# Lift-the-flap: what, where and when for context reasoning

**Mengmi Zhang**[1]    **Claire Tseng**[2]    **Karla Montejo**[3]    **Joseph Kwon**[4]    **Gabriel Kreiman**[1]

[1] Boston Children's Hospital, Harvard Medical School
[2] Harvard College
[3] Mayo Clinic Graduate School of Biomedical Sciences
[4] Yale University
Address correspondence to `gabriel.kreiman@tch.harvard.edu`

## Abstract

Context reasoning is critical in a wide variety of applications where current inputs need to be interpreted in the light of previous experience and knowledge. Both spatial and temporal contextual information play a critical role in the domain of visual recognition. Here we investigate spatial constraints (what image features provide contextual information and where they are located), and temporal constraints (when different contextual cues matter) for visual recognition. The task is to reason about the scene context and infer what a target object hidden behind a flap is in a natural image. To tackle this problem, we first describe an online human psychophysics experiment recording active sampling via mouse clicks in lift-the-flap games and identify clicking patterns and features which are diagnostic for high contextual reasoning accuracy. As a proof of the usefulness of these clicking patterns and visual features, we extend a state-of-the-art recurrent model capable of attending to salient context regions, dynamically integrating useful information, making inferences, and predicting class label for the target object over multiple clicks. The proposed model achieves human-level contextual reasoning accuracy, shares human-like sampling behavior and learns interpretable features for contextual reasoning.

## 1   Introduction

The tiny object on the table is probably a spoon, not an elephant. Objects do not appear in isolation. Instead, they co-vary with other objects, their sizes and colors usually respect regularities with respect to nearby elements, and objects tend to appear at specific locations within a scene. Humans exploit these contextual associations. Contextual analyses based on the statistical summary of object relationships, provide an effective source of information for perceptual inference tasks, such as object detection ([39, 34, 21, 40, 30]), scene classification ([19, 41, 47]), semantic segmentation ([47]), and visual question answering ([38]).

An example of how contextual information is incorporated during object recognition is lift-the-flap books, where a flap covers part of the page. Children make guesses about what is behind the flap based on the context and check their answers by lifting the flap (Figure 1a). Here we investigate *what* image features matter for contextual reasoning and *where* those features are with respect to the target object of interest. Furthermore, scene interpretation in humans involves a sequence of eye movements [49], each one of these image samples providing additional context to inform interpretation of the contents of the next location. Therefore, we also investigate *when* scene information matters for context reasoning.
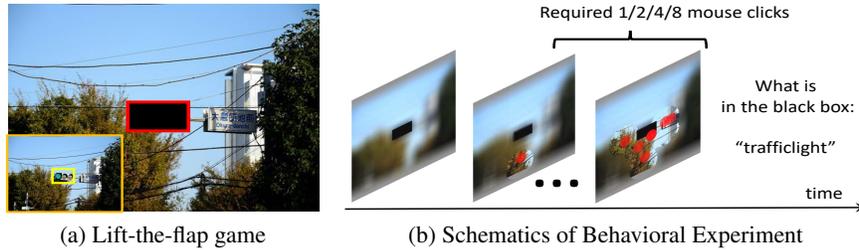
Preprint. Under review.

Required 1/2/4/8 mouse clicks

What is
in the black box:

"trafficlight"

time

(a) Lift-the-flap game

(b) Schematics of Behavioral Experiment

Figure 1: **Schematic of the lift-the-flap task and human behavioral experiment**. (a) The task requires subjects to capitalize on the scene context in a natural image to infer what is behind the black box (the hidden target). The original image (bottom left) reveals the target object ("traffic light"); this image was *not* shown in the actual experiments. (b) A blurred image with the hidden target was presented to the subject. To identify contextual areas of importance surrounding the hidden target, subjects used the computer mouse to click on image a pre-specified number of times (red dots). Upon clicking a certain location, a circle of fixed radius was revealed at high resolution. After the required number of clicks, subjects typed a noun describing the hidden target object. The experiments were conducted online using Amazon Mechanical Turk [42] on 100 subjects, 50 trials per subject. Figure 3 also shows results for a variation of the experiment conducted in the lab while tracking eye movements.

To tackle the problem of contextual reasoning, we introduce the lift-the-flap task and conduct online psychophysics experiments where subjects make mouse clicks while they explore important contextual cues to identify a hidden target (Figure 1b). We investigate the contextual reasoning strategies observed from human active sampling patterns. As a proof of concept, we propose a recurrent attention model (ClickNet), to automatically learn these contextual reasoning strategies. The model guides attention towards regions with informative context, decides where to sample the image, and makes inferences about the target behind the flap. The learnt sampling patterns and predicted class labels by ClickNet share remarkable similarities with human behavior.

## 2 Related Works

### 2.1 Role of Context in Human Vision

The cognitive science literature has shown that contextual information affects the efficiency of several visual processes [2, 6, 22, 4, 18, 1], such as object recognition [2], object detection [6, 22], visual working memory [16, 1] and visual search [20]. Objects appearing in a familiar background can be detected more accurately and processed more quickly than objects appearing in an incongruent scene. Here we focus on what visual features contribute to contextual reasoning, which parts of image regions attract humans' attention for making inferences, and the dynamic sequence of directed sampling needed for making inferences about a hidden target.

### 2.2 Role of Context in Computer Vision

Contextual reasoning about objects and relations is critical to machine vision. In fact, many object recognition studies using natural image datasets such as ImageNet, rely implicitly but strongly on contextual feature regularities [17, 8]. Several studies employ contextual information in order to improve object detection [34, 21, 40, 30]. The types of contexts can be exploited in the form of global scene context [40], ground plane estimation [34], geometric context [21], relative location [14], 3D layout [28], and spatial support and geographic information [15]. In [19, 47, 27], researchers proposed Conditional Random Field (CRF) models that reason jointly across multiple computer vision tasks in image labeling and scene classification. Additionally, [32] studies the role of context in both object detection and semantic segmentation tasks, demonstrating improved performance in both tasks compared to raw image features. Recently, several neural network architectures incorporating contextual information have been successfully applied in object priming [39], place and object recognition [44, 41], object detection [30], and visual question answering [38]. Here we focus on developing a biologically inspired computational model for contextual reasoning that can automatically and dynamically sample image regions of interest, integrating information in

2

space and time to make inferences about a hidden object. Additionally, we compare the model's performance against human behavior in the same task.

Several interesting approaches have combined graphical models with deep neural networks for structural inference, primarily in structured prediction tasks [31, 11, 10, 38, 23, 5, 45]. [23] designed a structured model to improve classification performance by leveraging relations among scenes, objects, and their attributes. A structured inference model is also used in [11, 13] to analyze relations in group activity recognition. Several works, like Structural-RNN [24] and Interaction Net [5], combine the power of spatiotemporal graphs and sequence learning, and evaluate the model from motion prediction to object interactions. These works assume full contextual information is available, while in our experiment we consider only partial contextual information that is sequentially revealed after a mouse click. [10] proposed DeepLab which inputs the response at the final layer of a deep neural network to a CRF model for semantic image segmentation. Subsequently, [35, 50] transformed the CRF model into a Recurrent Neural Network in an end-to-end fashion. Breaking away from this previous work where graph optimization is performed globally, our proposed model selects important visual features using an attention mechanism and integrates partial information over multiple steps, which is computationally more efficient and accurate in the current task (Section 5).

