

# Decision-making as an encoding-decoding process and its correlation with neuronal activity and behavior

Ramon Nogueira Mañas

---

TESI DOCTORAL UPF / 2017

Thesis Supervisor

Prof. Rubén Moreno Bote,  
Department of Information and Communication Technologies





A Gelasio y Daría,



# Agradecimientos

Me gustaría empezar esta tesis doctoral con la sección dedicada a los agradecimientos, ya que, si bien es cierto que durante los últimos cinco años he dedicado mucho tiempo y esfuerzo a la realización de este proyecto, nada de todo esto hubiera sido posible sin la inestimable ayuda de un gran número de personas. Pido disculpas por adelantado a todos aquellos de los que me pueda estar olvidando. Son todos los que están, pero probablemente no están todos los que son.

En primer lugar me gustaría agradecer profundamente a mis padres Daría y Gelasio por todo el apoyo que me han dado durante todos estos años; a Teresa por estar siempre a mi lado tanto en los buenos como en los malos momentos; a Isabel, Albert, Teresa y Pep por su apoyo y por los buenos ratos compartidos y a Àlex, Héctor, Carles, Joan y Ninja por todos los años de convivencia y por haber podido compartir con ellos una de las etapas más intensas, estimulantes y divertidas de mi vida.

También me gustaría darle las gracias a Rubén Moreno Bote por todo lo que me ha enseñado, por todas las conversaciones que hemos tenido durante estos cinco años y por haber confiado en mí desde el principio; a Julio Martínez Trujillo por haberme tratado como a un miembro más de su laboratorio y por haberme dado la confianza necesaria para mejorar como científico y a Iñigo Romero Arandia por todas las charlas científicas y no-científicas que hemos tenido y por haberme ayudado siempre que he tenido alguna dificultad.

Asimismo, me gustaría agradecer a todos mis amigos y amigas, tanto del ámbito académico como del personal, todo el apoyo y ayuda que me han dado durante estos años: Alex Hyafil, Josefina Cruzat, Vicente Pallarés, Ruggero Bettinardi, Víctor Saenger, Andrea Insabato, Iruñe Fernández, Gisela Pi, Alice Del Genovese, Joao Barbosa, Gabriela Mochol, Philipp Schustek, Dávid Samu, Belén Sancristóbal, Silvana Silva, Txema Esnaola, Marina Vegué, Ana Martín, Roberto Gulli, Matthew Leavitt, Lyndon Duong, Ben Corrigan, Rogelio Luna, Miguel Burgale-

ta, Adrián Ponce, Marina Cunquero, Laura Puigcerver, Claudia Caprile, Víctor Segura, Roberto Bailón, Albert Sidro, Brendan Golden, Monka, Joel López, Miquel Muntaner, Brasi, Fausto Altamirano, Òscar Rota, Tote Ruiz, David Pineda, Virginia Bru, Marçal Gabaldà y César Guillén. También me gustaría agradecer a Jan Drugowitsch, Gustavo Deco, Luca Bonatti, Mavi Sánchez-Vives, Juan Abolafia, Emili Balaguer-Ballester, Greg DeAngelis, Jordi Navarra, Jaime de la Rocha y Albert Compte ya que su ayuda en muchos casos ha sido fundamental para que esta tesis fuera posible.

Igualmente, me gustaría agradecer todo el apoyo recibido por las familias Nogueira y Mañas durante estos cinco años, así como a los González-Zapata, Miguel, Teresa y Elena, por haberme hecho sentir siempre como un miembro más de la familia.

Por último, me gustaría dedicar esta tesis a mis padres Daría y Gelasio por todo lo que me han dado en la vida. Gracias a la educación, cariño y soporte que he recibido durante todos estos años he podido desarrollarme como la persona que soy hoy en día y como el científico que me gustaría ser en el futuro. En especial me gustaría dedicar esta tesis a la memoria de mi padre, Gelasio Nogueira Esmorís (1947-2014), al que hubiera hecho muy feliz poder estar compartiendo este momento tan importante conmigo.

## **Abstract**

One of the most important goals in theoretical neuroscience is to determine what are the fundamental principles underlying the processing of information in the brain and ultimately characterize the link between neuronal activity and behavior. Even though many important steps have been done in this direction, we are still far from providing a clear and robust answer to this question. In this thesis I will present a set of results that will be analyzed under the encoding-decoding framework in decision-making, a fundamental part of cognition. In particular, I will present a set of electrophysiological, behavioral and mathematical results that have been used to study the encoding of information on the cortex of behaving monkeys and the integration of sensory with prior evidence on rats performing an outcome-coupled perceptual decision-making task.

**Keywords:** theoretical neuroscience, computational neuroscience, decision making, encoding-decoding, Bayesian inference, pairwise correlations, global activity, visual cortex, prefrontal cortex, orbitofrontal cortex, prior information, empirical artificialism.

## Resum

Un dels objectius més importants de la neurociència teòrica és determinar quins són els principis fonamentals subjacents en el processament de la informació al cervell i en última instància caracteritzar el nexa entre l'activitat neuronal i el comportament. Tot i que s'han produït avenços importants en aquesta direcció, encara estem lluny de poder proporcionar respostes clares i robustes a aquesta pregunta. En aquesta tesi presentaré un conjunt de resultats que han estat analitzats des del paradigma de codificació-decodificació en la presa de decisions, una part fonamental de la cognició. En particular, presentaré un conjunt de resultats electrofisiològics, comportamentals i matemàtics que han estat utilitzats per a estudiar la codificació d'informació a l'escorça de micos conductuals i en la integració de l'evidència prèvia amb la sensorial en rates realitzant una tasca perceptual de presa de decisions acoblada a la seva resposta.

**Paraules clau:** neurociència teòrica, neurociència computacional, presa de decisions, codificació-decodificació, inferència bayesiana, correlacions a parells, activitat global, escorça visual, escorça pre-frontal, escorça orbito-frontal, informació prèvia, artificialisme empíric.

## Resumen

Uno de los objetivos más importantes de la neurociencia teórica es determinar cuáles son los principios fundamentales subyacentes en el procesamiento de la información en el cerebro y en última instancia caracterizar el nexo entre la actividad neuronal y el comportamiento. Aunque se han producido importantes avances en esta dirección, aún estamos lejos de poder proporcionar respuestas claras y robustas para esta pregunta. En esta tesis voy a presentar un conjunto de resultados que han sido analizados desde el paradigma de codificación-decodificación en la toma de decisiones, una parte fundamental de la cognición. En particular, voy a presentar un conjunto de resultados electrofisiológicos, comportamentales y matemáticos que han sido usados para estudiar la codificación de información en la corteza de monos conductuales y en la integración de la evidencia previa con la sensorial en ratas realizando una tarea perceptual de toma de decisiones acoplada a su respuesta.

**Palabras clave:** neurociencia teórica, neurociencia computacional, toma de decisiones, codificación-decodificación, inferencia bayesiana, correlaciones a pares, actividad global, corteza visual, corteza pre-frontal, corteza orbito-frontal, información previa, artificialismo empírico.



## Prólogo

Fue por julio de 2012, poco antes de acabar el máster en física teórica, cuando me empecé a dar cuenta de que la cosmología y la gravitación eran campos de la ciencia apasionantes pero alejados del tipo de investigación que yo quería hacer durante los próximos cuatro o cinco años. Llegué a la conclusión de que necesitaba encontrar un campo donde sintiera que estaba haciendo ciencia más tangible y mundana, en comparación con la física teórica, donde a veces sentía que estaba haciendo más metafísica matemática que investigación basada en el método científico.

Estuve buscando doctorados durante algún tiempo y recuerdo que tuve la posibilidad de elegir algunos temas tan interesantes como dispares entre sí: simulación de agujeros negros, estudio de atmósferas de exoplanetas, modelización y análisis de la actividad de circuitos neuronales y hasta la realización de un doctorado en arqueoastronomía. Finalmente decidí dar un salto al vacío y hacer un doctorado en un campo del no había oído a hablar en mi vida, la neurociencia teórica. Fue una gran decisión, aunque totalmente atribuible a una mezcla entre la intuición y el azar.

Durante el primer año de doctorado tengo que reconocer que no entendía nada. Iba haciendo algunas cosas, pero mi trayectoria se asemejaba bastante a la de un invidente avanzando por un camino lleno de obstáculos donde su única guía se basa en las instrucciones dadas por su mentor en un idioma desconocido pero familiar. Con el tiempo empecé a abrir los ojos y descubrí un mundo que me fascinó por varias razones, pero principalmente por tratarse de un campo donde se mezclan las matemáticas, la física, la química, la biología, la informática, la psicología, la filosofía y la lingüística, entre otras. Desde el momento en que me di cuenta de lo mucho que me gustaba esta rama de la ciencia, he intentado aprender todo lo que he podido del mayor número de personas posibles.

Gracias a este doctorado he visitado muchos países distintos y he podido conocer a gente muy interesante y estimulante intelectualmente. Me gustaría destacar el verano de 2015 que pasé en el laboratorio del Dr. Julio Martínez Trujillo en London, Canadá. Durante ese tiempo aprendí muchas cosas sobre neurociencia experimental, pero esa estancia sobretodo

me sirvió para ganar confianza en mí mismo y para hacer un salto cualitativo en mi capacidad de generar contenido científico. Durante esa estancia creé, gracias a la ayuda de muchas personas, lo que ha acabado siendo el segundo capítulo de esta tesis.

También tengo que reconocer que he pasado por momentos bastante duros durante estos años. El proyecto principal de mi doctorado, lo que ahora es el tercer capítulo, fue un proyecto que duró más de tres años. La publicación de éste en marzo de 2017 en una revista científica prestigiosa, ha sido uno de los momentos más satisfactorios que he vivido durante mis años del doctorado, pero me costaría afirmar que todas las inseguridades y angustias vividas durante los últimos dos años previos a su publicación han valido la pena. En cualquier caso, lo pasado, pasado está, y me quedo con la experiencia vivida para poder gestionar mejor las situaciones futuras. Asimismo, la experiencia adquirida durante este primer proyecto me ha servido para poder desarrollar el segundo desde una posición mucho más madura científicamente, pero sobretodo personalmente.

También me gustaría mencionar que la enfermedad y posterior fallecimiento de mi padre en otoño del 2014 ha sido el momento más difícil de estos cinco años, por no decir que ha sido uno de los momentos más duros de mi vida. Con el tiempo me he dado cuenta de que inconscientemente pasé el duelo de su pérdida refugiándome en mi trabajo. De alguna forma me fue muy bien tener la mente ocupada durante los meses posteriores, pero sin duda alguna esa experiencia ha marcado un antes y un después en mi vida.

Independientemente de los buenos y malos momentos, analizando la situación con perspectiva, la realización de este doctorado ha sido una gran experiencia de vida. He adquirido mucho conocimiento científico, he conocido a mucha gente muy interesante y he visitado lugares que no sabía que existían. Siempre podría haber hecho algunas cosas mejor, pero también podría haber hecho muchas más cosas peor. Creo que por encima de todo destacaría lo afortunado que me siento de haber podido descubrir el mundo del *machine learning*, el análisis de datos y la inteligencia artificial. Estamos viviendo un momento histórico donde estas disciplinas han pasado a dominar tanto la ciencia como la tecnología mundial y estoy

muy contento de haber podido tener un contacto relativamente profundo y cuantitativo con estas áreas gracias a la realización de este doctorado.

Una anécdota muy curiosa y enriquecedora ocurrida durante este doctorado que me gustaría mencionar brevemente fue cuando, después de publicar el primer artículo, nos llamaron desde diferentes medios de prensa para preguntarnos sobre nuestra investigación ya que les había parecido interesante y querían publicarlo. Recuerdo especialmente un par de conversaciones de más de media hora con unos periodistas de *RTVE* y del *ABC*. Tengo que reconocer que, aparte de la satisfacción que experimenté por el hecho de que medios nacionales se hicieran eco de nuestro trabajo, no deja de sorprenderme como transformaron todo lo que les conté en un mensaje bastante más digerible para el público no-científico que sobretodo se caracterizaba por ser incorrecto y relativamente manipulado. Aún recuerdo la frase que me dijo de Rubén después de leer uno de los artículos publicados en un medio digital: ¿Pero qué les has contado? Aprendí mucho de cómo funciona la sociedad actual gracias a esa experiencia. En caso de que el lector de esta tesis esté interesado, fueron publicados dos artículos en dos periódicos: página 44 del *ABC* y página 15 de *El punt avui* del 25 de Marzo de 2017. También fuimos publicados online tanto en medios exclusivamente digitales como en medios que combinan la publicación en papel con la digital.<sup>1</sup>

Para finalizar me gustaría comentar que gracias a la realización de este doctorado he tenido la suerte de conocer a mi compañera de viaje Teresa, con la que he compartido muchos buenos momentos durante los últimos años y con la que ahora empiezo una nueva etapa en otro continente. Empezamos nueva vida como marido y mujer, y tal y como hemos hecho hasta ahora, estoy seguro de que vamos a saber exprimir al máximo todas y cada una de la experiencias que la vida nos tiene preparadas.

Barcelona, 27 de septiembre de 2017

---

<sup>1</sup>No tendría mucho sentido escribir las direcciones de las páginas en este prólogo, pero en el caso de que el lector tenga interés, una forma de acceder a ellas es buscando *ramon nogueira orbitofrontal cortex* en *google*.



# Contents

<b>Agradecimientos</b>	<b>v</b>
<b>Abstract</b>	<b>vii</b>
<b>Resum</b>	<b>viii</b>
<b>Resumen</b>	<b>ix</b>
<b>Prólogo</b>	<b>xi</b>
<b>List of figures</b>	<b>xix</b>
<b>1 INTRODUCTION</b>	<b>1</b>
1.1 Encoding and Decoding . . . . .	2
1.2 Encoding of information . . . . .	9
1.2.1 Common trial-by-trial variability . . . . .	12
1.2.2 Global activity as an internal state . . . . .	23
1.3 Prior Information . . . . .	32
1.3.1 Behavioral protocols for studying prior . . . . .	33
1.3.2 The OFC in decision-making . . . . .	40
<b>2 ENCODING OF INFORMATION AND ITS LINK WITH BEHAVIOR</b>	<b>49</b>
2.1 Introduction . . . . .	50
2.2 Results . . . . .	52

2.2.1	Tuning and noise features of the neural code affecting encoding of information . . . . .	52
2.2.2	Analytical Decoding Performance on <i>in vivo</i> recordings . . . . .	53
2.2.3	Perturbing the original dataset by the bootstrap method . . . . .	59
2.2.4	Encoded information is not affected by mean pairwise correlation or global activity . . . . .	61
2.2.5	Selectivity length and projected precision are the most important features for behavior . . . . .	65
2.2.6	A biologically-constrained neuronal model accounts for the experimental findings . . . . .	71
2.3	Discussion . . . . .	77
2.4	Methods . . . . .	80
2.4.1	Analytical expressions for encoded information . . . . .	80
2.4.1.1	Analytical DP for an arbitrary linear classifier . . . . .	81
2.4.1.2	Optimal linear classifier . . . . .	83
2.4.1.3	Analytical DP for the optimal linear classifier . . . . .	83
2.4.1.4	Analytical DP for shuffled neuronal recordings . . . . .	85
2.4.1.5	Analytical DP under suboptimal readouts . . . . .	86
2.4.1.5.1	Variability-blind classifier . . . . .	86
2.4.1.5.2	Correlation-blind classifier . . . . .	86
2.4.1.6	Analytical expression for DP and differential correlations . . . . .	87
2.4.2	Analytical vs fitted decoding performance . . . . .	91
2.4.3	Bootstrap analysis . . . . .	92
2.4.3.1	General Method . . . . .	92
2.4.3.2	Conditioning Method . . . . .	93
2.4.4	Network Model . . . . .	95
2.4.4.1	Generative model . . . . .	96

2.4.4.2	Information encoded by the model . . .	99
2.4.4.3	Analysis of surrogate data . . . . .	100
2.4.5	Experimental Methods . . . . .	101
2.4.5.1	Coarse random dot motion task recorded in MT . . . . .	101
2.4.5.1.1	Subjects and recordings . . .	101
2.4.5.1.2	Experimental task . . . . .	102
2.4.5.1.3	Neuronal data analysis . . .	103
2.4.5.2	Attentional task recorded in LPFC 8a .	104
2.4.5.2.1	Subjects and recordings . . .	104
2.4.5.2.2	Experimental task . . . . .	105
2.4.5.2.3	Neuronal data analysis . . .	106
2.4.5.3	Fine motion discrimination task recorded in MT . . . . .	108
2.4.5.3.1	Subjects and recordings . . .	108
2.4.5.3.2	Experimental task . . . . .	109
2.4.5.3.3	Neuronal data analysis . . .	109

### **3 INTEGRATION OF PRIOR WITH SENSORY INFORMATION 113**

3.1	Introduction . . . . .	114
3.2	Results . . . . .	116
3.2.1	Animals use task-contingencies to improve per- formance . . . . .	116
3.2.2	Single-cells encode upcoming choice and second- order prior . . . . .	121
3.2.3	OFC encodes immediate prior and anticipates fu- ture choices . . . . .	124
3.2.4	Build-up of choice-related neuronal signals . . .	129
3.2.5	Expected value and outcome representations . .	130
3.2.6	Behaviorally irrelevant prior is not represented in OFC . . . . .	133
3.2.7	Population decoding reveals a hierarchy of vari- ables in OFC . . . . .	135

3.3	Discussion . . . . .	137
3.4	Methods . . . . .	142
3.4.1	Behavioral task . . . . .	142
3.4.2	Presentation of acoustic stimuli . . . . .	144
3.4.3	Logistic regression of behavior . . . . .	144
3.4.4	Psychometric curve analysis . . . . .	146
3.4.5	Surgical procedure . . . . .	147
3.4.6	Tetrodes and microdrives . . . . .	148
3.4.7	Electrophysiological recordings from awake, freely moving rats . . . . .	149
3.4.8	Experimental setup . . . . .	149
3.4.9	Neural data . . . . .	150
3.4.10	ROC analysis . . . . .	151
3.4.11	Generalized linear model for neuronal activity . .	151
3.4.12	Correlation of regression weights . . . . .	156
3.4.13	Population decoding . . . . .	158
3.4.14	Conditioned population decoding . . . . .	159
3.4.15	Information ranking . . . . .	160
3.4.16	Correlation between behavior and neuronal activ- ity across rats . . . . .	160
3.4.17	Power Analysis . . . . .	161
<b>4</b>	<b>DISCUSSION</b>	<b>163</b>
4.1	Encoding of Information . . . . .	164
4.2	Prior Information . . . . .	168
4.3	Future research . . . . .	173
4.4	General perspective . . . . .	177
	<b>Bibliography</b>	<b>183</b>

# List of Figures

1.1	Perception as an encoding-decoding process . . . . .	4
1.2	Pairwise correlations limit information . . . . .	15
1.3	Differential correlations . . . . .	20
1.4	Attention modulates the tuning curve . . . . .	24
1.5	Global activity does not affect information . . . . .	30
1.6	Prior information and the psychometric curve . . . . .	36
1.7	Anatomy of the OFC . . . . .	41
1.8	OFC anticipates choices . . . . .	44
2.1	Selectivity length and projected precision . . . . .	54
2.2	Four datasets in this study . . . . .	56
2.3	Analytical expression is a good approximation . . . . .	58
2.4	Analytical expression is a good approximation 2 . . . . .	60
2.5	Analytical expression for the shuffled datasets . . . . .	62
2.6	LDA is the best classifier . . . . .	64
2.7	Perturbation method based on bootstrapping . . . . .	66
2.8	Only SL and PP affect information . . . . .	67
2.9	Behavioral performance for the different tasks . . . . .	69
2.10	SL and PP are the most influential on behavior . . . . .	70
2.11	Encoded information vs performance . . . . .	72
2.12	Robustness under other metrics and conditionings . . . . .	73
2.13	Model of the cortex . . . . .	74
2.14	Model of the cortex 2 . . . . .	76
3.1	Rats use the additional contingencies . . . . .	118

3.2	Outcome-coupled hidden Markov chain . . . . .	120
3.3	Psychometric curves . . . . .	121
3.4	Factors influencing choice. Logistic regression . . . . .	122
3.5	OFC neurons encode past information . . . . .	123
3.6	OFC neurons encode essential quantities . . . . .	125
3.7	Neurons in IOFC integrate information . . . . .	126
3.8	Encoding of information only after correct . . . . .	128
3.9	Stability of the encoding . . . . .	131
3.10	Temporal evolution of the encoding for each rat . . . . .	132
3.11	Temporal evolution of the encoding in passive task . . . . .	134
3.12	Population decoding of prior and upcoming choice . . . . .	136
3.13	Hierarchy of encoded variables . . . . .	138
3.14	Behavior vs encoded information . . . . .	139
3.15	Linear regression produces equivalent results . . . . .	153

# Chapter 1

## INTRODUCTION

Computational or theoretical neuroscience is the scientific discipline that aims to determine the fundamental principles of information processing in the brain and ultimately characterize the link between brain activity and behavior. It lies on the assumption that all behavior, from the most basic forms of interaction with the environment to the highest instances of rationality and cognition, is fully determined by the state, dynamics and interaction of the whole set of fundamental information-processing units composing the nervous system, the neurons.

Even though this dissertation just represents a tiny contribution to the field of theoretical neuroscience, trying to provide an answer to this question has been the main driving force of this doctoral thesis. In this section I will start by introducing the main theoretical framework where this work is embedded. Then, I will present a brief review of the set of studies that I find more relevant, which will naturally bring the reader to the main questions I have tried to answer during this dissertation. The rest of the thesis is dedicated to answer these questions by providing mathematical, behavioral and electrophysiological evidence.

## 1.1 Encoding and Decoding

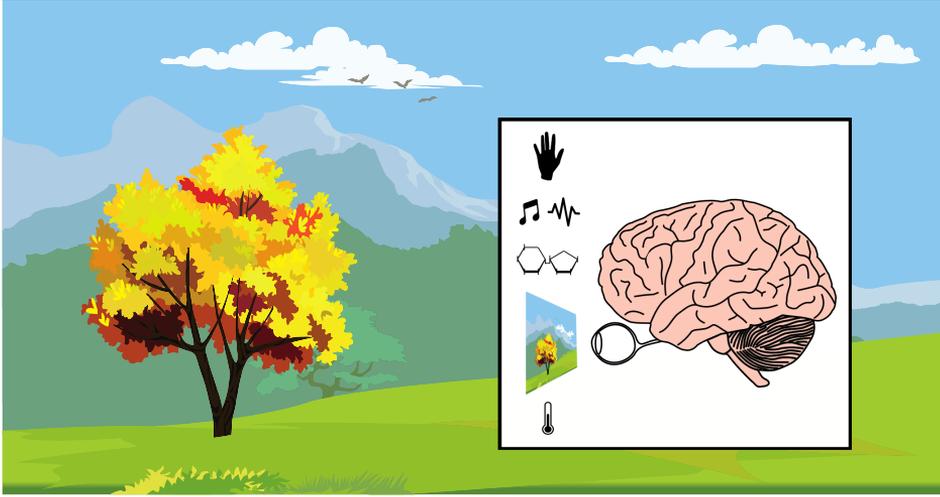
One of the most important goals of any living being is making sense of its environment and use this information to guide its behavior and ultimately maximize its chances for survival [1]. Information about the environment is gathered through a set of channels (senses) which provide the living being with the external-world information (sensory information). The amount of necessary information to fully characterize the external world is potentially infinite. Biological systems have to deal with finite resources and intrinsic noise and therefore the gathered sensory information can be understood as a corrupted or noisy representation of the real hidden state of the world. This leads to the general idea that the task of collecting environmental information and process it to get an accurate description of the real state of the world (perception) is an inference process, in particular a probabilistic inference process [2–5].

Each living being is optimized for its particular environment [1, 6]. The process of inferring the hidden state of the world through noisy samples of it, or perception, has undergone an optimization process as well, which has mostly been driven by the physical properties of this environment. The laws of physics, and in particular, the way matter and energy behave, can be very different among all possible worlds living being are embedded. For instance, the nervous system of humans has been optimized in an environment where the quantum properties of nature can be negligible, the speeds of moving objects correspond to a very small fraction of the speed of light and the surrounding masses are not large enough to produce gravitational interactions (besides the Earth, the Sun and the Moon). Humans have been optimized to mainly make sense of the electromagnetic interactions (which provide the structure of matter as we understand it) and the gravitational force that Earth is permanently performing on all the bodies around us. Our finite external-information gathering channels have evolved in such a way that most of the information we collect comes from the visible spectrum portion of the scattered electromagnetic radiation by the surrounding bodies, the energy perturbations of matter that propagate through space (and time) and the chemical,

thermal and mechanical identification of the different elements surrounding or entering our body [7]. In a different environment, for instance, in the world of the very small, the gathering of information could have been optimized to make sense of nuclear forces and we would perceive what for us are fundamental physical quantities like time, space, energy and velocity in a totally different way. On the contrary, if we would have been optimized for the world of the astrophysical masses, we would have a very different perception of the gravitational interaction and perhaps the relativity of time and space would be integrated in our inferential process in a more insightful way [6].

