

# A dual foveal-peripheral visual processing model implements efficient saccade selection

Emmanuel Daucé, Pierre Albiges, Laurent U Perrinet

## ▶ To cite this version:

Emmanuel Daucé, Pierre Albiges, Laurent U Perrinet. A dual foveal-peripheral visual processing model implements efficient saccade selection. Journal of Vision, 2020, 20 (8), pp.22. 10.1167/jov.20.8.22. hal-02947410

## HAL Id: hal-02947410 https://amu.hal.science/hal-02947410

Submitted on 24 Sep 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

# A dual foveal-peripheral visual processing model implements efficient saccade selection

## **Emmanuel Daucé**

**Pierre Albiges** 

## Laurent U. Perrinet

Institut de Neurosciences de la Timone (UMR 7289), Aix Marseille University, CNRS, Marseille, France

Institut de Neurosciences de la Timone (UMR 7289), Aix Marseille University, CNRS, Marseille, France

Institut de Neurosciences de la Timone (UMR 7289), Aix Marseille University, CNRS, Marseille, France



We develop a visuomotor model that implements visual search as a focal accuracy-seeking policy, with the target's position and category drawn independently from a common generative process. Consistently with the anatomical separation between the ventral versus dorsal pathways, the model is composed of two pathways that respectively infer what to see and where to look. The "What" network is a classical deep learning classifier that only processes a small region around the center of fixation, providing a "foveal" accuracy. In contrast, the "Where" network processes the full visual field in a biomimetic fashion, using a log-polar retinotopic encoding, which is preserved up to the action selection level. In our model, the foveal accuracy is used as a monitoring signal to train the "Where" network, much like in the "actor/critic" framework. After training, the "Where" network provides an "accuracy map" that serves to guide the eye toward peripheral objects. Finally, the comparison of both networks' accuracies amounts to either selecting a saccade or keeping the eye focused at the center to identify the target. We test this setup on a simple task of finding a digit in a large, cluttered image. Our simulation results demonstrate the effectiveness of this approach, increasing by one order of magnitude the radius of the visual field toward which the agent can detect and recognize a target, either through a single saccade or with multiple ones. Importantly, our log-polar treatment of the visual information exploits the strong compression rate performed at the sensory level, providing ways to implement visual search in a sublinear fashion, in contrast with mainstream computer vision.

## Introduction

#### **Problem statement**

The field of computer vision was recently recast by the outstanding capability of convolution-based deep neural networks to capture the semantic content of images and photographs. Human performance is now outreached by computer algorithms in numerous image categorization tasks (He, Zhang, Ren, & Sun, 2015). One of the reasons explaining this breakthrough is a strong reduction in the number of parameters used to train the network, through a massive sharing of weights in the convolutional layers. Reducing the number of parameters and/or the size of the visual data that need to be processed is a decisive factor for further improvements. Initially trained on energy-greedy, high-performance computers, these algorithms are now designed to work on more common hardware such as desktop computers with dedicated GPU hardware (Sandler, Howard, Zhu, Zhmoginov, & Chen, 2018). Despite much effort in hardware and software optimization, the processing of pixel-based images is still done at a cost that scales linearly with the image size: All pixels present in the image are systematically processed by the computer algorithm, even the ones that are useless for the task at hand. Current computer vision algorithms consequently manipulate millions of pixels and millions of variables with ensuing energy consumption, even in the case of downsampled images and with a still prohibitive cost for large images and videos. The need to detect visual objects at a glance while running on resource-constrained embedded hardware, for instance, in autonomous driving, introduces a necessary trade-off between efficiency and accuracy, requiring renewed mathematical treatment and computational implementations.

Citation: Daucé, E., Albiges, P., & Perrinet, L. U. (2020). A dual foveal-peripheral visual processing model implements efficient saccade selection. *Journal of Vision*, 20(8):22, 1–20, https://doi.org/10.1167/jov.20.8.22.

Received January 27, 2020; published August 20, 2020

ISSN 1534-7362 Copyright 2020 The Authors



Interestingly, things work differently when human vision is considered. First, human vision is still unsurpassable in the case of ecological real-time sensory flows. Indeed, object recognition can be achieved by the human visual system both rapidly—in less than 100 ms (Kirchner & Thorpe, 2006)—and at a low energy cost (< 5 W). On top of that, it is mostly self-organized, robust to visual transforms or lighting conditions. and can learn with few examples. If many different anatomical features may explain this efficiency, the main difference lies in the fact that its sensor (the retina) combines a nonhomogeneous sampling of the world with the capacity to rapidly change its center of fixation: On the one hand, the retina is composed of two separate systems: a central, high-definition fovea (a disk of about 6 degrees of diameter in visual angle around the center of gaze) and a large, lower-definition peripheral area (Strasburger et al., 2011). On the other hand, the human vision is active and dynamic: The retina is attached at the back of the eye, which is capable of low-latency, high-speed eye movements. In particular, saccades are stereotyped eye movements that allow for efficient changes of the position of the center of gaze: They take about 200 ms to initiate, last about 200 ms and usually reach a maximum velocity of approximately 600 degrees per second (Bahill, Clark, & Stark, 1975). The scanning of a full visual scene is thus not done in parallel but sequentially, and only scene-relevant regions of interest are scanned through saccades. This implies a decision process between each saccade that decides where to look next. This behavior is prevalent in biological vision with on average a saccade every 2 s, that is, almost a billion saccades in a lifetime. The interplay of peripheral search and focal inspection allows human observers to engage in an integrated action/perception loop that sequentially scans and analyzes the different parts of the visual scene.

Take, for instance, the case of an encounter with a friend in a crowded café. To catch the moment of his or her arrival, a face-seeking visual search is needed, possibly under heavy sensory clutter conditions. To do so, relevant parts of the visual scene need to be scanned sequentially with the gaze. Each saccade may potentially allow you to recognize your friend, provided it is accurately focused on each target face. The main feature of this task is thus the monitoring of a particular *class* of objects (e.g., human faces) in the periphery of the visual field before the actual eye displacement and the processing of the foveal visual data. Searching for *any* face in a peripheral and crowded display needs thus to precede the recognition of a specific face *identity*. Since biological vision is the result of a continual optimization under strong material and energy constraints via natural selection, it is important to understand both its ground principles and its specific computational and material constraints in order to implement effective biomimetic vision systems.

The problem we address is thus how to ground an artificial visual processing system on top of the material constraints found in human vision that is conforming to the structure of the visual input and to the capability of the visual apparatus to rapidly scan a visual scene through saccades in order to find and identify objects of interest. We thus start from an elementary visual search problem, which is how to locate an object in a large, cluttered image, and take human vision as a guide for efficient design.

#### State of the art

The visual search problem, that is, finding and identifying objects in a visual scene, is a classical task in computer vision, appealing as well to machine learning, signal processing, or robotics. Crucially, it also speaks to neuroscience, for it refers to the mechanisms underlying foveation and more generally to low-level attention mechanisms. When restricted to a mere "feature search" (Treisman & Gelade, 1980), many computational solutions are proposed in the computer vision literature. Notably, recent advances in deep learning have been proven efficient to solve the task with models such as faster-RCNN (Ren, He, Girshick, & Sun, 2017) or YOLO (Redmon, Divvala, Girshick, & Farhadi, 2016). Typical object search implementations predict in the image the probability of proposed bounding boxes around visual objects. While rapid, the potential number of boxes may significantly increase with image size, and the approach more generally necessitates dedicated hardware to run in real time (Feng, Jiang, Yang, Du, & Li, 2019). Under fine-tailored algorithmic and material optimization, the visual search problem can be considered in the best case as *linear* in the number of pixels (Strengert, Kraus, & Ertl, 2006), which still represents a heavy load for real-time image processing. This poses the problem of the *energy scaling* of current computer vision algorithms to large/high-definition visual displays. This scaling problem becomes even more crucial when considering a dynamical stream of sensory images.

Analogously to human visual search strategies, low-level attentional mechanisms may help guide the localization of targets. A sequence of saccades over a natural scene defines a scan-path that provides ways to define *saliency maps* (Itti & Koch, 2001). These quantify the attractiveness of the different parts of an image that are consistent with the detection of objects of interest. Essential to understand and predict saccades, they also serve as phenomenological models of attention. Estimating the saliency map from a luminous image is a classical problem in neuroscience that was shown to be consistent with a distance from baseline image statistics known as the "Bayesian surprise" (Itti & Baldi, 2009). Such an approach was extended in the



Figure 1. Problem setting: In generic, ecological settings, when searching for one target (from a class of targets) in a cluttered environment, the visual system is bound with an action selection problem. It is synthesized in the following virtual experiment: (A) After a fixation period FIX of 200 ms, an observer is presented with a luminous display DIS showing a single target from a known class (here digits) put at a random position within the field of view. The display is presented for a short period of 500 ms (light shaded area in B), which is enough to perform at most one saccade on the potential target (SAC, here successful). Finally, the observer has to identify the digit by a keypress ANS. *NB*: the target contrast is here enhanced to 100% for better readability. (B) Prototypical trace of a saccadic eye movement to the target position. In particular, we show the fixation window FIX and the temporal window during which a saccade is possible (green shaded area). (C) Simulated reconstruction of the visual information from the internal retinotopic map at the onset of the display DIS and after a saccade SAC, the dashed red box indicating the foveal region. The task does not consist of inferring the location of the target, but rather to infer an action that may provide relevant pixels at the center of fixation, allowing to identify the target's category. By comparison with the external display (see A), the action is processed from log-polar coefficients, representing a focal sample of the total visual field. Controlling the clutter and reducing the contrast of the digit allows modulating the task's difficulty.

