

# Predictable Effects of Visual Salience in Experimental Decisions and Games

By XIAOMIN LI AND COLIN F. CAMERER\*

*Bottom-up stimulus-driven visual salience is largely automatic, effortless, and independent of a person’s “top-down” perceptual goals; it depends only on features of a visual stimulus. Algorithms have been carefully trained to predict stimulus-driven salience values for each pixel in any image. The economic question we address is whether these salience values help explain economic decisions. Our first experimental analysis shows that when people pick between sets of fruits that have artificially induced value, predicted salience (which is uncorrelated with value by design) leads to mistakes. Our second analysis uses evidence from games in which choices are locations in images. When players are trying to cooperatively match locations, predicted salience is highly correlated with the success of matching ( $r=.57$ ). In competitive hide-seeker location games, players choose salient locations more often than predicted by the unique Nash equilibrium. This tendency creates a disequilibrium “seeker’s advantage” (seekers win more often than predicted in equilibrium). The result can be explained by level- $k$  models in which predicted stimulus-driven salience influences level-0 choices and thereby influences overall perceptions, beliefs, and choices of higher-level players. The third analysis shows that there is an effect of visual salience in matrix games, but it is small and statistically weak. Applications to behavioral IO, price and tax salience, nudges and design, and visually-influenced beliefs are suggested.*

*JEL: D91 - Role and Effects of Psychological, Emotional, Social, and Cognitive Factors on Decision Making, C91 - Laboratory, Individual Behaviors, C72 - Noncooperative Games*

\* Tuesday 12<sup>th</sup> October, 2021, Li: Div HSS, Caltech, xli2@caltech.edu. Camerer (corresponding author): Div HSS and Computational and Neural Systems, Caltech, camerer@caltech.edu. Support was provided by the Behavioral and Neuroeconomics Discovery Fund (PI Camerer), Tianqiao and Chrissy Chen Center for Social and Decision Neuroscience, Alfred P. Sloan Foundation (G201811259), and NIMH Conte Center P50MH094258. Thanks to audiences at the Caltech Graduate Proseminar, Sloan/NOMIS Conference on Decision and Cognition, IAREP/SABE (Middlesex), Columbia University, Peking University, UC-Berkeley Behavioral Economics, Virtual Process Tracing Conference, Bocconi IGIER, and the Pitt Behavioral and Experimental Economics seminar. We especially thank Elke Weber, Vince Crawford, Richard Thaler (for the Schelling list idea), Adam Sanjurjo (for “The Pearl” tip), the Editor and several anonymous referees for helpful comments, Luca Polonio and Giorgio Coricelli for sharing the valuable gaze data in games, Anne Karing for a valuable image, Gidi Nave, Xintong Han, and Nina Solovnyeva (SURF 2020). Extra special thanks go to Eskil Forsell, Milica Moorman, Alec Smith whose energetic prior research on this topic laid the foundation for this paper, although their own work was never published due to poor project management by the senior author. We maintain a visual saliency processing

## I. Introduction

Features of a stimulus that grab attention are called “salient”. Of the different types of externally triggered sensory salience, visual salience is the best understood and is clearly important given the amount of information that people process through the visual system. This investigation is about whether one type of visual salience can be predicted and can help explain choices in experimental economic decisions and games.

Many economists have studied attention and salience recently, as part of growth in the foundations of behavioral economics. Notable contributions include Salience Theory (Bordalo et al. (2012b, 2013a,b)), a related model of focusing (Kőszegi and Szeidl (2013), and theories of rational (Sims (2003, 2006); Caplin and Dean (2015); Caplin et al. (2019); Kőszegi and Matějka (2020); Caplin et al. (2020); Mackowiak et al. (2020)) and dynamic inattention (Schwartzstein (2014); Gagnon-Bartsch et al. (2018)) The SAM algorithm salience is different from these economic models in content and purpose. We defer the comparison of those models to the penultimate Section VII.

To begin with, there is an important distinction between “bottom-up” and “top-down” salience (e.g., Chun et al. 2011; Baluch and Itti 2011).<sup>1</sup>

Bottom-up salience is what the human visual system notices most quickly and automatically. Bottom-up salience is also called “stimulus-driven”– the term we will use from now on– because it depends only on the properties of a stimulus. Stimulus-driven properties can be further divided into low- and high-level fea-

toolbox including tutorials at <https://github.com/lixiaomin328/imageToolboxForEconomists.git>

<sup>1</sup>There is an ongoing debate in attention science about how sharp the bottom-up vs. top-down distinction is. Awh et al. (2012) gives the example of the history of selective attention to a feature, which seems to influence future attention. That influence is not purely stimulus-driven (because it depends on previous attentive behavior, not just the stimulus itself) nor is it accomplishing a goal. Another example is faces. Faces are considered to be bottom-up salient for humans but they also help achieve a variety of goals that are generally evolutionarily important (such as emotional communication, friend-foe detection, mate choice, and social learning). These goals might also be even more important in a particular domain, like decoding facial emotion while watching a dramatic movie. So a person watching a movie sees faces that have both automatic bottom-up salience, and additional top-down salience to achieve the goal of understanding the movie. In general, the two processes together can be thought of as a “family of filters” that have been adaptively shaped by forces ranging almost continuously from evolutionarily-conserved universal principles to others locally tuned by personal experience and valuation.

tures (Judd et al., 2009). “Low-level features” are independent of object identity, meaning, and categorization; they include intensity, orientation, color, and motion. “Higher-level features” combine low-level features to identify and categorize objects, and direct attention to objects that are familiar, semantically meaningful, and generally valued. Faces, people, and text are generally salient high-level features.

Many algorithms have been trained to predict stimulus-driven salience using large image sets and eyetracking data from people who are “freely gazing” at the images for 3-5 seconds. These algorithms produce “salience maps” which closely match the actual gaze patterns.

In contrast to stimulus-driven attention, top-down attention is directed to achieve specific goals. We will therefore refer to top-down attention as “goal-directed” attention.<sup>2</sup> Goal-directed attention includes “extra-retinal<sup>3</sup> information such as intrinsic expectations, knowledge and goals” (Baluch and Itti 2011).

To illustrate the distinction between stimulus-driven and goal-directed attention, consider the classic study by Yarbus (2013), done in 1967. He showed subjects a painting of people in a room. One group was told to freely gaze. Another group was told to “estimate the material circumstance” of the people in the painting. The third group was told to “estimate the ages of the people” in the painting. Eyetracking showed that each of the three groups looked at somewhat different parts of the images.<sup>4</sup> Their gaze differences were due to differences in goal-directed attention. However, there was also a substantial overlap in measured attention. For example, people in both the free gaze and the “estimate the ages” goal conditions looked at faces in a similar way. This overlap indicates that

<sup>2</sup>Stimulus-driven and goal-directed attention are also sometimes called “exogeneous” and “endogeneous” in attention psychology. While we will not use this terminology, it is useful to emphasize the difference between stimulus-driven and rational (endogeneous) attention models, discussed later in Section VII.

<sup>3</sup>“Extra-retinal” means that the information attended to because of goal-directed guidance is not input to the retina, but is instead represented in the visual cortex and other regions such as the superior colliculus (see Veale et al. 2017); that information is in the proverbial “mind’s eye” rather than coming from the retina.

<sup>4</sup>Reversing the order of inference in Yarbus’s early study, Haji-Abolhassani and Clark (2014) showed that perceptual goals could be inferred reliably from eye-gaze patterns.

the measured attention to faces was both stimulus-driven and goal-directed when the goal was age estimation.

The hypothesis tested in this paper is whether stimulus-driven salience influences incentivized choices in three experiments involving decisions and strategic games. This type of salience lies outside of popular rational inattention modeling, which is a specific mathematical derivation of optimal goal-directed attention (discussed further in Section VII). The results, therefore, provide evidence that goal-directed models (including rational inattention) are leaving out an important type of attention—stimulus-driven salience—which is important behaviorally.

A preview of the first experiment illustrates the conflict between stimulus-driven salience and goal-directed perception. Subjects saw two sets of fruits, on the left and right halves of their computer screen. The two fruit sets were constructed to have different stimulus-driven salience, and different induced monetary value. The subjects' goal was to choose the set with the highest induced value, which requires goal-directed perception. Under time pressure, stimulus-driven salience sometimes shifted choices toward the high-salience options, even if those were low-value choices (see also Milosavljevic et al. 2012; Towal et al. 2013).

The other two experiments test whether stimulus-driven salience influences strategic choices that are intended to accomplish goals of (1) either coordinated matching, or hiding and seeking, in location games and (2) maximizing payoffs in normal-form matrix games. Predicted salience helps explain choices in the first set of location games and is weakly associated with low-level thinker choices (as classified by eyetracking) in the second set of normal-form games.

The empirical analysis uses an algorithm called the Salience Attentive Model (SAM).<sup>5</sup> SAM takes any 2-D color image as an input and predicts stimulus-driven attention—what most people will look at—in the first few seconds. The SAM algorithm is general, so it can be applied to any economic or social decisions influenced

<sup>5</sup>SAM is the first of several acronyms we use repeatedly. They are summarized in Appendix Table II.

by images. Potential applications include: advertisements; visual design features of “nudges”; televised political debates; e-commerce websites; virtual house tours; retail tags showing prices, promotions, or taxes; point-of-purchase displays; social media; face-to-face interviewing; and graphical display of information.

Here is the structure of the paper: The next section II presents the SAM algorithm. Section III describes the choice experiment pitting stimulus-driven salience against goal-directed attention. The results from location game experiments are described in Section IV and explained with cognitive hierarchy and level-k modeling in Section V. Section VI is about matrix games. Section VII describes several recent economic models of salience and attention and contrasts them with our approach. Section VIII concludes by speculating about other economic applications.

## II. The Salience Attentive Model (SAM) algorithm

Algorithms that take images as inputs, and output predictions about where people will look, have been an active area of research in visual neuroscience since the 1990s. A brief history will help clarify what the algorithms do (see Appendix A for more details).

The earliest algorithms included only low-level features (Itti et al., 1998). Using these features as a starting point was motivated by decades of research on the cognitive neuroscience of perception, including animal and human neuroanatomy, and detailed understanding of functions and interaction of different parts of human visual cortex.<sup>6</sup> We note these facts as an indication for readers of how much is known about basic aspects of the neural circuitry underlying attention and its connection to behavior, including the ability to causally change attention and

<sup>6</sup>Veale et al. (2017) is an excellent review. An elegant recent example found that stimulus-driven salient features are associated with measured neural activity in a specific area of the visual cortex called V1 (Chen et al., 2016). (V1 got that label because it is activated by retinal input earlier in time than other regions and detects only the simplest low-level features, such as orientation and direction. Krasovskaya and MacInnes (2019) review other examples of how well algorithmic salience is associated with measured neural activity in the visual cortex.) Other studies show that microstimulating and lesioning specific regions of the brain (in non-human animals) can *causally* change goal-directed attention and behavior (Baluch and Itti 2011).

subsequent behavior.

The early low-level algorithms were steadily improved by adding features that are higher-level, and generally salient, such as faces (Cerf et al., 2008). In the hunt for better predictive accuracy, in 2014 state-of-the-art algorithms switched to a neural network structure in which there is less *a priori* specification of what salient features are (Vig et al., 2014). These neural networks consist of multiple “layers” of connected discrete nodes. Each node in one layer receives weighted inputs from nodes in an earlier layer, and contributes weighted output as an input to nodes in a later layer. The initial input layer is based on a stimulus, and the final output layer encodes or “sees” an approximation of the stimulus. The network is “trained” by inputting stimuli— such as images— and propagating weighted inputs and outputs to eventually create a stimulus-specific output layer. That predicted output layer is then compared to the objective stimulus, and the connecting weights linking the different layer nodes are adjusted to improve accuracy. The SAM algorithm uses several modern variants of these methods to improve accuracy and training speed.<sup>7</sup> The network structure is usually “pre-trained” using a borrowed “backbone” network that encodes low-level features. The network is then trained further to learn encoding of semantically meaningful objects which are commonly present in the image sets and are looked at by the training subjects (such as apples, prices, people, and text; see Cornia et al. (2018)). The images in the SAM training sets were highly varied, and most subjects were students or others recruited at American campuses (see Appendix Table A1 for details).

<sup>7</sup>In technical jargon, SAM is a convolutional neural network with a salience encoder using a long short-term memory structure. Convolution is a method that combines encoding at different spatial scales. Crudely speaking, if features are encoded at fine-grained spatial scales and also at supersets of those fine-grained scales the object is “big”. The “long short-term memory” LSTM property is a kludge to retain memory so that backpropagation algorithms that adjust hidden-layer weights based on prediction errors do not overreact and create “vanishing gradients”— which are bad. SAM uses ResNet as its “backbone” (there is also a version with a VGG backbone). The backbone is the earliest part of the network (i.e, the layers closest to stimulus input, encoding low-level features). That part of the network typically has many layers and is therefore the most computationally demanding. It is used for low-level feature extraction from the input image. People nowadays mostly use established backbones such as ResNet or VGG, much like using a standard set of code then adding further code by hand.

These algorithms have progressed quickly because researchers can try out new ideas on four popular open-access salience datasets (SALICON, MIT1003, MIT300, CAT2000). These are sets of images along with “ground truth” data on what people actually looked at in the first five seconds of free gaze, recorded using eye-tracking and other high-quality methods for measuring visual attention.

SAM and similar algorithms are now highly accurate. The reported performance of SAM on the website MIT-saliency is 0.88 using the AUC-Judd area-under-the-curve measure (Riche et al., 2013). An AUC of .50 is random and 1.0 is perfectly accurate. The SAM accuracy of .88 is a little better than earlier algorithms and approaches the accuracy of the best human-to-human benchmark, which is .92.<sup>8</sup>

Figure 1 shows an example image and its associated SAM saliency maps. The salience map assigns a salience value from zero to one to each pixel of the image. The salience map is typically shown as a “heatmap” in grayscale or in color, with warmer (redder) colors indicating higher salience.<sup>9</sup> We adopted the default parameters from the original approach and applied them to our image dataset. There are no additional free parameters.<sup>10</sup>

To illustrate salience and choice, Figure IIa shows the map drawn by Schelling (1960) in a famous discussion of focality and “psychological prominence”. The map shows small square houses, a pond in the lower left, two places marked x and y, and a river running horizontally through the lower third of the map. A bridge spans the river. Schelling wrote:

Two people parachute unexpectedly into the area shown, each with a map and knowing the other has one, but neither knowing where

<sup>8</sup>The best human benchmark indicates how strongly two different sets of human fixation maps correlate for the same image. Each of the two sets contains many different individuals. Human-human accuracy is less than 1.0 because of idiosyncratic individual differences in their fixations, which make predictions from one group to another less than perfect (Judd et al., 2012).

<sup>9</sup>We use the standard color protocol “jet” in Matlab for all the heatmaps in this paper.

<sup>10</sup>Note that this CNN model, or any simpler variations of it, could be retrained on new data to understand different kinds of salience. Two studies have coded abstract features of strategies in 2-person matrix games (e.g., minimax, equal payoffs, level-1) and fit machine learning models using those features to explain observed choices. Hartford et al. (2016) is a neural network and Fudenberg and Liang (2019) is a random forest.

the other has dropped nor able to communicate directly. They must get together quickly to be rescued. Can they study their maps and “coordinate” their behavior? (p. 56)

Schelling said seven of the eight people (87.5%) who saw the map chose to rendezvous at the bridge.

In a larger incentivized experiment, N=61 UCLA students earned \$1 if they matched. They chose the bridge 59% of the time (see Figure IIb).<sup>11</sup> The SAM algorithm predicts that the bridge area, and the upper left road fork, are the most salient features (Figure IIc).

Note that SAM does *not* predict that the “x spot” is salient, even though it was chosen by 25% of the subjects. For stimulus-driven algorithms, “x” is a special configuration of low-level features—two lines with diagonal orientation, that meet symmetrically in the middle. The x is also a high-level feature because it is a letter in many languages; that is, it is a recognized semantic object. However, the algorithm, as it was trained on other images, was not originally capable of learning that “x” is familiarly known (to many UCLA subjects) to sometimes indicate locations of buried treasure on a map. So the x has minimal stimulus-driven salience and SAM did not learn its goal-directed value for coordinating a meeting place on a map.

To distinguish the effects of purely visual salience and goal-directed attention further, we did an online experiment in which the 10 most prominent map locations were described in a verbal list. There was no accompanying visual map. Just as in the map experiment, matching the list choices of others gave a reward. The subjects’ list choices were *not* the same as the map-based choices. The most popular choices were “x on the map” and “small house near the pond” (49% and 14%). Only 5% chose “bridge” (see Appendix E). This discrepancy shows that the popularity of the bridge choice depends on visual salience rather than

<sup>11</sup>These data were collected in conjunction with Milica Moormann and Alec Smith.



its semantic content.<sup>12</sup>

#### A. *Explainable AI and the SAM black box*

Before proceeding, we note that the SAM algorithm is neither a model nor a mechanism, in the sense that economists typically use those terms. Neural network models (including SAM) are often called “black boxes” because the basis of their predictions is in “hidden layers” that are difficult to interpret. One cannot readily do the comparative statics analysis that is useful in economics: e.g., there is no simple mathematical way to easily compute how a change in an input image leads to a change in the outputted salience map.

However, an active area called “explainable AI” is concerned precisely with how to make opaque AI output more understandable (Belle and Papantonis, 2020; Hinton et al., 2015; Ras et al., 2018; Arrieta et al., 2020; Fan et al., 2020; Lipton, 2018).<sup>13</sup> Some progress has already been made in explainability for deep neural networks predicting visual salience. For example He et al. (2019) used an image set in which a neural network predicts visual salience in a set of images. The categorical features in each image were also laboriously annotated ‘by hand’. That is, people looked at the images and coded the locations of vehicles, plants, animals, etc. Then salience, as encoded at the middle-layer output of the neural network, was extracted (like examining a partially-finished manufactured product). They found that the hand-coded categorized features were often correlated with the middle-layer salience predictions at these features’ locations. That correlation means that much of what the hidden middle layers were doing is learning the semantic categories of image features. In order of importance, 12 categories of features— a person’s head, “other”, an object, a person’s body part, etc<sup>14</sup>— were

<sup>12</sup>Rihn et al. (2019) finds a related effect, that visual attention to a logo rather than text description of a type of plant changes valuation.

<sup>13</sup>Igami (2020) explains the connection between some high-profile neural net training methods and structural estimation approaches invented in economics. This equivalence does not, however, guarantee the explainability of the *content* of the resulting neural networks.

<sup>14</sup>The rest of the list is: food, plant, symbol, vehicle, drink, animal head, text.

most commonly encoded by the middle network layers.

The method just described is one way to measure the “feature relevance” of a predicted salience map. Feature relevance could be applied to all the images in our investigation as well, to improve explainability. For a set of maps like Schelling’s, each spatial location has one or more codeable features— the distance from the center, roads, forks in roads, ponds, rivers, houses, bridges, etc (which were elements of the list version of the experiment). If these features and their locations are hand-coded, regressing the SAM salience values at each location against that location’s features will measure how well the SAM salience values are approximated by a function of the coded features. A good fit means the black-box salience output is approximated by explainable features. The size and statistical strength of the regression coefficients indicate which features are most salient.

The Schelling map example sets up the empirical question in this paper: How well does stimulus-driven salience— as predicted by SAM— predict actual choices in decisions and games? Does stimulus-driven salience get partially or entirely inhibited when there is also goal-directed attention?

We describe three experimental applications. They are:

- 1) Choices between visual images of two sets of fruits: The sets varied in induced values and in predicted salience. These data measure how often people picked lower-value sets because they were higher in stimulus-driven salience.
- 2) Strategic choices of locations in visual images: In Schelling-style matching games, both players were rewarded if they matched by choosing the same location. In hider-seeker games, the hider wanted to mismatch and the seeker wanted to match. These data measure whether cognitive hierarchy or level-k structural models can fit data, and more ambitiously, make accurate cross-game predictions from the hider-seeker game to the matching game.