## 3 Lift-the-flap task

### 3.1 Human Behavioral Experiments

Subjects were presented with a natural image where one object was hidden by a rectangular black box and everything else was blurred. They were allowed a fixed number of mouse clicks between 1 and 8, each click revealed part of the image in high resolution. After the target number of clicks, they had to provide a single word to describe the object hidden behind the black box (Figure 1b). The clicking experiments were run on Amazon Mechanical Turk [42]. The stimulus set consisted of 573 images from the test set of the MSCOCO Dataset [29], spanning 80 object categories. This dataset has been widely used for object recognition and detection studies [29]. We constrained the stimulus set to have a uniform distribution of 6 - 8 target objects per category. To avoid any potential memory effects, subjects were only exposed to each image once. The trial presentation order was randomized.

### 3.2 Ground Truth Responses

In contrast to other experiments where subjects are forced to perform N-way categorization (e.g., [37]), here there were no constraints on what words subjects could use to describe the hidden object in the experiment. This probing mechanism was implemented for two reasons: first, it is difficult for humans to memorize 80 object classes in advance and there could be non-uniform memory effects impacting the results; second, we were concerned that presenting humans with an 80-choice question in each trial could introduce biases in their inference and decision processes.

We could not simply use the 80 category labels to evaluate performance because subjects could use other similar words or synonyms and we are interested in the context reasoning process rather than the subjects' language abilities. Therefore, to evaluate humans' performance, we separately collected a distribution of ground truth answers for each hidden target by presenting to 10 *other* subjects, who did not participate in the main task, the same set of images with the target objects highlighted by a bounding box (not hidden). During the lift-the-flap task, a response was considered to be correct if it matched any of the ground truth labels, allowing for plurals and misspellings.

### 3.3 Evaluation Metrics

We introduce several evaluation metrics to measure contextual reasoning accuracy and to compare the consistency of mouse clicking patterns between humans and computational models. We evaluated ClickNet on the MSCOCO Dataset using the typical classification accuracy measure. In Fig 5b, we report **top-1 classification accuracy** as a function of **context-object ratio**. The context-object ratio is defined as the total area of the image *excluding* the hidden target divided by the hidden target object size. For example, a context-object ratio of 1 implies that the size of the
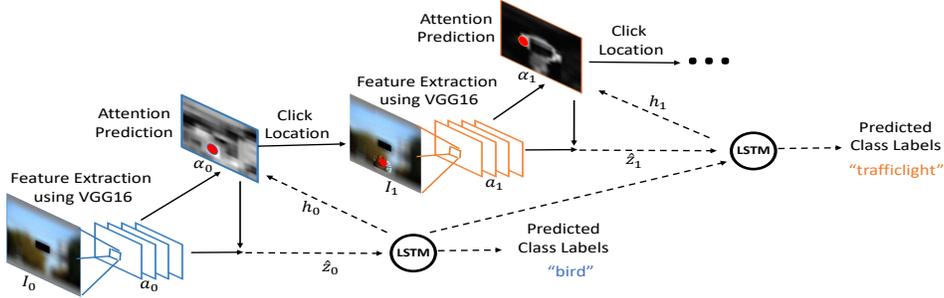
Figure 2: **Architecture overview of the ClickNet model**. The diagram depicts the iterative modular steps carried out by ClickNet for contextual reasoning over multiple clicks in the lift-the-flap task (Figure 1b). ClickNet consists of 3 main modules: feature extraction, attention, and recurrent memory. For illustrative purposes, only the first and second clicks in a trial are shown here. ClickNet performs feature extraction using the VGG16 network pre-trained on ImageNet and produces feature maps $a_0$. Conditioned on the hidden state $h_0$ and feature maps $a_0$, ClickNet produces an attention map $\alpha_0$, which is used to select the next click location (red dots) and to modulate the feature maps for contextual reasoning (Figure S1). The recurrent network in the LSTM module (Figure S2) integrates over time the attentionally modulated feature maps $\widehat{\mathbf{z}}_0$, and outputs a predicted class label after each click (here "bird" and "trafficlight" before the 1st and 2nd clicks). After the first click, the input image gets updated with parts at clicked locations revealed in high resolution. These three modular steps repeat until the specified number of clicks have been made. See supplementary figures S1 and S2 for implementation details on the attention and LSTM modules, respectively.

black box and the size of the contextual information is the same (see Figure 5b for example images with different context-object ratios).

To measure the degree of consistency between two mouse clicking patterns (human versus human; or human versus computational models), we computed the minimum Euclidean distance between the sequence of clicks in each trial, regardless of order. The smaller the median in the distribution of distances, the more similar the two mouse clicking patterns are.

## 4 ClickNet Architecture

We propose a recurrent neural network for context reasoning (ClickNet), extending previous work on image captioning [46]. ClickNet integrates attention-modulated context information over multiple clicks, makes a decision about the next click location based on the attention map, and infers the class label of the hidden target after every click (Figure 2).

As in the human psychophysics experiment (Figure 1b), ClickNet is first presented with a blurred image $I_0$, which is the original image $\mathbf{I}$ with uniform gaussian blur and where the target object is covered by a black box. ClickNet makes the first attempt to predict a class label $y_0$ out of a pre-defined set of $C$ object classes and decides its first click location $m_1$. In every trial, over a series of $T$ clicks, the input image $I_t$ to ClickNet gets updated with circular regions of constant radius $R$ centered at all previous click locations $M = \{m_1, ... m_T\}$, revealed in its original resolution in $\mathbf{I}$. The black box is constant and none of its content is ever revealed, even if the model opts to click within the box or if the circle centered on the click encompasses part of the box.

### 4.1 Convolutional Feature Extraction

At each time $t$ where $t \in \{0, ..., T\}$, ClickNet takes $I_t$ as input and uses a feed-forward convolutional neural network to extract feature maps $a_t$. We use the VGG16 network [36], pre-trained on ImageNet [12]. To focus on specific parts of the image and select features at those locations, we have to preserve the spatial organization of features; thus, ClickNet uses the output feature maps at the last convolution layer of VGG16.

A feature vector $\mathbf{a_{ti}}$ of dimension $D$ represents the part of the image $I_t$ at location $i$, where $i = 1, .., L$ and $L = W \times H$ and $W$ and $H$ are the width and height, respectively, of the feature map:

$$a_t = \{\mathbf{a_{t1}}, ..., \mathbf{a_{tL}}\}, \quad \mathbf{a_{ti}} \in \mathbb{R}^D \tag{1}$$

