asked ‘What did you notice?’ and if they answered ‘a sequence’ or ‘pattern’ they were asked to ‘report the sequence, as far as you noticed, verbally’.

Subjects were then given a cued generation task to assess their ability to generate the contingencies they had encountered during a block. The test involved presenting subjects with a test sequence of four trials, generated from the same transition matrix, after which they were asked which object was most likely to occur next. The last in the test sequence varied through the four possible targets (*i.e.* target numbers 1 to 4) to test contingencies associated with each target. Twelve blocks in total generated 48 responses per subject. Given that we knew the conditional and marginal distributions used in each block, we were able to ask whether there was any evidence, within subjects’ responses, in favor of explicit learning of the conditional probabilities. We assessed this using the likelihood ratio of their responses based on the conditional distribution  $p(E_t | E_{t-1})$ , relative to the marginal distribution  $p(E_t)$ . This odds-ratio provided a principled test of whether the subjects’ responses were informed explicitly by the contingencies to which they had been exposed.

### Subjects

Informed consent was obtained from 13 right-handed subjects (8 males; age range 22-35 years; mean age 27). Ethics approval was obtained from the joint ethics committee of the Institute of Neurology, University College London and National Hospital of Neurology and Neurosurgery, London. A behavioral study on 12 different subjects (7 male; age range 23-34; mean age 26) was performed to assess awareness of patterns within blocks of trials.

### Imaging

A 2T Siemens VISION system was used to acquire T1-weighted anatomical images and gradient-echo echo-planar T2\*-weighted MRI image volumes with blood oxygenation level dependent (BOLD) contrast. A total of 552 volumes were acquired per subject plus 6 initial ‘dummy’ volumes to allow for T1 equilibration effects. Volumes were acquired continuously every 2506ms. Each volume comprised 33 3.3mm axial slices, with an in-plane resolution of 3×3mm, positioned to cover the entire cerebrum. The imaging time series were realigned, slice-time corrected, normalized into the standard anatomical space defined by Montreal Neurological Institute (MNI) and smoothed with a Gaussian kernel of 6mm full width half maximum.

The data were analyzed using the software Statistical Parametric Mapping (SPM2, <http://www.fil.ion.ucl.ac.uk/spm>), employing an event-related model (Josephs et al., 1997) and a two-stage random effects procedure. A model of the BOLD response to trials (explanatory variables) was constructed by convolving a series of modulated delta-functions (one for each trial) with a canonical hemodynamic response function (HRF). The delta functions were modulated by the estimates of trial-by-trial surprise ( $\hat{h}$ ), entropy ( $\hat{H}$ ) trial-by-trial reduction in surprise ( $\hat{i}$ ) and mutual information ( $\hat{I}$ ). This meant that the effect of each information theoretic measure was modeled at the neuronal level and whose consequences on the BOLD response could be predicted. This predicted BOLD response for each information theoretic index was used in a general linear model to investigate the measured BOLD response. Maximum correlation among the regressors was -0.1728 (between mutual information and entropy). Nuisance variables included an exponentially decaying covariate (half life of 3 blocks) to model non-specific adaptation, response errors, low frequency drifts in signal (cut off 64 seconds) and movement parameters, calculated during realignment. The 5-second cue periods before each block were modeled using delta functions at the beginning of each block. Parameter estimates for each subject and regressors were calculated for each voxel (Friston, 1995). For the second stage of the random-effects analysis, subject-specific parameters for each of the four information theoretic measures were entered into four one-sample t-tests.

## Results

### Behavioral

**Reaction Times**—All incorrect responses were removed and the average reaction times, over all blocks and subjects calculated. These are shown in Figure 7 and demonstrate a large initial reduction followed by a gradual decrease in reaction times with trial number within a block. This indicates implicit or explicit learning. In addition an AnCova of the reaction times was performed using the four information theoretic indices as explanatory variables. These results are shown in Figure 8, which demonstrates a significant reduction in reaction times (ms/bit) for both mutual information (66 ms/bit;  $p < 0.001$ ) and reduction in surprise (21 ms/bit;  $p < 0.001$ ). There was a significant increase with entropy (169 ms/bit;  $p < 0.001$ ) and surprise (63 ms/bit;  $p < 0.001$ ). In short, subjects took less time to respond when the sequence was predictable and longer to less frequent (more surprising) events, irrespective of the sequences predictability.

It is important to appreciate that the reductions in reaction time predicted by changes in any one of the four measures cannot be explained by changes in the others. This is the nature of inference with the general linear model (in this instance, analysis of covariance); in which one regressor ‘explains away’ any effect that could be explained by another. This means that not only do surprising events incur a longer reaction time but also, in the context of sequences that are inherently unpredictable, there is a further increase above and beyond the average of the trial-bound increases. In other words, the reaction time appears to be sensitive, both to the probabilistic attributes of specific events and to the probabilistic context in which these events occur. This was in contra-distinction to the neurophysiologic responses (see functional imaging results) that seemed to be much more sensitive to the probabilistic context, as indexed by mutual information.

**Free Generation Task**—In the behavioral study subjects reported positively in 46% of blocks to the first question, ‘did you notice anything about the stimulus?’. These subjects were then questioned further. In 24% of blocks subjects were unable to describe what they noticed. In 22% of blocks they were able to give examples of what they had noticed, however, most of these were incorrect. Subjects were able to correctly identify simple

repetitions of a single object (3-10 trials) or alternations between two objects lasting for 2-4 cycles in 10% of blocks.

**Cued Generation Task**—We were interested in assessing whether responses were based on an explicit knowledge of the conditional distribution or just the marginal distribution. The logarithms of the odds-ratio for each subject are shown in Figure 9. For all subjects (except one) the evidence is in favor of the marginal distribution and significantly so for ten out of twelve of the subjects. This means that, despite faster reaction times, they were not able to use what they had learnt explicitly. We conclude from this that subjects did not have explicit access to the conditional probabilities acquired implicitly. The distinction between implicit and explicit learning is not central to our basic hypothesis that the brain represents sequential predictability. Indeed, our statistical model was based upon an ideal Bayesian observer that is indifferent to the cognitive mechanisms that mediate learning. However, from a cognitive neuroscience perspective, our results can be interpreted within the domain of implicit learning.

In what follows, we only discuss results that survived a correction for multiple comparisons, using the corrected p-value based on spatial extent (see Table 1). A height threshold of  $p < 0.001$  uncorrected defined the spatial extent.

### Functional imaging results

Activity in left hippocampus, bilateral parieto-occipital sulcus, left retrosplenial cortex and right anterior cingulate was positively correlated with mutual information. No significant effects were seen for the remaining information theoretic indices at this threshold. See Table 1 for coordinates and Z-scores of significant regions. Figure 10 shows orthogonal sections of a Statistical Parametric Map (SPM) centered on the local maxima (voxel coordinate [-30,-18,-24]) of the left hippocampus. This demonstrates the response of the left hippocampus to mutual information. The bottom right panel of this figure shows parameter estimates at the same local maxima to all information theoretic indices (taken from a one-sample t-test at the second level) and demonstrates the selective response to mutual information. Parameter estimates measured in parieto-occipital sulcus, retrosplenial cortex and anterior cingulate are shown similarly in Figure 11 (voxel coordinates of local maxima given in Table 1). All regional responses correlate selectively with mutual information, except in anterior cingulate, which also showed significant negative correlation with entropy.

### Discussion

This study was designed to engage the hippocampal system using a 1<sup>st</sup> order Markov sequence in a sequential reaction time task (SRTT). Recent reports (Rose et al., 2002; Schendan et al., 2003) have presented evidence for a relational memory account of learning. This calls on the hippocampal system to represent temporally distinct and novel relationships, regardless of whether the task is learned implicitly or explicitly. The MTL encodes context and its activity may reflect relationships among events. In particular the hippocampus may mediate expectancies and inferences (Eichenbaum et al., 1999) based on the probabilistic structure of past events, particularly the conditional uncertainty about what will happen next. We chose a 1<sup>st</sup> order Markov sequence as it contains contingencies (between consecutive trials) and has a precise mathematical structure from which information theoretic indices are easily calculated. The notion of an ‘ideal’ observer was introduced to estimate these measures, which were updated with each new sample. These quantities were then used to predict behavior and brain responses while subjects performed the task.

Reaction times decreased with both mutual information measures (reduction in surprise and its expectation). The stronger the dependence between consecutive events (*i.e.* high predictability) the lower the reaction time. As noted above, this is an extremely interesting result that suggests a behavioral facilitation, not only for frequently encountered events, but also conferred by the probabilistic context in which events occur. This might be explained by an increased reliance on prior expectations that speeds up reaction times irrespective of a particular event's probability.

