

**Title:** Hippocampal-entorhinal transformations in abstract frames of reference

**Authors:** Raphael Kaplan<sup>1,2</sup> & Karl J Friston<sup>1</sup>

**Affiliations:**

1-Wellcome Centre for Human Neuroimaging, University College London, United Kingdom

2- Egil and Pauline Braathen and Fred Kavli Centre for Cortical Microcircuits, Kavli Institute for Systems Neuroscience, Norwegian University of Science and Technology, Trondheim, Norway

**Correspondence:** raphael.s.m.kaplan@ntnu.no

**Abstract:** Knowing how another's preferences relate to our own is a central aspect of everyday decision-making, yet how the brain performs this transformation is unclear. Here, we ask whether the putative role of the hippocampal-entorhinal system in transforming first person and extra-personal spatial cues during navigation extends to transformations in abstract decision spaces. In our functional magnetic resonance imaging study, subjects learned a stranger's preference for an everyday activity – relative to a personally known individual – and subsequently decided how the stranger's preference relates to other familiar people's preferences. Across reference frames, we observed signals in the retrosplenial cortex and hippocampal body during decisions that require precise memories for preferences. In contrast, the entorhinal cortex/subiculum exhibited reference frame-sensitive responses to the relative distance between the ratings of the stranger and the familiar choice options. Taken together, these data implicate the hippocampal-entorhinal system in the assimilation of knowledge in an abstract metric space.

## Introduction

Learning other people's attributes is facilitated by expressing personal preferences *ordinally*—whether we prefer one thing to another—and *metrically*—how much more we prefer one thing over another. Ordinal and metric coding are particularly important when acquiring knowledge about new people. This type of learning involves relating a new person's attributes to prior beliefs about other people; either by adopting an egocentric or extra-personal frame of reference. For instance, imagine you are preparing dinner for a foreign visitor that says, “I like spicy food”. If they are Vietnamese, their preference for spicy food is probably greater than Germanic tastes, despite everyone declaring the same preference.

On one hand, progress has been made in linking the hippocampus to the maintenance of an ordinal sequence or 'hierarchy' of personal attributes (Eichenbaum, 2015; Schiller et al., 2015; Kumaran et al., 2012, 2016) and, similarly, category learning (Zeithamova et al., 2008; Mack et al., 2017, 2018). Yet, the neural representation of metrically coded knowledge remains elusive, even though metric coding affords the transformation of knowledge learned via egocentric/relative and extra-personal/absolute frames of reference.

Clues about the neural computations underlying the transformation of abstract knowledge among frames of reference may come from research on the role of the hippocampal formation in path integration: the process of calculating one's position by estimating the direction and distance one has travelled from a known point. During path integration, specific sub-regions of the hippocampal formation have been implicated in integrating environmental and first person representations of space, in order to reach a desired location (McNaughton et al., 2006). In particular, grid cells in entorhinal/subicular areas are selectively active at multiple spatial scales when an animal encounters periodic triangular locations covering the entire environment (Hafting et al., 2005; Boccara et al., 2010); while hippocampal place cells code specific locations in an environment (O'Keefe & Dostrovsky, 1971). Working together with boundary vector cells in entorhinal/subicular areas – that code an environmental boundary at a particular direction and distance (O'Keefe & Burgess, 1996; Hartley et al., 2000; Burgess et al., 2000) – and head direction cells (Taube et al., 1990), spatially-

modulated neurons in the hippocampal formation are thought to serve collectively as a cognitive map of the environment (O'Keefe & Nadel, 1978; McNaughton et al., 2006; Hartley et al., 2013).

Notably, recent findings have extended the idea of map-like coding in the entorhinal cortex and subiculum to humans. Human entorhinal/subicular regions respond to both the distance of goal locations (Howard et al., 2014; Chadwick et al., 2015) and discrete abstract relations (Garvert et al., 2017), suggesting that entorhinal/subicular areas might represent abstract knowledge metrically along multiple dimensions. Taken together, these results imply that neural computations in the hippocampal formation related to spatial exploration, may also underlie the integration of information learned in different reference frames during more abstract types of memory-guided decision-making (Kaplan et al., 2017a).

We investigated whether specific brain regions, including sub-regions of the hippocampal formation – like the entorhinal/subicular area – facilitate switching between relative and absolute reference frames during memory-guided decisions. To test this hypothesis, we developed a novel experimental task, where healthy volunteers were first asked to rate 1-9, on a 0-10 scale, how likely (likelihood) they (*self*), a close friend (*friend*), and a typical person (*canonical*) were to partake in a variety of everyday scenarios (e.g., eat spicy food, read a book, cycle to work; Fig. 1A). To allow for strangers with more extreme ratings than the familiar individuals in the fMRI paradigm, subjects were restricted to rating between 1-9. Subsequently, subjects performed a forced-choice fMRI task (Fig. 1B), where they judged the proximity of a *stranger's* rating relative to the ratings of the self, friend, and canonical individuals, for a given everyday scenario (Fig. 1B). Specifically, on each (self-paced) trial, subjects were shown the preference of a *stranger* on a 0-10 scale, relative to a known person (*anchor*) and subjects had to select which of the two remaining (*non-anchor*) familiar individuals were closest to the *stranger*. After making a decision, a jittered (mean=2.43s) intertrial interval (ITI) period started, where a white fixation point overlaid on a black background appeared on the screen.

In our task, it was crucial that the *anchor* was always placed in the middle of the scale, to ensure that subjects used their prior knowledge in order to infer the *stranger's* absolute preference and form an appropriate mental number line of personal preferences (i.e., remembering the absolute positions of the different preferences on the number line, instead of the visible relative position of the

anchor in the middle of the screen with a stranger presented somewhere between 0-10; see Fig. 1C-D). In other words, subjects had to infer the *stranger's* position in relation to the *anchor's* 'true' rating for that scenario (i.e., the rating they gave for that anchor and scenario during the training period) and mentally rescale this information relative to the closest boundary (0 or 10). Note that this is a non-trivial task because the preferences of *strangers* were not conserved over attributes.

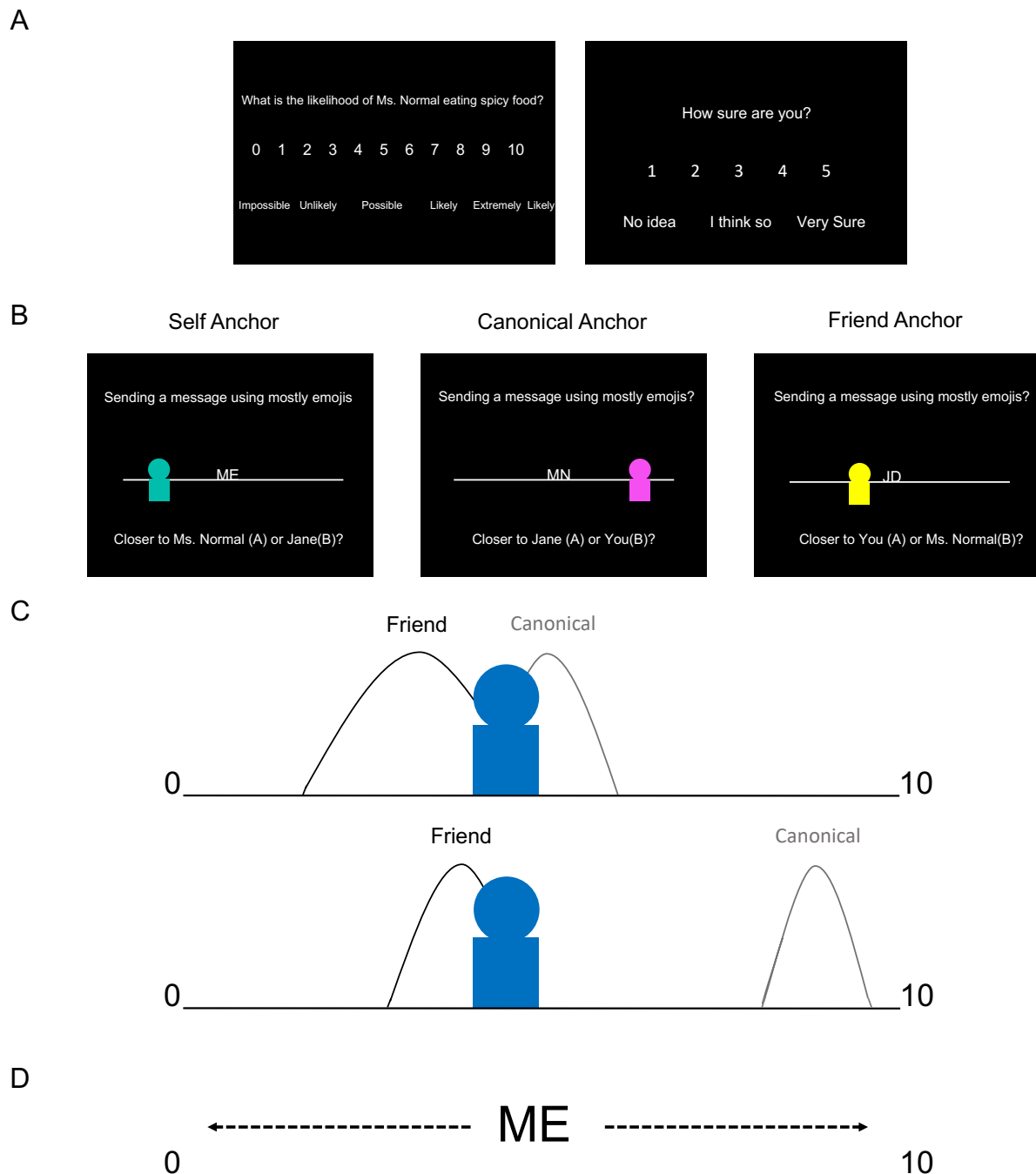


Figure 1. Experiment. *A*. Right before fMRI scanning, subjects were instructed to choose a friend with a different personality of the same gender. Subsequently, subjects rated from 1-9, on a 0-10 scale, how likely (likelihood) they (*self*), a close friend (*friend*), and the typical person (*canonical*) were to partake

in a variety of everyday scenarios (e.g., eat spicy food, read a book, cycle to work). To allow for strangers with more extreme ratings than the familiar individuals in the fMRI paradigm, subjects were restricted to rating between 1-9. Subjects also reported their confidence, on a 1-5 scale, for each scenario. *B. fMRI paradigm.* During a forced-choice task, subjects made a decision on the relative proximity of a *stranger's* likelihood rating for an everyday scenario relative to the likelihood ratings for the *self*, *friend*, and *canonical* individuals for that same scenario. On each self-paced trial (max. allowed response time 9s), subjects viewed a personal preference for a new *stranger* presented relative to one of the known individuals' initials (*anchor*) on a number line. Subjects had to determine which one of the two remaining (i.e., non-anchor) familiar individuals was closer to the *stranger's* rating. Crucially, the anchor individual (e.g., ME=self; MN=canonical 'Mr./Ms. Normal'; JD=friend's initial) was always placed in the middle of the scale, ensuring that subjects had to use memory of their ratings made before fMRI scanning to infer the *stranger's* absolute preference relative to the *anchor's* true rating and the ends of the number line. Additionally, subjects were instructed that the number line ranged from 0 and 10. For example, if the stranger was  $\frac{3}{4}$  to the right of the anchor with a rating of 9, the participant would infer that the stranger's rating was  $\frac{3}{4}$  of the way between 9 and 10 (the right boundary of the scale). Consequently, the participant would indicate that the stranger's rating would be approximately near 10. Notably, the numbers on the scale were not visible during the task. After making a decision, an intertrial interval (mean=2.43s) screen with a white fixation point in the center of a black screen appeared. *C. Illustration of the behavioral model.* Top illustration shows an ambiguous, less discriminable choice, while the bottom illustration shows a straightforward, highly discriminable choice. We quantified the difficulty of discriminating a particular choice by fitting a formal signal detection model, based on the relative distance between the two choice individuals on the scale and how confident subjects were in their ratings (e.g., comparing the *stranger* rating, represented by the blue avatar, with their rating for their *friend* and the *canonical* individual). Subjective confidence was represented by the standard deviation for each rating (e.g., curve width in the illustrations), where lower confidence entails higher standard deviations, and helped account for the influence of memory on choice behavior. *D. Anchor rescaling.* Illustration of how the stranger's rating is inferred by mentally rescaling the anchor individual's rating from its perceived relative position on the screen (5), to its absolute position (participant's rating) on the preference scale.