Regardless of the particular properties of the environment, the inference of the hidden state of the world performed by any agent in an arbitrary world can be understood as an encoding-decoding process (Fig. 1.1) [4, 8–10]. Under this framework, the real state of the world  $s$  generates a representation on the activity of the sensory channels  $r$ . As stated before, the biological constraints together with the intrinsic noise present in any biological system, make  $r$  to be a corrupted version of  $s$ , being  $r$  compatible with many different and potentially exclusive real states of the world  $s$ . The high-dimensional real state of the world  $s$  will be compressed to a lower dimensional representation and perturbed with this intrinsic noise to produce  $r$ . The process of representing the external state of the world by the biological sensory devices is known as the encoding process. Once the information is encoded by the information-gathering units of the sensory system, it can be read-out by other units to create percepts or estimates of the hidden state of the world an ultimately guide behavior.

Even though there is not a general agreement on how the inference process is performed, it exists a large set of experimental evidence that suggest that humans and other animals behave as optimal Bayesian observers [11–14]. Under this framework sensory information about the external world is represented by a conditional probability density function, the posterior probability distribution. The posterior distribution provides the probability for each state of the world  $s$  given the corrupted representation  $r$  by the sensory system. For instance, the perceived height of a tree is not going to be represented as a single number  $h$  but as a



**Figure 1.1:** Perception as an encoding-decoding inference process. The real hidden state of the world  $s$  generates a corrupted representation  $r$  on the activity of the sensory system (encoding). The encoded information will be read-out by the other units in order to produce percepts (decoding) and ultimately guide behavior.

probability distribution over all the possible real tree heights given the information provided by the retina  $p(h|r)$ . By representing information in this way, the system is able to integrate information over space and time, to combine it from different sources and to propagate it across all the information processing stages efficiently. Bayesian inference consists on four main terms: the posterior probability function, the likelihood, the prior and the cost function. The posterior distribution, as stated before, represents the knowledge of the observer about the underlying parameter to be inferred  $p(h|r)$  (belief). The likelihood function, also known as the generative model, represents the probability distribution for the different response instances of the encoding system given a particular value of the hidden parameter  $p(r|s)$ . From now on, under the Bayesian framework, I will refer to the likelihood function, the encoding stage and the generative model indistinctly, as they all represent the probability of a particular

sensory representation  $r$  given the real state of the world  $s$ . The prior probability accounts for the *a priori* probability distribution of the underlying parameter  $p(s)$ . The prior probability can be hardwired in the biological system or it can be acquired through the interaction with the environment. This interaction can span from a whole life experience to a few seconds. [15–17]. The relationship between these three quantities is captured by the Bayes’ theorem

$$p(s|r) = \frac{p(r|s)p(s)}{p(r)}, \quad (1.1)$$

or, if we consider  $p(r)$  to be just a normalization constant, then

$$p(s|r) \propto p(r|s)p(s). \quad (1.2)$$

Once the probability distribution for the hidden state of the world has been fully characterized, the next step for the observer is to choose the optimal value  $\hat{s}$ , which will be determined by the loss function. The optimal decision-making rule for an agent will be

$$\operatorname{argmin}_{\hat{s}} \int p(s|r) L(\hat{s}, s) ds, \quad (1.3)$$

or equivalently

$$\operatorname{argmin}_{\hat{s}} E_{s|r}[L(\hat{s}, s)], \quad (1.4)$$

where  $L(\hat{s}, s)$  is the loss function. This is an optimization problem that depends on the loss function. For instance, if the loss-function is defined as

$$L = (s - \hat{s})^n, \quad (1.5)$$

then the optimal  $\hat{s}$  will be the mode, the median and the mean of  $p(s|r)$  for  $n = 0$ ,  $n = 1$  and  $n = 2$  respectively.

The Bayesian approach corresponds to the optimal approach, that is, it provides an upper bound for the discrimination accuracy given a noisy representation of the state of the world by the encoding units [18]. Even

though it can be generalized to more classes or to the continuous case, for the sake of simplicity from now on I will consider a binary classification task, where the task is to correctly distinguish between two different classes ( $s_1$  and  $s_2$ ) given the sensory response. This is so because most of the electrophysiological experiments performed on behaving animals are generally characterized by a decision-making protocol where animals have to decide between two alternative options. Although there is a mathematical framework that extends to the multiple and continuous options, there are several experimental limitations for such protocols, like the limitation on the number of available trials per option, or the difficulty for some animal species to learn more complicated tasks that require higher cognitive capabilities. Moreover, from now on I will refer to information as the accuracy of an inference process to properly estimate the encoded variable. Even though the inferential process can be performed in many different ways, in this thesis I will refer to linear inferential processes, or in other words, parameter estimations performed by a linear transformation of the sensory representation  $r$ . For instance, when the parameter to be inferred is continuous, information is characterized by the linear Fisher information, the variability of an optimal linear estimator. When the parameter to be inferred is discrete, I will refer to information as the percentage of correct classifications (decoding performance; DP) performed by the optimal linear estimator.

As stated before, one of the most important parts when performing such a task is to characterize the posterior distribution  $p(s|r)$ . This can be achieved by a generative model or by a discriminative model [19]. In the generative model our aim is to first characterize the likelihood function  $p(r|s)$  (encoding) and then find the posterior distribution using Bayes' theorem and the prior distribution. On the contrary, when taking the discriminative approach, our aim is to model the posterior directly. Regardless of the approach we take, by the Bayes theorem, the posterior on  $s_1$  can be expressed as

$$p(s_1|\mathbf{r}) = \frac{1}{1 + \exp(-a)}, \quad (1.6)$$

where

$$a = \ln \left( \frac{p(\mathbf{r}|s_1)p(s_1)}{p(\mathbf{r}|s_2)p(s_2)} \right), \quad (1.7)$$

and where

$$p(s_2|\mathbf{r}) = 1 - p(s_1|\mathbf{r}). \quad (1.8)$$

When taking the generative approach our aim is to model  $p(\mathbf{r}|s_1)$  and  $p(s_1)$ . If the inferring system knows the generative model and the prior distribution, modeling the posterior distribution is straightforward. This is the desired scenario because characterizing the posterior distribution on  $s$  by the combination of the generative model  $p(r|s)$  and the prior  $p(s)$  provides the inferring agent with explicit representations of all sources of information. However, this is not always possible because a full characterization of both  $p(r|s)$  and  $p(s)$  can be very costly in terms of acquisition, storage and processing of all the necessary data, especially when the sensory representation is multidimensional.

On the contrary, in the discriminative approach the aim is to characterize the posterior distribution directly, regardless of the exact shape the likelihood and the prior distributions can take. Although the desired approach is the former, very good approximations or even the optimal solution can be achieved when modeling directly the posterior distribution. For instance, we can consider in eq. (1.6)  $a = \boldsymbol{\omega}\mathbf{r} + \omega_0$ , which corresponds to a linear read-out of the sensory representation  $\mathbf{r}(s)$ . The set of linear weights can be found by maximizing the probability of obtaining the experimental data under this model [19]. This approach makes no assumption about the mapping between  $r$  and  $s$ , so finding a close approximation to the optimal set of weights is just a matter of how large is the dataset used to fit the model. In particular, if the underlying generative model is Gaussian with equal covariance matrices for  $s_1$  and  $s_2$ , then (1.7) becomes

$$a = \boldsymbol{\omega}\mathbf{r} + \omega_0, \quad (1.9)$$

where

$$\boldsymbol{\omega} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad (1.10)$$

$$\omega_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \left( \frac{p(C_1)}{p(C_2)} \right). \quad (1.11)$$

If the underlying generative model is indeed gaussian and the optimization process is provided with enough data, the discriminant approach would be able to reach eqs. (1.10) and (1.11) through the optimization process. Moreover, as stated in eq. (1.11), in the discriminant approach prior information could be potentially incorporated to the model implicitly. The decision-making agent can be taking into account this extra source of information by just a bias term in the argument of the expression for the posterior distribution. This is in contrast with the generative approach, where the prior knowledge has to be explicit in order to generate a representation of the posterior distribution out of the likelihood function (eqs. (1.2) (1.7)).

Nevertheless, it exists a third possible approach, the so-called Discriminant function [19] or naive approach. Under this framework a simple rule of thumb or phenomenological algorithm is taken, without any explicit reference to the posterior distribution, the likelihood or the prior. Depending on the particular problem, a very simple rule can achieve very high performances because it can correspond to very good approximations for the generative or discriminative models. The perceptron algorithm [20] or the fisher discriminant [21, 22] would correspond to examples of discriminant functions that can achieve very high performances out of relatively simple non-probabilistic data-based algorithms. For some particular underlying generative structures of the data, these algorithms are equivalent to the optimal probabilistic approaches, for instance, if the sensory representation  $r$  follows the gaussian distribution, the optimal set of weights  $\omega$  given by the generative approach (eqs. (1.10) and (1.11)) are equivalent to those given by the fisher discriminant analysis [19, 21, 22]. Naive decision-making algorithms can also take into account prior information, both in an implicit and explicit way. Indeed, it exists behavioral evidence that in some perceptual decision-making task humans follow sub-optimal decision rules where both sensory and prior information are being used [23].

Therefore, either from the Bayesian or from a more general not necessarily optimal approach, perception and decision-making can be understood as a process where the current sensory information provided by the encoding units is combined with the prior information to adaptively guide behavior. This will be the general guideline of this thesis. In the first part of the thesis (chapter 2) I will present a set of both theoretical and experimental results where I will analyze the role of neuronal ensemble tuning and trial-by-trial variability on the amount of encoded information by a network and its effects on behavior. On the second part of the thesis (chapter 3) I will present an study performed on behaving rats where we analyze the role of prior information on behavior and their correlates with neuronal activity on the orbitofrontal cortex (OFC). In order to smoothly guide the reader to both studies, in the following subsections I am providing a specific introduction to each of them. I will review the most important works on each field and introduce the questions I will try to answer for the next chapters.

## 1.2 Encoding of information

Information from the external world is encoded by the brain in the firing rate of sensory neurons. A neuron is encoding a particular external parameter if it has a different mean firing rate for each different value the external parameter can take. This concept is known as the tuning curve and it is defined as

$$f(s) = \int rp(r|s)dr. \quad (1.12)$$

It is important to note that in this expression the noise of the response is averaged out, so when we talk about tuning curve we actually talk about the mean neuronal response for a each different value of the parameter to be encoded. The number of spikes in a particular time window is known as the spike count ( $n$ ), and the number of spikes per time unit is known as the firing rate of a neuron ( $r$ ). There is a linear relationship between them  $r = n/\Delta t$ . From now on, when talking about neuronal activity I

will be referring to any of them indistinctly. There are many examples of cortical neurons encoding parameters of the external world. One of the earliest and most important one corresponds to an study performed on the primary visual cortex of anesthetized cats [24] where it was found that neurons fired more vigorously to particular orientations of light bars in a dark background. This way it was shown for the first time that neurons encode the external world through their firing rate, in this particular case the orientation of bars in the visual field of the animal. Many other examples have been found since then on other sensory modalities like hearing [25, 26], touch [27, 28], vestibular [29, 30], olfactory [31, 32] and others.

Tuning curves are defined as the mean activity per stimulus condition (see eq. (1.12)) because it is a widely tested experimental evidence that neuronal responses are variable both in the timing and number of spikes when presented with identical stimuli [33, 34]. More formally, the trial-by-trial variability for neuron  $i$  on condition  $s$  is defined as

$$\sigma_i^2(s) = E[(n_i - E[n_i])^2] \simeq \frac{1}{M-1} \sum_{j=1}^M (n_{ij}(s) - \langle n_i(s) \rangle)^2, \quad (1.13)$$

where  $E[\cdot]$  is the expectation,  $M$  is the number of trials when presented stimulus  $s$ ,  $n_{ij}(s)$  is the number of spikes depicted by the neuron  $i$  on trial  $j$  and  $\langle n_i(s) \rangle$  is the mean number of spikes across trials when presented stimulus  $s$ , or equivalently  $\langle n_i(s) \rangle = \Delta t f_i(s)$ .

When inferring what is the underlying parameter governing the activity of a neuron, the problem becomes non-trivial when we obtain different responses for identical presentations of external (or internal) world variables. Many authors consider spiking variability just as a form of noise that harms signal processing [34, 35], even though it is not clear yet if that is the case or it just represents deterministic coding of experimentally uncontrolled variables. Assuming spiking variability arises as a form of noise, two main different explanations have been proposed to account for its origins. It has been proposed to come from the intrinsic noise of the

system, such as ion channels or stochastic synaptic release [36, 37] or/and from the complex dynamics of the neural networks [38–40].

It has also been shown experimentally that the trial-by-trial variability depicted by single neurons is approximately proportional to its mean activity  $\sigma_i^2(s) \simeq \langle n_i(s) \rangle$  [34, 41], a typical feature of a Poisson point process. Although it is more or less accepted that the neuronal activity is governed by this probability distribution, some recent studies have been able to fit with great accuracy the neuronal statistics by including an additional term of common variability before the Poisson step [42–44].

Based on these experimental findings a very useful tool in computational neuroscience consists in modeling the noisy activity of neurons as a stochastic process where each trial or stimulus presentation corresponds to a particular realization of this process. This way the encoding process  $p(r|s)$  for single cells can be naturally formalized as

$$E[n] = g^{-1} \left( \sum_{i=1}^k \omega_i x_i + \omega_0 \right), \quad (1.14)$$

where  $E[n]$  is the expectation on the number of spikes across trials,  $g(\cdot)$  is the link function and  $\sum_{i=1}^k \omega_i x_i + \omega_0$  is a weighted sum of the factors that could be influencing the spike count of a neuron in a particular trial (linear projection of  $\{x_i\}$  to a one-dimensional space). The trial-by-trial variability can potentially follow many different probability distributions like the Gaussian, Poisson, Binomial, etc., depending on the particular features of the problem. This family of models is known as Generalized Linear Models (GLMs) and in the methodology of chapter 3 a detailed description of the Poisson GLM can be found.

Inference on the presented stimulus  $s$  from the neuronal activity can only be performed when  $f(s)$  is not flat. In particular, the larger the modulation of  $f(s)$  with respect to  $s$ , the easier the inference on the stimulus presented. This can be formalized by the derivative of the tuning curve  $f'(s)$  with respect to  $s$ . Additionally, as stated before, if neurons were fully deterministic, the inference problem would be trivial. However, neurons do not present the same activity pattern when presented with

identical stimuli. The amount of information encoded by a single cell (accuracy when inferring the underlying parameter) can be fully determined by these two factors: how large is the change of  $f(s)$  with respect to  $s$  and how variable is the response of that neuron for identical stimulus presentations. However, neuronal information is not only encoded at the single cell level but at the network level. Therefore the concepts of tuning and variability need to be extrapolated to populations of neurons. The tuning of a neural population becomes a vector  $\mathbf{f}(s)$ , being the  $i$ -th entry the tuning curve of neuron  $i$ . The trial-by-trial variability becomes a matrix  $\Sigma(s)$ , where the diagonal elements account for the single cell variability  $\Sigma_{ii}(s) = \sigma_i^2$  and the off-diagonal terms account for the covariance between neuron  $i$  and neuron  $j$ ,  $\Sigma_{ij} = \text{cov}(n_i, n_j)$ .

After these definitions, my purpose for the upcoming section is to present a brief review on the different studies aiming to identify the effects of the common trial-by-trial variability (section 1.2.1) and global activity (section 1.2.2) on the amount of information encoded by a network. This review will lead naturally to chapter 2 of my thesis. In this chapter I will present an study that aims to identify what are the exact features of the neural code affecting the amount of information encoded by a population of neurons and their effects on the performance of behaving animals.

## 1.2.1 Common trial-by-trial variability

Trial-by-trial fluctuations are not independent among neurons but it exists a small shared component, the so-called pairwise correlations. Pairwise correlations are often quantified with the Pearson correlation of the spike count of the neuronal pair to repeated presentations of the same stimulus. From now on I will refer to the Pearson correlation between the spike count of neuron  $i$  and neuron  $j$  as  $\rho_{ij}$ , and to the mean Pearson correlation across all pairs of a neuronal population as  $r_{sc}$  (or  $\langle r_{sc} \rangle$ ). The relationship between the covariance and the Pearson correlation between neuron  $i$  and  $j$  is  $\Sigma_{ij} = \sigma_i \sigma_j \rho_{ij}$ . Pairwise correlations have usually mean values of 0.01 - 0.2 for nearby neurons, although the magnitude depends on a broad range of factors, like response strength, the time window used to measure

correlations, spike sorting conventions and internal states [45]. The role of pairwise correlations on the encoding of information has been widely studied [35, 46–48] and in the following I am briefly reviewing the most important studies addressing this question. It is important to recall that it exists a larger literature than the presented below, but I have specifically chosen the studies that in my opinion are more relevant to the field and in particular to my thesis.

In 1994 E. Zohary and colleagues [35] published a study where they aimed to identify what was the role of pairwise correlations on the amount of information encoded by a network. Their motivation was based on the experimental findings that the behavioral accuracy of monkeys in a perceptual discrimination task [49] was not considerably better than the discrimination accuracy of single neurons recorded in their middle temporal cortex (MT or V5). Neuronal discrimination accuracy here is referred to the amount of information encoded by the neuron about the underlying parameter ruling its activity, or in other words, the performance shown by an ideal observer on inferring the identity of the stimulus presented to the cell.

Even though many possible explanations have been proposed for this problem during the last twenty years, in my opinion this is still one of the most interesting unsolved problems in computational neuroscience. Neurons in MT are generally characterized for encoding motion features of the visual field such as direction of motion and velocity [49, 50]. If information in the brain is mainly processed in a feedforward way, downstream neurons could read out the activity of MT neurons to infer what was the stimulus (direction of motion) presented to the monkey. As stated before, the response of neurons to identical stimulus presentations is variable, hence by reading out simultaneously from a pool of MT neurons the noise can be averaged out and get a statistically reliable signal that can accurately guide behavior [34]. By pooling out from many neurons, the variability at the inferential stage could potentially vanish and produce infinitely accurate behavior. Let's express this concept formally. As stated in a previous section, information on the encoding of a continuous variable can be characterized by the variability of the estimate across tri-

als. If the task of a downstream neuron in a given trial  $j$  is to infer the mean parameter ( $\lambda$ ) underlying a set of identical and independent Poisson processes for a particular time window  $T$ , then the optimal strategy is to average out the spike count of all neurons

$$\hat{\lambda}_j = \frac{1}{NT} \sum_{i=1}^N n_{ij}. \quad (1.15)$$

The variability on the estimate of  $\lambda$  is then

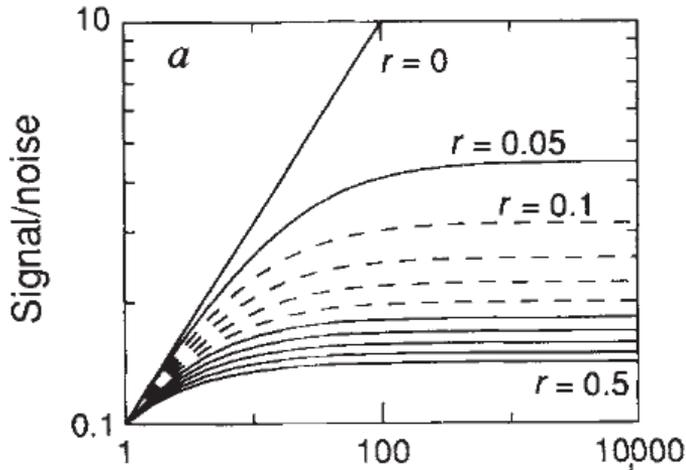
$$\text{Var}(\hat{\lambda}) = E[\hat{\lambda}^2] - E[\hat{\lambda}]^2 = \left( \frac{\lambda}{NT} + \lambda^2 \right) - \lambda^2 = \frac{\lambda}{NT}, \quad (1.16)$$

where we have made use of the identity  $\text{Var}(n) = E[n] = \lambda T$  for a Poisson process. This expression converges to zero for large ensemble sizes. In other words, the accuracy that can be obtained when estimating the encoded parameter by a large population of independent and identical neurons tends to be infinite for large ensemble sizes. The question then is, if sensory percepts are obtained by averaging out the activity of neuronal populations, how is it possible that the accuracy of a whole monkey is not outstandingly superior than the accuracy of a single MT neuron?

Their explanation for this apparent paradox was based on pairwise correlations among sensory neurons. If the estimation is performed by a set of identical but uniformly correlated set of neurons, the variability on the estimation (encoded information) can be expressed by (analogous to eq. (1.16) under the presence of correlations and for large ensemble sizes; see [34])

$$\text{Var}(\hat{\lambda}) \simeq \tilde{r}\lambda, \quad (1.17)$$

where  $\tilde{r} = \langle r_{sc} \rangle$ . In Fig. 1.2 they showed that indeed the information that could be extracted from a network by pooling the activity of its units strongly depended on the strength of the correlation. When all units are independent ( $\langle r_{sc} = 0 \rangle$ ) information grows linearly with the size of the read-out neuronal ensemble as expected by eq. (1.16). For this particular



**Figure 1.2:** Pairwise correlations limit the accuracy of an ideal observer to estimate the underlying parameter of a set of Poisson neurons. The larger the correlation among neurons, the lower will be the information plateau for large ensemble sizes. Extracted from [35].

situation it is very surprising that the sensitivity of the whole monkey, which has access to large ensemble sizes (potentially  $10^4$  or  $10^5$  neurons), does not exceed significantly the sensitivity of a single neuron. However, when pairwise correlations start increasing, the amount of information reaches a plateau for a particular ensemble size. The larger the  $\langle r_{sc} \rangle$ , the faster the saturation regime is achieved. In my opinion, this is a very graceful solution for the apparent paradox. If pairwise correlations are present in the sensory cortex (which is an experimentally well tested truth [45, 46]), pooling from many neurons would not necessarily imply large behavioral accuracy and the sensitivity of single neurons and the whole monkey could be of the same order of magnitude.

Even though this was a very elegant study that solved an important problem in neuroscience, some years later it was shown it was not such a general result. The main concerns about this study were that the neuronal ensemble consisted in a set of  $M$  homogeneous neurons and that the information metric was defined as the pooled activity divided by the sum

of all the terms in the covariance matrix, what they called the signal-to-noise ratio. A set of  $M$  homogeneous neurons refers to a collection of neurons that have identical tuning curves, only differing on the position of the preferred stimulus, or the stimulus that depicts the higher neuronal response. Neurons in the cortex, however, have been shown to have heterogeneous tuning curves, each one with a different baseline, width and preferred stimulus [51].

When inferring the underlying parameter driving a stochastic process  $\hat{s}$ , it exists a theoretical upper bound on the accuracy of the estimate. This quantity is known as the Fisher information and in particular it provides us with an upper bound for the inverse of the standard deviation (squared error) of the parameter to be inferred assuming the estimate is unbiased ( $\langle \hat{s} \rangle = s$ ). If the estimate is performed using a linear read-out of the noisy sources

$$\hat{s} = \boldsymbol{\omega}(\mathbf{r} - \mathbf{f}(s_0)) + s_0, \quad (1.18)$$

then the expression for linear Fisher information ( $I_F$ ) is

$$I_F = \mathbf{f}'^T \boldsymbol{\Sigma}^{-1} \mathbf{f}', \quad (1.19)$$

where  $s_0$  is the reference stimulus around which we are inferring  $s$ ,  $\mathbf{f}'$  is the derivative of the tuning curve with respect to the stimulus and  $\boldsymbol{\Sigma}$  is the covariance matrix of the ensemble's activity.

In general, linear Fisher is not equivalent to the full Fisher information, but it is the case for any response distribution of the exponential family with linear sufficient statistics [52]. From now when talking about Fisher information I will be referring to linear Fisher information and the parameter estimate will be also assumed to be performed using a linear read-out (linear decoder). Linear decoders are very relevant in neuroscience because they can be learned and implemented easily in a biological circuit [53, 54] and because the read-out weights can be easily associated with the synaptic weights connecting neurons.

The next relevant paper I wanted to review was published in 1999 by Larry Abbott and Peter Dayan [48]. It is a theoretical study where

realistic tuning and noise properties of the neuronal ensemble are used to derive the role of pairwise correlations on information. They showed that in general pairwise correlations are not harmful, however when considering very particular cases of the tuning and noise properties of the network, they would be harmful for the neural code. It is an important paper because, in contrast with the idea proposed in [35], it claimed for the first time that pairwise correlations were not harmful in general, but it depended on the relationship between the tuning curves and the noise structure. Three correlation matrices were defined in this study

$$\Sigma_{ij} = \sigma^2(\delta_{ij} + c(1 - \delta_{ij})) \quad (1.20)$$

$$\Sigma_{ij} = \sigma^2(\delta_{ij} + c(1 - \delta_{ij}))f_i(s)f_j(s) \quad (1.21)$$

$$\Sigma_{ij} = \sigma^2 e^{-\Delta/L}, \quad (1.22)$$

where  $\sigma^2$  is the trial-by-trial variability of each neuron,  $\delta_{ij}$  is the Kronecker delta,  $c$  is the correlation between pairs of neurons,  $\Delta$  is the difference between the peaks of two different neurons' tuning curves and  $L$  is a parameter that controls the general range of correlations. Matrices (1.20), (1.21), (1.22) are known as additive (or uniform), multiplicative and limited range-correlations respectively. It is important to note that limited-range correlations seem to be a very good approximation for the correlations profile found in many regions of the cortex [35, 55–58]. In this study it was shown that for the additive noise model (eq. (1.20)) information would saturate for large ensemble sizes if all neurons shared the same tuning dependency with  $s$  with a free additive term

$$f_i(x) = p(x) + q_i, \quad (1.23)$$

for any function  $p$  and number  $q_i$  (additive separability). For the multiplicative case the tuning curves family that would limit information for large ensemble sizes would be

$$f_i(x) = p(x)q_i + r(x) + s_i, \quad (1.24)$$

for any function  $p$  and  $r$  and numbers  $q_i$  and  $s_i$  (multiplicative separability). Finally for the limited-range correlations information would only not grow without limits as a function of the ensemble size when  $L$  tends to zero as ensemble size tends to infinite. The most important conclusion that can be extracted from this study is that information in general would grow with the ensemble size in the presence of pairwise correlations unless the tuning curves of the neuronal ensemble present additive separability (uniform correlations) or multiplicative separability (multiplicative correlations), which have not been reported in any cortical recording.