- 3) Two-player 2x2 matrix games: These data measure whether stimulus-driven salience biases— which happen to predict looking at the top row and the left column in the matrices— can potentially explain strategy choices. This is a tough challenge for stimulus-driven theories because the experimental participants had a clear goal, to choose payoff-maximizing rows or columns. They may have ignored stimulus-driven salience entirely.

### III. Decisions: Fruit displays

#### A. Study 1: Salience and Induced Value in Visual Fruit Displays

The first experiment measured the empirical importance of visual salience in a simple setting that is lifelike. Subjects were shown two fruit sets presented on the left and right parts of an image, as shown in Figure IIIa. Each fruit type (e.g., apples or oranges) had a unique, pre-determined induced monetary value (Smith, 1976) that subjects learned before making choices. The induced values artificially created value, so that there is an objectively best choice, and we can clearly judge if people are making mistakes.<sup>15</sup>

N = 97 participants did this study on Prolific (a European online data collection platform), following a pre-registration process on the Open Science Foundation website (OSF)<sup>16</sup>. All the participants were pre-screened to have a prior approval rate of at least 70% based on their previous participation. Each subject was only allowed to participate in one experimental session (including pilot studies). Participation from mobile phones and tablets was not allowed in order to control for possible display effects.<sup>17</sup> There were five questions to check subjects' comprehension after the instruction session. We exclude individuals who failed more than one question. See a full description of the experiment block design in Appendix

<sup>15</sup>We also hope that the induced monetary values swamped minor differences in intrinsic subjective value from personal or aesthetic preferences for fruits.

<sup>16</sup><https://osf.io/>

<sup>17</sup>Even though computer screens also differ in size, phones and tablets have more variation in screen sizes.

F.F1 and Figure F1.

The total value of a fruit set is the simple sum of the values of all fruits in that set. The everyday analogue to this task is a retail vendor who is buying fruits at a wholesale market to resell, and has in mind a retail price for each fruit. The retail price of the fruit induces value to the vendor. Subjects learned the induced values of different fruits before the main session of 20 choices.<sup>18</sup>

While the vendor should optimally be computing resale value, the visually salient properties of fruit (such as color, intensity, and orientation), are hypothesized to influence stimulus-driven perception. The salience and value properties are independently controlled in the design.<sup>19</sup> In the choice sets, visual salience and fruit value were either positively or negatively correlated. The empirical question is whether subjects can ignore, or inhibit, visual salience, which is not generally correlated with induced value and could therefore lead to mistakes.

The main experiment included 20 images like those in Figure III. Choices were made with a 10s time limit. Trials were balanced across induced values, numbers of fruits in the two sets, and whether the more salient set was on the left or right (see Appendix F.F2). Subjects earned money based on the induced value of the sets they chose in an incentive-compatible design (a 10% chance of earning the value of what they chose on one randomly selected trial).

The average difference between the most salient peaks in the two fruit sets was 0.23 on the 0-1 scale of salience. More ambitious designs could obviously covary the size of the salience difference and the size of value difference between the two sets. In half of the trials, SAM-salience and induced value are “congruent” – one set is higher in both salience and induced value. In the other half of the trials, they are “incongruent” – the high-salience set has a lower induced value or vice versa.

<sup>18</sup>They experienced an untimed, but incentivized session before the main session. More experimental details are in Appendix F.F1.

<sup>19</sup>It is possible that stimulus-driven salience of fruits is correlated with their subjective value in the natural ecology– e.g., brightness might be visually salient, and also correlate with ripeness and fruit taste or nutrition. However, even if this is the case, by design stimulus-driven salience is uncorrelated with induced value, which is the only type of value a payoff-maximizing agent should attend to.

The dependent variable is 0-1 choice accuracy— did they choose the most highly-valued set? With a 10-second time limit, choice accuracies were 85% and 79% in the congruent and incongruent conditions. This drop in accuracy, when salience conflicts with valuation, is highly significant (p-value = 0.002, two-sided t-test).

We test for the effect of salience, controlling for the value gain from choosing correctly, using a logistic regression of the form:

$$(1) \ y_{ij} = \beta_1(S_j^H - S_j^L) + \beta_2 \text{abs}(V_j^L - V_j^R) + \beta_3(S_j^H - S_j^L) \text{abs}(V_j^L - V_j^R) + \beta_4 X_i + \epsilon_{ij}$$

with robust standard errors clustered at the subject level. The variable  $y_{ij}$  is accuracy (a 0-1 dummy variable, for person  $i$  at image  $j$ );  $V_j^L$  and  $V_j^R$  are the monetary values of the left and right sets in image  $j$ , and  $\text{abs}(V_j^L - V_j^R)$  is the absolute induced value difference ( $\text{abs}(\text{valueDiff})$  in Table 1). The congruency variable defined earlier is  $S_j^H - S_j^L$ , the difference in salience of the high- and low-valued sets. We are therefore regressing choice accuracy on congruency, absolute value difference, their interaction, and covariates.<sup>20</sup> The results are summarized in Table I. The induced value difference and congruency variables are both significantly associated with choice, with comparably large t-statistics (around 3-4).

There are two boundary conditions in which the effect of salience disappears. When the value difference is large the accuracy is 94% for both congruent and incongruent conditions (p=0.91 for the test for a difference). When the value difference is small, the accuracy is lower and salience-value incongruence does have an effect (78% vs. 69%, p = 0.01). (The Table I results pooled both types of images).

The second boundary condition is endogenous time allocation: When there is no time limit (N=22)<sup>21</sup>, participants in both conditions are near the ceiling of perfect accuracy (congruent 94% and incongruent 96%).

<sup>20</sup>“Covariates”: is yes when the current model contains covariates of education, gender, income, and self-reported fruit preference.

<sup>21</sup>An additional batch of subjects collected on Prolific did only the unlimited time experiment.

At this point, readers may be curious why subjects don't just ignore the stimulus-driven salience of the fruits. The reason is that, in economics jargon, perceptions are not freely disposable. The visual perceptual system is highly evolved to distill a huge amount of visual input into a much smaller amount of useful information, and to not waste the small amount that seems useful. The fastest parts of that process occur implicitly (without conscious awareness) in less than a second. Inhibiting any rapid highly-evolved implicit behavior is mentally difficult. One type of evidence about inhibition difficulty is that exogenous manipulation of attention— adding more ‘involuntary’ attention to a choice object— increases later choice of that object (albeit by a small amount; see Shimojo et al. (2003); Armel et al. (2008); Pachur et al. (2018) and see Mormann and Russo (2021) for a contradictory view).<sup>22</sup>

A mechanistic explanation for why irrelevant salience affects choices comes from a popular class of psychological models for how attention and decision time influence choice. These “accumulators” (or diffusion drift) models assume that over time perceptions and memory cumulate a running value of a latent numerical “evidence” variable (Ratcliff, 1978; Ratcliff et al., 2016; Fudenberg et al., 2018). A choice is made when the variable level crosses a mental threshold or barrier. In these models, if stimulus-driven initial perceptions enter the accumulator variable, there is not a known mechanism that will fully erase their effect. If the time to decision can be endogenously chosen by the decision-maker, then a very high threshold can be set which will dilute the early effect of stimulus-driven perceptions, but will not always fully inhibit that effect. (This is consistent with the absence of a salience effect in untimed trials.)

A different way to model why stimulus-driven perception influences choice comes from the signal-extraction model of Cunningham (2013). In that model, an “upstream” sensory system sends information to a more “informed” downstream

<sup>22</sup>A related phenomenon is called the “mere exposure” effect in psychology. Mere exposure means that repeated presentation of one unfamiliar stimulus tends to slightly increase expressed likings for that stimulus, compared to similar stimuli with less exposure (see Zajonc 1968 and Bornstein 1989 for meta-analytic review).

system and the two kinds of information are integrated. However, the downstream system only has partial information about input to the sensory system. Intuitively, the brain may partially accumulate the stimulus-driven perception into a decision variable as if it *might* have come from value-driven attention.<sup>23</sup>

#### IV. Study 2: Matching and Hider-Seeker Location Games

This section reports new experimental data from location games. Schelling’s map game is an example of a location game. In our general location games, two players saw a common visual image and simultaneously choose a location— a pixel. A circle was created around the pixels (with a radius of 108 pixels). The circle was about 1/5 of the screen width. The baseline circle size was chosen so that if players were choosing pixels randomly, they would match 7.1% of the time. (One experimental treatment below varied the circle size.)

In matching games, both players wanted to match by choosing locations that had overlapping circles. In hider-seeker games, seekers wanted to match and hiders wanted to mismatch. Interactions of the hider-seeker kind include predator-prey relations in nature. Human examples include choosing passwords to outwit hackers, other “coded” language and signals used in sports, gangs, and in other rivalries to coordinate action with teammates and avoid detection by the other side. Industries such as fashion can have follower-leader dynamics (e.g., fashion leaders want to “hide” by choosing unique new designs, and outsiders want to “seek” by matching those designs which induce hider-seeker structure). Visual salience might conceivably play a role in some of these games.

The experiment had three blocks of games (Figure G1): matching, the hider-seeker game in the role of seeker or hider, and the hider-seeker game in the opposite role of the one in the second block. The matching block always came first, followed by the hider and seeker blocks in a randomized order between

<sup>23</sup>This kind of upstream-downstream integration is likely to be common in the brain, leading to illusions like the “atmosphere” illusion (people do not fully undo the effects of unusual foggy or clear days on distance perception).

subjects. During each block, there was a “feedback” sequence in which the choice the other player made was revealed to a player right after both choices, by showing the circle around the other player’s pixel choice and the player’s own circle. In a “no feedback” sequence those results were not revealed.

The matching block had two sets of 20 images for each of the two feedback treatments (40 images in total). The hider-seeker game used a different set of 19 images for each of the two feedback treatments (38 images in total). For each image, subjects played once as a hider and once as a seeker. An additional short session of hider-seeker games followed in the last block (16 images) with a bonus payment 10x higher than in the baseline, to test for effects of higher incentives.

There was unlimited time to read instructions but only 6s to make a choice. Subjects got no payoff if they didn’t respond before the known time limit (see the instructions in Appendix G). The results shown to subjects in the feedback condition were drawn from previous choices of actual subjects (using different previous subjects for each image).

N=151 subjects participated, excluding a pilot dataset for power analysis. Of these 151 subjects, N=29 subjects (13 males, 16 females) participated in the lab, one at a time, in a small testing room where their eye movements were recorded. N=15 of those subjects were from the Caltech community and N=14 from the neighboring community (there were no differences in results between the two groups). The bonus payments were \$0.2, \$0.1, and \$0.4 in matching, hiding, and seeking games respectively, for each “win” per trial (image). Participants were paid the cumulative monetary amount at the end of the experiment. In the lab experiments, all the visual images were displayed on a computer screen in 1920x1080 resolution. The other (N=122) subjects participated online through Amazon Mechanical Turk (“MTurk”).<sup>24</sup> Images were randomly selected from a large image pool (273) with five categories (abstract art, city, faces, social scenes,

<sup>24</sup>Online experiments have the same instructions and block orders as the in-lab version, except that now everything is shown in a web browser. This study was pre-registered on the Open Science Framework (<https://osf.io/yuqjg/>) during data collection and before analysis. The sample size was pre-determined before the data collection process, based on a pilot study (N= 29) carried out in March 2017.



nature). The image set contained images with only one obvious salience center and more complex images which have multiple salience centers (Judd et al., 2009).

There were some behavioral differences between choices in the feedback and no-feedback conditions.<sup>25</sup> The largest effect is that the matching rate is higher with feedback than with no feedback (64% vs. 35%). However, the seeker win rate in hider-seeker games is the same in both conditions (9%) and most other differences are not substantial. We therefore report only data from the feedback condition in this main text. The corresponding no-feedback results are in an Appendix C.

Figure IV shows examples of result screens that subjects saw during the experiment.

#### A. Analysis and Results

Equilibrium game theory generates a statistical benchmark for what people might do.<sup>26</sup> In location games, strategies are pixels in x-y space (and their resulting circles).

For the matching coordination game, choices by the two players of any two pixels which create overlapping circles constitute a pure strategy Nash equilibrium. One image contains about two million ( $=1920 \times 1080$ ) pixels. Since any pixel match is a pure equilibrium, there are an enormous number of equilibria. There are also many mixed equilibria. So standard equilibrium theories do not rule out any of the location choices.<sup>27</sup>

<sup>25</sup>Both feedback and no-feedback blocks were included because each one answers a different question of interest. To help ensure increased subject comprehension in learning-by-doing, and especially in testing equilibrium concepts, the standard practice in experimental economics is to provide feedback. However, whether salience is predictive even with no feedback is an interesting question too. That is why we did both.

<sup>26</sup>A game-theoretic idea which might help explain how salience influences choices is “correlated equilibrium” (Aumann, 1974). When both players receive a common public signal and a strategy is conditioned on the signal values, a correlated equilibrium occurs when nobody wants to deviate from recommended strategies. Stop signs and green-yellow-red traffic lights, for example, act as correlating devices (also enforced by law) to create a commonly-observed visual signal which coordinates traffic and reduces accidents. In these terms, our study is about whether the stimulus-driven visual salience of image locations works as a “correlating device” in matching games.

<sup>27</sup>Note that if players have a personal utility from picking a specific location or a type of image feature, such preferences might conceivably reduce the set of equilibria, particularly if a selection principle such as payoff-dominance is applied (see Bacharach (1993); Bacharach and Bernasconi (1997)). However, such results would likely be sensitive to whether such preferences were commonly known.

For the hider-seeker game, there is a *unique* Nash equilibrium in which all locations are chosen equally often.<sup>28</sup> The fact that equal randomization over all strategies is the unique hider-seeker equilibrium is an example of how game theory logic conflicts with the result of human biology. We are so good at quickly noticing salient information, while amateurs at rapidly choosing what is *unsalient* in order to hide.<sup>29</sup> The last thing the brain is equipped to do is to *ignore* salient differences among many objects and choose them equally often.<sup>30</sup>

### B. Matching games

To analyze the behavioral data, we first test whether subjects are playing an equal random mixture across all pixels and their associated salience levels. To compare results from different images, all salience values in this section refer to the normalized levels, which are the rank percentiles of raw measures from the algorithm, ranked within each image. We calculated the normalized salience value for each chosen pixel and then compared these salience values against the baseline of equal randomization independent of salience. Kolmogorov-Smirnov tests reject the hypothesis of equal randomization for all treatment conditions ( $p < 10^{-4}$ ). Subjects' choices are not independent of salience.

To see examples of how salience affects choices, the choices from all the subjects

<sup>28</sup>For those unfamiliar with game theory, intuition can be gained by a simplified example. Suppose there are just two locations and the hider chooses them with probabilities  $p$  and  $1 - p$ . If the seeker matches those probabilities she has a  $p^2 + (1-p)^2$  chance of winning. This sum is always lower if the seeker chooses the most likely spot (i.e., the location with  $p > 0.5$ ) because if  $p > 0.5$ , then  $p > p^2 + (1-p)^2$ . To defend against this, the hider should mix equally, so  $p = 0.5$ . Every new location that is added should also have a  $\frac{1}{n}$  chance of being chosen (if there are  $n$  locations) by an iterated logic. A special design that, if a circle touches any boundary, it wraps around from the opposite boundary, guarantees the equilibrium.

<sup>29</sup>A similar conflict between logic and biology occurs in the games "rock, paper, scissors" (e.g., Crawford et al., 2013). When players display the three choices with their hands, there is a slight tendency to match an opponent's choice (e.g., playing rock against rock) more often than predicted in equilibrium. The explanation is that imitation of another person's body movements is such a highly-adapted automatic behavior, that the brain cannot inhibit the response, even though it reduces performance (e.g., you should play paper rather than imitating rock).

<sup>30</sup>The difficulty of inhibiting certain kinds of perception is illustrated by Steinbeck (2011). In "The Pearl" the protagonist Kino has hidden a valuable pearl that everyone in the small town covets. An unscrupulous doctor comes to treat Kino's baby, hoping to find out about the pearl. "The doctor shrugged, and his wet eyes never left Kino's eyes. He knew the pearl would be buried in the house, and he thought Kino might look toward the place where it was buried. "It would be a shame to have it stolen before you could sell it", the doctor said, and he saw Kino's eyes flick involuntarily to the floor near the side post of the brush house."

are plotted on two specific images in Figure V.

The salience heat map is in the middle column. The right column shows, using red dots, the subjects' actual location choices. The predicted salience in the middle column and the observed choice maps in the right column are highly overlapping, especially for the top panel image.

Statistically, the mean salience level of the pixel locations chosen in the coordination game is 0.87. This is far above the chance level of 0.5 ( $p < 10^{-4}$ ).

### C. How predictable is the matching rate across images?

Intuitively, the matching rate for an image should be affected by how dispersed salience is. When salience is highly concentrated, then the rate of choosing the same pixels, and matching, should increase. And if salience is not highly concentrated, the matching rate should be lower.

Dispersion of salience throughout an image can be measured by the number of local salience centers.<sup>31</sup> Figure VI shows that indeed, the matching rate<sup>32</sup> is strongly negatively correlated with the number of salience centers (Pearson  $r = -0.57$ ,  $p < 10^{-4}$ ,  $df = 38$ ).<sup>33</sup>

The matching rates span a range from a high rate of about 75%, for one salience center, to just above random (20%) for seven salience centers. These results suggest that for *any* image, the matching rate could be predicted *ex-ante* with substantial accuracy from the salience map, before any data are collected. Put the other way around, it is possible to find images with salience distributions

<sup>31</sup>The typical raw salience map has flat local maxima with many adjacent pixels with nearly equal salience. To detect salience centers we first apply a Gaussian smoothing (with [300pixel,300pixel] window size and standard deviation  $\sigma=75$  pixels) to the entire image to smooth hyperlocal spikes in salience. Then we simply take the number of local maxima for the salience distribution using the Matlab function `imregionalmax()` with default settings. That function takes the local maximum inside each 3pixel\*3pixel patch. If the original image has two local maxima that are close enough together, the Gaussian filter combines them.

<sup>32</sup>This result is based on all images from both the feedback session and the no-feedback session using the in-lab dataset (image N = 40).

<sup>33</sup>At a reader's suggestion, we also calculated whether the number of salience centers was correlated with the seeking win rate in hider-seeker games (across the N=38 images). This is an interesting question because if there are many strategically naive hidiers, the correlation will be positive. However, there is no correlation (Pearson  $r=-0.10$ ,  $p=0.23$ ).

that will predictably yield either near-perfect matching or near-random matching. This could be a useful tool for designers who are trying to either enhance shared attention or undermine it.

#### D. Hider-seeker games

For the hider-seeker game, we start with an example image and data. Figure VII shows that subjects' choices are more spread out than in the matching game examples (c shows hiding data and d shows seeking data.)

In Figure VIIc, there is no distinct peak of the hider choice distribution, and few choices are in the most salient area.

The direction of effects suggested by these two examples holds more generally. The mean salience levels of hider and seeker click points were 0.53 and 0.61, close to the chance level of 0.50.<sup>34</sup> The same in-lab group ( $N = 29$ ) with payoffs ten times higher had very similar results, averaging salience levels of 0.51 and 0.64 for hidere and seekers.<sup>35</sup> A paired t-test showed this difference in choice salience between hidere and seekers is highly significant ( $p < 10^{-4}$ ), reflecting what is suggested by the Figure VII example. The no-feedback results had a similar difference (see Appendix C).