## 4.2 Attentional Modulation

We use a "soft-attention" mechanism as introduced by [3] to compute "the context gist" $\widehat{\mathbf{z_t}}$ on $I_t$ (Figure S1). For each location $i$ in $a_t$, the attention mechanism generates a positive scalar $\alpha_{ti}$, representing the relative importance of the feature vector $\mathbf{a_{ti}}$ for context reasoning. This relative importance $\alpha_{ti}$ depends on the feature vectors $\mathbf{a_{ti}}$, combined with the hidden state at the previous step $\mathbf{h_{t-1}}$ of a recurrent network described below:

$$e_{ti} = A_h \mathbf{h_{t-1}} + A_a \mathbf{a_{ti}}, \quad \alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{i=1}^{L} \exp(e_{ti})} \tag{2}$$

where $A_h \in \mathbb{R}^{1 \times n}$ and $A_a \in \mathbb{R}^{1 \times D}$ are weight matrices initialized randomly and to be learnt. Because not all attended regions might be useful for context reasoning, the soft attention module also predicts a gating vector $\beta_t$ from the previous hidden state $h_{t-1}$, such that $\beta_t$ determines how much the current observation contributes to the context vector at each location: $\beta_t = \sigma(W_\beta \mathbf{h_{t-1}})$, where $W_\beta \in \mathbb{R}^{L \times n}$ is a weight matrix and each element $\beta_{ti}$ in $\beta_t$ is a gating scalar at location $i$. As also noted by [46], $\beta_t$ helps put more emphasis on the salient objects in the images. Once the attention map $\alpha_t$ and the gating scale $\beta_t$ are computed, the model applies the "soft-attention" mechanism to compute $\widehat{\mathbf{z_t}}$ by summing over all the $L$ regions in the image:

$$\widehat{\mathbf{z_t}} = \sum_{i=1}^{L} \beta_{ti} \alpha_{ti} \mathbf{a_{ti}} \tag{3}$$

The next click location $m_{t+1}$ corresponded to the maximum on the attention map:

$$m_{t+1} = \arg\max_i \alpha_{ti} \tag{4}$$

The attentional module is smooth and differentiable and ClickNet can learn all the weight matrices in an end-to-end fashion via back-propagation.

## 4.3 Recurrent Connections using long short-term memory (LSTM)

We use a long short-term memory (LSTM) network to output a predicted class label $y_t$ based on the previous hidden state $\mathbf{h_{t-1}}$ and the context gist vector $\widehat{\mathbf{z_t}}$ for $I_t$ (Figure S2). Our implementation of LSTM closely follows [48] where $T_{s,t} : \mathbb{R}^s \to \mathbb{R}^t$ defines a linear transformation with learnable parameters. The variables $\mathbf{i_t}, \mathbf{f_t}, \mathbf{c_t}, \mathbf{o_t}, \mathbf{h_t}$ represent the input, forget, memory, output and hidden state of the LSTM respectively:

$$\begin{pmatrix} \mathbf{i_t} \\ \mathbf{f_t} \\ \mathbf{o_t} \\ \mathbf{g_t} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} T_{D+n,n} \left( \widehat{\mathbf{z_t}}, \mathbf{h_{t-1}} \right) \tag{5}$$

$$\mathbf{c_t} = \mathbf{f_t} \odot \mathbf{c_{t-1}} + \mathbf{i_t} \odot \mathbf{g_t}, \quad \mathbf{h_t} = \mathbf{o_t} \odot \tanh(\mathbf{c_t}) \tag{6}$$

where $n$ is the dimensionality of LSTM, $\sigma$ is the logistic sigmoid activation, and $\odot$ indicates element-wise multiplication.

To cue ClickNet about the location of the hidden target, we initialize the memory state $\mathbf{c_0}$ and hidden state $\mathbf{h_0}$ of the LSTM based on a binary mask that contains zeros everywhere and ones in the hidden target location. Specifically, $\mathbf{c_0}$ and $\mathbf{h_0}$ are predicted by an average of all feature vectors $a_0$ over all $L$ locations with two separate linear transformations $W_{c0} \in R^{n \times D}$ and $W_{h0} \in R^{n \times D}$:

$$\mathbf{c_0} = W_{c0} \left( \frac{1}{L} \sum_i^L \mathbf{a_{0i}} \right), \quad \mathbf{h_0} = W_{h0} \left( \frac{1}{L} \sum_i^L \mathbf{a_{0i}} \right) \tag{7}$$
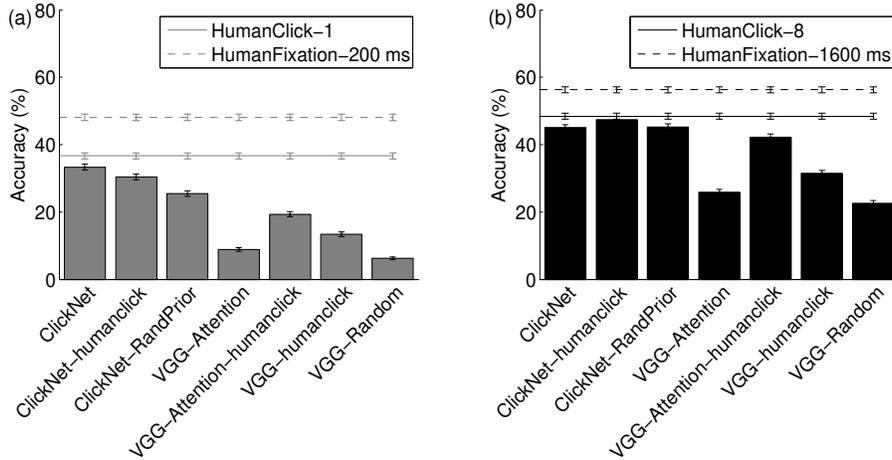
Figure 3: **Contextual reasoning accuracy of humans and models**. Performance for humans (horizontal lines) and models (bars) for **(a)** 1 click (gray), and **(b)** 8 clicks (black) (Fig S4 shows results for 2 and 4 clicks, and additional comparison models). Section 3.3 defines the evaluation metric and Section 4.5 describe each model. Error bars denote SEM across images.

To predict the class label $y_t$ of the hidden target, the LSTM computes a classification vector where each entry denotes a class probability given the hidden state:

$$y_t = \arg\max_c p(y_c), \quad p(y_c) \propto L_h \mathbf{h_t} \tag{8}$$

where $L_h \in \mathbb{R}^{C \times n}$ is a matrix of learnt parameters initialized randomly.

### 4.4 Training and Implementation Details

In training the model, we introduced a regularization to constrain $\sum_t \alpha_{ti} = 1$. This is to encourage ClickNet to acquire as much context information as possible by exploration. This regularization term was empirically important to improve the context reasoning accuracy. We trained ClickNet end-to-end by minimizing the cross entropy loss between the predicted label $y_t$ at each time step $t$ and the ground truth label $x$, and a regularization term for exploration:

$$LOSS = \sum_{t=1}^{T} (-\log(P(y_t|x))) + \lambda \sum_{i}^{L} (1 - \sum_{t}^{T} \alpha_{ti})^2 \tag{9}$$

We used all images from the MSCOCO training set for training and validating all models. On every training image, each object can be blocked out as the hidden target. The input image size to ClickNet was $400 \times 400$ pixels. We used a Gaussian filter of size $51 \times 51$ with variance 64 pixels to blur the images. The radius $R$ of the circular region revealed by each click was 55 pixels. The hyper-parameter for the regularization term in the loss functino was $\lambda = 2$. As in the human psychophysics experiments (Fig 1b), in each trial, we set the total number of time steps $T = 8$ for training ClickNet (ClickNet predicts the label after the 1st click at $T = 1$). The dimension of the LSTM module was $n = 512$. The feature maps extracted from the last convolution layer was of size $2048 \times 28 \times 28$, and the total number of locations was $L = 28 \times 28 = 784$. The Adam optimizer [25] was used with a learning rate of $10^{-4}$ to fine-tune the VGG16 network, and a learning rate of $4 \times 10^{-4}$ to train the attentional module and the LSTM module. The network was developed in Pytorch, based on [46]. All source code for our proposed architecture, and the data from the psychophysics experiments will be released publicly upon publication.

