

Learning where to look for a target before seeing it

\begin{abstract}

A realistic system for object categorization should be able to find the target in possibly large images, independently of its position in visual space. Current solutions leverage this solution by processing the different hypothesis (classes) at all possible spatial configuration. However, this can be costly in terms of computing time especially without dedicated parallel hardware. % We explore here a solution inspired by the anatomy of the human visual system, that is, the combination of a foveated sensor with the capacity of rapidly moving the center of fixation using saccades. Indeed, the position and category of objects in images are a priori independent and we hypothesize that the retinotopic map overlays a central area dedicated to object categorization and a peripheral area dedicated to

Using this hypothesis, we formalize this problem in a probabilistic setting which allows us to build two parallel but interactively connected system: a classical image classification algorithm assuming that gaze is centered on the object on one side and a system learning to infer the position of the target on the other. Until the classification is confident, the system performs a saccade to the most likely position in the image. Overall, the computational cost of this strategy is less than that in holistic methods. %

We tested this framework on a simple task of finding digits in a large, cluttered image. Results demonstrate that it is possible to correctly learn the position of a target from a given class, and this before actually seeing a foveated image of the target. We compare the results of this model with classical psychophysical results in visual search. This provides evidence of the importance of such strategies in computer vision and we highlight some predictions of our model.

\end{abstract}

Notes :

- differential processing
 - periphery : low resolution in space
high in time
 - fovea : high spatial resolution
low temporal
- saccadic suppression
 - Ziad Hafez : elevation of perceptual thresholds at the time of saccade
 - frequency tuning in SC

Notes existantes

2017-09-28

- manu → papier frontié
- coarse periphery
- Hermann grid ?
- invariances

2018-03-15 sujet pierre

- def sujet / protocole
- def méthode avec carte accès
- oral

2018-12-21

- outline paper
- Figures
- yaka

① INTRODUCTION

1a Active Inference in Machine Learning

The prominence of automatic methods to identify objects in natural images is ever increasing. The performance of such systems recently reached the same performance as human observers [ref needed]. Moreover, these systems which were trained on energy greedy high performance computers are now designed to work on more common hardware such as desktop computers with a decent GPU. However, such methods are not yet available for mobile devices, as will be necessary for instance for the fast detection of visual objects in autonomous driving, as when has to identify a pedestrian from a sign pole. More importantly, the robustness of such methods is still lower than that of humans. Indeed, it is still difficult for such a system to learn to categorize a particular object-class given all its possible spatial configurations and the respective geometrical, visual transformations. This explosion

of combinations is currently handled by increasing accordingly the number of parameters, hence the energy consumption of such methods. As a consequence, state-of-the-art classification architectures contain many millions of parameters while still handling relatively small images.

On the contrary, the human visual system is able to perform such a feat very rapidly (less than 100 ms) [Kirchner, -] and at a low energy cost ($\sim 5\text{W}$). On top of that the system is mostly autonomous, robust to transforms or lighting conditions and can learn with a few examples. If many different features may explain such efficiency, a main difference of the human visual system with classical computer vision approaches is the fact that its sensor (the retina) combines

a non homogeneous sampling of the world with the capacity to rapidly change its center. Indeed, on the one hand the retina is composed of two separate systems: a central, high definition area and a large peripheral area. On the other hand, the retina is attached on the back of the eye and the eye is capable of low latency, high speed ($>500^{\circ}/s$) eye movements. In particular, saccades allow for efficient changes of the position of the center of gaze. The interplay of those 2 properties allow to engage observers in an action perception cycle which sequentially scans parts of the image. This behavior is prevalent during our lifetime (2/3 saccades per second = 2 billions per life). This behavior is one type of active inference [Friston] and we will envision herein how to incorporate it to classical computer vision schemes.

To explain and take advantage of this visual behaviour, it is of particular importance to understand its computational and biological (neurophysiological) principles. One main hypothesis regarding this active vision is that visual scenes most often consist of a single visual object of interest. Take for instance the case of a conversation with a friend in a noisy café. To ease the understanding of his voice and emotion you will track his face despite all the remaining visual clutter. Such a visual experience can be simplified in a manner reminiscent of psychophysical experiments. An observer is asked to classify digits (for instance as taken from the MNIST database) as they are shown on a computer display. However these digits can be placed at a random position on the display, and visual clutter is added as a background to the image (see Figure 1). This defines more precisely our problem: how do we

identify a small object in a large image while not knowing its position?

This joint problem of localization and identification found many solutions in computer vision. Notably, recent advances in deep-learning have provided with efficient solution such as Faster-RCNN [Ren 2015] or Yolo [Redmon 2015]. This last implementation is particularly interesting as it predicts in the image the probability of proposed bounding boxes around the visual object. While rapid, the amount of such boxes greatly increases with image size and necessitates a dedicated hardware.

When limiting our problem to few objects of interest in the image, this strategy amounts to a classical problem in neuroscience, that is, the transformation of a luminous image into a saliency map [Itti 2001]. Such a computation

is essential to understand and predict saccades but also as models of attention. Recently, deep learning methods have extended this model by learning the computation of saliency maps over large databases of natural images [Kummerer 2014].

While these methods are efficient at predicting the probability of fixation, they miss an essential point in the action perception cycle: They operate on the full image while the retina operates on the non-uniform, foveated sampling of visual space (Figure 1-B). As such we believe that this is an essential step to reproduce and understand this active vision process.