It is apparent from these reaction time results that subjects were able to represent the contingencies between consecutive trials. We had hypothesized that the activity of those brain regions involved in this representation would vary with an information theoretic measure of temporal association. Indeed our neuroimaging results demonstrate responses, within an interconnected network involving the left hippocampus, bilateral parieto-occipital sulcus, left retrosplenial cortex and right anterior cingulate, that correlates with the mutual information between consecutive trials. Responses in these regions increased when events became more predictable (reduction in uncertainty or relative increase in certainty), *i.e.* measured as an increase in mutual information, irrespective of how surprising the actual event was. This result supports the notion that specific brain regions, critically including the hippocampal system and its connected structures, may be sensitive to uncertainty within Markov sequences.

Subjects were unable to use explicit knowledge to reproduce contingencies. However, reaction times were sensitive to dependencies within a block. This indicates that implicit learning had occurred. However, establishing that learning is truly implicit is difficult (Shanks, 1994). Within the implicit SRTT learning literature, evidence implicating the hippocampal system, striatum and cortical components of fronto-striatal pathways has been reported (Berns et al., 1997; Rose et al., 2002; Schendan et al., 2003). In particular, Schendan *et al* measured MTL responses during implicit and explicit learning of second order sequences (*i.e.* contingencies exist among more than 2 events). This was motivated by a study (Curran, 1997) that demonstrated impaired implicit learning of higher order associations compared to first order (or pairwise association) in patients with anterograde amnesia.

Results from navigation research are relevant as navigation involves processing sequential information, with many reports in the literature of activity within the network connecting parieto-occipital sulcus through retrosplenial cortex to MTL (Burgess et al., 2001). Maguire (Maguire, 2001) reports functional imaging and patient studies implicating the retrosplenial cortex in navigation. Evidence from patient studies suggests that only the right hippocampus is necessary for navigation, while the left may have a more general function in episodic memory. The parieto-occipital sulcus and retrosplenial cortex provide input to the hippocampal system, consistent with its involvement in encoding predictability. Lateralisation of hippocampal function has been reported in context-dependent episodic memory involving the left hippocampus, whilst the right is associated with spatial navigation (Burgess et al., 2001). Given the contextual nature of Markov sequences, it is interesting to note that we found left hippocampal activity was correlated with mutual information. Several studies have reported activations in the AC *e.g.* (Berns et al., 1997) too. This is reasonable as a loop of reciprocal connections exists between the frontal lobes and basal ganglia (Seger, 1994) and patients with striatal damage are typically impaired on implicit SRTT tasks.

No significant effects (corrected for multiple comparisons) were detected for the remaining indices. Strange *et al* (Strange, 2005) reported left hippocampal activity correlated with an estimate of the entropy of an independently sampled sequence in a SRTT. In addition, they



detected an extensive bilateral cortico-thalamic network correlated with stimulus-bound surprise. There are several factors that would explain the superficial discrepancy between the results of Strange *et al* and those presented above. First, the fact we failed to demonstrate a significant effect of entropy does not mean that this effect was absent (we used a very conservative correction procedure for our imaging analysis). A second and more compelling reason relates to the motivation for the current study. If the hippocampus is specifically interested in the relational or temporal structure of sequences, it might respond selectively to the conditional entropy of the current stimulus given the preceding one. As indicated in Figure 5 the conditional entropy has two components; the entropy *per se* and the mutual information. In the Strange *et al* study, the sequence was random and the conditional entropy was exactly the same as the entropy. In our study we deliberately introduced variations in mutual information that represented the pre-dominant changes in conditional entropy. In short a parsimonious explanation for the positive results of the two studies is that the hippocampus shows selective responses to changes in conditional entropy. This speaks to a specific role in temporal sequencing and the encoding of conditional uncertainty.

Our model of implicit learning was based upon a Bayesian update scheme and touches on an active area of current research, the learning of probability distributions. Recent work has shown that there are intelligent Bayesian priors that dramatically reduce the bias in the calculation of entropy from small data-sets (Nemenman and Bialek, 2002;Paninski, 2003). It is possible that one could use fMRI responses to disambiguate among different models of density learning that would be expressed primarily in the dynamics or temporal evolution of reaction times and event-related responses. This report limits itself to a simple model based upon the assumption of an ideal observer. Clearly, this may not be the best model but was sufficient to disclose predictability-related responses in the hippocampus. Our model is sufficient in the sense that had it not predicted the observed physiological responses sufficiently accurately, we would not have obtained significant results. However, it should be noted that other observer models might also have been significant.

There is a growing interest in the role of predictability in reward processing. Indeed, basic models of reinforcement learning are predicated on temporal difference models that encode the predicted future reward (Sutton and Barto, 1981). These models have been refined and examined from a cognitive neuroscience perspective, placing predictability in a central position, not only for reinforcement and emotional learning but also for perception itself. Indeed, our own work on hierarchical processing, of a Bayesian sort, provides an algorithmic perspective on the central role of parameters encoding uncertainty about the causes of sensory input and uncertainty about the input itself (Friston, 2003). The main neurobiological insight that obtains from this study is that the hippocampus may not only be involved in sequence learning (Schendan et al., 2003) but may also be involved in the representation of how learnable sequences are. In machine learning, this learnability or predictability is critical for estimation and inference: It balances the relative weight assigned to prior expectations and the likelihood of obtaining subsequent data or sensory input. It is therefore possible that the hippocampus and related structures encode uncertainty to finesse representational learning and perceptual inference.

In summary, this study provides a quantitative functional anatomic basis for learning contextual relationships engendered by conditional dependence among consecutive events. The notion of an 'ideal' observer was used to calculate the mutual information as a measure of conditional uncertainty. Regions whose activity correlated with this index were the hippocampal system, parieto-occipital sulcus, retrosplenial cortex and anterior cingulate. These regions have been implicated in many sequence-learning and navigation studies, suggesting that they may be involved in encoding the expected uncertainty of temporal events as they unfold.

## Acknowledgments

We thank the Wellcome Trust for funding this work. This work was presented in provisional form at the Human Brain Mapping conference 2003, New York, United States.

