**TASK**

Each subject performs three tasks, in the following order:
1. In scanner: 144 trials of the pleasurable image viewing task described below (divided into two blocks of 72 trials each.
2. Outside scanner: 48 trials of a fractal choice task.
3. Outside scanner: valence rating of all the images stimuli used in the task.

*Viewing task*. The basic structure of each trial is as depicted in the following figure:



Note:
- Each trial starts with a fixation cross of random duration, Unif(1,5) seconds
- Then a pie stimulus is shown depicting the probability $p$ of getting a rewarding image to view in the trial. With probability 1-p, the subject receives a "neutral stimulus", such as a solid frame of the same size and constant color.
- The subject then waits either 1, 2, 3, or 4 secs, with equal probability for the outcome phase. Screen indicates the waiting period with "..." in the middle of the screen.
- Finally the outcome (image or neutral) is revealed and shown for two seconds.
- A different rewarding image is shown in each trial (see more on Stimuli below). The same neutral stimulus is shown in non-reward trials.

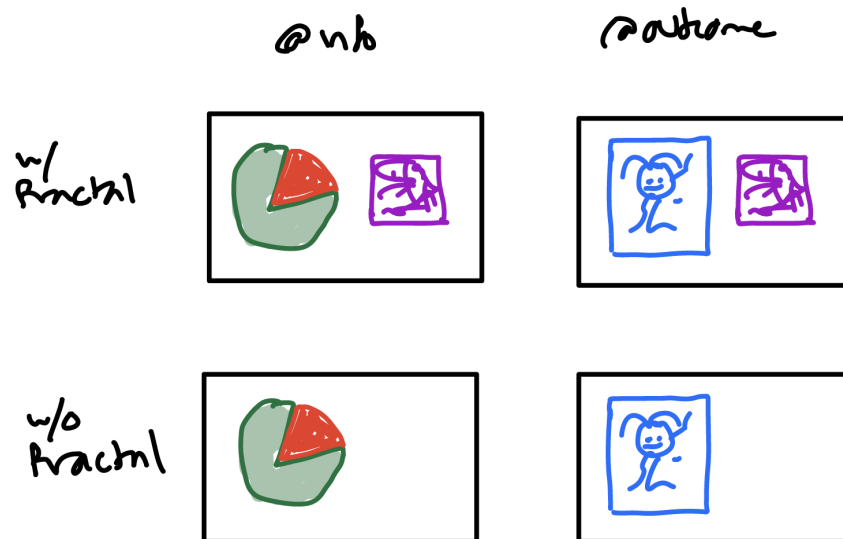The trial conditions are as follows:

| P = | 1 | $\frac{2}{3}$ | $\frac{1}{3}$ | $\emptyset$ |
|---|---|---|---|---|
| **Block H** | | | | |
| Freq | $\frac{1}{3}$ | $\frac{1}{3}$ | $\frac{1}{3}$ | — |
| # trials image | 24 | 16 | 8 | — |
| # trials $\emptyset$ | | 8 | 16 | — |
| **Block L** | | | | |
| Freq | — | $\frac{1}{3}$ | $\frac{1}{3}$ | $\frac{1}{3}$ |
| # trials image | — | 16 | 8 | — |
| # trials $\emptyset$ | — | 8 | 16 | 24 |