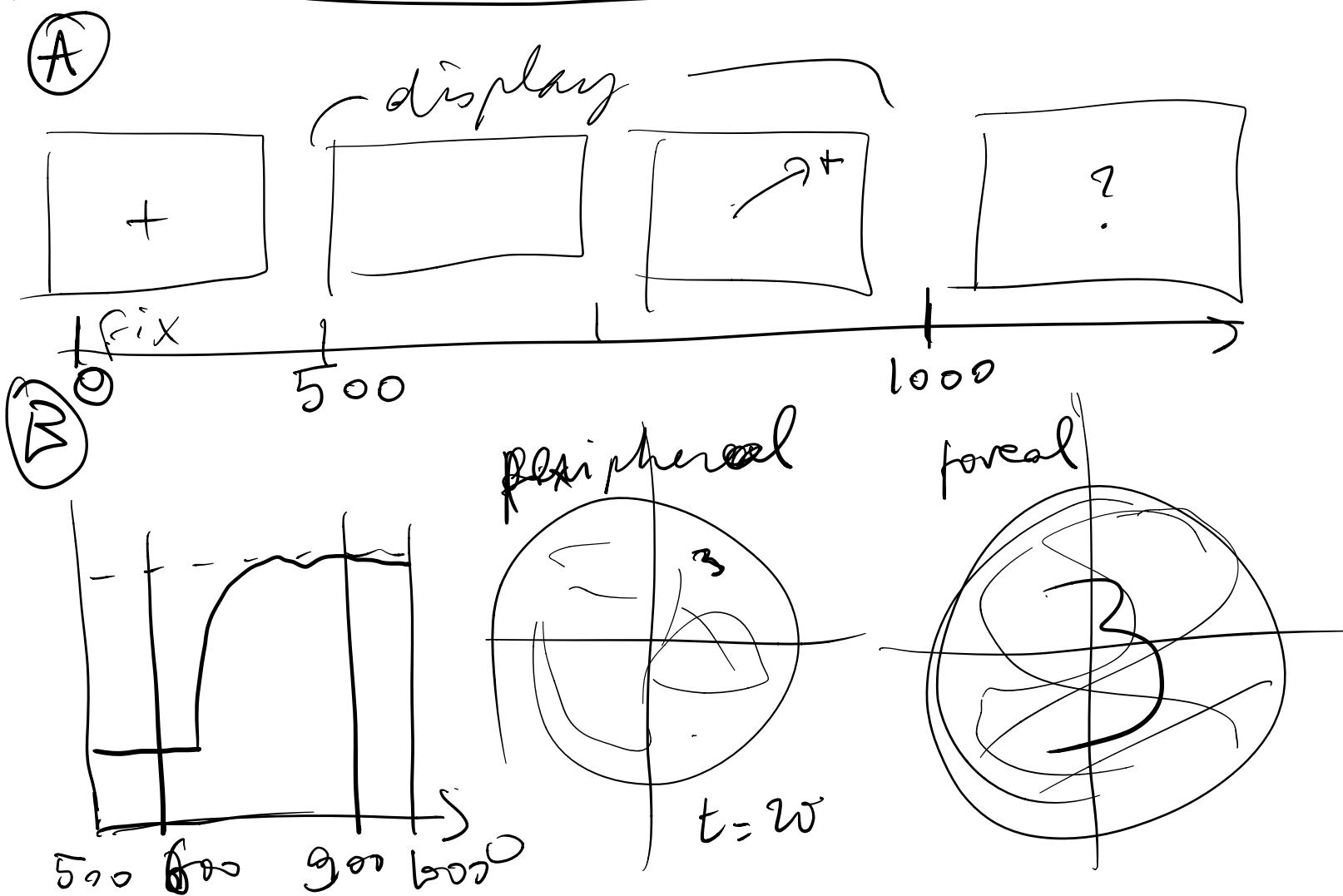
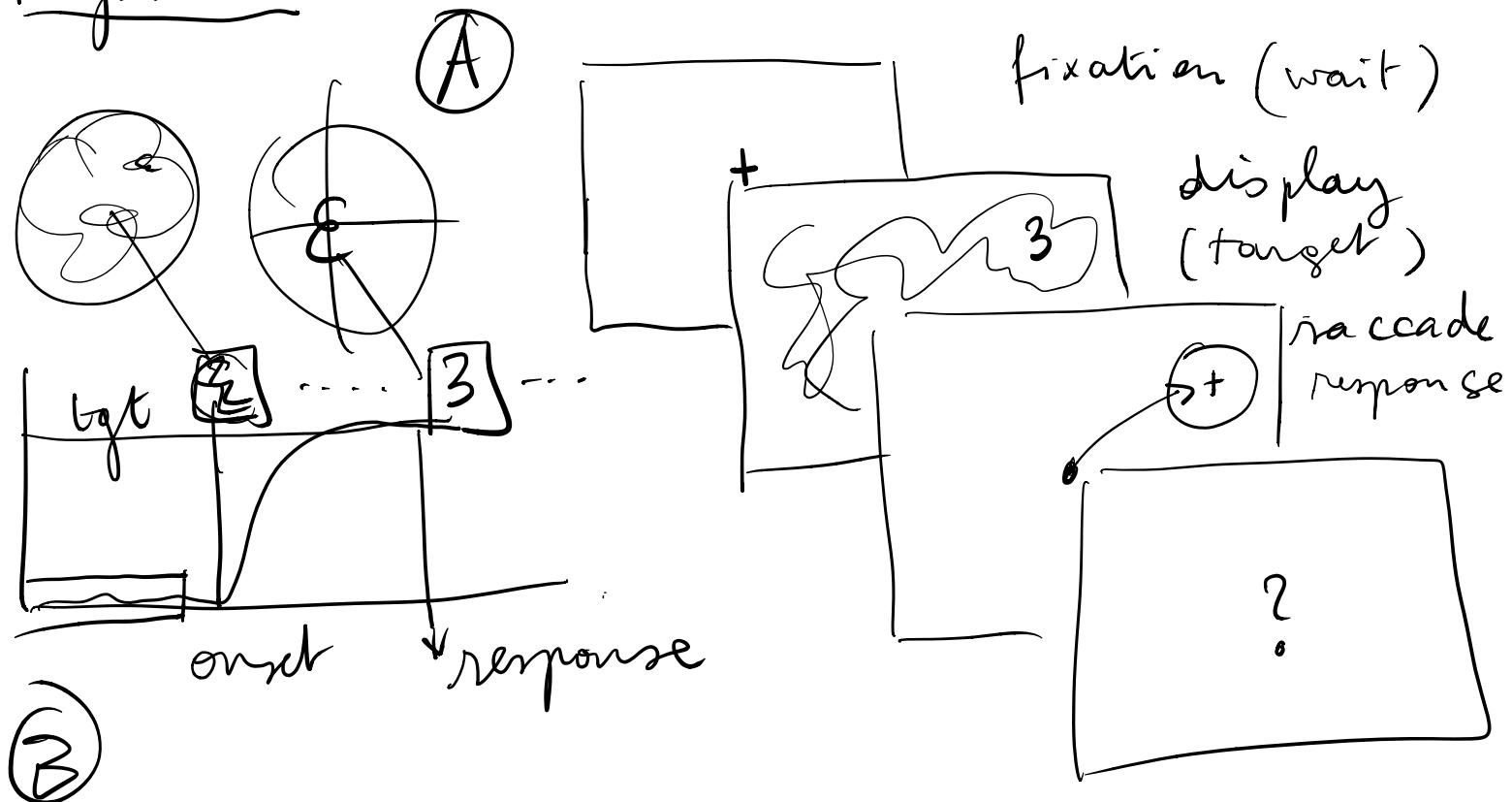
R.Botics