### 4.5 Variations of Proposed Network Architecture and Comparative Methods

Previous work has shown that it is possible to augment vision systems with human perceptual supervision on many difficult computer vision tasks, such as [43, 26]. One central goal in our study
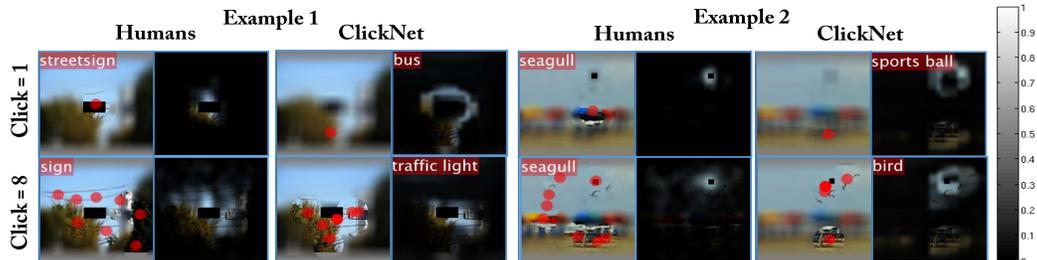
Figure 4: **Example visualization for humans and ClickNet**. Two example trials (first four columns is example 1, last four columns is example 2), either with 1 click (rows 1) or 8 clicks (rows 2), for one human (columns 1, 2, 5, 6) or ClickNet (column 3, 4, 7, 8) (Fig S3 in Appendix shows results for 2 and 4 clicks). Red dots denote clicked locations. Top-left corner shows output labels after the required number of clicks. Column 2 and 6 show the mouse click maps aggregated over subjects. Brighter regions denote more mouse clicks (see scale bar on right). Column 4 and 8 show the attentional map predicted by ClickNet. Brighter regions denote higher attentional values.

is to investigate what, where, and when matter for human contextual reasoning in the lift-the-flap game and whether these factors could help improve current machine learning algorithms. We now introduce two variations of ClickNet with human inputs at the **testing** stage:

**ClickNet-humanclick.** Instead of clicking at the location with highest activation value on the attention map predicted by ClickNet, we substitute the input with human clicking images.

**ClickNet-RandPrior.** We observe strong spatial bias in the human clicking patterns where most of the clicks tend to be nearby the hidden target (see Sec 5.4 for more discussions). To test this is a useful clicking strategy, we generate random clicks along either side of the black bounding box and use these clicking images with strong spatial prior as inputs to ClickNet.

To study the role of attention and recurrent connections, we introduced two ablated models.

**Variations of VGG16.** One intuitive way of solving the context reasoning problem is to use a feed-forward object recognition network pre-trained on ImageNet, e.g. VGG16 [36], and fine-tune it to classify the hidden target on MSCOCO dataset. During training, the input to the network was an image where one object on the image was randomly covered with a black bounding box. We tested the performance of this alternative model on the 573 images selected for human psychophysics experiments with different input variations: human clicking images (**VGG-humanclick**), the blurred images (**VGG-Blur**), the full-resolution images (**VGG-Fullres**) and images with random clicks (**VGG-Random**).

**VGG-Attention.** Previous work has demonstrated the efficiency of attention in computer vision tasks [33], such as question answering and image captioning [46]. To study the effect of attention in contextual reasoning, we added an attention module to the end of VGG16. To make the complexity of the architecture comparable with ClickNet, we added the same number of fully connected layers as in the LSTM module. As in ClickNet, we used the location with the highest activation value on the attention map to predict the next click. VGG-Attention takes the updated image as input and iteratively predicts the hidden target label. In contrast to ClickNet, the network is feed-forward and there is no incorporation of past information integrated over clicks. We also test VGG-Attention with human clicks (**VGG-Attention-humanclick**) and randomly generated click locations with strong spatial priors (**VGG-Attention-RandPrior**).

We considered several competitive baselines and existing methods of modeling temporal dynamics.

**Human-fixations.** We were concerned about the variable viewing conditions in the MTurk experiments. Therefore, we conducted in-lab psychophysics measurements as a benchmark. In the in-lab experiment, after 500 ms fixation, a bounding box with a fixation cross in the center presented for 1,000 ms indicated the target position and cued subjects to attend to the hidden target location. To ensure that in-lab subjects paid attention to the hidden target location, we recorded their eye movements using an EyeLink D1000 system (SR Research, Canada). The image with the black box was shown for 200, 400, 800, or 1600 ms. Subjects freely moved their eyes; after stimulus offset, subjects said a single noun describing what the hidden target was. We recruited 4 naive subjects (22 to 24 years old, 2 female), each one participating in 573 trials.

**SVM-category.** To study the effect of object co-occurrences, we used a binary vector of size $1 \times C$ as input to a classifier, where the $i$th entry is 1 if there was an object from category $i$th in the image
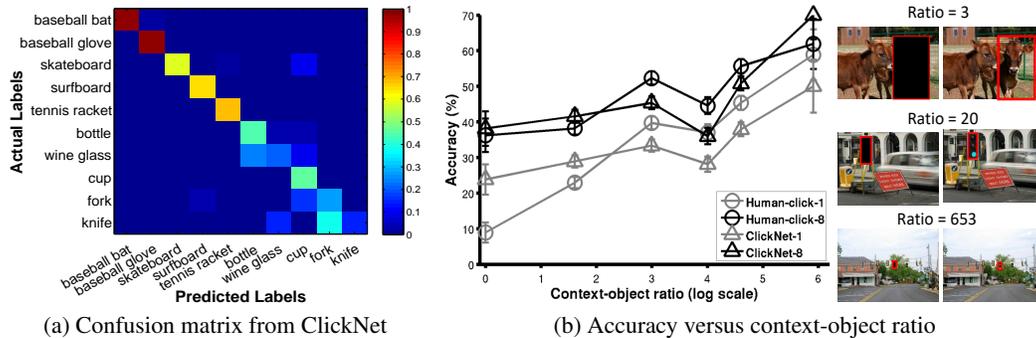
(a) Confusion matrix from ClickNet

(b) Accuracy versus context-object ratio

Figure 5: **Improvement in contextual reasoning accuracy with context-object ratio and patterns of mistakes**. (**a**) Partial view for illustrative purposes of confusion matrix showing the mistakes made by ClickNet among 10 of the 80 object categories in MSCOCO (Fig S5 in Appendix shows the complete confusion matrix with all 80 categories). The element in row $i$, column $j$ denotes the probability that ClickNet predicted label $j$ while the ground truth label was $i$ (see scale bar at right). The sum of probabilities in a row in the full confusion matrix (but not here) equals 1. (**b**) Human (circle) and model (triangle) accuracy for 1 click (gray) and 8 clicks (black) as a function of context-object ratio, shown in logarithmic scale (Sec. 3.3). Right: 3 example images with different context-object ratios. Only the images on the first column were shown (the second column is shown here for illustrative purposes only). Error bars indicate SEM across images.

and 0 otherwise. We assume that the model had perfect information about *all* the object labels (all objects were visible except for the hidden target). A multi-class support vector machine (SVM) classifier was used to predict the hidden target based on this vector.