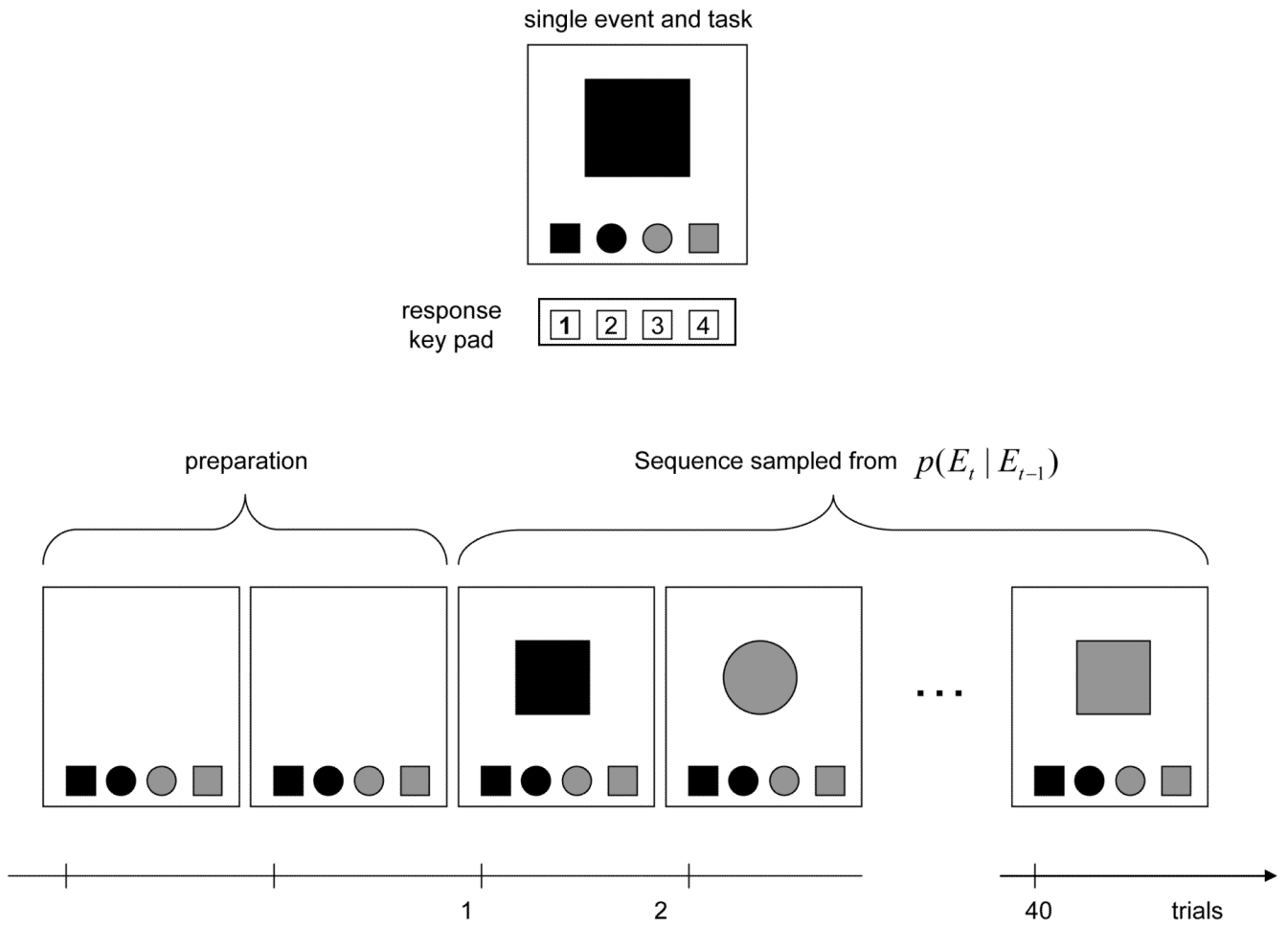
## Reference List

1. Aron AR, Shohamy D, Clark J, Myers C, Gluck MA, Poldrack RA. Human midbrain sensitivity to cognitive feedback and uncertainty during classification learning. *J.Neurophysiol.* 2004; 92:1144–1152. [PubMed: 15014103]
2. Berns GS, Cohen JD, Mintun MA. Brain regions responsive to novelty in the absence of awareness. *Science.* 1997; 276:1272–1275. [PubMed: 9157889]
3. Burgess N, Maguire EA, Spiers HJ, O'Keefe J. A temporoparietal and prefrontal network for retrieving the spatial context of lifelike events. *Neuroimage.* 2001; 14:439–453. [PubMed: 11467917]
4. Chun MM, Phelps EA. Memory deficits for implicit contextual information in amnesic subjects with hippocampal damage. *Nat.Neurosci.* 1999; 2:844–847. [PubMed: 10461225]
5. Cohen, N.; Eichenbaum, H. *Memory, Amnesia and the Hippocampal System.* MIT Press; Cambridge, MA: 1993.
6. Cover, TM. *Elements of Information Theory.* Wiley-Interscience; 1991.
7. Cox, DR.; Miller, HD. *The Theory of Stochastic Processes.* Methuen; London: 1965.
8. Curran T. Higher-order associative learning in amnesia: evidence from the serial reaction time task. *The Journal of Cognitive Neuroscience.* 1997; 9:522–533.
9. Eichenbaum H, Dudchenko P, Wood E, Shapiro M, Tanila H. The hippocampus, memory, and place cells: is it spatial memory or a memory space? *Neuron.* 1999; 23:209–226. [PubMed: 10399928]
10. Fiorillo CD, Tobler PN, Schultz W. Discrete coding of reward probability and uncertainty by dopamine neurons. *Science.* 2003; 299:1898–1902. [PubMed: 12649484]
11. Friston K. Learning and inference in the brain. *Neural Netw.* 2003; 16:1325–1352. [PubMed: 14622888]
12. Friston KJ. Statistical parametric maps in functional imaging: a general linear approach. *Hum.Brain Mapp.* 1995; 2:189–210.
13. Hasselmo ME. Neuromodulation and the hippocampus: memory function and dysfunction in a network simulation. *Prog.Brain Res.* 1999; 121:3–18. [PubMed: 10551017]
14. Honda M, Deiber MP, Ibanez V, Pascual-Leone A, Zhuang P, Hallett M. Dynamic cortical involvement in implicit and explicit motor sequence learning. A PET study. *Brain.* 1998; 121(Pt 11):2159–2173. [PubMed: 9827775]
15. Josephs O, Turner R, Friston KJ. Event-related fMRI. *Hum.Brain Mapp.* 1997; 5:243–248. [PubMed: 20408223]
16. Maguire EA. The Retrosplenial contribution to human navigation: A review of lesion and neuroimaging findings. *Scand.J.Psychol.* 2001; 42:225–238. [PubMed: 11501737]
17. Nemenman I, Bialek W. Occam factors and model independent Bayesian learning of continuous distributions. *Phys.Rev.E.Stat.Nonlin.Soft.Matter Phys.* 2002; 65:026137. [PubMed: 11863617]
18. Paninski. Estimation of Entropy and Mutual Information. *Neural Comp.* 2003; 15:1191–1253.
19. Poldrack RA, Rodriguez P. Sequence learning: what's the hippocampus to do? *Neuron.* 2003; 37:891–893. [PubMed: 12670418]
20. Rose M, Haider H, Weiller C, Buchel C. The role of medial temporal lobe structures in implicit learning: an event-related FMRI study. *Neuron.* 2002; 36:1221–1231. [PubMed: 12495634]
21. Schendan HE, Searl MM, Melrose RJ, Stern CE. An FMRI study of the role of the medial temporal lobe in implicit and explicit sequence learning. *Neuron.* 2003; 37:1013–1025. [PubMed: 12670429]
22. Seger CA. Implicit learning. *Psychol.Bull.* 1994; 115:163–196. [PubMed: 8165269]
23. Shanks DR, St.J MF. Characteristics of dissociable human learning systems. *Behavioral and Brain Sciences.* 1994; 17:367–447.

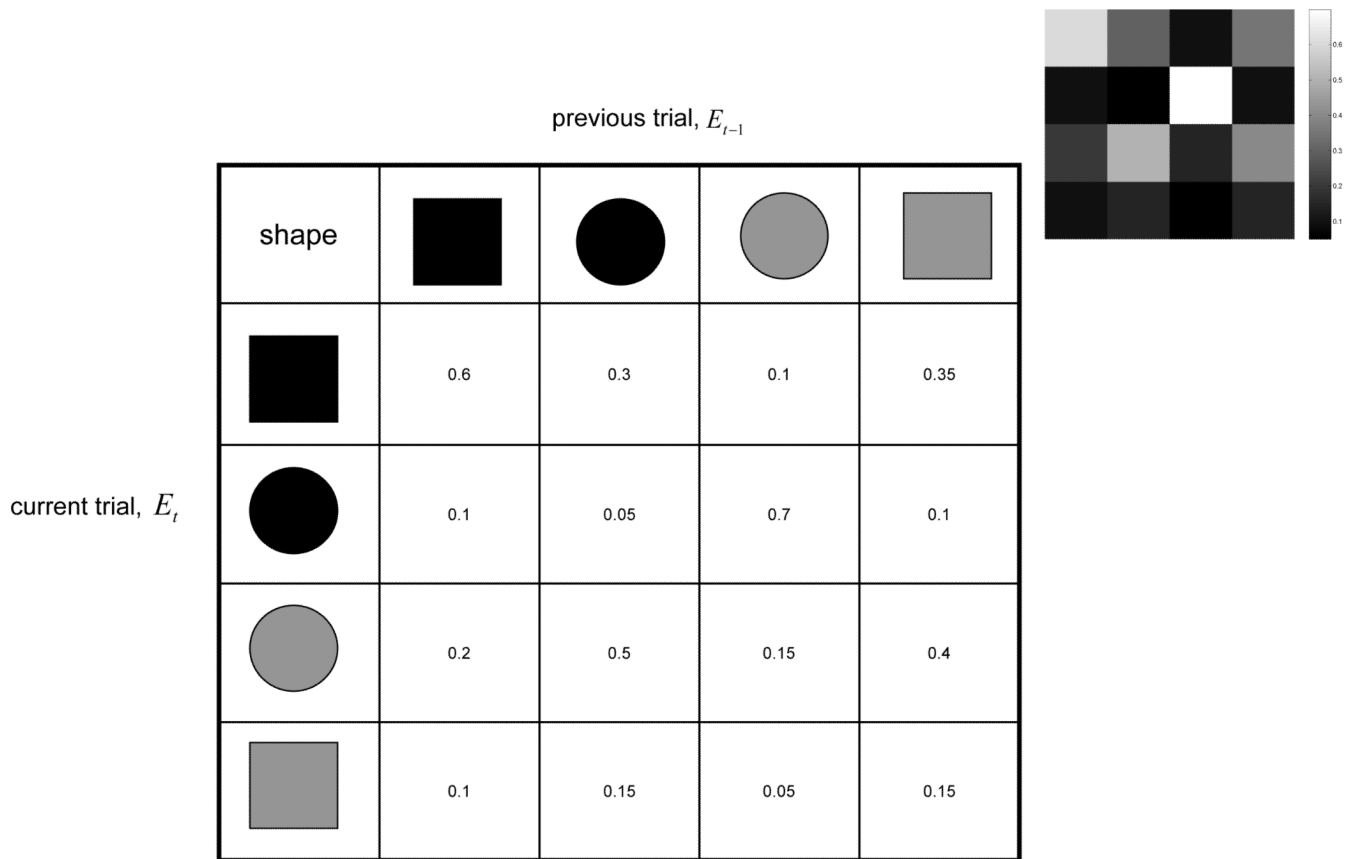
24. Strange B. Information theory, novelty and hippocampal responses: unpredicted or unpredictable? *Neural Networks*. 2005; 18:225–230. [PubMed: 15896570]
25. Sutton RS, Barto AG. Toward a modern theory of adaptive networks: expectation and prediction. *Psychol.Rev.* 1981; 88:135–170. [PubMed: 7291377]
26. Wallenstein GV, Eichenbaum H, Hasselmo ME. The hippocampus as an associator of discontinuous events. *Trends Neurosci.* 1998; 21:317–323. [PubMed: 9720595]
27. Willingham DB. Some are more equal than others. *J.Am.Dent.Assoc.* 1980; 100:997–998. [PubMed: 6929851]
28. Yu AJ, Dayan P. Acetylcholine in cortical inference. *Neural Netw.* 2002; 15:719–730. [PubMed: 12371522]