→ identify gap

↓
= "category before seeing"

An interesting perspective is given with previous modeling of foveated sensors. Indeed the non-uniform sampling of visual space is often modeled as a log-polar conformal mapping [Javier Traver 2010] which has a long history in computer vision and robotics. A first property of this mapping is the separation between the foveal and the peripheral areas as we defined above. This transformation has also other notable properties, such as the correspondence (by way of translations) in the radial and angular directions to rotations and scalings (respectively) in the visual domain. However, this sensor is not coupled to an action (but see [ref needed]). This paper aims at addressing the fragmentation of these studies respective to their fields (Machine learning, neuroscience, robotics) to propose a novel computational model of this active vision behaviour.

Figure 1:



(A) Problem setting: After a fixation period, an observer is presented with a luminous display which shows a target (here a digit) at a random position. The display is presented for a short period but enough to perform a saccade on the potential target. In particular, the configuration of the display is such that by adding clutter and reducing the size of the digit it may become necessary to perform a saccade to be able to see the digit.

Finally, the observer identifies the digit.

(B) We show a prototypical trace of a saccadic eye movement to the target position. In particular, we show the fixation window used to ensure fixation during that window (green shaded area). Overlaid is a simulation of the retinotopic map at the onset of the display and after a (successful) saccade. This demonstrates that the position of the target has to be inferred from a degraded (sampled) image and that a correct identification is mediated by the action to the location of the target * before seeing it *.

<http://www.cnbc.cmu.edu/cns/papers/Najemnik-Geisler-05-N.pdf>

<https://cyber.sci-hub.se/MTAuMTAzOC9zNDE1OTMtMDE4LTAyNTUtNQ==/samonds2018.pdf>

<http://www.cs.colorado.edu/~mozer/Teaching/syllabi/ProbabilisticModels/readings/NajemnikGeisler2005Supplemental.pdf>

"the eye movement is presented briefly such that no eye movement is possible."

16

State-of-the-art

You ARE HERE

Several studies are relevant to our endeavor. First, one can consider optimal strategies to resolve the problem of the visual search of a target [Nayemik 2005]. In a setting similar to that presented in Figure {fig:intro}, where the target is an oriented edge and the background pink noise, authors show first that a bayesian ideal observer provides with an optimal strategy and second that human observers are close to that optimal performance. Though they can predict a sequence of saccades in this perception/action loop, this model is limited by the simplicity of the display (elementary edges and a finite number of locations on a triangular grid) and by the abstract level of (Bayesian) modelling. Though its apparent simplicity, this important study could predict some characteristics of visual scanning such as the trade-off between memory content and rapidity.

Looking more closely at neurophysiology, the study of [Samonds 2016] allows to go further in our understanding of the interplay between saccadic behaviour and the statistics of the input. In this study, authors were able to manipulate the size of saccades by manipulating key properties of the presented (natural) images. In particular, smaller images generate smaller saccades.

Interestingly, they also explained the size of saccades for different species, including mice which lack a foreal region, by the size of receptive fields. One key prediction of this study for our problem is the fact that saccades seem optimal to a priori decorrelate the visual input, that is, to minimize redundancy, knowing statistics of natural inputs.