In 2014 a very important paper regarding the role of pairwise correlations on the amount of information encoded by a network was published by Ruben Moreno-Bote and colleagues [47]. In this study it was fully characterized for the first time the exact pattern of correlations that limited the amount of information that could be linearly extracted from a network. Interestingly, the study was grounded on a theoretical scheme that went beyond pairwise correlations per se: the amount of information that can be extracted can never exceed the amount of information entering the network.

In order to further understand this study it is important to be more precise when referring to input information. In the statistical inference framework information is defined as the accuracy on estimating a particular hidden variable of a stochastic process. As stated before Fisher information is a quantity that provides us with an upper bound on the inverse squared error when inferring a parameter in an unbiased way from a noisy source. Now let's consider a set of  $N$  neurons, each one with a tuning curve  $f_i(s)$ . The activity in each trial for the set of neurons will be

$$\mathbf{r}_i = \mathbf{f}(s) + M\mathbf{z}_i, \quad (1.25)$$

where  $\mathbf{z}_i$  is a  $N$ -dimensional vector with each entry being a random number drawn from a zero-mean and unit-variance Gaussian distribution in each trial. In general, the stimulus presented in a particular trial is assumed to be a controlled parameter but this might not be the most general case. The stimulus presented to a network could be itself noisy from one trial to the next, for instance if the presented stimulus  $s$  is the pooled sig-

nal coming from the retina or the cochlea. In this case in each trial the effective stimulus presented to a network will be

$$s_i = s + \delta s_i, \quad (1.26)$$

where the  $\delta s_i$  is the fluctuation term for a particular trial of stimulus real stimulus presented. If we introduce equation (1.26) in (1.25) and assume  $\delta s_i$  is small

$$\mathbf{r}_i = \mathbf{f}(s) + \mathbf{f}'(s)\delta s_i + M\mathbf{z}_i. \quad (1.27)$$

We can evaluate the covariance matrix of the trial-by-trial neuronal activity

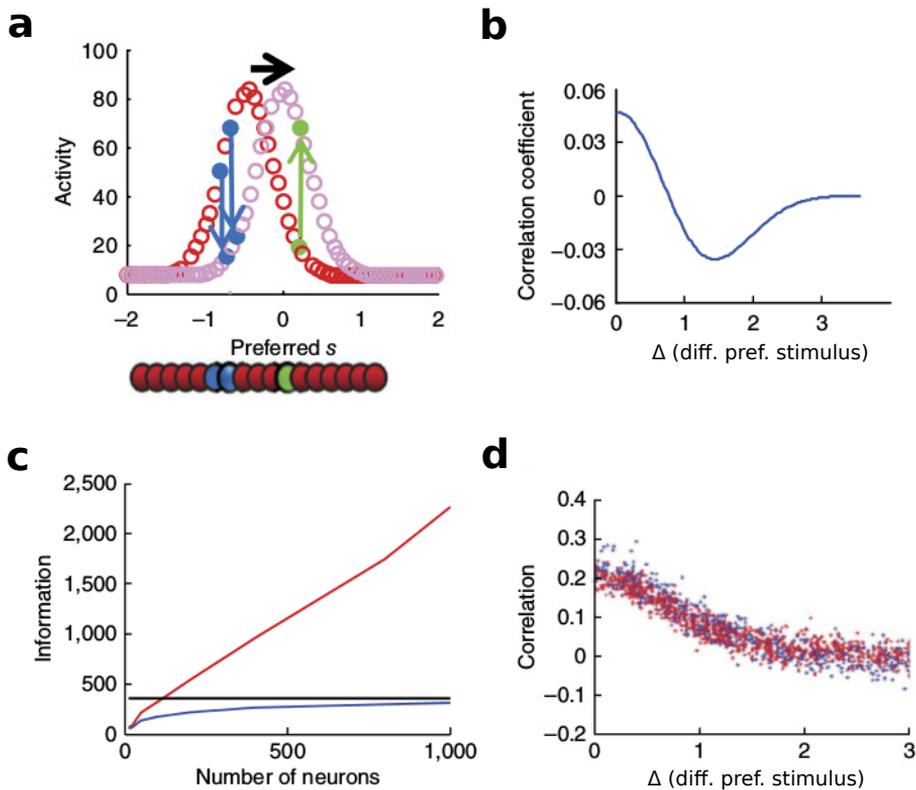
$$\Sigma = \Sigma_0 + \epsilon \mathbf{f}' \mathbf{f}'^T, \quad (1.28)$$

where  $\epsilon = \text{Var}(\delta s)$  and where  $\Sigma_0 = MM^T$  is the original covariance of the population. When including this term, the linear fisher information can be re-evaluated and it takes the form [47]

$$I = \frac{I_0}{1 + \epsilon I_0}, \quad (1.29)$$

where  $I_0 = \mathbf{f}'^T \Sigma_0^{-1} \mathbf{f}'$ . The intuition behind these equations is simple: when the stimulus  $s$  is noisy in a trial-by-trial basis, the joint activity of the ensemble will fall within the line defined by  $\mathbf{f}'$  and the inference process will tend to make an estimation closer to  $s + \delta s$  than  $s$  itself. Because we ignore that this is the case, from our point of view the noise will be just larger (eq. (1.28)). Going back to the starting point, the input information is therefore  $1/\epsilon$ , as this quantity is the upper bound an optimal inference agent would be able to obtain about the real value of  $s$  just by reading from the noisy stimulus introduced in the network.

The correlations induced by the trial-by-trial presentations of the input signal are known as differential correlations and they have a very interesting graphical interpretation. In Fig. 1.3a it is depicted how neurons would fluctuate from one trial to the next when slightly different stimuli



**Figure 1.3:** Differential correlations limit the amount of information encoded by a neural network. (a) Noise at the sensory layer produces positive and negative pairwise correlations when pairs of neurons are similarly and dissimilarly tuned (blue and green dots) respectively. (b) Pattern of differential correlations as a function of tuning similarity. (c) In the absence of differential correlations the amount of information encoded by a neuronal ensemble grows without limits (red line), in contrast with information encoded in the presence of differential correlations (blue line). (d) Differential correlations (red dots) can be masked by non-differential correlations (blue dots) and still have a dramatic effect on information. Extracted from [47].

are presented. The plot shows a population hill of activity where neurons are plotted according to their preferred stimulus. When  $s + \delta s$  is

presented, two similarly tuned neurons (blue) will both present a negative fluctuation in their activities, producing a positive pairwise correlation. However, this same presentation of a corrupted stimulus would generate opposite activity fluctuations on pairs of dissimilarly tuned neurons, producing negative pairwise correlations. Because of this reason the pattern of differential correlations in Fig. 1.3b reminds to a limited-range correlation pattern, where similarly tuned neurons present a larger correlation coefficient than dissimilarly tuned neurons.

In Fig. 1.3c it is shown how information grows with the ensemble size ( $N$ ) in two different situations. The red curve has been calculated using a model where the covariance matrix does not present any additional term parallel to  $\mathbf{f}'$  (eq. (1.28)), while the blue curve does. In both cases information grows with  $N$  but in the presence of differential correlations, the amount of information that can be extracted from the network regarding  $s$  reaches a plateau. The plateau will correspond to the input information, which corresponds to  $1/\epsilon$ . On the contrary, when the input signal is fully deterministic (red curve; infinite information on the input signal) information grows linearly with  $N$ . As stated at the beginning, one of the most interesting parts of this study is that it is shown that pairwise correlations are not the limiting factor per se but noise on the input stage.

If differential correlations are a sufficient and necessary condition for the information to saturate as a function of the ensemble size, is there any way they can be detected at the experimental level? The question is yes, but it is not easy. One possible option would be plotting the correlation coefficient as a function of the pair difference on preferred stimulus orientations and check if it resembles Fig. 1.3b, the pattern of differential correlations. In order to answer this question Moreno-Bote and colleagues plotted correlation coefficients as a function of the difference in preferred stimulus orientation in the presence of differential correlations (blue) and without differential correlations (red) (Fig. 1.3d). The two patterns are indistinguishable because differential correlations so to say are masked by regular non-limiting pairwise correlations. Even though their contribution to the total correlation could be considered to be negligible, their impact on information can be dramatic (Fig. 1.3c).

However, there is yet another possible strategy in order to experimentally detect and characterize differential correlations. If we were able to compute in a reliable manner the real Fisher information of a neuronal ensemble (I will come to this point later), when plotting it against the ensemble size we could claim we have found differential correlations if eventually the relationship starts deviating from the linear case (sublinearly; see Fig. 1.3c (blue)). The question then is: at what ensemble size we would find such a deviation? The answer to this question is not trivial and even though it is not included in this dissertation, it is a current study we are performing to be submitted in the near future.

In my opinion the presented study is the only one that really provides a clear and complete solution to the role of pairwise correlations on information. Pairwise correlations are not the problem but the finite information entering the network, which produces a very particular pattern of trial-by-trial shared noise.

It also exists a set of studies where the different values of mean pairwise correlations are compared to the performance of behaving animals [58–60]. This is in my opinion the best approach to take, as the ultimate goal of theoretical neuroscience is to characterize the link between neuronal activity and behavior. In studies [58, 60] it was found that attention improved perceptual accuracy by reducing the mean pairwise correlations present in V4 neurons. Later on, Ruff and Cohen [59] showed that attention, and consequently perceptual accuracy, affected differently similarly than dissimilarly tuned neurons. Under attention, dissimilarly tuned neurons experienced an increase while similarly tuned neurons experienced a decrease on their mean pairwise correlations. This was an interesting experimental finding as it was clearly aligned with the theoretical prediction stated in [61].

Even though these studies systematically explored the role of mean pairwise correlations on behavior, they made it through an indirect path, the modulation of attention by the spatial allocation of attention. However, studying how correlated activity can influence the performance of behaving animals in a more explicit and direct manner provides a more accurate characterization of the link between behavior and neuronal activ-

ity. In chapter 2 we characterized what is the role of correlations in both the encoding of information and behavioral performance by the combination of a parametric with a non-parametric approach based on bootstrapping from the original dataset.

## 1.2.2 Global activity as an internal state

So far I have reviewed what is the role of individual and shared trial-by-trial variability on the reliability of the neural code or in other words, on the capability of a network to encode information about the hidden state of the world. However, these are not the only magnitudes that could affect information. Another important quantity is the global activity understood as an internal state of the neural network. From now on when I will refer to global activity as the mean activity across the neuronal population for a particular time window

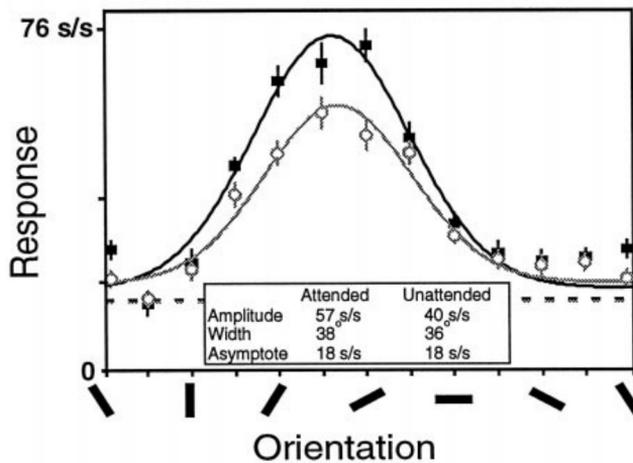
$$g(\Delta t) = \frac{1}{N} \sum_{i=1}^N n_i(\Delta t), \quad (1.30)$$

where  $N$  is the population size.

There is not a common agreement on how this quantity might affect the amount of encoded information, but it exists a large literature where this question has been addressed from different perspectives. In the following I will present what in my opinion have been the most important or influential studies, explaining them in more or less detail depending on their complexity and/or relevance.

In 1999 Carrie McAdams and John Maunsell published a seminal study [62] where they showed that the tuning curves of monkey V4 neurons were modulated with attention. In particular they designed a match-to-sample task where monkeys had to report whether two consecutively presented Gabor patches had the same orientation or not. In some trials monkeys attended the Gabor patches while in the rest of trials they did not. For each neuron, a tuning curve was computed for the attended (Fig 1.4; black curve) and the unattended condition (gray curve). At the population level they found that the stronger effect on the neurons' tuning

curve was multiplication. A smaller component of tuning displacement was found as well. No changes were found on the width of the tuning. When the monkey attended the Gabor stimulus, the relationship between mean firing rate and orientation of the patch was enhanced (multiplication) with respect to the unattended condition. A multiplicative scaling of the tuning curve has also been reported when increasing the contrast of the presented stimuli [62]. Because of this similarity between contrast and attention on the neural code, the authors suggest that perhaps both mechanisms are related in terms of information processing.



**Figure 1.4:** Attention increases single neuron selectivity by scaling the tuning curve multiplicatively. The mean response (vertical axis) for each stimulus orientation (horizontal axis) is multiplicatively enhanced when the presented stimuli are attended. Extracted from [62].

Even though the content of this paper is highly related to our approach in chapter 2, it is not identical. First of all their definition of multiplicative scaling is highly related with our definition of global activity  $g$  (see eq. (1.30)) but not equivalent. A multiplicative scaling of the tuning curves of all neurons would imply an increase in global activity, but an increase of global activity does not imply a multiplicative scaling of the tuning curves. For instance, if tuning curves undergo an additive increase, global

activity would increase as well but no scaling would have occurred. Second, they are just characterizing a multiplicative scaling of neuronal tuning curves under attention but not showing how this multiplicative scaling could affect the amount of information encoded in a network. We know that information is not encoded at the single cell level but across the network. Not taking into account the pairwise covariability can lead to important biases or even wrong estimates of information encoded in a network (see previous section and [47]). Even though in the derivations of chapter 2 a very similar reasoning could be applied to the effect of a tuning multiplicative scaling in the amount of information encoded in a network, we will characterize explicitly the exact quantities of the neural code that play a role on information. Finally, in this study they do not show any effect of this multiplication on behavior. Even though it has been shown that attention enhances behavioral performance in general [58, 60] and they show that attention modulates multiplicatively neuronal tuning curves (and therefore information at the single cell level), they do not provide any clear link besides this indirect relationship. It is important to provide a direct mapping between a particular quantity of the neural code and behavior in order to claim for a causal relationship, even though a direct relationship could still arise from the connection to a third real causal quantity. In this thesis I will try to provide such a direct link and also I will present a method where we are able to disentangle what are the quantities that are directly affecting behavior from those that are just linked indirectly.

In 2006 Ma and colleagues published an interesting paper [63] that, even though a bit orthogonal to our claims, it is worth explaining it here briefly. The main claim of this study is that the joint activity of a population of neurons represents the posterior distribution of an encoded parameter (see section 1.1). They propose a neural implementation for such a representation called Probabilistic Population codes (PPC). During the encoding phase, the trial-by-trial variability can be understood as noise. In this paper it is proposed that noise could represent uncertainty about the inferential process in a natural way. Instead of being the noise a sub-product of a biological system harming the reliability of the neural code,

it is the way the brain encodes for uncertainty, or which is the same, the degree of belief for a particular probabilistic parameter estimate. The global activity of the network would determine how narrow the posterior distribution is, or in other words, the larger the global activity, the larger the information (reliability) of the neural code.

In 2014 Ecker and colleagues published another interesting study [43] where they tried to account for the difference in Fano factors and pairwise correlations between the awake and the anesthetized state of monkeys. Their main claim was that during anesthesia there was a common global modulation of the network that was playing a very important role on each neuron's firing rate. This common modulation would produce the experimentally observed larger Fano factors and pairwise correlations in anesthesia when compared to awake. In particular they modeled the firing rate of each neuron  $i$  as a function of time by the equation

$$r_i(t) = f_i(s(t)) + c_i x(t) + \xi_i, \quad (1.31)$$

where  $f_i(s(t))$  was the tuning curve for neuron  $i$ ,  $s(t)$  was the stimulus presented to the network at time  $t$ ,  $c_i$  was the coupling of each neuron to the time-varying global modulation variable  $x(t)$  and  $\xi_i$  accounted for the model's assumption of Gaussian additive independent noise. After fitting the model and subtracting the effect of the global variable  $x(t)$  on the anesthetized data, the population's obtained values for Fano factor and pairwise correlation became statistically equivalent to those values obtained on the awake dataset. In my opinion this study represents a great contribution to the scientific community because with a very simple explanation they were able to explain a rather controversial and unexplained experimental evidence. In this paper the authors do not mention what might be the role of this common global modulation in terms of information processing nor refer to possible implications on behavior. This is a key difference between this study and our approach in the chapter 2.

In a similar line of research in 2014 Goris et al. published a study [44] where they tried to characterize many of the experimentally found statistics of the neuronal activity by adding a gain factor before the Poisson step. In particular they proposed a model where in each trial and neuron the

spike count was obtained from a Poisson process

$$p(N|\mu, \Delta t) = \frac{(\mu\Delta t)^N}{N!} \exp(-\mu\Delta t), \quad (1.32)$$

where  $\mu = f(s)g$  is the combination of the tuning curve and a global multiplicative gain. If the trial-by-trial distribution of the multiplicative gain is a gamma distribution then eq. (1.32) becomes a negative binomial distribution. This way the supralinear relationship between the variance and the mean firing rate can be explained as well as some of the relationships between pairwise correlations and distance or tuning similarities. In this study the global modulation factor is equivalent to those described in [62]. It arises as a multiplicative scaling in the tuning curve of the neurons. A multiplicative scaling always produces an increase in global activity as defined by eq. (1.30), but an increase in global activity does not necessarily imply increase in multiplicative scaling in general. In this study it is not addressed the role of this multiplicative scaling on the amount of information processed by the neural network nor how it might ultimately affect behavior. As in some of the previous presented studies, in my opinion it is a very valuable scientific study but their approach is rather different to the one we took in chapter 2. The intersection with this thesis lies solely on the concept of a global modulation, which has been defined similarly but not equivalently as eq. (1.30) accounts for a general definition regardless of the exact mechanism underlying a change in global activity.

So far I have reviewed studies that in my opinion are important for this part of the introduction, however there are in particular two papers that deserve a deeper analysis due to their relevance with this thesis: [42] and [51].

The starting point for the study performed by Lin et al. 2015 [42] is about explaining trial-by-trial individual and shared variability in a neuronal population. To do so they propose a model for the neuronal activity that includes both an additive and a multiplicative term, the so-called affine model. The affine model would correspond to a mixture between the additive [43] and the multiplicative [44] model. In particular they

assumed that the expected spike count of neuron  $i$  on trial  $k$  was

$$E[n_{ik}] = g_k f_i(s_k) + a_i h_k, \quad (1.33)$$

where  $g_k$  is a common multiplicative gain,  $f_i(s)$  is the tuning curve for neuron  $i$ ,  $h_k$  is the common additive modulation and  $a_i$  is the coupling from each neuron to the additive term. Even though in the affine model each neuron is differently coupled to the additive modulation, the multiplicative gain is shared. The trial-by-trial variability was modeled using a negative binomial distribution. By this model they were able to explain most of the individual and shared trial-by-trial variability. Interestingly this model outperforms the pure multiplicative and pure additive described in studies [43] and [44] respectively. One of the most interesting parts of this study was that they further analyzed what was the role of the additive and multiplicative terms on the amount of information encoded in a network. To do so they generated surrogate population activity to simulate an orientation and contrast discrimination task. They first found that the larger the mean of the multiplicative gain the larger the information in the network, and the larger the mean of the additive term, the lower the information. Additionally they found that the information was very weakly affected by the variability of the additive term whereas variability on the multiplicative term had a major impact on the discrimination capabilities of the network. In my opinion this paper is very interesting in terms of providing a very good description of cortical variability by the affine model. With a natural extension of a very simple idea they outperform the pure multiplicative and pure additive models in terms of explanatory power. The encoding of information is very interesting but it should be extended and studied in more detail in order to make robust claims about the exact role of the additive and multiplicative terms on information encoding. As in the rest of studies herein presented, it lacks a direct test with behavioral performance.

The last study I wanted to review was published by Iñigo Arandia-Romero and colleagues in 2016 [51]. Its strength lies on the fact that their aim was not to characterize the statistical properties of the neuronal activity per se, but how the different types of global modulations could

affect the encoding of information, a much more interesting question in my opinion. The starting point for this study is the multi-gain model for the neuronal response

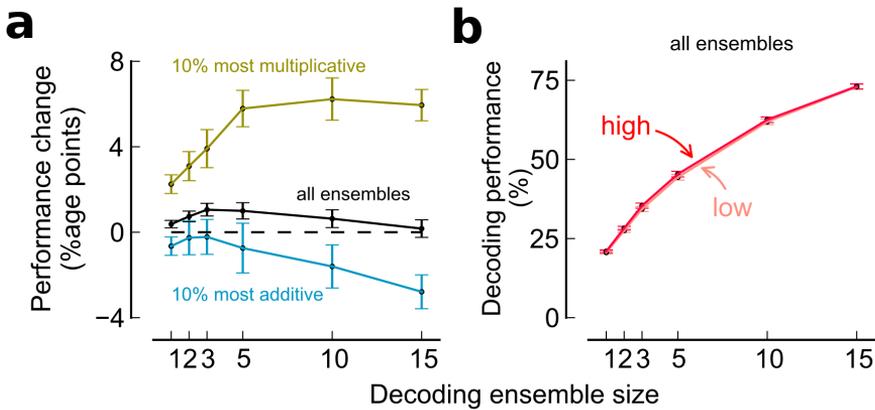
$$f_i(s, g) = (1 + \alpha_i g)h_i(s) + \beta_i g, \quad (1.34)$$

where  $f_i(s)$  is the experimentally-measured tuning curve of neuron  $i$ ,  $h_i(s)$  is the original tuning curve of neuron  $i$  before any global modulation,  $g$  is the global modulation term, and  $\alpha_i$  and  $\beta_i$  are the multiplicative and additive coupling of neuron  $i$  to the global fluctuation parameter. In this study the global parameter  $g$  was defined in the same manner as in equation (1.30). The amount of information (Fisher information) encoded by each neuron under this model is

$$I_i(s, g) = \frac{(1 + \alpha_i g)^2 h_i'^2(s)}{(1 + \alpha_i g)h_i(s) + \beta_i g}, \quad (1.35)$$

where  $h'(s)$  accounts for the derivative of the original tuning curve  $h(s)$  with respect to the stimulus  $s$ . From this expression it can be seen that neurons that tend to be multiplicatively coupled to the global gain will increase their information with larger values of  $g$ , whereas neurons that tend to be additively coupled will decrease their discriminability for larger values of  $g$ .

In Fig. 1.5 they show the effect of encoded information when the neuronal ensemble is constructed using multiplicative-like neurons and additive-like neurons. Following the intuition provided by eq. (1.35), those sets of trials with larger global activity produced larger and fewer encoded information on multiplicative and additive ensembles respectively. (Fig. 1.5a). Moreover, the amount of information encoded by the network as a whole did not change when comparing high vs low global activity of the population. I think this is a very interesting result: global activity in the population defined as eq. (1.30) has nothing to do with the amount of information encoded in a network, but it depends on whether global activity comes from a multiplicative of an additive modulation. In chapter 2 I will further confirm this results in different datasets, task and



**Figure 1.5:** The amount of information encoded by a network is independent from its global activity. **(a)** Changes in global activity produce and increase and a decrease in the amount of encoded information for multiplicative and additive subpopulations of neurons respectively. **(b)** Encoded information does not change as a function of global activity when considering the whole network. Extracted from [51].

brain regions and additionally I will link it with the behavioral performance of the decision-making agents.

It is important to mention here that it exists a fundamental difference between the set of presented studies and our approach to global activity as an internal state of the neural network (eq. (1.30)). In these studies they propose different modulatory mechanisms for the global activity that principally account for the single neuron and population activity statistics: an additive [43], a multiplicative [44, 62] and both an additive and a multiplicative common terms [42, 51]. Our approach does not consider the underlying mechanisms for such a modulation but just considers how the ensemble’s global activity might affect the amount of information encoded by the neural network. A narrower characterization of how the different mechanisms underlying global fluctuations might be affecting the encoding of information is scientifically very interesting but out of the scope of this chapter.

I also find important to clarify that most of the studies aiming to an-

swer what is the role of pairwise correlations on the amount of information encoded by a neural network base their analysis on approaches that are difficult to test experimentally. As stated before, Fisher information is defined as the variability on the estimate performed by an optimal agent about a continuous latent variable. Even though mathematically very convenient, it presents difficulties when implemented on experimentally-realistic electrophysiological recordings. The most important problem is based on the fact that on real experiments involving both behavior and electrophysiology the presented stimuli or experimental conditions constitute a discrete set of values. To accurately assess Fisher information one needs precise estimations of  $f'(s)$ , the derivative of the tuning curve with respect to the stimulus  $s$ . Under discrete sets of stimulation values such estimates tend to be inaccurate and consequently Fisher information estimates tend to be inaccurate as well. The natural extension of Fisher information to discrete sets of values is the percentage of correct classifications or decoding performance (DP). This metric accounts for the number of correctly classified patterns of activity over the total.