a

### Sequential reaction time task

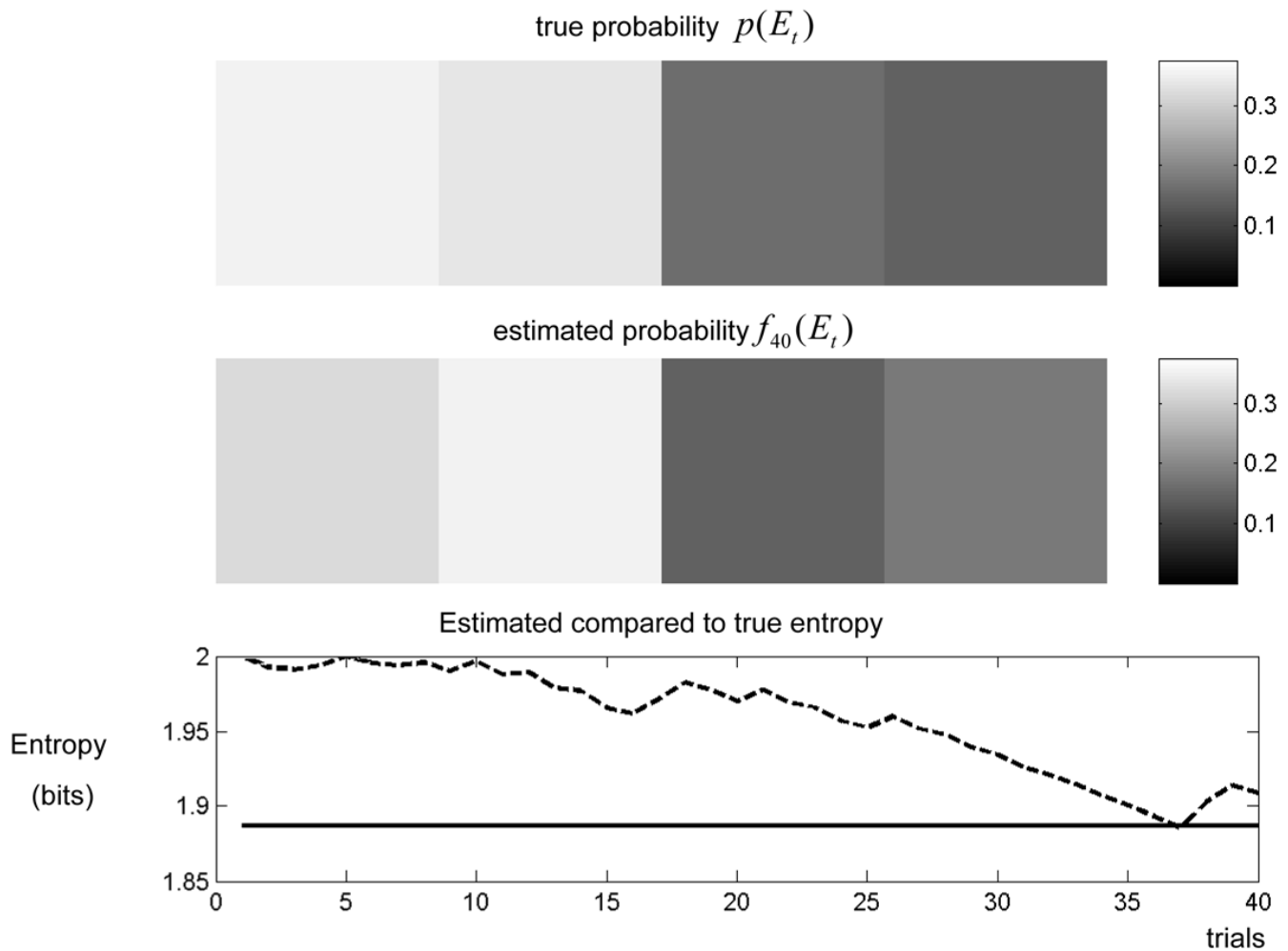


b

Transition matrix,  $p(E_t | E_{t-1})$ **Figure 1.**

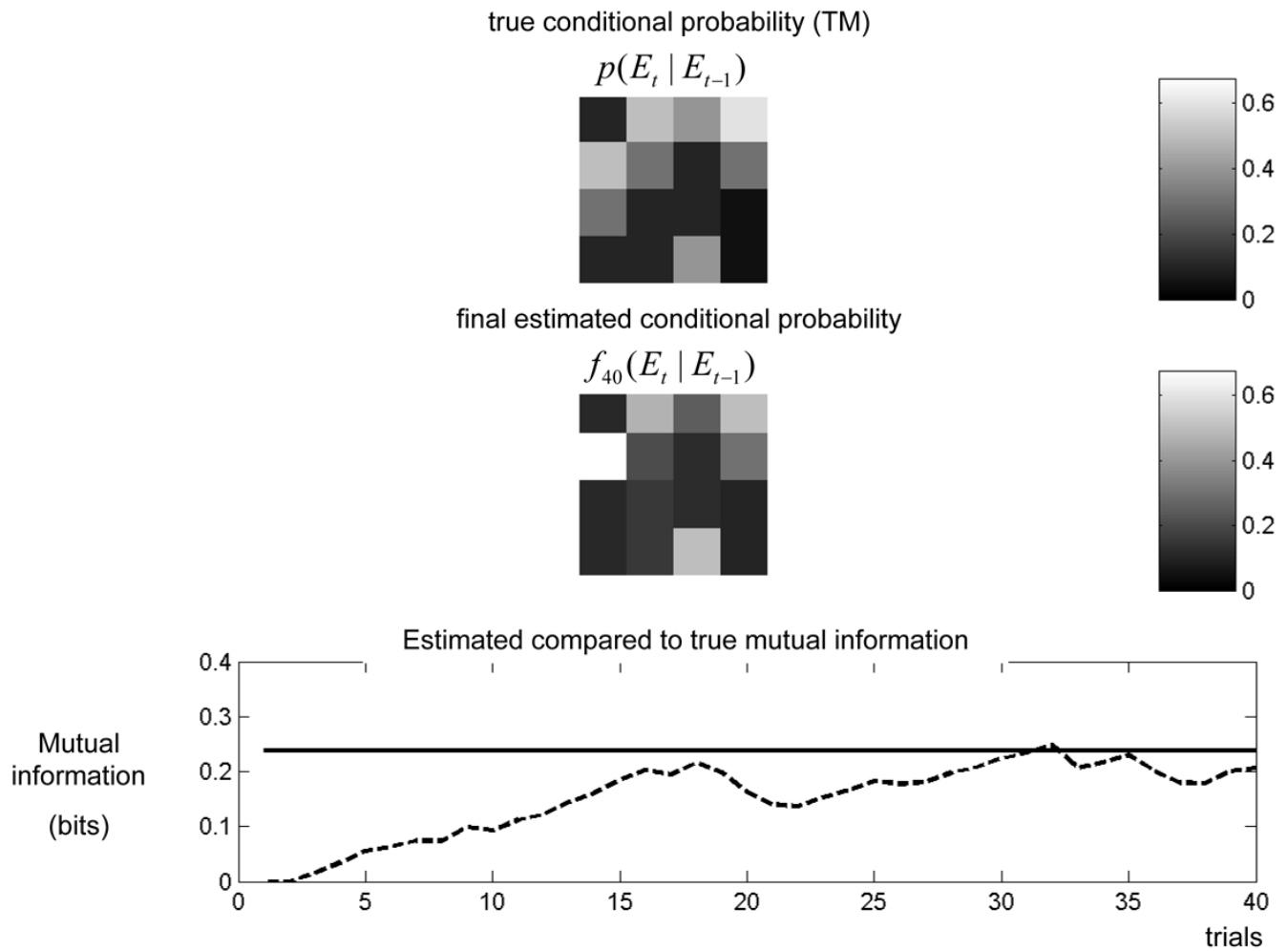
(a) The four alternative choice reaction time paradigm. Subjects were required to press the key indicating the position of the target in the row beneath (indicated in the top figure by '1' in bold type). Below is a schematic of samples displayed over one block. A row of possible targets was shown for 5 seconds before trials began. Each block consisted of 40 trials (first two and last one shown) with 12 blocks over the experiment. (b) An example of a transition matrix quantifying dependence among consecutive trials. The sequence of events produced by sampling from this distribution is an example of a 1<sup>st</sup> order Markov sequence. Gray-scale plot (top right) represents conditional probabilities.