Note:
- There are two blocks of 72 trials each.
- The blocks differ on their distribution over the probability of getting a rewarding image.
- Trials are sampled by semi-randomizing from a list, to make sure that we get the trial counts in the table
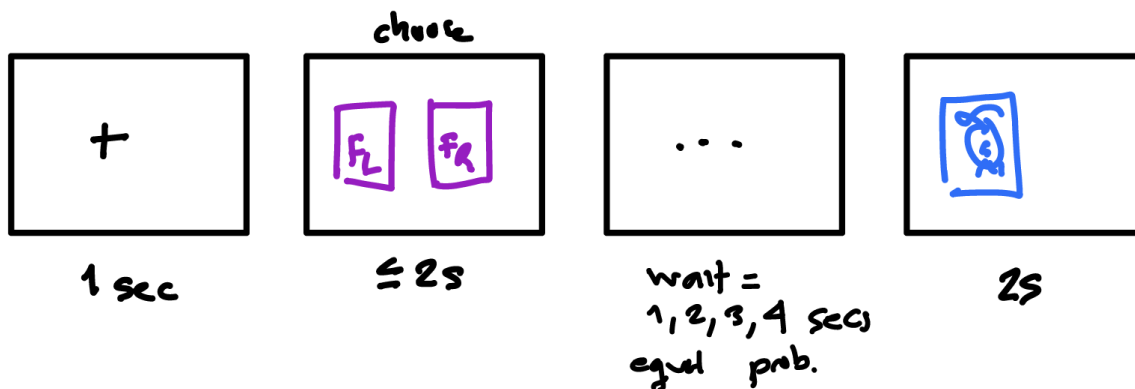
Fractals:
- In some of the screens, the subjects are also shown one of 7 different fractals while they see the stimuli.
- Let F(O,p) denote the fractal that is shown at the outcome screen next to the image, when the outcome is O and the probability of reward in the trial $p$.
- Let F(B,p) denote the fracthat that is shown at the info screen next to the pie chart, when the block is B and the probability of reward in the trial is $p$.
- Fractals shown at outcome screen: F(Rew,1), F(noRew,0), F(Rew, ⅓), F(noRew, ⅓)
- Fractals shown at info screen: F(H,1), F(H,⅓), F(L,⅓)
- The mapping of fractals to screens stays constant for the entire task, but is randomized for each subject.
- Figure below shows the arrangement of the screens (both for fractals and non-fractal trials).
- Left right location of the image randomized at outcome screens.
- To keep subjects engaged, they need to respond the location of the image/box (left-right) at outcome presentation. RT and response recorded, but the screen still shown for two seconds. Image border added upon response with feedback.

|            | @ info | @ outcome |
|------------|--------|-----------|
| w/ fractal | | |
| w/o fractal | | |

Other comments about the task:
- Duration:
  - Average trial is 8.5 s
  - Block of 72 trials is 612 secs = 10.2 min
- Wait delays are semi-randomized so that there is the same number of delays for each of the types of trials listed in the table above.

*Fractal choice task*. Every trial the subject makes a choice between a pair of the fractals shown in the viewing task. The trials are structured as follows:
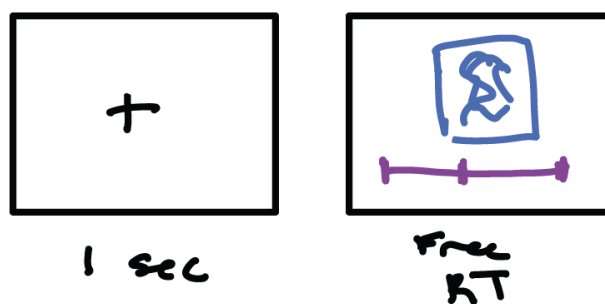


Note:

- Subjects have up to 2s to choose Left/Right, and choice and RT recorded.
- There are 48 trials.
- 18 of the trials involve a binary choice between these pairs, 3 times each:
  - F(Rew,⅓) vs F(Rew, 1)
  - F(noRew,0) vs F(noRew, ⅓)
  - F(H,1) vs F(H,⅓)
  - F(L,⅓) vs F(H,⅓)
  - F(L,⅓) vs F(noRew,0)
  - F(H,⅓) vs F(noRew,⅓)
- 30 of the trials involve a choice between the other 15 remaining pairs, 2 times each.
- Order of trials is randomized with the constraint that no pair is repeated within two trials.
- Left-right location of fractals also randomized.
- At outcome, the image associated with the chosen fractal is shown. For fractals associated with the info screen, the image is chosen with the probabilities associated with the info screen.
- Wait times are uniformly drawn 1,2,3,4 each trial.

Note:
- Mean trial duration is 7.5 secs
- Block of 48 trials takes ~ 360 sec = 6 min

*Image rating task*. Subjects are asked to provide a rating for each of the images shown in the previous two tasks. Depending on the subject choices, this involves between 72 & 72 + 48 = 120 images.