Another important feature of many of the studies presented above that is difficult to test on experimentally-realistic datasets is that they aim to study how correlations affect information for very large ensemble sizes. As recording techniques nowadays just allow to record tens or hundreds of neurons simultaneously, this approach is very challenging when generating direct testable predictions on real experiments. Even though exploring the large ensemble size regime is a fundamental step for characterizing the link between neuronal activity and behavior, it is also very important to derive theories that state what is the role of pairwise correlations for information metrics and ensemble sizes that can be tested experimentally. This is, in my opinion, one of the main contributions of chapter 2.

In chapter 2 I will present a set of theoretical and experimental results where we aimed to identify what are the most important features affecting the amount of information encoded by a neuronal ensemble. We started by deriving an analytical expression for the amount of encoded information as expressed by the decoding performance of a cross-validated linear

classifier on experimentally-realistic ensemble sizes. Then I will present a novel non-parametric method based on bootstrapping trials with replacement that allowed us to generate perturbations on different features of the neural code. By the combination of these two strategies we were able to validate consistently across four datasets involving different brain regions and tasks that neither mean pairwise correlations nor global activity of the network had any influence at all on the amount of encoded information by the neural network. Moreover, when assessed how perturbations on the different features of the neural code affected the behavior we found that selectivity length and projected precision were the most influential factors on the performance of the behaving animals. A simple neuronal model with experimentally-realistic covariance matrices, global networks dynamics and activity statistics was able to reproduce qualitatively all the findings reported by the experimental datasets.

### **1.3 Prior Information**

The field of decision-making comprises a vast bulk of literature, ranging from studies involving humans in economic tasks to models of spiking units resembling neurons during visual discrimination tasks. Even though I find all the different approaches for studying decision-making really interesting, in chapter 3 I will focus on studying the effect of prior information on the behavior of rats performing an outcome-coupled perceptual decision-making task and how this information is encoded in the orbitofrontal cortex (OFC). Therefore, the principal goal for this part of the introduction is to guide the reader towards the main question I want to address in chapter 3. I will do so by stating the general theoretical framework the study is embedded in and by reviewing the most important studies of the field. Unfortunately I will not be able to cover all the relevant studies but those I am going to present here are in my opinion very important contributions to the field, in particular to how prior information affects decision-making.

### 1.3.1 Behavioral protocols for studying prior

A widely accepted and fruitful framework when studying decision-making is understanding it is a process that consists on the integration of current sensory information with previous experience or prior information. Current sensory information, or the encoding process, is performed by the sensory neurons in the brain and its efficiency and reliability can potentially depend on many different quantities of the neural code (see section 1.2 and chapter 2). The term prior information encompasses a large set of different concepts but in general it can be understood as the set of regularities in the environment that have been acquired by of a decision-making agent that are used to produce faster and more accurate estimation of the state of the world than when only using sensory information. Importantly, the environmental regularities can be of many different natures and time scales. For example, it has been reported that baby humans possess important intuitions about very general physical laws of our environment which can only be explained as genetically-encoded prior information [64–66] (so to say long-term 'evolutionary' priors).

At the short time scale it is also clear that living beings make use of the temporal dependency of events in nature to anticipate upcoming situations and optimize their behavior. In this thesis I will focus on the short term prior information. In particular our goal is to study prior information that is linked to previous rewards and choices of the behaving agent. In my opinion this is one of the most interesting types of prior because the combination of recent choices and their associated outcomes is one of the key factors predicting upcoming events in nature, and therefore crucial for adaptative behavior. Keeping track of this information and using it optimally is a fundamental feature of living beings. At this point I think it is important to remark that the term prior information has a very clear and well defined conceptual and mathematical form in the Bayesian framework, yet it is not the only framework where we can talk about prior information. In suboptimal probabilistic frameworks and even in non-probabilistic frameworks we can still talk about prior information in the way that has been defined above.

The study of how decision-making agents commit to a particular option has been classically addressed using environments where the experimentalist can control most of the variables that can affect the choice of the decision-maker. Actually this methodology is not exclusive to cognitive decision-making studies but a fundamental idea of the scientific method for understanding nature. The common scheme of these experiments is the following: an agent is presented with two or more options and then asked to commit to a decision. By controlling the features of the presented options, a mapping can be inferred between the options presented to the observer and the decisions it makes. Even though these experiments are highly controlled, there will always be external and internal uncontrolled variables that will affect the decision. This uncontrolled variability can in principle be averaged out if the protocol is repeated many times, either by performing the same experiment on one subjects many times or by performing the experiment on many subjects few times.

One of the most important experimental protocols in cognitive neuroscience is the random dot motion task (RDM) [49, 50]. The general protocol of this task is the following: the subject is presented with a set of moving dots and it has to report what direction these dots are moving. Most of the times the task is binarized so the subject has to choose one of two possible options. As stated before, several repetitions of the stimulus will be presented to the observer so that the trial-by-trial variability can be averaged out and statistical inference can be performed on behavior and neuronal recordings. The difficulty of the task can be controlled in different ways, for instance by changing the exposure of the observer to the stimulus. The longer the presentation of the stimulus the easier the task. Another way of controlling the difficulty of the stimulus is by the coherence parameter. This parameter determines the percentage of coherently moving dots for the presented stimulus. When the coherence parameter is 100% all the dots will move coherently to one particular direction, whereas when the coherence is 0%, all the dots will move in a random direction. Primates reach high performances at around 5% - 10% coherence levels [49, 67]. This protocol and variations of it have been widely used to study how the middle temporal cortex (MT/V5) [49, 67],

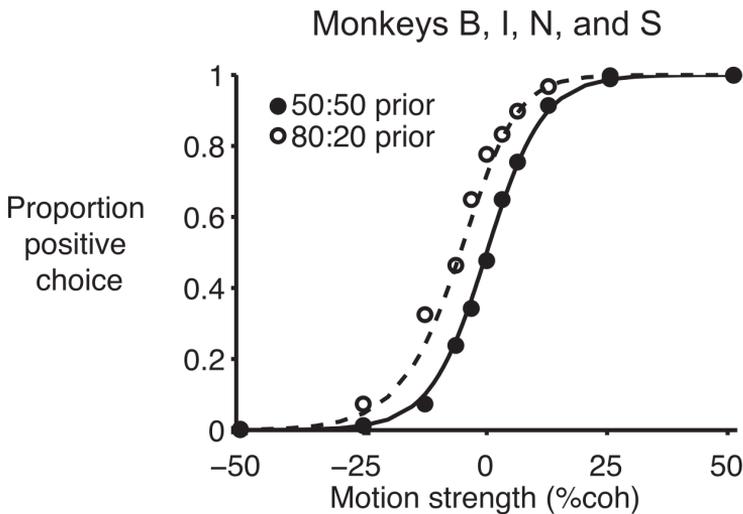
the lateral intra-parietal cortex (LIP) [68], or the lateral prefrontal cortex (LPFC) [69] are involved in perceptual decision-making as well as working memory, among others.

In most of the studies involving the RDM protocol, the stream of stimulus presentations has no statistical regularities: the percentage of the two possible presented motion directions is counterbalanced per experiment and independent from one trial to the next. This way the observer cannot make use of any additional source of information when performing the task, and therefore its decisions should have to be based purely on the sensory information presented. Because events in nature are rarely time independent and living beings have been optimized in such an environment, this is a rather unnatural condition and often signals of temporal dependency can be found in the stream of decisions performed by the decision-making agents.

Another very important experimental protocol used in cognitive neuroscience is the tactile frequency discrimination task [27]. This task, also known as the flutter discrimination task, consists mainly in the consecutive presentation of two vibration frequencies to a somatosensory receptor, generally a finger. The task of the decision-making agent is to report whether the second presented frequency is higher or lower than the first one. As most of the decision-making tasks it is based on instrumental conditioning, and therefore the agent is only rewarded when it makes a correct response. This task was designed to study perceptual decision-making on the somatosensory channel and how sensory information is represented in the somatosensory areas of the brain. As stated before, these classical decision-making paradigms were not designed to study how prior information shapes the stream of decisions and therefore the two conditions are counterbalanced. Additionally, it does not exist any temporal correlation between the stream of presented stimulus conditions and therefore, the optimal strategy for the decision-making agent is to report her percept based solely on the acquired sensory information. It exists a vast literature of classical protocols for perceptual decision-making tasks for most of the sensory channels, namely the auditory [25, 26], the olfactory [32, 70] or the gustatory system [71], besides the visual

and the somatosensory just explained.

Even though these classical decision-making protocols were not designed to study the effects of prior information on behavior, some years ago Hanks and colleagues [72] performed an experiment involving a variation of the RDM task to study prior information on both human and non-human primates. The protocol is equivalent to the classical RDM but they changed the prior probability for the two possible stimulus conditions (direction of motion right or left) to be presented. In particular, they presented 80% of the times one particular direction of motion. By embedding decision-making agents in this particular environment they were able to incorporate this source of information and bias their decisions accordingly. This is what we see in Fig 1.6, monkeys presented a systematic shift of their psychometric curve towards the more frequent direction of motion.



**Figure 1.6:** Perceptual responses on the RDM task experience a shift towards the most presented stimulus condition. The probability of making a positive choice (vertical axis) increases alongside the probability of presenting of a positive motion stimulus. Extracted from [72].

Their main goal however was to model how prior information was incorporated in a drift-diffusion model. Explaining the exact details of their model is beyond the scope of this thesis but the main idea is that they modeled prior information as a dynamic additive term in the decision variable that pushes the momentary evidence to the option with a larger prior. The alternative model, prior information incorporated as a static offset on the drift-diffusion model was ruled out by finding that the dynamical bias model could explain better the psychophysical data collected on two humans and four monkeys. Neurons in the LIP region of the behaving monkeys were also recorded and they report that the features of their firing rate suggest an encoding of the dynamical bias term. In my opinion this is a very important work in order to understand how prior information might be incorporated to the decision-making process. The way prior information is manipulated in this study is by breaking the symmetry on the percentage of stimuli presented to the subjects. However, there are many other possibilities in the way prior information can be manipulated. Departing from the 50%-50% condition is a first step in my opinion but there are richer and more ecologically-realistic options.

Another very interesting approach could be achieved by generating a serial dependency on the stream of presented stimuli while keeping the fraction of stimulus counterbalanced on each experiment. As the stream of natural stimuli any living being is exposed to generally presents a component of temporal auto-correlation, this represents in my opinion a better protocol for studying the effect of prior information on decision-making agents. This would be a very interesting experimental design because the probability of presenting a particular stimulus on the upcoming trial can be controlled by the experimentalist and therefore analyze how this might affect both the behavior and the neuronal activity of the decision-making agent.

Natural stimuli are in general temporally correlated and therefore having a sense of the external statistical regularities of the environment can be very advantageous in adaptive behavior. However, previous choices and their associated outcomes can also be of great importance on anticipating upcoming events or assisting the inference of ambiguous sensory

stimuli. Therefore, the set of recent choices and outcomes has to be understood as a form of prior information. This is the main scope of chapter 3 in the thesis: characterizing how the set of recent choices and outcomes can affect the behavior of decision-making agents as well as its neuronal representation.

I think that it is important here to mention that it has been found recently that the inactivation of some areas of the prefrontal cortex in rats can produce performance enhancements in a perceptual decision-making task [73]. Before the inactivation rats based their decisions on the combination of sensory information with the set of previous choices and rewards. This elicited suboptimal behavior because the experimental design of this task did not incorporate previous choices and rewards as valuable information sources. This is a very interesting and confirmatory result about the importance of previous choices and rewards as a form prior information because rats have been optimized in an environment where this sort of information is fundamental for guiding behavior optimally.

That said, before contextualizing in more detail chapter 3, I think it is important to analyze some additional studies that provide further evidence on the relevance of previous choices and rewards as a form of prior information as well as their neuronal representations.

In 2004 Dominic Barraclough and colleagues [74] published a very interesting paper where they analyzed the behavior and neuronal activity of two monkeys while they played a game analogous to the matching pennies against a computer. Describing the exact details of the paper is out of the scope of this thesis but the general idea is that monkeys were trained to choose one of two targets so that the machine could not predict their choice. There were three levels of 'machine intelligence': in the first one the machine selected choices in a random manner, in the second one it used the information of monkeys' previous choice to predict the upcoming response and in the last level the machine used both the history of monkeys' rewards and outcomes to predict their upcoming response. The more information the machine was using the more difficult for the monkey was to be unpredictable. While performing the task neurons in the prefrontal cortex were recorded and they found that some of them were

encoding for task relevant quantities like previous trial choice or previous trial outcome. This is a very interesting result, and even though these results were not interpreted in the prior information framework, they provide clear evidence that during the decision-making process monkeys can take previous choices and rewards into account to maximize performance and that this information is encoded by neurons in the prefrontal cortex. It is important to remark though that this result is not embedded into a perceptual decision-making framework and therefore it is difficult to interpret it from the prior information point of view.

Another study I found important to mention here was published in 2012 by Moonsang Seo and colleagues [75]. As in the previous study, the interpretation of their results was not embedded in the framework of prior information integrated with sensory information but I think that the behavioral protocol presented in this study is very close to the one described in chapter 3. The task consisted in reporting whether there was a larger fraction of red vs blue pixels in a presented circle on a screen. The larger the fraction difference, the easier was for the monkey to perform the task. The flow of responses the monkey had to make was not independent from one trial to the next but it had to follow a predetermined sequence of responses (out of 8 different sequences). The response of the monkey therefore could be based on integrating the perceptual information gathered by the visual system with the non-sensory information about the particular state in the sequence of responses the monkey was embedded.

Just from sensory information the task could be performed with a perfect performance, however for difficult trials the additional source of information could be very useful. Additionally, when the monkey made a mistake, it was brought back to the previous stage of the sequence which provided further assistance for the task. As mentioned above the behavioral results were not analyzed from the integration of prior with sensory information but they could have been. Moreover, the type of prior information provided to the monkey in this task is very similar to the prior information in our behavioral protocol, where previous choice and their associated outcome provide the decision making agent with all the necessary non-sensory information. In particular, based on the previous choice

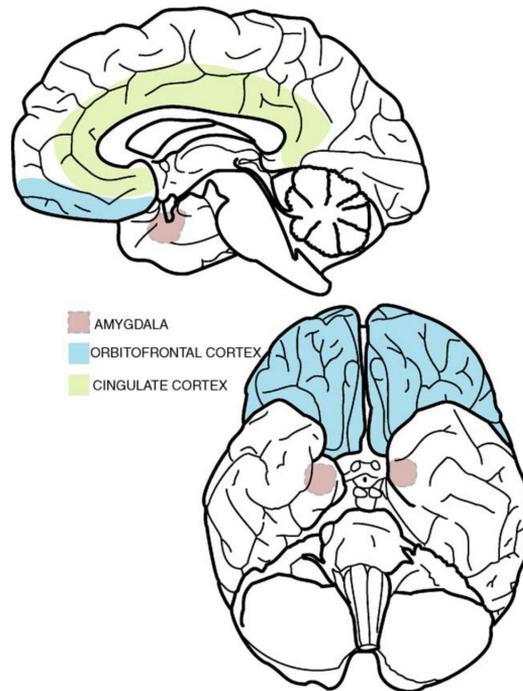
and outcome the monkey will be able to figure out in what response sequence he is embedded in so that it can assist its sensory percept. When it makes a mistake, as he is brought back to the previous stage in the response sequence, he can use this information to make a responses opposite to its previous mistake, regardless of the sensory information it has been provided.

### **1.3.2 The OFC in decision-making**

As stated earlier in this thesis, computational neuroscience lies on the assumption that all behavior arises from the state, dynamics and interactions of the whole set of neurons present in the nervous system. Even though it is still a topic of debate among scientists, it exists a large set of experimental evidence that shows functional localization in the brain [7]. Whereas primary sensory areas have been principally found in the occipital (visual) and temporal lobes (auditory and olfactory), some parts of the parietal and frontal lobes have been associated with motor planning and execution and with higher cognitive functions like decision-making, working-memory, attention, reasoning and language production. In particular, it has been found in some regions of the prefrontal and parietal cortex signals that are fundamental for the process of decision-making [67, 76–81].

The orbitofrontal cortex (OFC) is a region of the prefrontal cortex that has been traditionally thought to play a fundamental role in adaptive and goal directed behavior [82]. In primates it is located in the ventral surface of the frontal lobe (Fig. 1.7) and it can be defined as the region of the prefrontal cortex that receives projections from the magnocellular medial nucleus of the mediodorsal thalamus [83]. It receives inputs from all sensory modalities: gustatory, olfactory, somatosensory, auditory and visual [84].

The OFC has been shown experimentally to encode a myriad of decision-related variables, principally those involved in economic decision-making [85]. G. Schoenbaum and colleagues presented a study in 1998 [86] where they reported signals of expected outcome in rats' OFC and amyg-



**Figure 1.7:** The orbitofrontal cortex (OFC) is located in the ventral surface of the frontal lobe. Extracted from [84].

dala. Rats were trained to discriminate between two different odors. One of the odors was associated with a rewarding fluid and the other was associated with an aversive fluid. Neurons recorded in their OFC and amygdala showed a significant encoding of the nature of the fluid to be delivered during the presentation of the olfactory cues. In 1999 L. Tremblay and W. Schultz published another seminal study [87] where they showed that neurons in the macaques' OFC modulated their firing rates with respect to visual cues that predicted different types of rewards. In conjunction with the behavioral response to the task, the authors concluded that these set of neurons were encoding the upcoming rewards associated with each of the different visual cues. Some years later, in 2006, C. Padoa-Schioppa and J.A. Assad reported that neurons in macaques' OFC

encoded the value of presented and chosen goods, regardless of the visual and spatial traits of the visual stimuli and the motor action performed by the monkey [81]. Another important article regarding the functional role of OFC in decision-making was published in 2009 by S.W. Kennerley and colleagues [81]. In this study they recorded from the anterior cingulate cortex (ACC), the OFC and the lateral prefrontal cortex (LPFC) while monkeys performed a task that involved decision-making variables such as potential payoff, probability of success and cost in terms of time and effort. All three areas were found to encode the outcome values that were represented by each choice. Together with other experimental findings on humans [88–90], the traditional view on the functional role of the OFC is that it is responsible for the encoding of the expected values associated with the different choice options during decision-making [82, 85].

With the exception of some studies [91–94] the OFC has not been generally considered to be related with action selection and initiation. In the study performed by C. Feuerstein and colleagues [91] a group of rats was trained on an odor discrimination task while neurons were recorded in their medial and lateral orbitofrontal cortex. During stimulus sampling rats were presented with one of two odors in a central port. The identity of the odor was associated with a drop of water on one of the two lateral ports. They found that some neurons encoded stimulus identity while other neurons modulated their firing rate with respect to the direction of motion the rat was about to take. Due to their experimental finding on the encoding of choice-related signals in the OFC, in this study they claimed for a re-evaluation of the classical view where the OFC is exclusively associated with value-related variables in economic decision-making.

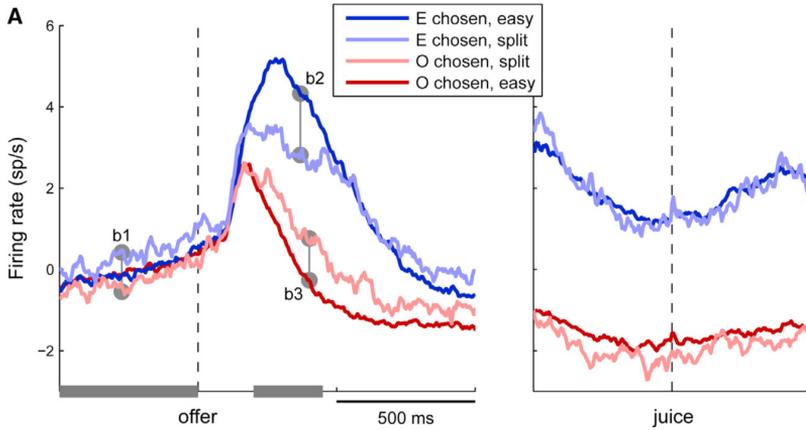
In that same year a study conducted at the G. Schoenbaum lab [92] found that neurons in the rats' OFC were encoding for the time delay before the delivery of a reward regardless of its absolute value. Some years later another study [93] found that neurons in rats' OFC modulated their firing rate with respect to the choice and the expected outcome on an odor discrimination task. The instrumental task was designed such that choice related signals could be distinguished from those related to the expected value on a trial-by-trial basis. On each trial a particular odor

was pseudorandomly chosen from a set four different odors: *a*, *b*, *c* or odor *d*. Odors *a-b* and odors *c-d* were associated with left and right response ports respectively and odors *a-c* and odors *b-d* were associated by small and large reward respectively. This way they could dissociate the encoding of both signals in the OFC of the behaving rats.

These three studies performed on rats reported heterodoxal views for the functional role of the OFC on decision-making. Even though at first glance this could indicate a different functional role between rodents and primates, in 2013 C. Padoa-Schioppa published a paper where it was described for the first time choice-related signals on monkeys' OFC even before the presentation of the different options to be chosen [94]. In this study he recorded from neurons in the OFC while monkeys performed an economic task where they had to choose between two different offered juices. He found that different populations of neurons encoded for the offered value, chosen value and chosen juice during the course of the trial. Interestingly, he found that a particular set of neurons in the OFC encoded information about the upcoming choice on easy trials before the presentation of the two offered juices (see Fig. 1.8), consistently with similar findings on LIP neurons and perceptual decision-making tasks [95, 96].

Even though the results reported in this study are very well aligned with those presented in [15] (see chapter 3) and provide evidence for a unified role of the OFC in action selection and initiation across species, in my opinion it includes some methodologies regarding the data analysis that might invalidate their main conclusions.

Recently a new hypothesis has been proposed for a unified functional role of the OFC in the process of decision-making [97]. They propose that the role of the OFC is to represent the current state within an abstract cognitive space of the task. This information would be used elsewhere in the brain to support learning and behavior, in particular, reinforcement learning (RL). Reinforcement Learning is an area of machine learning that models the behavior of a virtual agent in a particular environment under the constraint of maximizing a cumulative reward function [98]. It is out of the scope of this thesis to explain this theory in detail but the RL framework basically consists of a virtual behaving agent that changes its



**Figure 1.8:** Before the presentation of the different offers this neuron shows predictive activity for the upcoming choice of the monkey. The firing rate averaged across the set of trials where the monkey chose option E is larger than when the monkey chose option O before the presentation of the offer. The effect is only visible for difficult trials. Extracted from [94].

state  $s$  within an environment through actions  $a$ . The mapping between each state and the action to take is guided by the aim of maximizing a cumulative reward function, which can be expressed by a discounted sum of the present and future rewards associated with each state of this environment. Formally RL is defined by the terms:

- $S$ : set of states,  $s \in S$ .
- $A$ : set of actions,  $a \in A$ .
- $T(s'|s, a)$ : transition probability from state  $s$  to  $s'$  under action  $a$ .
- $R(s, a, s') \rightarrow r$ : reward obtained in  $s'$  when transitioning from  $s$  under action  $a$ .

The goal is to find the policy  $\pi(s) \rightarrow a$  that maximizes the future discounted sum of rewards, or in other words, to find the  $a^*$  that maximizes

$$Q(s, a) = E \left[ \sum_{t=0}^{\infty} \gamma^t r_t | s, a \right], \quad (1.36)$$

where  $0 < \gamma < 1$  controls for the relative importance of temporally distant rewards with respect to the close ones. In the RL framework two main approaches are generally considered: the model-based and the model-free RL. In the model-based it is assumed and explicit knowledge of the transition function  $T$ , the reward function  $R$  and the state and action spaces  $S$  and  $A$ . Equation (1.36) can be therefore expressed as

$$Q(s, a) = \sum_{s'} T(s', a, s) [R(s, a, s') + V(s')], \quad (1.37)$$

where

$$V(s') = \max_{a'} Q(s', a'), \quad (1.38)$$

which are known as the Bellman equations. The optimal policy  $\pi^*$  is then defined by the action  $a$  that maximizes eq. (1.37). In the model-free approach the behaving agent does not know explicitly  $T$  or  $R$ . What they do is to build approximations of eq. (1.36) as they interact with the environment. The total amount of expected reward can be written as the sum of the immediate reward and the future reward

$$\begin{aligned} V(s) &= E \left[ \sum_{t=0}^{\infty} \gamma^t r_t | s \right] = E \left[ r_0 + \sum_{t=1}^{\infty} \gamma^t r_t | s \right] = \\ &= E[r_0 + V(s') | s] = E[r_0 | s] + E[V(s') | s]. \end{aligned} \quad (1.39)$$

While interacting with the environment the agent samples the actions, states, transition probabilities and rewards to approximate its expectations. It will use the difference between its expectations and the reality ( $\delta_t$ ) to update them in an iterative manner

$$\delta_t = r_t + V_t(s') - V_t(s) \quad (1.40)$$

$$V_{t+1}(s) \leftarrow V_t(s) + \epsilon \delta_t, \quad (1.41)$$

where  $0 < \epsilon < 1$  is the learning rate. This is also known as temporal difference reinforcement learning (TDRL) and as stated above it is based on updating predictions from the interaction with the environment. While model-based RL can evaluate the optimal policy  $\pi^*$  due to the explicit knowledge of the full environment rewards and actions, model-free RL learns it implicitly by sampling the environment.