## Estimating entropy



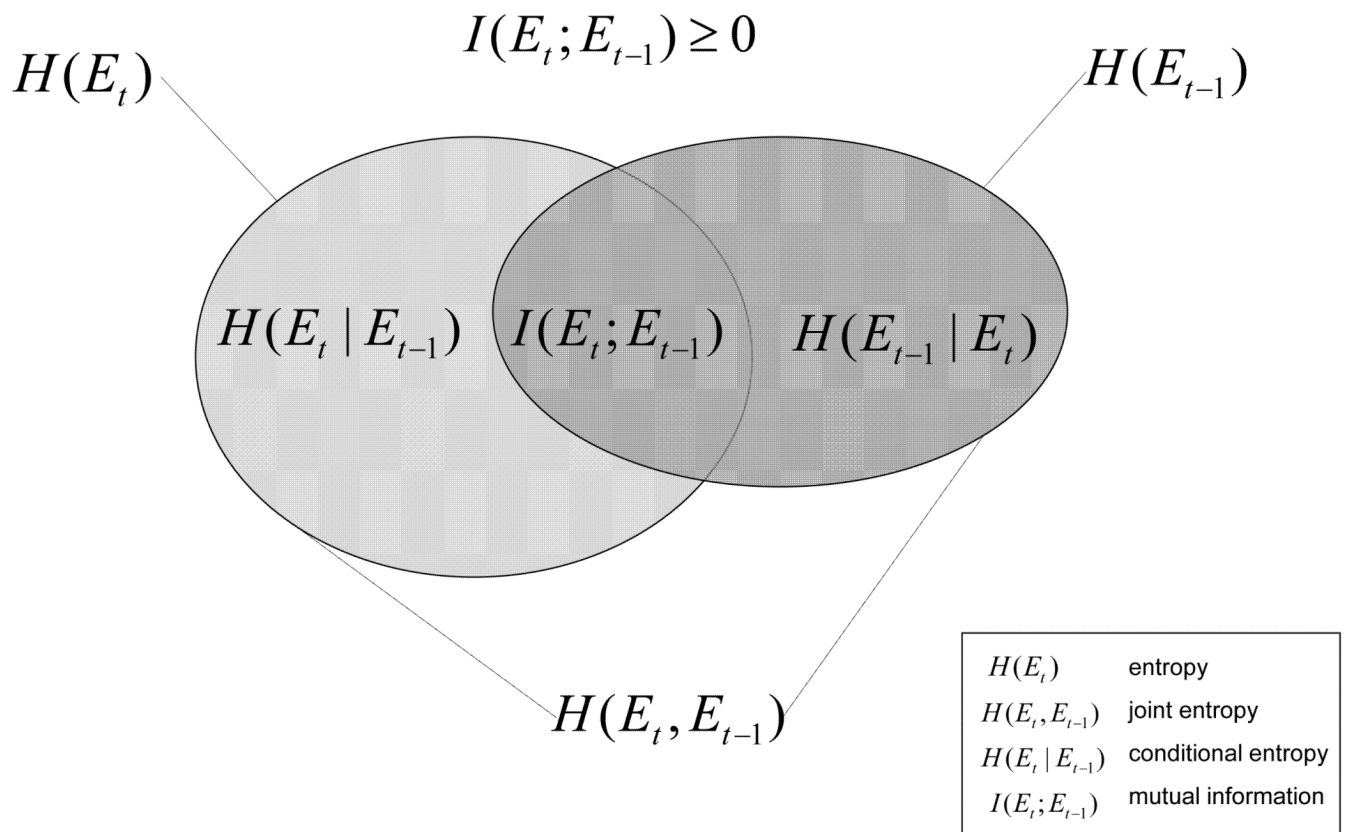
**Figure 2.** Estimating the marginal probability distribution. Top figure: grayscale representation of the distribution to be estimated,  $p(E_t)$ . Middle figure: the last estimate,  $f_{40}(E_t)$  of the marginal distribution. Lower figure: the estimated entropy (dashed line) decreasing towards the true value (solid line) as sampling increases.

## Estimating mutual information



**Figure 3.** Estimating  $p(E_t | E_{t-1})$ . Top figure: the conditional distribution to be estimated. Middle figure: the final estimate,  $f_{40}(E_t | E_{t-1})$  of this distribution given 40 trials. Lower figure: the mutual information of the estimated transition matrix (dashed line) rising towards the true value (solid line) as sampling increases.

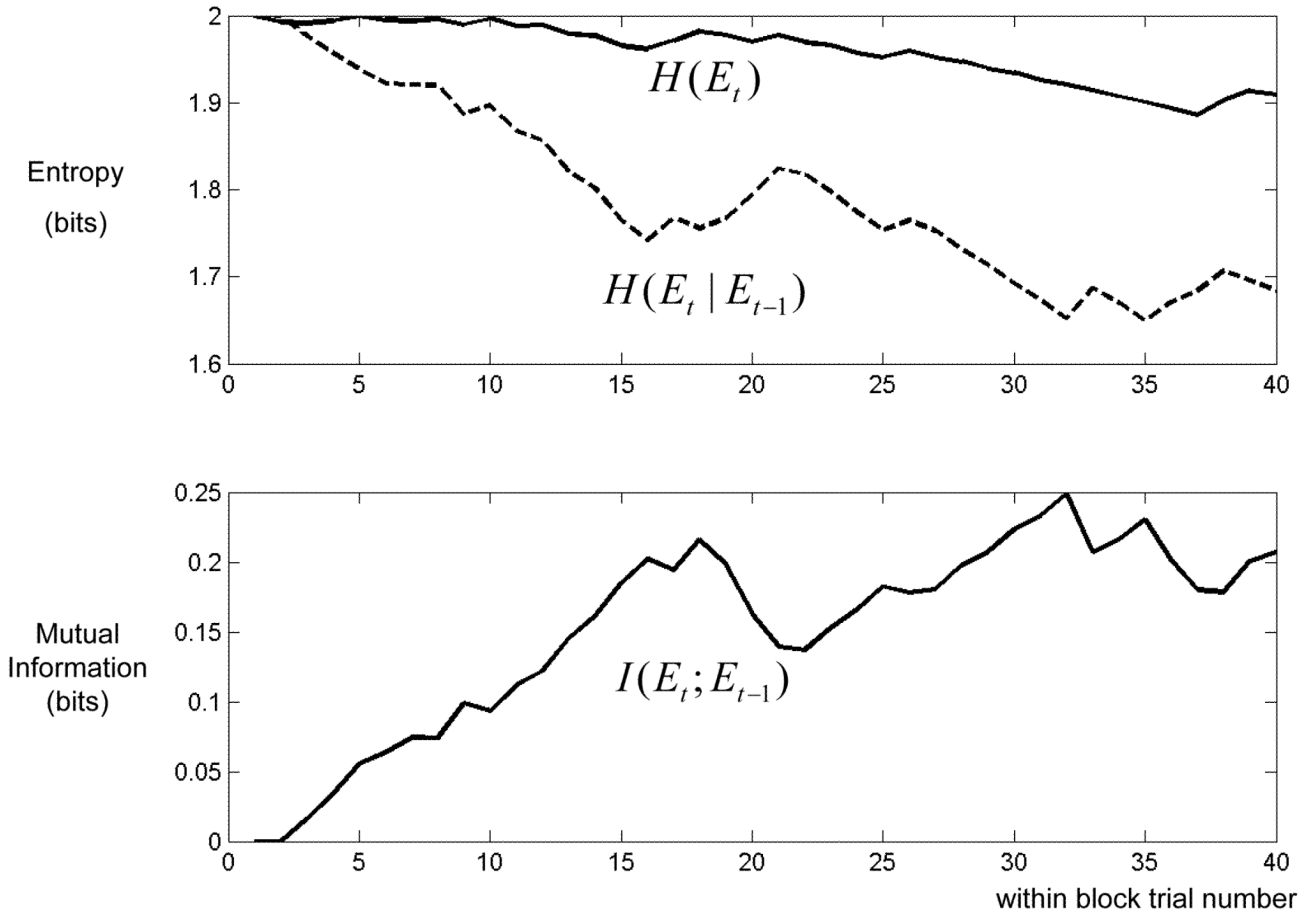
## Relationship among entropies and mutual information



**Figure 4.** Relationship between entropy, conditional entropy and mutual information illustrated using a Venn diagram. The degree to which ‘surprise’ is reduced in  $E_t$  conditional on  $E_{t-1}$  is measured by the mutual information.