**Seeker's advantage:** Recall that the theoretical frequency with which two randomly chosen location circles will match is 0.071. Table II presents the realized matching probability in each specific game condition.<sup>36</sup>

To check robustness, the hider-seeker game experiments were replicated in two other conditions: A high-payoff condition with payments 10 times as large ( $N=29$ )<sup>37</sup> and a between-subjects condition where subjects played only one of the

<sup>34</sup>p-value = 0.02, t-test CI: [0.51, 0.56]

<sup>35</sup>Hiding: p-value for test against null of .50 salience = 0.59, CI: [0.48,0.54], seeking: p-value  $< 10^{-4}$

<sup>36</sup>Tests to compare the matching rates with random baseline were carried out by bootstrapping a person's hiding data and a different person's seeking data (or two data points from matching game) for 1000 batches (batch size is a total number of different pairs). We get the empirical distribution for the matching rate and statistical significance against baseline 0.071 from that bootstrap. Specifically, each sample is drawn by matching two random users (different ones). The batch-seeking win rate is calculated accordingly. All values were calculated from the average of 500 iterations of randomly matching two data points from the data set if two subjects were in the same sub-block, same image.

<sup>37</sup>They did this session at the end of the in-lab group experimental session. See the full batch descrip-

two hider or seeker roles across all their trials ( $N=53$ )<sup>38</sup>. In both conditions the seeker win rate was 9.0%, the same as in the baseline experiment. All differences from the equilibrium prediction of 7.1% were highly significant.<sup>39</sup>

To test whether the seeker advantage is only present under time pressure,  $N=46$  people from MTurk participated in the same hider and seeker experiment, but without a time limit. The seeker’s win rate was again 9% (p-value = 0.002 for comparison with Nash benchmark 7.1%). Subjects spent on average of 3.14s, 4.61s, and 6.44s in matching, hiding, and seeking conditions, respectively, when there was no time limit.<sup>40</sup>

The seeker’s advantage could depend on the size of the circle that is drawn to surround the chosen pixel. To explore this possibility, in another experiment the circle size was enlarged to be 1.5x as large as in the original experiments. Then the chance/equilibrium matching rate is about twice as high, 16%. The seeker win rate was 18%, so there is still a small seeker’s advantage exactly equal in absolute size (+2%) to the benchmark circle results (p = 0.003,  $N = 66$ ).

The seeker’s advantage must be due to a correlation between the hidiers’ choices and the seekers’ choices, which should not happen in equilibrium (except for sampling error).<sup>41</sup> We have already shown that both hidiers and seekers choose slightly higher salience locations, but at different frequencies. But how exactly do those biases lead to the seeker’s advantage? Figure VIII presents the seeking win rates conditional on different salience levels for hidiers and seekers. The seeker’s advantage is mainly due to the concentration of wins when both players choose

tion in Table G1.

<sup>38</sup>This was an mTurk separate sample, see Appendix Table G2

<sup>39</sup>The 9% win-rate for seekers does not seem to be much larger than the equilibrium prediction of 7%. However, under the null hypothesis of Nash equilibrium, this win rate should be identically distributed for all images, and for all people. This null hypothesis supplies a lot of statistical power. A more conservative approach averages all data within an image and tests whether the image-wise matching rates are above 7% ( $N=19$ ,  $p= 0.0005$ ). A different conservative approach averages win rates for individuals and tests whether the average individual seeker win rate is different than the Nash 7% ( $N = 29$ ,  $p= 0.002$ ).

<sup>40</sup>The standard deviations were 7.10s, 15.54s, and 19.49s for matching, hiding, and seeking. These large standard deviations are not unusual for an online experiment with unlimited time because some subjects take much longer time than others.

<sup>41</sup>We know that people are capable of approximate equal randomization in these games because when they play a random computer opponent their choices are approximately equally random (Heinrich and Wolff, 2012).

locations that are in the top 10% in salience.

## V. A salience-influenced cognitive hierarchy model (SCH)

This section describes a parametric behavioral model meant to explain choices and their salience-sensitivity, closely following Crawford and Iriberri (2007a). It uses the level-k model of Stahl and Wilson (1994) and Nagel (1995), later extended by Camerer et al. (2004).

The SCH model combines cognitive hierarchy levels, a quantal response function (softmax) and a salience-influenced level 0 assumption.

### A. General model description

The population consists of different levels of players starting from level zero. The proportion of level  $k$  players is  $f(k)$ , with  $f(k)$  assumed to be Poisson distributed with parameter  $\tau$ .

For all levels of players, there is randomness which will be described using a conventional logit softmax function  $\frac{e^{\lambda x_n}}{\sum_m e^{\lambda x_m}}$  with parameter  $\lambda$ . Higher  $\lambda$  corresponds to more sensitivity to  $x_n$ .

In this SCH specification, the nonstrategic level zero players weakly prefer salient choices. The probability of choosing strategy/pixel  $n$  depends on the direct salience value <sup>42</sup>  $S_n$  of that pixel from SAM according to:

$$P_{0n} = \frac{e^{\lambda(1+\mu S_n)}}{\sum_m e^{\lambda(1+\mu S_m)}}$$

If  $\mu = 0$ , salience is ignored and level 0 types choose randomly among all points. We assume that  $\lambda$  and the salience weight  $\mu$  are common across subjects, although heterogeneous versions could be used (e.g., Rogers et al. 2009).

All levels of players above zero behave in the same way as in a standard cognitive

<sup>42</sup>Just as before, the salience values refer to the normalized ranking with respect to each image. This way, we can use data from different images and salience distributions in a common specification.

hierarchy model. Level  $k$  players assume that all other players are only of lower levels (0 to  $k - 1$ ), using normalized Poisson frequencies  $f(k)$ . A level  $k$  player calculates the expected payoffs of choosing  $n$ , denoted as  $EU_{kn}$ . The probability of a level- $k$  player  $i$  choosing option  $n$  is:

$$P_{kn} = \frac{e^{\lambda EU_{kn}}}{\sum_m e^{\lambda EU_{km}}}$$

Note that salience only enters *directly* into the value calculations of level 0 players. This assumption tests whether a model in which salience only enters  $k \geq 1$  level players through beliefs (and hence uses goal-directed attention) is a good approximation.<sup>43</sup>

### B. Model fitting results

Besides the SCH above, there are many other ways to specify models of limited strategic thinking, which have been mixed and matched in previous research. We therefore fit six model specifications to the hide-seeker data (see Appendix D).

Some specifications restrict the frequency of actual level 0 types to be zero,  $f(0) = 0$ , as if level 0 players are only a figment of the imagination of higher-level types (though see Wright and Leyton-Brown, 2019). Restricting  $f(0) = 0$  in this way clearly degrades fit (Table A2). We therefore focus only on  $f(0) > 0$ .

A close relative of SCH is the “Level- $k$ ” model, in which level  $k$  types believe all others are level  $k-1$  (rather than distributed from 0 to  $k-1$  as in SCH) (Crawford and Iriberri, 2007a,b). Level- $k$  is usually estimated non-parametrically, allowing all frequencies  $f(k)$  (up to some maximum  $k$ ) to be estimated separately.

Both SCH and Level- $k$  specifications with role-specific level frequencies fit the overall data about equally well by the AIC criterion (although SCH is a little better by BIC). These games are not an ideal testing ground for comparing such

<sup>43</sup>This is similar to Mehta et al., 1994a for matching games, in which “secondary salience” is derived from primary salience.

differences. The goal, instead, is to see if either SCH or level-k variants can explain both matching and hider-seeker games, which have different goal-directed attentional demands.

We first focus on the preferred specification of SCH. It has four free parameters— $\mu$ , the salience weight parameter;  $\lambda$ , the softmax parameter; and two role-specific parameters  $\tau_s$ , and  $\tau_h$  which are the Poisson distribution parameters of strategic levels for hidiers and seekers separately. (Allowing different  $\lambda$  and  $\mu$  parameters for hidiers and seekers fits worse due to the large BIC penalty for extra parameters).

We used a standard training-testing separation to avoid over-fitting. Recall that each subject did two sessions.<sup>44</sup> We use the first session data as a training set to estimate parameters. The parameter values are then fixed and used to predict data from the second session test set (see Appendix). The best fitting parameter values and measures of fit are shown in Table III.

Figure IX compares the actual choice density (frequency) function and best-fit model predicted density functions for the hider-seeker game. Training data are shown in the top Figures IXab and test data are shown on the bottom Figures IXcd. In the choice data, there is a sharp density increase starting around 0.9 salience for both roles (although note that the y-axes are different, so the actual increase is about half as big for hidiers as for seekers). There is also a smaller trend of slightly *decreasing* choice from the very lowest salience to medium salience levels for hidiers (but not for seekers). This small dip reflects the fact that some hidiers did manage to strategically choose the lowest-salience locations. SCH can roughly fit these two major features of the data.

However, the best-fit values of  $\tau$ , 0.4 and 0.1 for hidiers and seekers, are much lower than typical estimates around  $\tau = 1.5$  (e.g. Camerer, Ho and Chong, 2004; see also Riche et al., 2013, although Fudenberg and Liang (2019) find minimal prediction error in a large interval (0, 1.25) including low  $\tau$  values).

<sup>44</sup>Two sessions contain different image sets. A first session of normal payment trials including feedback and no-feedback trials and a second session of high payment trials.



The low values of  $\tau$  estimated for SCH result from the fact that the ability to identify  $\tau$  is limited in these visual choice games. A single-peaked SCH with Poisson  $f(k)$  does not meet the calibration challenge well. Level-1 hidiers should anticipate high-salience choices by level-0 seekers and move sharply to anti-salient locations. But there are not that many low-salience choices in the hider data (as Figure IXd shows). The SCH distribution explains the infrequency of low-salience hiding the only way it can, by simply estimating few level-1 types through a low value of  $\tau$ .

The Level-k model gives better insight here about plausible level frequencies.<sup>45</sup> Compared to SCH, the best Level-k specification estimates lower frequencies of level 0 ( $\hat{f}_s(0) = .17$  and  $\hat{f}_h(0) = .29$ ) for seekers and hidiers, and a higher salience weight  $\hat{\mu} = .18$  for level 0 types. Level-k also estimates larger frequencies of level 2 and 3 types ( $\hat{f}_s(3) = .66$ ,  $\hat{f}_h(2) = .61$ ). While the overall Level-k fit is just a little less accurate than SCH, this type distribution is more consistent with experimental results than the SCH estimates of low  $\tau$  (see Appendix D). So while it is clear that both specifications fit the salience-choice profiles adequately (as seen in all the Figures including Appendix Figure D1), they suggest different evidence of level frequencies. These games were chosen to investigate the effect of predictable salience, but were not ideal to recover levels accurately. Better methods can be developed.

### C. Cross-game predictive validation

To further test generalizability of SCH, parameters estimated from fitting the SCH model to hider-seeker data will now be used to predict choice behavior in the matching game. There is no guarantee that this cross-game portability will work at all (see Hargreaves Heap et al., 2014). Identification of the salience weight  $\mu$  in hider-seeker games comes purely from the level 0's choices and from higher-

<sup>45</sup>A better way to identify  $\tau$  is by creating games in which different level types choose distinct strategies (such as in the matrix games pioneered by Stahl II and Wilson, 1994, and see Nagel, 1995; Ho et al., 1998; Costa-Gomes and Crawford, 2006; Kneeland, 2015; Fragiadiakis et al., 2017).

level player beliefs and choices. In the matching games, all higher-level types are similarly guided by goal-directed attention since they are all trying to match the lower-level types. The strength of salience-sensitivity that is estimated in the two cases could easily be different. Furthermore, matching and hiding are completely opposite in strategic motives.

Figure X compares predictions of the salience-frequency profile on the test set of matching game data. The left graph shows predictions based on using hider-seeker training— that is, the free parameters are trained on the hider-seeker data, then fixed and used to predict (“test on”) the matching game results. The right graph shows predictions of matching test-set data using matching data for training (i.e., using the two-session train-test cross-validation described above). Of course, training on the matching data and then predicting matching test data should be more accurate than training on a different type of game, and it is (LL = -1943). However, training on the hider-seeker data and testing on matching is only about 10% worse (LL=-2176). Comparing Figures Xab shows that the main difference is that the hider-seeker trained parameters underestimate how sharply matching-game test data respond to the highest salience.<sup>46</sup>

The hider-seeker structure is a good example of how stimulus-driven and goal-directed salience can be combined. Level-0 players are only influenced by stimulus-driven salience (from the SAM algorithm) because they do not have a strategic goal. Higher-level types need to compute expected values of strategies, which requires goal-directed attention. But they also form beliefs about level-0’s which requires simulating the stimulus-driven attention of level-0’s. Therefore, both types of attention need to be combined to make good choices. The fact that hiders lose more often than expected in equilibrium is associated (via the structural model) with the fact that they are choosing too many locations that have stimulus-driven salience. Their goal of hiding, which should guide perception to low-

<sup>46</sup>We did not do the opposite analysis, predicting hider-seeker data based on parameters estimated from matching game. The meaning of doing this opposite analysis is limited due to the identification problem. Using matching game data only is not enough to identify the strategic level parameters because all level players are using similar strategies of choosing salient locations.

saliency locations, does not appear to sufficiently inhibit stimulus-driven saliency.

## VI. Study 3: Matrix games

Location game experiments are unusual. Most game theory experiments, following visual conventions in textbook game theory, use normal-form games in a matrix format (or occasionally game trees). To establish boundaries of where visual saliency is predictive and where it is not, it is therefore useful to ask whether SAM saliency can help explain choices in the common matrix game format.

First, note that the SAM training set does not contain images that resemble matrices of payoffs. Subjects in matrix game experiments also have a clear attentional goal, which is to look at numbers in a matrix to make a high-payoff choice. These goals are likely to create a complicated visual search to compute beliefs and implement decision rules, which is different than the rapid stimulus-driven attention that SAM is designed to predict.

In fact, many studies using Mouselab and eye-tracking stretching back three decades have shown patterns of search consistent with goal-directed perception for strategic thinking (Camerer et al., 1993; Costa-Gomes et al., 2001; Johnson et al., 2002; Arieli et al., 2011; Brocas et al., 2014; Polonio et al., 2015; Devetag et al., 2016). Furthermore, most of the behavioral studies about coordination and hide-seeker games have aimed at establishing general principles of focality or psychological prominence from strategic goals and set-theoretic properties of strategies (see Appendix B for a review). So it is already known that goal-directed allocation of attention is evident in choices from matrix payoff games. An unanswered new question is whether stimulus-driven SAM saliency has any *additional* predictive power or not.

The possible influence of visual saliency is tested here using data from Polonio et al. (2015). In their experiment, N=56 people played 32 normal form games with different strategic structures. Eye-tracking was used to record visual attention. These data are especially useful because actual gaze maps can then be compared

with both SAM predictions and with actual choices.<sup>47</sup>

Figure XIa shows one example of the type of matrix that subjects see on their computers (it’s a prisoner’s dilemma in structure). Row player payoffs are in the lower left of each matrix cell, and column player payoffs are in the upper right of each matrix cell.

Figure XIb is the average prediction from the SAM algorithm about where people look, averaged over all 32 games. There is a predicted bias toward looking more at the top row and the left column, as well as a row-player payoff bias (even for column players). Figure XIc is the average measured attention map calculated from eye-tracked gaze data over different types of games (filtering out gazes that are away from payoffs). The comparison between Figures XIb (algorithm) and c (gaze data) suggest that the algorithm does predict the actual attention allocation during game play rather well. This visual impression is supported by conventional statistics used in visual science.<sup>48</sup> Much to our surprise, the actual human gaze data are also quite similar for row and column players (as is the SAM salience map, because it does not use vary with player roles). This is surprising because higher-level strategic thinkers need to direct attention to different payoffs.

The main question is whether there is a congruency effect (as in the fruits experiment 1): That is, does salience affect how often people choose the equilibrium strategy? We look at the 24 games which contain a unique equilibrium strategy for both players. We also use Polonio et al. (2015)’s own classification of subjects into three groups based on strategic levels of thinking from 0 to 2, using the Cognitive Hierarchy model.<sup>49</sup>

Figure XIId shows that level 0 and 1 types do choose the salient strategy more

<sup>47</sup>See Appendix H for more details.

<sup>48</sup>In the computer vision field, two validation scores, AUC and CC are commonly used metrics to evaluate how closely salience algorithm predictions are correlated with actual human gazes. AUC: area under the receiver operating characteristics curve and CC: Pearson Correlation (see Kummerer et al. 2018). The Appendix Table H1 shows these statistics.

<sup>49</sup>They type-classify players based on their gaze patterns on matrix games. The level-0s only focus on the payoff property itself (intra-cell). Level-1 players compare their own payoffs (own focused). Level-2 players also look at others’ payoffs (distributed attention). This classification from gaze data was then correlated with predictions about what choices the three types should make.

often when it is an equilibrium, and level 2’s go slightly in the opposite direction. This is consistent with the idea that level 0’s are not using goal-directed attention, and level 1’s and 2’s use more goal-directed attention.<sup>50</sup>

Table 4 tests whether the likelihood of choosing the equilibrium strategy depends on salience Congruency. There is no general effect when all level types are pooled together (Model 1). However, Model (2) shows that there is a substantial effect of Congruency, but only for Level-0 players. (Note that Level-2 is the omitted level category so that the Congruency main effect estimates the Level-2 effect, which is negative). However, the significance of the Level-0 effect is only  $p=0.12$  when Bonferroni-corrected for multiple comparisons.

Thus, the evidence for an influence of stimulus-driven salience is suggestive but not statistically strong. It is also a surprise that the salience map and gaze data are so similar. Future experiments could explicitly manipulate salience (guided by SAM predictions) of particular payoffs to see if stronger effects can be created.

## VII. Comparison with other salience and attention approaches

This section briefly reviews recent economic theories which have analyzed salience and attention, and describes the relations of those theories to our approach.

### A. *Salience theory*

Salience Theory is a theory of salience that has been widely applied for the last 10 years in economics and finance, and in other areas (Bordalo et al., 2012b, 2013a,b). It was the first economic theory to specify exactly how salience is generally derived from attributes, and affects choice, in order to make clear predictions testable from observable data. The goal of this section is to describe how salience is computed in that theory and compare it to stimulus-driven SAM algorithmic salience.

<sup>50</sup>Note that we did not pre-register this prediction, so our conclusions should rightly be taken as exploratory and not a planned test of a hypothesis.

In Saliency Theory, attribute values of choice objects which are relatively farther from a reference point (such as the average attribute value<sup>51</sup>) are judged to be more salient. We’ll use the notation from analysis of multi-attribute choice (Bordalo et al., 2013b), to see how Saliency Theory works. A choice  $k$  has attribute level  $a_k$  along a particular attribute dimension. The average level across the entire choice set is  $\bar{a}$ .

The saliency function is defined by  $\sigma(a_k, \bar{a})$ . This function is assumed to obey two properties called ordering and diminishing sensitivity.

Ordering means that increasing the magnitude of the attribute level  $a_k$  by  $\epsilon$  from  $\bar{a}$ , while decreasing the reference point in the opposite direction by  $\epsilon'$ , increases saliency.<sup>52</sup> Kőszegi and Szeidl (2013) proposed a similar “focusing” model in which all values of an attribute are weighted more heavily when an attribute has more wide-ranging utilities (see Bordalo et al. 2013b (p. 815-16) for comparison). Diminishing sensitivity means that increasing the level of both  $a_k$  and  $\bar{a}$  by the same positive amount reduces the saliency of  $a_k$ . Although ordering and diminishing sensitivity are enough for most of the applications to work, a more strict version further assumes homogeneity of degree zero (i.e.,  $\sigma(a_k, \bar{a}) = \sigma(\alpha a_k, \alpha \bar{a})$  for  $\alpha > 0$ ). A simple saliency function which satisfies all these properties is  $\frac{|a_k - \bar{a}|}{|a_k| + |\bar{a}|}$ . We now make two remarks about Saliency Theory.