AIM (Bruce & Tsotsos, 2009) and SUN (Zhang, Tong, Marks, Shan, & Cottrell, 2008) models. Recently, the saliency approach was updated using deep learning to estimate saliency maps over large databases of natural images (Kummerer, Wallis, Gatys, & Bethge, 2017). While efficient at predicting the probability of fixation, these methods miss an essential component in the action-perception loop: They operate on the full image while the retina operates on the nonuniform, foveated sampling of visual space (see Figure 1C). Herein, we believe that this constitutes an essential factor to reproduce and understand the active vision process.

Foveated models of vision have been considered for a long time in robotics and computer vision as a way to leverage the visual scene scaling problem. Focal computer vision relies on a nonhomogeneous compression of an image that maintains the pixel information at the center of fixation and strongly compresses it at the periphery, including pyramidal encoding (Kortum & Geisler, 1996; Butko & Movellan, 2010), local wavelet decomposition (Daucé, 2018) and log-polar encoding (Traver & Bernardino, 2010). A recent deep-learning-based implementation of such compression shows that in a video flow, a log-polar sampling of the image is sufficient to provide a reconstruction of the whole image (Kaplanyan et al., 2019). However, this particular algorithm lacks a system predicting the best saccadic action to perform. In summary, though focal and multiscale encoding is now largely considered in static computer vision, sequential implementations have not been shown effective enough to overtake static object search methods. Several implementations of a focal sequential search in visual processing can be found in the literature, with various degrees of biological realism (Mnih, Heess, Graves, et al., 2014; Fu, Zheng, & Mei, 2017) that often rely on a simplified focal encoding, long training procedures, and bounded sequential processing. More realistic attempts to combine foveal encoding and sequential

In contrast with the phenomenological (or "bottom-up") approaches, active models of vision (Najemnik & Geisler, 2005; Butko & Movellan, 2010; Daucé, 2018) provide the ground principles of saccadic exploration. In general, they assume the existence of a generative model from which both the target position and category can be inferred through active sampling. This comes from the constraint that the visual sensor is foveated but can generate a saccade. Several studies are relevant to our endeavor. First, one can consider optimal strategies to solve the problem of the visual search of a target (Najemnik & Geisler, 2005). In a setting similar to that presented in Figure 1A, where the target is an oriented edge and the background is defined as pink noise, authors show first that a Bayesian ideal observer comes out with an optimal strategy and second that human observers are close to that optimal performance. Though it is well predicting sequences of saccades in a perception action loop, this model is limited by the simplicity of the display (elementary edges added on stationary noise, a finite number of nonoverlapping locations on a discrete grid) and by the abstract level of modeling. Despite these (inevitable) simplifications, this study could successfully predict some key characteristics of visual scanning such as the trade-off between memory content and speed. Looking more closely at neurophysiology, the study by Samonds, Geisler, and Priebe (2018) allows one to go further in understanding the interplay between saccadic behavior and the statistics of the input. In this study, authors were able to manipulate the size of saccades by monitoring key properties of the presented (natural) images. For instance, smaller images generate smaller saccades. Interestingly, they also predicted the size of saccades from the size of visual receptive fields for different species, including mice, which lack a foveal region. One key prediction of this study that is relevant for our problem is the fact that saccades seem optimal to a priori decorrelate the visual input, that is, to minimize redundancy in the sequence of generated saccades, knowing the statistics of the visual inputs.

A further modeling perspective is provided by Friston, Adams, Perrinet, and Breakspear (2012). In this setup, a full description of the visual world is used as a generative process. An agent is completely described by the generative model governing the dynamics of its internal beliefs and is interacting with this image by scanning it through a foveated sensor, just as described in Figure 1. Thus, equipping the agent with the ability to actively sample the visual world allows one to interpret saccades as optimal experiments, by which the agent seeks to confirm predictive models of the (hidden) world. One key ingredient to this process is the (internal) representation of counterfactual predictions, that is, the probable consequences of a possible hypothesis as they would be realized into actions (here, saccades). Following such an active inference scheme, numerical simulations reproduce sequences of eye movements that fit well with empirical data (Mirza, Adams, Mathys, & Friston, 2018). Compared to previous studies (Najemnik & Geisler, 2005), saccades are not the output of a value-based cost function such as a saliency map but are the consequence of an active strategy by the agent to minimize the uncertainty about his or her beliefs, knowing his or her priors on the generative model of the visual world.

#### Outline

Daucé, Albiges, & Perrinet

So far, few models in active vision come with an integrated processing of the visual scene, from early visual treatment toward saccade selection. The difficulty lies in combining object hypothesis (feature) space along with their spatial mapping. As pointed out earlier, the system needs to guess where the interesting objects lie in space before actually knowing what they are. Establishing the position of the objects in space is thus crucial, for it resorts to the capability of the eye to reach them with a saccade, so as to finally identify them. Inferring the target's position in the peripheral visual field is thus an essential component of focal visual processing, and the acuity of such target selection ultimately conditions the capability to rapidly and efficiently process the scene. Stemming from the active vision principles, we thus address the question of the interplay of the location and identity processing in vision and provide an artificial vision setup that efficiently implements those principles. Moreover, despite refined generative models, the processing of the visual data found in biomimetic models should resort to a combination of local features to build posterior beliefs. Herein, our framework is made as general as possible, with minimal mathematical treatment, to speak largely to fragmented domains, such as machine learning, neuroscience, and robotics.

The article is organized as follows. After this introduction, the principles underlying accuracy-based saccadic control are defined in the second section. We first define notations, variables, and equations for the generative process governing the experiment and the generative model for the active vision agent. Complex combinatorial inferences are here replaced by separate pathways, that is, the spatial ("Where") and categorical ("What") pathways, whose output is combined to infer optimal eye displacements and subsequent identification of the target. Our agent, equipped with a foveated sensor, should learn an optimal behavior strategy to actively scan the visual scene. Numerical simulations are presented in the Results section, demonstrating the applicability of this framework to tasks with different complexity levels. The Discussion section finally summarizes the results, showing their relative advantages in comparison with other frameworks and providing ways toward possible improvements. Implementation details are provided in the Methods section, giving ways to reproduce our results, showing in particular how to simplify the learning using accuracy-driven action maps.

## Setup

#### **Experimental design**

In order to implement our visual processing setup, we provide a simplified visual environment toward which a visual agent can act on. This visual search task is formalized and simplified in a way reminiscent to classical psychophysical experimentation: An observer is asked to classify digits (for instance, as taken from the MNIST data set, as introduced by Lecun, Bottou, Bengio, & Haffner, 1998) as they are shown with a given size on a computer display. However, these digits can be placed at random positions on the display, and visual clutter is added as a background to the image (see Figure 1A). In order to vary the difficulty of the task, different parameters are controlled, such as the target eccentricity, the background noise period, and the signal/noise ratio (SNR). The agent initially fixates the center of the screen. Due to the peripheral clutter, he or she needs to explore the visual scene through saccades to provide the answer. He or she controls a foveal visual sensor that can move over the visual scene through saccades (see Figure 1B). When a saccade is actuated, the center of fixation moves toward a new location, which updates the visual input (see Figure 1C). The lower the SNR and the larger the initial target eccentricity, the more difficult the identification. There is a range of eccentricities for which it is impossible to identify the target from a single glance, so that a saccade is necessary to reduce the relative eccentricity and issue a proper response. This setup implies also that the position of the object may be detected in the peripheral clutter *before* being properly identified.

This setup provides the conditions for a separate processing of the visual information. On the one side, the detailed information present at the center of fixation needs to be analyzed to provide specific environmental cues. On the other side, the full visual field, that is, mainly the low-resolution part surrounding the fovea, needs to be processed in order to identify regions of interest that deserve fixation. This basically means making a choice of what is interesting next. The actual content of putative peripheral locations does not need to be known in advance, but it needs to look interesting enough and of course to be reachable by a saccade. This is reminiscent of the What/Where visual processing separation observed in primates' ventral and dorsal visual pathways (Mishkin, Ungerleider, & Macko, 1983).