Each trial looks as follows:



Note:
- Images shown in random order.
- Subjects enter ratings by moving a slider bar up or down with 7 possible discrete locations (1-7) and pressing enter. Initial location of the slider randomized every trial.
- Rating and RT saved.

**STIMULUS SET**

We will use the images from Crockett et al, 2013, Neuron, "Restricting Temptations: Neural mechanisms of precommitment"), which was shown to activate reward regions like vmPFC upon viewing. This stimulus set was also used by Iigaya et al 2016 eLife and Iigaya et al 2020 Scientific Advances.

The stimulus set has 379 images of models in lingerie, but no nudity. Similar to the ads for Victoria Secret that are shown in newspapers and TV.

We will select the 120 images with highest ratings for use in the experiment, based on pre-existing rankings mentioned in the Iigaya paper. This same set of images will be used for all subjects, but in random order.

**SUBJECTS AND DATA COLLECTION**

Subjects will have to be heterosexual males, to minimize variance across subjects on the extent to which images are rewarding.

We will carry out data collection in three phases:
1. Prolific behavioral study. N = 50.
   >> Task is identical except that initial fixation cross has duration of 1 sec.
2. FMRI exploratory sample. N = 30
3. FMRI confirmatory sample. N = 30.

**SCANNING NOTES**

***TBC
Structural scans at the end (might not be needed if subject in Conte set)

**LOGIC OF THE EXPERIMENT**

The goal of the experiment is to look for evidence that experienced utility is affected by both consumption experiences (cEU) and good/bad news (sEU), as proposed by the Koszegi-Rabin model. Critically, sEU is hypothesized to be proportional to reward prediction errors (sEU = a * RPE).

The logic of our test relies on a combination of the BOLD responses in the viewing task and the choices in the fractal choice task.

The following table describes the EU signal at each screen of the viewing task, under the assumption that cEU(Reward) = 1 and cEU(noReward) = 0. This is without loss of generality, since these cEU can be renormalized:

**@info screen**

| | H,1 | H,$\frac{2}{3}$ | H,$\frac{1}{3}$ | L,$\frac{2}{3}$ | L,$\frac{1}{3}$ | L,0 |
|---|---|---|---|---|---|---|
| EU = cEU | 0 | 0 | 0 | 0 | 0 | 0 |
| EU = cEU + sEU | $\frac{1}{3}$ | 0 | $-\frac{1}{3}$ | $\frac{1}{3}$ | 0 | $-\frac{1}{3}$ |

**@outcome screen**

| | R,1 | R,$\frac{2}{3}$ | R,$\frac{1}{3}$ | nR,$\frac{2}{3}$ | nR,$\frac{1}{3}$ | nR,0 |
|---|---|---|---|---|---|---|
| EU = cEU | 1 | 1 | 1 | 0 | 0 | 0 |
| EU = cEU + sEU | 1 | $1\frac{1}{3}$ | $1\frac{2}{3}$ | $-\frac{2}{3}$ | $-\frac{1}{3}$ | 0 |

We will localize areas on vmPFC that encode EU with the contrast Rew > noRew at outcome, pooling all trials.

Then we will extract the signals of the responses in this area, separately at info and outcome screens, using a ROI analysis to test that it reflects both cEU and RPE, consistent with the table.

The hypothesis is that we will find evidence for cEU + sEU experience utility in both screens. (Note that we have previous data from a related old unpublished Bushong study that found clean evidence of this at outcome).

Problem: This is consistent with both the KR model of experience utility, and with this area containing units that respond to a mixture of cEU and RPE, but not to cEU + sEU, as KR hypothesizes.

The choice data from the fractals can then be used to test between the two models.