First, attributes—such as product quality, or endowment states—do not have to be numbers to be judged as salient. They could be perfume aromas or restaurant noise levels. However, attributes are assumed to have subjective estimated values, so that saliency can be computed and used to weight attributes in computing decision values. A salient thinker will overweight the salient attributes

<sup>51</sup>In some applications, it is plausible that an external reference point which is not part of a choice set influences saliency. For example, the explanation of endowment effects works with goods that have two attributes, and the consideration set includes having nothing (0,0) (Bordalo et al., 2012a). Including this null state makes the best quality of the initially-endowed good salient, which creates a valuation that is inflated (compared to a no-saliency benchmark). For example, in Thaler (1985), when people are asked about their willingness-to-pay for a beer on a hot day, most people will value hotel cans more than the cans from a normal corner shop, even though they are identical goods.

<sup>52</sup>Formally, define a sign function by  $\mu(a_k - \bar{a}) = 1$  iff  $a_k - \bar{a} \geq 0$  and  $\mu(a_k - \bar{a}) = -1$  iff  $a_k - \bar{a} < 0$ . Ordering is the property that  $\sigma(a_k + \mu(a_k - \bar{a})\epsilon, \bar{a} - \mu(a_k - \bar{a})\epsilon') > \sigma(a_k, \bar{a})$ , for  $\epsilon, \epsilon' \geq 0$  and  $\epsilon + \epsilon' > 0$ .

and underweight the unsalient ones.

Second, like our work, Saliency Theory was clearly motivated by ideas and evidence in psychology and neuroscience. Ancestors of context-sensitivity and the ordering property are common in both historical and modern psychology. (For example, we have repeatedly noted the importance of low-level contextual contrast in Itti et al. (1998) and later algorithms.) Diminishing sensitivity is also a ubiquitous psychophysical (Weber-Fechner) principle of perception. William James’s (1863) speculative list of things that engage “passive immediate sensorial attention” included “strange things” which can be translated as context-deviating attributes or objects. In modern neuroscience, saliency is often defined as absolute magnitude (deviation from zero) and is known to be encoded in the brain (Litt et al., 2011; Armel et al., 2008; McCoy and Platt, 2005).

Recent perceptual judgment experiments (Kunar et al. (2017)) illustrate one way that saliency of extreme values impacts judgment. Participants saw sequences of 12 two-digit numbers, presented rapidly ( $< 100\text{msec}$ ) one at a time. Judgments reflected more attention to the highest and lowest numbers in each stream (which are those with the highest BGS saliency; see also Tsetsos et al. 2012).<sup>53</sup> Larger differences are also more salient when people are looking for one target object out of many (including “distractors”). The target is easier to find when it is more different than distractors on features— such as searching for an X in a group of O’s rather than in a group of Y’s. The target-distractor differences should be expressible as numbers similar to normalized values of  $|a_k - \bar{a}|$  (Wolfe and Horowitz 2017, p. 2), as in Saliency Theory, but we do not know of direct equivalences of this sort.<sup>54</sup>

<sup>53</sup>Kunar et al. (2017) also found that when people were instructed to report whether they saw a specific target number, they missed that number more often when it was preceded by the highest or lowest number in the sequence. This is consistent with the joint hypothesis that people were more attentive to the extreme numbers, and exhibit a typical “attentional blink” in which attention lapses a bit after the high attention paid to extreme numbers.

<sup>54</sup>Wolfe and Horowitz (2017) compile a list of visual properties of features that robustly “guide” attention. In vision science jargon, a variable X guides attention if a target having property X increases the accuracy and speed of finding that target. Relative size and higher subjective value are two guiding variables in Wolfe and Horowitz (2017).

Because of this generality, Saliency Theory has been used to explain or interpret phenomena and empirical evidence in finance, lottery choices (including drug trafficking), legal judgment, price-quality markets and cross-game attention (Bordalo et al., 2013a; Cosemans and Frehen, 2020; Spitmaan et al., 2019; Dertwinkel-Kalt and Köster, 2020; Magliocca et al., 2019; Bordalo et al., 2015; Dertwinkel-Kalt et al., 2017; Avoyan and Schotter, 2020).

Saliency Theory and stimulus-driven salience (as defined and applied above) focus on different aspects of salience and their implications. In most applications the two theories do not make competing predictions, without additional specialized assumptions. The experiment 1 fruit sets design is an example. SAM salience predicts visual salience of images, then investigates whether that special type of salience affects choices. In contrast, Saliency Theory is about salience of valued attributes, regardless of how they are displayed or described, so it does not have a natural role for aspects of visual salience that are unrelated to attribute values.

Both theories are simplifications which have advantages and limits. Saliency Theory has the advantage of portability to many familiar microeconomic and social science applications. It benefits from the simplicity which comes from ignoring details of visual perception. Algorithmic SAM-type salience has the advantage of predicting rapid stimulus-driven visual attention for all possible images, but applying the theory to familiar domains such as price-quality competition is not straightforward (as noted in our discussion of explainable AI) and stimulus-constrained.

### *B. Rational inattention*

“Rational inattention” (RI) models assume that people optimally trade off the benefits and costs of paying closer attention. In more technical terms, endogeneously-allocated attention creates a subjective perception of objective factors. More accurate subjective perception is more costly but also improves expected decision



value.<sup>55</sup> These models are goal-directed because there is a clear goal— better perception is chosen to improve decision value.

RI models often start with a prior belief distribution  $\mu$  over a set of states  $\{\omega|\omega \in \Omega\}$ . In our fruit experiment, each  $\omega$  is a possible image. For each image, there is an optimal action  $a \in \{L, R\}$  (left or right, depending on which has the higher induced value). Denote the optimal actions by  $a^*(\omega)$ . There is also a pair of numbers  $S^L(\omega), S^R(\omega)$  which are the predicted SAM saliencies in the L and R halves of an image  $\omega$ .<sup>56</sup>

In RI, attention creates a set of latent “signals”  $\gamma(\omega)$  from a mapping  $\pi : \Omega \rightarrow \Delta(\Gamma)$  (Caplin and Dean, 2015; Caplin et al., 2019). In the fruit example,  $\gamma(\omega)$  could be the subjective belief probability of image  $\omega$  after all the learning processes. The “rationality” in RI comes from the assumption that the signal structure is chosen to maximize a gross decision value minus a cost of attention. The key term in the decision value is  $\max_{a \in A} \sum_{\omega \in \Omega} \gamma(\omega)u(a, \omega)$ . Since the saliencies  $S^L(\omega), S^R(\omega)$  do not enter the utility function  $u(a, \omega)$  and do not provide information about the optimal action  $a^*(\omega)$ , an RI agent should ignore them.<sup>57</sup> However, the results from the fruit experiment show that stimulus-driven salience can interfere with goal-directed RI and moves decisions away from RI optimality.

### C. *Dynamic channeled inattention and Bayesian surprise*

Some economic models seek to understand the dynamic effects of limited attention. This is different than our use of predicted salience to understand static choices.

Schwartzstein (2014) studies a problem of forecasting a binary variable  $y$  which depends on  $x$  and a subjectively encoded variable  $z$ . When  $z$  is expected to be

<sup>55</sup>For more detail see Sims (2003, 2006); Caplin and Dean (2015); Caplin et al. (2019); Kőszegi and Matějka (2020); Caplin et al. (2020); Mackowiak et al. (2020).

<sup>56</sup>To be clear, the fruits experiment is not an ideal proper test of RI. To do so would require controlling the set  $\Omega$  more carefully, and assuming, measuring, or inducing a prior belief that salience and induced value are uncorrelated, which was not done.

<sup>57</sup>It would be useful to figure out precisely how to integrate the effect of stimulus-driven salience into RI, to explain examples like the fruits experiment. Li (2020) p. 81 provides a saliency-sensitive state separation that can explain the saliency effect in simple choices.

important enough in forecasting  $y$ , with an expected value above a “busyness” threshold  $b$ ,  $z$  is accurately encoded. Otherwise,  $z$  is ignored and if  $z$  is ignored, no missing value is imputed.

Gagnon-Bartsch et al. (2018) proposed a similar idea of “channeled attention” during learning, in which people do not always recognize the results of their inattention. For example, a person who often forgets to take her medicines, but does not have a strong prior belief that she might forget, does not notice or keep track of her forgetting. She won’t pay for a reminder technology. They refer to these missed data as “statistical gorillas” (from the famous attention-blindness experiment of Simons and Chabris 1999). They derive dynamic conditions under which statistical gorillas will be noticed or not.

In dynamic image sequences, such as movies, one property of images that is known (from eyetracking) to grab attention strongly is called “Bayesian surprise”. This concept begins with a prior belief over “models” in model space  $\mathcal{M}$ . Itti and Baldi (2009) used an example in which a person turns on her TV, not knowing what channel was last watched and will pop up first.  $\mathcal{M}$  is the set of possible TV channels.  $P(D|M)$  are the likelihoods of perceptual data  $D$  conditional on a model  $M$  (a TV channel). For example, if blonde women are more common on  $M = \{\text{Fox News}\}$  than other channels, then  $P(\text{blonde women}|\{\text{Fox News}\}) > P(\text{blonde women}|\mathcal{M})$ .

“Surprise” for a given  $(D, M)$  combination is defined as  $S(D, M) \equiv \log \frac{P(M)}{P(M|D)}$ .<sup>58</sup> A person might be greatly surprised, for example, by seeing a blonde woman on the sports channel ESPN if  $P(\text{ESPN}) \gg P(\text{ESPN}|\text{blonde woman})$ . The ratio  $\frac{P(\text{ESPN})}{P(\text{ESPN}|\text{blonde woman})}$  and its logarithm will then be much greater than one, measuring how surprising that data-model combination is. Experienced surprise from data  $D$ , averaged over model posteriors, is a measure of overall experienced

<sup>58</sup>There is a loose relation between the ratio  $\frac{P(M)}{P(M|D)}$  and a concept of representativeness as relative likelihood  $P(D|M_1)/P(D|M_2)$  (see Tenenbaum et al. 2001 and Bordalo et al. 2016 for stereotypes, where  $D$  is a social type and models  $M$  are groups). The surprise ratio for a particular  $M$  is a measure of how *unrepresentative* or *anomalous*  $D$  is, and the summation adds up the total degree of unrepresentativeness of  $D$  for all models  $M$ .

surprise from perceptual data  $D$ :

$$\sum_M P(M|D)S(D, M)$$

Note that Bayesian surprise does not fit into the stimulus-driven vs. goal-directed dichotomy. It depends on a perceiver’s prior beliefs so it is not purely stimulus-driven. But surprise-detection is also highly general and is therefore not typically considered a perceptual goal like, say, searching for a familiar face in a crowd or for a high resale value fruit.<sup>59</sup>

Bayesian surprise is not used in the types of experiments in this paper because the presented images were not deliberately linked in a dynamic sequence (as in a movie). However, in typical experiments prior perceptual beliefs are induced by short exposures to each of a large number of images, so that what is surprising in a subsequent image (relative to those priors) can be quantified. This could easily be done in the fruits experiment. For example, if many images in a row included no apples, then in a new image with an apple, the apple would be Bayesian-surprising and is predicted to be salient and attract attention.

The Bayesian surprise model is well-supported experimentally (Itti and Baldi (2009)) and has the advantage that some analytical results are available for the class of conjugate priors (Baldi and Itti, 2010). Potential economic applications include a sequential visual presentation of price changes in a time series, or testing for salience from a new advertising campaign, product design, or logo change.<sup>60</sup>

<sup>59</sup>Prof. Pierre Baldi said in a personal communication that “...Bayesian surprise is agnostic with respect to any bottom-up or top-down considerations.”

<sup>60</sup>Note that there is an apparent opposition between ignoring statistical gorillas in Gagnon-Bartsch et al. (2018) and Bayesian surprise. A gorilla on a basketball court is typically very high in Bayesian surprise and hence predicted to be quite salient; then why don’t people notice the gorilla? The answer is that scarce attention is focused on one mentally taxing goal— counting basketball passes (the instructed goal in the seminal study)— so that a Bayesian-surprising object is ignored. Magic tricks work the same general way: Skillful “misdirection” draws attention away from the sneaky sleight of hand (Macknik et al., 2008; Wiseman and Nakano, 2016). In economic settings, Bayesian surprise and other goal-directed attention will be productive substitutes, an hypothesis which can be tested by phenomena like timing and reaction to unusual corporate earnings announcements (e.g., DeHaan et al. (2015)).

*D. Relative attention  $m(x)$*

It is useful to have a simple measure of inattention, as revealed by choices, to compare across domains. A good one is summarized by Gabaix (2019). Define both rational and behavioral actions, as a function of a perceived normative variable  $x$  (such as a price), by  $a^r(x) = \operatorname{argmax}_a u(a, x)$  (rational) and  $a^b(x) = a^r(mx)$  (behavioral). The behavioral model is assumed to maximize, but underperceives or underweights the true variable value  $x$ , shrinking it toward zero to a degree measured by a parameter  $m < 1$ . (A canonical example is paying too little attention to a hidden component of price, such as taxes, where  $m < 1$  measures the degree of tax underweighting.) Gabaix (2019) shows that  $m(x)$  can be recovered from the ratio of marginal effects of the  $x$  variable on actions in different attention treatments,  $\frac{a_x^b}{a_x^r}$ . His Table 1 summarizes numerical estimates from several field experiments and datasets.

A version of  $m(x)$  can also be computed from the fruit experiment data based on an *ad hoc* assumption. Suppose choice under time pressure is designated as the “behavioral” condition and choice with unlimited time is designated as the “rational” condition. The intuition is that in the behavioral condition stimulus-driven salience is not fully inhibited (even though it is irrelevant), which reduces the influence of goal-directed attention ( $m < 1$ ) to the induced value  $x$  variable. From Table 1 the marginal effect of the normative variable (the induced value difference) on choice accuracy is  $\hat{a}_x^b = .795$ . The value of  $\hat{a}_x^r$  can be computed from the same regression as in Table 1, using data from the unlimited time treatment. That value turns out to be  $\hat{a}_x^r = 2.249$ . The ratio of the behavioral and rational coefficient estimates is therefore  $\frac{\hat{a}_x^b}{\hat{a}_x^r} = \frac{.795}{2.249}$ , which is .35. This figure is close to the mean  $m(x) = .44$  reported in Gabaix (2019) Table 1. This numerical exercise shows how the effect of stimulus-driven salience as a behavioral condition can be compared numerically to other kinds of limited attention.

## VIII. Discussion and Conclusion

Our study leads to two new conclusions:

- 1) Stimulus-driven salience can be predicted by an underlying neuro-computational theory (SAM) of which features of an image or information display most people look at first. SAM-estimated salience has a small but significant effect in visualized binary set choices (fruit sets) and in matrix games. These effects are not always strong because in both cases stimulus-driven attention competes with goal-directed attention in a way that SAM the algorithm does not attempt to predict.
- 2) In the main set of experiments with location matching games, salience is a good predictor of which location people choose, and how often their choices match ( $r=-.57$ ). In hider-seeker games, a salience-influenced cognitive hierarchy model (and a similar level-k model) can account for the small, but robust, seeker’s advantage in hider-seeker games. Parameters fit to hider-seeker data can also “portably” predict the salience-choice relation in matching games, even though the hider-seeker game is strictly competitive and matching is cooperative.

### A. *Where else in economics could salience be useful?*

Before proceeding to further visual salience speculation, note that vision is only one of five senses; other sensory systems have salience structures too. Auditory (sound) attention is also driven by both goal-directed and stimulus-driven processes. One can attend to an important conversation to achieve a social goal while tuning out background noises at a party. But the stimulus-driven system will hijack attention if a champagne glass shatters with a loud crash.<sup>61</sup> Research parallel

<sup>61</sup>A general example is the stimulus-driven salience of human screams (which have an unpleasant power spectrum quality called “roughness”). Screams are rated more quickly as fear-inducing, are more accurately localized, and activate the amygdala and primary auditory cortex more strongly (Arnal et al., 2015). Kaya and Elhilali (2014) proposed a salience map based on five features (envelope, harmonicity, spectrogram, bandwidth, and modulation) and tested it.

to ours could explore auditory salience in domains like advertising, business communication, security analyst earnings calls, open-outcry auctions, negotiations, etc.

This last section speculates how an empirical understanding of stimulus-driven salience might improve other economic studies.

- *Behavioral IO*: The fruits experiment is a paradigm that invites thinking and future exploration about the supply-side response to consumer psychology, a subfield called behavioral industrial organization (Heidhues and Kőszegi 2018, and others). A central concept in behavioral IO is whether product attributes are “shrouded” (Gabaix and Laibson, 2006)– that is, deliberately hidden by sellers. Measuring whether attributes are low in stimulus-driven salience is one scientific measure of shrouding, which is perhaps useful for consumer policy regulators.

By understanding stimulus-driven salience, a retailer could create a product display with the goal of maximizing profit margin. High-margin items would be displayed to maximize their stimulus-driven salience. An open and interesting question is whether consumers can recognize and ignore such supply-side salience manipulations.

- *Tax and price salience in consumer markets*: Price and value components that are presented to sensory systems, such as explicit price tags that the eye can see, seem to receive more decision weight than equivalent components that need to be imagined and computed. This effect was first shown for unit-cost price tags by (Russo, 1977) and has been shown carefully in many recent studies (Ott and Andrus, 2000; Hossain and Morgan, 2006; Min Kim and Kachersky, 2006; Finkelstein, 2009; Taubinsky and Rees-Jones, 2017). In principle, SAM could be applied to visual images of store price tags or e-commerce websites, as was done in the fruit-valuation Study 1, to guess the visual salience of explicit and hidden prices. These measures could be

compared to salience as measured from behavior in these papers, and as summarized in the  $m(x)$  measure in Gabaix’s (2019) Table 1.

- *Nudges and design*: “Nudges” are changes in design and choice architecture, which do not drastically change information content or incentives, but can make information processing simpler and improve decisions.<sup>62</sup> Many nudge experiments have been done and are ongoing. But their effects are often unpredictable (Milkman et al., 2021; DellaVigna and Linos, 2020).

Predictions about what nudges are visually salient might help us understand what has worked and create better designs. If a financial regulator is trying to design a form to nudge goal-directed attention toward particular information, for example, their design will probably work better if the targeted information also has stimulus-driven salience (e.g. Hilchey et al. 2021).

- *Beliefs*: Besides influencing choices, visual salience can influence what information is processed and what *beliefs* result.<sup>63</sup> Padilla et al. (2017) showed a striking example of an effect of stimulus-driven salience on beliefs about hurricanes. The National Hurricane Center currently shows potential geospatial paths with a “cone of uncertainty”, a 2D confidence interval forecasting a range of areas that a hurricane might conceivably reach. The cone becomes wider, spreading out geographically, for forecasts projecting more days ahead (which are typically more uncertain). An alternative visualization is an “ensemble plot” which shows many distinct possible individual paths and does not draw a cone around them (cf. “spaghetti plots”). Padilla et al. (2017) apply the Itti et al. (1998) algorithm (a precursor to SAM) to these two different visualizations. The algorithm predicts that cone plots will focus attention on the center and on the furthest boundaries of the cone, where the cone is widest. This perception biases actual human judgments of whether the hurricane will grow in storm size and intensity (e.g., wind

<sup>62</sup>See Goldin (2015); Thaler and Sunstein (2009); Luo et al. (2021).

<sup>63</sup>See Padilla et al. 2018; Itti et al. 1998; Mackowiak et al. 2020.

speed) in the future. These subjective beliefs reflect a cognitive mistake: People think the growing size of the *cone* predicts that the size of the *storms* and their intensity will grow.

The ensemble plot has different predicted salience and a different effect on beliefs. Predicted salience is highest at the location where different paths are clustered before they diverge into different paths. Attention is widely dispersed over the ending points of the different trajectories (rather than concentrated at the cone plot boundary). As a result, judgments about future storm size and intensity are not infected by a size bias (as they are from cone plots). Thus, the cone plot leads to mistaken beliefs and the ensemble plot does not. The salience algorithm accurately predicted the direction of that effect.