In summary, training subjects to think of personal preferences in the form of a mental number line allowed us to probe different everyday scenarios at various levels of discriminability and use a well-characterized signal detection model of decision-making (Tanner & Swets, 1954).

Discriminability was determined by the relative distance between individuals on a scale and how confident subjects were about a particular preference. Furthermore, this modeling approach allowed us to account for differing mnemonic demands related to the subject knowing their ratings better for themselves (*self*) than their *friend* and *canonical* ratings (Fig. 1C). Casting our paradigm in terms of coordinate transforms, we investigated how personal preferences are represented in the brain in two complementary ways. First, how the stranger's rating is inferred by mentally rescaling the anchor individual's rating from its perceived relative position on the screen, to its absolute position (the participant's actual rating) on the preference scale (i.e., anchor rescaling). Second, uncertainty in preference discrimination based on the relative distance between the stranger's rating and the ratings of the personally familiar individuals that the stranger is compared with (i.e., choice discriminability).

We subsequently refer to these two elements of representing personal preferences in our task as anchor rescaling (Fig. 1D) and choice discriminability (Fig. 1C), respectively.

## Results

### Behavior

Subjects' ratings of personally known individuals followed a consistent structure. As expected, subjects tended to rate the *canonical* individual towards the middle of scale, 5, with the *self* and *friend* ratings being more evenly dispersed between 1-9 (Fig. 2A). Subjects' ratings for the *canonical* exemplar and their *friend* were generally consistent before and after scanning; after scanning, subjects were on average within  $\pm 1$  of their original ratings for 67.7% (SD=8.15%) of trials and only made large deviations from their original ratings  $> \pm 3$  on 8.54% (SD=6.2%) of occasions (Fig. 2B). There were overlapping ratings (the same rating for two individuals on the same scenario) for a small subset of preferences, and there was a significant effect of known person (*self*, *canonical*, or *friend*) in the number of overlapping ratings ( $F(2,22)=6.74$ ;  $p=.005$ ), such that there was higher overlap between *self* and *friend* ratings than between other pairings (Fig. 2C). Crucially, trials with overlapping ratings were always considered correctly answered and not included in subsequent behavioral and fMRI analyses. When relating subjective confidence to rating consistency, confidence ratings significantly correlated with rating consistency ( $t(23)=-4.11$ ;  $p<.001$ ), where higher confidence ratings had more rating consistency. These results suggest a metacognitive validity of the subjective ratings, relating to memory for specific personal preferences.

During the fMRI experiment, subjects made 72.6% of decisions correctly (SD=1.13%;  $n=24$ ), after eliminating trials for which answer was correct (i.e., same or equidistant ratings). The mean reaction time (RT) was 4.07 s (SD=.803 s). There was a significant effect of condition ( $p<.05$ ) for RT ( $F(2,22)=4.57$ ;  $P=.022$ ; Fig. 2D), but not performance ( $F(2,22)=2.73$ ;  $P=.087$ ). The RT effect was driven by significantly higher RT for *self* versus *friend* anchor conditions ( $t(23)=2.71$ ;  $P=.012$ ; Fig. 2D). Notably, there was a negative correlation between trial by trial RT and accuracy ( $t(23)=-5.51$ ;  $p<.001$ ), i.e., quicker RT for accurate choices.

We then related RT and accuracy to the relative distance between the stranger's rating and the ratings of the non-anchor individuals presented as choice options for each trial (i.e., choice

discriminability without the use of a model and confidence ratings). As expected, we found that this distance negatively correlated with RT ( $t(23)=-5.42; p<.001$ ) and positively correlated with accuracy ( $t(23)=15.8; p<.001$ ); i.e., larger distances related to quicker and more accurate decisions. Highlighting how the metric nature of these relative distances influenced behavior, we observed a steady increase in the quickness (Fig. S1) and accuracy (Fig. 2E) of the participant's responses when the stranger's rating was closer to one choice option versus another. We also found that mean deviations in rating consistency positively correlated with RT ( $t(23)=2.69; p=.013$ ) and negatively correlated with accuracy ( $t(23)=-5.23; p<.001$ ); i.e., scenarios with inconsistent ratings/poor memory related to slower and less accurate decisions. Further relating these relative distances to subjects' behavior, we then used a signal detection model fit to subjects' performance that was based on relative rating distances and subjective confidence ratings. Given the relationship between subjective confidence and rating consistency, subjective confidence ratings also helped account for the different mnemonic demands induced by the different conditions.

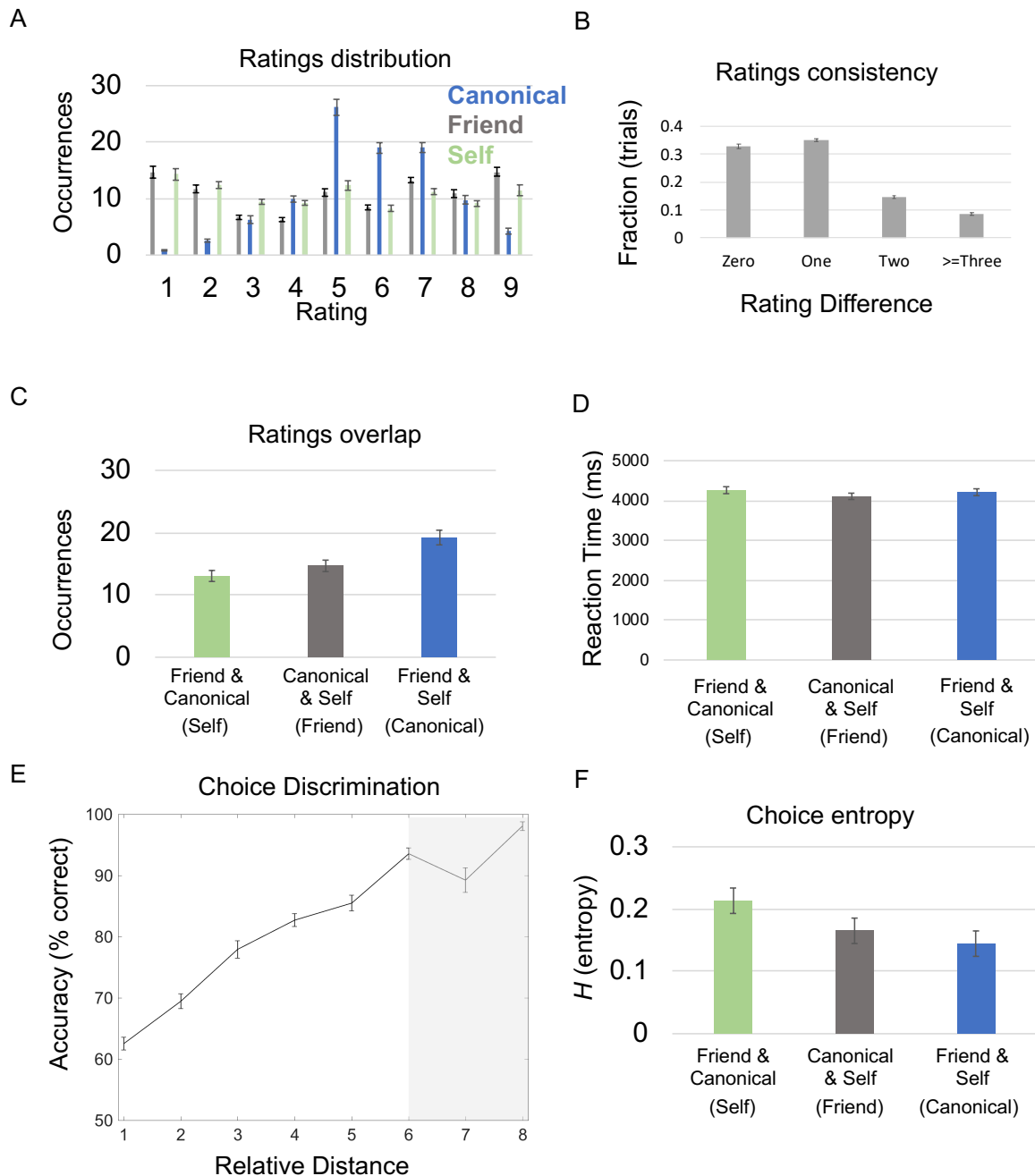


Figure 2. Behavioral Results. *A.* Mean ratings across subjects for each familiar individual and every scenario. Occurrences are out of the 100 total trials per condition. *B.* Rating Consistency. Difference in ratings pre- and post-fMRI scanning for *friend* and *canonical* individuals *C.* Ratings and Overlap. Significant effect of condition for ratings overlap (the same ratings) between the two individuals ( $p < 0.001$ ). Individuals being compared are listed below each bar with the corresponding anchor/condition name listed in parentheses. Occurrences are out of the 100 trials per condition. *D.* Reaction time: Significant effect of condition for reaction time/decision speed ( $p = .022$ ). Individuals being compared are listed below each bar with the corresponding anchor condition listed in parentheses. *E.* Relative distances and performance: Significant relationship ( $p < .001$ ) between accuracy and the relative distance between strangers' ratings and the non-anchor individuals for each trial. 50% represents chance level of accuracy. The absolute value of relative distances were rounded to the closest integer and plotted from 1 to 8. 7 & 8 on the x-axis are shaded in gray because 18/24 and 13/24 of subjects had trials with relative distances of 7 & 8, respectively. *F.* Choice entropy: Significant effect of condition for choice entropy ( $p < 0.001$ ). Individuals being compared listed below each bar, with corresponding anchor listed in parentheses. All error bars showing mean  $\pm$  SEM.



### *A computational model of latent preferences*

We quantified the difficulty of discriminating a particular preference by fitting a formal signal detection model (Tanner & Swets, 1954) based on the relative distance between the two choice (*non-anchor*) individuals on a scale and subjective confidence. Specifically, we characterized discriminability using the entropy of decision probabilities based on a softmax function of likelihood and confidence ratings (Fig. 1C). Using this modeling approach allowed us to characterize the relative distance between the ratings, along with the subjects' confidence ratings, in a single measurement of trial-specific choice discriminability. Importantly, high entropy corresponds to lower choice discriminability that is induced by similar ratings. To estimate the requisite softmax (sensitivity or precision) parameter, we modeled performance in terms of entropy ( $H$ ) over trials (and subjects) using a simple linear regression model. The ensuing behavioral model provided an estimate of the precision parameter ( $\beta$ ) and associated measure of trial-specific choice discriminability for each subject.