#### **Computational implementation**

In order to show it is possible to learn such a task, it is sufficient to demonstrate the existence of a simple "deep learning" neural network that would implement it, for instance, through the effective success of its training. More specifically, this class of modern parametric classifiers is composed of many layers (hence the terminology) that can be trained through gradient descent over arbitrary input and output feature spaces. For our specific problem, the anatomy of the agent is made of two separate pathways for which a different processing is realized by two different neural networks (see Figure 2). The proposed computational architecture is connected in a closed-loop fashion to the visual environment, with the capacity to produce saccades whose effect is to shift the visual field from one visual position to another. By analogy with biological vision, the identification of the target is assumed to rely on the very central part of the retina (the fovea) that comes with higher density of cones and thus higher spatial precision. In contrast, the saccade planning should rely on the full visual field, with peripheral regions having a lower density of sensors and thus a lower sensitivity to high spatial frequencies.

In a stationary condition, where the target's position and identity do not change over time, each saccade thus provides a new viewpoint over the scene, allowing one to form a new estimation of the target identity. Following the active inference setup (Najemnik & Geisler, 2005; Friston et al., 2012), we assume that, instead of trying to detect the actual position of the target, the agent tries to maximize the counterfactual benefit in scene understanding that would be gained by any potential saccade. The focus is thus put on action selection metric rather than spatial representation. This means in short estimating how accurate a categorical target classifier will be after moving the eye. In a full setup, predictive action selection means first predicting the future visual field, denoted x', which is obtained at the center of fixation, and then predicting how good the estimate of the target identity, denoted y, that is, p(y|x'), will be at this location. In practice, predicting a future visual field over all possible saccades is too computationally expensive. Better off instead is to record, for every context x, the improvement obtained in recognizing the target after different saccades  $a, a', a'', \dots$  If a is a possible saccade and x' the corresponding future visual field, the result of the central categorical classifier over x' can either be correct (1) or incorrect (0). If this experiment is repeated many times over many visual



Figure 2. Computational graph. Based on the general anatomy of the visual pathways, we define two streams of information, one stream for identifying a target ("What"?) and the other for categorizing it in visual space ("Where"?). (A) The visual display is constructed the following way: first a  $128 \times 128$  natural-like background noise is generated, characterized by noise contrast, mean spatial frequency, and bandwidth (Sanz-Leon, Vanzetta, Masson, & Perrinet, 2012). Then, a sample digit is selected from the MNIST data set (of size  $28 \times 28$ ), rectified, multiplied by a contrast factor, and overlaid on the background at a random position (see another example in Figure 1A, DIS). Last, a circular mask is put on. (B) The visual display is then transformed in two sensory inputs: (i) A  $28 \times 28$  central foveal-like snippet is fed to a classification network ("What" pathway). (ii) A log-polar set of oriented visual features is fed to the "Where" pathway. This log-polar input is generated by a bank of filters whose centers are positioned on a log-polar grid and whose radius increases proportionally with the eccentricity. (C) The "What" network is implemented using the three-layered LeNet neural network (Lecun et al., 1998). This network outputs a vector predicting the accuracy of detecting the correct digit. In parallel, the "Where" network is implemented by a three-layered neural network consisting of the retinal log-polar input, two hidden layers (fully connected linear layers combined with a ReLU nonlinearity), and a collicular-like accuracy map at the output. This map has a similar log-polar (retinotopic) organization and predicts the accuracy at the hypothetical position of a saccade. Both networks learn to associate the output with the ground truth through back-propagation. (D) For a given display, the network provides two accuracy outputs. The two streams converge toward a decision layer that compares the central and the peripheral accuracy, which are predicted by both pathways, in order to decide whether to issue a saccadic or a categorical response. If the predicted accuracy in the output of the "What" network is higher than that predicted in the "Where" network, we interrupt the visual search and classify the foveal image using the "What" pathway such as to give the answer (ANS). In the other case, the position of maximal activity in the "Where" pathway serves to generate a saccade that shifts the center of gaze. For each saccade realized, the center of vision is displaced and the process is repeated.

scenes, the probability of correctly classifying the future visual field x' from a is a number between 0 and 1 that reflects the frequency of correct classifications. The putative effect of every saccade can thus be condensed in a single number, the *accuracy*, that quantifies the final benefit of issuing saccade a from the current observation x. Extended to the full action space A, this forms an accuracy map that should monitor the selection of saccades. This accuracy map can be trained by trials and errors, with the final classification success or failure used as a teaching signal. Our main assumption here is that such a *predictive accuracy map* is at the core of a realistic saccade-based vision system. From the active inference standpoint, the separation of the scene analysis in two independent tasks relies on a simple "naïve Bayes" assumption (see Method). Each processing is assumed to be realized in parallel through different pathways by analogy with the ventral and dorsal pathways in the visual pathways (see Figure 2). A first classifier is thus assigned to process only the pixels found at the center of fixation, while a second one processes the full visual field with a retina-mimetic central log-polar magnification. The first one is called the "What" network, and the second one is the "Where" network (see Figure 7 for details). This combination of a scalar drive with action selection is reminiscent of the actor/critic principle proposed for a long time in the reinforcement learning community (Sutton & Barto, 1998). In biology, the ventral and the dorsolateral division of the striatum has been suggested to implement such an actor-critic separation (Joel, Niv, & Ruppin, 2002; Takahashi, Schoenbaum, & Niv, 2008). Consistently with those findings, our central accuracy drive and peripheral action selection map can respectively be considered the "critic" and the "actor" of an accuracy-driven action selection scheme, with foveal identification/disambiguation taken as a "visual reward."

The operations that transform the initial visual data should preserve the initial retinotopic organization, so as to form a final retinotopic accuracy map. Accordingly with the visual data, the retinotopic accuracy map may thus provide more detailed accuracy predictions in the center and coarser accuracy predictions in the periphery. Finally, each different initial visual field may bring out a different accuracy map, indirectly conveying information about the retinotopic position of the target. A final action selection (motor map) should then overlay the accuracy map through a winner-takes-all mechanism (see Figure 2D), implementing the saccade selection in a biologically plausible way, as it is thought to be done in the superior colliculus, a brain region responsible for oculomotor control (Sparks & Nelson, 1987). With the saccadic motor output showing a similar log-polar compression than the visual input, the saccades should be more precise at short than at long distance, and several saccades may be necessary to precisely reach distant targets.

In practice, the "What" and 'Where" networks are both implemented in PyTorch (Paszke et al., 2019) and trained with gradient descent over multiple layers. Each network is trained and tested separately. Because the training of the "Where" pathway depends on the accuracy given by the "What" pathway (and not the reverse), we trained the latter first, though a joint learning also yielded similar results. Finally, these are evaluated in a coupled, dynamic vision setup.

## Results

#### One-saccade setup

After training, the "Where" pathway is now capable of predicting an accuracy map (see Figure 3), whose maximal argument drives the eye toward a new viewpoint with one saccade. There, a central snippet is extracted, that is processed through the "What" pathway, allowing one to predict the digit's label. Examples of this simple sequence with one saccade are presented in Figure 3, where the digits contrast parameter is set to 70% and the digits eccentricity varies between 0 and 40 pixels. The presented examples correspond to strong eccentricity cases, when the target is hardly visible on the display (Figure 3a) and almost invisible on the reconstructed input (Figure 3b). The radial maps (Figure 3c,d) respectively represent the actual and the predicted accuracy maps. The final focus (foveal area at the location of the selected saccade) is represented in Figure 3e, with cases of classification success (Figure 3A,B) and cases of classification failures (Figure 3C–E). In the case of successful detection (Figure 3A,B), the accuracy prediction is not perfect and the digit is not perfectly centered on the fovea. This "close match" still allows for a correct classification, as the digit's pixels are fully present on the fovea. The case of Figure 3B,C is interesting for it shows two cases of a bimodal prediction, indicating that the network is capable of doing multiple detections in a single pass, that is, at a glance. The case of Figure 3C corresponds to a false detection, with the true target detected still, though with a lower intensity. The case of Figure 3D is a "close match" detection that is not precise enough to correctly center the visual target. With some pixels of the digit being invisible on the fovea, the label prediction is mistaken. The last failure case (Figure 3E) corresponds to a correct localization that is harmed by a wrong label prediction, which is due to the "What" classifier inherent error rate.