An economic example of a similar kind is the visualization of regression discontinuity effects. Korting et al. (2020) show that axis-scaling, x-axis bin width, and spacing all influence the perceptions people have about causal effects when shown different graphs based on the exact same data. SAM or other salience algorithms could be applied to these data, to learn more about how stimulus-driven processes affect what scientific consumers think a graph is telling them.

## REFERENCES

- Arieli, A., Y. Ben-Ami, and A. Rubinstein (2011). Tracking decision makers under uncertainty. *American Economic Journal: Microeconomics* 3(4), 68–76.
- Armel, C., A. Beaumel, and A. Rangel (2008). Biasing simple choices by manipulating relative visual attention. *Judgment and Decision making* 3(5), 396–403.
- Arnal, L. H., A. Flinker, A. Kleinschmidt, A.-L. Giraud, and D. Poeppel (2015). Human screams occupy a privileged niche in the communication soundscape. *Current Biology* 25(15), 2051–2056.



- Arrieta, A. B., N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al. (2020). Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion* 58, 82–115.
- Aumann, R. J. (1974). Subjectivity and correlation in randomized strategies. *Journal of Mathematical Economics* 1(1), 67–96.
- Avoyan, A. and A. Schotter (2020). Attention in games: An experimental study. *European Economic Review*, 103410.
- Awh, E., A. V. Belopolsky, and J. Theeuwes (2012, August). Top-down versus bottom-up attentional control: a failed theoretical dichotomy. *Trends in Cognitive Sciences* 16(8), 437–443.
- Bacharach, M. (1993). 13 variable universe games. *Frontiers of Game Theory*, 255.
- Bacharach, M. and M. Bernasconi (1997). The variable frame theory of focal points: An experimental study. *Games and Economic Behavior* 19(1), 1–45.
- Baldi, P. and L. Itti (2010). Of bits and wows: A bayesian theory of surprise with applications to attention. *Neural Networks* 23(5), 649–666.
- Baluch, F. and L. Itti (2011). Mechanisms of top-down attention. *Trends in neurosciences* 34(4), 210–224.
- Bar-Hillel, M. (2015). Position effects in choice from simultaneous displays: A conundrum solved. *Perspectives on Psychological Science* 10(4), 419–433.
- Belle, V. and I. Papantonis (2020). Principles and practice of explainable machine learning. *arXiv preprint arXiv:2009.11698*.
- Bordalo, P., K. Coffman, N. Gennaioli, and A. Shleifer (2016). Stereotypes. *The Quarterly Journal of Economics* 131(4), 1753–1794.

- Bordalo, P., N. Gennaioli, and A. Shleifer (2012a). Salience in experimental tests of the endowment effect. *American Economic Review* 102(3), 47–52.
- Bordalo, P., N. Gennaioli, and A. Shleifer (2012b). Salience theory of choice under risk. *The Quarterly Journal of Economics* 127(3), 1243–1285.
- Bordalo, P., N. Gennaioli, and A. Shleifer (2013a). Salience and asset prices. *American Economic Review* 103(3), 623–28.
- Bordalo, P., N. Gennaioli, and A. Shleifer (2013b). Salience and consumer choice. *Journal of Political Economy* 121(5), 803–843.
- Bordalo, P., N. Gennaioli, and A. Shleifer (2015). Salience theory of judicial decisions. *Journal of Legal Studies* 44(S1), S7–S33.
- Bornstein, R. F. (1989). Exposure and affect: overview and meta-analysis of research, 1968–1987. *Psychological bulletin* 106(2), 265.
- Brocas, I. and J. D. Carrillo (2021). The development of randomization and deceptive behavior in mixed strategy games.
- Brocas, I., J. D. Carrillo, S. W. Wang, and C. F. Camerer (2014). Imperfect choice or imperfect attention? Understanding strategic thinking in private information games. *Review of Economic Studies* 81(3), 944–970.
- Camerer, C. F., T.-H. Ho, and J.-K. Chong (2004). A cognitive hierarchy model of games. *The Quarterly Journal of Economics* 119(3), 861–898.
- Camerer, C. F., E. Johnson, T. Rymon, and S. Sen (1993). Cognition and framing in sequential bargaining for gains and losses. *Frontiers of game theory* 104, 27–47.
- Caplin, A., D. Csaba, J. Leahy, and O. Nov (2020). Rational inattention, competitive supply, and psychometrics. *The Quarterly Journal of Economics* 135(3), 1681–1724.

- Caplin, A. and M. Dean (2015). Revealed preference, rational inattention, and costly information acquisition. *American Economic Review* 105(7), 2183–2203.
- Caplin, A., M. Dean, and J. Leahy (2019). Rational inattention, optimal consideration sets, and stochastic choice. *The Review of Economic Studies* 86(3), 1061–1094.
- Cerf, M., J. Harel, W. Einhäuser, and C. Koch (2008). Predicting human gaze using low-level saliency combined with face detection. In *Advances in Neural Information Processing Systems*, pp. 241–248.
- Chen, C., X. Zhang, Y. Wang, T. Zhou, and F. Fang (2016). Neural activities in v1 create the bottom-up saliency map of natural scenes. *Experimental Brain Research* 234(6), 1769–1780.
- Chun, M. M., J. D. Golomb, and N. B. Turk-Browne (2011). A taxonomy of external and internal attention. *Annual Review of Psychology* 62, 73–101.
- Cornia, M., L. Baraldi, G. Serra, and R. Cucchiara (2018). Predicting human eye fixations via an lstm-based saliency attentive model. *IEEE Transactions on Image Processing* 27(10), 5142–5154.
- Cosemans, M. and R. Frehen (2020). Saliency theory and stock prices: empirical evidence. *Journal of Financial Economics*.
- Costa-Gomes, M., V. P. Crawford, and B. Broseta (2001). Cognition and behavior in normal-form games: An experimental study. *Econometrica* 69(5), 1193–1235.
- Costa-Gomes, M. A. and V. P. Crawford (2006). Cognition and behavior in two-person guessing games: An experimental study. *American Economic Review* 96(5), 1737–1768.
- Crawford, V. P. (2014). A comment on how portable is level-0 behavior? a test

of level-k theory in games with non-neutral frames by heap, rojo-arjona, and sugden.

Crawford, V. P., M. A. Costa-Gomes, and N. Iriberry (2013). Structural models of nonequilibrium strategic thinking: Theory, evidence, and applications. *Journal of Economic Literature* 51(1), 5–62.

Crawford, V. P. and N. Iriberry (2007a). Fatal attraction: Saliency, naivete, and sophistication in experimental “hide-and-seek” games. *American Economic Review* 97(5), 1731–1750.

Crawford, V. P. and N. Iriberry (2007b). Level-k Auctions: Can a Nonequilibrium Model of Strategic Thinking Explain the Winner’s Curse and Overbidding in Private-Value Auctions? *Econometrica* 75(6), 1721–1770.

Cunningham, T. (2013). Biases and implicit knowledge.

DeHaan, E., T. Shevlin, and J. Thornock (2015). Market (in) attention and the strategic scheduling and timing of earnings announcements. *Journal of Accounting and Economics* 60(1), 36–55.

DellaVigna, S. and E. Linos (2020). Rcts to scale: Comprehensive evidence from two nudge units. Technical report, Working Paper, UC Berkeley.

Dertwinkel-Kalt, M., K. Köhler, M. R. Lange, and T. Wenzel (2017). Demand shifts due to saliency effects: Experimental evidence. *Journal of the European Economic Association* 15(3), 626–653.

Dertwinkel-Kalt, M. and M. Köster (2020). Saliency and skewness preferences. *Journal of the European Economic Association* 18(5), 2057–2107.

Devetag, G., S. Di Guida, and L. Polonio (2016). An eye-tracking study of feature-based choice in one-shot games. *Experimental Economics* 19(1), 177–201.

- Falk, R., R. Falk, and P. Ayton (2009). Subjective patterns of randomness and choice: Some consequences of collective responses. *Journal of Experimental Psychology: Human Perception and Performance* 35(1), 203.
- Fan, F., J. Xiong, and G. Wang (2020). On interpretability of artificial neural networks. *arXiv preprint arXiv:2001.02522*.
- Finkelstein, A. (2009). E-ztax: Tax salience and tax rates. *The Quarterly Journal of Economics* 124(3), 969–1010.
- Fragiadiakis, D., A. Kovaliukaite, and D. R. Arjona (2017). Testing cognitive hierarchy assumptions. Technical report, Working paper, 2017. 4, 13, 14, 15, 17.
- Fudenberg, D. and A. Liang (2019). Predicting and understanding initial play. *American Economic Review* 109(12), 4112–41.
- Fudenberg, D., P. Strack, and T. Strzalecki (2018). Speed, accuracy, and the optimal timing of choices. *American Economic Review* 108(12), 3651–84.
- Gabaix, X. (2019). Behavioral inattention, handbook of behavioral economics: Applications and foundations. *Chapter 4*, 261–343.
- Gabaix, X. and D. Laibson (2006). Shrouded attributes, consumer myopia, and information suppression in competitive markets. *The Quarterly Journal of Economics* 121(2), 505–540.
- Gagnon-Bartsch, T., M. Rabin, and J. Schwartzstein (2018). *Channeled attention and stable errors*. Harvard Business School.
- Goldin, J. (2015). Which way to nudge: Uncovering preferences in the behavioral age. *Yale Law Journal* 125, 226.
- Haji-Abolhassani, A. and J. J. Clark (2014). An inverse yarbus process: Predicting observers’ task from eye movement patterns. *Vision Research* 103, 127–142.

- Harel, J., C. Koch, and P. Perona (2007). Graph-based visual saliency. In *Advances in neural information processing systems*, pp. 545–552.
- Hargreaves Heap, S., D. Rojo Arjona, and R. Sugden (2014). How Portable Is Level-0 Behavior? A Test of Level-k Theory in Games With Non-Neutral Frames. *Econometrica* 82(3), 1133–1151.
- Hartford, J. S., J. R. Wright, and K. Leyton-Brown (2016). Deep learning for predicting human strategic behavior. In *Advances in Neural Information Processing Systems*, pp. 2424–2432.
- He, S., H. R. Tavakoli, A. Borji, Y. Mi, and N. Pugeault (2019). Understanding and visualizing deep visual saliency models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10206–10215.
- Heidhues, P. and B. Köszegi (2018). Behavioral industrial organization. *Handbook of Behavioral Economics: Applications and Foundations 1* 1, 517–612.
- Heinrich, T. and I. Wolff (2012). Strategic reasoning in hide-and-seek games: A note.
- Henderson, J. M. and T. R. Hayes (2017). Meaning-based guidance of attention in scenes as revealed by meaning maps. *Nature Human Behaviour* 1(10), 743.
- Hilchey, M. D., M. Osborne, and D. Soman (2021). Does the visual salience of credit card features affect choice? *Behavioural Public Policy*, 1–18.
- Hinton, G., O. Vinyals, and J. Dean (2015). Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Ho, T.-H., C. Camerer, and K. Weigelt (1998). Iterated dominance and iterated best response in experimental “p-beauty contests”. *The American Economic Review* 88(4), 947–969.

- Hossain, T. and J. Morgan (2006). ... plus shipping and handling: Revenue (non) equivalence in field experiments on ebay. *Advances in Economic Analysis & Policy* 6(2).
- Igami, M. (2020). Artificial intelligence as structural estimation: Deep blue, bonanza, and alphago. *The Econometrics Journal* 23(3), S1–S24.
- Itti, L. and P. Baldi (2009). Bayesian surprise attracts human attention. *Vision Research* 49(10), 1295–1306.
- Itti, L., C. Koch, and E. Niebur (1998). A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20(11), 1254–1259.
- James's, W. (1863). Principles of psychology.
- Johnson, E. J., C. Camerer, S. Sen, and T. Rymon (2002). Detecting failures of backward induction: Monitoring information search in sequential bargaining. *Journal of Economic Theory* 104(1), 16–47.
- Judd, T., F. Durand, and A. Torralba (2012). A benchmark of computational models of saliency to predict human fixations.
- Judd, T., K. Ehinger, F. Durand, and A. Torralba (2009). Learning to predict where humans look. In *Computer Vision, 2009 IEEE 12th international conference on*, pp. 2106–2113. IEEE.
- Kaya, E. M. and M. Elhilali (2014). Investigating bottom-up auditory attention. *Frontiers in human neuroscience* 8, 327.
- Kneeland, T. (2015). Identifying Higher–Order Rationality. *Econometrica* 83(5), 2065–2079.
- Korting, C., C. Lieberman, J. Matsudaira, Z. Pei, and Y. Shen (2020). Visual inference and graphical representation in regression discontinuity designs.

- Kőszegi, B. and F. Matějka (2020). Choice simplification: A theory of mental budgeting and naive diversification. *The Quarterly Journal of Economics* 135(2), 1153–1207.
- Kőszegi, B. and A. Szeidl (2013). A model of focusing in economic choice. *The Quarterly Journal of Economics* 128(1), 53–104.
- Krasovskaya, S. and W. J. MacInnes (2019). Saliency models: A computational cognitive neuroscience review. *Vision* 3(4), 56.
- Kummerer, M., T. S. Wallis, and M. Bethge (2018). Saliency benchmarking made easy: Separating models, maps and metrics. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 770–787.
- Kunar, M. A., D. G. Watson, K. Tsetsos, and N. Chater (2017). The influence of attention on value integration. *Attention, Perception, & Psychophysics* 79(6), 1615–1627.
- Lewis, D. (2008). *Convention: A philosophical study*. John Wiley & Sons.
- Li, X. (2020). *Attention, Strategy, and the Human Mind*. Ph. D. thesis, Division of Humanities and Social Science, Caltech.
- Lipton, Z. C. (2018). The mythos of model interpretability. *Queue* 16(3), 31–57.
- Litt, A., H. Plassmann, B. Shiv, and A. Rangel (2011). Dissociating valuation and saliency signals during decision-making. *Cerebral cortex* 21(1), 95–102.
- Luo, Y., D. Soman, and J. Zhao (2021). A meta-analytic cognitive framework of nudge and sludge.
- Macknik, S. L., M. King, J. Randi, A. Robbins, J. Thompson, S. Martinez-Conde, et al. (2008). Attention and awareness in stage magic: turning tricks into research. *Nature Reviews Neuroscience* 9(11), 871–879.



- Mackowiak, B., F. Matejka, M. Wiederholt, et al. (2020). Rational inattention: A review. *CEPR Discussion Papers* (15408).
- Magliocca, N. R., K. McSweeney, S. E. Sesnie, E. Tellman, J. A. Devine, E. A. Nielsen, Z. Pearson, and D. J. Wrathall (2019). Modeling cocaine traffickers and counterdrug interdiction forces as a complex adaptive system. *Proceedings of the National Academy of Sciences* 116(16), 7784–7792.
- McCoy, A. N. and M. L. Platt (2005). Risk-sensitive neurons in macaque posterior cingulate cortex. *Nature neuroscience* 8(9), 1220–1227.
- Mehta, J., C. Starmer, and R. Sugden (1994a). Focal points in pure coordination games: An experimental investigation. *Theory and Decision* 36(2), 163–185.
- Mehta, J., C. Starmer, and R. Sugden (1994b). The nature of salience: An experimental investigation of pure coordination games. *The American Economic Review* 84(3), 658–673.
- Milkman, K., D. Gromet, H. Ho, J. Kay, T. Lee, P. Pandiloski, Y. Park, Y. Rai, M. Bazerman, J. Beshears, L. Bonacorsi, C. Camerer, E. Chang, E. Chapman, R. Cialdini, H. Dai, L. Eskreis-Winkler, A. Fishbach, J. Gross, A. Horn, A. Hubbard, J. SJ, D. Karlan, T. Kautz, E. Kirgios, E. Klusowski, A. Kristal, R. Ladhania, G. Loewenstein, J. Ludwig, B. Mellers, S. Mullainathan, S. Saccardo, J. Spiess, G. Suri, J. Talloen, J. Taxer, Y. Trope, L. Ungar, K. Volpp, A. Whillans, J. Zinman, and A. Duckworth (2021). A mega-study approach to applied behavioral science. *Nature*, *in press*.
- Milosavljevic, M., V. Navalpakkam, C. Koch, and A. Rangel (2012). Relative visual saliency differences induce sizable bias in consumer choice. *Journal of Consumer Psychology* 22(1), 67–74.
- Min Kim, H. and L. Kachersky (2006). Dimensions of price salience: a conceptual framework for perceptions of multi-dimensional prices. *Journal of Product & Brand Management* 15(2), 139–147.

- Mormann, M. and J. E. Russo (2021). Does attention increase the value of choice alternatives? *Trends in cognitive sciences*.
- Nagel, R. (1995). Unraveling in guessing games: An experimental study. *The American Economic Review* 85(5), 1313–1326.
- Ott, R. L. and D. M. Andrus (2000). The effect of personal property taxes on consumer vehicle-purchasing decisions: a partitioned price/mental accounting theory analysis. *Public Finance Review* 28(2), 134–152.
- Pachur, T., M. Schulte-Mecklenbeck, R. O. Murphy, and R. Hertwig (2018). Prospect theory reflects selective allocation of attention. *Journal of Experimental Psychology: General* 147(2), 147.
- Padilla, L. M., S. H. Creem-Regehr, M. Hegarty, and J. K. Stefanucci (2018). Decision making with visualizations: a cognitive framework across disciplines. *Cognitive research: principles and implications* 3(1), 29.
- Padilla, L. M., I. T. Ruginski, and S. H. Creem-Regehr (2017). Effects of ensemble and summary displays on interpretations of geospatial uncertainty data. *Cognitive research: principles and implications* 2(1), 1–16.
- Polonio, L., S. Di Guida, and G. Coricelli (2015). Strategic sophistication and attention in games: An eye-tracking study. *Games and Economic Behavior* 94, 80–96.
- Ras, G., M. van Gerven, and P. Haselager (2018). Explanation methods in deep learning: Users, values, concerns and challenges. In *Explainable and Interpretable Models in Computer Vision and Machine Learning*, pp. 19–36. Springer.
- Ratcliff, R. (1978). A theory of memory retrieval. *Psychological review* 85(2), 59.
- Ratcliff, R., P. L. Smith, S. D. Brown, and G. McKoon (2016). Diffusion decision

- model: Current issues and history. *Trends in Cognitive Sciences* 20(4), 260–281.
- Riche, N., M. Duvinage, M. Mancas, B. Gosselin, and T. Dutoit (2013). Saliency and human fixations: state-of-the-art and study of comparison metrics. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, pp. 1153–1160. IEEE.
- Rihn, A., X. Wei, and H. Khachatryan (2019). Text vs. logo: Does eco-label format influence consumers’ visual attention and willingness-to-pay for fruit plants? an experimental auction approach. *Journal of Behavioral and Experimental Economics* 82, 101–452.
- Rogers, B. W., T. R. Palfrey, and C. F. Camerer (2009). Heterogeneous quantal response equilibrium and cognitive hierarchies. *Journal of Economic Theory* 144(4), 1440–1467.
- Rubinstein, A., A. Tversky, and D. Heller (1997). Naive strategies in competitive games. In *Understanding Strategic Interaction*, pp. 394–402. Springer.
- Russo, J. E. (1977). The value of unit price information. *Journal of Marketing Research*, 193–201.
- Schelling, T. C. (1960). *The Strategy of Conflict*. Harvard university press.
- Schwartzstein, J. (2014). Selective attention and learning. *Journal of the European Economic Association* 12(6), 1423–1452.
- Shimojo, S., C. Simion, E. Shimojo, and C. Scheier (2003). Gaze bias both reflects and influences preference. *Nature neuroscience* 6(12), 1317–1322.
- Simons, D. J. and C. F. Chabris (1999). Gorillas in our midst: Sustained inattentive blindness for dynamic events. *Perception* 28(9), 1059–1074.
- Sims, C. A. (2003). Implications of rational inattention. *Journal of monetary Economics* 50(3), 665–690.