In detail, we want to score the difficulty of each trial. This difficulty is the entropy ( $H$ ) of the choice probability ( $p$ ), based upon a softmax function of the log odds ratio of a target (*stranger*) preference ( $T$ ) being sampled from the distributions of (*non-anchor*) reference subjects A and B. Let the preference of individual  $i$  have a mean  $\mu_i$  and standard deviation  $\sigma_i$ , then the log probability of sampling  $T$ , from the preference distribution of the  $i$ -th reference subject is (ignoring constants):

$$L_i = -\frac{1}{2} \frac{(T - \mu_i)^2}{\sigma_i^2}$$

Equation 1

If we assume a precision of  $\beta$ , then the probability of choosing subject A over B is a softmax function of the log odds ratio:

$$P_A = \frac{\exp(\beta \cdot L_A)}{\sum_{i \in \{A, B\}} \exp(\beta \cdot L_i)}$$

Equation 2

and the difficulty is the entropy:

$$H = \sum_{i \in \{A, B\}} -P_i \ln P_i$$

### Equation 3

This will give difficulties between 0 and  $\ln(2)$  i.e., 0.6931. Given the mean and standard deviation reported by subjects, we can evaluate the probability of choosing each (*non-anchor*) individual on each trial, given the *stranger's* target preference.

This model of choice behavior provided trial-specific measures of choice uncertainty ( $H$ ) that enabled us to identify its fMRI correlates. Following previous work (Daw et al., 2006; Glascher et al., 2010; Daw, 2011), we used the mean precision parameter ( $\beta$ ), evaluated over subjects to compute trial-specific choice entropies as a predictor for our fMRI responses. As expected, choice entropy negatively correlated with performance over trials ( $t=-10.9$ ;  $p < .001$ ). As might be expected, we observed a significant effect of condition for mean choice entropy ( $F(2,22)=17.6$ ;  $P < .001$ ; see Fig. 2F), where choice entropy was significantly higher for *self* than *canonical* ( $t(23)=6.06$ ;  $P < .001$ ) or *friend* ( $t(23)=4.23$ ;  $p < .001$ ) anchor trials (Fig. 2D). We assume that this was likely due to lower confidence (i.e., memory demands) for *friend* and *canonical* ratings, relative to ratings of *self* preferences.

## fMRI Analyses

### *Choice Discrimination*

To test how strangers' preferences are compared to personally known people, we characterized the brain's response to uncertainty in preference discrimination. Specifically, we were interested in what brain regions related to choice discrimination in different reference frames. In a whole-brain analysis, we investigated which brain regions responded to choice entropy (i.e., how discriminable the choices were). We observed significant effects in the bilateral retrosplenial cortex (RSc:  $x=6, y=-37, z=16$ ;  $Z$ -score= 4.57; cluster-level FWE  $p=.042$  at  $p < .0001$  uncorrected; Fig. 3A-B), the right hippocampal body ( $x=36, y=-19, z=-10$ ;  $Z$ -score=3.67; small-volume corrected (SVC) peak-voxel  $p=.036$ ; Fig. 3A-B), and a cluster in the body of the hippocampus that peaked in the left entorhinal/subicular region ( $x=-21, y=-22, z=-18$ ;  $Z$ -score 3.27; SVC peak-voxel  $p=.049$ ) for increasing choice entropy (i.e., a positive correlation between the fMRI signal and choice entropy/ambiguity). We did not find any significant regions responding to increasing choice entropy elsewhere.

Furthermore, we did not observe any correlation between the fMRI signal and decreasing choice entropy (i.e., a negative correlation between the fMRI signal and choice entropy).

In a whole-brain analysis – comparing differences in entropy effects between the three conditions (*self*, *friend*, *canonical*) – the only significant effect of condition was in a right entorhinal/subicular area ( $x=24, y=-28, z=-22$ ; F-stat:9.79; Z-score=3.44; SVC peak-voxel  $p=.035$ ; Fig. 3C), extending into anterior parahippocampal cortex. In parallel, the left subiculum/entorhinal cluster that exhibited a significant effect of entropy, also exhibited an interaction between condition and entropy ( $t(23)=2.54$ ;  $p=.018$ ). In other words, the effect of discriminability depended upon the *anchor* condition.

Subsequent t-tests on the bilateral entorhinal/subicular regions exhibiting choice entropy effects indicated that the positive correlation with entropy was higher for *canonical* and *friend* anchor trials. The bilateral entorhinal/subicular effect (right: Z-score=3.98; peak voxel  $p=.005$ ; left: Z-score=3.68; peak voxel  $p=.014$ ; p-values are SVC within region; Fig. S2) was the strongest effect of any area in the brain for this contrast. Note that *canonical* and *friend* anchor trials involved choices with *self* preferences. Subjects were either deciding whether a *stranger's* rating for a scenario was closer to the *self* versus *canonical* individual (*friend* anchor), or the *self* versus the *friend* (*canonical* anchor). Conversely, *self* anchor trials involved deciding whether a *stranger's* rating was closer to the *friend* versus the *canonical* individual. In other words, entorhinal/subicular areas responded to more fine-grained/ambiguous choices involving self-comparisons, but readily discriminable choices otherwise.

One alternative explanation for the observed choice entropy effects was that they could be purely explained by memory demands, instead of the relative distance between the stranger's rating and the non-anchor individuals. Hence, we conducted an analysis to test whether our hippocampal formation—including the right entorhinal/subicular region—and retrosplenial choice entropy effects were primarily explained by covariance due to memory demands. In this second general linear model (GLM2), we accounted for covariance due to subjective confidence ratings and RT (i.e., memory demands) in order to examine whether there were any residual effects related to the relative distance (i.e., choice discriminability) between the strangers' ratings and the ratings of the non-anchor

individuals for each trial. After accounting for memory-related factors, there remained a significant effect of condition ( $F(2,22)=6.67;p=.005$ ; Fig. 3D) in the right entorhinal/subicular region, where this effect was primarily driven by responses to greater relative distances for friend versus canonical anchor trials ( $t(23)=3.01;p=.006$ ). In contrast, we did not observe any significant effect in the right hippocampal body (main effect:  $t(23)=1.16$ ;  $p=.258$ ); interaction: ( $F(2,22)=.069;p=.933$ ), nor bilateral RSc (main effect:  $t(23)=.94;p=.357$ ; interaction:  $F(2,22)=.304;p=.741$ ) related to choice discrimination after accounting for memory demands. In other words, choice entropy effects in the right entorhinal/subicular area of the hippocampal formation primarily related to the relative distances between the stranger's rating and the non-anchor individuals, while the choice entropy effect in the right hippocampal body was primarily related to the increased memory demands of ambiguous choices.

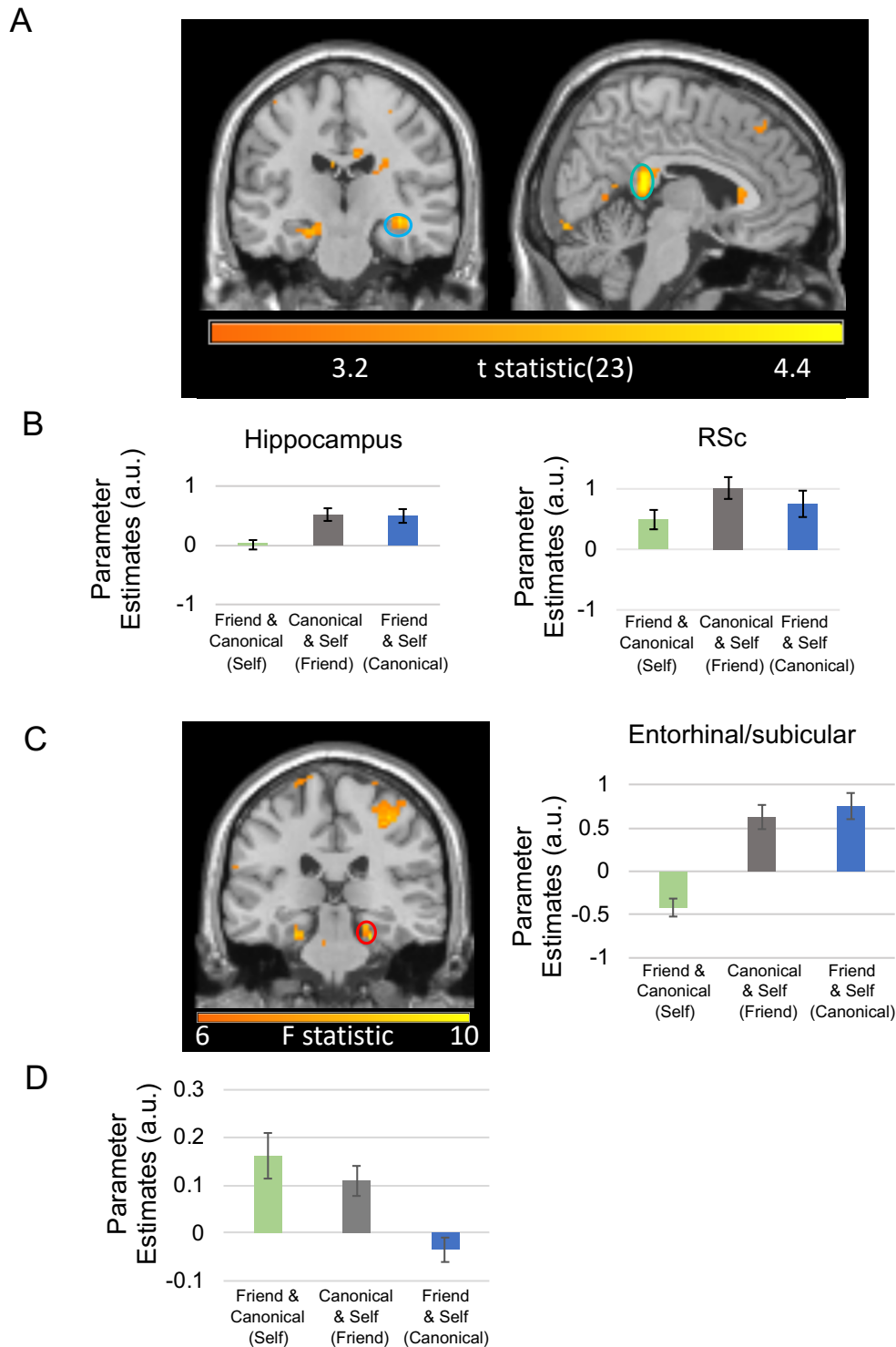


Figure 3. Choice Discrimination Effects. *A*. Regions significantly responding to high choice entropy. Left: Coronal image showing right hippocampus circled in light blue. A portion of the left subicular cluster is also visible. Right: Sagittal image showing retrosplenial cortex (RSc) circled in turquoise. *B*. Effect size for a 10-mm sphere around the right hippocampal and RSc peaks (mean  $\pm$  SEM). A positive effect size indicates a positive BOLD correlation with choice entropy (i.e., more ambiguous choices), whereas a negative effect size indicates a negative BOLD correlation with choice entropy (i.e., straightforward choices). Individuals being compared listed below each bar, with corresponding anchor individual (condition name) provided in parentheses. *C*. Left: Coronal image of right entorhinal/subicular region exhibiting effect of choice entropy by condition circled in red. Portion of left entorhinal/subicular region showing same effect is also visible. Right: Effect size for a 10-mm

sphere around right entorhinal/subicular region exhibiting effect of choice entropy by condition (mean  $\pm$  SEM). *D.* Effect size for a 10-mm sphere around the same right entorhinal/subicular region, which also displays an effect of relative distance (choice discrimination) by condition (mean  $\pm$  SEM), after accounting for covariance due to memory demands. A positive effect size indicates a positive BOLD correlation with the relative distance between the stranger's rating and the closest non-anchor individual's rating versus the other (i.e., straightforward choices), whereas a negative effect size indicates a negative BOLD correlation with smaller relative distances in the same comparison (i.e., ambiguous choices). All highlighted regions survived FWE correction for multiple comparisons at  $p < 0.05$  and are displayed at an uncorrected statistical threshold of  $p < .005$  for display purposes.