A further modelling perspective is given by [Friction12]. In this model, we first have a full description of the visual world as a generative process, here of the presentation of faces, knowing they are constituted of their components: mouth, eye, etc...
An

agent is completely described by giving 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 we described in Figure `ref{fig:intro}`. Equipping the agent with the ability to actively sample the visual world enable us to explore the idea that actions (saccadic eye movements) are 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, of the probable consequences of possible hypothesis as they would be realized into actions.

Such a model constitutes an Active Inference scheme (cite of Mirza 18) and simulations of the resulting optimization scheme reproduce sequential eye movements which fit well with empirical data. Compared to Najemnik (2005), saccades are not the output of a value-based cost function but a consequence of the agency for the agent to minimize the uncertainty about his beliefs, knowing his priors on the generative model of the visual world. Such an approach applies well to our settings as described in Fig {fig:intro}. In particular, we will similarly include a generative process of the visual world as a digits at a random position and embedded in a cluttered noise. Then, we will equip the agent with a pre-activated sensor and with the ability to actively scan the visual image. Knowing such priors, we will optimize the behaviour of this agent and explore its properties.

1c outline of the paper : the uncoupling hypothesis

The goal of this work is to emulate a model of active vision which is able to find a visual target from a known target, while the position is unknown. This comes with the constraint that the visual sensor is foveated. We will use this constraint as an asset to minimize the overall computational cost of finding the target. This generic visual search problem is of broad general interest in machine learning, computer vision and robotics, but also to neuroscience, as it speaks to the mechanisms underlying foveation and more generally to low-level attention mechanisms.

Our aim is to produce a generic principled model and to implement a fast computer implementation. As such, we will start as in \cite{Friston12} by a probabilistic formulation and use the fundamental hypothesis outlined in

Figure \ref{fig:intro}: The position of an object is independent from its identity. 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 [Melissa Le-Houérou and Wolfe, 2012] for more complex visual scene grammars). From this independence hypothesis, we may separate the two tasks (identification vs localization) in two separate pathways with different morphologies (respectively foreal and peripheral). Note that from the retinotopic projection of the visual information, this independence is conditional on action: These pathways should update their beliefs upon decisions made in each respective pathway.

This paper is organized as follows. After this introduction, we will define the notations, variables and equations for the generative process governing the experiment and the generative model for the active vision agent. In particular, we will derive our method to simplify the learning of an optimal agent given these definitions. In section \backslash ref{sec: results}, we will then show results of numerical simulations of this agent. We will first demonstrate some applications of this framework to different levels of complexity of the problem. This will allow us to derive some limits of this agent and, as in \backslash citet{Najemnik05}, we will draw some analogies with biologically observed eye movements. Finally, in section \backslash ref{sec: discussion}, we will summarize these results in comparison with other similar schemes. We will conclude by showing the relative advantages of using this active inference approach.

Theory (methods)

~~1 Introduction~~

~~1.1 Issue~~

~~1.2 State of the art~~

~~1.3 Outline~~

1.3.1 Notations

- \mathbf{x} : visual field (image)
- \mathbf{y} : target category (categorical)
- \mathbf{u} : target position (real coordinates or categorical, retinocentric referential)

Generative model :

$$\mathbf{x} \sim P(X|\mathbf{y}, \mathbf{u})$$

Full inference (posterior):

$$P(Y, U|\mathbf{x}) \propto P(\mathbf{x}|Y, U)$$

Independence assumptions :

$$P(Y, U) = P(Y)P(U) \quad (\text{toujours vrai}) \quad (1)$$

$$P(Y, U|X) = P(Y|X)P(U|X) \quad (\text{faux s'il y a plusieurs cibles}) \quad (2)$$

Partial inference on object category:

$$P(Y|\mathbf{x}, \mathbf{u}) \propto P(\mathbf{x}|Y, \mathbf{u})$$

Partial inference on object position:

$$P(U|\mathbf{x}, \mathbf{y}) \propto P(\mathbf{x}|U, \mathbf{y})$$

Marginals:

- $P(Y|\mathbf{x}) = \int P(Y|\mathbf{x}, \mathbf{u})d\mathbf{u}$
- $P(U|\mathbf{x}) = \int P(U|\mathbf{x}, \mathbf{y})d\mathbf{y}$

1.3.2 What we did so far...

Consider a view \mathbf{x} that contains a single target \mathbf{y} at unknown retinocentric position \mathbf{u} . The brain needs to guess both \mathbf{y} and \mathbf{u} with limited computational resources.

We assume here that the brain adopts independence assumption (2), making a separation between the “Where” and the “What” pathways, forming separate (and cheaper) inferences :

- $p(Y|\mathbf{x})$
- $p(U|\mathbf{x})$

Another assumption is that the category \mathbf{y} is *translationally invariant*: given a transformation \mathcal{T} ,

$$\mathcal{T}(\mathbf{u}, \mathbf{y}) = (\mathcal{T}(\mathbf{u}), \mathbf{y})$$

Now, given \mathbf{x} and the separation assumption, it is sensible to change the viewpoint to better estimate \mathbf{y} , because \mathbf{y} is invariant to the viewpoint transformation.