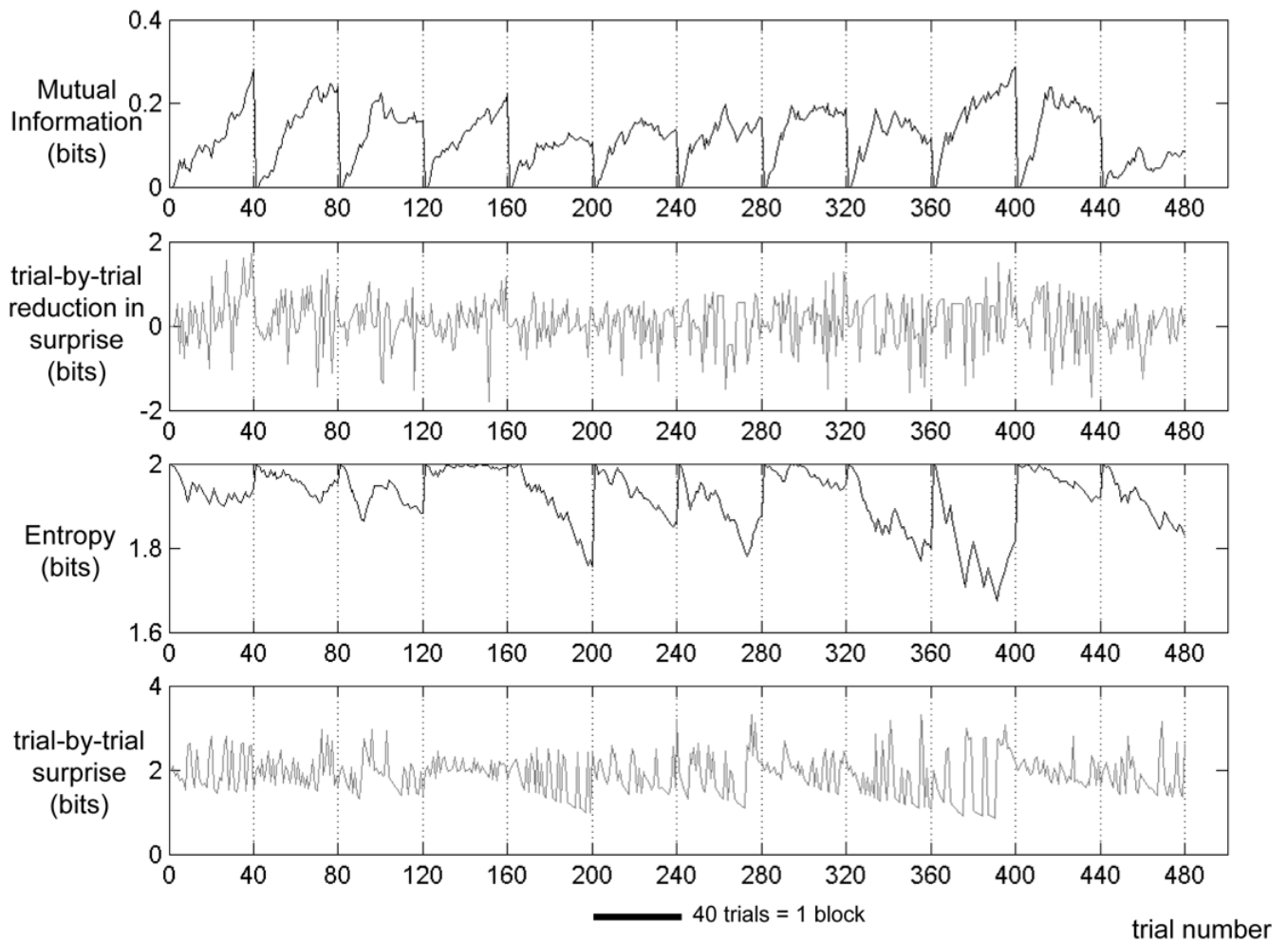
## Comparing entropy, conditional entropy and mutual information



$$H(E_t) - H(E_t | E_{t-1}) = I(E_t; E_{t-1})$$

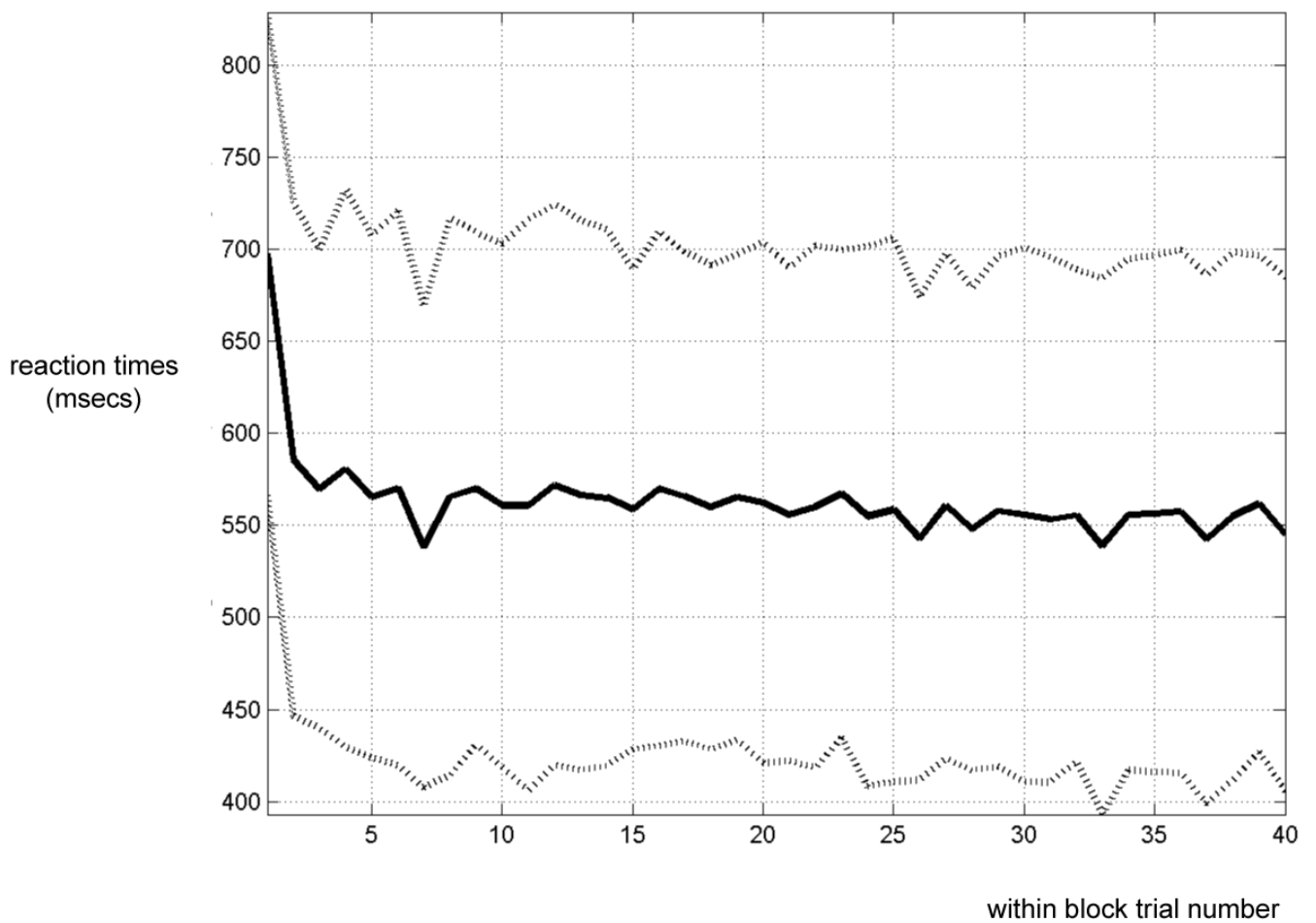
**Figure 5.** Example of entropy, conditional entropy and mutual information calculated over a single block. Dependence between consecutive trials, embedded within the transition matrix, is evidenced in the top graph by the reduction in conditional entropy (dashed line) compared to the entropy (solid line). The difference is the mutual information, which is a measure of the average contingency among consecutive events (bottom graph).

