

AIX-MARSEILLE UNIVERSITÉ

MÉMOIRE DE MASTER

---

# Computational modelling of visual object localization in the magnocellular pathway

---

*Auteur:*

Pierre ALBIGÈS

*Superviseur:*

Emmanuel DAUCÉ

*Une thèse présentée à*

ECOLE DOCTORALE DES SCIENCES DE LA VIE ET DE LA SANTÉ

*en vue de l'obtention du diplôme de*

MASTER DE NEUROSCIENCES, SPÉCIALITÉ INTÉGRATIVES ET COGNITIVES

*et réalisée au sein de*

Institut de Neurosciences des Systèmes

*Durant la période : 04/12/2017 - 02/03/2018*

*“Wolves have no Kings”*

Robin Hobb

AIX-MARSEILLE UNIVERSITÉ

# *Abstract*

Faculté des Sciences, département de Biologie

Master de Neurosciences

Master de Neurosciences, spécialité Intégratives et Cognitives

**Computational modelling of visual object localization in the magnocellular pathway**

by Pierre ALBIGÈS

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...



# *Acknowledgements*

The acknowledgments and the people to thank go here, don't forget to include your project advisor...



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Vision naturelle . . . . .	1
1.2	Vision artificielle . . . . .	3
<b>2</b>	<b>Matériel et méthodes</b>	<b>5</b>
2.1	Support physique et numérique . . . . .	5
2.2	Modèle POMDP . . . . .	5
2.3	Champs rétinien . . . . .	6
2.4	Apprentissage supervisé . . . . .	6
	<b>Bibliography</b>	<b>9</b>
<b>A</b>	<b>Figures</b>	<b>11</b>
<b>B</b>	<b>Code source</b>	<b>15</b>





# 1 Introduction

## 1.1 Vision naturelle

Tous les êtres vivants utilisent la vision à un degré ou à un autre, et pour de nombreuses espèces -y compris la notre- elle est même la modalité perceptive principale. Elle est alors primordiale pour appréhender l'environnement et interagir avec celui-ci, que ce soit dans une optique de survie de l'individu ou dans la construction de relations sociales.

Chez les vertébrés la vision débute à la surface de la rétine, où les cellules photovoltaïques (cônes et bâtonnets) réalisent la transduction des signaux lumineux qui les atteignent en signaux électriques, transmissibles au réseaux nerveux en aval.

Les cônes et les bâtonnets sont différenciables par un certain nombre de caractéristiques, notamment leur sensibilité aux longueurs d'ondes lumineuses et leur distribution au sein de la rétine. Ces différences permettent à notre rétine de rester fonctionnelle dans de nombreuses situations, y compris lorsque la luminance est très faible (le seuil absolu de la rétine humaine correspondant à 70 photons) et lui permet donc de nous fournir des informations pertinentes dans une grande variété de contextes. Le champ visuel peut être divisé en deux parties. La vision centrale (environ  $2^\circ$  chez l'humain) est soutenue par la fovea, une région rétinienne comprenant uniquement des cônes. On y observe l'acuité visuelle la plus importante (la région présente les champs récepteurs les plus petits de la rétine).

La vision périphérique comprends majoritairement des bâtonnets, présente une faible acuité et une perception des couleurs très faible (voir nulle). Elle est par contre très sensible à des variations de luminance et de fréquence spatiale (donc aux mouvements).<sup>2</sup>

L'acuité visuelle diminue en fait avec l'excentricité par rapport à la fovéa. Autrement dit, les champs récepteurs visuels grandissent avec cette excentricité.

Lors de l'exploration de son environnement visuel, un agent va pouvoir détecter des stimuli dans sa vision périphérique mais ne va pas y présenter assez d'acuité pour réaliser une description précise.

En conséquences, l'agent va réaliser des saccades oculaires (mouvements brefs (20-60ms) des globes oculaires) afin de placer l'image de la cible (ou tout du moins sa position prédite dans l'espace) au niveau de la fovéa, permettant ainsi de traiter les informations en provenant avec la plus grande précision possible.