In the study by Wilson et al. 2014 [97], the authors proposed that OFC would be critical for representing the state  $s$  where the behaving agent would be located in this abstract cognitive space. This abstract space would consist on the combination of the multisensory information gathered by the behaving agent with the additional sources of information like previous experiences and working memories. They also propose a combination between model-free and model-based RL as the basis for behavior in animals. While OFC would be responsible for encoding the current location within the state space, ventral striatum (VS) and dorso-lateral striatum (DLS) would encode state and action values respectively for the model-free RL and model-based RL would occur on dorsomedial striatum (DMS) as well as VS. Their hypothesis is very well aligned with the set of results reported in [15] (see chapter 3), where neurons in the OFC were found to encode and integrate current sensory information with previous rewards and choices while rats were performing an outcome-coupled perceptual decision-making task.

In summary, in chapter 3 I will present a set of electrophysiological and behavioral results that aim to shed light on how prior information can be incorporated on a perceptual decision-making task as well as how this process is encoded in the OFC. Three rats were trained on an auditory time-interval categorization task where the previous choice and reward were predictive for the upcoming stimulus. We show that at the behavioral level rats do combine both sensory and prior information when committing to a decision on a trial-by-trial basis. We also found that populations of neurons in the rats' OFC encoded prior information in the form of previous choices and rewards as well as current stimulus and difficulty,

among others. Interestingly, OFC activity was predictive of rats' upcoming choice even before stimulus presentation. The encoding time profiles for the different task variables suggest that the OFC could be responsible for integrating current sensory information with prior information to guide behavior on a trial-by-trial basis.

My aim in this part of the introduction was to shortly review what have been the most important studies addressing the question of how prior information can affect perceptual decision-making as well as how this information can be represented in the brain. I also found convenient to present what is the general view on the functional role of OFC in decision-making by providing a set of brief descriptions of the most relevant studies of the field.



## Chapter 2

# ENCODING OF INFORMATION AND ITS LINK WITH BEHAVIOR

*The study presented in this chapter corresponds to the article in progress Nogueira, R. et al. Dissecting the most influential features of the neural code on information encoding and behavior.*

**Identifying the features of the neuronal code that are involved in the encoding of information is crucial to characterize the link between neuronal activity and behavior. Here, we dissect the neuronal code and find that only two factors, selectivity length and projected precision, affect the amount of information encoded in neuronal populations. If these two factors are controlled for, other features such as mean pairwise correlations or global activity do not contribute to information across four datasets involving different brain areas and tasks. Further, we show that selectivity length and projected precision are the most influential factors on behavioral performance, while mean pairwise correlations or global activity do not have consistent**

**effects. A biologically-constrained model of the cortex that follows close-to-optimal readout of information to guide choices is consistent with our results**

## **2.1 Introduction**

Understanding the neuronal code means understanding what are the statistical features of neuronal activity that affect the encoding of information as well as their link with behavior [99, 100]. Neurons are tuned to a large set of external and internal variables by eliciting different responses for different values of the encoded variable [24, 25, 49, 70, 101–104]. However, neuronal responses are typically largely variable both on the number and timing of spikes when presented with identical stimuli from one trial to the next [33]. Additionally, this neuronal variability is not independent but it is correlated across neurons, the so-called ‘noise correlations’ [34, 35, 48]. Noise correlations, even if they are extremely weak, can dramatically reduce the precision of the network’s encoding capabilities [47], and therefore have a major impact on behavior.

Experimental evidence supports the hypothesis that noise correlations can harm the neural code, as removing them produces a significant increase on the amount of information that can be extracted from neuronal ensembles [105]. However, a similar study also reported apparently contradictory results, where information was weakly improved when considering shared trial-by-trial variability [106]. It has also been shown that attention can have a positive impact on a visual discrimination task by reducing the overall noise correlations in the visual cortex [58, 60] or by selectively increasing or reducing pairwise correlations depending on the tuning similarity of the cortical neurons [59]. Perceptual learning has also been shown to be linked to a reduction of mean pairwise correlations [107]. These results suggest that mean pairwise correlations play a role on behavioral performance.

Global modulations of activity induced by attentional or motivational changes can also affect information and behavior due to modulations of

the tuning curves [62, 108]. Indeed, in the form of a multiplicative gain modulating neuronal tuning, global modulations have been hypothesized as the mechanism by which information processing abilities of the network increase under attention [62]. Also, gain modulations affect overall mean pairwise correlations [43, 44], thus implying its potential role on information. The fact that there is a relationship between global fluctuations and the distribution of information in neuronal ensembles [51], also suggests that global modulations of activity indeed play a role in information processing.

However, the reported dependencies of tuning and mean pairwise correlations with information, behavior, attention and perceptual learning, could be the result of third, untested variables, that explain those modulations and that they are themselves correlated with tuning and mean pairwise correlations. To aid at a coherent picture of the reported modulations, we aim at uncovering what are those hidden statistical variables of the neuronal code that more directly affect information and behavior. Using recent developments of encoding theory [47], we develop an analytical expression for the decoding performance of optimal linear decoders, which provides very tight predictions on the performance of trained optimal decoders in four different datasets spanning two monkey areas and three different tasks. From this analytical expression, we identify the two main statistical features of the neuronal code that affect decoding performance, called here ‘selectivity length’ (SL) and ‘projected precision’ (PP). We designed a perturbation technique that allows creating virtual instances of a particular experiment so that we could test and confirm the effect of SL and PP on information. Other statistical features of the neuronal code, such as mean pairwise correlations (MPC) or global activity gain (GA) did not play any role on the amount of encoded information by the neuronal population when SL and PP were controlled for. Further, we found that SL and PP, but not or inconsistently so MPC and GA, were modulated with changes in behavioral performance, suggesting that behavior was generated by an optimal or close-to-optimal readout of the population activity. Finally, we built a biologically-constrained model of cortex that follows optimal readout of information to guide choices and

that shows qualitatively the SL and PP dependencies reported in the data. Overall, our results consistently show that the same features of the neural code that are relevant for information are relevant for behavior.

## 2.2 Results

### 2.2.1 Tuning and noise features of the neural code affecting encoding of information

Identifying what are the most important features of the neuronal code on the amount of information encoded by a neural network is crucial for understanding the link between neuronal activity and behavior. The amount of information encoded by a neural network can be characterized by the percentage of correctly classified patterns of activity performed by a linear readout neuron downstream in the information processing pathway. This quantity is known as the decoding performance (DP) of the linear classifier, and its use has been widely extended recently due to its suitability as an information metric when evaluated on real experiments involving behavioral protocols and recordings up to tens of neurons [15, 51, 100, 105]. In Fig. 2.1a it is shown the trial-by-trial joint set of activities  $\mathbf{r} = (r_1, r_2)$  of an example network of two-neurons when presented with stimulus 1 ( $s_1$ ; green dots) and stimulus 2 ( $s_2$ ; blue dots). A readout neuron aim is to perform a weighted integration of the ensemble's activity  $\hat{s} = \boldsymbol{\omega}^T \mathbf{r} + \omega_0$  and compare it to a reference threshold to decide whether  $s_1$  or  $s_2$  was presented to the network in a trial-by-trial basis. We derived an analytical expression for the percentage of correct classifications performed by a linear read-out neuron of an arbitrarily large neuronal ensemble on a binary task (see section 2.4)

$$DP = \Phi \left( \frac{1}{2} \frac{\boldsymbol{\omega}^T \Delta \mathbf{f}}{\sqrt{\boldsymbol{\omega}^T \Sigma \boldsymbol{\omega}}} \right) \quad (2.1)$$

where  $\Phi()$  is the zero-mean and unit-variance cumulative Gaussian and  $\Sigma$  is the trial-by-trial covariance matrix among neurons in the population.

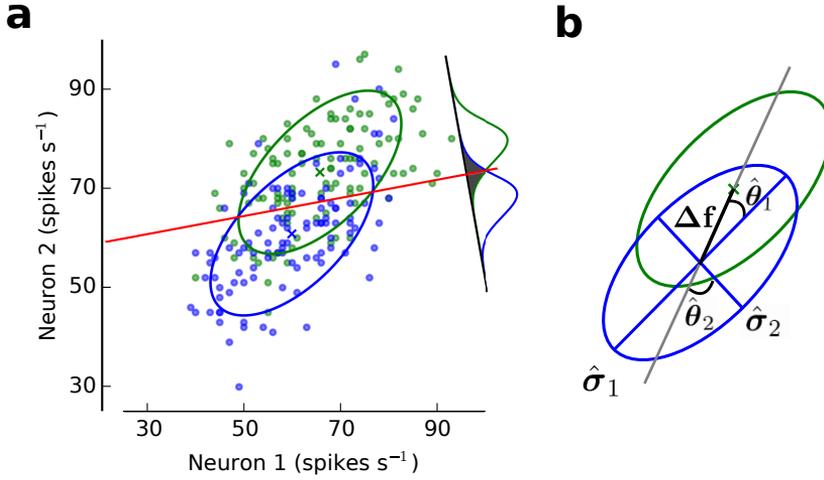
The term  $\Delta\mathbf{f}$  is the selectivity vector and it represents the tuning of the network to the stimulus  $s$  and  $\omega$  is the set of read-out integration weights used by the downstream neuron. Based on the tuning and noise properties of the neuronal ensemble, a particular set of read-out weights that maximize the performance of the read-out neuron can be derived (optimal classifier) [109]. When the read-out is optimal the above expression becomes

$$DP = \Phi \left( \frac{1}{2} |\Delta\mathbf{f}| \sqrt{\sum_{i=1}^N \frac{\cos^2 \hat{\theta}_i}{\hat{\sigma}_i^2}} \right) \quad (2.2)$$

The first term  $|\Delta\mathbf{f}|$  is the selectivity length (SL) and it represents the norm of the selectivity vector  $\Delta\mathbf{f}$  (Fig. 2.1b). The second term  $\sqrt{\sum_{i=1}^N \frac{\cos^2 \hat{\theta}_i}{\hat{\sigma}_i^2}}$  is the projected precision (PP) and it is defined as the projection of the inverse covariance ellipsoid (precision matrix) on the stimulus axis  $\mathbf{u}_{\Delta\mathbf{f}}$ , the direction of the vector  $\Delta\mathbf{f}$  (Fig. 2.1b). The amount of encoded information by a neuronal ensemble can be therefore fully characterized by two independent factors of the neural code: the magnitude of the population’s tuning to the stimulus  $s$  (SL) and the relative orientation of the covariance matrix with respect to the stimulus axis (PP). Is it important to note that the term inside the cumulative Gaussian in Eq. (2.2) is equivalent to  $d' = \sqrt{\Delta\mathbf{f}\Sigma^{-1}\Delta\mathbf{f}} = SL \times PP$  [110], however by rotating the original reference framework when evaluating this term, a novel and insightful view is obtained as we can fully detach the contributions from the magnitude of the ensemble’s tuning and the trial-by-trial variability on the relationship between the neural code and information.

## 2.2.2 Analytical Decoding Performance on *in vivo* recordings

We tested our analytical expression for the encoded information (Eq. (2.2)) on four different datasets consisting on simultaneously recorded units in monkeys (2 to  $\sim 50$  neurons) from different brain areas and tasks:



**Figure 2.1: Information encoded by a neuronal ensemble can be fully characterized by the selectivity length and the projected precision.** (a) The trial-by-trial joint activity of a network consisting on  $N$  neurons can be characterized by an ellipsoid embedded in an  $N$ -dimensional space. The covariance matrix and mean activity of the population will determine the shape and location of the ellipsoid on the neural space for stimulus 1 ( $s_1$ ; green) and or stimulus 2 ( $s_2$ ; blue). Information extracted by a downstream linear read-out neuron will be potentially different from chance (50%) if the joint activity corresponding to  $s_1$  and  $s_2$  is linearly separable on a one-dimensional space defined by the integration weights  $\omega$  (red line) (b) Information depends only on the selectivity length (SL) and the projected precision (PP) (see section 2.4). The norm of the selectivity vector  $\Delta \mathbf{f}$  (SL), corresponds to the distance between the mean activity of the population when presented with  $s_1$  and with  $s_2$  (distance between the centers of the green and blue ellipsoids). The PP is calculated from the angle ( $\hat{\theta}_i$ ) between each eigenvector of the covariance matrix and the selectivity vector ( $\Delta \mathbf{f}$ ) as well as from their eigenvalues ( $\hat{\sigma}_i^2$ ; length of each axis). Larger amounts of information will be encoded when the longest axis of the ellipsoid are orthogonal to  $\Delta \mathbf{f}$ .

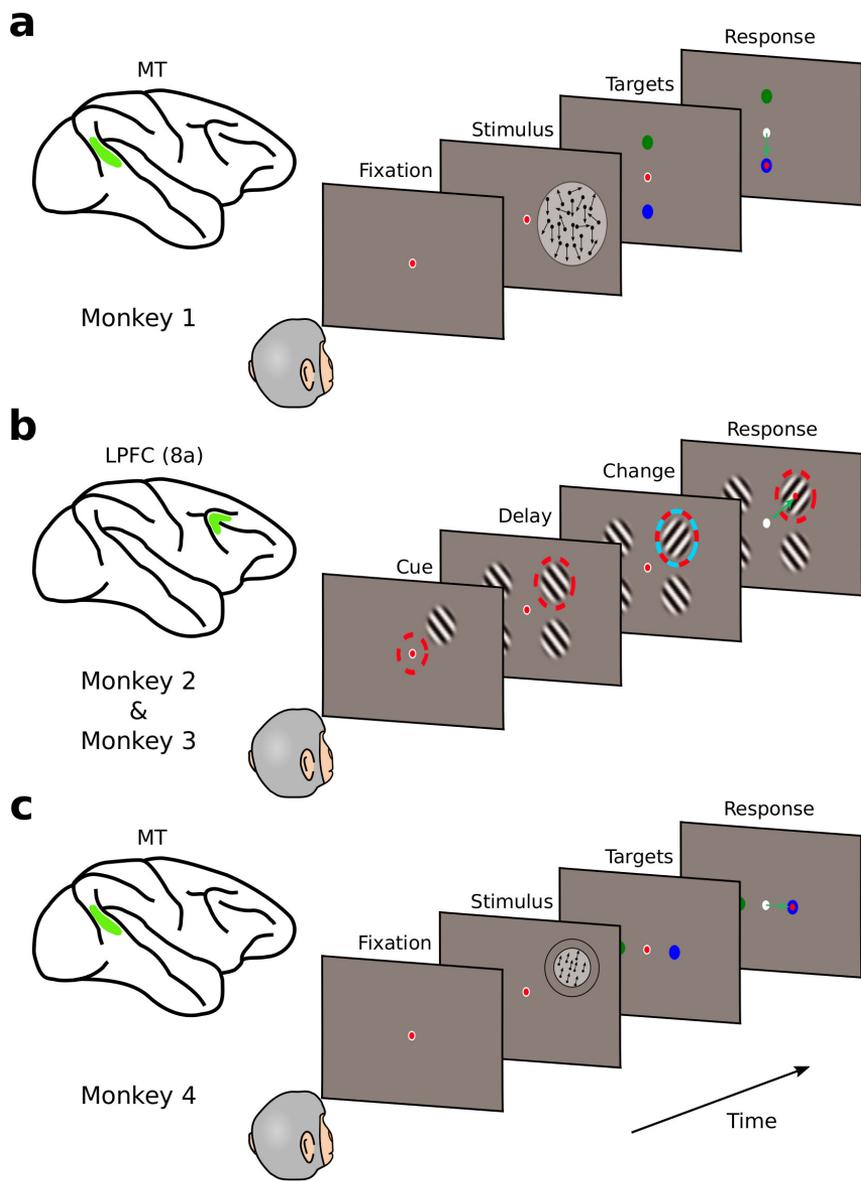
middle temporal (MT) neurons while monkeys performed a coarse motion discrimination task (monkey 1) [35] (Fig. 2.2a); lateral prefrontal cortex (LPFC, area 8A) neurons while monkeys performed an attentional

task (monkey 2 and monkey 3) [105] (Fig. 2.2b); and MT neurons while monkeys performed a fine motion discrimination task (monkey 4) (Fig. 2.2c) (see section 2.4 for a detailed description of the datasets and the classifier’s task for each monkey).

In Fig. 2.3a four examples corresponding to the four different datasets used in this study are shown. For each panel, on the horizontal axis it is plotted the DP of a Linear Discriminant Analysis (LDA) evaluated on a hold out fraction of trials (test set) after training the optimal set of read-out weights ( $\omega, \omega_0$ ) on the rest of trials ( $DP_{cv}$ ) (train set; 5-fold cross-validation). On the vertical axis it is plotted the analytical DP (Eq. (2.2)) using the whole dataset ( $DP_{th}$ ). For each panel a different ensemble size is used for illustration purposes and each dot represents a particular neuronal ensemble of the specified size (see section 2.4).

If Eq. (2.2) is a good approximation for the real amount of information encoded by the neuronal ensemble, it should be able to explain a large fraction of the whole variance depicted by the DP of a cross-validated linear classifier ( $DP_{cv}$ ). We performed a linear fit of the analytical  $DP_{th}$  against the cross-validated  $DP_{cv}$  for all ensemble sizes and monkeys (see section 2.4; Fig. 2.3b). For all datasets and ensemble sizes, the percentage of variance on  $DP_{cv}$  that was explained by  $DP_{th}$  was, in all cases, above 96.4%, (mean across monkeys and ensemble sizes =  $97.8\% \pm 0.6\%$ ). In Fig. 2.4 summary statistics for the goodness-of-fit (% explained variance), slope and intercept on the test and train set for different linear classifiers are provided. Even though Eq. (2.2) was obtained by assuming stimulus-independent covariances and Gaussian noise on the population activity, these results can be thought as a first experimental proof that the amount of information linearly encoded by a neuronal ensemble can be fully characterized by the SL and the PP. In the following sections we will present a novel technique based on bootstrap that will allow us to further confirm our hypothesis and ultimately test it on the performance of behaving animals.

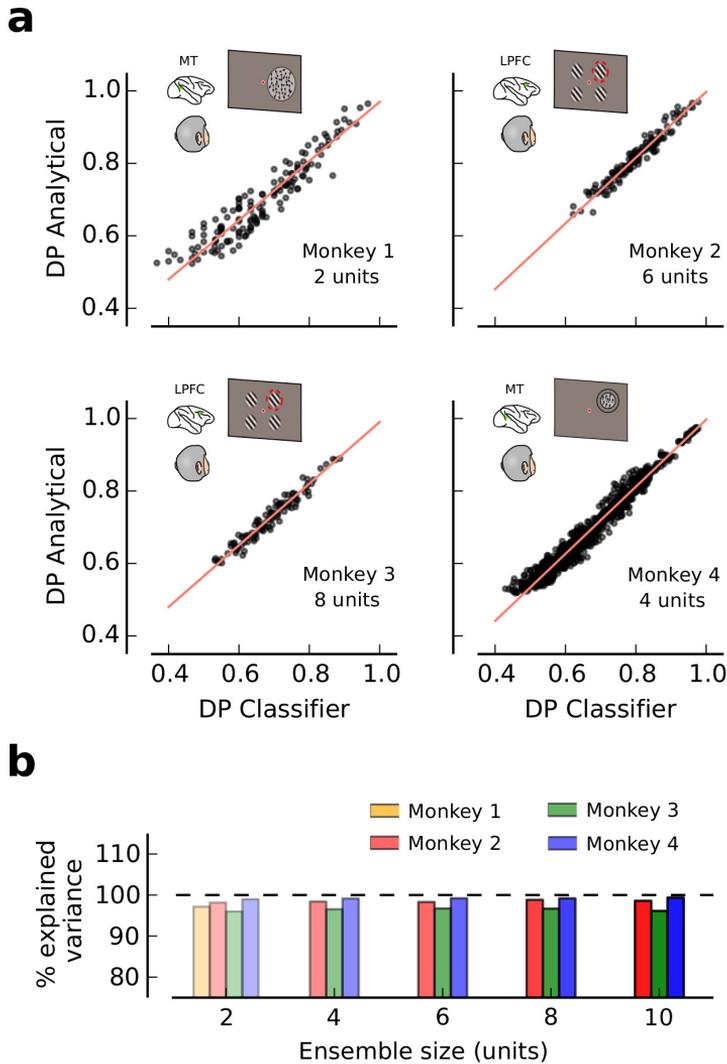
In order to rule out the possibility that the results on Figs. 2.3 and 2.4 could be obtained when using a naïve classifier, we performed the same analysis but using two suboptimal classifiers instead: the correlation-



**Figure 2.2: Four datasets involving different brain areas and tasks were used in this study. (a) One monkey performed a coarse motion discrimination task (monkey 1) while pairs of units were recorded in middle temporal cortex (MT) [35].**

**Figure 2.2:** (cont.) After stimulus presentation (random dot kinematogram) the monkey had to report the direction of motion by a saccadic movement to one of the two targets. Difficulty was controlled by the percentage of coherently moving dots in the stimulus. **(b)** Two monkeys performed an attentional task (monkeys 2 and 3) while units were simultaneously recorded ( $\sim 50$ ) in the lateral prefrontal cortex (LPFC, area 8A) [105]. Four Gabor filters were presented on the screen and the task was to make a fast saccadic movement to the attended location after a change in orientation of the cued filter was produced. **(c)** One monkey performed a fine motion discrimination task (monkey 4) while units were simultaneously recorded ( $\sim 25$ ) in middle temporal cortex (MT). After stimulus presentation (100% coherent random dot kinematogram) the monkey had to report whether dots were moving upward-leftwards or upward-rightwards by a saccadic movement to one of two targets. Difficulty was controlled by the angle defined by the direction of motion of the moving dots with respect to the vertical (see section 2.4).

blind classifier [111] and the variability-blind classifier (see section 2.4). We plotted (Fig. 2.5a) the ratio between the percentage of variance explained by Eq. (2.2) over the percentage of variance explained by the equivalent expression when using the correlation-blind and the variability-blind classifier (top and bottom panel respectively; see section 2.4). For all datasets and ensemble sizes, Eq. (2.2) has a larger explanatory power than both the correlation-blind and the variability-blind classifiers. Interestingly, when performed the same analysis after destroying the correlated structure of the dataset (shuffling method; see [61, 105]), the analytical expression for the correlation-blind classifier outperformed the optimal (before shuffling) and the variability-blind classifier as the best approximation for the amount of encoded information (Fig. 2.5b). Despite the set of assumptions needed for mathematical tractability, these results provide further robustness the SL and PP being the only factors affecting the encoding of information (Eq. (2.2); see section 2.4) as the reported goodness-of-fits cannot be achieved when using classifiers that do not accurately represent the correlated structure of the neuronal population. It is important to remark at this point that the amount of encoded infor-



**Figure 2.3: The analytical expression provides very tight predictions for the actual amount of encoded information.** (a) For all datasets the analytical expression for the decoding performance of a linear classifier ( $DP_{th}$ ; vertical axis) was a very good approximation for the real amount of encoded information as revealed by the decoding performance of a cross-validated linear classifier ( $DP_{cv}$ ; horizontal axis). Each panel corresponds to a different dataset. For illustrative purposes, a different ensemble size is used from each dataset for both the  $DP_{cv}$  and the  $DP_{th}$ .

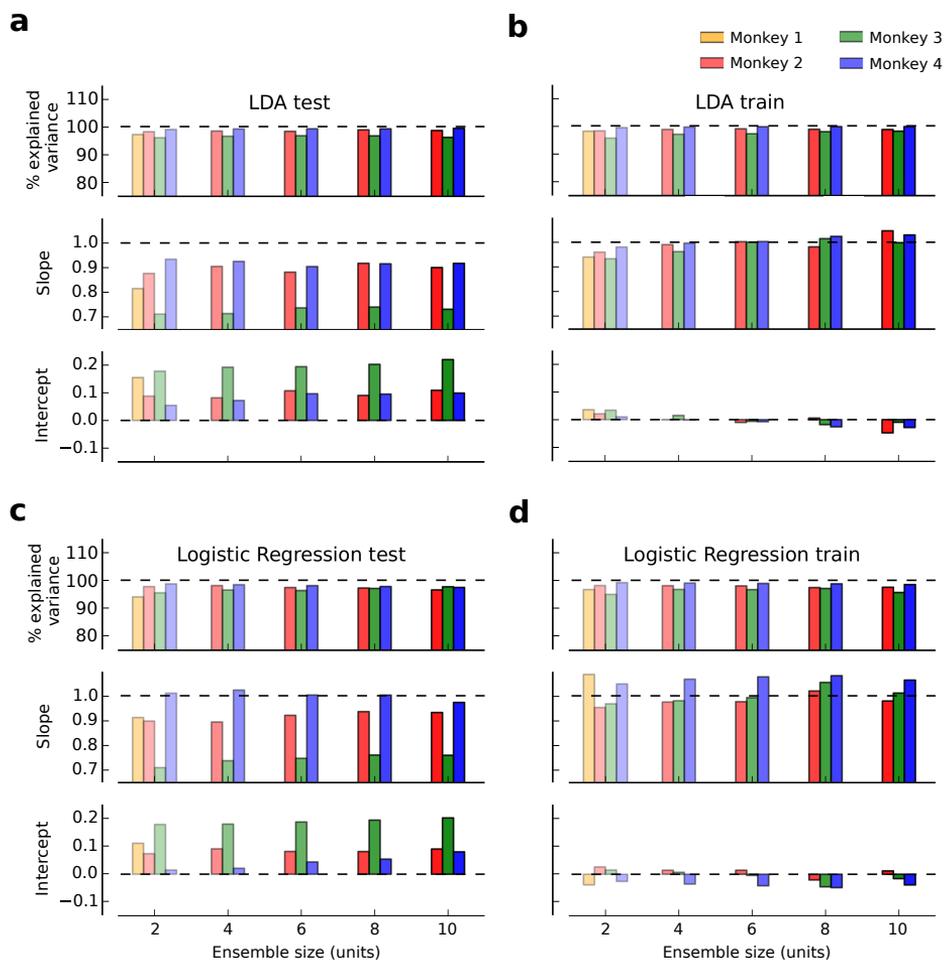
**Figure 2.3:** (cont.) **(b)** Summary statistics of the analytical approximation. In all datasets and ensemble sizes the percentage of variance explained by the analytical expression was a least 96.4% of the total variance depicted by the cross-validated linear classifier (mean across monkeys and ensemble sizes =  $97.8\% \pm 0.6\%$ ).

mation  $DP_{cv}$  has been characterized throughout this study by the performance of a cross-validated LDA because it showed depicts the largest DP when compared to the logistic regression (LR; linear classifier) and to the quadratic discriminant analysis (QDA; non-linear classifier) (Fig. 2.6).