**SVM-category-instances.** Extending SVM-category, we constructed a vector of size $1 \times C$ where the $i$th entry contained the number $n$ of instances of the $i$th category in the image. A multi-class SVM classifier was used to predict the hidden target based on this vector.

**Hidden Markov Model (HMM).** To study the temporal dynamics over multiple clicks, we considered a Hidden Markov Model where we used all the training images in MSCOCO dataset to calculate the co-occurence matrix of size $C \times C$ as the transition probability matrix. We use normalized uniform vector of size $1 \times C$ as the initial probability. We fine-tuned VGG16 on the MSCOCO dataset and used it for classifying the cropped region at the human clicked locations where the classification vector contributes to emission matrix. The Viterbi decoding algorithm [7] was used for making inferences about the hidden target.

**DeepLab-Conditional Random Field (CRF)** One interesting solution to reason about the hidden target is to run state-of-the-art semantic segmentation algorithms and use majority voting on the predicted labels over all pixels in the bounding box. We used the instantiation in DeepLab-CRF [9].

## 5 Results

### 5.1 What: Importance sampling via attention and prior information

Subjects inferred the identity of the hidden target object with $36.7 \pm 0.9\%$ accuracy after 1 click (Fig 3, gray horizontal line). Performance showed a small, but significant improvement when allowing subjects 8 clicks, reaching $48.4 \pm 0.9\%$ (Fig 3, black horizontal line, $p < 10^{-9}$, two-tailed t-test, t=-6, df=2593). In-lab experiments corroborated these results showing accuracies of $48.1 \pm 0.9\%$ after 200 ms exposure and $56.3 \pm 0.9\%$ after 1,600 ms exposure to the images.

The same images that subjects saw were used to evaluate ClickNet (Figure 3). The ClickNet model showed a close approximation to human performance in the mturk experiments, reaching a top-1 classification performance of $33.3 \pm 0.9\%$ for 1 click and $45.0 \pm 0.9\%$ for 8 clicks. In both cases, performance was only slightly lower than human performance in 1 click ($p = 0.17$, two-tailed t-test, t=1.4, df=1828) and 8 clicks ($p = 0.17$, two-tailed t-test, t=1.4, df=1909). Performance for intermediate numbers of clicks is shown in (Figure S4). For all the computational models, random guessing would yield accuracy = $1.25\%$.

The worst performing model, VGG-Random, yielded performance above chance levels, emphasizing that even small amounts of high-resolution contextual data at arbitrary locations can help solve the problem. Yet, VGG-Random was well below ClickNet's performance ($p < 10^{-5}$,

two-tailed t-test, t=4, df=1144). Adding attention to the model (VGG-Attention) yielded only minimal improvement. An important ingredient missing in VGG-Random and VGG-Attention is the informed location of the clicks. Humans and ClickNet do not sample the image randomly, but rather explore informative locations. Figure 4 shows qualitative examples of clicking patterns from humans and ClickNet. Both humans and ClickNet attend to salient regions on the images. For example, clicks often occur near traffic signs in the first example and near chairs and birds in the second example. Accordingly, substituting the random clicks for the human clicks into the VGG models (VGG-humanclick and VGG-Attention-humanclick) yielded a large performance boost. Conversely, substituting the ClickNet clicks with random clicks leads to large drop in performance when there is only 1 click, even when we artificially try to boost performance by constraining the clicks to be near the hidden target object (ClickNet-RandPrior). This effect is also evident with 2 clicks and 4 clicks (Fig. S4), but disappears with 8 clicks because there is already a lot of high resolution information in the image surrounding the hidden target, and ClickNet can integrate information over time to capitalize on it. Interestingly, the ClickNet sampling clicks are sufficiently close to human clicks that substituting the ClickNet clicks with human clicks does not improve performance (ClickNet-humanclick).

We considered several other comparative models (Fig. S4). Interestingly, using just a few clicks, ClickNet reaches performance that is essentially equivalent to that of VGG using a full resolution version of the entire image. Other comparative models (VGG-Blur, SVM-category, SVM-category-instances, HMM, DeepLab-CRF) showed above chance performance but their accuracies were well below ClickNet.

## 5.2 What: the more, the merrier

To investigate how much context information is needed to enhance recognition, we evaluated accuracy as a function of context-object ratio (Fig 5b, Sec 3.3). Images with higher context-object ratio contained more context information for inference, and yielded higher accuracy both for humans and models. Similarly, accuracy improved with increasing numbers of clicks (Fig 3 and Fig 5b).

It is not just the quantity of context, but also the specific quality of contextual information that matters. In the real world, objects do not tend to appear in isolation but rather they co-vary with other objects. As ClickNet explores more regions on the image, it integrates information at previous clicked locations and learns associations of objects. The pattern of mistakes made by ClickNet is indicative of those associations (Fig 5a and Fig S5). ClickNet often makes "reasonable" wrong guesses when there is ambiguity in context reasoning, as humans do. For example, knife tends to be associated and therefore confused with spoon, fork, and wine glass, but knife seldom co-occurs with baseball bat or skateboard in these images.

## 5.3 Where: consistency of human and model clicks

We hinted at the presumed similarity in the clicking patterns between humans and ClickNet based on the accuracy of the ClickNet-humanclick and ClickNet-RandPrior comparative models in Fig 3. To more directly assess whether ClickNet learned to sample the image to gather information about areas of contextual relevance, we directly quantified the similarity in clicking patterns (Figure 6). To interpret the distances between human clicks and model clicks, we computed the degree of human-human consistency in the clicking patterns. The clicking patterns of ClickNet were overall similar to those made by humans. The model clicks were still different from the consistency between two humans (for 8 clicks, $p < 10^{-15}$, two-tailed t-test, t=-30.4, df=29542); yet, the model clicks were much more similar to human clicks than random clicks (for 8 clicks, $p < 10^{-15}$, two-tailed t-test, t=-50.3, df=35310).

## 5.4 Where: tendency of clicking nearby the target

There was a strong spatial bias towards clicking near the target for both humans and ClickNet (see examples in Figure 4). To quantify this spatial bias, we computed the Euclidean distance between the clicked locations and the center of the bounding box, normalized by the diagonal of the bounding box (Figure 6c). Humans tended to click within approximately one diagonal distance of the target box. Interestingly, although ClickNet does not take any human supervisory signal during training, ClickNet still learned to capture the tendency of clicking near the target.