- Sims, C. A. (2006). Rational inattention: Beyond the linear-quadratic case. *American Economic Review* 96(2), 158–163.
- Smith, V. L. (1976). Experimental economics: Induced value theory. *The American Economic Review* 66(2), 274–279.
- Spitmaam, M., E. Chu, and A. Soltani (2019). Salience-driven value construction for adaptive choice under risk. *Journal of Neuroscience* 39(26), 5195–5209.
- Stahl II, D. O. and P. W. Wilson (1994). Experimental evidence on players' models of other players. *Journal of Economic Behavior & Organization* 25(3), 309–327.
- Steinbeck, J. (2011). *The pearl*. Penguin UK.
- Taubinsky, D. and A. Rees-Jones (2017). Attention variation and welfare: theory and evidence from a tax salience experiment. *The Review of Economic Studies* 85(4), 2462–2496.
- Tenenbaum, J. B., T. L. Griffiths, et al. (2001). The rational basis of representativeness. In *Proceedings of the 23rd annual conference of the Cognitive Science Society*, pp. 1036–1041. Citeseer.
- Thaler, R. (1985). Mental accounting and consumer choice. *Marketing science* 4(3), 199–214.
- Thaler, R. H. and C. R. Sunstein (2009). *Nudge: Improving decisions about health, wealth, and happiness*. Penguin.
- Towal, R. B., M. Mormann, and C. Koch (2013). Simultaneous modeling of visual saliency and value computation improves predictions of economic choice. *Proceedings of the National Academy of Sciences* 110(40), E3858–E3867.
- Tsetsos, K., N. Chater, and M. Usher (2012). Salience driven value integration explains decision biases and preference reversal. *Proceedings of the National Academy of Sciences* 109(24), 9659–9664.

- Veale, R., Z. M. Hafed, and M. Yoshida (2017). How is visual salience computed in the brain? insights from behaviour, neurobiology and modelling. *Philosophical Transactions of the Royal Society B: Biological Sciences* 372(1714), 20160113.
- Vig, E., M. Dorr, and D. Cox (2014). Large-scale optimization of hierarchical features for saliency prediction in natural images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2798–2805.
- Wiseman, R. J. and T. Nakano (2016). Blink and you’ll miss it: the role of blinking in the perception of magic tricks. *PeerJ* 4, e1873.
- Wolfe, J. M. and T. S. Horowitz (2017). Five factors that guide attention in visual search. *Nature Human Behaviour* 1(3), 1–8.
- Wright, J. R. and K. Leyton-Brown (2019). Level-0 models for predicting human behavior in games. *Journal of Artificial Intelligence Research* 64, 357–383.
- Yarbus, A. L. (2013). *Eye movements and vision*. Springer.
- Zajonc, R. B. (1968). Attitudinal effects of mere exposure. *Journal of personality and social psychology* 9(2p2), 1.

## Tables

TABLE I—INFLUENCE OF SALIENCE-VALUE CONGRUENCY IN A SIMPLE CHOICE PROBLEM (FRUIT SETS)

<i>Dependent variable: Accuracy (0,1)</i>					
	(1)	(2)	(3)	(4)	(5)
Congruency	0.83*** (0.32)	0.90*** (0.29)	0.89*** (0.33)	0.97*** (0.31)	1.26*** (0.41)
abs(valueDiff)		0.80*** (0.23)		0.80*** (0.23)	0.77*** (0.23)
Interaction: Congruency*abs(valueDiff)					-0.55 (0.63)
Constant	1.54*** (0.10)	0.78*** (0.17)	1.99*** (0.61)	1.25** (0.57)	1.26** (0.57)
Covariates	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>
Observations	1,382	1,382	1,307	1,307	1,307
Log Likelihood	-644.7	-591.8	-607.5	-556.7	-556.3
Akaike Inf. Crit.	1,293	1,189	1,239	1,139	1,141

*Note:*

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

*Note:* The congruency variable is the difference between the the maximum salience level of the more valuable set and the maximum salience level of the less valuable set (between 0-1). This variable will be positive if one option is both more salient and more valuable. abs(ValueDiff) is the absolute value of the induced-value difference between left and right sets. Standard deviations are clustered on per subject level. "Covariates" denotes whether the current model contains covariates: education, gender, income and self-reported fruit preference (we ask them which fruit they prefer in everyday consumption: apples, oranges or equal preference). The main effect estimates are not sensitive to these covariates, as is evident comparing specifications (k-2) to (3-4).

TABLE II—REALIZED MATCHING RATE

	Matching rate	N of observations
Nash mixed prediction	0.071	
Matching game	0.64 (0.006)	559
Hider-seeker game	0.09 (0.002)	1060(531(H),529(S))
Hider-seeker game (between-subjects)	0.09 (0.002)	1325(600(H),725(S))
Hider-seeker high payoff (10x)	0.09 (0.003)	892(446(H),446(S))

*Note:* Statistical tests against the null hypothesis that the seeker win rate is the baseline level and choices are independently and identically distributed across subjects (which is the Nash benchmark prediction). The number in the bracket is the standard error of the seeking win-rate in each condition.

TABLE III—ESTIMATION DETAILS, ROLE-SPECIFIC SCH

-	$\lambda$	$\mu$	$\tau_h$	$\tau_s$
Best fit parameters	100	0.06	0.4	0.1
Number of observations	1096 for hiders and 1090 for seekers			
95% CI	[72.3,100]	[0.05,0.08]	[0.32,0.47]	[0.08,0.13]

*Note:* The parameters  $\mu$  and  $\lambda$  are constrained to be the same for both hiders and seekers. The confidence interval in the table is calculated using the bootstrap method with data batch size 1096 for hider, 1090 for seeker and the number of iterations is 100.



TABLE IV—THE EFFECT OF OF SALIENCE-EQUILIBRIUM CONGRUENCY IN MATRIX GAMES

	<i>Dependent variable:</i>	
	Whether the choice is an equilibrium strategy	
	(1)	(2)
Congruency (whether equilibrium strategy is salient)	0.008 (0.082)	-0.208 (0.142)
Congruency*Level-0		0.465** (0.189)
Congruency*Level-1		0.073 (0.197)
Constant	0.240* (0.140)	0.348** (0.143)
Observations	1,323	1,323
Log Likelihood	-910.061	-908.221
Akaike Inf. Crit.	1,834.122	1,834.443

*Note:*

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

*Note:* The dependent variable is (0-1) whether the chosen strategy in a matrix game is the equilibrium strategy. (All games in the dataset have a unique Nash strategy.) “Congruency” indicates whether the equilibrium option in that particular game is also more salient (which is the top row/left column). Covariates (coefficients not reported) are: game types (DSS, PD, DSO), Role (Row, Column), Levels (Level-0,Level-1,Level-2). Standard errors are clustered at the individual level.

## Figures

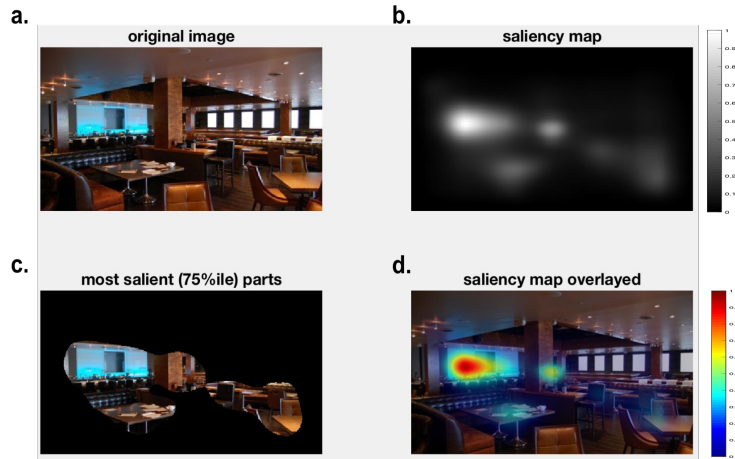


FIGURE I. A SALIENCE ALGORITHM EXAMPLE

*Note:* a: An original image. b: The SAM saliency map, in which the brightness indicates the saliency level. c: The area of the original image which is 75% most salient. This area is generated from ranking all saliency values of each pixel. d: The original image with the saliency heatmap overlaid onto it (the colorbar on the right indicates the corresponding saliency values).

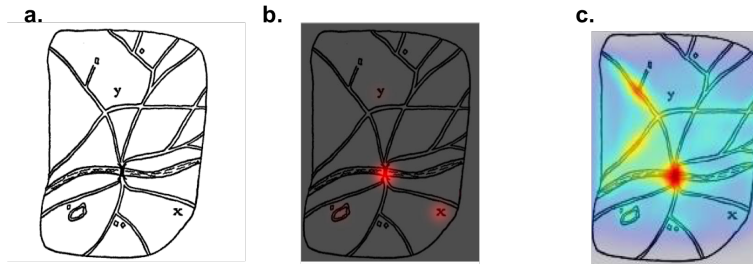


FIGURE II. SCHELLING'S MAP REVISITED

*Note:* (a) Original map; (b) Choice frequencies heatmap, where redness indicates choice frequency; and (c) The SAM algorithm predicted salience heatmap.

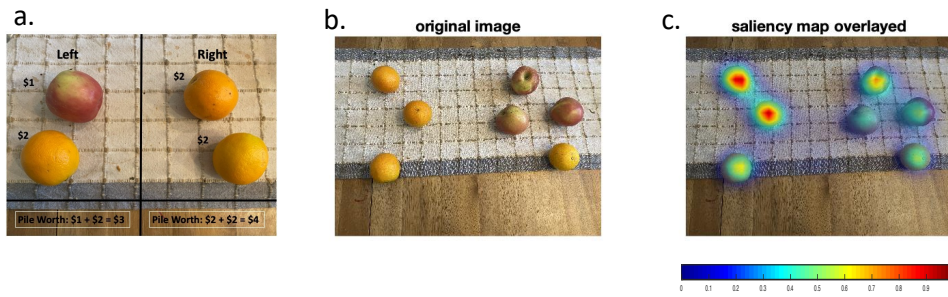


FIGURE III. FRUIT EXPERIMENT IMAGES

*Note:* a) Illustrates the rules of this task. Each fruit was worth a certain amount of dollars. The value of a set was the sum of all fruit values in that set. b) Presents a sample image of an actual trial in this task, as subjects saw it (the dollar values were not shown). c) Presents the SAM saliency map for the sample image in b). The left set was more salient than the right set in this example. All images used in this task had a saliency distribution similar to this example, in that the saliency peak is only distributed in one of the two sets. At the saliency peak, the value of saliency was 1 (the peak is located in the middle orange of the left set). In test images, the difference between the left and right saliency peaks had an average difference of 0.23.

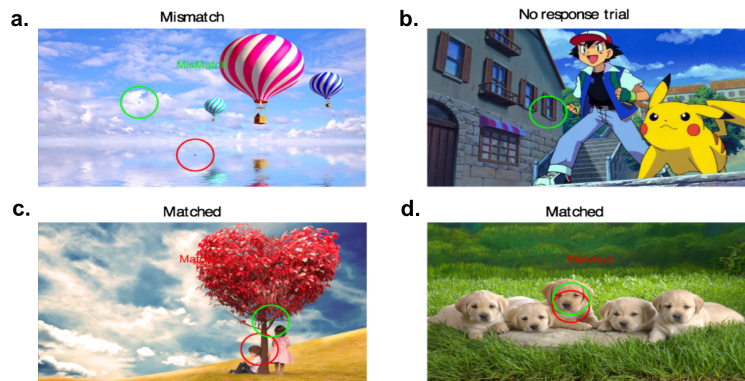


FIGURE IV. EXAMPLES OF TRIAL OUTCOMES WITH FEEDBACK, SHOWING CIRCLED PIXEL CHOICES

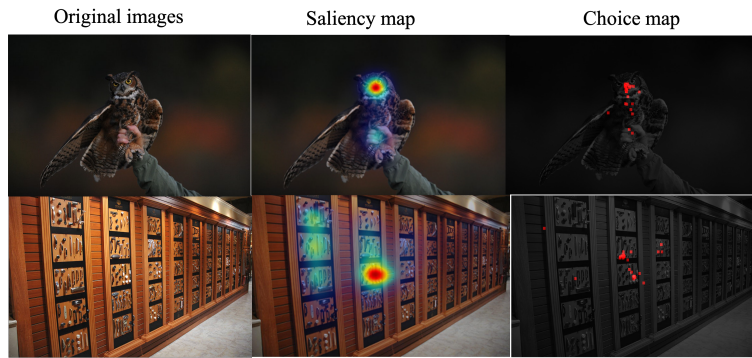


FIGURE V. TWO MATCHING GAME IMAGES, SALIENCE HEATMAPS, AND CHOICES (RED)

*Note:* (Left column) The original images. (Middle column) The original images overlaid by the SAM saliency maps. (Right column) The grayscale original image overlaid with the actual empirical choice distributions (each red dot is one choice).

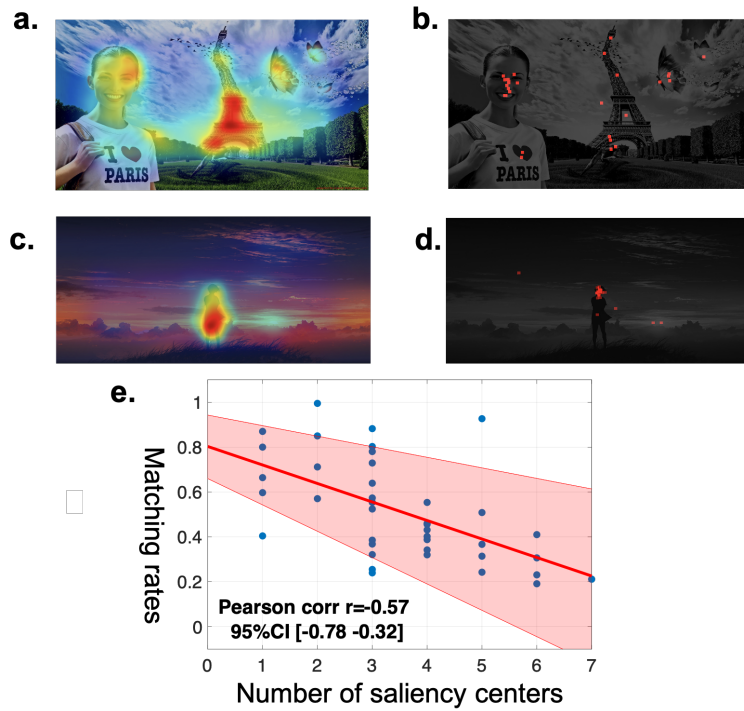


FIGURE VI. CORRELATION ACROSS IMAGES BETWEEN MATCHING RATE AND NUMBER OF SALIENCE CENTERS

*Note:* (a) is an image with seven saliency centers; (c) is an image with one saliency center. (b,d) are corresponding maps (red dots) of actual choice data in each matching game. The choice map in (b) is more dispersed because the saliency centers in (a) are more numerous. (e) plots the correlation between the number of saliency centers and the matching rate using both the feedback session and the no-feedback session to get a larger image pool ( $N = 40$  images).

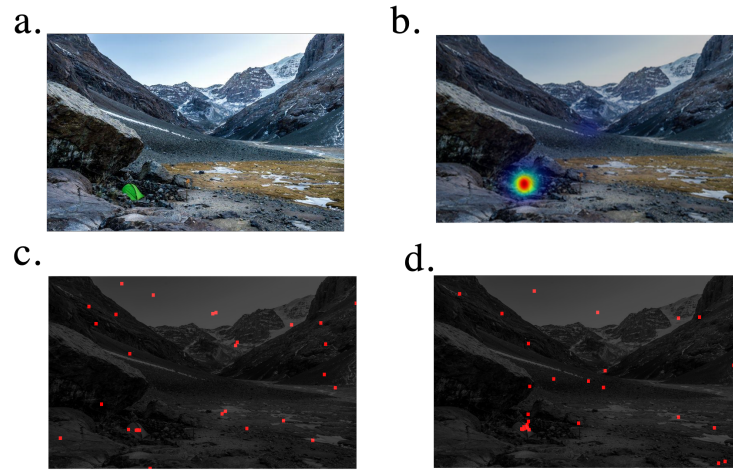


FIGURE VII. A HIDER-SEEKER GAME IMAGE, SALIENCE MAP, AND CHOICES

*Note:* a: The original image. b: The original image overlaid by the saliency map. c,d: The grayscale original image overlaid with the actual empirical choice distributions (each red dot represents an actual choice from one person). c is for hider choices and d is for seeker choices.



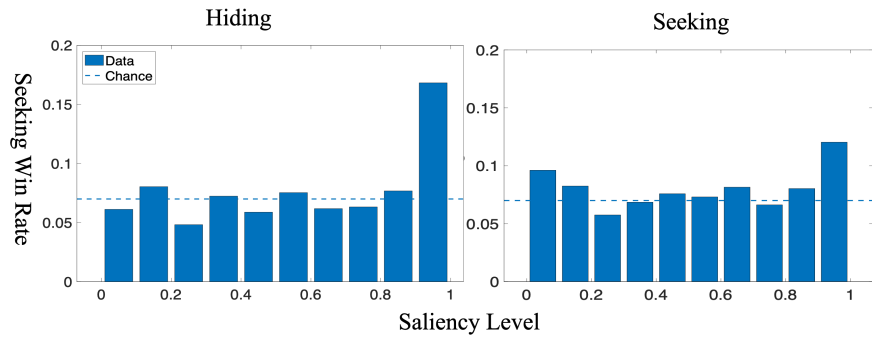


FIGURE VIII. SEEKING WIN RATES AS A FUNCTION OF DIFFERENT SALIENCY LEVELS

*Note:* This figure shows the average seeking win rate of hiders and seekers separately, at each saliency level bin from 0 to 1 (with bin size 0.1). This conditional seeking win rate looks a little different between hiders and seekers mainly because their choices are distributed differently across saliency levels.

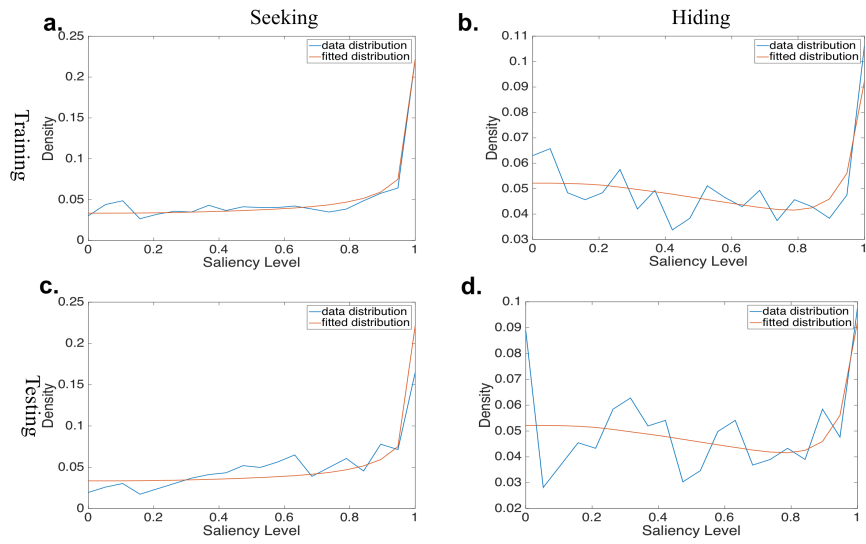


FIGURE IX. FREQUENCY OF CHOICE BY SALIENCE LEVEL WITH MODEL FITTED DISTRIBUTIONS

*Note:* The graphs indicate what percentage of choices were made for locations with the saliency of those locations on the x-axis. a: Choice data and model prediction in the training dataset seeking condition. b: Choice data and model prediction in the training dataset hiding condition. c: Choice data and model prediction in the testing dataset seeking condition. d: Choice data and model prediction in the testing dataset hiding condition.