L'activité rétinienne est transmise le long des voies nerveuses visuelles jusqu'au cortex visuel, où sera réalisé la majorité du traitement des informations -notamment haut niveau- qu'elle code. Entre la rétine et le cortex existe un certain nombre d'étapes, mais tout au long de ces voies la distribution rétinienne de l'information (la rétinotopie) est conservée.

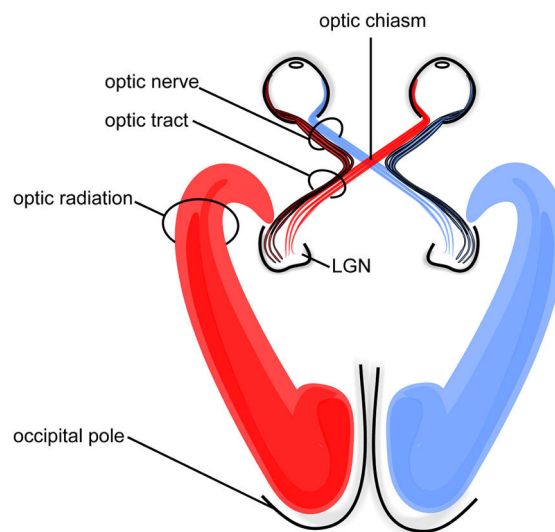


FIGURE 1.1: Schéma des voies visuelles précorticales humaines (adapté de Hofer S. et al., 2010 via Wikimedia Commons [CC BY 3.0])

Dans leurs travaux de 1962 et 1977, Hubel et Wiesel émettent l'hypothèse des courants visuels. Cette hypothèse définit trois voies nerveuses naissant dans le **corps genouillé latéral** (LGN, où sont présents les corps des types cellulaires donnant leurs noms aux voies) et projetant sur le **cortex visuel primaire** (V1) : **magnocellulaire** (M), **parvocellulaire** (P) et **koniocellulaire** (K). Chacune de ces voies supporte le transport d'informations codant pour des caractéristiques différentes des stimuli visuels. L'activité des cellules M ne distingue pas les couleurs mais est sensible à des différences fines de luminosité, de contraste et de fréquence spatiale. Ces caractéristiques semblent lier la voie M notamment au traitement de la luminosité et des mouvements.<sup>2</sup>

Cette multiplicité de voies visuelles réalisant des traitements différents des stimuli est conservée

au delà de V1, où l'on décrit deux courant sortants : la **voie ventrale** (transportant et traitant majoritairement pour des informations provenant de la voie P) et la **voie dorsale** (transportant et traitant majoritairement pour des informations provenant de la voie M).

La voie ventrale communique ainsi principalement avec les aires cérébrales du lobe temporal, l'activité de son réseau étant primordiale pour la reconnaissance et l'identification des objets visuels. La voie dorsale quant à elle communique principalement avec les aires du lobe pariétal, l'activité de ce réseau étant primordiale pour le traitement des relations spatiales entre les objets visuels ainsi que pour le guidage attentionnel et physique vers eux.<sup>2</sup>

Parmi ce réseau dorsal, on trouve l'**aire intrapariétale latérale** (LIP), qui reçoit en partie des informations directement depuis V1 et V2 (contournant donc le traitement d'aires en amont, dont l'aire MT) codant pour des stimuli dans le champs visuel périphérique. Des travaux ont d'ailleurs relié l'activité des neurones du LIP à la planification des saccades oculaires et à la représentation spatiale des objets visuels.

## 1.2 Vision artificielle

La vision représentant notre modalité perceptive principale et les aires dédiées à traiter ses informations occupant une part significative de notre système nerveux central (environ 50% chez certains primates),<sup>3</sup> l'étudier permet de mieux comprendre le fonctionnement général de notre système nerveux. De nombreux domaines d'étude s'intéressent donc au fonctionnement du système visuel. Parmi eux, les neuromathématiques se basent sur les données expérimentales (anatomique, physiologique et comportementale) pour émettre des modèles mathématiques sur le fonctionnement d'une partie ou de l'ensemble de la modalité visuelle. Idéalement, ces théories doivent pouvoir expliquer son activité dans l'ensemble des contextes observables, mais aussi pouvoir prédire son comportement dans de nouveaux contextes.<sup>3</sup>