### *Anchor Rescaling*

In our task, reference frame transformations rely on flexibly relating different known individuals' ratings to each other. We wanted to capture the initial demands of mentally rescaling the anchor and, consequently, the corresponding stranger from its observed (relative) position on the screen to its actual rating (absolute position) on the scale. To test this, we investigated which regions responded to how far the *anchor's* preference deviated from the middle of the scale (Anchor Rescaling;  $|\text{anchor rating} - 5|$ ). In a whole-brain analysis, the only significant main effect of increased anchor rescaling towards the limits of the scale was in the superior parietal lobule (SPL) bilaterally (left:  $x = -57, y = -52, z = 38$ ;  $Z\text{-score} = 4.70$ ; FWE cluster  $p < .001$ ; right:  $x = 33, y = -61, z = 52$ ;  $Z\text{-score} = 4.46$ ; FWE cluster  $p < .001$ ; Fig. S3). Critically, the anchor rescaling effect in right SPL could not be explained by mnemonic demands (i.e., subjective confidence ratings and RT), because after accounting for these mnemonic demands, there was still a significant effect of greater anchor rescaling in the right SPL ( $t(23) = 2.21; p = .037$ ), though not left SPL ( $t(23) = 1.77; p = .089$ ).

However, we did not observe any significant effects for decreased anchor rescaling (maintaining the rating in the middle of the scale). Likewise, we did not observe any regions showing an interaction with anchor condition.

### *Accuracy*

Further characterizing the functional contribution of different brain regions, we asked if regional responses during memory-guided decisions depended on whether subjects made a correct, or incorrect choice for each trial. In a whole-brain analysis, we observed significant activation for correct trials in the right ventral striatum ( $x = 24, y = 8, z = -6$ ;  $Z\text{-score} = 4.31$ ; FWE cluster  $p = .004$ ; Fig. 4A-B),

extending into ventromedial prefrontal cortex, and another separate activation in left posterior parietal cortex (PPC) in the depth of intraparietal sulcus ( $x=-24,y=-61,z=40$ ;  $Z\text{-score}=4.37$ ; FWE cluster  $p=.005$ ; Fig. S4). Additionally, testing whether any regions exhibiting entropy effects also showed performance correlates, we found that both bilateral RSc ( $t(23)=2.54;p=.019$ ) and left subiculum/entorhinal area ( $t(23)=3.34;p=.003$ ) activity related to correct choices. We did not observe any significant activation that preceded incorrect choices, nor significant interactions between anchor condition and accuracy.

Finally, in a follow-up analysis to test whether the right ventral striatum accuracy effect was driven by choices where highest/lowest judgment between the two choice individuals needed to be made. We split decision-making trials in terms of whether a choice required determining which choice option was highest/lowest, instead of closest. Ignoring the three anchor conditions, we split the data into two conditions; one where the target was between the two choice individuals, and the second where the target was between either higher or lower than both choice individuals. We asked whether subjects who performed better for the latter scenario, also recruited the ventral striatum when making highest/lowest versus proximity judgments. We found a between-subject correlation between our right ventral striatal effect for highest/lowest judgments and subject performance for highest/lowest versus proximity judgments ( $r=.561;p=.004$ ) (Fig. 4C).

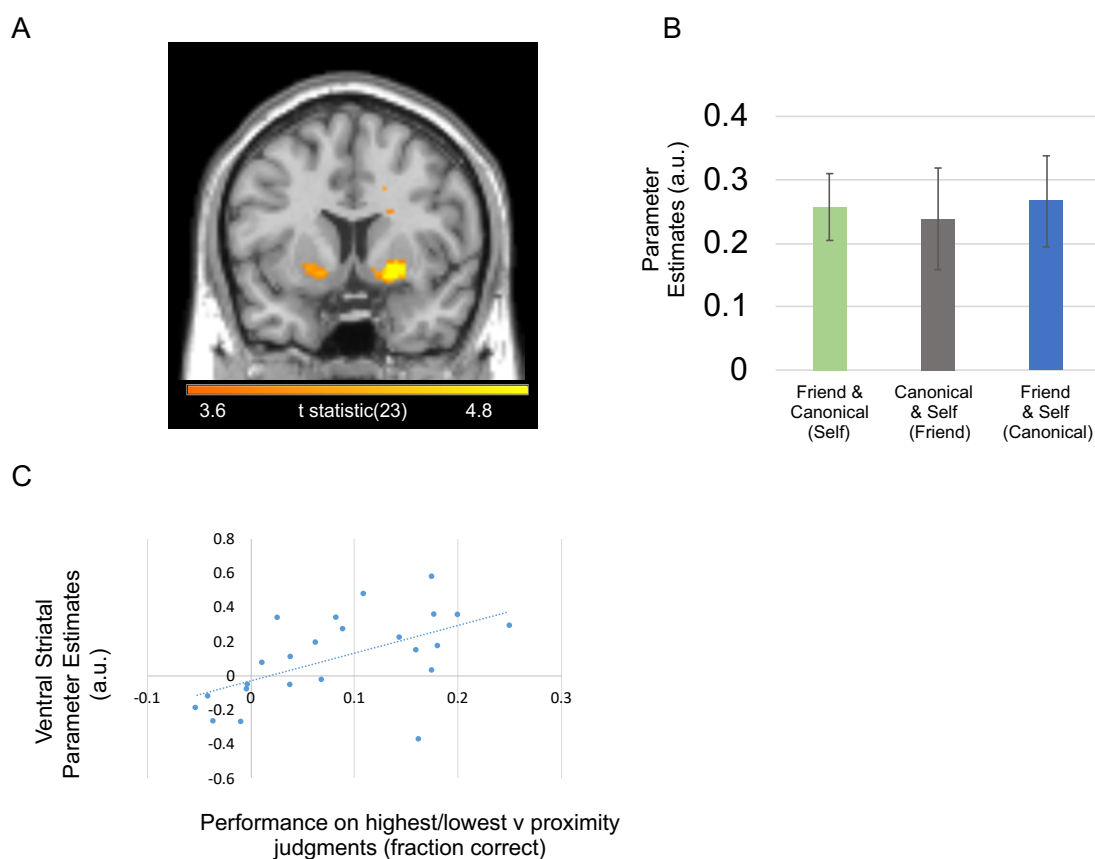


Figure 4. Striatum and Decision Accuracy. *A*. Ventral striatal activity related to correct versus incorrect choices. The ventral striatal cluster survived cluster-level FWE correction at  $p < 0.05$  and is displayed at an uncorrected statistical threshold of  $p < .001$ . *B*. Effect sizes for a 10-mm sphere around right ventral striatum peak voxel (mean  $\pm$  SEM). Individuals being compared list below with anchor in parentheses. A positive effect size indicates a positive BOLD correlation with correct choices. *C*. Plot showing between-subject correlation for subjects' ventral striatal fMRI signals with behavioral performance. Ventral striatal effects were for highest/lowest judgments versus proximity judgment trials (GLM2). Behavioral performance was taken from subjects' performance on trials, where highest/lowest judgments could be used, versus their performance on trials where a proximity judgment could be used. Subjects who performed better for highest/lowest judgments exhibited increased ventral striatal activity for the same contrast.

## Discussion

Using fMRI, a signal detection model, and a novel memory-guided decision-making paradigm, we asked how different brain regions integrate knowledge of personal attributes within different reference frames (Fig. 1). We observed hippocampal and retrosplenial cortex (RSc) responses to fine-grained choices involving any frame of reference that was primarily due to the increased memory demands of such trials (Fig. 3A-B). When focusing on reference frame-sensitive responses, we identified an entorhinal/subicular region that related to the relative distance between the stranger's rating and the ratings of choice options (Fig. 3C-D). Notably, subiculum/entorhinal and RSc responses were higher leading up to correct choices. In parallel, superior parietal lobule activity



increased when the anchor required a mental shift from the middle of the scale towards the periphery, which was partially due to the extra time it took to rescale the anchor's position (Fig. S3). Finally, we found that striatal responses also preceded correct choices, which were partially driven by decisions about which individual had the highest or lowest preference rating (Fig. 4). In what follows, we relate these memory-guided choice findings to the wider literature – and to the different hippocampal formation responses we observed. We then speculate on potential neural computations induced by our task.

### *The role of the hippocampus and retrosplenial cortex during option discrimination*

Our results build on previous results in rodents and humans showing metric coding of learned non-explicitly spatial representations in the hippocampus (Tavares et al., 2015; Aronov et al., 2017). Modeling knowledge-dependent ratings with a confidence measure, representing the uncertainty of a rating (Fig. 1C), allowed us to capture mnemonic influences on choice discriminability in more detail. Hippocampal fMRI signal increases were primarily associated with deliberation and subjective confidence, suggesting that hippocampal responses observed in our task may reflect increased sampling of personal experience during memory-guided decision-making (Shadlen & Shohamy, 2016; Bornstein & Norman, 2017; Bakkour et al., 2018). Notably, the hippocampal responses observed here closely resemble putative pattern separation signals in episodic memory (Marr, 1971). Future work may determine whether hippocampal responses during less discriminable knowledge/memory-guided choices parallel non-linear hippocampal attractor dynamics during human spatial navigation (Steemers et al., 2016) and how these responses arise in hippocampal subfields during pattern separation (Leutgeb et al., 2007; Chen et al., 2011; Yassa & Stark, 2011, Duncan et al., 2012; Kumaran, 2012; Kyle et al., 2015; Stokes et al., 2015).

In parallel, we observed retrosplenial responses when successfully resolving more ambiguous choices. It is important to note that our task always requires the subject to reference themselves in some aspect, either as the anchor used to infer the *stranger's* preference, or as one of the choice options. Consequently, there is always a first person to extra-personal (or *vice versa*) transformation required by our task. On a related note, our results parallel effects related to heading representations

during spatial navigation that are tied to local reference frames in RSc and SPL (Marchette et al., 2014). Notably, previous work has implicated the RSc in updating first-person information when progressing through an unfolding event (Summerfield et al., 2009; Vann et al., 2009). Given the RSc's known contribution to spatial navigation, autobiographical memory, and theory of mind (Spreng et al., 2009; Vann et al., 2009), our results highlight a role for the RSc in maintaining precise memories for different personal attributes during social decision-making.