This is where *active inference* comes into the play:

- Consider that the true target is $\hat{\mathbf{y}}$
- Consider that the target current retinocentric position is \mathbf{u}
- Then, for any translation $\delta\mathbf{u}$, the future posterior on the true target is estimated by: $\mathbb{E}_{\mathbf{x}' \sim p(X|\hat{\mathbf{y}}, \mathbf{u} + \delta\mathbf{u})} p(\hat{\mathbf{y}}|\mathbf{x}')$
- And the optimal translation is: $\underset{\delta\mathbf{u}}{\operatorname{argmax}} \mathbb{E}_{\mathbf{x}' \sim p(X|\hat{\mathbf{y}}, \mathbf{u} + \delta\mathbf{u})} p(\hat{\mathbf{y}}|\mathbf{x}')$

If now \mathbf{u} is unknown and needs to be guessed from \mathbf{x} , the optimal translation is:

$$\underset{\delta\mathbf{u}}{\operatorname{argmax}} \mathbb{E}_{\mathbf{u} \sim p(U|\mathbf{x})} \mathbb{E}_{\mathbf{x}' \sim p(X|\hat{\mathbf{y}}, \mathbf{u} + \delta\mathbf{u})} p(\hat{\mathbf{y}}|\mathbf{x}')$$

with :

- $p(U|\mathbf{x})$ the inferred target position
- and $\mathbb{E}_{\mathbf{x}' \sim p(X|\hat{\mathbf{y}}, \mathbf{u} + \delta\mathbf{u})} p(\hat{\mathbf{y}}|\mathbf{x}')$ the expected inference on the actual target.

1.3.3 Accuracy maps

In practice, it is computationally impossible to make exact guesses about the future observation \mathbf{x}' . Our second assumption is that instead of predicting future inferences on true target, the brain trains a *parametric accuracy map* by experience (trial and error).

In a model-based approach, the *accuracy maps* can be calculated using a parametric classifier :

- Given a training set $\{(x_1, u_1, y_1), \dots, (x_n, u_n, y_n)\}$:
 - Train a classifier p_θ that estimates $p(Y|\mathbf{x})$.
- Then, for each class $\hat{\mathbf{y}}$, taking $\tilde{\mathbf{y}} \sim p_\theta(Y|\mathbf{x})$, the classification rate $r_\theta(\mathbf{u})$ is an estimator of the posterior expectation :

$$\begin{aligned} r_\theta(\mathbf{u}) &= \mathbb{E}_{\mathbf{x} \sim p(X|\hat{\mathbf{y}}, \mathbf{u})} \mathbb{E}_{\tilde{\mathbf{y}} \sim p_\theta(Y|\mathbf{x})} \delta_{\hat{\mathbf{y}}=\tilde{\mathbf{y}}} \\ &= \mathbb{E}_{\mathbf{x} \sim p(X|\hat{\mathbf{y}}, \mathbf{u})} p_\theta(\hat{\mathbf{y}}|\mathbf{x}) \\ &\simeq \mathbb{E}_{\mathbf{x} \sim p(X|\hat{\mathbf{y}}, \mathbf{u})} p(\hat{\mathbf{y}}|\mathbf{x}) \end{aligned}$$

that forms an *accuracy map* for each target position \mathbf{u} .

1.3.4 Parametric transformation (Colliculus?) map

One can now select $\delta\mathbf{u}$ with the parametric estimator:

$$\begin{aligned} \widehat{\delta\mathbf{u}} &\simeq \underset{\delta\mathbf{u}}{\operatorname{argmax}} \mathbb{E}_{\mathbf{u} \sim p(U|\mathbf{x})} r_\theta(\mathbf{u} + \delta\mathbf{u}) \\ &= \underset{\delta\mathbf{u}}{\operatorname{argmax}} Q(\delta\mathbf{u}|\mathbf{x}) \end{aligned}$$

with $Q(\delta\mathbf{u}|\mathbf{x})$ the *transformation map*, given the view \mathbf{x} and the marginal posterior estimate $p(U|\mathbf{x})$.

It must be noticed that, given $\hat{\mathbf{u}} = \underset{\mathbf{u}}{\operatorname{argmax}} r_\theta(\mathbf{u})$, the transformation map is maximal at $\delta\mathbf{u} = \hat{\mathbf{u}} - \mathbf{u}$. Each initial \mathbf{u} provides a different transformation map, that is a shift of the original accuracy map (*Ergodic assumption??*).

We assume in the following that a parametric action value map Q_ψ can be trained on top of the parametric classifier p_θ and its accuracy map r_θ . The training set is $\{(\mathbf{x}_1, \mathbf{u}_1), \dots, (\mathbf{x}_n, \mathbf{u}_n)\}$ and the accuracy map classifier learns to associate each \mathbf{x} with its full transformation map $Q(\cdot|\mathbf{x})$.

1.3.5 Algorithms

Once p_θ and Q_ψ are trained, the recognition algorithm is straightforward:

Single saccade algorithm:

1. Read the view \mathbf{x}
2. Choose $\delta\mathbf{u}$ according to $Q_\psi(\cdot|\mathbf{x})$
3. Move the eye
4. Update the view \mathbf{x}'
5. Identify the target with $\tilde{\mathbf{y}} \sim p_\theta(Y|\mathbf{x}')$

Multi saccades algorithm:

1. $q(Y) \leftarrow$ uniform distribution
2. Read the view \mathbf{x}
3. Choose $\delta\mathbf{u}$ according to $Q_\psi(\cdot|\mathbf{x})$
4. Repeat several times up to some posterior confidence threshold:
 - (a) Move the eye
 - (b) Read \mathbf{x}
 - (c) $q(Y) \leftarrow q(Y) \times p_\theta(Y|\mathbf{x})$
 - (d) normalize q
 - (e) Choose $\delta\mathbf{u}$ according to $Q_\psi(\cdot|\mathbf{x})$ (with some inhibition of return mechanism)
5. Identify the target with $\tilde{\mathbf{y}} \sim q(Y)$