The fractals should acquire a valence proportional to the EU at each moment when they are shown. This leads to different predictions of their value, as shown in the table above, which should be reflected in the choices.
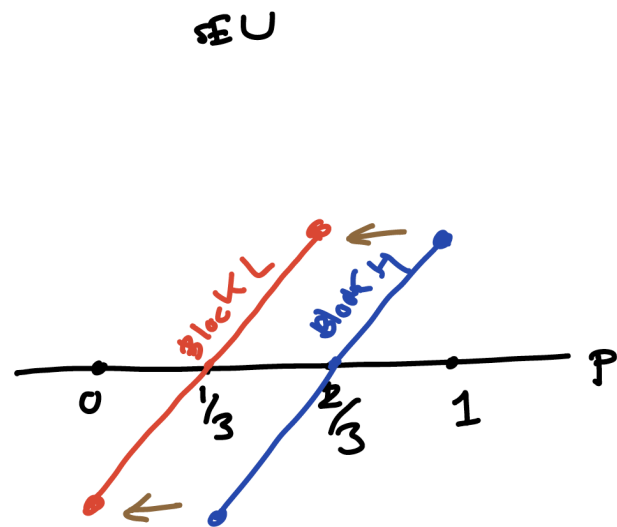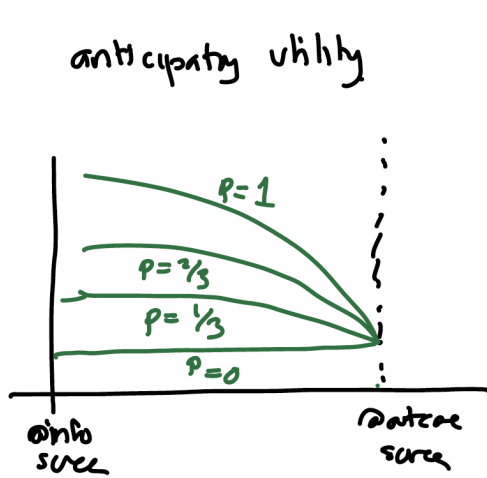
We will test this in two different ways:
- Direct predictions for specific binary pairs:
    - $F(Rew, \frac{1}{3}) > F(Rew, 1)$ under surprise utility, but not under cEU utility. This also provides evidence for positive utility of positive surprise.
    - $F(noRew, 0) > F(noRew, \frac{1}{3})$ under surprise utility, but not under cEU utility. This also provides evidence for negative utility of negative surprise.
    - $F(H, 1) > F(H, \frac{1}{3})$ under surprise utility, but = if under cEU.
    - $F(L, \frac{1}{3}) > F(H, \frac{1}{3})$ under surprise utility, but = under cEU. Also provides direct evidence that sEU changes with prior believes.
    - $F(L, \frac{1}{3}) = F(noRew, 0)$ & $F(H, \frac{1}{3}) = F(noRew, \frac{1}{3})$, which provides direct evidence of same sEU a the info and outcome choice screens.
- Estimation of the fractal values using logistic regression:
    - Using the choice data we can estimate a hierarchical logistic regression model of the fractal values.
    - We can then compare those with the values in the table above predicted by both theories.


***About anticipatory utility.*** The image viewing task has a delay between the information and outcome screens, which introduces the possibility that anticipatory utility might be at work. In fact, using a related task and the same stimulus set, igaya et al 2016 eLife and Iigaya et al 2020 Scientific Advances have argued for the existence of anticipatory utility in this task.

Although the main goal of our task is to test for surprise EU, the design will also allow us to look for anticipatory utility and separate it from any existing sEU signals.

Here is the logic of the test.

The figure below shows the predictions for the anticipatory utility and sEU signals at the info and wait screens. Critically, note that the sEU signals arise at the info screen, but are not sustained through the wait period, and have a different very specific predicted pattern across the two blocks. In contrast, the aEU signals are identical across both blocks (since they only depend on the probability of getting a rewarding image on the trial), are present over the entire wait period, and are decreasing for at least the last few seconds of wait.

anticipatory utility

$P=1$

$P=\frac{2}{3}$

$P=\frac{1}{3}$

$P=0$

@info
succ

@outcome
succ

$\mathscr{E}U$

Block L

Block H

0   $\frac{1}{3}$   $\frac{2}{3}$   1

$P$

Note also that the existence of the two blocks decorrelates significantly the anticipatory and RPE signals at the info screen.