We asked whether this spatial bias in sampling behavior is sufficient to explain performance in this task in a modified version of ClickNet. We removed the clicks dictated by the attention module and instead forced the clicks to be randomly distributed while still respecting the
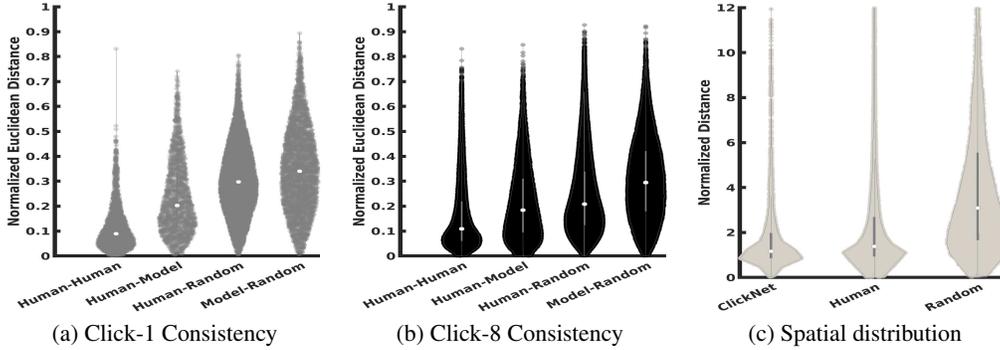
(a) Click-1 Consistency     (b) Click-8 Consistency     (c) Spatial distribution

Figure 6: **Model click locations were similar to human sampling.** (**a-b**) Consistency of click patterns between human subjects (human-human), consistency between humans and model (human-model), consistency between humans and random (human-random) and consistency between models and random (model-random) for 1 click (**a**) and 8 clicks (**b**) measured by the distribution of normalized Euclidean distances with respect to the diagonal of the image between any pairs of clicks by humans and ClickNet or random clicks. In each trial, we permute the sequence of mouse clicks between pairs of human and model clicks such that their sum of Euclidean distance is minimized across all clicks. The white circles denote the median of the distribution and the light grey bar denote the 1st and 3rd quartiles. (**b**) Euclidean Distance between click locations and center of the hidden target bounding box normalized by the diagonal of the hidden target bounding box.

spatial distribution in Fig 6c (ClickNet-RandPrior). Both ClickNet and ClickNet-humanclick surpassed ClickNet-RandPrior by 17% and 10% respectively (Figure 3). Similar results were obtained when using only the VGG architecture: VGG-Attention-humanclick was 16% better than VGG-Attention-RandPrior (Fig S4 in Appendix). Therefore, the spatial bias in clicking behavior is not sufficient to explain performance in this task. Sampling for context reasoning involves more than clicking near the target.

## 5.5  When: role of recurrent connections

Several lines of evidence support the importance of the recurrent network in the LSTM module in ClickNet. ClickNet outperformed the competitive baselines and state-of-the-art comparative methods to make inferences (Figure 3 and Fig S4 in Appendix). We tested whether the co-occurrence of object categories or the number of objects per category present in the image would be sufficient for context reasoning (SVM-category and SVM-category-instances). Even though these alternative models were exposed to full contextual information on the image and assumed perfect labeling of all objects in the image (except for the hidden target object), there was still a large overall performance drop in performance with respect to ClickNet. Moreover, graphical models for inference, such as Hidden Markov Model and DeepLab with Conditional Random Field (Sec 4.5) failed to reach ClickNet's accuracy in this task (Fig S4). The ablation studies eliminating the LSTM module further support the role of integrating information over multiple clicks in this task, as evidenced by the observation that ClickNet outperformed the VGG-Attention model.

## 6  Discussion

Here we quantitatively studied the role of contextual information in visual recognition in human observers and computational models in a task that involved inferring the identity of a hidden target object. Contextual influenced recognition based on the amount of context, the specific location of contextual cues, and the dynamic sampling among salient visual features. We introduced a recurrent neural network model that combines a feed-forward visual stream module that extracts image features in a dynamic fashion, combined with an attention module to prioritize different image locations and select the next sampling step, and a recurrent LSTM module that integrates information over time and produces a label for the hidden object. Surprisingly, even though the model lacks the expertise that humans have in interacting with objects in their context, the model adequately predicts human sampling behavior and reaches almost human-level performance in this contextual reasoning task. The model opens the doors to examine more complex form of reasoning about scenes and how to integrate sequential sampling with prior knowledge.

# References

[1] E. Aminoff, N. Gronau, and M. Bar. The parahippocampal cortex mediates spatial and nonspatial associations. *Cerebral Cortex*, 17(7):1493–1503, 2006.

[2] M. E. Auckland, K. R. Cave, and N. Donnelly. Nontarget objects can influence perceptual processes during object recognition. *Psychonomic bulletin & review*, 14(2):332–337, 2007.

[3] J. Ba, V. Mnih, and K. Kavukcuoglu. Multiple object recognition with visual attention. *arXiv preprint arXiv:1412.7755*, 2014.

[4] M. Bar and E. Aminoff. Cortical analysis of visual context. *Neuron*, 38(2):347–358, 2003.

[5] P. Battaglia, R. Pascanu, M. Lai, D. J. Rezende, et al. Interaction networks for learning about objects, relations and physics. In *Advances in neural information processing systems*, pages 4502–4510, 2016.

[6] I. Biederman, R. J. Mezzanotte, and J. C. Rabinowitz. Scene perception: Detecting and judging objects undergoing relational violations. *Cognitive psychology*, 14(2):143–177, 1982.

[7] P. Blunsom. Hidden markov models. *Lecture notes, August*, 15(18-19):48, 2004.

[8] W. Brendel and M. Bethge. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. *arXiv preprint arXiv:1904.00760*, 2019.

[9] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.

[10] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2018.

[11] W. Choi and S. Savarese. A unified framework for multi-target tracking and collective activity recognition. In *European Conference on Computer Vision*, pages 215–230. Springer, 2012.

[12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[13] Z. Deng, A. Vahdat, H. Hu, and G. Mori. Structure inference machines: Recurrent neural networks for analyzing relations in group activity recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4772–4781, 2016.

[14] C. Desai, D. Ramanan, and C. C. Fowlkes. Discriminative models for multi-class object layout. *International journal of computer vision*, 95(1):1–12, 2011.

[15] S. K. Divvala, D. Hoiem, J. H. Hays, A. A. Efros, and M. Hebert. An empirical study of context in object detection. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 1271–1278. IEEE, 2009.

[16] A. Friedman. Framing pictures: The role of knowledge in automatized encoding and memory for gist. *Journal of experimental psychology: General*, 108(3):316, 1979.

[17] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.

[18] J. O. Goh, S. C. Siong, D. Park, A. Gutchess, A. Hebrank, and M. W. Chee. Cortical areas involved in object, background, and object-background processing revealed with functional magnetic resonance adaptation. *Journal of Neuroscience*, 24(45):10223–10228, 2004.

[19] J. M. Gonfaus, X. Boix, J. Van de Weijer, A. D. Bagdanov, J. Serrat, and J. Gonzalez. Harmony potentials for joint classification and segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 3280–3287. IEEE, 2010.

[20] J. M. Henderson, P. A. Weeks Jr, and A. Hollingworth. The effects of semantic consistency on eye movements during complex scene viewing. *Journal of experimental psychology: Human perception and performance*, 25(1):210, 1999.

[21] D. Hoiem, A. A. Efros, and M. Hebert. Geometric context from a single image. In *Computer Vision, 2005. ICCV 2005. Tenth IEEE International Conference on*, volume 1, pages 654–661. IEEE, 2005.

[22] A. Hollingworth. Does consistent scene context facilitate object perception? *Journal of Experimental Psychology: General*, 127(4):398, 1998.

[23] H. Hu, G.-T. Zhou, Z. Deng, Z. Liao, and G. Mori. Learning structured inference neural networks with label relations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2960–2968, 2016.

[24] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5308–5317, 2016.