2 Methods

2.1 Visual transformation

2.1.1 Wavelets

2.1.2 Log Gabor

2.2 Accuracy map

2.3 Network architecture

3 Results

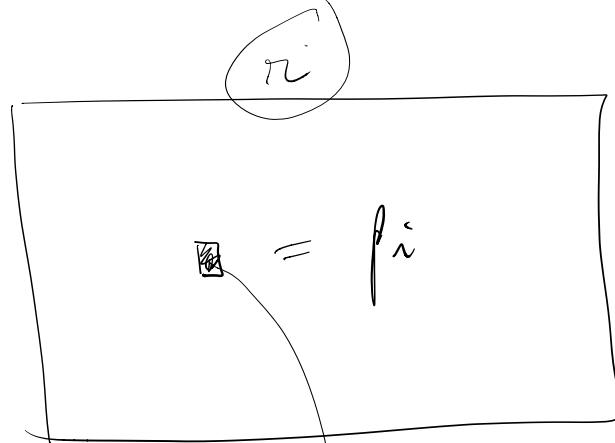
[https://github.com/laurentperrinet/WhereIsMyMNIST/blob/master/2018-11-13-Where%20recap%20\(clut](https://github.com/laurentperrinet/WhereIsMyMNIST/blob/master/2018-11-13-Where%20recap%20(clut)

+ Parker de

Mirza 2018

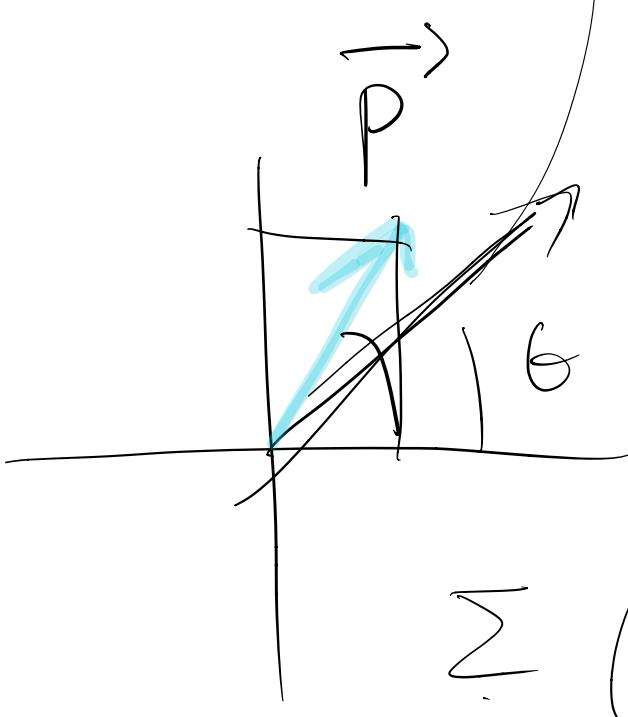
→ MDP

STATS IN
CIRCULAR
MAPS



cf. Kaplan 13
(Eq. 10)

<https://laurentperrinet.github.io/publication/kaplan-13/>



1. make a 2D vector from the polar coord
2. make the mean

$$\sum_i \left(p_i \cos \theta_i, p_i \sin \theta_i \right)$$

$$\bar{\theta} = \text{ArcTan} 2 (p_y, p_x)$$

Similar for radius

$$\bar{r} = \frac{\sum_i p_i r_i}{\sum_i p_i}$$

$$\log \bar{r} = \sum_i p_i \log(r_i)$$

4 Discussion

4.1 Summary

4.2 Limits

4.3 Perspectives

Predictions of the model → discussion

- Reaction times versus proba of detection
- can we get fooled by "fake" digits?
- bigger digits \Rightarrow bigger saccades.
↓
the type of database \Rightarrow context influences size of saccades of Yarbus

Interaction What and Where



$$P(u, y | \text{oc})$$

?

"what" avant la saccade

$$P(y | \text{oc}, u=0)$$

en sortie on a un vecteur

$$p_{y=y} = [\dots] \text{ de probas}$$

en fonction d'un critère (entropie, seuil)

on peut choisir

$$q_0 = \max_y (P_{Y=y})$$

- on accepte STAY \rightarrow évaluation

- sinon, signal Go

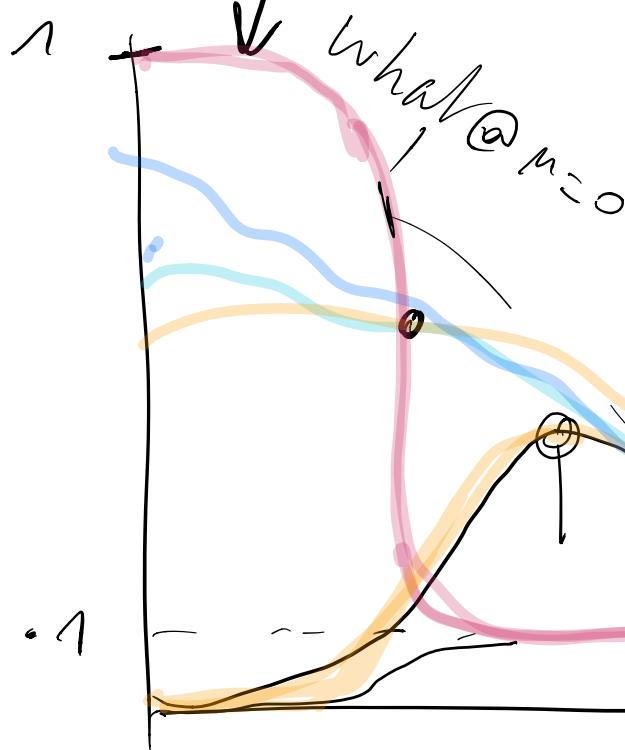
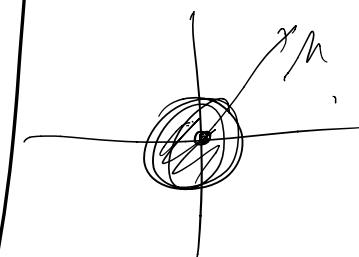
prediction locale du what sur le where :