To test for the robustness of our framework, the same experiment was repeated at different SNRs of the input images. Both pathways being interdependent, it is indeed crucial to disentangle the relative effect of both sources of errors in the final accuracy. By manipulating the SNR and the target eccentricity, one can precisely monitor the network detection and recognition capabilities, with a detection task ranging from "easy" (small shift, strong contrast) to "highly difficult" (large shift, low contrast). The digit recognition capability is systematically evaluated in Figure 4 for different eccentricities and different SNRs. We test the final accuracy of the system for three target contrasts conditions ranging from 30% to 70% of the maximal contrast and 10 different eccentricities ranging from 4 to 40 pixels. It is averaged over 1,000 trials both on the initial central snippet and the final central snippet (that is, at the landing of the saccade). The (transparent) orange bars provide the initial classification rate (without saccade) and the blue bars provide the final classification rate (after saccade): see Figure 4. As expected, the accuracy decreases in both cases with the eccentricity, for the targets become less and less visible in the periphery. The decrease is rapid in the presaccadic case: The accuracy drops to the baseline level for a target distance of approximately 20 pixels from the center of gaze. The postsaccadic accuracy has a much wider range, with a slow decrease up to the border of the visual display (40 pixels away from the center). When varying the target contrast, the

Daucé, Albiges, & Perrinet

8



Figure 3. (A–E) Representative samples of active vision after training the "Where" network: (A, B) classification success samples and (C–E) classification failure samples. Digit contrast set to 70%. From left to right: (a) The initial 128  $\times$  128 visual display, with blue cross giving the initial center of gaze. The visual input is retinotopically transformed and sent to the multilayer neural network implementing the "Where" pathway. (b) Magnified reconstruction of the visual input, as it shows off from the primary visual features through an inverse log-polar transform. (c, d) Color-coded radial representation of the output accuracy maps, with dark violet for the lower accuracies and yellow for the higher accuracies. The network output ("Predicted") is visually compared with the ground truth ("True"). (e) The foveal image as the 28  $\times$  28 central snippet extracted from the visual display after doing a saccade, with label prediction and success flag in the title.

presaccadic accuracy profile is scaled by the reference accuracy (obtained with a central target), whose values are approximately 92%, 82%, and 53% for contrasts of 70%, 50%, and 30%. The postsaccadic accuracy profile undergoes a similar scaling at the different contrast values, indicating the critical dependence of the global setup to the central processing reliability.

The high-contrast case (see Figure 4) provides the greatest difference between the two profiles, with an accuracy approaching 90% at the center and 60% at the periphery. This allows one to recognize digits after

one saccade in a majority of cases, up to the border of the image, from very scarce peripheral information. With decreasing target contrast, a general decrease of the accuracy is observed, both at the center and at the periphery, with about 10% decrease with a contrast of 50% and 40% decrease with a contrast of 30%. In addition, the proportion of false detections also increases as contrast decreases. At 40 pixels away from the center, the false detection rate is approximately 40% for a contrast of 0.7, 60% for a contrast of 0.5, and 80% for a contrast of 0.3 (with a recognition close to the



Figure 4. Effect of contrast and target eccentricity. The active vision agent is tested for different target eccentricities (in pixels) and different contrasts. The final classification rate is plotted as transparent orange and blue bars, which correspond respectively to the presaccadic accuracy from the central classifier (0 saccade) and the postsaccadic accuracy (1 saccade). These are plotted with respect to the target's eccentricity and averaged over 1,000 trials per eccentricity.

baseline at the periphery in that case). The difference between the initial and the final accuracies is maximal for eccentricities ranging from 15 to 30 pixels. This optimal range reflects a proportion of the visual field around the fovea where the target detection is possible, but not its identification. The visual agent knows *where* the target is, without exactly knowing *what* it is.

#### Analysis

This full covering of the  $128 \times 128$  image range is done at a much lesser cost than what would be done by a systematic image scan, as in classic computer vision. Taking *n* the number of pixels in the original image (in our case,  $n = 128 \times 128 = 16,384$ ), our log-polar encoding provides  $O(\log n)$  log-polar visual features by construction. The total visual data processed is the addition of the C pixels processed at the fovea and the  $O(\log n)$  log-polar visual features processed at the periphery. The total processing cost is thus  $O(C + \log n)$ . Taking C as a constant, the total processing cost can be  $O(\log n)$  (for constant processing times do not change the order). In the case of multiple saccades (see next section), the total cost is  $O(k \times (C + \log n))$  with k the number of saccades. If the number of saccades k is bounded by a constant K, this allows one to estimate the processing cost as  $O(K \times (C + \log n))$  in the worst case, which also resumes to  $O(\log n)$ . This is to be contrasted, for instance, with the linear cost obtained with a full convolutional scan with a window of size C and a stride of 1, which is precisely  $O(C \times n)$ . Various optimizations can of course be considered, one of which is the well-known max-pooling principle used in deep learning, but image processing (without compression loss) is generally considered linear in the size of the visual data processed (Strengert et al., 2006).

Our sublinear processing time thus justifies a strategy that may have been chosen in a variety of natural vision systems. The compromise between the urgency to detect and the need to be accurate may justify the different balances that may exist in different species. In particular, this may justify the differences observed between prey (with a less sparse cone density at the periphery) and predators (with a tendency toward denser foveal regions).

#### Multiple-saccades setup

In our simulation results, the postsaccadic accuracy is found to overtake the presaccadic accuracy *except* when the target is initially close to the center of gaze. When closely inspecting the 1- to 10-pixels eccentricity range in our first experiment (not shown), a decision frontier between a positive and a negative information gain is found at 2 to 3 pixels away from the center. Inside that range, no additional saccade is expected to be produced, and a categorical response should be given instead. It is crucial here to understand that this empirical accuracy difference can be predicted. by construction, as the difference of the maximal outputs of the "Where" and the "What" pathways. This difference-of-accuracies prediction can serve as a decision criterion before actuating the saccade, like a GO/NOGO signal. It is moreover interpretable as an approximation of the information gain provided by the "Where" pathway, with the true label log-posterior seen as a sample of the posterior entropy (see Equation 1).

After the first saccade, while the decision criterion is not attained, additional saccades may be pursued in order to search for a better centering. In the case of a false detection, for instance, the central accuracy estimate should be close to the baseline and may





Figure 5. Multiple saccade setup. (A) Example of a trial with a sequence of three saccades. The subjective visual field is reconstructed from the log-polar visual features, with the red square delineating the  $28 \times 28$  foveal snippet after zero, one, two, and three saccades (from left to right). After the first saccade, the accuracy predicted by the "Where" network is higher than that predicted by the "What" network, and a corrective saccade is realized to center the target. After this second saccade, the foveal accuracy is higher than that predicted in the periphery, and the answer ANS is given. (B) Average classification accuracies measured for different target eccentricities (in pixels) and a different number of saccades. Target contrast set to 70%. (Transparent) light orange bars: presaccadic central accuracy (0 saccade) with respect to eccentricity, averaged over 1,000 trials per eccentricity. Blue bars: final classification rate after one, two, and three saccades (from left to right, respectively).

allow to "explain away" the current center of gaze and its neighborhood, encouraging one to actuate long-range saccades toward less salient peripheral positions, making it possible to escape from initial prediction errors. This incitement to select a saccade "away" from the central position is reminiscent of a well-known phenomenon in vision known as the "inhibition of return" (Itti & Koch, 2001). Combining accuracy predictions from each pathway may thus allow one to refine saccade selection in a way that complies with the sequential processing observed in biological vision. Note that extended to a multitarget case, the information gain maximization principle still holds as a general measure of the improvement of scene understanding through multiple saccades. It is uncertain, however, whether biologically realistic implementations would be possible in that case. In particular, we predict that such a mechanism is dependent on the class of inputs and would be, for instance, different when searching for faces as compared to digits.

Some of the most peripheral targets are thus difficult to detect in just one saccade, resulting in degraded performances at the periphery (see Figure 4). Even when correctly detected, our log-polar action maps also preclude precise centering. As a consequence, peripheral targets are generally poorly centered after the first saccade, as shown, for instance, in Figure 3D, resulting in classification errors. The possibility to perform a sequential search using more saccades is thus crucial to allow for a better recognition. Results on multisaccade visual search results are presented in Figure 5.

An example of a trial with a sequence of three saccades is shown in Figure 5A. A hardly visible peripheral target (digit) is first approximately shifted to the foveal zone thanks to the first saccade. Then, a new retinal input centered at the new point of fixation is computed, such that it generates a novel predicted accuracy map. The second saccade allows one to improve the target centering. As the predicted foveal accuracy given by the "What" network is higher than the peripheral one given by the "Where" network, a third saccade would not improve the centering: The stopping criterion is met. In practice, one or two saccades were sufficient in most trials to reach the actual target. Another behavior was also observed for some "bad start," which exhibited a false localization (as in Figure 3C, for instance), when the target is shifted away in the opposite direction and the agent cannot recover from its initial error. From Figure 5B, this case can be estimated at about 15% of the cases for the most peripheral targets.