### *Entorhinal-subicular representations of social knowledge along multiple scales*

We observed entorhinal/subicular responses that were sensitive to larger relative distances when comparing the stranger's rating with the *canonical individual* and the *self (friend anchor)*, versus comparing the stranger to the *self* and *friend (canonical anchor)*. One potential contributor is the greater ratings overlap for the *self* and *friend*, which induced smaller relative distances when making comparisons with the stranger's rating. Despite this difference potentially influencing our entorhinal/subicular choice discrimination results, similar effects were not observed in the hippocampal body, which primarily related to mnemonically demanding decisions. Notably, in rodents, hippocampal place representations are typically self-referenced and continuous. In contrast, entorhinal/subicular grid and boundary vector representations use multiple reference frames, and can either be continuously or discretely coded (Hartley et al., 2013; Diehl et al., 2017). Our data potentially extend this functional dissociation to memory-guided decision-making, where the body of the hippocampus generally related to greater sampling of personal experience (Bornstein & Norman, 2017; Bakkour et al., 2018; Shadlen & Shohamy, 2016), while the entorhinal/subicular region related to discriminating the relative distance between the stranger's rating and choice options differently, depending on the familiar individuals being compared (i.e., reference frame). Further support for this distinction stems from the putative role of the entorhinal cortex (Buzsaki & Moser, 2013; Maas et al., 2015; Navarro et al., 2015) and subicular region (Dalton & Maguire, 2017) in integrating incoming sensory input with learned hippocampal representations. Building on recent studies of macaque and human entorhinal cortex in spatial and non-spatial tasks, our results potentially speak to the organizational principles governing metrically coded maps (Doeller et al., 2010; Chadwick et al.,

2015; Constantinescu et al., 2016; Lositsky et al., 2016; Vass et al., 2017; Julian et al., 2018; Meister et al., 2018; Nau et al., 2018), along with the implicit encoding of discrete graphs (Garvert et al., 2017), which may be located within the same brain region. This architecture portends a domain general role for entorhinal cortex in memory-guided decision-making involving multiple reference frames.

### *Striatal involvement in knowledge-guided social decision-making*

We observed a ventral striatal signal corresponding to correct choices, which also related to choices where highest/lowest—as opposed to proximity—judgments could be made. This distinction between hippocampal-entorhinal versus striatal responses to different knowledge-guided decision-making strategies closely follows spatial navigation strategies; where the subjects align themselves to a single landmark, instead of a boundary, in order to infer the direction of a goal location (Doeller & Burgess, 2008). The implicit functional neuroanatomy may parallel the dissociation between hippocampal and striatal responses observed in the current study. This follows since boundary-oriented navigation is linked to the hippocampus, whereas landmark-based navigation is related to the striatum (Doeller et al., 2008). More generally, this striatal versus hippocampal dissociation highlights a potential mechanism for how reinforcement/procedural learning could reduce metrically-coded relational knowledge into more efficiently usable heuristics (Simon & Daw, 2011; Gershman & Daw, 2017), paralleling categorical versus coordinate-based judgments in spatial cognition (Kosslyn, 1987; Baumann & Mattingley, 2014).

### *Allocentricity and the dimensionality of preferences*

Our findings highlight the role of the hippocampal formation, namely the entorhinal cortex and subiculum, in the functional anatomy of how we transform egocentric and allocentric reference frames during decision-making. However, since all of our conditions involve a transformation, it is unclear whether humans ever preferentially use an extra-personal or ‘allocentric’ frame of reference during decision-making without drawing upon a first-person reference. Even in spatial and episodic memory this question is difficult to resolve (see Filimon, 2015); since self-referencing helps us relate

past experience to our current environment (Vann et al., 2009). Useful clues about how we flexibly transform egocentric and allocentric reference frames come from the social psychology literature. A subset of social psychology has focused on how individualized representations of others' preferences are generated by anchoring to a known preference, and then adjusting accordingly, a phenomenon known as 'anchoring and adjustment' (Epley & Gilovich, 2001; Epley et al., 2004; Tamir & Mitchell, 2013). Given their similarities, jointly studying 'anchoring and adjustment' with boundary-oriented navigation may offer a promising avenue for translational research across species and levels of analysis.

We tested metrically-coded preferences along a single dimension (namely the likelihood of people preferring things), but the true dimensionality of personal preferences is less clear. It is plausible that we learn about others' preferences within a multi-dimensional trait space (Tamir & Thornton, 2018). Yet we do not know the number of dimensions that support this trait space (Tavares et al., 2015). Further work may help relate what we know about navigating physical space with the mental exploration of more abstract spaces (Kaplan et al., 2017a).

Our study induces a spatial (metric) strategy, preventing us from asking whether humans must use map-like coding in order to relate others' preferences to their own. Furthermore, decisions in our task involve an abstract one-dimensional spatial discrimination, where subjects determine the relative proximity of a stranger's rating to the non-anchor individuals. Future work can implicitly test discrimination of metrically coded decision variables in different reference frames more implicitly (i.e., without an explicit spatial element) to determine whether the spatial element is required for hippocampal-entorhinal involvement. Still, theoretical work has provided support for the use of a spatial strategy, where non-human primate research has investigated social cognition in terms of coordinate transforms (Chang et al., 2013). Social coordinate transformation experiments have captured how social variables are encoded in individual neurons – and how encoding changes over different computational stages during social decision-making (Chang, 2013), like current versus long-term frames of reference (Boorman et al., 2013). Extending such work to abstract knowledge can potentially help us understand how coordinate transformations can generally guide everyday decisions.

## *Conclusion*

Metric coding of decision variables informs decisions by providing coordinates and boundaries that can be translated between different frames of reference. We provide evidence that neural computations – that integrate relative and absolute coordinates during spatial navigation – also extend to relating others’ personal attributes to our own. Consequently, these data provide important clues about how hippocampal-entorhinal map-like coding may facilitate memory-guided decision-making in a domain general manner.

## **Materials and Methods**

### *Subjects*

Twenty-four healthy adult subjects were studied and compensated (16 female; mean age in 25.5 y; SD of 5.38 y) and gave informed written consent to participate. This study was approved by the local research ethics committee at University College London. The study was conducted in accordance with Declaration of Helsinki protocols. All subjects were right-handed had normal or corrected-to-normal vision and reported good health with no prior history of neurological disease.

### *Task*

Stimuli were presented using the Cogent (<http://www.vislab.ucl.ac.uk/cogent.php>) toolbox running in MATLAB (Mathworks, Natick, MA, USA). Subjects performed a self-paced (max 9s), forced-choice, social decision making-task featuring 100 different personal preferences for 3 personally familiar individuals, which included themselves. Subjects viewed a personal preference for a novel stranger, represented by an avatar (Fig. 1), that was presented on a scale relative to one of three known individuals (*anchor*). Immediately prior to scanning, subjects were first trained to infer the stranger’s rating and then look below the scale in order to decide which of the two remaining familiar individuals was closer to that rating. All of the information needed to perform the task was presented at once on the screen and the only aspects of the task that varied from trial to trial were the anchor and choice options.

All ratings were obtained from subjects approximately 45 minutes prior to fMRI scanning. Subjects gave likelihood ratings about 110 everyday scenarios (e.g. eating spicy food; see Table S1 for all scenarios) from 1-9 on a 0-10 scale (0 being impossible to 10 being extremely likely) for themselves and two other familiar individuals. For one of the familiar individuals (*friend*), subjects were asked to choose their closest friend of the same gender that had the most distinct personality from the subject and – unlike the subject – ideally lived outside of London, in order to induce different ratings between individuals. The third individual was a typical/*canonical* individual of the same gender named Mr/Ms. Normal. Mr/Ms. Normal was supposed to correspond with what they thought a normal person/exemplar at their stage of life would do for each scenario (Fig. 1A). During the rating period, subjects were instructed to keep track of how their ratings related to each other for a given scenario (e.g., the participant being a bit more likely to have spicy food than their friend). Crucially, subjects also gave confidence ratings for the *canonical* exemplar and their *friend* on every scenario on a scale from 1-5 (No Idea to Very Sure) in order to assess the consistency of their rating (Fig. 1A). As a second confirmation of rating consistency, subjects also rated the familiar individuals for every scenario after the fMRI task.

Subjects then performed a brief practice version of the fMRI social decision-making task outside of the scanner, where subjects viewed a personal preference for a *stranger* that was presented on a 0-10 scale relative to one of the known individuals (*anchor*). Subjects had a maximum of 9 seconds to decide whether the *stranger's* rating was closer to the two remaining (*non-anchor*) individuals. The *anchor* was indicated by presenting their initials, which would either be the friend's initials, ME (Self), or MN, which was an abbreviation for Mr/Ms. Normal (*Canonical*). All information related to the decision was presented at once, with the scenario being listed above the stranger and anchor individual on the scale (Fig. 1B). Directly below the scale “closer to self/friend/canonical individual (A) or self/friend/canonical individual (B)” was written. After making a decision, there was then a jittered intertrial interval (ITI) period (mean=2.43s; range=0.25-9s) where a white fixation point overlaid on a black background was presented.

Crucially, the anchor individual was always placed in the middle of the scale, so that subjects needed to use prior social knowledge in order to infer the *stranger's* absolute preference and form a

mental number line of the *non-anchor* preferences. In other words, subjects had to infer the stranger's true (absolute) position in relation to the anchor's rating for that scenario and the closest boundary (0 or 10). For example, if the stranger was  $\frac{3}{4}$  to the right of an anchor individual that was rated a 9, the participant would infer that the stranger's rating was  $\frac{3}{4}$  of the way between 9 and 10 (the right boundary of the scale). Consequently, the subject would indicate the stranger's rating would approximately be near 10 (actual rating=9.75). We assumed that subjects represented people's preferences on this number line, similar to a mental number line (Dehaene et al., 1993), thereby enabling them to choose the *individual* (A or B) that was closest to the *stranger*. If both individuals in the choice were the same or equidistant from the stranger, subjects were instructed that either answer was counted as correct. Same/equidistant trials were not used in either behavioral or fMRI analyses.

Subjects then practiced the task using the last 10 scenarios they rated, which were repeated three times each. On the first few practice trials, subjects verbally rehearsed transforming the anchor's rating from the middle of the screen to its actual (absolute) position and then inferring the stranger's rating with verbal feedback from the experimenter until they performed the transformation correctly.

Then during fMRI scanning, subjects performed the task for the remaining 100 scenarios in three different self-paced runs (once for each anchor individual) each lasting approximately 10 minutes (a maximum of 15 minutes). From a relative viewpoint, the stranger was either  $\pm 1/4$  or  $3/4$  of the way between the anchor individual and the boundary ratings of 0 and 10, which was presented equally by run and condition. Subjects were instructed that it was a different stranger on each trial, so the stranger didn't conserve any preferences. To further emphasize the lack of conserved preferences, the stranger avatar randomly alternated between five colors (red, green, yellow, magenta, and cyan). Once again after completing the fMRI task, subjects gave likelihood ratings for the friend and canonical individuals outside of the scanner.