### 2.2.3 Perturbing the original dataset by the bootstrap method

To further confirm that the only features of the neural code affecting the amount of information encoded by a neural network were the SL and the PP we used a novel technique based on bootstrap. For each dataset we were able to produce virtual instances of the same experiment by sampling trials with replacement from based on the original dataset. By using this method we were able to generate distributions of the different features of the neural code and therefore perturbing them and evaluating their effects on the amount of encoded information on the network and ultimately on monkeys' behavioral performance.

From the original dataset we first calculated the amount of encoded information as expressed by the cross-validated DP of an LDA ( $DP_{cv}$ ), the selectivity length (SL), the projected precision (PP), the mean pairwise correlations (MPC), the global activity of the network (GA) and the monkeys' behavioral performance (B) (see section 2.4). Then, by sub-sampling trials with replacement from the original dataset, we were able to generate a distribution for the amount of encoded information ( $DP_{cv}$ ) and all other of quantities (Fig. 2.7a). Each element of the distribution can be therefore thought as a perturbation with respect to the original dataset's value for the amount encoded information.



**Figure 2.4: The analytical expression provides excellent fits also when using different linear classifiers and goodness-of-fit metrics** (a) The slope and intercept parameter of the linear fit between the analytical and the cross-validated decoding performance (DP) are close to 1.0 and 0.0 respectively when using a Linear Discriminant Analysis (LDA) and the test set on all ensemble sizes and datasets. Deviations from 1.0 (slope) and 0.0 (intercept) can be explained by the limited number of trials used to train and test the LDA which produce in some cases  $DP < 0.5$ . (b) All goodness-of-fit metrics are enhanced on all datasets and ensemble sizes when using the train set for testing the performance of the LDA.

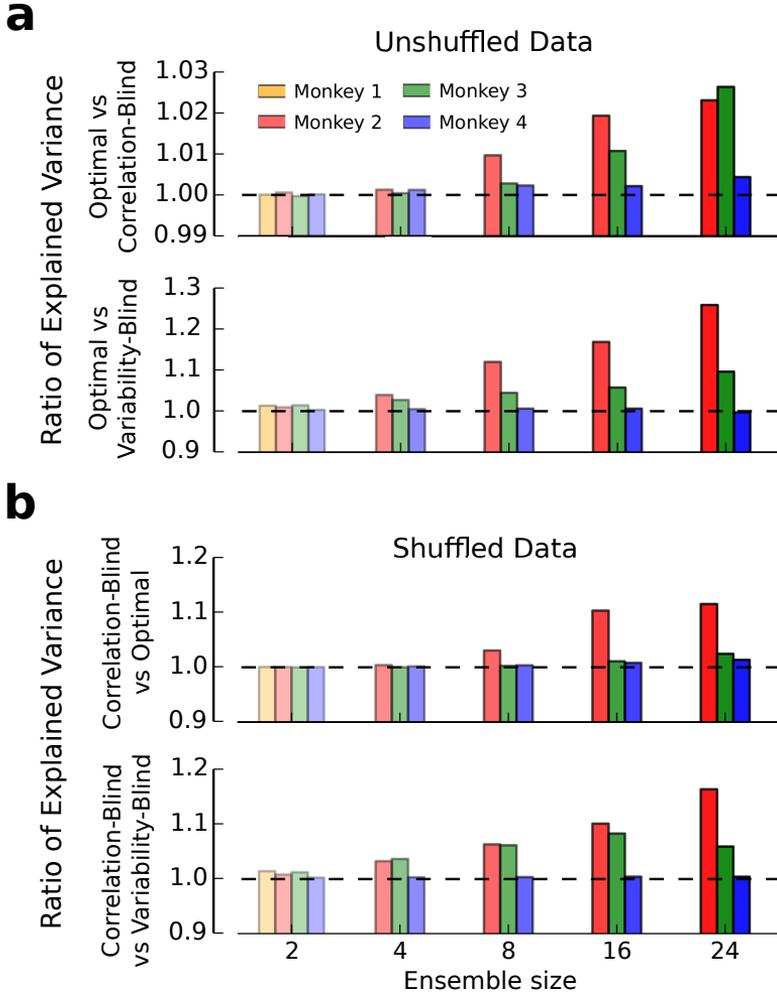
**Figure 2.4:** (cont.) This enhancement is produced because evaluating a model on the train set will in general produce better results than the same model evaluated on a hold-out partition of the dataset (test set) (c) When using a Logistic Regression (LR) instead of an LDA as the amount of encoded information by the neural network, the fitting results are qualitatively equivalent to panel (a). (d) As when comparing panel (a) and (b), when using the train set to evaluate the performance of the LR, all goodness-of-fit metrics get enhanced with respect to panel (c).

Because different features of the neural code might be correlated, perturbations on these different features might not be independent. To avoid spurious results we used a conditioning strategy where, when studying how each feature affected information (and behavior; see below), we only used bootstrap iterations that produced perturbations of this particular feature while fixing the rest. For instance, when studying the effects of MPC on the amount of information encoded by the neural network ( $DP_{cv}$ ), we only used bootstrap iterations for which both SL and PP were roughly fixed (see section 2.4; Fig. 2.7b).

## 2.2.4 Encoded information is not affected by mean pairwise correlation or global activity

By perturbing the different features of the neural code (bootstrap and conditioning method; see Fig. 2.7 and section 2.4) we were able to identify what were the quantities that mostly affected the amount of information encoded in a neural network. In particular we show how perturbations in SL (red), PP (green), MPC (blue) and GA (orange) produced a change in the amount of encoded information ( $DP_{cv}$ ) (Fig. 2.8) for all the different datasets used in this study and ensemble sizes (see Fig. 2.2 and section 2.4).

In the vertical axis it is plotted the percentage change in  $DP_{cv}$  when evaluated on those bootstrap iterations that elicited positive perturbations with respect to those that elicited negative perturbations (and keeping the

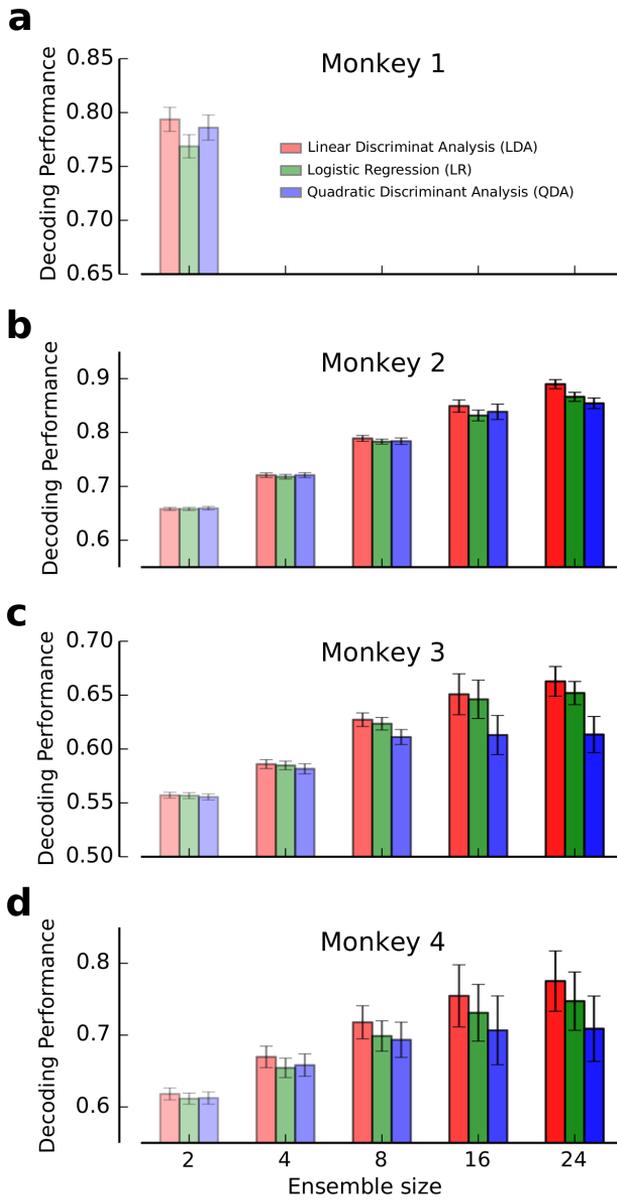


**Figure 2.5: The optimal and the correlation-blind are the best expressions for the unshuffled and the shuffled datasets respectively.** (a) Ratio between the goodness-of-fit corresponding to the optimal classifier and the correlation-blind classifier (top panel) and the variability-blind classifier (bottom panel). The goodness-of-fit is assessed as the percentage of variance depicted by the cross-validated linear classifier ( $DP_{cv}$ ) that can be explained by the analytical expression (see Fig. 2.3; and section 2.4). In all monkeys and ensemble sizes (2, 4, 8, 16 and 24 units) the analytical expression for the optimal classifier is the best approximation for the amount of information encoded by the neural population ( $DP_{cv}$ ).

**Figure 2.5:** (cont.) **(b)** Ratio between the goodness-of-fit corresponding to the correlation-blind classifier and the optimal classifier (top panel) and the variability-blind classifier (bottom panel). The amount of encoded information ( $DP_{cv}^{sh}$ ) is better approximated by the correlation-blind classifier than the optimal classifier (before shuffling) and the variability-blind classifier when pairwise correlations are removed (see section 2.4 for the correlation-blind and variability-blind analytical expression).

rest of features unperturbed, see section 2.4) and in the horizontal axis it has been plotted the different ensemble sizes used in this study (2, 4, 6, 8, and 10 units). For monkey 1 (Fig. 2.8a) perturbations on SL and PP produced a significant positive change in encoded information of  $9.45\% \pm 0.68\%$  (Wilcoxon signed-rank test,  $P = 4.2 \times 10^{-37}$ ) and  $5.56\% \pm 0.43\%$  ( $P = 3.2 \times 10^{-38}$ ) respectively. On the contrary, perturbations on mean pairwise correlations and global activity produced a change on  $DP_{cv}$  of  $0.24\% \pm 0.19\%$  and  $-0.20\% \pm 0.19\%$ , and they were not significantly different from zero ( $P = 0.38$  and  $P = 0.13$  respectively). For monkeys 2, 3 and 4 (Figs. 2.8b,c,d) the obtained results were qualitatively equivalent to those obtained in Fig. 2.8a. Consistently with Eq. (2.2) the amount of information is fully determined by the selectivity length and the projected precision of the neural code, while neither mean pairwise correlations nor global activity of the network play any role on the amount of encoded information by a neuronal ensemble. Another useful way of understanding these results is that on average, on those trials where SL and PP were larger by chance, the amount of encoded information was also larger.

It exists a widely accepted approach when characterizing the role of pairwise correlations on the amount of information encoded by a neural ensemble that consists on shuffling across trials the activity depicted by each neuron for a fixed stimulus condition [51, 61, 105, 112]. Eq. (2.2) allows also predicting beforehand the effect of shuffling across trials (removing pairwise correlations) on the neural code by just comparing the projected precision of the original dataset with the equivalent term calculated after removing the off-diagonal terms in the covariance matrix (see



**Figure 2.6: Linear discriminant analysis outperforms logistic regression and quadratic discriminant analysis.** (a) Mean decoding performance (DP) for a linear discriminant analysis (LDA), a logistic regression (LR) and a quadratic discriminant analysis (QDA) (5-fold cross-validation) for an ensemble size of 2 units.

**Figure 2.6:** (cont.) The LDA shows the largest mean DP across files and coherences. **(b)-(c)** Equivalent to **(a)** for monkeys 2 and 3 and for larger ensemble sizes (2, 4, 8, 16 and 24 units). As in **(a)**, the LDA depicts the largest mean DP across files and independent subgroups of neurons for all ensemble sizes. **(d)** Equivalent to **(a-c)** for monkey 4. As in **(a-c)**, the LDA depicts the largest mean DP across angle of motion, files and independent subgroups of neurons for all ensemble sizes. In all panels error bars correspond to the s.e.m.

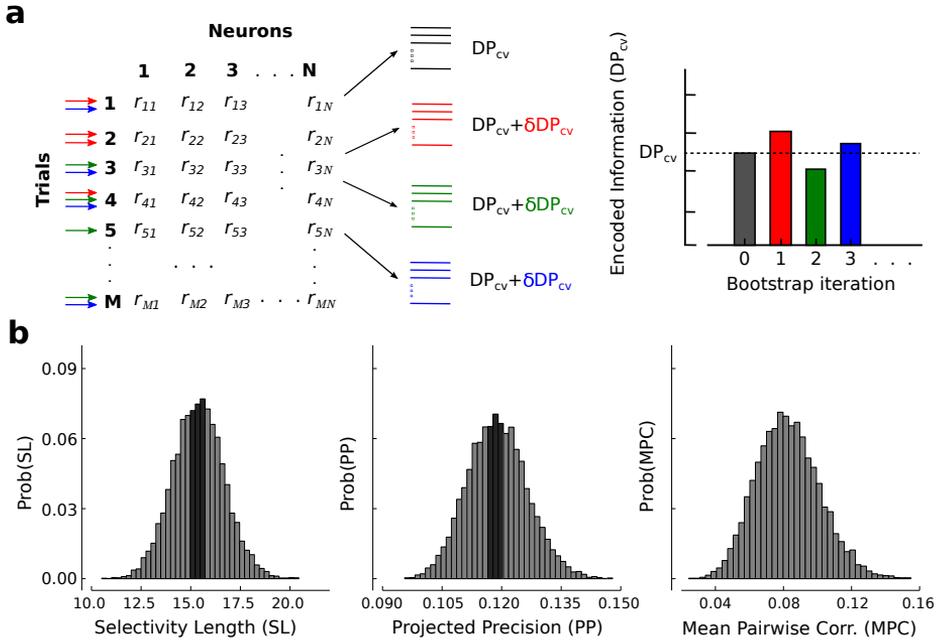
section 2.4).

## 2.2.5 Selectivity length and projected precision are the most important features for behavior

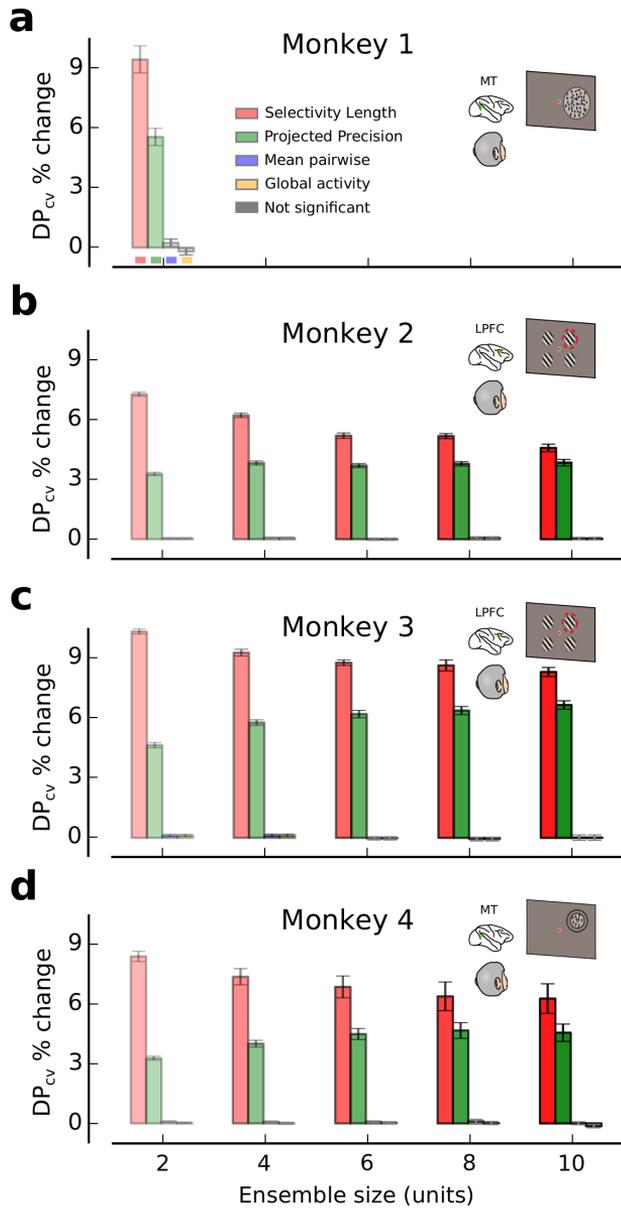
Once it has been shown that SL and PP both play a fundamental role on the amount of encoded information at the neuronal level, it naturally arises the question of what is the role of these neural code's features on the task performance of behaving monkeys. If behavior is produced by downstream neurons reading out from brain areas responsible for encoding the task-relevant variables, behavioral performance should be therefore linked to the amount of encoded information by these neural ensembles (see Fig. 2.9 and section 2.4 for a detailed description of behavioral performance for each task).

To answer this question we proceeded in the same way as in Fig. 2.8, but instead of focusing on encoded information ( $DP_{cv}$ ), we evaluated the percentage change in monkeys' performance ( $B$ ) when evaluated on those bootstrap iterations that elicited positive perturbations with respect to those that elicited negative perturbations for each feature of the neural code (SL, PP, MPC and GA) while keeping the rest fixed (see section 2.4). As in Fig. 2.8, the results are reported separately for all the different ensemble sizes used in this study (2, 4, 6, 8 and 10 units).

The percentage change produced in the performance of monkey 1 when evaluated on positive perturbations of SL and PP with respect to their negative perturbations (conditioned to iterations that fixed the rest



**Figure 2.7: The amount of encoded information, the features of the neural code and the performance of the animals can be perturbed by producing virtual instances of the original dataset. (a)** By subsampling trials with replacement (bootstrap) from the original dataset, perturbations on the amount of encoded information ( $DP_{cv}$ ), a range of values for the features of the neural code and animal performance (B) can be generated. On each bootstrap iteration we can calculate the amount of information encoded by the network and evaluate how these perturbations are affected by perturbations of the selectivity length (SL), projected precision (PP), mean pairwise correlations (MPC) and global activity (GA) (see section 2.4). **(b)** Distributions of SL, PP and MPC generated by the bootstrap method for a dataset (Fig. 2.2c; ensemble size = 10 units). Because perturbations on the different features of the neural code can be correlated we used a conditioning method where only a particular subset of bootstrap iterations were used. When studying for instance how perturbations of MPC affected the amount of encoded information, we used only those bootstrap iterations that produced negligible perturbations of both SL and PP (dark grey region) (see section 2.4).

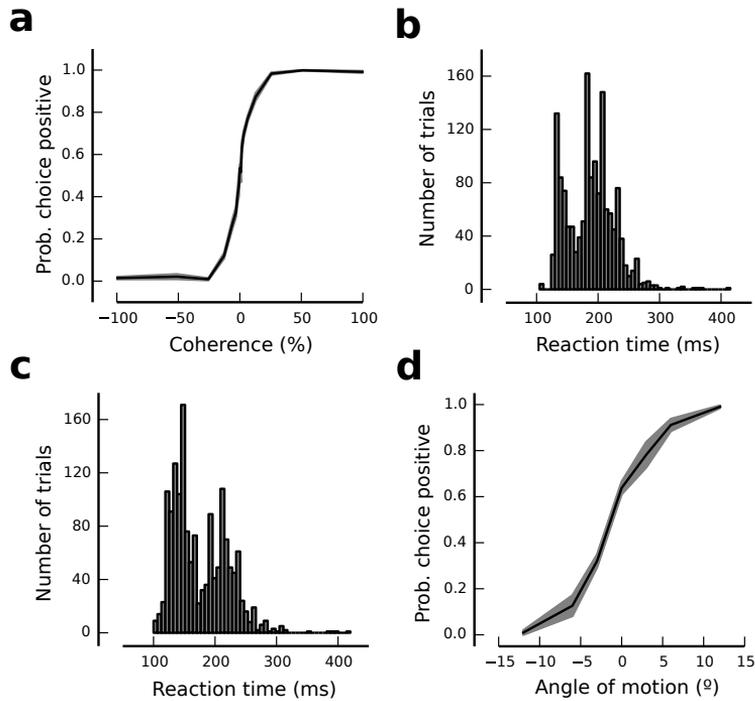


**Figure 2.8: Encoded information is not affected by mean pairwise correlation or global activity if SL and PP are accounted for.** (a) Percentage change in the amount of information encoded by the network ( $DP_{cv}$ ) when selectively

**Figure 2.8:** (cont.) perturbing the different features of the neural code (selectivity length (SL): red; projected precision (PP): green; mean pairwise correlations (MPC): blue; global activity (GA): orange) as a function of the different ensemble sizes used in this study (2, 4, 6, 8 and 10 units). Only SL and PP play a role on the amount of information encoded by neural populations. Results corresponding to pairs neurons recorded in MT during a coarse discrimination motion task [35] (monkey 1; see Fig. 2.2a and section 2.4). **(b)-(c)** Analogous to (a) but on ensembles of neurons simultaneously recorded ( $\sim 50$ ) in LPFC 8a during an attentional task [105] (monkeys 2 and 3; see Fig. 2.2b and section 2.4). **(d)** Analogous to (a-c) but on ensembles of units simultaneously recorded ( $\sim 25$ ) in middle temporal cortex (MT) during a fine discrimination motion task (monkey 4; see Fig. 2.2c and section 2.4). In all panels errorbars correspond to s.e.m. and significant deviations from zero are calculated by a Wilcoxon signed rank test (not significant if  $P > 0.05$ ).

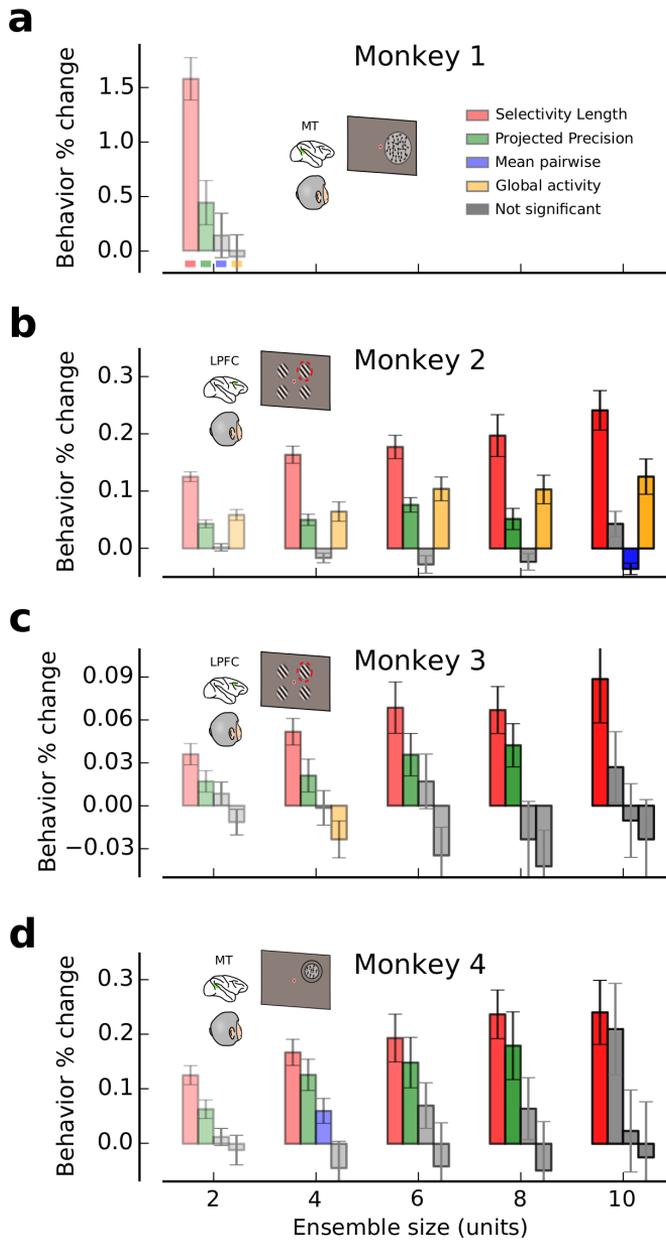
of features; see section 2.4) was  $1.58\% \pm 0.19\%$  (Wilcoxon signed-rank test,  $P = 3.3 \times 10^{-12}$ ) and  $0.44\% \pm 0.20\%$  ( $P = 2.1 \times 10^{-3}$ ) respectively (see Fig. 2.10a). When assessed on mean pairwise correlations and global activity of the network we obtained a non-significant percentage change of  $0.14\% \pm 0.20\%$  (Wilcoxon signed-rank test,  $P = 0.52$ ) and  $-0.06\% \pm 0.20\%$  ( $P = 0.66$ ). For monkey 2, 3 and 4 (Fig. 2.10b,c,d) we obtained qualitatively similar results, where SL and PP were in general the most influential factors on the monkeys' task performance across all ensemble sizes. For monkey 2 it was found an effect of GA on the performance of the monkey across ensemble sizes, but it was not consistent with the rest of datasets. Our results indicate that, on average, on those trials where SL and PP were larger by chance, the performance of the task was larger because monkeys had access to a larger amount of encoded information by the corresponding brain areas. As SL and PP are the most important parts of the code under optimality (see section 2.4), results shown in Fig. 2.10 further support the idea that encoded information is optimally read out to guide behavior [111].