Overall, as shown in Figure 5B, the corrective saccades implemented in this multiple saccade setup provide a significant improvement in the classification accuracy. Except at the center, the accuracy increases by about 10% both for the midrange and the most peripheral eccentricities. Most of the improvement, however, is provided by the first corrective saccade. The second corrective saccade only shows a barely significant improvement of about 2%, which is only visible at the periphery. The following saccades would mostly implement target tracking, without providing additional accuracy gain. A three-saccade setup finally allows a wide covering of the visual field, providing a close to central recognition rate at all eccentricities, with the residual peripheral error putatively corresponding to the "bad start" target miss cases.

## Discussion

#### Summary

In summary, we have proposed a visuomotor action-selection model that implements a focal accuracy-seeking policy across the image. Our main modeling assumption here is an *accuracy-driven* monitoring of action, stating in short that the ventral classification accuracy drives the dorsal selection in building an extrafoveal accuracy map. The comparison of both accuracies amounts either to select a saccade or to keep the eye focused at the center, so as to identify the target. The predicted accuracy map has, in our case, the role of a value-based action selection map, as it is the case in model-free reinforcement learning. However, it also owns a probabilistic interpretation, making it possible to combine concurrent accuracy predictions, such as the ones done through the "What" and the "Where" pathways. This allows one in particular to explain more elaborate aspects of the whole decision-making processes, such as the inhibition of return (Itti & Koch, 2001), without further specific heuristic.

Moreover, one crucial aspect highlighted by our model is the importance of centering objects in recognition. Despite the robust translation invariance observed on the "What" pathway, a small tolerance radius of about 4 pixels around the target's center needs to be respected to maximize the classification accuracy. The translation invariance is in our case an effect of the max-pooling operations in the convolutional layers, built in at the core of the "What" layer. This relates to the idea of finding an absolute referential for an object, for which the recognition is easier. If the center of fixation is fixed, the log-polar encoding of an object has the notable properties to map object rotations and scalings toward translations in the radial and angular directions of the visual domain (Traver & Bernardino, 2010). Extensions to scale and rotation invariance would in principle be feasible through central log-polar encoding, with little additional computational cost. This prospect is left for future work.

11

#### Comparison with other models

Daucé, Albiges, & Perrinet

A lot of computer models found in the literature reflect to some degree the foveal/sequential visual processing principles developed here. Since the question of a normative and quantitative comparison with them is important, no specific or unified data set is proposed at present to address this specific case. Every model found uses a different retinal encoding, different computing methodologies, and different training data sets. We thus provide here a qualitative comparison with the more prominent computer-based focal vision models proposed in the literature.

First, active vision is of course an important topic in mainstream computer vision. In the case of image classification, it is considered a way to improve object recognition by progressively increasing the definition over identified regions of interest, referred as "recurrent attention" (Mnih et al., 2014; Fu et al., 2017). Standing on a similar mathematical background, recurrent attention is, however, at odds with the functioning of biological systems, with a mere distant analogy with the retinal principles of foveal-surround visual definition. Phenomenological models, such as the one proposed in Najemnik and Geisler's seminal study (Najemnik & Geisler, 2005), rely on a rough simplification, with foveal center-surround acuity modeled as a response curve. Despite providing a bio-realistic account of sequential visual search, the model owns no foveal image-processing implementation. Stemming on Najemnik and Geisler's principles, a trainable center-surround processing system was proposed in Butko and Movellan (2010), with a sequential scan of an image in a face-detection task. However, the visual search task relies there on a systematic scan over a dynamically blurred image, with all the visual processing delegated to standard feature detectors.

In contrast, the Akbas and Eckstein (2017) model (foveated object detector) uses an explicit bio-inspired log-polar encoding for the peripheral processing, with trainable local features. With a focus put on the processing effectiveness provided by this specific compression, the model approaches the performance of state-of-the-art linear feature detectors, with multiscale template matching (bounding box approach). However, the use of a local/linear template matching processing makes here again the analogy with the brain oversimplistic.

Denil et al.'s (2012) article is probably the one that shows the closest correspondence with our setup. It owns an identity pathway and a control pathway, in a What/Where fashion, just as ours. Interestingly, only the "What" pathway is neurally implemented using a random foveal/multifixation scan within the fixation zone. The "Where" pathway, in contrast, mainly implements object tracking, using particle filtering with a separately learned generative process. The direction of gaze is here chosen so as to minimize the target's position, speed, and scale uncertainty, using the variance of the future beliefs as an uncertainty metric. The control part is thus much similar to a dynamic region of interest (ROI) tracking algorithm, with no direct correspondence with foveal visual search or with the capability to recognize the target

#### Perspectives

We have thus provided here a proof of concept that a log-polar retinotopy can efficiently serve object detection and identification over wide visual displays. Despite its simplicity, the model used to generate our visual display allowed to assess the effectiveness and robustness of our learning scheme, which should be extended in the future to more complex displays and more realistic closed-loop setups. In particular, the restricted  $28 \times 28$  input used for the foveal processing is a mere placeholder that should be replaced by more elaborate computer vision frameworks, such as Inception (Szegedy et al., 2015) or VGG-19 (Simonyan & Zisserman, 2014), that can handle more ecological natural image classification setups.

The main advantage of our peripheral image processing is its cost-efficiency. Our full log-polar processing pathway consistently conserves the high compression rate performed by retina and V1 encoding up to the action selection level. The organization of both the visual filters and the action maps in concentric log-polar elements, with radially exponentially growing spatial covering, can thus serve as a baseline for a future sublinear (logarithmic) complexity for visual search in computer vision. Our work thus illustrates one of the main advantages of using a focal/sequential visual processing framework, which is providing a way to process large images with a sublinear processing cost. This may allow one to detect an object in large visual environments, which should be particularly beneficial when the computing resources are under constraint, such as for drones or mobile robots.

If the methodology and principles developed here are clearly intended to deal with real images, an important contribution of the study is providing principles that justify the separation between a ventral and a dorsal stream in the early visual pathways. If some forms of "dual-pathway models" have been proposed in the past through separating the central and the peripheral processing, as in Denil et al. (2012), and also in one instance of the Akbas and Eckstein (2017) model, their guiding principles stem on computing efficacy rather than biological fidelity. We thus think that our principled ventral/dorsal concurrent processing, rooted on dorsal accuracy map predictions, is both important and novel.

Finally, our approach relies on a strong idealization, assuming the presence of a unique target. This is well adapted to a fast-changing visual scene as is demonstrated by our ability to perform as fast as five saccades per second to detect faces in a cluttered environment (Martin, Davis, Riesenhuber, & Thorpe, 2018). However, some visual scenes—such as when looking at a painting in a museum—allow for a longer inspection of their details. The presence of many targets in a scene should be addressed, which amounts to sequentially select targets, in combination with implementing a more elaborate inhibition of return mechanism to account for the trace of the performed saccades. This would generate more realistic visual scan-paths over images. Actual visual scan-paths over images could also be used to provide priors over action selection maps that should improve realism. Identified regions of interest may then be compared with the baseline bottom-up approaches, such as the low-level feature-based saliency maps (Itti & Koch, 2001). Maximizing the information gain over multiple targets needs to be envisioned with a more refined probabilistic framework extending previous models (Friston et al., 2012), which would include phenomena such as mutual exclusion over overt and covert targets. Next, scan-paths could be generated on actual dynamical scenes, despite existing oculomotor delays (Perrinet et al., 2014), extending such a framework in the temporal domain. How the brain may combine and integrate these various probabilities dynamically is still an open question that amounts to the fundamental binding problem.

## Methods

#### Image generation

We first define here the generative model for input display images as shown first in Figure 1A (DIS) and

as implemented in Figure 2A. Following a common hypothesis regarding active vision, visual scenes consist of a single target embedded in a large image with a cluttered background.

#### Targets

We use the MNIST data set of handwritten digits introduced by Lecun et al. (1998): Samples are drawn from the data set of 60,000 grayscale,  $28 \times 28$  pixel images. There are 10 categories (from "zero" to "nine"). These are separated between a training and a validation set (see below the description of the "Where" network). Note that to simplify the task, there is one and only one target per image.

#### Full-scale images

We call "full-scale images" input images that correspond to a discretized, rectangular sampling (pixels) of the visual field, which are usually used in computer vision. These input images are set to a size of  $128 \times 128$  pixels in which we embed the target. Each target location is drawn at random in this large image. To enforce isotropic generation (at any direction from the fixation point), a centered circular mask covering the image (of radius 64 pixels) is defined. Also, the target's location is such that the embedded sample fits entirely into that circular mask.

#### Background noise setting

To implement a realistic background noise, we generate synthetic textures (Sanz-Leon et al., 2012) using a bidimensional random process. The texture is designed to fit with the statistics of natural images. We chose an isotropic setting where textures are characterized by solely two parameters, one controlling the median spatial frequency of the noise, the other controlling the bandwidth around the central frequency. Equivalently, this can be considered the band-pass filtering of a random white noise image. The spatial frequency is set at 0.1 pixel<sup>-1</sup> to fit that of the original digits. This specific spatial frequency occasionally allows one to generate some "phantom" digit shapes in the background. Finally, these images are rectified to have a normalized contrast.