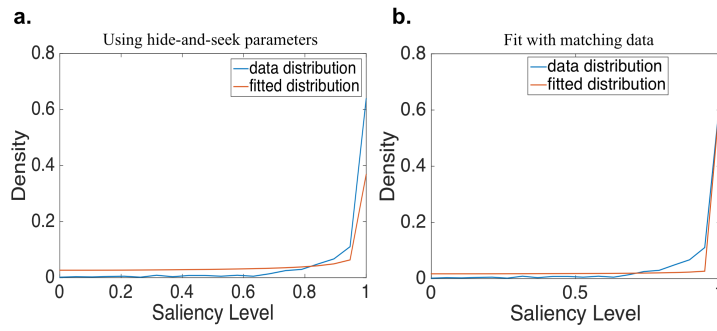


FIGURE X. THE SCH MODEL CALIBRATED ON HIDER-SEEKER GAME DATA CAN PREDICT MATCHING GAME CHOICES.

*Note:* The comparison between the matching data distribution and the two fitted matching game distributions. (a) Parameter estimates from the hider-seeker game are used to predict matching game results. log-likelihood: -2176 (b) Parameter estimates from the training matching game data are used to predict test matching game data. log-likelihood: -1943

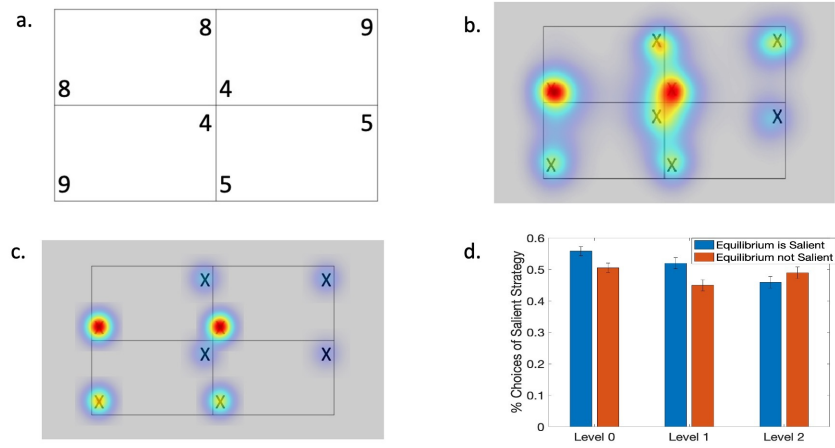


FIGURE XI. SALIENCY AND CHOICES IN MATRIX GAMES

*Note:* a) One example (Prisoners' Dilemma) of the games used in the experiment. b) The average SAM prediction of all games. c) The ground truth gaze density map generated by gaze data. d) Percentage of choices choosing the most SAM salient strategy grouped by levels (strategic thinking levels classified by gaze and behavior data by Polonio et al). (N: level 0 =551, level 1 =402, level 2 =371) Source: Polonio et al. (2015)

## Appendix

### A. HISTORY AND DETAILS OF SALIENCY ALGORITHMS

The SAM algorithm takes one image as an input and outputs its predicted saliency map. The saliency map is a saliency value from zero to one (least salient to most salient) assigned to each pixel on an image. Figure I in the text is a specific example of the SAM saliency map from one of the pictures we used in our experiments. A little history of saliency mapping may be useful here to convey

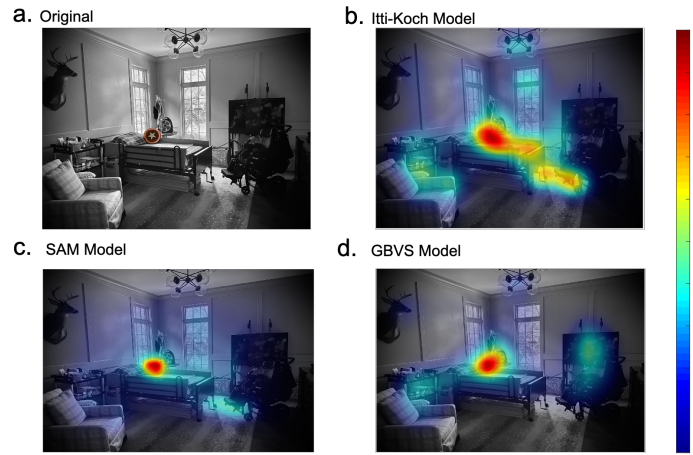


FIGURE A1. COMPARISONS BETWEEN THREE SALIENCY MODELS

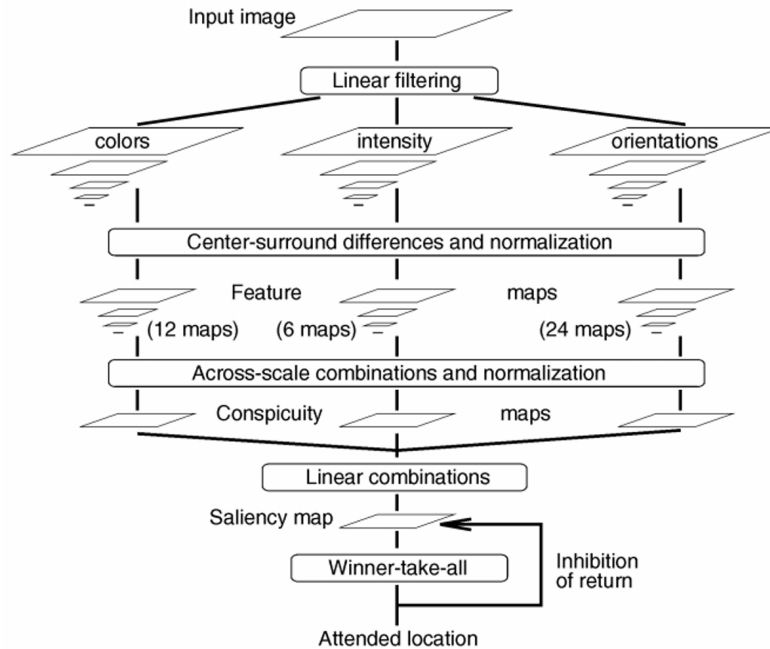
*Note:* We show the result of different saliency model on the same image (a). Both the Itti-Koch model and GBVS are fully interpretable (shown in b and d). Also in this example GBVS (Harel et al., 2007) and SAM have very similar results. All three outputs are color-plot in the same standard colormap function in matlab (type: “jet”).

how well-founded these algorithms are.

Inspired by a deep understanding of how the human visual system prioritizes attention, a series of progressively improving algorithms were developed to use visual images as inputs, and output predictions about where people will look in the first 1-2 sec of processing (Itti et al., 1998; Harel et al., 2007; Judd et al., 2009). Figure A1 compares two early algorithms and SAM in one example image.

Figure A2 shows an early algorithm from Itti et al. (1998). These early algorithms used a combination of handcrafted features to extract information about contrast, color, and orientation. Dark-light contrast is special because it marks boundaries between objects. Color and orientation are also thought to have adaptive value in parsing images in ways that are ecologically useful.

FIGURE A2. ITTI-KOCH MODEL

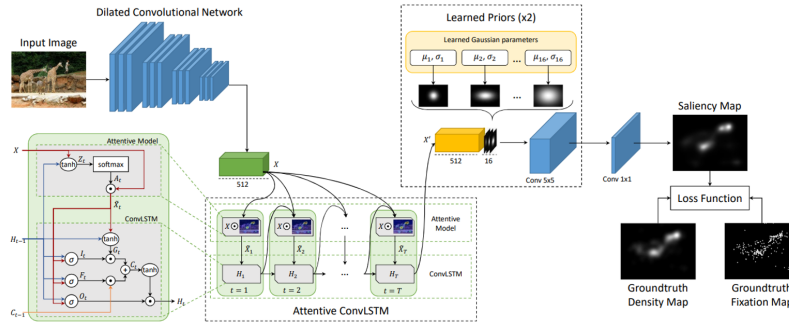


*Note:* presents the model of Itti-Koch (1998).

Consider the stick figure “I”. The bottom-up perception is a black vertical line of a certain length, with slightly extended top and bottom horizontal lines on top of the vertical line, surrounded by contrast with a white background.

A more abstract high-level concept of bottom-up is that (depending on the subject population) the I may be more familiar, valuable, or behaviorally useful. An English speaker will perceive “I” as a marker of first-person communication; a student just learning Roman numerals will perceive “I” as the number one;

FIGURE A3. SAM MODEL



*Note:* presents the deep neural network model structure of the state-of-art SAM model framework.

and an architecture aficionado may perceive it as an Iconic column, a part of a building. All the latter forms of salience use semantic knowledge— which is local and acculturated— about the world to inform the perception of what “I” means and what to do with that information. Which features or objects are personally relevant, valued, familiar, and novel will also be trained into a bottom-up algorithm, but will depend on the image set and characteristics of the subject population.<sup>64</sup>

The SAM algorithm (see Figure A3) we use was tuned using human free gaze data on a large number of images, without any special goals or incentives. The subject are just told to look. These algorithms were **not** designed to predict active choices in games with specific goals, such as matching and hider-seeker games. The matching goal, for instance, is to choose a location another person is also likely to choose. This is a goal-directed influence on perception which is likely to produce visual fixations that are different from free viewing. Thus, the extent to which SAM can predict the influence of predicted salience is probably a lower bound on how well better models, incorporating top-down goals, will do. To help clarify, the origins of SAM and previous algorithms, Table A1 describes

<sup>64</sup>See studies on the effects on perception of recent choice history (Awh et al., 2012), familiarity and novelty (Itti and Baldi, 2009), value for consumer goods (Towal et al., 2013), and self-reported “meaning” (Henderson and Hayes, 2017).

TABLE A1—SUBJECTS AND SOURCE IMAGE INFORMATION OF SAM TRAINING

Dataset	Subjects restriction	Image Categories
Salicon	Age: 19-28	91 categories including persons, vehicle, outdoor, animal, and etc.
MIT1003	Age: 18-35, We guess from Boston/MIT area. <sup>65</sup>	natural indoor and outdoor scenes
MIT300	Age: 18-50 We guess from Boston/MIT area	natural indoor and outdoor scenes
CAT2000	age 18-27, observers were undergraduates at USC from different majors and from mixed ethnicities.	Action, Affective, Art, Black & White, Cartoon, Fractal, Indoor, Inverted, Jumbled, Line Drawing, Low Resolution, Noisy, Object, Outdoor Man-made, Outdoor Natural, Pattern, Random, Satellite, Sketch, and Social.

the sets of images used and some characteristics of the subjects whose free gaze data were used to train SAM. The original papers are not crystal clear on who the subjects were, which is an indication that the authors think of the perceptual processes they are studying as rather homogeneous across people.

## B. FOCALITY IN PREVIOUS GAME EXPERIMENTS

There is a substantial, interesting series of experimental studies about focality in matching games. These studies are quite different from our approach but are described here for completeness.

There was a long lag between Schelling’s early 1960 discussion and later bursts of careful experimentation on focality.

Mehta et al. (1994b) proposed an important contrast between “secondary salience” and “Schelling salience”. Following Lewis (1969, pp. 24-36), they suggested that when players are not sure what to choose, they choose according to “primary salience”, which is “some (possibly stochastic) process that brings one of the labels to the player’s mind” (p. 660). *Secondary* salience is the belief about what creates primary salience for *others*. This process can obviously be iterated further.

Their experiments supported this distinction. In “picking” conditions people just picked an object from a choice set (e.g., a set of flowers). In “matching”



conditions their choices were matched with randomly chosen others and rewarded if they matched. The hypothesis is that picking measures primary salience and matching measures secondary salience. Indeed, the most common modal choice in the picking condition was usually chosen much more often when matching.

Note that this primary-secondary distinction is instantiated naturally in the SCH model (although that model was developed to explain behavior in a much wider range of games). In SCH, the process that brings one of the labels to the player’s mind— its primary salience— is predicted *ex ante* from the bottom-up SAM model. In the Mehta et al. (1994b) paradigms primary salience has to be *measured* by having people choose objects in the picking condition. Using SAM a primary salience prediction is delivered for all images; no new data or free parameters are needed.

In contrast to primary and secondary salience, an object has “Schelling salience” if it is unique or is chosen by a rule that leads to unambiguous results, which improves matching. Schelling salience need not arise from primary or secondary salience. For example, in a list of historical figures including Adolf Hitler, Hitler could be Schelling-salient even though few people would pick Hitler (primary salience) or think others would pick Hitler (secondary salience). Indeed, Mehta et al. (1994b,a) find evidence for both secondary and Schelling salience in their data.

More ambitiously, Bacharach (1993); Bacharach and Bernasconi (1997) proposed general principles underlying focality in matching choices from sets of objects, essentially trying to unpack Schelling salience into component parts. Their idea was that if people know their goal is coordination, they will try to naturally categorize objects into subsets and chose from more distinctive— e.g., smaller— subsets. However, subjects’ actual choices were not always consistent with the most non-obvious of their principles. Their experiments are elegant and careful. They were held back by the fact that a key element of the theory— “noticing” set-theoretic features— is measured only crudely (by self-report), whereas we now

have eyetracking to measure noticing directly.

Focality is also likely to work differently in hider-seeker games (HS). Studies by Mehta et al. (1994b) Bacharach (1993); Bacharach and Bernasconi (1997) were focused on coordination; at that time in the research history, there was no ambition to create theories of focality that would span games of different competitive structures. Understanding matching was difficult enough.

In a separate strand of cumulated regularity, an early study by Rubinstein et al. (1997) (RTH) used a four-choice hider-seeker game. Their canonical example is a choice between four letters ordered from left-to-right, where one letter is a singleton subset, like so:

A B A A

RTH hypothesized that the left and right A letters are avoided (because of “extremity-aversion”; cf. Bar-Hillel 2015). They hypothesize that the single B is clearly focal because it is both visually and semantically unique, and it will therefore be avoided by hidere. That leaves the second “interior” A from the right, which is least focal when compared to other choices (and therefore uniquely non-focal, giving it an ironic strategic focality due to uniqueness).

In these early studies, extremity-aversion and B-focality are simply hypothesized intuitions; they were not guided by data or visual perception principles. On this basis, RTH predicted that the third A would be chosen most often. Indeed, in their experiments, the third A is chosen most frequently both by hidere (40%) and seekers (45%). As a result, there is a “seeker advantage” because the seekers win more often than Nash equilibrium prediction of 25%. However, our replications in Caltech and UCLA subjects found much lower rates of the choice of the third “inner” A, around 29%, closer to the Nash 25% prediction (unpublished data).

Falk et al. (2009) used visual hider-seeker games similar to the four-letter choice. One game required choosing 3 cells out of the 25 locations in a 5x5 matrix. They

observe both an edge aversion and a seeker advantage.<sup>66</sup> There is a lot of other interesting data and psychology in their paper. In a recent study Brocas and Carrillo (2021) targeting at development and social choices, seeker advantages in hider-seeker games were discovered early in the life stages (second grade kids), and were still present in the adolescents control group (with effect size increasing with age).

In the main text we noted that our modeling builds upon Crawford and Iriberri (2007a) (hereafter CI), they advanced a novel analysis of games like ABAA, based on level-k modeling. They hypothesized that behavior could be consistent with a level-k approach<sup>67</sup>, in which level-0 behavior is influenced by salience. Rather than using an algorithm to predict salience, salience is parameterized by the frequencies of the outer A's and the central A. CI also assumed that level-k types only best respond to level k-1 types and that the population didn't contain any actual level zero types. Under this framework, they estimated both level zero players' preferences towards different options (saliency biases) and population frequencies of level types. The general approach fits behavior well. Our paper expands on this approach by predicting saliency independently of choice, using no new data, in location games.

Hargreaves Heap et al. (2014) questioned the strength of the CI conclusions on the grounds that the salience of the extreme A's and the central A were estimated parametrically and not constrained across game structures. They created choice sets with a single "oddity" that is visually or semantically unique (e.g. a list of words which are all diseases plus the word "fitness".) They test whether the oddity is equally salient for level 0 players in three types of games— coordination (matching), discoordination (players win if they both choose something different), and hider-seeker. They reject the hypothesis that level 0 salience is the same across games. Crawford (2014) commented on their paper.

<sup>66</sup>Based on data reported in their paper, the seeking win rate in this experiment is 10.37% while the chance level is only 6.25%, implying a seeker advantage of +4.12%. These numbers are rather close to our own, which are about 7% and 9%, although the paradigms differ a lot.

<sup>67</sup>See Stahl II and Wilson 1994; Nagel 1995 and see Crawford et al. (2013) for a thorough review.

We find better “portability” of saliency across matching and hider-seeker games. Specifically, we are able to predict the saliency-choice correspondence in matching games from SCH hider-seeker estimation.

### C. RESULTS FROM NO-FEEDBACK TRIALS

The realized matching rates when there is no feedback are as follows. Note that the hider-seeker matching rate (9%) is the same as in trials with feedback:

TABLE C1—REALIZED MATCHING RATE

	No feedback	Number of observations
Nash mixed prediction	0.071	
Matching game	0.35(0.004)	550
Hider-seeker game	0.09(0.002)	523(s)+527(h)

*Note:* Statistical tests are against the null hypothesis that the seeker win rate is the baseline level and choices are independently and identically distributed across subjects (which is the Nash benchmark prediction).

Figure C1 below is a quantile to quantile plot, plotting the percentage rank of saliency for each location against the percentage rank of choice frequencies for those locations in matching games. To get the Q-Q plot, we first mapped all users’ choice data (not only click points, but all points which fell into the circle) onto a one-dimensional saliency value, normalized from zero to one. (The highest saliency point in each entire image is one, and the lowest is zero). Then we ranked all these realized saliency values for all choices in the targeted sub-block. We also transformed the rank of the choice frequencies across all subjects into rank percentages. We plotted the normalized saliency value, which was also the percentage of saliency, against the percentage of points chosen with the same saliency ranking. The Q-Q plot below shows that all quantiles of choice data are above the same quantiles of saliency level, and hence above the diagonal dashed line that would result if people were choosing independently of saliency.

Figure C2 presents both Q-Q plots and density maps in the hider-seeker game. Figures C2 a-b indicate that seekers’ choices are more biased towards salient

locations than hiders' choices are, and both are much less saliency-biased than in the matching games (recall Figure C1). Keep in mind, however, that the hiders should be choosing locations as low in saliency as they can perceive (i.e., a best-response Q-Q curve would be underneath the 45-degree identity line).

The density maps in Figures C2 c-d take every location in every game, and assign each one a saliency level (0-1 normalized within each image), and computes the frequency with which “strategies” (=locations) were chosen across all games and subjects. For hider-seeker games, these should be flat horizontal lines in equilibrium(except for sampling error). However, there are a disproportionate number of choices of high-saliency locations (that is, the densities turn up sharply at the right end of the scale). Seekers choose the highest-saliency locations about three times as often, and hiders choose them about two times as often. There is a slightly disproportionate tendency to choose the lowest saliency locations (near zero at the left end of the scale), especially for hiders.

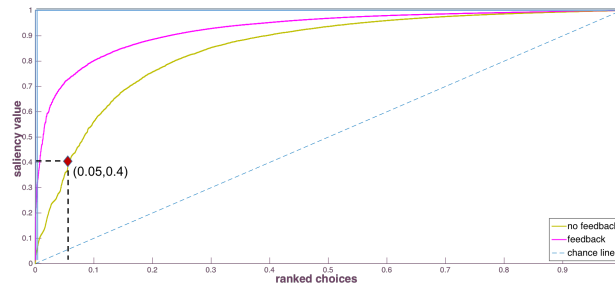


FIGURE C1. MATCHING GAME Q-Q PLOT OF CHOICE FREQUENCY(X-AXIS) AND SALIENCY RANKS (Y-AXIS)

*Note:* The red-diamond point (0.05,0.4) indicates that only 5 percent of choice points were made at the locations at or below 40% saliency. Equivalently, 95% of the points fall within the top 60% most salient points. Choices generated by chance thus correspond to a diagonal line of this plot from (0,0) to (1,1). The maximal accuracy is the blue line:  $y = 1$  for all  $x > 0$ , which would occur only if all choices fall on exactly the most salient point.