[25] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[26] A. Kovashka, O. Russakovsky, L. Fei-Fei, K. Grauman, et al. Crowdsourcing in computer vision. *Foundations and Trends® in computer graphics and Vision*, 10(3):177–243, 2016.

[27] L. Ladicky, C. Russell, P. Kohli, and P. H. Torr. Graph cut based inference with co-occurrence statistics. In *European Conference on Computer Vision*, pages 239–253. Springer, 2010.

[28] D. Lin, S. Fidler, and R. Urtasun. Holistic scene understanding for 3d object detection with rgbd cameras. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1417–1424, 2013.

[29] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[30] Y. Liu, R. Wang, S. Shan, and X. Chen. Structure inference net: Object detection using scene-level context and instance-level relationships. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6985–6994, 2018.

[31] K. Marino, R. Salakhutdinov, and A. Gupta. The more you know: Using knowledge graphs for image classification. *arXiv preprint arXiv:1612.04844*, 2016.

[32] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 891–898, 2014.

[33] T. V. Nguyen, Q. Zhao, and S. Yan. Attentive systems: A survey. *International Journal of Computer Vision*, 126(1):86–110, 2018.

[34] D. Park, D. Ramanan, and C. Fowlkes. Multiresolution models for object detection. In *European conference on computer vision*, pages 241–254. Springer, 2010.

[35] A. G. Schwing and R. Urtasun. Fully connected deep structured networks. *arXiv preprint arXiv:1503.02351*, 2015.

[36] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[37] H. Tang, M. Schrimpf, W. Lotter, C. Moerman, A. Paredes, J. O. Caro, W. Hardesty, D. Cox, and G. Kreiman. Recurrent computations for visual pattern completion. *Proceedings of the National Academy of Sciences*, 115(35):8835–8840, 2018.

[38] D. Teney, L. Liu, and A. van den Hengel. Graph-structured representations for visual question answering. *arXiv preprint*, 2017.

[39] A. Torralba. Contextual priming for object detection. *International journal of computer vision*, 53(2):169–191, 2003.

[40] A. Torralba, K. Murphy, and W. Freeman. Using the forest to see the trees: Ob-ject recognition in contex. *Comm. of the ACM*, 2010.

[41] A. Torralba, K. P. Murphy, and W. T. Freeman. Contextual models for object detection using boosted random fields. In *Advances in neural information processing systems*, pages 1401–1408, 2005.

[42] A. M. Turk. Amazon mechanical turk. *Retrieved August*, 17:2012, 2012.

[43] C. Vondrick, H. Pirsiavash, A. Oliva, and A. Torralba. Learning visual biases from human imagination. In *Advances in neural information processing systems*, pages 289–297, 2015.

[44] K. Wu, E. Wu, and G. Kreiman. Learning scene gist with convolutional neural networks to improve object recognition. In *Information Sciences and Systems (CISS), 2018 52nd Annual Conference on*, pages 1–6. IEEE, 2018.

[45] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 2, 2017.

[46] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015.

[47] J. Yao, S. Fidler, and R. Urtasun. Describing the scene as a whole: Joint object detection, scene classification and semantic segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 702–709. IEEE, 2012.

[48] W. Zaremba, I. Sutskever, and O. Vinyals. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*, 2014.

[49] M. Zhang, J. Feng, K. T. Ma, J. H. Lim, Q. Zhao, and G. Kreiman. Finding any waldo with zero-shot invariant and efficient visual search. *Nature communications*, 9(1):3730, 2018.

[50] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1529–1537, 2015.

# A Appendix

We provide supplementary figures and materials here. All labels in supplementary figures and tables are pre-fixed with letter S in front.
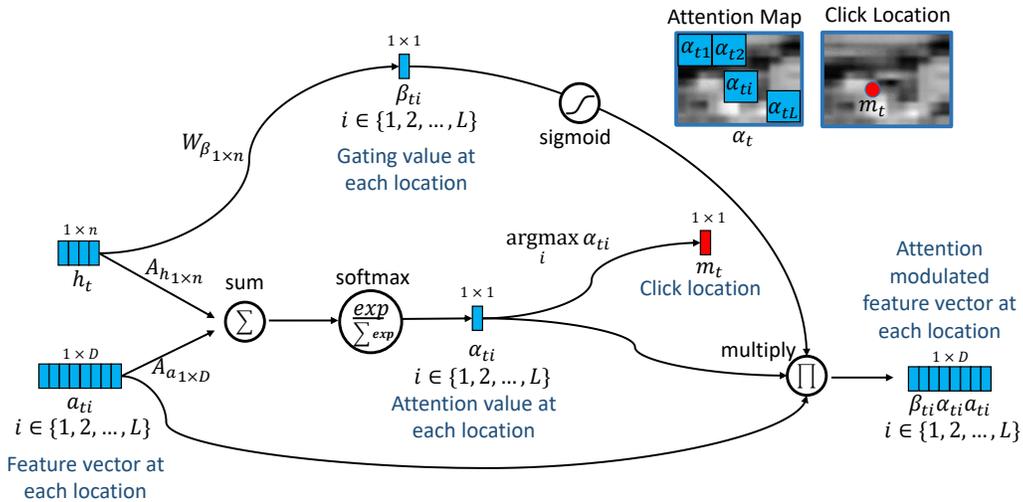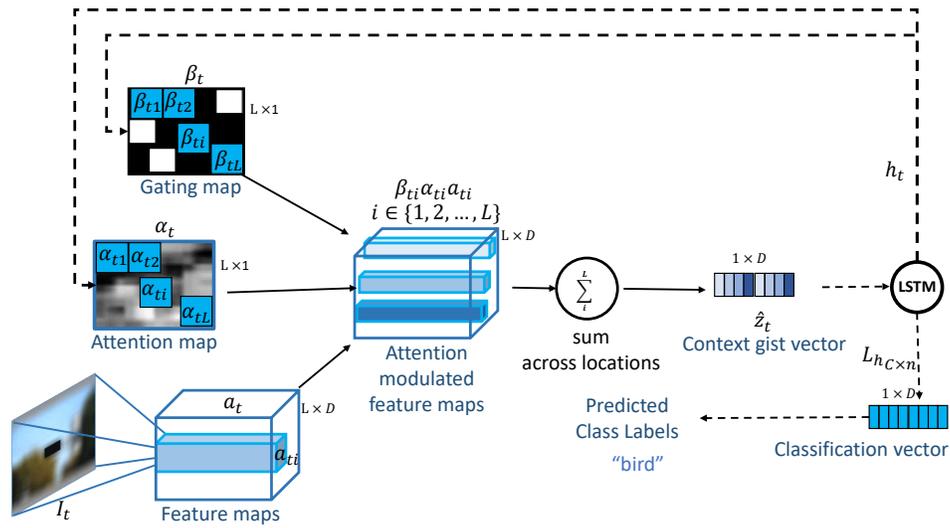


Figure S1: **Schematic illustration of the attention module implementation**. Expanding on the overall ClickNet architecture shown in Fig 2, here we zoom into the attention module. The attention module takes as inputs the features at each location $a_{ti}$ and the output of the LSTM module $h_t$ and selects the next click location $m_t$ and a map that modulates the features at each location (see Section 4 for a description of all the variables).