## Typical information theoretic indices (1 subject) over an experiment



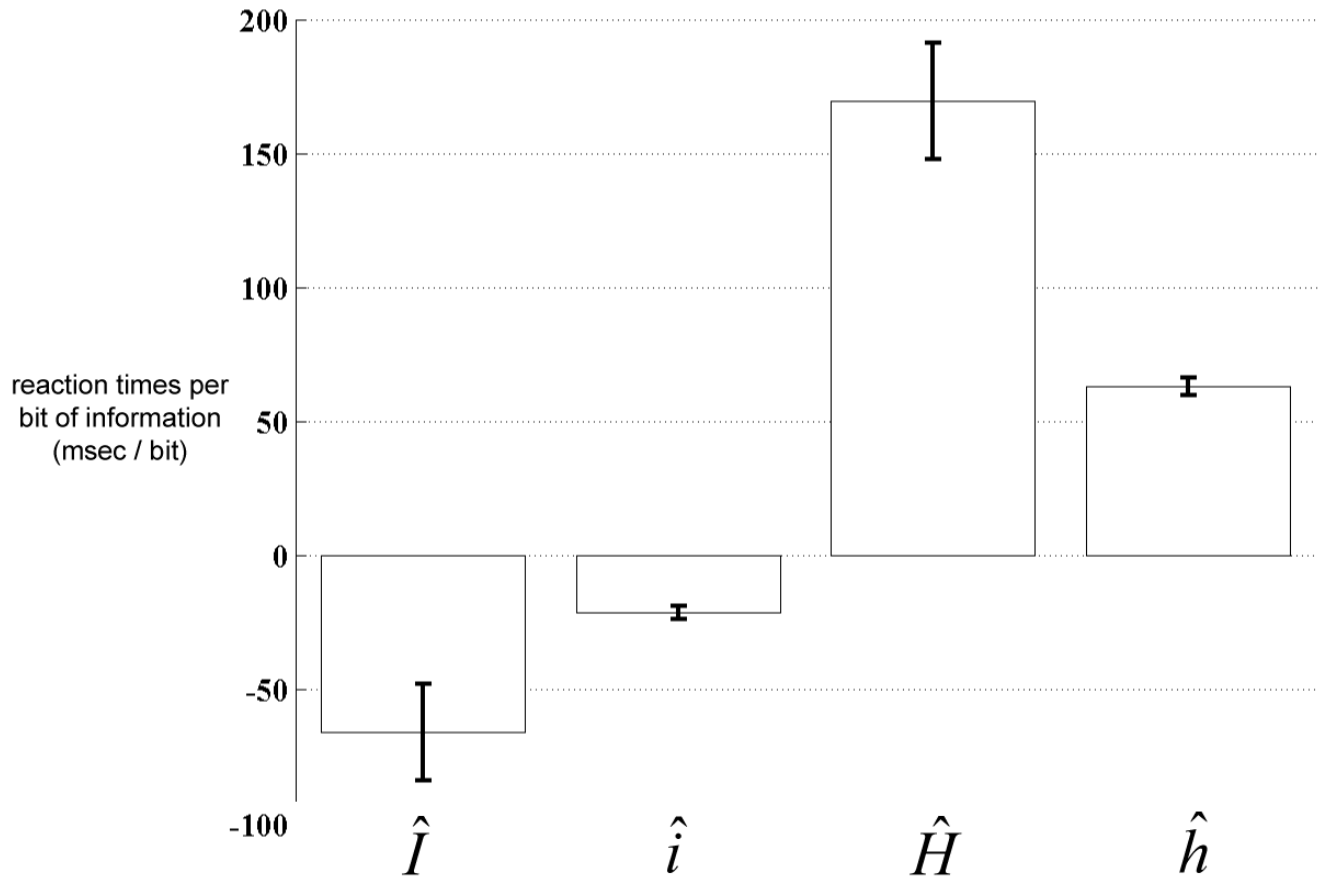
**Figure 6.** Plots of surprise, entropy, reduction in surprise and mutual information over 12 blocks during one experiment. These time series were calculated for each subject and used as regressors in a general linear model (SPM2 software) of the BOLD time series.

## Average reaction times over blocks and subjects



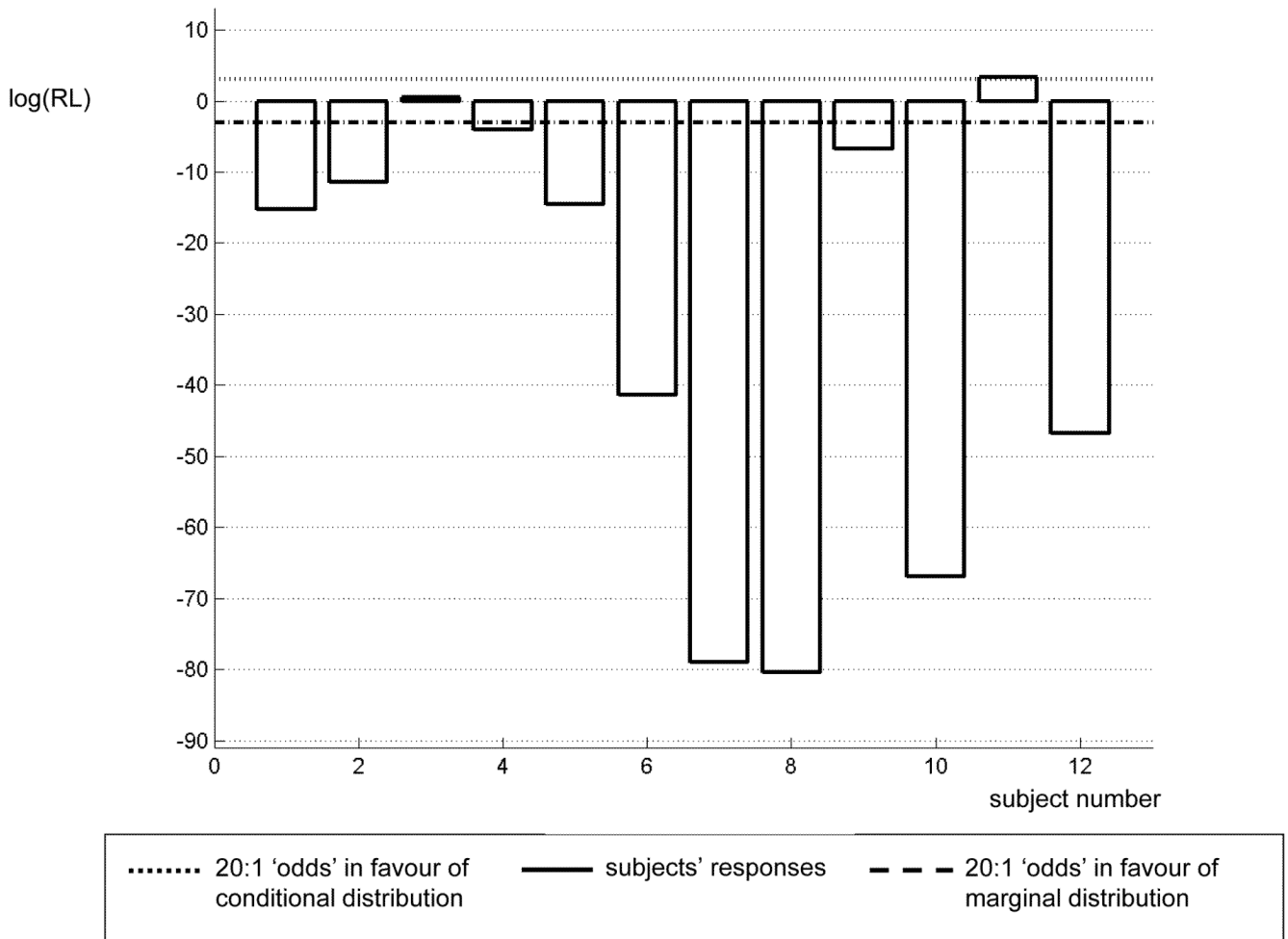
**Figure 7.** Temporal dynamics of learning. Average reaction times ( $\text{ms} \pm \text{standard deviation}$ ) for all subjects (25 in total) are shown. Reaction times decrease, on average, as sampling increases.