By directly correlating how perturbations on the amount of encoded information ( $DP_{cv}$ ) affected perturbations on behavioral performance, we



**Figure 2.9: Behavioral performance definitions for each task.** (a) Psychometric curve for monkey 1. Percentage of positive choices as a function of motion coherence. On difficult trials the percentage of correct choices is lower than on easy trials (grey region depicts the s.e.m.). (b)-(c) Distribution of reaction times for monkeys 2 and 3. In the attentional task behavioral performance is defined as mean time (across trials in a particular dataset) elicited from the change in orientation until the saccade to the cued Gabor filter (saccadic movement). (d) Psychometric curve for monkey 4. Percentage of positive choices as a function of the angle of motion with respect to the vertical. On difficult trials the percentage of correct choices is lower than on easy trials (grey region depicts the s.e.m.).

also confirmed that the encoded and the decoded features of the neural code matched (Fig. 2.11). So to say, we tested whether encoded information correlated with behavior before splitting the former into its most important features. In all monkeys and ensemble sizes we found a sig-



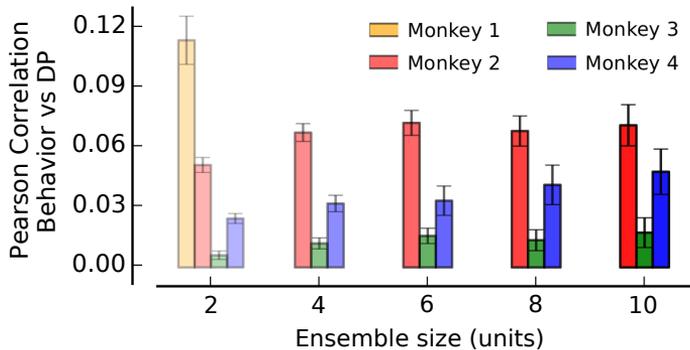
**Figure 2.10: Selectivity length and projected precision are the most important features on behavior.** (a) Percentage change in behavioral performance when selectively perturbing the different quantities of the neural code

**Figure 2.10:** (cont.) (selectivity length (SL): red; projected precision (PP): green; mean pairwise correlations (MPC): blue; global activity (GA): orange) as a function of the different ensemble sizes used in this study (2, 4, 6, 8 and 10 units). SL and PP are the most important factors influencing monkey’s performance in a coarse discrimination motion task while pairs of neurons were recorded in MT [35] (monkey 1; see see Fig. 2.2a and section 2.4). (b)-(c) Analogous to (a) but on ensembles of neurons simultaneously recorded ( $\sim 50$ ) in LPFC 8a during an attentional task [105] (monkeys 2 and 3; see Fig. 2.2b and section 2.4). (d) Analogous to (a-c) but on ensembles of units simultaneously recorded ( $\sim 25$ ) in middle temporal cortex (MT) during a fine motion discrimination task (monkey 4; see Fig. 2.2c and section 2.4). SL and PP are the most influential factors on behavioral performance consistently across all datasets and ensemble sizes while MPC and GA did, or inconsistently, relate to behavior. In all panels errorbars correspond to s.e.m. and significant deviations from zero are calculated by a Wilcoxon signed rank test (not significant if  $P > 0.05$ ).

nificant correlation between these two quantities (for ensemble size = 2 units, monkey 1:  $\rho = 0.11 \pm 0.01$  (s.e.m.), Wilcoxon signed-rank test,  $P = 4.5 \times 10^{-14}$ ; monkey 2:  $\rho = 0.051 \pm 0.004$ ,  $P = 4.8 \times 10^{-16}$ ; monkey 3:  $\rho = 0.006 \pm 0.002$ ,  $P = 0.011$ ; monkey 4:  $\rho = 0.024 \pm 0.002$ ,  $P = 2.1 \times 10^{-10}$ ). Results obtained in Fig. 2.10 were qualitatively equivalent when using different metrics to quantify the relationship between perturbations of the neural code with perturbations on behavioral performance and when using different conditioning windows (Fig. 2.12; see section 2.4).

## 2.2.6 A biologically-constrained neuronal model accounts for the experimental findings

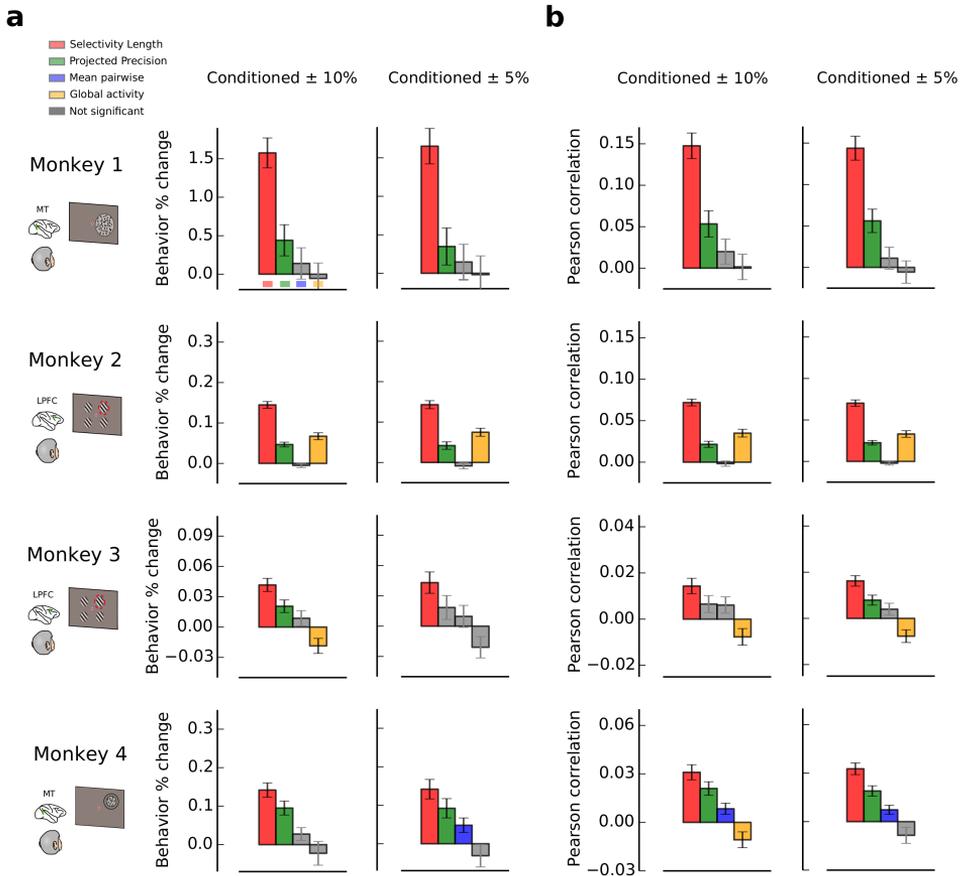
The results depicted by Figs. 2.10 and 2.11 show a relatively weak coupling between the most important features of the neural code or the amount of encoded information with behavior. If in all monkeys the recorded areas are responsible for encoding the task-relevant variables of the experiment (MT for coarse and fine motion discrimination task and LPFC 8a



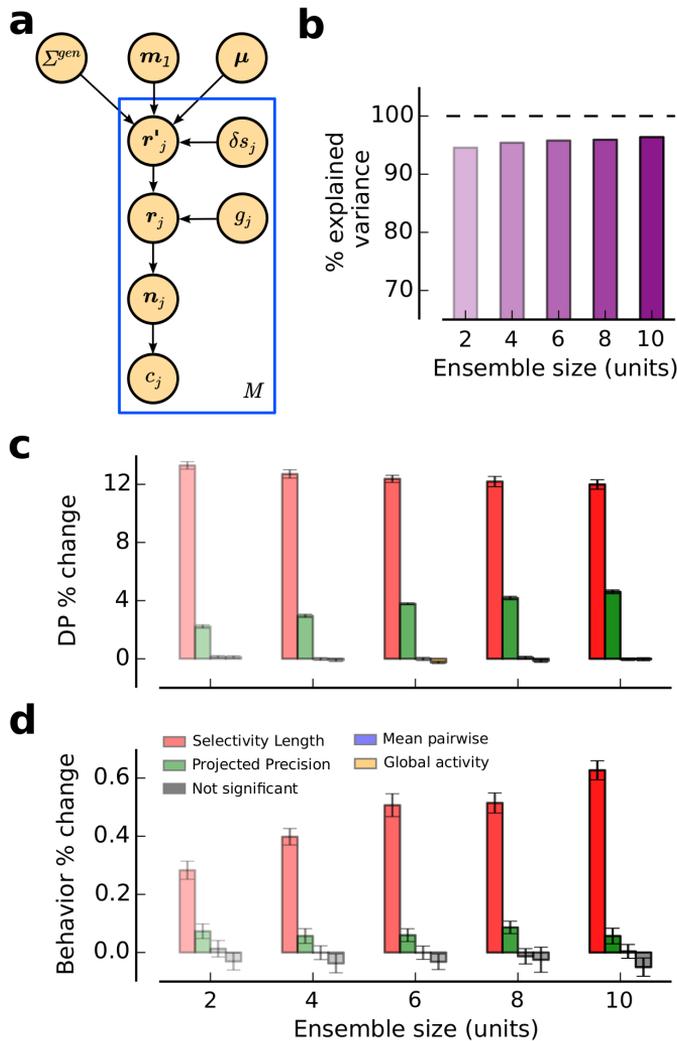
**Figure 2.11: Amount of encoded information correlates with behavioral performance in all datasets and ensemble sizes.** Pearson correlation between perturbations on the amount of information encoded in the network and perturbations on monkeys’ performance for all monkeys and ensemble sizes (bootstrap method; see Fig. 2.7 and 2.4). Across all datasets and ensembles sizes trials associated with larger encoded information are also associated with better task performances. Errorbars correspond to s.e.m. and significant deviations from zero are calculated by a Wilcoxon signed rank test (not significant if  $P > 0.05$ ).

for an attentional task), a larger coupling could be expected.

To understand the reasons how this weak correlation could be obtained, we built a simple neuronal model used to generate a set of simulated ensemble activity on a trial-by-trial basis that we could fit later in the analysis we described above (Fig. 2.13a). In each trial, a stimulus was presented to the ensemble and the population activity ( $N = 1000$  units) was drawn from a multivariate Gaussian distribution. Before applying the Poisson step on a particular trial, the activity of the population was transformed by a common additive and multiplicative global modulation [51]. Limited-range correlations were used [46, 51, 55–57, 113] and an additional small term of differential correlations was added to the covariance matrix to control the information in the network [47] as the network size increased. Behavior was also simulated in each trial by reading out optimally from the entire neuronal ensemble. Surrogate neural code’s features, encoded information and behavior were calculated using the same



**Figure 2.12: Results obtained in Fig. 2.10 are robust under different conditionings and metrics.** (a) Equivalent to Fig. 2.10 averaged across ensemble sizes. On the vertical it is plotted Behavior % change (see section 2.4) and for the conditioning it has been used those bootstrap iterations that produced perturbations departing  $\pm 10\%$  (left column) and  $\pm 5\%$  (right column) from the median value of the distribution. (b) Pearson correlation between perturbations of the different features of the neural code and behavioral performance averaged across ensemble sizes. For the conditioning it has been used those bootstrap iterations that produced perturbations departing  $\pm 10\%$  (left column) and  $\pm 5\%$  (right column) from the median value of the distribution.



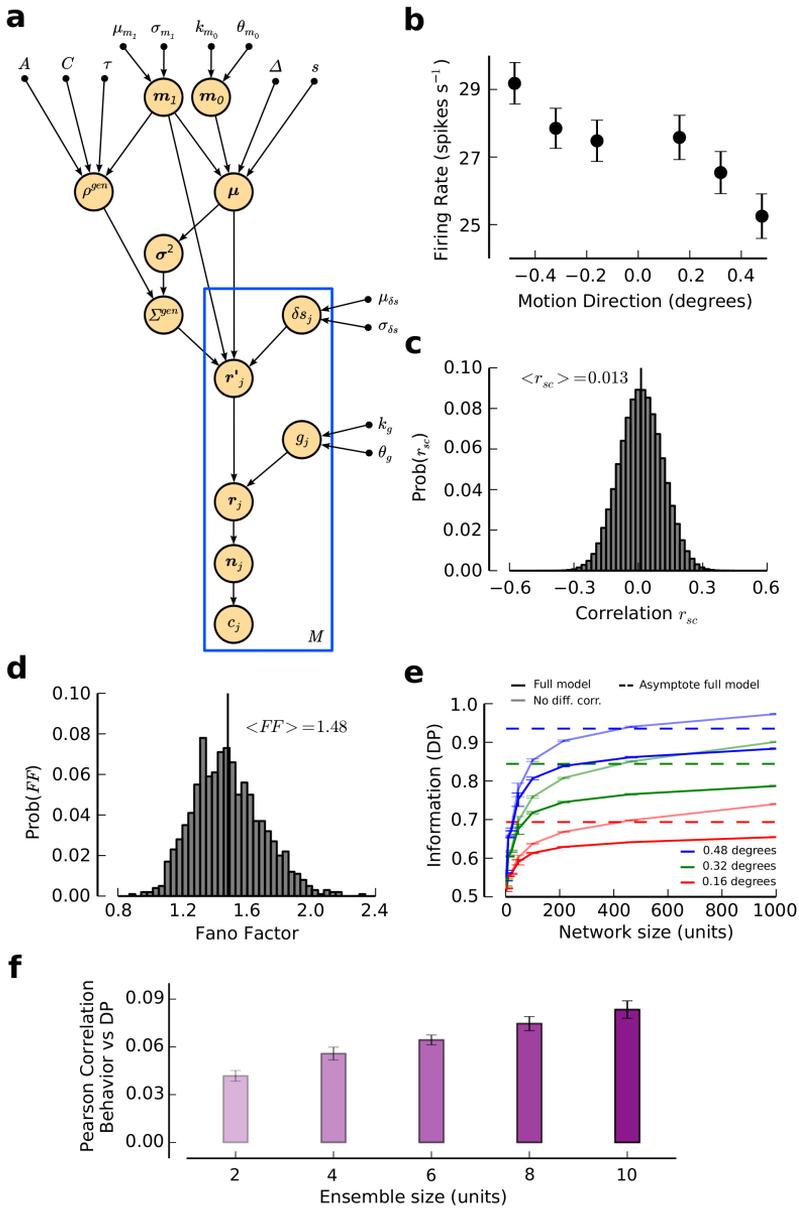
**Figure 2.13: A biologically constrained neuronal model of MT and LPFC accounts for the all the experimental findings. (a)** A preliminary activity pattern  $r'_j$  was obtained by drawing an  $N$ -dimensional sample from a multivariate Gaussian distribution and corrupting it with sensory noise ( $\delta s_j$ ) on trial  $j$ . Then, a homogeneous modulation and a Poisson step were applied to produce the population spike count ( $n_j$ ). The choice of the virtual agent ( $c_j$ ) was obtained by reading-out optimally the population activity pattern on each trial.

**Figure 2.13:** (cont.) **(b)**  $DP_{th}$  was a very good approximation of  $DP_{cv}$  on the surrogate data as well (see section 2.4 and Fig. 2.3). **(c)** Percentage change in the amount of information encoded by the network  $DP_{cv}$  when selectively perturbing the different features of the neural code as a function of the different ensemble sizes. Only SL and PP play a role on the amount of information encoded by neural populations (see Fig. 2.8). **(d)** Equivalent to (c) but on the behavioral performance of the virtual agent. SL and PP are the most important factors influencing the behavioral performance (see Fig. 2.10).

time window (1 sec) (Figs. 2.13a and 2.14; see section 2.4 for a full description of the model).

In Fig. 2.13b it is shown the percentage of variance depicted by a cross-validated linear classifier trained and tested on the surrogate dataset ( $DP_{cv}$ ) that can be explained by Eq. (2.2). As in Fig. 2.3b, Eq. (2.2) is a very good approximation for the amount of encoded information in the neural network (mean explained variance = 95.6%). Consistently with Fig. 2.8a-d (monkey 1-4), in Fig. 2.13c it is shown that SL and PP are the only features of the neural code affecting information while mean pairwise correlations and global activity played no role in the amount of encoded information by the surrogate population.

As surrogate trial-by-trial behavior is fully defined from the optimal read-out of the whole ensemble's activity, a high correlation between perturbations in  $DP_{cv}$  and behavioral performance could be expected. We show that (Fig. 2.13f) a model of surrogate behavior determined by an optimal read-out of a whole population can produce the counterintuitive weak coupling between encoded information and behavior (see Fig. 2.11 for the comparison with *in vivo* recordings). This weak coupling is mainly determined by the small proportion of units used to assess ensemble's encoded information ( $DP_{cv}$ ) when compared to the full network driving behavior. Finally in Fig. 2.13d it is shown that qualitatively equivalent results to Fig. 2.10 can be found under this model. When reading-out optimally from the information-encoding ensemble, SL and PP are the most important features of the neural code affecting behavior.



**Figure 2.14: A biologically-constrained model of the cortex can account for the weak coupling between encoded information and behavioral performance**

**Figure 2.14:** (cont.) **(a)** Graphical model depicting in full detail the generative process underlying the surrogate dataset (see 2.4 for a detailed description of the different terms). **(b)** Tuning curve of an example neuron used in the model. This neuron is encoding the direction of motion presented to the network ( $s_1$  or  $s_2$ ) on its firing rate. The larger the angle of motion with respect to the vertical, the larger the difference between the firing rate associated with  $s_1$  and  $s_2$  (more information). **(c)** Pearson correlation on the spike count of all the neuronal pairs in an example recording session (distribution of pairwise correlations). The mean value of the distribution is in agreement with the empirical values of mean pairwise correlations found in in vivo recordings [45]. **(d)** Distribution of Fano Factors for all the units in an example dataset. The mean value of the distribution is in agreement with the empirical values of fano factors found in in vivo recordings **(e)** Amount of encoded information as a function of the network size ( $DP_{cv}$ ) for the three different stimulus intensities. When not including differential correlation in the model information does not saturate before reaching the natural boundary  $DP_{cv} = 1.0$  (light colored lines) for any of the stimulus intensities. When including differential correlations on the model, information is decreased (colored lines) and an asymptote on  $DP_{cv}$  is produced for very large network sizes (dashed colored line). **(f)** Pearson correlation between the perturbations on the amount of encoded information and the behavioral performance of the virtual agent. A simple model of the cortex can account for the counterintuitive weak coupling between these two quantities (see Fig. 2.11).

## 2.3 Discussion

Identifying what are the most influential features of the neural code on the amount of information encoded by a neuronal ensemble is a fundamental step for fully characterizing the link between neuronal activity and behavior. By deriving an analytical expression for the performance of a downstream neuron linearly reading-out a population neurons (Eq. (2.2)), we have been able to identify the tuning and noise factors determining the amount of encoded information. Even though some assumptions have been used for the analytical derivation, Eq. (2.2) is a very good approximation for information as expressed by the decoding performance of a

cross-validated linear classifier on four different datasets involving different brain areas and tasks ( $DP_{cv}$ ).

To further confirm the results obtained by Eq. (2.2), we developed a novel non-parametric method where we generated virtual instances of the original datasets by bootstrapping trials with replacement and treated them as perturbations. By selectively perturbing the different features of the neural code we were able to study their effects on the amount of encoded information. Consistently with the parametric approach (Eq. (2.2)) we found that only selectivity length (SL) and projected precision (PP) played a role at all, whereas neither mean pairwise correlation nor global activity of the network had any effect on the amount of encoded information by a neuronal ensemble. Moreover, when we evaluated how neural code's features affected behavior, we found that SL and PP were the features affecting the most the behavioral performance of four different monkeys involving different brain regions and tasks. Finally, we were able to reproduce all the experimental findings by a simple neuronal model of the cortex.

It is important to remark that previous works have already identified the exact shape of pairwise correlations that limit information [46, 47, 113], however the present study represents an important contribution in the field because it extends these results to experimentally-realistic information measures and ensemble sizes and tests the predictions on *in vivo* recordings for both neuronal encoding of information and behavior. Assessing the amount of encoded information using fisher information (FI) on experimental datasets [47, 113] can be in many situations misleading because the locality assumption is very challenging to fulfill and because of the difficulty to get an accurate estimate of the tuning curve's derivative  $f'$  in real recordings. Moreover, in most of the previous theoretical studies, the aim is to characterize how information behaves for large neuronal populations, whereas a theoretical scheme was missing for small and experimentally-realistic ensemble sizes. Nevertheless, it is important to remark that important contributions have been made in this direction (Kanitscheider 2015 Plos). Here we claim that by evaluating the amount of encoded information as expressed by the DP of a linear

classifier we naturally extend the concept of FI to discrete inference tasks and therefore to realistic electrophysiological datasets as the training of a linear classifier is agnostic to the explicit knowledge of the tuning curves and covariance matrices. Consistently with these studies [47, 113], our results conclusively show that mean pairwise play no role at all on the amount of information encoded by a network. Shuffling across trials for a fixed stimulus condition is a highly popular technique as well when studying the role of mean pairwise correlations on the neural code's information capabilities [51, 61, 105, 112]. The bootstrap method presented in this study complements the shuffling method as both represent non-parametric techniques for modifying the original dataset. However, Eq. (2.2) provides an analytical generalization of the shuffling method because whether pairwise correlations are harmful or beneficial for the neural code can be calculated by comparing the original projected precision with its equivalent after removing the off-diagonal terms in the covariance matrix. Interestingly our expression is generalizable to any linear classifier (Eq. (2.1)) and therefore robust predictions could be made about the role of shared and individual variability under sub-optimal information-extracting strategies.

Another important feature of the neural code is the global activity of the network, defined as the mean activity across the neuronal population for a particular time window. Even though it has been shown previously that global activity plays no role on the amount of encoded information [51], Eq. (2.2) provides a general explanation for this result that can be extrapolated to most of the nowadays experimental datasets. Even though it has not been explicitly stated in this study, if the global activity is split into an additive [42, 43, 51] and a multiplicative term [42, 44, 51, 62, 63], the former does not have any effect on either SL nor PP while the latter implies solely an elongation of the selectivity vector (increase in SL), as experimentally confirmed on V1 neurons in monkeys encoding the orientation of moving gratings [51]. Moreover, this study is the first one providing experimental evidence about the weak link between both global activity of the network and mean pairwise correlations with the performance of behaving animals. It is important to note the existence

of previous studies reporting a link between mean pairwise correlations and behavioral performance [58–60], however our claim here is that the real correlate with performance is the projected precision. Mean pairwise correlations and projected precision are two highly correlated variables and therefore an apparent modulation of performance with mean pairwise correlations could arise if not controlling for projected precision.

Overall in this study we claim for a re-evaluation of the amount of resources devoted to studying the role of pairwise correlations on both encoding of information and behavior. Mean pairwise correlations are a very appealing feature of the neural code from the experimental perspective because obtaining an accurate estimate of them is relatively easy. As shown in this study both parametrically (Eq. (2.2)) and non-parametrically (perturbations by bootstrap), this feature is completely decoupled from the amount of information encoded by a neuronal ensemble and very weakly coupled with behavioral performance. Even though some previous studies already pointed in this direction [47, 113], this study serves as a link between the electrophysiological and the theoretical studies by showing the important features of the neural code on experimentally-realistic ensemble sizes and information measures like the decoding performance of a cross-validated linear classifier.

## **2.4 Methods**

### **2.4.1 Analytical expressions for encoded information**

The amount of information encoded by a neural network can be characterized by the percentage of correct classifications (decoding performance) that can be performed when linearly reading-out the activity of a population of neurons on a binary classification task. On the following we are going to derive a set of analytical expressions for the decoding performance (DP) of a set of linear classifiers.