### *Computational Model*

We used a trial-by-trial computational model of subjects' performance in order to quantify how subjects' individual ratings and confidence judgments determined subjects' choice discriminability during social decision-making (see Fig. 1C). Choice discriminability (or inverse

precision) was measured using Shannon entropy (Shannon, 1948). Shannon entropy ( $H$ ) was calculated by estimating the distribution of choice probabilities using a softmax function of individual ratings. This included the standard deviation ( $\sigma_i$ ) derived from subject's confidence ratings ( $C_i$ ) for each scenario ( $\sigma_i=1$  for all self ratings), where the constant term accommodated the assumed stability (reduced memory demands) of self ratings, as well as the dynamic range of confidence ratings.

$$\sigma_i = 1 + \frac{1}{2}(5 - C_i) : C_i \in \{1, \dots, 5\}$$

Equation 4

Furthermore, the inverse temperature parameter ( $\beta$ ) of the softmax function was optimized using each subject's trial-by-trial performance. We constructed predictions of neuronal responses in terms of the Shannon entropy ( $H$ ) of the choice probability in each trial for fMRI (Kaplan et al., 2017b).

We optimized subject-specific parameters across trials using maximum likelihood estimation and the optimization toolbox in MATLAB (MathWorks, Inc). We then calculated trial-by-trial parameter estimates of  $H$  (choice discriminability) using the group average softmax (precision or inverse temperature) parameter. The ensuing trial-by-trial choice entropy measures were then used to predict BOLD responses in our neuroimaging analyses.

### *fMRI Acquisition*

Functional images were acquired on a 3T Siemens Trio scanner. BOLD T2\*-weighted functional images were acquired using a gradient-echo EPI pulse sequence acquired obliquely at 45° with the following parameters: repetition time, 3,360 ms; echo time, 30 ms; slice thickness, 2 mm; inter-slice gap, 1 mm; in-plane resolution, 3 × 3 mm; field of view, 64 × 72 mm<sup>2</sup>; 48 slices per volume. A field-map using a double echo FLASH sequence was recorded for distortion correction of the acquired EPI (Weiskopf et al., 2006). After the functional scans, a T1-weighted 3-D MDEFT structural image (1 mm<sup>3</sup>) was acquired to co-register and display the functional data.

### *fMRI Analysis*



Functional images were processed and analyzed using SPM12 ([www.fil.ion.ucl.ac.uk/spm](http://www.fil.ion.ucl.ac.uk/spm)). The first five volumes were discarded to allow for T1 equilibration. Standard preprocessing included bias correction for within-volume signal intensity differences, correction for differences in slice acquisition timing, realignment/unwarping to correct for inter-scan movement, and normalization of the images to an EPI template (specific to our sequence and scanner) that was aligned to the T1 Montreal Neurological Institute (MNI) template. Finally, the normalized functional images were spatially smoothed with an isotropic 8-mm full-width half maximum Gaussian kernel. For the model described below, all regressors, with the exception of six movement parameters of no interest, were convolved with the SPM hemodynamic response function. Data were also high-pass filtered (cut-off period = 128 s). Statistical analyses were performed using a univariate GLM with an event-related experimental design.

GLM1: There were two periods of interest, the self-paced (9s maximum) social decision-making and jittered baseline inter-trial interval (ITI), which were modeled as boxcar functions and convolved with a canonical hemodynamic response function (HRF) to create regressors of interest. For each social decision-making regressor (*self*, *friend*, and *canonical* anchor trials), there were parametric regressors based on choice entropy (discriminability), anchor rescaling, and accuracy (whether the choice was correctly answered; 1=incorrect choice; 2=correct choice). Inferences about these effects were based upon t- and F-tests using the standard summary statistic approach for second level random effects analysis.

GLM2: In a control analysis to determine whether choice discriminability effects were explained by covariance in mnemonic demands, we included additional parametric regressors based on the sum of subjective confidence ratings for the non-anchor individuals ( $\sigma_i$ ) in our signal detection model; Equation 1) and RT. In place of choice entropy to measure choice discriminability, we used the absolute value of the relative distance between the strangers' ratings and the non-anchor individuals for each trial. Anchor rescaling and accuracy parametric regressors were also still included in this model. GLM2 used the same conditions and periods of interest as GLM1.

GLM3: In a follow-up analysis to assay whether the ventral striatum accuracy effect was driven by choices where highest/lowest judgment between the two choice individuals needed to be made, we split decision-making trials in terms of whether a choice required determining which choice option was highest/lowest, instead of closest. Ignoring the three conditions, we split the data into two conditions; one where the target was between the two choice individuals, and the second where the target was either higher or lower than both choice individuals.

All initial analyses were whole-brain analyses. Subsequently, post hoc statistical analyses were conducted using 10-mm radius spheres in MarsBar toolbox (Brett et al., 2002) within SPM12 around the respective peak voxel specified in the GLM analysis. This allowed us to compare the effects of different parametric regressors of interest (e.g., to determine whether an accuracy effect was present in a region defined by an orthogonal main effect of choice discriminability). This ensured we did not make any biased inferences in our post hoc analyses. We used 10-mm radius spheres for post hoc statistical analyses—instead of clusters—since our hippocampal-entorhinal effects were corrected for multiple comparisons at the peak-voxel level.

Given the previously hypothesized role of the entorhinal/subicular region and hippocampus in absolute versus relative coding of environmental cues, we report whether peak-voxels in these regions survive small-volume correction for multiple comparisons ( $p < 0.05$ ) based on bilateral ROIs in the entorhinal/subicular region (mask used in Chadwick et al., 2015; Garvert et al., 2017) and the hippocampus (mask created using Neurosynth, Yarkoni et al., 2011). For all analyses outside of the ROIs, we report activations surviving an uncorrected statistical threshold of  $p = 0.001$  and correction for multiple comparisons at the whole-brain level (FWE  $p < 0.05$ ), unless otherwise noted. Coordinates of brain regions are reported in MNI space.

## **Acknowledgments**

The authors would like to thank Asaf Gilboa, Jochen Michely, Philipp Schwartenbeck, and Geert-Jan Will for helpful discussion, along with Martin Chadwick and Mona Garvert for providing an entorhinal/subicular mask. We thank Jacob Bellmund and Joshua Julian for helpful comments on a previous version of this manuscript. The authors would also like to thank Megan Creasey and Clive

Negus for help with scanning and the Wellcome Centre for Human Neuroimaging for providing facilities. This work was funded by a Wellcome Trust grant to KJF (Ref: 088130/Z/09/Z), and a Sir Henry Wellcome Postdoctoral Fellowship awarded to RK (Ref: 101261/Z/13/Z). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

### **Author Contributions**

R.K. conceived the experiment, collected and analyzed the data, and wrote the draft with input from K.J.F. K.J.F. also helped with the study design and model implementation.

### **Declaration of Interests**

The authors declare no competing interests.

### **References**

- Aronov, D., Nevers, R., and Tank, D.W. (2017). Mapping of a non-spatial dimension by the hippocampal-entorhinal circuit. *Nature* 543, 719-22.
- Bakkour, A., Zylberberg, A., Shadlen, M.N., and Shohamy, D. (2018). Value-based decisions involve sequential sampling from memory. *bioRxiv*, 269290.
- Baumann, O., Mattingley, J.B. (2014) Dissociable roles of the hippocampus and parietal cortex in processing of coordinate and categorical spatial information. *Front Hum Neurosci* 8, 73
- Boccaro, C.N., Sargolini, F., Thoresen, V.H., Solstad, T., Witter, M.P., Moser, E.I., and Moser, M.B. (2010). Grid cells in pre- and parasubiculum. *Nat Neurosci.* 13, 987-94.
- Boorman, E.D., Rushworth, M.F., Behrens, T.E. (2013) Ventromedial prefrontal and anterior cingulate cortex adopt choice and default reference frames during sequential multi-alternative choice. *J Neurosci*, 33:2242-53.
- Bornstein, A.M., Norman, K.A. (2017) Reinstated episodic context guides sampling-based decisions for reward. *Nat Neurosci*, 20:997-1003.
- Brett, M., Anton, J.L., Valabregue, R., and Poline, J.B. (2002). Region of interest analysis using an SPM toolbox; Sendai, Organization for Human Brain Mapping.
- Burgess, N., Jackson, A., Hartley, T., and O'Keefe, J. (2000). Predictions derived from modelling the hippocampal role in navigation. *Biological cybernetics* 83, 301-12.
- Buzsaki, G. and Moser, E.I. (2013). Memory, navigation and theta rhythm in the hippocampal-entorhinal system. *Nat Neurosci.* 16, 130-8.

- Chadwick, M.J., Jolly, A.E., Amos, D.P., Hassabis, D., and Spiers, H.J. (2015). A goal direction signal in the human entorhinal/subicular region. *Curr Biol.* 25, 87-92.
- Chang, S.W. (2013). Coordinate transformation approach to social interactions. *Front Neurosci.* 7, 147.
- Chang, S.W., Gariépy, J.F., and Platt, M.L. (2013). Neuronal reference frames for social decisions in primate frontal cortex. *Nat Neurosci.* 16, 243-50.
- Chen, J., Olsen, R.K., Preston, A.R., Glover, G.H., and Wagner, A.D. (2011). Associative retrieval processes in the human medial temporal lobe: hippocampal retrieval success and CA1 mismatch detection. *Learn Mem.* 18, 523-8.
- Constantinescu, A.O., O'Reilly, J.X., and Behrens, T.E.J. (2016). Organizing conceptual knowledge in humans with a gridlike code. *Science* 352, 1464-68.
- Dalton, M.A., and Maguire, E.A. (2017). The pre/parasubiculum: a hippocampal hub for scene-based cognition? *Curr Opin Behav Sci.* 17, 34-40.
- Daw, N.D., O'Doherty, J.P., Dayan, P., Seymour B., Dolan, R.J. (2006) Cortical substrates for exploratory decisions in humans. *Nature.* 441:876-9.
- Daw, N.D. in *Decision Making, Affect, and Learning: Attention and Performance XXIII*, eds:Delgado, M.R., Phelps, E.A., Robbins, T.W. (Oxford Univ Press, Oxford) (2011).
- Diehl, G.W., Hon, O.J., Leutgeb, S., and Leutgeb, J.K. (2017). Grid and Nongrid Cells in Medial Entorhinal Cortex Represent Spatial Location and Environmental Features with Complementary Coding Schemes. *Neuron* 94, 83-92.
- Dehaene, S., Bossini, S., and Giraux, P. (1993). The mental representation of parity and number magnitude. *Journal of Experimental Psychology: General.* 122, 371-96.
- Doeller, C.F. and Burgess, N. (2008). Distinct error-correcting and incidental learning of location relative to landmarks and boundaries. *Proc Natl Acad Sci USA.* 105, 5909-14
- Doeller, C.F., King, J.A., and Burgess, N. (2008). Parallel striatal and hippocampal systems for landmarks and boundaries in spatial memory. *Proc Natl Acad Sci USA.* 105, 5915-20
- Doeller, C.F., Barry, C., and Burgess, N. (2010). Evidence for grid cells in a human memory network. *Nature* 463, 657-61.
- Duncan, K., Ketz, N., Inati, S.J., and Davachi, L. (2012). Evidence for area CA1 as a match/mismatch detector: a high-resolution fMRI study of the human hippocampus. *Hippocampus* 22, 389-98.
- Eichenbaum, H. (2015). The Hippocampus as a Cognitive Map ... of Social Space. *Neuron*, 87, 9-11.
- Epley, N. and Gilovich, T. (2001). Putting adjustment back in the anchoring and adjustment heuristic: Differential processing of self-generated and experimenter-provided anchors. *Psychol Sci.* 12, 391-96.
- Epley, N., Keysar, B., Van Boven, L., and Gilovich, T. (2004). Perspective tasking as egocentric anchoring and adjustment. *J Pers Soc Psychol.* 87, 327-39.
- Filimon, F. (2015). Are all spatial reference frames egocentric? Reinterpreting evidence for allocentric, object-centered, or world-centered reference frames. *Front Hum Neurosci.* 9, 648.