#### Mixing the signal and the noise

Finally, both the noise and the target image are merged into a single image. Two different strategies are used. A first strategy emulates a transparent association, with an average luminance computed at each pixel, while a second strategy emulates an opaque association, choosing for each pixel the maximal value. The quantitative difference was tested in simulations but proved to have a marginal importance, and results shown here are the result of the opaque association.

#### Active inference and the Naïve Bayes assumption

Saccade selection in visual processing can be captured by a statistical framework called a partially observed Markov decision process (POMDP) (Najemnik & Geisler, 2005; Butko & Movellan, 2010; Friston et al., 2012), where the cause of a visual scene is made up from the couple of independent random variables of the viewpoint and of the scene elements (here a digit). For instance, changing the viewpoint will lead to a different scene rendering. A generative model tells how the visual field should look knowing the scene elements and a certain viewpoint. In general, active inference assumes a hidden external state e, which is known indirectly through its effects on the sensor. The external state corresponds to the physical environment. Here the external state is assumed to split into two (independent) components, namely, e = (u, y), with u the interoceptive body posture (in our case, the gaze orientation, or "viewpoint") and y the object shape (or object identity). The visual field x is the state of the sensors, that is, a partial view of the visual scene, measured through the generative process:  $x \sim p(X|e)$ .

Using the Bayes rule, one may then infer the scene elements from the current viewpoint (model inversion). With the real physical state *e* being hidden, a parametric model  $\theta$  is assumed to allow for an estimate of the cause of the current visual field through model inversion thanks to the following Bayes formula:  $p(E|x) \propto p(x|E;\theta)$ . It is also assumed that a set of motor commands  $A = \{..., a, ...\}$  (here saccades) may control the body posture (gaze orientation) but not the object's identity, so that y is independent of a. Actuating a command a changes the viewpoint to u', which feeds the system with a new visual sample  $x' \sim p(X|u', y)$ . The more viewpoints you have, the more certain you are about the object identity through a chain rule sequential evidence accumulation.

In an optimal search setup, however (Najemnik & Geisler, 2005), you need to choose the next viewpoint that will help you to disambiguate *at best* the scene. In a predictive setup, the consequence of every saccade should be analyzed through model inversion *over the future observations*, that is, predicting the effect of every action to choose the one that may optimize future inferences. The benefit of each action should be quantified through a certain metric (future accuracy, future posterior entropy, future variational free energy, ...) that depends on the current inference p(U, Y|x). The saccade *a* that is selected thus provides a new

visual sample from the scene statistics. If well chosen, it should improve the understanding of the scene (here the target position and category). However, estimating in advance the effect of every action over the range of every possible object shape and body posture is combinatorially hard, even in simplified setups, and thus infeasible in practice.

The predictive approach necessitates in practice to restrain the generative model in order to reduce the range of possible combinations. One such restriction, known as the "naïve Bayes" assumption, considers the independence of the factors that are the cause of the sensory view. The independence hypothesis allows considering the viewpoint u and the category y being independently inferred from the current visual field, that is, p(U, Y|x) = p(U|x)p(Y|x). This property is strictly true in our setting and is very generic in vision for simple classes (such as digits) and simple displays (but see Võ & Wolfe, 2012, for more complex visual scene grammars).

#### Foveal vision and the "What" pathway

At the core of the vision system is the identification module, that is, the "What" pathway (see Figure 2). It consists of a classic convolutional classifier for which we will show some translation invariance in the form of a shift-dependent accuracy map. Importantly, it can quantify its own classification uncertainty, which may allow comparisons with the output of the "Where" pathway.

The foveal input is defined as the  $28 \times 28$  grayscale image cropped at the center of gaze (see dashed red box in Figure 1C). This image is passed unmodified to the agent's visual categorical pathway (the "What" pathway), which is realized by a convolutional neural network, here the well-known "LeNet" classifier (Lecun et al., 1998). The network structure that processes the input to identify the target category is made of three convolution layers interleaved with max-pooling layers, followed by two fully connected layers as provided (and unmodified) by its PyTorch library implementation (Paszke et al., 2019). Each intermediate layer's output is rectified, and the network output uses a sigmoid operator to predict the likelihood of detecting each of the 10 digits. The index of one of the 10 output neurons with maximum probability provides the image category. It is first trained over the (centered) MNIST data set after approximately 20 training epochs. This strategy achieves an average 98.7% accuracy in the center on the validation data set (Lecun et al., 1998).

To achieve an even more generic "What" pathway, a specific data set is constructed to train the network. It is made of randomly shifted digits overlaid on a randomly generated noisy background, as defined above. Both the shift and the contrast relative to the



Figure 6. (A) Input samples from the "What" training set, with randomly shifted targets using a Gaussian bivariate spatial shift with a standard deviation of 15 pixels. The target contrast is randomly set between 30% and 70%. (B)  $55 \times 55$  shift-dependent accuracy map, measured for different target eccentricities on the test set after training.

background noise make the task more difficult than the original MNIST categorization. The relative contrast of the digit is randomly set between 30% and 70% of the maximal contrast. The network is trained incrementally by progressively increasing the shift variability (of a bivariate central gaussian) and by increasing the standard deviation from 0 to 15 (with a maximal shift set at 27 pixels). The network is trained on a total of 75 epochs, with 60,000 examples generated at each epoch from the original MNIST training set, using the cross-entropy loss. The shifts and backgrounds are regenerated at each epoch. The shifts' standard deviation increases one unit every five epochs such that at the end of the training, many digits fall outside the center of the fovea, so that many examples are close to impossible to categorize, either because of a low contrast or a too large eccentricity. At the end of the training process, the average accuracy is thus 34% and a maximum accuracy of 91% at the center.

After training, this shift-dependent accuracy map is validated by systematically testing the network accuracy on every horizontal and vertical shift, each on a set of 1,000 cluttered target samples generated from the MNIST test set and within the range of  $\pm 27$ pixels (see Figure 6). This forms a  $55 \times 55$  accuracy map showing higher accuracy at the center and a slow decreasing accuracy with target eccentricity (with an accuracy plateau over 70% showing a relative shift invariance on around a 7-pixel eccentricity radius). This shift invariance is a known effect of convolutional computation. Note that the categorization task is here harder by construction and the accuracy that is obtained here is lower (with a central recognition rate of around 80%). The accuracy sharply drops for eccentricities greater than 10 pixels, reaching the baseline 10% chance level at shift amplitudes at around 20 pixels.



Figure 7. Implementing the "Where" pathway: (A) A visual display is transformed by a feature vector, in which elements compute the similarity of the full image with a bank of oriented filters placed at positions defined by a log-polar grid. This defines a linear transform of the  $128 \times 128 = 16,384$  input pixels into 2,880 coefficients. It is possible to represent this information in visual space by using the pseudo-inverse of the linear transform (see, for instance, Figure 1C). (B) The "Where" network consists of two hidden layers composed with a RELU operator transforming the retinal feature vector. A sigmoid operator ensures that this output vector is a distribution of predicted likelihoods in log-polar space. (C) Similar to (A), any full accuracy map computed by shifting the pre-computed shift-dependent accuracy map of the "What" pathway (see Figure 6) can be transformed into a distribution in log-polar space, similar to a collicular representation. As the full accuracy map is itself a distribution, this can be implemented by a linear (matrix) transform. In practice, one can use the inverse of this linear transform to project any collicular representation into the visual space, for instance, to predict for the position with maximal accuracy (red cross).

#### "Where" pathway: Transforming log-polar feature vectors to log-polar action maps

Here, we assume the "Where" implements the following action selection: where to look next in order to reduce the uncertainty about the target identity? The "Where" pathway is thus devoted to choosing the next saccade by predicting the location of the target in the (log-polar) visual field. This implies moving the eye to increase the "What" categorization accuracy. For a given visual field, each possible future saccade has an expected accuracy that can be trained from the "What" pathway output. To accelerate the training, we use an equivalent strategy by training the network on a translated accuracy map (see below for details). The output is thus an accuracy map that tells for each possible visuomotor displacement the value of the future accuracy.