sachant G_0 ,

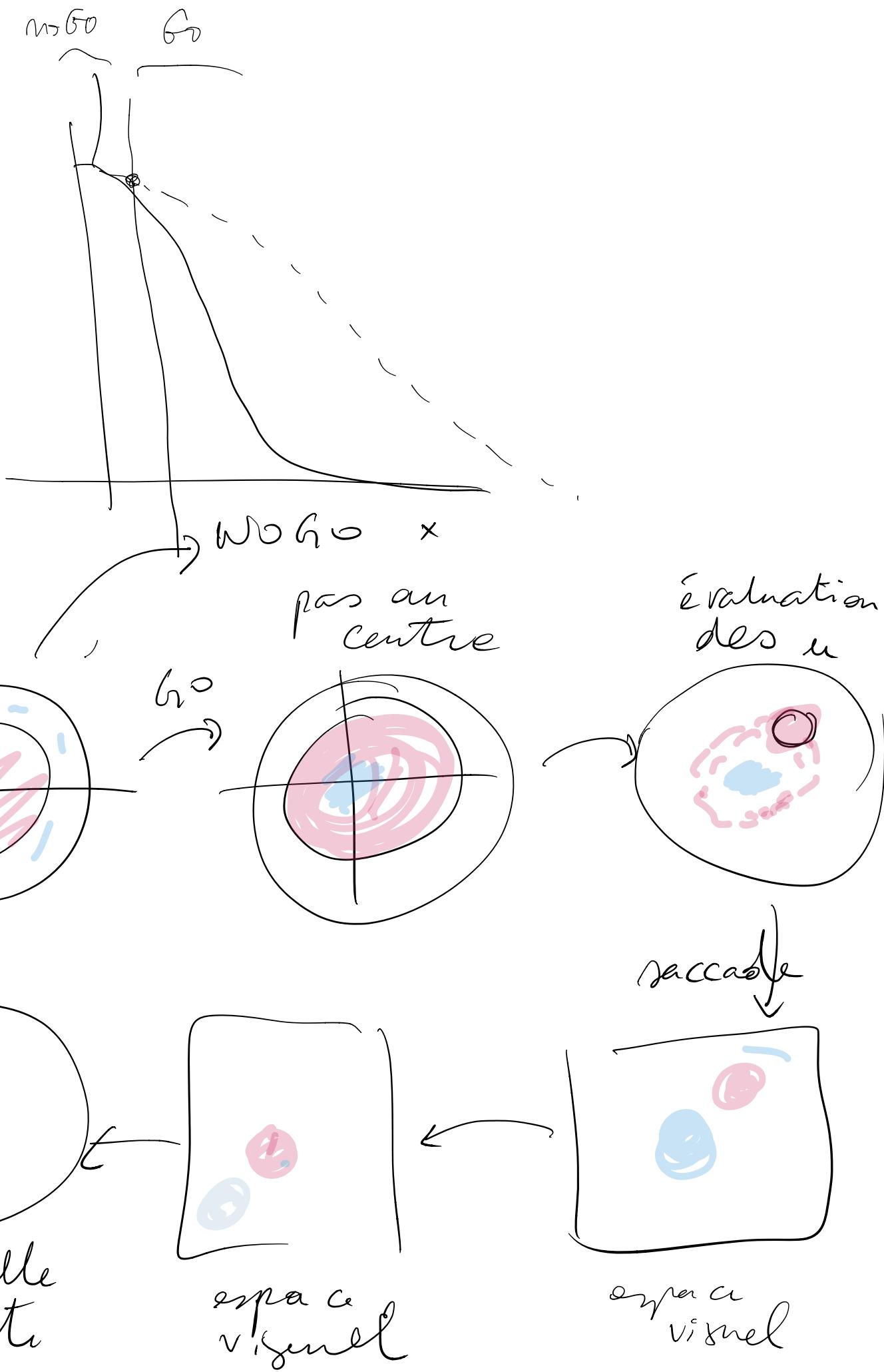
$$P(a | x, u) = \underset{y}{\operatorname{Max}} [P(y | x, u)]$$

$$P(a | x, u=0)$$

$$P(a_u | a_{u=0})$$



avant de savoir si on y va



General case: Visual information gain maximization

Consider a view \mathbf{x} generated from a target \mathbf{y} viewed at retinocentric position \mathbf{u} .

Consider first that :

- The generative model $p(X|\mathbf{y}, \mathbf{u})$ is known
- The retinocentric position \mathbf{u} is known.
- The view \mathbf{x} is known.
- The target category \mathbf{y} is unknown.

The question comes how to choose the new retinocentric position \mathbf{u}' in order to maximize the *mutual information* between $\mathbf{x}|\mathbf{u}$ (current view) and $\mathbf{x}'|\mathbf{u}'$ (future view).