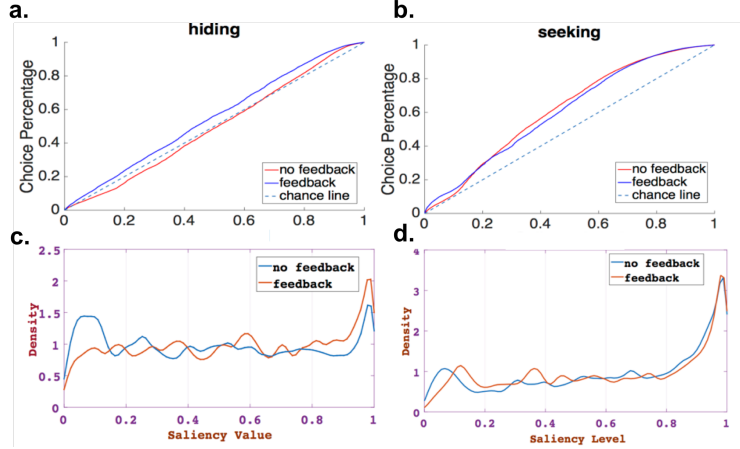


FIGURE C2. HIDER-SEEKER GAME Q-Q PLOT OF CHOICE FREQUENCY(X-AXIS) AND SALIENCY RANKS (Y-AXIS)

*Note:* a, b: Q-Q plots for hiding role (a) and seeking role (b). c, d: Kernel pdf density map of the choice frequency as a function of location saliency ranks. The x- axis is the rank of the saliency values and the y-axis is the probability density. Note: The kernel is Gaussian. The bandwidth is calculated using the formula:  $\sigma \times \frac{4}{3N}^{0.2}$ , in which  $\sigma$  is the standard deviation of the samples and  $N$  is the number of observations.

#### D. SCH MODEL COMPARISON WITH DIFFERENT SPECIFICATIONS

In this subsection, we are going to compare four different sub-models. We choose the Bayesian information criterion (BIC) to be the criterion for model selecting, since it balances the goodness of fit and the possibility of overfitting.

In all cases we restricted the softmax sensitivity parameter  $\lambda$  from 0 to 100. Larger values carry little extra information since  $\lambda = 100$  is close to best response. Constraining  $\lambda$  also makes it easier to create a bootstrapped confidence interval, which is useful due to the non-smoothness of the target function (likelihood function).

Here are descriptions of models we are going to test (and see Table D1):

- Model 1: There are only two types of players: 1) naive players who play as level zero players described in the main text. 2) equilibrium players who do pure randomization. Both the proportion of naive players,  $p_s$  and  $p_h$  serve as free parameters.

- Model 2: There is no level zero player in the real population, but higher level types believe there is. The hidiers and the seekers have different  $\tau$ s but the same  $\mu$  and  $\lambda$ .
- Model 3: Same as model 2, except that level zero players exist both in the belief structure and in the population.
- Model 4: This model fits hiding data and seeking data separately using two sets of parameters. Each game has three parameters:  $\mu$ ,  $\lambda$ , and  $\tau$ . The best fit model of it dominates the best fit of model 3 since model 3 is a special case of model 4. However, model 4 allows more free parameters, which the BIC value will penalize.
- Model 5: This model fits hiding data and seeking data using a common  $\mu$ ,  $\lambda$ , but uses the level-k belief framework<sup>68</sup> rather than CH, assuming the population consists of players whose level ranging from one to four (no level 0's).
- Model 6: This is the same as model 5, except it allows level 0 types.

Table D1 lists the best fit results of each model. Both BIC and AIC indicate that model 3 is the best performing model. Model 2 performs worst for the reason that without level zero types, the model structure will over predict the frequency of pure anti-salient hidiers, which is not seen in the data..

The Level-k Model 6 is almost as accurate by AIC and BIC, and we commented on what can be learned from it in the text. Figure D1 plots predictions of that model and the data, for comparison to Figure IX.

#### E. WORD LIST VERSION OF MAP COORDINATION TASK

One way to see how important visual saliency is for coordination is to test how people behave when they face the coordination problem in a non-image en-

<sup>68</sup>See Nagel, 1995; Crawford and Iriberry, 2007a,b

TABLE D1—MODEL COMPARISONS FOR HIDER-SEEKER GAME

Model	Description	Free parameters (Estimated)	AIC	BIC
1	Level 0+equilibrium	$p_s, p_h$ [1, .3]	12716	12728
2	Role-specific $\tau_x, f(0)=0$	$\mu, \lambda, \tau_s, \tau_h$ [.004, 99, .46, .002]	12780	12803
<b>3</b>	Role-specific $\tau_x, f(0) \neq 0$	$\mu, \lambda, \tau_s, \tau_h$ <b>[.06, 100, .40, .10]</b>	<b>12650</b>	<b>12673</b>
4	Role-specific $\tau_x, \mu_x, \lambda_x$	$\mu_s, \lambda_s, \tau_s, \mu_h, \lambda_h, \tau_h$ [.01, 90, .40, .07, 90, .50]	12646	12680
5	Level-k role-specific $f(k), f(0)=0$	$\mu, \lambda, f_s(1), f_s(2), f_s(3)$ $f_h(1), f_h(2), f_h(3)$ [1, 99, .22, 0, .78, .83, .05, .12]	12681	12738
<b>6</b>	Level-k role-specific $f(k), f(0) \neq 0$	$\mu, \lambda, f_s(0), f_s(1), f_s(2), f_s(3),$ $f_h(0), f_h(1), f_h(2), f_h(3)$ <b>[.18, 99, .29, .05, 0,</b> <b>.66, .17, .22, .61, 0]</b>	<b>12652</b>	<b>12709</b>

*Note:* Each model in the table is specified in the text list. BIC is defined as  $-2 \cdot \log L + \text{numParam} \cdot \log(\text{numObs})$  and AIC is  $-2 \cdot \log L + 2 \cdot \text{numParam}$

vironment. The SAM algorithm does not apply to such a scenario. We tried a non-visual version of Schelling’s location game, in which subjects were asked to coordinate on ten locations described in Figure IIa but only in a word list. The options were: house at the bottom of the map, bridge, small house near the pond, house at the top of the map, pond, two houses together, creek, fork in the road, X on the map, and Y on the map. The questions were presented in a randomized order. N=37 people participated the survey on Prolific. Each of them only answered the question once. Most of them choose the option “x on the map” (49%) while none of them chooses “y on the map” (see table E1). Only 5% people choose the bridge, which was the most popular option when the question was presented in an image format.



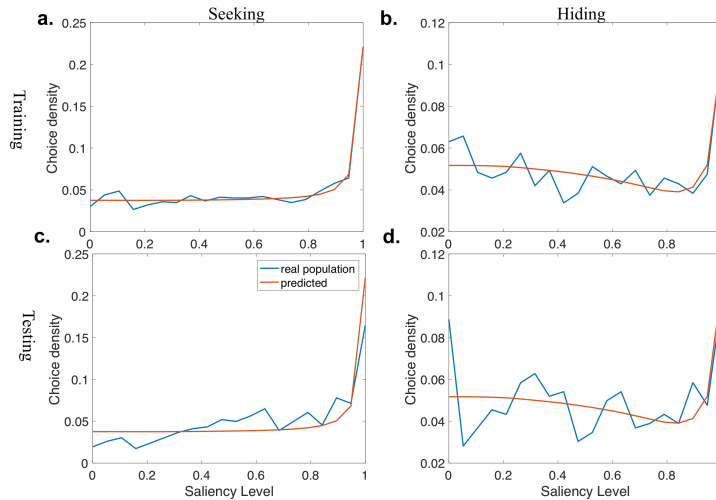


FIGURE D1. LEVEL-K MODEL 6 TRAINING-TESTING COMPARISON

*Note:* The x-axis is the saliency values of all click points. Each point on a graph indicated what percentage of choices were made for locations within images based on the saliency of those locations. (a): Choice data and model prediction in the training dataset seeking condition. (b): Choice data and model prediction in the training dataset hiding condition. (c): Choice data and model prediction in the testing dataset seeking condition. (d): Choice data and model prediction in the testing dataset hiding condition. Can be compared to Figure IX in the text.

## F. FRUIT EXPERIMENT

### F1. Fruit Experiment -Data

$N = 75$  participants took part in this study on Prolific, a European online data collection platform, following a pre-registration process on the Open Science Foundation website (OSF). All the participants were pre-screened to have a prior approval rate of at least 70% based on their previous participation. Each subject was only allowed to participate once for all types of batches (including pilot studies). Participation from mobiles and tablets were not allowed in order to control for attention effects.

The experiment design in timeline is shown in Figure F1. Subjects first read instructions freely until they fully understood. They were then asked to answer five comprehension questions as a check. Subjects who made more than one mistake are excluded. Then they played a session with unlimited time to get

TABLE E1—THE CHOICE PERCENTAGE OF ALL CHOICES IN FIGURE II

	Percentage
X on the map	0.49
House at the bottom of the map	0.08
Bridge	0.05
Small house near the pond	0.14
House at the top of the map	0.08
Pond	0.03
Two houses together	0.05
Creek	0.05
Fork in the road	0.03
Y on the map	0.00

*Note:* This table represents the percentage of people (N=37) playing the map game based on a list of verbal description rather than the visual map in text Figure II. Each participants played once.

familiarized with the rules (this part was incentivised also but was only for training purposes, as is not counted in the reported dataset). After that, they entered the main task session, where they would encounter 20 new images in a randomized order.

### *F2. Stimuli Properties and Selection Mechanism*

We took 72 photos of different combinations of real fruits displayed on a dining table. In the text, Figure III showed examples of SAM predictions. Each image contains two sets of fruits and each set contains three to five fruits. We flipped all the images in the horizontal direction so that we got another 72 images with the same content, but with the set locations flipped horizontally.<sup>69</sup> These 144 images are our image pool.

We selected 20 images from the image pool and all of the selected images satisfy four conditions:

- 1) **One-side salience centered** All of the selected images are strictly one-side salience centered, which means that the most salient locations only appear in one fruit set. Figure IIIc represents an example of a one-side

<sup>69</sup>This procedure is to avoid any left-right biases when taking images. It is done using a matlab function `flipimg()`.

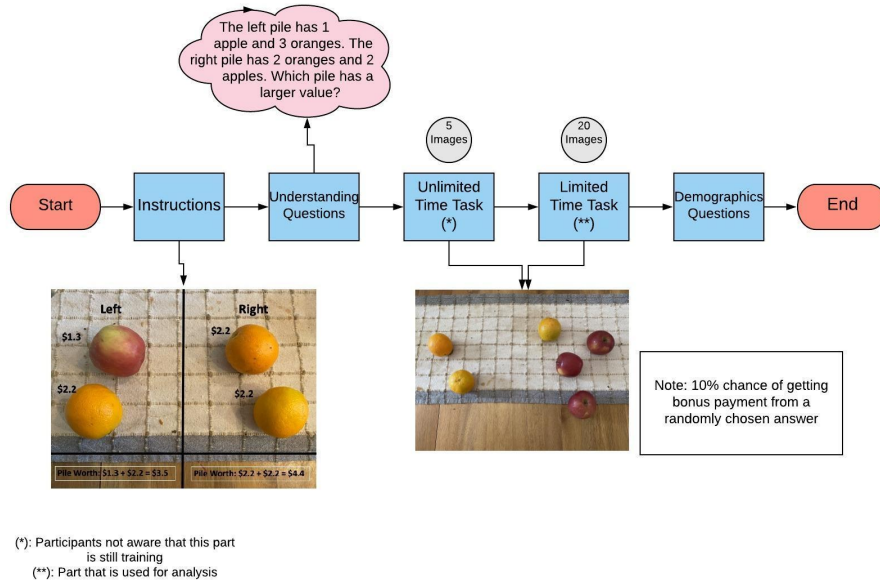


FIGURE F1. EXPERIMENT FLOW.

*Note:* Subjects first experienced an introduction session explaining the basic tasks followed by a testing session asking questions about the rules. They then experienced an unlimited time session with  $N=5$  images which will count into payment but only for training purposes. Afterwards, they will do a session with time limit. The experiment ends with demographic questions and payment realizations.

salient image, while Figure IIIa shows an image that is not one-side salience centered. Formally, consider two sets of pixels constituting the left set and the right set,  $P_l$  and  $P_r$ . Function  $s$ , the salience model, maps the union of  $P_l$  and  $P_r$  to  $[0, 1]$ . The most salient location of an image consists of a set of pixels  $S_h : \{x | s(x) > 0.99\}$ .<sup>70</sup> An image is one-side salience centered, if and only if exactly one of the two conditions hold true:  $S_h \cap P_l = \emptyset$  or  $S_h \cap P_r = \emptyset$ .

- 2) **Balanced salience center locations:** The selected image set has salience centers equally located on the left side or right side. Half of the images have salience centers on the left and the other half have them on the right.

<sup>70</sup>Since salience is a relative measure, there will always be at least one pixel with salience value one.

- 3) **Balanced valuation distribution:** There are only two types of fruits: oranges and apples. Each apple is worth 1.3 dollars and each orange is worth 2.2 dollars.<sup>71</sup> The total value differences between two sets range from 0.4 dollars to 4 dollars. There are exact 50% of rounds with the more rewarding option located on the left and 50% of rounds with the more rewarding option on the right.
- 4) **Balanced congruences:** An image will be called “congruent” if the more rewarding option is also the more salient option. Among all images, there are 50% congruent images and 50% incongruent images. No image contains two sets with the same amount of values.
- 5) **Balanced number of fruits:** Among the total 20 images, in 18 images have the number of fruits only differ by one. The other two images differ by two. 11 images have more fruits on the left and 9 images have more fruits on the right.

#### G. EXPERIMENTAL PROCEDURES OF LOCATION GAMES

**Screen1:** You are now going to do a series of short games. In each one, you will see a series of pictures and you must choose a location on the picture by clicking with the mouse.

The rules of each game are slightly different, so read them carefully before you start! (You cannot go back and reread them.)

**Screen 2:** You’ll start with a few practice items to help you get familiar with the basic set-up.

Use the mouse to click a location anywhere on picture. Notice that your selection is the entire area within the circle.

You will have 6 seconds to make your selection before the picture disappears. If you do not make a selection within 6 seconds, you will not get credit for that

<sup>71</sup>We did a pilot experiment with integer unit values. It turned out to be that integer values were too easy for the subjects so we didn’t see any variation in choice accuracy.

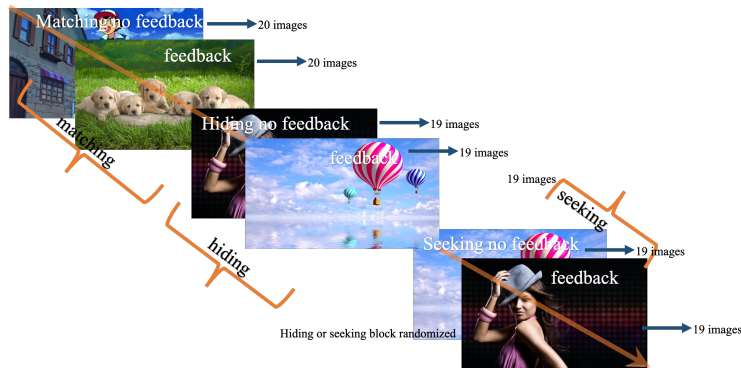


FIGURE G1. BLOCK DESIGN OF THE LOCATION GAME EXPERIMENT

*Note:* This figure shows the block design of the main experiment (location game). Each participants experienced matching game first, then hiding or seeking game in a randomized order. Under each game, there are two sub-block, the first one is always without feedback and the second one is with feedback.

picture.

**Screen before each session depending on games:** Matching: Now you are playing a matching game with several other research participants like you.

For each image, you will play against a randomly selected opponent. If you and your opponent choose the same location in the picture, you both win \$x.<sup>72</sup> If there is any intersection between your location and your opponent's location, it will count as a "match".

You won't find out how much you won in this phase until the end of the game.

As before, you will only have 6 seconds to make your choice for each image.

<sup>72</sup>The value of x changes with games, we pay \$0.2, \$0.1, \$0.4 for a success in matching, hiding and seeking.

TABLE G1—SUMMARY OF DATASETS

Type of datasets	Platform	Whether no-feedback <sup>73</sup>	N of subjects	Games	Time limit	Between vs Within
Main Batch	In lab	Yes	29	M,H,S	6s	Within
Main Batch Online	mTurk	Yes	38	M,H,S	6s	Within
Big Circle	mTurk	No	67	M,H,S	6s	Within
High Reward	mTurk	No	29	H,S	6s	Within
No time-limit	mTurk	No	49	M,H,S	Inf	Within
Between-Subject	mTurk	No	53	H,S	6s	Between
Time pressure	mTurk	Yes	31	M,H,S	2s	Within

*Note:* This table summarized seven different datasets collected at different times. Only the high reward group and the main batch group are the same group of participants. All other batches are completed by a new group of people. Repeated participation is not allowed in all other batches.

TABLE G2—LOCATION GAMES: DATASET USAGE SUMMARY

Analysis Names	Dataset Used	Observations
Seeking Win Rates (Seeker’s advantage)	The main results: in-lab dataset. Also reported this percentage for other robustness checks.	M:559,H:529,S:531 M:458,H:441,S:452 (Main Batch Online)
SCH model: training	In-lab dataset with both feedback group and no feedback group.	H:1096,S:1090
SCH model: testing	In-lab dataset, high reward group.	H=446,S=446
Choice saliency level analysis	In-lab dataset with both feedback and no feedback group (in footnote and appendix).	M:1139,H:1096,S:1090
Matching rate/Saliency center	In-lab dataset both feedback group and no-feedback group.	M:1139

*Note:* This table summarizes the dataset we used for each part of the analysis. We mainly and consistently use the dataset we collected in lab for all the analysis. For the seeker’s advantage part, we also tested different conditions for robustness checks. The “observation” column denotes the total number of observations under each game. M,S,H denotes for matching, seeking and hiding, separately. We omit all the missing data which happens rarely in the in lab sessions and more commonly in online sessions.

## H. MATRIX GAME

N=56 people played 32 normal form games with different strategic structures: Dominant Solvable Self (DSS), Dominant Solvable Other (DSO), Prisoner’s Dilemma (PD), and Stag Hunt (SH).

In the 32 games, 24 games contain a unique equilibrium (SH has two). Each player is either assigned as a row player or a column player. Both roles saw the original games without any transpose.

TABLE H1—EVALUATION OF SAM ON MATRIX GAME EXPERIMENT

	AUC	CC
SAM vs fixations(games)	0.96	0.47
Chance level	0.5	0
Range	(0,1)	(0,1)

*Note:*

The table reports two common evaluation metrics for the matrix game experiment in Section VI. It reports area under the receiver operating characteristics(AUC) and Pearson Correlation Coefficient (CC)(Kummerer et al., 2018) . We show SAM’s performance on human eye-fixations for matrix games. The results on both metrics show that SAM predicted human fixations far better than chance.

## I. ADDITIONAL ANALYSIS

TABLE II—SUMMARY OF ACRONYMS

<b>SAM</b>	The saliency model we used, Saliency Attentive Model
<b>CI</b>	Crawford and Iriberry (2007)
<b>ABAA</b>	Hider- seeker game using these four letters
<b>CH</b>	Cognitive hierarchy model
<b>SCH</b>	Saliency Cognitive hierarchy model
<b>BGS</b>	Bordalo, Gennaioli and Shleifer (BGS) saliency theory

*Note:* If readers have difficulty keeping track of all the acronyms, this table may help.