Mais ces modèles ne sont pas une finalité en soi dans la compréhension du système. Ils peuvent permettre d'identifier dans les théories de son fonctionnement des défauts ou des zones obscures à notre compréhension et donc diriger les études expérimentales vers ces points.<sup>3</sup>

L'identification de ces points et la démonstration de ces théories peut passer par le domaine des neurosciences computationnelles, qui va tenter d'appliquer ces modèles mathématiques dans des modélisations du système nerveux. Au delà de la possible validation ou invalidation des modèles, les neurosciences computationnelles permettent de résoudre des problèmes d'ingénierie (puissance de calcul disponible, vitesse de traitement, adaptabilité à l'environnement, ...) en s'inspirant des systèmes biologiques, très optimisés, et donc de créer des systèmes artificiels neuromimétiques plus performants et intégrables dans des systèmes embarqués ou des interfaces cerveau-machine.

Dans cette étude, nous avons tenté de construire un modèle simple de localisation de cible visuelle dans un champs rétinien : le modèle possède une vision centrale où son acuité est maximale et une vision périphérique dont l'acuité diminue avec l'excentricité.

Le modèle doit être capable de détecter dans sa vision périphérique une cible visuelle aux caractéristiques simples (représentée par un stimulus provenant le base de données MNIST), de prédire précisément sa position et de réaliser une saccade oculaire afin de la placer dans sa fovea, ce qui lui permet alors de l'identifier avec une certitude élevée.

Pour cela plusieurs méthodes (comprenant elles-mêmes plusieurs sous-méthodes) s'offrent à nous. Les **modèles de saillance** où l'on tente de décrire une image en fonction des régions (ou pixels) qui présentent la plus grande probabilité de fournir des informations pertinentes et donc qui devraient attirer les saccades. Ces modèles permettent donc de créer des histogrammes de probabilité de fixations oculaires mais présentent un certaines limites, dont le fait qu'ils ne vont pas cibler d'objets partiellement ou entièrement cachés.<sup>1</sup>

Les **modèles de contrôle**, tentent quant à eux de prédire quelles règles le système devra suivre pour répondre au mieux à une tâche. A chaque saccade, de nouvelles informations sur l'environnement sont collectées et permettent de changer l'opinion de l'agent sur le monde.<sup>1</sup>

## 2 Matériel et méthodes

### 2.1 Support physique et numérique

L'ensemble des modélisations ont été réalisés sur une ordinateur portable hébergeant une machine virtuelle. Leurs caractéristiques sont rassemblées dans le tableau suivant :

	Identifiant	Système d'exploitation	Processeur	Mémoire vive	Carte graphique
Machine physique	ASUS ROG G75VW	Windows 7 64-bit SP1	Intel Core i7-3610QM 2,30GHz (8CPU)	8 GB (DDR3)	NVIDIA GeForce GTX670M
Machine virtuelle (ressources allouées)	VirtualBox v.5.2.6	Ubuntu 16.04	4 CPU, 90% des ressources	5298 Mo	Support GPU non-utilisé

Les modélisation ont été réalisées à l'aide du langage de programmation **Python** (version 3.6.4) renforcé de la librairie **TensorFlow** (version 1.4) et de l'interface graphique **Jupyter**.

La base de données **MNIST** a été utilisée pour l'apprentissage et l'évaluation du modèle. Elle contient 70.000 images de chiffres manuscrits (60.000 pour l'entraînement, 10.000 pour l'évaluation), centrés et dont la taille a été normalisée. Chaque image est accompagnée d'un label décrivant quel chiffre elle contient.

### 2.2 Modèle POMDP

Le problème de recherche d'information dans un contexte d'exploration de l'environnement visuel peut être formulé comme un **processus de décision Markovien partiellement observable** (POMDP).<sup>1</sup> Dans un POMDP (figure A.1), l'agent perçoit partiellement l'**état de l'environnement**  $S$  à un temps  $t$  (dans ce travail l'environnement visuel, perçut au travers d'un champs rétinien) et peut réaliser des **actions**  $A$  (ici des saccades oculaires) qui peuvent avoir des conséquences l'environnement et sa perception  $O$ . L'agent va ainsi construire un **état de croyance**  $B$  (ici la catégorie prédite du stimulus) en fonction des observations et des actions réalisées jusqu'ici.<sup>1</sup>

Un tel système doit satisfaire la **propriété de Markov**, qui décrit que la distribution de probabilité des futurs états ne dépends que de l'état précédent et pas de toute la séquence d'états les précédents.