- Garvert, M.M., Dolan, R.J., and Behrens, T.E. (2017). A map of abstract relational knowledge in the human hippocampal-entorhinal cortex. *Elife* 6.
- Gershman, S.J. and Daw, N.D. (2017). Reinforcement learning and episodic memory in humans and animals: an integrative framework. *Ann Rev Psychol.* 68, 101-128.
- Glascher, J., Daw, N., Dayan, P., O'Doherty, J.P. (2010) States versus rewards: dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron.* 66:585-95.
- Hafting, T., Fyhn, M., Molden, S., Moser, M.B. and Moser, E.I. (2005). Microstructure of a spatial map in the entorhinal cortex. *Nature* 436, 801-6.
- Hartley, T., Burgess, N., Lever, C., Cacucci, F. and O'Keefe, J. (2000). Modeling place fields in terms of the cortical inputs to the hippocampus. *Hippocampus* 10, 369-79
- Hartley, T., Lever, C., Burgess, N. and O'Keefe, J. (2013). Space in the brain: how the hippocampal formation supports spatial cognition. *Philos Trans R Soc Lond B.* 369, 20120510.
- Howard, L.R., Javadi, A.H., Yu, Y., Mill, R.D., Morrison, L.C., Knight, R., Loftus, M.M., Staskute, L., and Spiers, H.J. (2014). The hippocampus and entorhinal cortex encode the path and Euclidean distances to goals during navigation. *Curr Bio.* 24, 1331-40.
- Julian, J.B., Keinath, A.T., Frazzetta, G., and Epstein, R.A. (2018). Human entorhinal cortex represents visual space using a boundary-anchored grid. *Nat Neurosci.* 21, 191-94.
- Kaplan, R., Schuck, N.W., and Doeller, C.F. (2017a). The Role of Mental Maps in Decision-Making. *Trends Neurosci.* 40, 256-9.
- Kaplan, R., King, J., Koster, R., Penny, W.D., Burgess, N., and Friston, K.J. (2017b). The neural representation of prospective choice during spatial planning and decisions. *PLoS Biol.* 15, e1002588.
- Kosslyn, S.M. (1987) Seeing and imagining in the cerebral hemispheres: A computational approach. *Psychological Review*, 94:148-75.
- Kumaran, D. (2012). What representations and computations underpin the contribution of the hippocampus to generalization and inference? *Front Hum Neurosci.* 6, 157
- Kumaran, D., Melo, H.L., and Duzel, E. (2012). The emergence and representation of knowledge about social and non-social hierarchies. *Neuron* 134, 653-8.
- Kumaran, D., Banino, A., Blundell, C., Hassabis, D., and Dayan, P. (2016). Computations Underlying Social Hierarchy Learning: Distinct Neural Mechanisms for Updating and Representing Self-Relevant Information. *Neuron* 92, 1135-47.
- Kyle, C.T., Stokes, J.D., Lieberman, J.S., Hassan, A.S., and Ekstrom, A.D. (2015). Successful retrieval of competing spatial environments in humans involves hippocampal pattern separation mechanisms. *Elife* 4, e10499.
- Leutgeb, J.K., Leutgeb, S., Moser, M.B., and Moser, E.I. (2007). Pattern separation in the dentate gyrus and CA3 of the hippocampus. *Science* 315, 961-6.
- Lositsky O, Chen J, Toker D, Honey CJ, Shvartsman M, Poppenk JL, Hasson U, Norman KA (2016) Neural pattern change during encoding of a narrative predicts retrospective duration estimates. *Elife* 5.

- Maass, A., Berron, D., Libby, L.A., Ranganath C, Duzel, E (2015). Functional subregions of the human entorhinal cortex. *Elife* 4.
- Mack ML, Love BC, Preston AR (2016) Dynamic updating of hippocampal object representations reflects new conceptual knowledge. *Proc Natl Acad Sci USA*, 113:13203-13208.
- Mack ML, Love BC, Preston AR (2017) Building concepts one episode at a time: The hippocampus and concept formation. *Neurosci Lett*
- Marchette SA, Vass LK, Ryan J, Epstein RA (2014) Anchoring the neural compass: coding of local spatial reference frames in human medial parietal lobe. *Nat Neurosci*. 17:1598-606
- Marr, D. (1971). Simple memory: a theory for archicortex. *Philos Trans R Soc Lond B Biol Sci*. 262, 23–81
- McNaughton, B.L., Battaglia, F.P., Jensen, O., Moser, E.I., and Moser, M.B. (2006). Path integration and the neural basis of the ‘cognitive map’. *Nat Rev Neurosci*. 7, 663-78,
- Meister, M.L.R. and Buffalo, E.A. (2018). Neurons in primate entorhinal cortex represent gaze position in multiple spatial reference frames. *J Neurosci*. 38, 2430-41.
- Nau, M., Navarro Schroeder, T., Bellmund, J.L.S., and Doeller, C.F. (2018). Hexadirectional coding of visual space in human entorhinal cortex. *Nat Neurosci*. 21, 188-90.
- Navarro Schroeder, T., Haak, K.V., Zaragoza Jimenez, N.I., Beckmann, C.F., and Doeller C.F. (2015). Functional topography of the human entorhinal cortex. *Elife* 4.
- O’Keefe, J. and Dostrovsky, J. (1971). The hippocampus as a spatial map: Preliminary evidence from unit activity in the freely-moving rat. *Brain Res*. 34, 171-5.
- O’Keefe J and Nadel L. *The Hippocampus as a Cognitive Map*. 1978:114–52 (Oxford Univ Press)
- O’Keefe, J., and Burgess, N. (1996). Geometric determinants of the place fields of hippocampal neurons. *Nature* 381, 425-28.
- Schiller, D., Eichenbaum, H., Buffalo, E.A., Davachi, L., Foster, D.J., Leutgeb, S., and Ranganath, C. (2015). Memory and Space: Towards an Understanding of the Cognitive Map. *J Neurosci*. 35, 13904-11.
- Shadlen, M.N. and Shohamy, D. (2016). Decision Making and Sequential Sampling from Memory. *Neuron* 90, 927-39.
- Shannon, C.E. (1948). A mathematical theory of communication. *The Bell System Technical Journal* 27, 379-423.
- Simon, D.A. and Daw, N.D. (2011). Neural correlates of forward planning in a spatial decision task in humans. *J Neurosci*. 31, 5526-39.
- Spreng, R.N., Mar, R.A. and Kim, A.S. (2009). The common neural basis of autobiographical memory, prospection, navigation, theory of mind, and the default mode: a quantitative meta-analysis. *J Cogn Neurosci*. 21, 489-510.

Stemers, B., Vicente-Grabovetsky, A., Barry, C., Smulders, P., Schroeder, T.N., Burgess, N., and Doeller, C.F. (2016). Hippocampal Attractor Dynamics Predict Memory-Based Decision Making. *Curr Biol.* *26*, 1750-7.

Stokes, J., Kyle, C., and Ekstrom, A.D. (2015). Complementary roles of human hippocampal subfields in differentiation and integration of spatial context. *J Cogn Neurosci.* *27*, 546-9.

Summerfield, J.J., Hassabis, D., and Maguire, E.A. (2009). Cortical midline involvement in autobiographical memory. *Neuroimage* *44*, 1188-1200.

Tamir, D.I. and Mitchell, J.P. (2013). Anchoring and adjustment during social inferences. *J Exp Psychol Gen.* *142*, 151-62.

Tamir, D.I. and Thornton, M.A. (2018). Modeling the Predictive Social Mind. *Trends Cogn Sci.* *22*, 201-12.

Tanner Jr., W.P. and Swets, J.A. (1954) A decision-making theory of visual detection. *Psychological Review*, *61*,401-9.

Taube, J.S., Muller, R.U., and Ranck, J.B. (1990). Head-Direction Cells Recorded from the Postsubiculum in Freely Moving Rats. I. Description and Quantitative Analysis. *J Neurosci.* *10*, 420-35.

Tavares, R.M., Mendelsohn, A., Grossman, Y., Williams, C.H., Shapiro, M., Trope, Y., and Schiller, D. (2015). A Map for Social Navigation in the Human Brain. *Neuron* *87*, 231-43.

Vann, S.D., Aggleton, J.P., and Maguire, E.A. (2009). What does the retrosplenial cortex do? *Nat Rev Neuro.* *10*, 792-802.

Vass, L.K. and Epstein, R.A. (2017). Common Neural Representations for Visually Guided Reorientation and Spatial Imagery. *Cereb Cortex.* *27*, 1457-71.

Weiskopf, N., Hutton, C., Josephs, O., and Deichmann, R. (2006). Optimal EPI parameters for reduction of susceptibility-induced BOLD sensitivity losses: A whole-brain analysis at 3 T and 1.5 T. *Neuroimage* *33*, 493–504.

Yarkoni, T., Poldrack, R.A., Nichols, T.E., Van Essen, D.C., and Wager, T.D. (2011). Large-scale automated synthesis of human functional neuroimaging data. *Nat Methods.* *8*, 665-70.

Yassa, M.A. and Stark, C.E. (2011). Pattern separation in the hippocampus. *Trends Neurosci.* *34*, 515-25.

Zeithamova, D., Maddox, W.T., Schnyer, D.M. (2008) Dissociable prototype learning systems: evidence from brain imaging and behavior. *J Neurosci.* *28*:13194-201.