# Primary visual representation: log-polar orientation filters

In order to reduce the processing cost, and in accordance with observations (Connolly & Van Essen, 1984; Sparks & Nelson, 1987), a similar log-polar compression pattern is assumed to be conserved from the retina up to the visuomotor layers. The nonuniform sampling of the visual space is adequately modeled as a log-polar conformal mapping, as it provides a good fit with observations in mammals (Traver & Bernardino, 2010) and has a long history in computer vision and robotics. Both the visual features and the output accuracy map are to be expressed in retinal coordinates. On the visual side, local visual features are extracted as oriented edges as a combination of the retinotopic transform with primary visual cortex filters (Fischer, Sroubek, Perrinet, Redondo, & Cristobal, 2007); see Figure 7A. The centers of these first- and

second- order orientation filters are radially organized around the center of fixation, with small and tightened receptive fields at the center and more large and scarce receptive fields at the periphery. The size of the filters increases proportionally to the eccentricity. The filters are organized in 10 spatial eccentricity scales (respectively placed at around 2, 3, 4.5, 6.5, 9, 13, 18, 26, 36.5, and 51.3 pixels from the center) and 24 different azimuth angles, allowing them to cover most of the original  $128 \times 128$  image. At each of these positions, six different edge orientations and two different phases (symmetric and antisymmetric) are computed. This finally implements a (fixed) bank of linear filters, which models the receptive fields of the primary visual cortex.

To ensure the balance of the coefficients across scales, the images are first whitehed and then linearly transformed into a retinal input as a feature vector x. The length of this vector is 2,880, such that the retinal filter compresses the original image by about 83%, with high spatial frequencies preserved at the center and only low spatial frequencies conserved at the periphery. In practice, the bank of filters is precomputed and placed into a matrix for a rapid transformation of input batches into feature vectors. This matrix transformation allows also the evaluation of a reconstructed visual image given a retinal activity vector thanks to a pseudo-inverse of the forward transform matrix. In summary, the full-sized images are transformed into a primary visual feature vector, which is fed to the "Where" pathway.

#### Visuomotor representation: "Collicular" accuracy maps

The output of the "Where" pathway is defined as an accuracy map representing the recognition probability after moving the eye, independently of its identity. Like the primary visual map, this target accuracy map is also organized radially in a log-polar fashion, making the target position estimate more precise at the center and fuzzier at the periphery. This modeling choice is reminiscent of the approximate log-polar organization of the superior colliculus (SC) motor map (Sparks & Nelson, 1987). To ensure that this output is a distribution function, we use a sigmoid operator at the output of the "Where" network. In ecological conditions, this accuracy map should be trained by sampling, that is, by "trial and error," using the actual recognition accuracy (after the saccade) to grade the action selection. For instance, we could use corrective saccades to compute (a posteriori) the probability of a correct localization. In a computer simulation, however, this induces a combinatorial explosion that does render the calculation not amenable.

In practice, as we designed the generative model for the visual display, the position of the target (which is hidden to the agent) is known. Combining this translational shift and the shift-dependent accuracy map of the "What" classifier (Figure 6B), the full accuracy map at each pixel can be thus predicted for each visual sample under an ergodic assumption by shifting the central accuracy map on the true position of the target (see Figure 7C). Such a computational shortcut is allowed by the independence of the categorical performance with position. This full accuracy map is a probability distribution function defined on the rectangular grid of the visual display. We project this distribution on a log-polar grid to provide the expected accuracy of each hypothetical saccade in a retinotopic space similar to a collicular map. In practice, we used Gaussian kernels defined in the log-polar space as a proxy to quantify the projection from the metric space to the retinotopic space. This generates a filter bank at 10 spatial eccentricities and 24 different azimuth angles, that is, 240 output filters. To ensure keeping a distribution function, each filter is normalized such that the value at each log-polar position is the average of the values that are integrated in visual space. Applied to the full-sized ground truth accuracy map computed in metric space, this gives an accuracy map at a different location of a retinotopic motor space.

#### **Classifier training**

The "Where" pathway is a function transforming an input retinal feature vector x into an output log-polar retinotopic vector *a* representing for each area of the log-polar visual field a prediction of the accuracy probability. Following the active inference framework, the network is trained to predict the likelihood  $a_i$  at position *i* knowing the retinal input *x* by comparing it to the known ground truth distribution computed over the motor map. The loss function that comes naturally is the binary cross-entropy. At each individual position *i*, this loss corresponds to the negative term of Kullback-Leibler divergence for a binomial random variable  $a_i$  given by the predicted map and the ground truth (see Figure 7B). The total loss is the average over all positions *i*. This scalar measures the distance between both distributions; it is always positive and null if and only if they are equal.

The parametric neural network consists of a primary visual input layer, followed by two fully connected hidden layers of size 1,000 with rectified linear activation and a final output layer with a sigmoid nonlinearity to ensure that the output is compatible with a likelihood function (see Figure 7B). An improvement in convergence speed was obtained by using batch normalization. The network is trained on 60 epochs of 60,000 samples, with a learning rate equal to  $10^{-4}$  and the Adam optimizer (Kingma & Ba, 2014) with standard momentum parameters. One full training takes about 1 hr on a laptop. The code is written in Python (Version 3.7.6) with PyTorch library (Paszke et al., 2019) (Version 1.1.0). The full



Figure 8. Quantitative role of parameters: We tested all parameters of the presented model, from that controlling the architecture of image generation to the parameters of the neural network implementing the "Where" pathway (including meta-parameters of the learning paradigm). We show here the results that show the most significative impact on average accuracy. The accuracy is given by a blue line, with the red line giving the rate of errors. The black dashed line gives the chance level (10%), while the blue box gives the 99% confidence interval as estimated over eight repetitions of the learning. (A) First, we tested some properties of the input, respectively from left to right: noise level (Noise), standard deviation of the distance of the target with respect to the fixation (Offset\_std), mean spatial frequency of clutter (Sf\_0), and bandwidth (B\_sf) of the clutter noise. This shows that average accuracy evolves with noise (see also Figure 4 for an evolution as a function of eccentricity) but also the characteristics of the noise clutter. In particular, there is a drop in accuracy whenever noise is of similar wavelength as digits, but it becomes less pronounced as the bandwidth increases. (B) Finally, we scanned parameters of the deep learning neural network. We observed that accuracy quickly converged after approximately 25 epochs (Epochs\_adam). We then tested different values for the dimension of respectively the first (Dim1) and second (Dim2) hidden layers, showing weak changes in accuracy. (C) The accuracy also changes with the architecture of the foveated input, as shown here by changing the number N\_azimuth of azimuth directions, which necessitates a longer training phase, such that the optimal number is around 24 azimuth directions.

scripts for reproducing the figures and exploring the results to the full range of parameters are available at https://github.com/laurentperrinet/WhereIsMyMNIST.

#### Quantitative role of parameters

In addition, we controlled that the training results are robust to changes in an individual experimental or network parameter from the default parameters (see Figure 8). From the scan of each of these parameters, the following observations were remarkable. First we verified that accuracy decreased when noise increased and while the bandwidth of the noise imported weakly, the spatial frequency of the noise was an important factor. In particular, final accuracy was worst for a clutter spatial frequency of  $\approx 0.07$ , that is, when the characteristic texture elements were close to the characteristic size of the objects. Second, we saw that the dimension of the "Where" network was optimal for a dimensionality similar to that of the input but that this mattered weakly. The dimensionality of the log-polar map is more important. The analysis proved that an optimal accuracy was achieved when using a number of 24 azimuthal directions. Indeed, a finer log-polar grid requires more epochs to converge and may result in an overfitting phenomenon, hindering the final accuracy. Such fine-tuning of parameters may prove to be important in practical applications and to optimize the compromise between accuracy and compression.

#### **Concurrent action selection**

Finally, when both pathways are assumed to work in parallel, each one may be used concurrently to choose

the most appropriate action. Two concurrent accuracies are indeed predicted through separate processing pathways, namely, the central pixel recognition accuracy through the "What" pathway and the log-polar accuracy map through the "Where" pathway. The central accuracy may thus be compared with the maximal accuracy, as predicted by the "Where" pathway.

From the information theory standpoint, each saccade comes with fresh visual information about the visual scene that can be quantified by a conditional *information gain*, namely:

$$IG_{max} = \max_{x'} \log p(y|x, x') - \log p(y|x)$$

with the left term representing the future accuracy (after the saccade is realized) and the right term representing the current accuracy as it is obtained from the "What" pathway. Estimating the joint conditional dependence in the first term being once again out of reach for computational reasons, the following approximative estimate is used instead:

$$\tilde{IG}_{max} \simeq IG_{max} = \max_{x'} \log p(y|x') - \log p(y|x), \quad (1)$$

which is a simple difference between the log accuracy after the saccade minus the log accuracy before the saccade. To provide a reliable estimate, the information gain may be averaged over many saccades and many target eccentricities (so that the information gain may be close to zero when the target eccentricity is close to zero). Since the saccade is subject to prediction errors and execution noise, the saccade landing position may be different from the initial prediction. The final accuracy, as instantiated in the accuracy map, contains this intrinsic imprecision and is thus necessarily lower than the optimal one. The consequence is that in some cases, the approximate information gain may become negative, when the future accuracy is actually lower than the current one. This is, for instance, the case when the target is exactly positioned at the center of the fovea.