Ainsi lors de l'évolution du système dans le temps, on considère que l'état suivant de l'environnement est uniquement influencé par son état actuel et l'action (éventuelle) réalisée par l'agent (équation 2.1).<sup>1</sup>

$$p(s_{t+1}|s_{1:t}, a_{1:t}, o_{1:t}) = p(s_{t+1}|s_t, a_t) \quad (2.1)$$

De même, les observations actuelles de l'agent ne dépendent que de l'état actuel de l'environnement et de l'action (éventuelle) qu'il réalise (équation 2.2).<sup>1</sup>

$$p(o_t|s_{1:t}, a_{1:t}) = p(o_t|s_t, a_t) L'obser \quad (2.2)$$

$$B_t^i = p(S_t = i|A_{1:t}, O_{1:t}) \quad (2.3)$$

## 2.3 Champs rétinien

Avant d'être utilisée par le modèle, l'image provenant de la base MNIST subit un certain nombre de transformations.

Chaque image présente originellement des niveaux de gris sur 28x28 pixels. Cette image est placée au hasard (mais de façon normocentrée) sur une image de 128x128 pixels (figure A.2) dont le fond est blanc (ne contient pas d'informations).

A cette image est ensuite appliqué un *filtre Wavelets* ou un *filtre LogGabor*, deux méthodes permettant d'obtenir un champs rétinien, c'est à dire un filtre visuel dont l'acuité est maximale en son centre (simulant la fovea) et diminuant avec l'excentricité (simulant la vision périphérique).

Les effets du *filtre Wavelets* sont visibles sur la figure A.3.

## 2.4 Apprentissage supervisé

Afin d'obtenir un modèle à la fois performant et adaptable, nous l'avons soumis à un apprentissage supervisé sous la forme d'une **régression linéaire multivariée** optimisée par **descente de gradient**.

Pour cela, nous calculons une hypothèse  $h_\theta$  (équation 2.4) sur la répartition des stimuli en multipliant chacune des valeurs  $x$  de l'entrée à un poids  $\theta$  qui lui est spécifique. Ces poids sont ensuite optimisés par descente de gradient, où ils sont comparés aux labels  $y$  des entrées (équation 2.5) pour un nombre d'exemples et d'itérations fixées. L'optimisation du paramètre d'apprentissage  $\alpha$  est nécessaire

au bon déroulement de l'apprentissage et sa valeur doit être adaptée pour éviter un sous- ou un sur-apprentissage (révélant respectivement une valeur trop faible ou trop importante).

$$h_{\theta}(x) = \theta^T x + b \quad (2.4)$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x_j^i \quad (2.5)$$

L'évolution du coût  $J(\theta)$  (équation 2.6) est un indicateur de l'efficacité de l'apprentissage, sa valeur devant décroître au fil de l'entraînement.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 \quad (2.6)$$

Nous construisons un modèle en deux couches. La première, *détecteur*, est entraînée à prédire la position du stimulus dans l'image, perçu au travers du champs rétinien et permet de réaliser une saccade jusqu'à ces coordonnées prédites, permettant d'approcher le stimulus de la vision fovéale. La seconde, *classifier*, est entraînée à prédire la catégorie du stimulus dans l'image, perçu au travers du champs rétinien et permet d'arrêter l'exploration de l'image lorsque sa prédiction présente une certitude assez élevée.





# Bibliography

- [1] Nicholas J Butko and Javier R Movellan. “Infomax control of eye movements”. In: *Autonomous Mental Development, IEEE Transactions on* 2.2 (2010), pp. 91–107.
- [2] John S. Werner and Leo M. Chalupa, eds. *The new visual neurosciences*. MIT Press. 2014, p. 1675. ISBN: 9780262019163.
- [3] Li Zhaoping. *Understanding vision : theory, models and data*. Oxford Uni. 2014, p. 383. ISBN: 9780199564668.