## Supplemental Figures

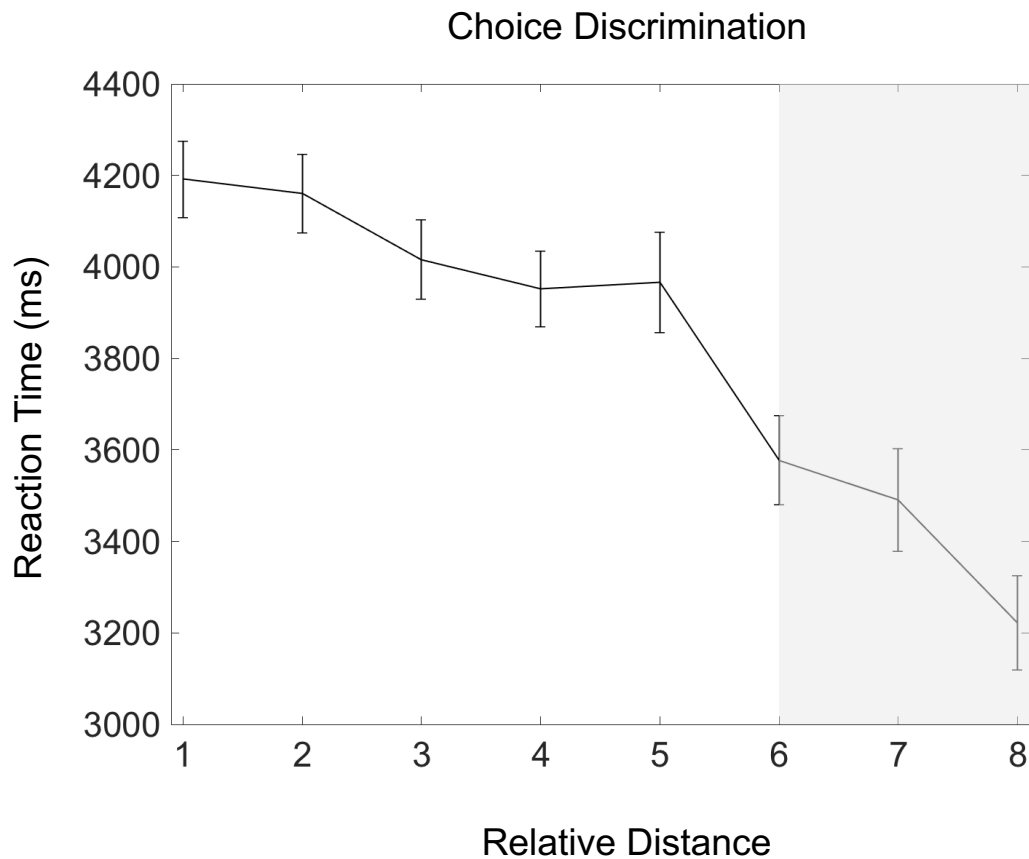
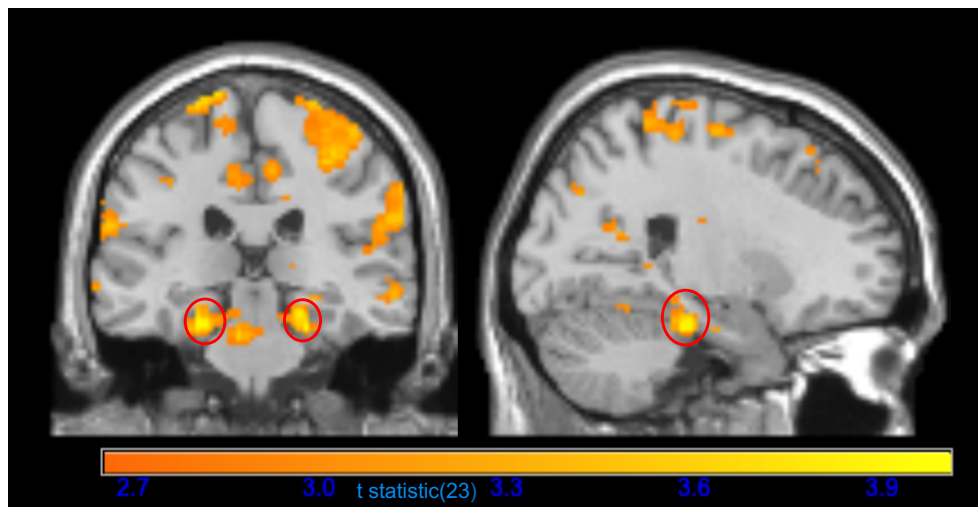


Figure S1. Relative distances and reaction time: Significant relationship ( $p < .001$ ) between decision speed (reaction time) and the relative distance between strangers' ratings and the non-anchor individuals for each trial. The absolute value of relative distances were rounded to the closest integer and plotted from 1 to 8. 7 & 8 on the x-axis are shaded in gray, because only 18/24 and 13/24 subjects had trials with relative distances of 7 & 8 respectively. Error bars showing mean  $\pm$  SEM.



A



B

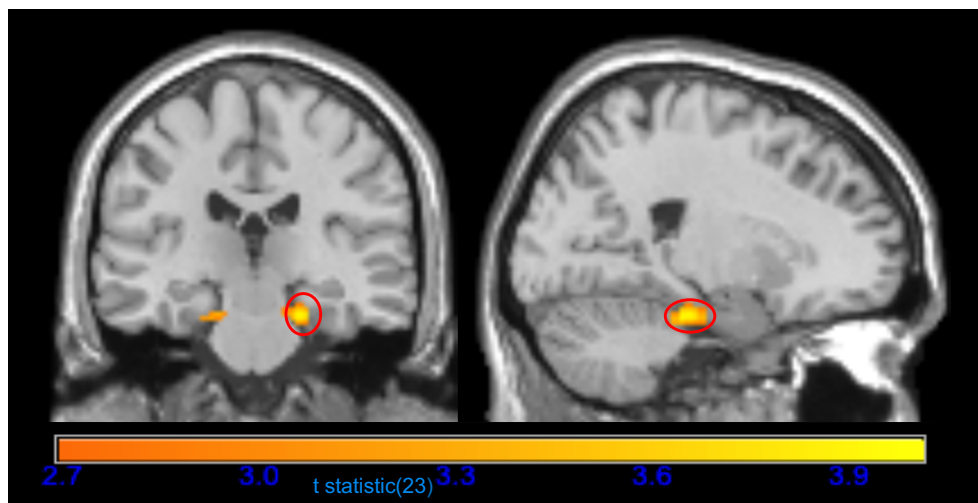


Figure S2. Anchor condition comparisons related to choice entropy and relative distance effects. A. Coronal and sagittal images showing bilateral entorhinal/subicular (circled in red) correlation with choice entropy for friend and canonical anchor trials (i.e., choices involving self-comparisons) versus self anchor trials (i.e., trials comparing canonical versus friend ratings). Figure 3C shows significant effect of condition. B. Coronal and sagittal images showing right entorhinal/subicular (circled in red) correlation with relative distance for self and friend anchor trials (i.e., choices involving comparisons with the canonical individual) versus canonical anchor trials (i.e., trials comparing self versus friend ratings). Figure 3D shows significant effect of condition. Highlighted regions survived peak-voxel FWE correction at  $p < 0.05$  and images are displayed at an uncorrected statistical threshold of  $p < .005$  for visualization purposes.

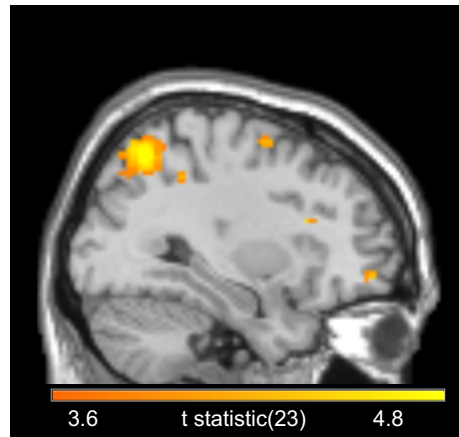


Figure S3. Superior Parietal Lobe Rescaling Effect. Sagittal image showing superior parietal lobule effect of mentally rescaling anchor towards the periphery in either direction, as seen in Figure 1D. Highlighted region survived cluster-level FWE correction at  $p < 0.05$  and image is displayed at an uncorrected statistical threshold of  $p < .001$ .

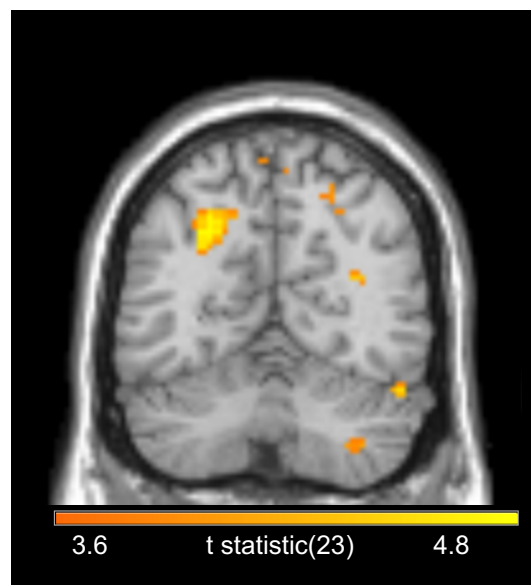


Figure S4. Coronal image of increased intraparietal sulcus activity for correct versus incorrect choices. Highlighted region survived cluster-level FWE correction at  $p < 0.05$  and image is displayed at an uncorrected statistical threshold of  $p < .001$ .

## Supplemental Table

1. Sending a message using mostly emojis 2. Eating spicy food 3. Following a TV series 4. Falling asleep in the car 5. Cycling 6. Meeting friends at the pub 7. Eating fast food 8. Downloading a movie 9. Speaking on the phone in a foreign language 10. Eating sushi 11. Wearing boots 12. Reading a novel 13. Having a house party 14. Growing vegetables 15. Taking an evening class 16. Going to a nightclub 17. Wearing a suit jacket 18. Sending a postcard when on holiday 19. Eating chocolate 20. Drinking coffee in the morning 21. Taking a dance class 22. Joining a political march 23. Volunteering 24. Jogging 5 miles 25. Smoking a cigarette 26. Needing caffeine to wake up in the morning 27. Catching a taxi 28. Travelling abroad for work 29. Listening to music on their way to work 30. Looking at new cars 31. Going to the drycleaners 32. Watching cartoons 33. Fixing things around the house 34. Reading a newspaper (online or print) 35. Having a meeting with a co-worker of a different nationality 36. Eating a kebab 37. Wearing a leather jacket 38. Hiring a person to clean their house 39. Hosting a dinner for friends 40. Visiting a children's museum 41. Sending a message via SMS 42. Dying hair 43. Painting their home 44. Leave a voicemail after a call 45. Drinking tea 46. Visiting an art museum 47. Communicate via email 48. Chat via instant messaging 49. Speak via video chatting 50. Having a conversation lasting more than an hour 51. Speaking with family 52. Ordering take away pizza 53. Cooking dinner 54. Making lunch 55. Going out to dinner with friends 56. Visiting a library 57. Making a new post on a blog 58. Sharing gossip 59. Joining a book club 60. Following politics 61. Sun tanning 62. Playing a board game 63. Watching a romantic movie 64. Holidaying in a tropical location 65. Going to a concert 66. Waking up early during the week 67. Attending a musical 68. Watching sport on TV 69. Playing videogames 70. Wearing a hat 71. Going on a road trip 72. Driving to work 73. Taking more than an hour to get ready in the morning 74. Wearing a sleeveless shirt 75. Walking to work 76. Going to a comedy show 77. Riding the bus 78. Spending a weekend afternoon on the sofa 79. Buying clothes from a secondhand shop 80. Buying designer clothes 81. Owns a pair of running shoes 82. Buying a new piece of furniture 83. Buying things from a moving sale 84. Decorating their place 85. Going to an antiques store 86. Hanging out with their neighbours 87. Going to a yoga class 88. Spending a holiday lying down on the beach 89. Going for a hike 90. Sleeping more than 7 hours a night 91. Playing football with friends 92. Meditating 93. Doing sudoku 94. Taking a nap 95. Staying up late during the week 96. Falling asleep in front of the television 97. Sleeping through their alarm 98. Sleeping on a flight 99. Taking the train 100. Going to the airport 101. Using a motorbike 102. Going to a work conference 103. Going out for lunch during the week 104. Speaking multiple languages at work 105. Working from home 106. Supervising other people 107. Changing jobs 108. Receiving a work-related phone call. 109. Checking work-related email at home 110. Taking a coffee break at work

Table S1. List of Rated Scenarios. Last 10 were used during practice trials and first 100 were used for the fMRI task.