Figure S2: **Schematic illustration of the LSTM module implementation**. Expanding on the overall ClickNet archietcture shown in Fig 2, here we zoom into the LSTM module. The LSTM module takes as input context gist vector $\widehat{\mathbf{z}}_t$ and integrates the information with the previous state to inform the attention module in the next time step via $h_t$ and to predict a class label (see Section 4 for a description of all the variables).
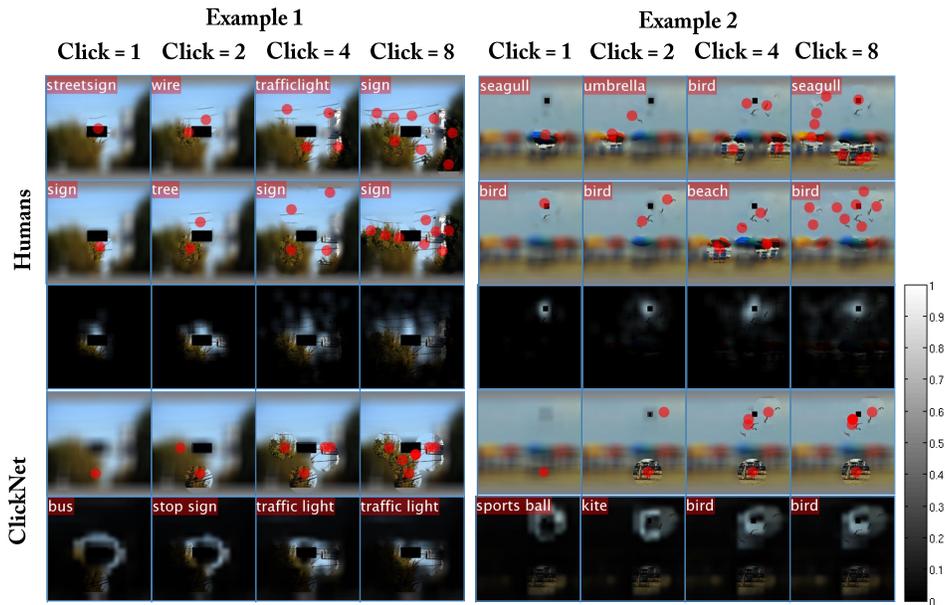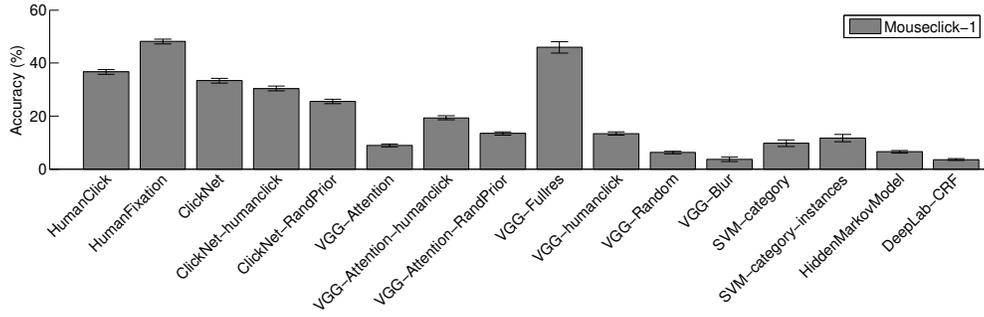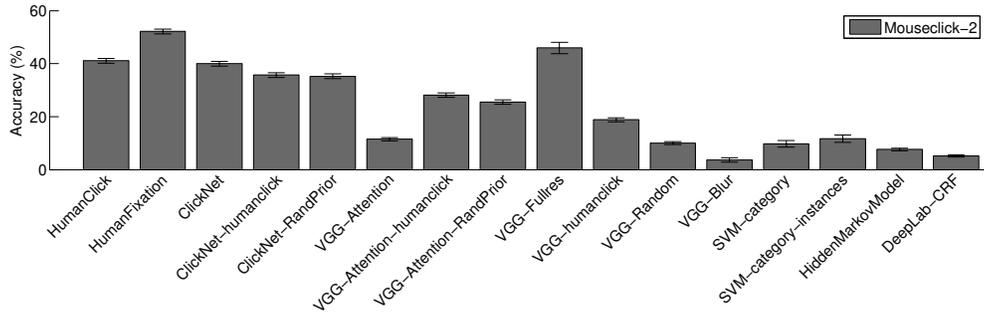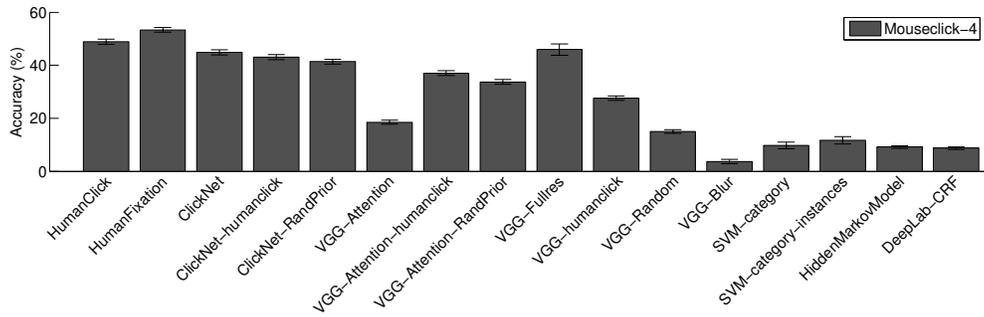


Figure S3: **Example visualziation for humans and ClickNet.** This figure shows click locations and attention maps using the same format as Fig 4, here adding results for 2 clicks and 4 clicks.
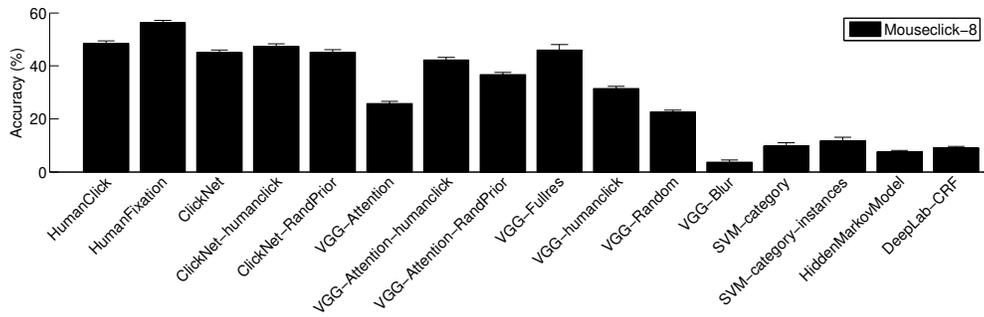
(a) Click-1



(b) Click-2



(c) Click-4



(d) Click-8

Figure S4: **Contextual reasoning accuracy of humans and models**. Expanding on the results in Fig. 3a-b, here we add the results for 2 clicks and 4 clicks, as well as additional comparative models (Section 4.5) describe each model).

Figure S5: **Confusion matrix for ClickNet with all click conditions**. The format is the same as in Figure 5a, except showing all 80 categories here. The element in row $i$, column $j$ denotes the probability that ClickNet predicted label $j$ while the ground truth label was $i$ (see scale bar on right). The sum of all probabilities in a row equals 1.