# A Figures

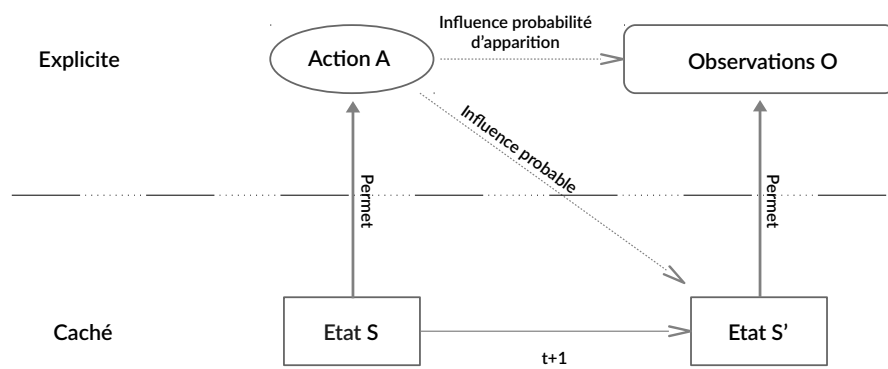


FIGURE A.1: Schéma des interactions entre l'agent et son environnement au cours du temps dans un modèle POMDP

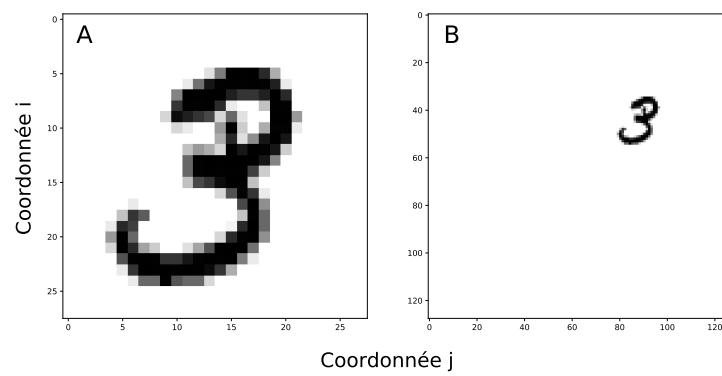


FIGURE A.2: **A.** Image originale tirée de MNIST ; **B.** Image après transformation géométrique et placé aux coordonnées ( $i = -20, j = 25$ )

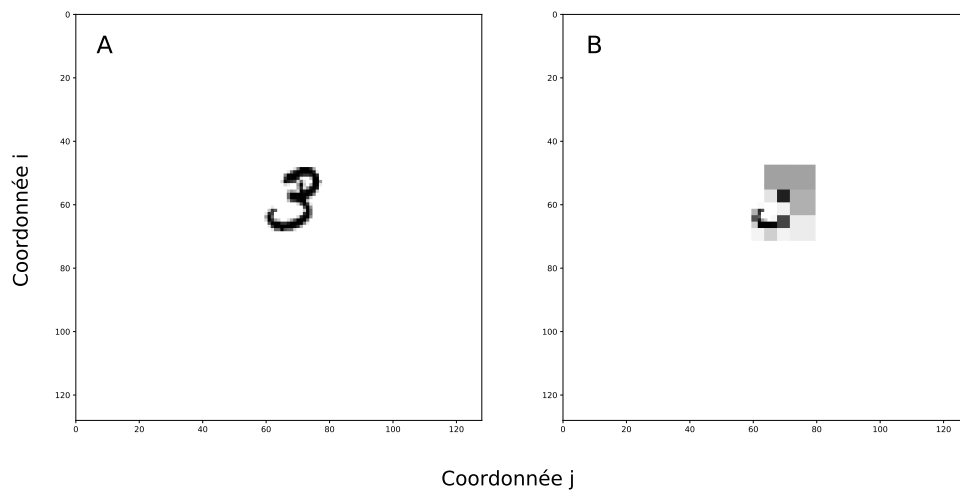


FIGURE A.3: **A.** Image avant application d'un filtre ; **B.** Image après transformation par vaguelettes (*filtre Wavelets*) et placé aux coordonnées ( $i = -6, j = 6$ )



## B Code source

L'ensemble du code source du modèle, ainsi que de ce rapport et d'autres documents complémentaires (dont les notes personnelles) sont entièrement disponibles **en ligne**.