In general, the visual Information Gain between two visual fields $\mathbf{x}|\mathbf{u}$ and $\mathbf{x}'|\mathbf{u}'$ is:

$$\text{IG}(\mathbf{x}|\mathbf{u}; \mathbf{x}'|\mathbf{u}') = -\log p(\mathbf{x}|\mathbf{u}) + \log p(\mathbf{x}|\mathbf{u}, \mathbf{x}', \mathbf{u}')$$

Information Gain Lower Bound Consider now that given \mathbf{x} and \mathbf{u} , the target category \mathbf{y} can be *inferred* using Bayes rule, i.e.:

$$P(Y|\mathbf{x}, \mathbf{u}) \propto P(\mathbf{x}|Y, \mathbf{u})$$

Then, it can be shown (see [?]) that :

$$\text{IG}(\mathbf{x}|\mathbf{u}; \mathbf{x}'|\mathbf{u}') \geq \mathbb{E}_{\mathbf{y} \sim p(Y|\mathbf{x}, \mathbf{u})} [\log p(\mathbf{y}|\mathbf{x}', \mathbf{u}') - \log(\pi(\mathbf{y}))]$$

with $\pi(\mathbf{y})$ the prior over the \mathbf{y} 's . When the prior is uniform, the information gain lower bound (IGLB) simplifies to $\mathbb{E}_{\mathbf{y} \sim p(Y|\mathbf{x}, \mathbf{u})} [\log p(\mathbf{y}|\mathbf{x}', \mathbf{u}')] + c$, with c a constant.

Predictive approach One can adopt a *predictive* approach to choose the new eye orientation \mathbf{e}' :

- First choose a new retinocentric position \mathbf{u}' that will maximize the information gain.
- Then choose \mathbf{e}' such that

$$\mathbf{z} - \mathbf{e}' = \mathbf{u}'$$

i.e.

$$\mathbf{e}' = \mathbf{e} + \mathbf{u} - \mathbf{u}'$$

The predictive approach needs three predictive steps:

- $p(Y|\mathbf{x}, \mathbf{u})$ is the current posterior over the target category inferred from the current observation,
- $\mathbf{x}' \sim p(X|\mathbf{y}, \mathbf{u}')$ is the predicted view generated by the model assuming that the target \mathbf{y} is seen from from \mathbf{u}' ,
- and $p(\mathbf{y}|\mathbf{x}', \mathbf{u}')$ is the predicted posterior for assumption \mathbf{y} , given \mathbf{x}' and \mathbf{u}' .

Then the optimal new retinocentric position is:

$$\hat{\mathbf{u}}' = \operatorname{argmax}_{\mathbf{u}'} \mathbb{E}_{\mathbf{y} \sim p(Y|\mathbf{x}, \mathbf{u})} [\mathbb{E}_{\mathbf{x}' \sim p(X|\mathbf{y}, \mathbf{u}')} [\log p(\mathbf{y}|\mathbf{x}', \mathbf{u}')]]$$

Taking $\delta\mathbf{e} = \mathbf{u} - \mathbf{u}'$, the optimal eye displacement is:

$$\widehat{\delta\mathbf{e}} = \operatorname{argmax}_{\delta\mathbf{e}} \mathbb{E}_{\mathbf{y} \sim p(Y|\mathbf{x}, \mathbf{u})} [\mathbb{E}_{\mathbf{x}' \sim p(X|\mathbf{y}, \mathbf{u} - \delta\mathbf{e})} [\log p(\mathbf{y}|\mathbf{x}', \mathbf{u} - \delta\mathbf{e})]]$$

(TODO : Attention il faudrait à partir de maintenant une carte qui moyenne les log posteriors car l'espérance du log n'est pas égale au log de l'espérance, i.e. $r_\theta^{\log}(\mathbf{u}|q) = \mathbb{E}_{\mathbf{y} \sim q(Y)} [\mathbb{E}_{\mathbf{x} \sim p(X|\mathbf{y}, \mathbf{u})} \log p_\theta(\mathbf{y}|\mathbf{x}, \mathbf{u})]$).

Plan Theory

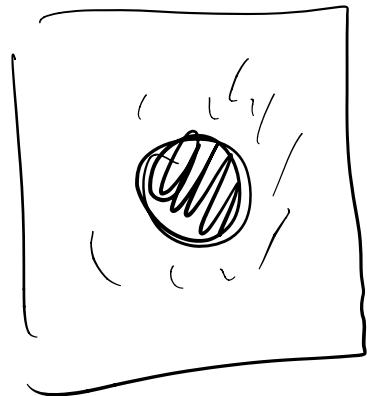
{ notations: proba,
convolution, ... }

a) generative process -

- $x = \text{digit}_g * \mathcal{D}_n + \text{noise}$
 - ↓
 - in $[0, 1]$
 - ↓
 - motion cloud.
- retinal transform
 $\mathcal{R}_n(x) \rightarrow \begin{array}{c} \boxed{\text{grid}} \\ \pi \end{array} \theta$

b) classifier = LeNet

- 98.1% accuracy $a = \% \text{ correct}$
- accuracy map



c) where detection

given $\mathcal{I}_0(x)$, can we estimate the accuracy in the visual space

$$IA(u) = P(a \mid \mathcal{I}_0(x)) = \text{function of coordinates}$$

- This is supervised by trial and error (corrective saccades)
- we can have an ergodic hypothesis and use the accuracy map
- it's a fitting problem = deep learning
- generative model = N-layered model
- "cost function" = $KL(p \parallel q)$ = free energy
(cf ITI 10g)
- Action $u^* = \text{ArgMin } F(u)$

④ inhibition of return