Keywords: object detection, active inference, visual search, visuomotor control, deep learning

## Acknowledgments

Commercial relationships: none. Corresponding author: Laurent U. Perrinet. Email: laurent.perrinet@univ-amu.fr. Address: Institut de Neurosciences de la Timone (UMR 7289) Aix Marseille Université, CNRS Faculté de Médecine - Bâtiment Neurosciences, 27, Bd Jean Moulin, Marseille, PACA 13385 Marseille Cedex 05.

#### References

- Akbas, E., & Eckstein, M. P. (2017). Object detection through search with a foveated visual system. *PLoS Computational Biology*, 13, e1005743, doi:10.1371/journal.pcbi.1005743.
- Bahill, A. T., Clark, M. R., & Stark, L. (1975). The main sequence, a tool for studying human eye movements. *Mathematical Biosciences*, 24, 191–204, doi:10/b655hq.01092.
- Bruce, N. D. B., & Tsotsos, J. K. (2009). Saliency, attention, and visual search: An information theoretic approach. *Journal of Vision*, 9(3), 5, doi:10/bzzq2x.
- Butko, N. J., & Movellan, J. R. (2010). Infomax control of eye movements. *IEEE Transactions on Autonomous Mental Development*, 2, 91–107.
- Connolly, M., & Van Essen, D. (1984). The representation of the visual field in parvicellular and magnocellular layers of the lateral geniculate nucleus in the macaque monkey. *Journal of Comparative Neurology, 226*, 544–564.
- Daucé, E. (2018). Active fovea-based vision through computationally-effective model-based prediction. *Frontiers in Neurorobotics*, 12, 76.
- Denil, M., Bazzani, L., Larochelle, H., & de Freitas, N. (2012). Learning where to attend with deep architectures for image tracking. *Neural Computation, 24*, 2151–2184.
- Feng, X., Jiang, Y., Yang, X., Du, M., & Li, X. (2019). Computer vision algorithms and hardware implementations: A survey. *Integration*, 69, 309–320.
- Fischer, S., Sroubek, F., Perrinet, R., Redondo, L.U., & Cristobal, G. (2007). Self-invertible 2D log-Gabor wavelets. *International Journal of Computer Vision*, 75, 231–246, doi:10.1007/s11263-006-0026-8.
- Friston, K. J., Adams, R. A., Perrinet, L. U., & Breakspear, M. (2012). Perceptions as hypotheses: Saccades as experiments. *Frontiers in Psychology*, 3, doi:10.3389/fpsyg.2012.00151.
- Fu, J., Zheng, H., & Mei, T. (2017). Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 4438–4446). Honolulu, Hawaii: IEEE.
- He, K., Zhang, X., Ren, S., & Sun, J. (2015, February). Delving deep into rectifiers: Surpassing humanlevel performance on ImageNet classification. arXiv:1502.01852 [cs], 04208.
- Itti, L., & Baldi, P. (2009). Bayesian surprise attracts human attention. *Vision Research*, 49, 1295–1306, doi:10.1016/j.visres.2008.09.007.

- Itti, L., & Koch, C. (2001). Computational modelling of visual attention. *Nature Reviews Neuroscience*, 2, 194–203, doi:10/chw2bk.
- Joel, D., Niv, Y., & Ruppin, E. (2002). Actor–critic models of the basal ganglia: New anatomical and computational perspectives. *Neural Networks*, 15, 535–547.
- Kaplanyan, A., Sochenov, A., Leimkühler, T., Okunev, M., Goodall, T., & Rufo, G. (2019). DeepFovea: Neural reconstruction for foveated rendering and video compression using learned statistics of natural videos. ACM Transactions on Graphics (TOG), 38 (pp. 730–742). New York, NY, USA: Association for Computing Machinery, doi:10.1145/3355089.3356557.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv. arXiv:1412.6980.
- Kirchner, H., & Thorpe, S. J. (2006). Ultra-rapid object detection with saccadic eye movements: Visual processing speed revisited. *Vision Research*, 46, 1762–1776, doi:10.1016/j.visres.2005.10.002.
- Kortum, P., & Geisler, W. S. (1996). Implementation of a foveated image coding system for image bandwidth reduction. *In Human Vision and Electronic Imaging*, 2657, 350–360.
- Kummerer, M., Wallis, T. S. A., Gatys, L. A., & Bethge, M. (2017). Understanding low- and high-level contributions to fixation prediction. In: *Proceedings* of the IEEE International Conference on Computer Vision (pp. 4789–4798). Venice, Italy: IEEE.
- Lecun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86, 2278–2324, doi:10/d89c25.
- Martin, J. G., Davis, C. E., Riesenhuber, M., & Thorpe, S. J. (2018). Zapping 500 faces in less than 100 seconds: Evidence for extremely fast and sustained continuous visual search. *Scientific Reports*, 8, 1–12, doi:10/gd5x5d. 00006.
- Mirza, M. B., Adams, R. A., Mathys, C., & Friston, K. J. (2018). Human visual exploration reduces uncertainty about the sensed world. *PLoS One*, 13, e0190429, doi:10.1371/journal.pone.0190429.
- Mishkin, M., Ungerleider, L. G., & Macko, K. A. (1983). Object vision and spatial vision: Two cortical pathways. *Trends in Neurosciences*, 6, 414–417.
- Mnih, V., Heess, N., Graves, A., & Kavukcuoglu, K. et al. (2014). Recurrent models of visual attention. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, & K. Q. Weinberger (Eds.), Advances in neural information processing systems (pp. 2204–2212). Curran Associates, Inc.
- Najemnik, J., & Geisler, W. S. (2005). Optimal eye movement strategies in visual search. *Nature, 434*, 387–391, doi:10/bcbw2b.

- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., & Chanan, G., ... Chintala, S. (2019). PyTorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett (Eds.), Advances in neural information processing systems 32 (pp. 8024–8035). Long Beach Convention Center, Long Beach: Curran Associates, Inc.
- Perrinet, L. U., Adams, R. A., & Friston, K. J. (2014). Active inference, eye movements and oculomotor delays. *Biological Cybernetics*, 108, 777–801, doi:10.1007/s00422-014-0620-8.
- Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 779–788). Las Vegas, Nevada, USA: IEEE.
- Ren, S., He, K., Girshick, R., & Sun, J. (2017). Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence, 39*, 1137–1149, doi:10/gc7rmb.08051.
- Samonds, J. M., Geisler, W. S., & Priebe, N. J. (2018). Natural image and receptive field statistics predict saccade sizes. *Nature Neuroscience*, 21, 1591, doi:10/gfgt3k.00002.
- Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L.-C. (2018, January). *MobileNetV2: Inverted residuals and linear bottlenecks*. arXiv:1801.04381 [cs], 00241.
- Sanz-Leon, P., Vanzetta, I., Masson, G. S., & Perrinet, L. U. (2012). Motion clouds: Model-based stimulus synthesis of natural-like random textures for the study of motion perception. *Journal of Neurophysiology*, 107, 3217–3226, doi:10.1152/jn.00737.2011.
- Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv. arXiv:1409.1556.
- Sparks, D. L., & Nelson, I. S. (1987). Sensory and motor maps in the mammalian superior colliculus. *Trends in Neurosciences*, 10, 312–317.
- Strasburger, H., Rentschler, I., & Jüttner, M. (2011). Peripheral vision and pattern recognition: A review. *Journal of Vision, 11*, 13, doi:10/fx28dx. 00489.
- Strengert, M., Kraus, M., & Ertl, T. (2006). Pyramid methods in gpu-based image processing. In *Proceedings vision, modeling, and visualization 2006* (pp. 169–176). Aachen, Germany: IOS Press.
- Sutton, R. S., & Barto, A. G. (1998). Reinforcement learning: An introduction (Vol. 1). Cambridge, MA: MIT.

- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1–9). Boston, Massachusetts, USA: IEEE.
- Takahashi, Y., Schoenbaum, G., & Niv, Y. (2008). Silencing the critics: Understanding the effects of cocaine sensitization on dorsolateral and ventral striatum in the context of an actor/critic model. *Frontiers in Neuroscience*, 2, 14.
- Traver, V. J., & Bernardino, A. (2010). A review of log-polar imaging for visual perception in robotics. *Robotics and Autonomous Systems*, 58, 378–398.

- Treisman, A. M., & Gelade, G. (1980). A featureintegration theory of attention. *Cognitive Psychology*, 12, 97–136, doi:10.1016/0010-0285(80) 90005.11957.
- Võ, M. L.-H., & Wolfe, J. M. (2012). When does repeated search in scenes involve memory? Looking at versus looking for objects in scenes. *Journal of Experimental Psychology: Human Perception and Performance, 38*, 23–41, doi:10.1037/a0024147.
- Zhang, L., Tong, M. H., Marks, T. K., Shan, H., & Cottrell, G. W. (2008). SUN: A Bayesian framework for saliency using natural statistics. *Journal of Vision*, 8(7), 32, doi:10/fcdh2d.