## AnCova of reaction times with information theoretic indices



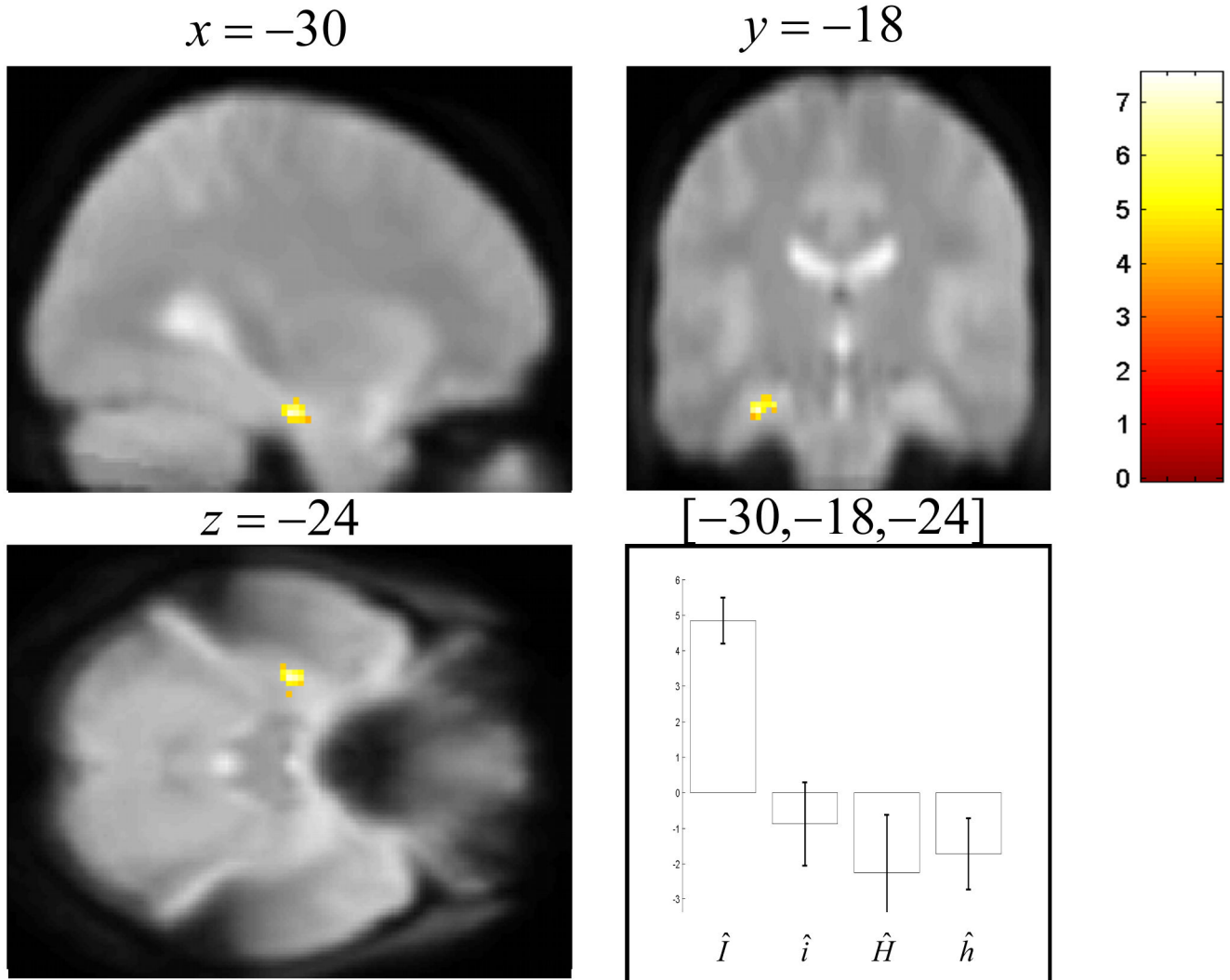
**Figure 8.** AnCova of reaction times (ms/bit  $\pm$  standard deviation). Less time is required to respond correctly to contingent events and response time increases with surprise.

### Relative likelihoods (RL) of cued generation task results



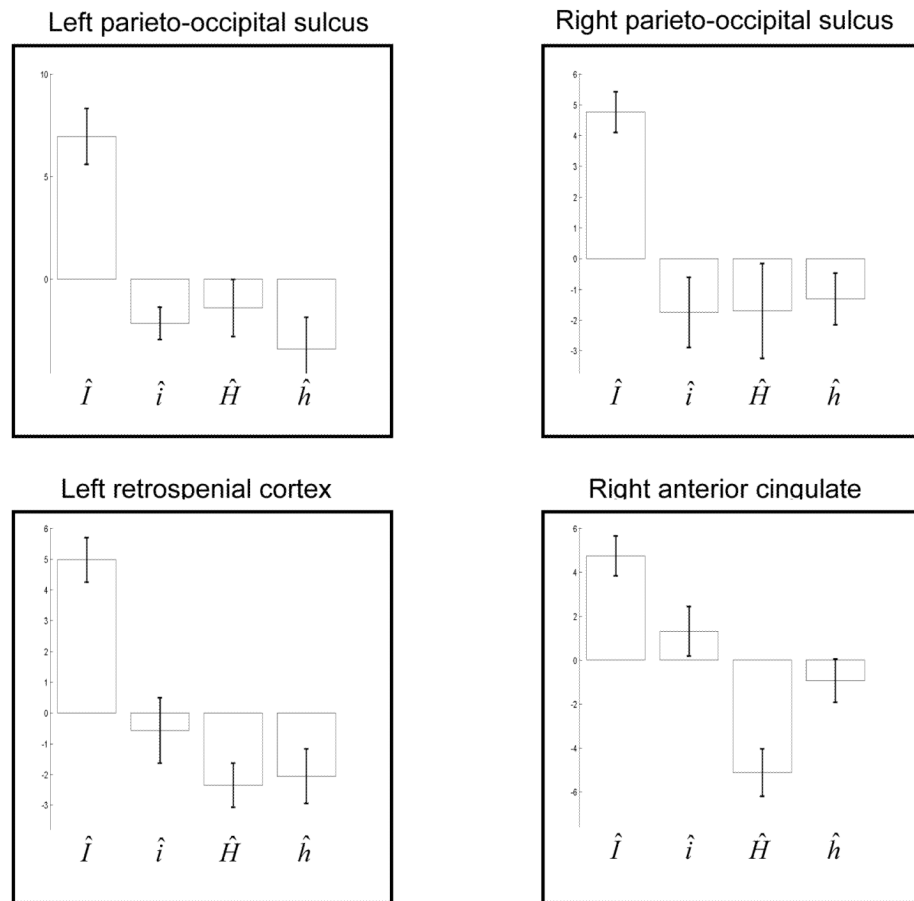
**Figure 9.** Cued generation task. Bar plot of the logarithm of relative likelihood (RL) for each subject. The RL compares the likelihoods of 2 models (conditional and marginal) of subject responses. The upper and lower bounds correspond to 20 and 0.05 (*i.e.*  $\pm \log(20) \approx \pm 3$ ). This corresponds to 20:1 'odds' in favor of the conditional and marginal model respectively.

## Random effects analysis: hippocampal response (Statistical Parametric Maps)



**Figure 10.** Random effects analysis of left hippocampal  $[-30, -18, -24]$  response to mutual information. The SPM (corrected at the cluster level to retain clusters at  $p < 0.05$ ; height threshold  $p < 0.001$ ) is overlaid on sections of a subject mean echo planar image. Parameter estimates ( $\pm$  standard error) from all 4 indices ( $\hat{I}$ ,  $\hat{i}$ ,  $\hat{H}$  and  $\hat{h}$ ) are shown in the bottom right panel. A significant effect was detected for mutual information only.

## Random effects analysis: parieto-occipital sulcus, retrosplenial cortex and anterior cingulate response



**Figure 11.** Parameter estimates ( $\pm$  standard error), from a random effects analysis, of all indices at local maxima within left (a) and right (b) parieto-occipital sulcus [-12, -62 16 and [22, -60, 12], (c) left retrosplenial cortex [-6, -48, 22] and (d) right anterior cingulate [4, 32, 16].

**Table 1**  
**Local maxima from random effects analysis: coordinates, Z scores and corrected p values**

List of regions (including coordinates, Z scores and corrected p-values) sensitive to mutual information

Brain Region	x,y,z maxima	Z score	p-values (corrected)
left hippocampus	-30 -18 -24	4.36	0.004
right parieto-occipital sulcus	22 -60 12	4.49	0.000
left parieto-occipital sulcus	-12 -62 16	3.9	0.000
left retrosplenial cortex	-6 -48 22	4.11	0.000
right anterior cingulate	4 32 16	3.87	0.001