### 2.4.1.1 Analytical DP for an arbitrary linear classifier

A binary classifier aims to classify an  $N$ -dimensional instance of a random variable  $\mathbf{x}$  as belonging to one of two possible classes  $C_1$  or  $C_2$ . The percentage of correctly classified instances is the Decoding Performance (DP) and, given the conditional distributions  $p(\mathbf{x}|C_1)$  and  $p(\mathbf{x}|C_2)$ , it can be expressed as

$$P[\text{correct}] = \int_{\Omega_1} p(\mathbf{x}|C_1)p(C_1)d\mathbf{x} + \int_{\Omega_2} p(\mathbf{x}|C_2)p(C_2)d\mathbf{x}, \quad (2.3)$$

where  $\Omega_1$  and  $\Omega_2$  are the regions of the  $\mathbf{x}$  space corresponding to  $C_1$  and  $C_2$ . We will focus on linear classifiers, where the boundary that delimits the two regions is linear on  $\mathbf{x}$ . A linear boundary on an  $N$ -dimensional space is an  $N - 1$  hyperplane characterized by the  $N$ -dimensional linear constraint  $\boldsymbol{\omega}^T \mathbf{x} + \omega_0 = 0$ . Eq. (2.3) can be re-expressed in terms of the decision variable, a 1-dimensional variable defined as  $z = \boldsymbol{\omega}^T \mathbf{x} + \omega_0$

$$P[\text{correct}] = \int_{Z_1} p(z|C_1)p(C_1)dz + \int_{Z_2} p(z|C_2)p(C_2)dz, \quad (2.4)$$

where  $Z_1$  and  $Z_2$  correspond to the regions defining  $C_1$  and  $C_2$  respectively. From now on we will refer to  $\boldsymbol{\omega}$  and to the classifier indistinctly. By setting  $\omega_0$  conveniently, we can make  $Z_1 = [0, \infty)$  and  $Z_2 = (-\infty, 0)$ . The above expression becomes then

$$\begin{aligned} P[\text{correct}] &= \int_0^\infty p(z|C_1)p(C_1)dz + \int_{-\infty}^0 p(z|C_2)p(C_2)dz \\ &= p(z > 0|C_1)p(C_1) + p(z < 0|C_2)p(C_2). \end{aligned} \quad (2.5)$$

If we assume the random variable  $\mathbf{x}$  (features) follows a multivariate Gaussian distribution and the covariance matrices are stimulus independent, then a linear combination of  $N$  normally distributed variables

will also produce a normally distributed variable, with mean and variance given by

$$\begin{aligned}\mu_{1z} &= \boldsymbol{\omega}^T \boldsymbol{\mu}_1 + \omega_0 \\ \mu_{2z} &= \boldsymbol{\omega}^T \boldsymbol{\mu}_2 + \omega_0 \\ \sigma_z^2 &= \boldsymbol{\omega}^T \boldsymbol{\Sigma} \boldsymbol{\omega}\end{aligned}\quad (2.6)$$

where  $E[\mathbf{x}|C_1] = \boldsymbol{\mu}_1$ ,  $E[\mathbf{x}|C_2] = \boldsymbol{\mu}_2$  and  $E[\mathbf{x}\mathbf{x}^T|C_1] = E[\mathbf{x}\mathbf{x}^T|C_2] = \boldsymbol{\Sigma}$ . Equation (2.5) can then be written as

$$P[\text{correct}] = \Phi\left(\frac{\mu_{1z}}{\sigma_z}\right) p(C_1) + \Phi\left(\frac{-\mu_{2z}}{\sigma_z}\right) p(C_2), \quad (2.7)$$

where  $\Phi(\cdot)$  is the cumulative standard normal distribution, defined as  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2} dt$ .

To find the desired value of  $\omega_0$  we start by calculating  $E[z]$

$$\begin{aligned}E[z] &= E[p(z|C_1)p(C_1) + p(z|C_2)p(C_2)] \\ &= \frac{1}{2}\boldsymbol{\omega}^T \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\omega}^T \boldsymbol{\mu}_2 + \omega_0,\end{aligned}\quad (2.8)$$

where we have assumed a uniform prior over the classes. By imposing  $E[z] = 0$  we ensure the classifier is unbiased, and therefore

$$\omega_0 = -\frac{1}{2}\boldsymbol{\omega}^T (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2). \quad (2.9)$$

Substituting Eqs. (2.9) and (2.6) into Eq. (2.7) we obtain

$$P[\text{correct}] = \Phi\left(\frac{\frac{1}{2}\boldsymbol{\omega}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\sqrt{\boldsymbol{\omega}^T \boldsymbol{\Sigma} \boldsymbol{\omega}}}\right) p(C_1) + \Phi\left(\frac{\frac{1}{2}\boldsymbol{\omega}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\sqrt{\boldsymbol{\omega}^T \boldsymbol{\Sigma} \boldsymbol{\omega}}}\right) p(C_2). \quad (2.10)$$

If  $p(C_1) = p(C_2) = 0.5$ ,  $\Delta \mathbf{f} \equiv (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$  and  $P[\text{correct}] = DP$ , the decoding performance (DP) of a binary arbitrary linear classifier can be written as

$$DP = \Phi \left( \frac{1}{2} \frac{\boldsymbol{\omega}^T \Delta \mathbf{f}}{\sqrt{\boldsymbol{\omega}^T \Sigma \boldsymbol{\omega}}} \right). \quad (2.11)$$

This expression therefore represents the percentage of correct classifications that would ideally be achieved by linearly reading out from a neuronal ensemble using an arbitrary set of weights  $\boldsymbol{\omega}$ . If we define encoded information as the upper bound on the amount of information that can be extracted from the population of neurons and the decoded information as the amount of information that is actually extracted from this population by an arbitrary classifier [61], the optimal classifier ( $\boldsymbol{\omega}_{opt}$ ) corresponds to the set of weights that decode all the encoded information in the network.

#### 2.4.1.2 Optimal linear classifier

To obtain the optimal classifier  $\boldsymbol{\omega}_{opt}$  we differentiate Eq. (2.11) with respect to the weights and set it to zero

$$\begin{aligned} \nabla DP &= \frac{1}{2} \Phi' \nabla \left( \frac{\boldsymbol{\omega}^T \Delta \mathbf{f}}{\sqrt{\boldsymbol{\omega}^T \Sigma \boldsymbol{\omega}}} \right) \\ &= \frac{\Phi'}{2\sqrt{\boldsymbol{\omega}^T \Sigma \boldsymbol{\omega}}} \left( \Delta \mathbf{f} - \frac{\boldsymbol{\omega}^T \Delta \mathbf{f} \Sigma \boldsymbol{\omega}}{\boldsymbol{\omega}^T \Sigma \boldsymbol{\omega}} \right) = 0 \quad (2.12) \\ &\Rightarrow \Delta \mathbf{f} - \frac{\boldsymbol{\omega}^T \Delta \mathbf{f} \Sigma \boldsymbol{\omega}}{\boldsymbol{\omega}^T \Sigma \boldsymbol{\omega}} = 0 \end{aligned}$$

The solution to the above equation is  $\boldsymbol{\omega}_{opt} \propto \Sigma^{-1} \Delta \mathbf{f}$ . This solution corresponds to a maximum for Eq. (2.12), as it can be shown that  $\nabla^2 DP|_{\boldsymbol{\omega}_{opt}} < 0$ . It is worth mentioning that the length of the vector  $\boldsymbol{\omega}_{opt}$  is a free parameter and therefore any classifier parallel to  $\boldsymbol{\omega}_{opt}$  will be solution for the above equation.

#### 2.4.1.3 Analytical DP for the optimal linear classifier

If  $\boldsymbol{\omega}_{opt}$  is introduced in Eq. (2.12), the analytical DP for the optimal classifier becomes

$$DP = \Phi \left( \frac{1}{2} \sqrt{\Delta \mathbf{f} \Sigma^{-1} \Delta \mathbf{f}} \right). \quad (2.13)$$

The term inside the cumulative is also known as  $d' = \sqrt{\Delta \mathbf{f}^T \Sigma^{-1} \Delta \mathbf{f}}$  [110]. It is a scalar quantity and therefore it remains invariant under unitary rotations of the reference frame. If we rotate the original neuron-based orthogonal basis so that  $\Sigma$  becomes diagonal (and also  $\Sigma^{-1}$ ), the above expression becomes

$$DP = \Phi \left( \frac{1}{2} |\Delta \mathbf{f}| \sqrt{\sum_{i=1}^N \frac{\cos^2 \hat{\theta}_i}{\hat{\sigma}_i^2}} \right). \quad (2.14)$$

(see [47] for a similar derivation for the case of linear Fisher information). Let's now analyze each of the following terms in detail. We refer to  $|\Delta \mathbf{f}|$  as Selectivity Length (SL) and it is the norm of the selectivity vector  $\Delta \mathbf{f}$ , which measures the overall modulation of the activity of the neuronal population as a function of the stimulus conditions. If the neurons in the population are strongly tuned, SL will be large and it will be easier to separate the ellipsoids defining the neuronal activity for each condition (Fig. 2.1b). The second term is called the Projected Precision (PP), and it is equal to the square root of the sum over all the ellipsoids' main axes (number of neurons) of  $\frac{\cos^2 \hat{\theta}_i}{\hat{\sigma}_i^2}$ . Each  $\hat{\theta}_i$  is the angle between the  $i$ -th eigenvector of the covariance matrix  $\Sigma$  and the stimulus axis  $\mathbf{u}_{\Delta \mathbf{f}}$ , the unitary direction defined by the stimulus vector  $\Delta \mathbf{f}$ . Each  $\hat{\sigma}_i^2$  is the  $i$ -th eigenvalue of the  $\Sigma$  matrix (Fig. 2.1b). Equation (2.14) can therefore be written as  $DP = \Phi \left( \frac{1}{2} SL \times PP \right)$ .

Despite the Gaussian and the stimulus-independent noise approximations, the analytical DP (Eq. (2.14)) is a very good predictor of the DP of a cross-validated linear classifier trained on *in vivo* electrophysiological recordings (Fig. 2.3). Moreover, this expression accounts for a larger amount of the variability depicted by a cross-validated linear classifier than the equivalent expressions for the variability-blind and the correlation-blind suboptimal classifiers (see below; see Fig. 2.5a).

#### 2.4.1.4 Analytical DP for shuffled neuronal recordings

Previous studies have compared the performance of a linear classifier on shuffled and unshuffled data to understand the role of noise correlations on DP [61, 105]. Eq. (2.14) provides an answer to the role of noise correlations on the DP of a linear classifier, but it also shows how destroying correlations affects DP. When noise correlations are destroyed by shuffling trials per neuron, the covariance matrix  $\Sigma$  is transformed into  $\Sigma_{sh}$ , which is close to a diagonal matrix except for finite data size effects. In order to find an analytical expression for the shuffled DP we just need to substitute  $\Sigma$  by  $\Sigma_{sh}$  in Eq. (2.11)

$$DP = \Phi \left( \frac{1}{2} \frac{\boldsymbol{\omega}^T \Delta \mathbf{f}}{\sqrt{\boldsymbol{\omega}^T \Sigma_{sh} \boldsymbol{\omega}}} \right). \quad (2.15)$$

The optimal classifier is therefore  $\boldsymbol{\omega}_{opt} \propto \Sigma_{sh}^{-1} \Delta \mathbf{f}$  and Eq. (2.14) becomes

$$DP = \Phi \left( \frac{1}{2} |\Delta \mathbf{f}| \sqrt{\sum_{i=1}^N \frac{\cos^2 \theta_i}{\sigma_i^2}} \right), \quad (2.16)$$

where the selectivity length is the same as in Eq. (2.14) and  $\frac{\cos^2 \theta_i}{\sigma_i^2}$  is the Shuffled Projected Precision (SPP). In the SPP,  $\sigma_i^2$  is the variability of each neuron and  $\theta_i$  is the angle between the stimulus axis  $\mathbf{u}_{\Delta \mathbf{f}}$  and each vector of the neuron-based basis, that is, the angle between  $\mathbf{u}_{\Delta \mathbf{f}}$  and each of the axis defining the original  $N$ -dimensional space of the neuronal activity. As stated before, due to the finite size effects, it is important to keep in mind that this is only true when considering the mean across shuffling iterations or in the limit of very large datasets (many trials). The above equation can also be expressed as

$$DP = \Phi \left( \frac{1}{2} \sqrt{\sum_{i=1}^N \frac{(\Delta f_i)^2}{\sigma_i^2}} \right), \quad (2.17)$$

where  $\sqrt{\sum_{i=1}^N \frac{(\Delta f_i)^2}{\sigma_i^2}}$  can be understood as the sum of the signal-to-noise ratio (SNR) of all the neurons in the ensemble. It is important to note that Eq. (2.16) is a better approximation for the DP of a cross-validated linear classifier than Eq. (2.14) when pairwise correlations are removed (Fig. 2.5b).

### 2.4.1.5 Analytical DP under suboptimal read-outs

As stated earlier,  $\omega_{opt}$  represents the set of weights that match encoded with decoded information (optimal). However, information could also be extracted sub-optimally. In the following we are providing the set of analytical expressions for decoded information when using the variability-blind classifier and the correlation-blind classifier [111].

#### 2.4.1.5.1 Variability-blind classifier

The variability-blind classifier only takes into account the neuronal selectivity and it is blind to any of the elements of the covariance matrix. Thus, it assumes that the readout vector is such that  $\omega_{opt} \propto \Delta \mathbf{f}$ . Introducing this expression in Eq. (2.11), we find

$$DP = \Phi \left( \frac{1}{2} \frac{|\Delta \mathbf{f}|}{\sqrt{\sum_{i=1}^N \hat{\sigma}_i^2 \cos^2 \hat{\theta}_i}} \right), \quad (2.18)$$

where  $\hat{\sigma}_i^2$  and  $\hat{\theta}_i$  are defined in the same way as in Eq. (2.14). For this sub-optimal classifier, as for the optimal, the stronger the neuronal selectivity (or selectivity length) the higher the classification performance. Moreover, if the smallest principal axis of the covariance matrix is aligned with the stimulus axis  $\mathbf{u}_{\Delta \mathbf{f}}$ , a very large DP could be achieved.

#### 2.4.1.5.2 Correlation-blind classifier

The correlation-blind classifier [111] takes into account the signal-to-noise ratio of each neuron, but ignores the pairwise correlation pattern

of the ensemble. This classifier will be determined by the set of weights  $\omega_i \propto \frac{(\Delta f_i)^2}{\sigma_i^2}$ , which can be also written as  $\omega \propto \Sigma_{sh}^{-1} \Delta \mathbf{f}$ . By substituting this expression in Eq. (2.11) we obtain

$$DP = \Phi \left( \frac{1}{2} |\Delta \mathbf{f}| \frac{\sum_{i=1}^N \frac{\cos^2 \theta_i}{\sigma_i^2}}{\sqrt{\sum_{i=1}^N \frac{\cos^2 \tilde{\theta}_i}{\tilde{\sigma}_i^2}}} \right), \quad (2.19)$$

where, as in Eq. (2.16)  $\sigma_i^2$  is the variability across trials of each neuron's activity and  $\theta_i$  is the angle between the stimulus axis  $\mathbf{u}_{\Delta \mathbf{f}}$  and the  $i$ -th axis of the original neuronal space. The terms  $\tilde{\sigma}_i^2$  and  $\tilde{\theta}_i$  correspond to the eigenvalues and the angle between  $\mathbf{u}_{\Delta \mathbf{f}}$  and the eigenvectors of the matrix  $\Sigma_{sh}^{-1} \Sigma \Sigma_{sh}^{-1}$  respectively. As in previous section, it is important to note that this is only true for the mean  $\Sigma$  across shuffling iterations or for the case when the number of trials is very large. The correlation blind classifier, even though suboptimal when pairwise correlations are not removed, is the optimal classifier when each neuron's firing rate is shuffled across trials (see Fig. 2.5b).

#### 2.4.1.6 Analytical expression for DP and differential correlations

The task of our binary classifier is to assign one of the two possible labels to a multidimensional pattern of activity. In most of the experimental situations, the two discrete labels are just particular instances of a continuous variable. For example, in a motion direction task the two directions to be discriminated for the subject are just two opposite values of the angle of motion in the screen, a continuous variable. Therefore, even though in the experiment we are just measuring two particular instances ( $s_1$  and  $s_2$ ;  $\Delta s \equiv (s_1 - s_2)$ ) of the population's tuning curve ( $\mu_1$  and  $\mu_2$ ), in reality the tuning curve is a continuously varying function of the parameter to be inferred  $\mathbf{f}(s)$ . This mapping between the one-dimensional  $s$  space to the  $N$ -dimensional space  $\mathbf{f}(s)$  is assumed to be differentiable, and we will refer to it as  $\mathbf{f}'(s)$ . It is important to remark that if  $s_1 - s_2$  is very small, then  $\mathbf{f}'(s_1) = \mathbf{f}'(s_2) \simeq \frac{\Delta \mathbf{f}}{\Delta s}$ .

In [47] it is shown that when the input signal itself is noisy in a trial-by-trial basis, then a very particular type of correlations are created on the neuronal population: the differential correlations. These correlations can be very weak but their impact on information can be dramatic. Differential correlations set an upper bound for the amount of information that can be linearly extracted from a network, and therefore increasing the size of the network does not always imply increasing the amount of information, as the information entering the network is already finite.

We aimed to characterize what is the effect of sensory noise in both  $s_1$  and  $s_2$  on the DP of a linear classifier. Even though differential correlations are just defined for a neighborhood of  $s$ , we are going to characterize how sensory noise affects the DP when both stimuli are corrupted with noise from the input layer. Let's define  $\mathbf{f}'(s_1) \equiv \mathbf{f}'_1$  and  $\mathbf{f}'(s_2) \equiv \mathbf{f}'_2$  as the derivatives of the population's tuning curve in  $s_1$  and  $s_2$ . When the sensory input is noisy, the covariance matrix incorporates a new term  $\Sigma = \Sigma_0 + \epsilon \mathbf{f}' \mathbf{f}'^T$ . From now on, in order to follow the notation introduced in [47] we will refer to the covariance matrix without differential correlations as  $\Sigma_0$ . Because  $\mathbf{f}'_1 \neq \mathbf{f}'_2$  in general, we are first going to write Eq. (2.5) in the most general way possible

$$DP = \Phi \left( \frac{1}{2} \frac{\boldsymbol{\omega}^T \Delta \mathbf{f}}{\sqrt{\boldsymbol{\omega}^T \Sigma_1 \boldsymbol{\omega}}} \right) p(C_1) + \left( \frac{1}{2} \frac{\boldsymbol{\omega}^T \Delta \mathbf{f}}{\sqrt{\boldsymbol{\omega}^T \Sigma_2 \boldsymbol{\omega}}} \right), \quad (2.20)$$

where  $\Sigma_1 = \Sigma_{10} + \epsilon \mathbf{f}'_1 \mathbf{f}'_1{}^T$  and  $\Sigma_2 = \Sigma_{20} + \epsilon \mathbf{f}'_2 \mathbf{f}'_2{}^T$ . If we assume equal priors and equal covariance matrices  $\Sigma_1 = \Sigma_2 = \Sigma = \Sigma_0 + \epsilon \mathbf{f}' \mathbf{f}'^T$ , then Eq. (2.20) becomes

$$\begin{aligned} DP &= \Phi \left( \frac{1}{2} \frac{\boldsymbol{\omega}^T \Delta \mathbf{f}}{\sqrt{\boldsymbol{\omega}^T \Sigma_0 \boldsymbol{\omega} + \epsilon (\boldsymbol{\omega}^T \mathbf{f}')^2}} \right) \\ &= \Phi \left( \frac{1}{2} \frac{\mu_z}{\sqrt{\sigma_{0z}^2 + \epsilon \sigma_{\epsilon z}^2}} \right) \\ &= \Phi \left( \frac{1}{2} \frac{\mu_z}{\sigma_z} \right) \end{aligned} \quad (2.21)$$

where we have defined  $\mu_z = \boldsymbol{\omega}^T \Delta \mathbf{f}$ ,  $\sigma_{0z}^2 = \boldsymbol{\omega}^T \Sigma_0 \boldsymbol{\omega}$  and  $\sigma_{\epsilon z}^2 = \epsilon (\boldsymbol{\omega}^T \mathbf{f}')^2$ , and where  $\sigma_z^2 \equiv \sigma_{0z}^2 + \epsilon \sigma_{\epsilon z}^2$ . When introducing sensory noise on  $s$ , we increase the variance of the decision variable  $z$  and therefore the classifier is less precise.

To find the optimal weights for Eq. (2.21) we can use the solution already known for Eq. (2.11)

$$\begin{aligned}
\boldsymbol{\omega}_{opt} &\propto \Sigma^{-1} \Delta \mathbf{f} \\
&= \Sigma_0^{-1} \Delta \mathbf{f} - \epsilon \frac{\Sigma_0^{-1} \mathbf{f}' \mathbf{f}'^T \Sigma_0^{-1} \Delta \mathbf{f}}{1 + \epsilon \mathbf{f}'^T \Sigma_0^{-1} \mathbf{f}'} \\
&= \boldsymbol{\omega}_{0,opt} - \epsilon \frac{(\mathbf{f}'^T \boldsymbol{\omega}_{0,opt})}{1 + \epsilon (\mathbf{f}'^T \boldsymbol{\omega}_\epsilon)} \boldsymbol{\omega}_\epsilon \\
&= \boldsymbol{\omega}_{0,opt} - \gamma \boldsymbol{\omega}_\epsilon
\end{aligned} \tag{2.22}$$

where  $\boldsymbol{\omega}_{0,opt} \equiv \Sigma_0^{-1} \Delta \mathbf{f}$ ,  $\boldsymbol{\omega}_\epsilon = \Sigma_0^{-1} \mathbf{f}'$ ,  $\gamma \equiv \epsilon \frac{(\mathbf{f}'^T \boldsymbol{\omega}_{0,opt})}{1 + \epsilon (\mathbf{f}'^T \boldsymbol{\omega}_\epsilon)}$  and where we have made use of the analytical expression for the inverse of the covariance matrix in the presence of differential correlations (Eq. (33) in [47]). It is important to make two remarks at this point (i) Eq. (2.11) is a general solution regardless of the exact shape of  $\Sigma$ , therefore it can be used to find the optimal linear classifier without explicitly deriving it and (ii) by assuming  $\Sigma_1 = \Sigma_2$  we are basically assuming that  $\Sigma_{10} = \Sigma_{20}$  and  $\mathbf{f}'_1 = \mathbf{f}'_2 \equiv \mathbf{f}'$ , which implies losing generality. Interestingly, when the optimal classifier without sensory noise  $\boldsymbol{\omega}_{0,opt}$  is perpendicular to  $\mathbf{f}'$ , then the second term in Eq. (2.22) vanishes and  $\boldsymbol{\omega}_{opt} = \boldsymbol{\omega}_{0,opt}$ .

It is worth studying in more detail the case where  $\mathbf{f}(s)$  is linear on  $s$  because they will be useful for deriving analytical expressions for the network model

$$\mathbf{f}'(s) = \frac{\Delta \mathbf{f}}{\Delta s}. \tag{2.23}$$

In this case the optimal classifier becomes

$$\begin{aligned}\boldsymbol{\omega}_{opt} &= \boldsymbol{\omega}_{0,opt} - \gamma \boldsymbol{\omega}_\epsilon \\ &= \left(1 - \frac{\gamma}{\Delta s}\right) \boldsymbol{\omega}_{0,opt}.\end{aligned}\quad (2.24)$$

where  $\gamma$  is now  $\gamma = \epsilon \frac{\Delta s (\Delta \mathbf{f}^T \boldsymbol{\omega}_{0,opt})}{\Delta s^2 + \epsilon (\Delta \mathbf{f}^T \boldsymbol{\omega}_{0,opt})}$ . When  $\mathbf{f}'(s) = \frac{\Delta \mathbf{f}}{\Delta s}$ , then the optimal classifier is proportional to  $\boldsymbol{\omega}_{0,opt}$  and therefore  $\boldsymbol{\omega}_{0,opt}$  is itself a solution of the optimal classifier when introduced sensory noise to the system.

The DP in this case can be expressed as

$$\Sigma^{-1} = \Sigma_0^{-1} - \frac{\epsilon}{1 + \epsilon \mathbf{f}'^T \Sigma_0^{-1} \mathbf{f}'} \Sigma_0^{-1} \mathbf{f}' \mathbf{f}'^T \Sigma_0^{-1}, \quad (2.25)$$

and then, by defining  $d'_0 = \sqrt{\Delta \mathbf{f}^T \Sigma_0^{-1} \Delta \mathbf{f}}$

$$\begin{aligned}d' &= \sqrt{\Delta \mathbf{f}^T \left( \Sigma_0^{-1} - \frac{\epsilon}{1 + \epsilon \mathbf{f}'^T \Sigma_0^{-1} \mathbf{f}'} \Sigma_0^{-1} \mathbf{f}' \mathbf{f}'^T \Sigma_0^{-1} \right) \Delta \mathbf{f}} \\ &= \sqrt{\frac{d'_0}{1 + \frac{\epsilon}{\Delta s^2} d'_0}}.\end{aligned}\quad (2.26)$$

Introducing it back to Eq. (2.13) we get

$$DP = \Phi \left( \frac{1}{2} \sqrt{\frac{d'_0}{1 + \frac{\epsilon}{\Delta s^2} d'_0}} \right). \quad (2.27)$$

When  $N$  is very large the expression becomes

$$DP = \Phi \left( \frac{1}{2} \frac{\Delta s}{\sqrt{\epsilon}} \right). \quad (2.28)$$

In the presence of noise at the input stage the decoding performance of a linear classifier will be constraint by the upper bound  $DP = \Phi \left( \frac{1}{2} \frac{\Delta s}{\sqrt{\epsilon}} \right)$ . The further apart the stimuli are, the higher the upper bound, and the larger the sensory noise, the lower the upper bound.

## 2.4.2 Analytical vs fitted decoding performance

We evaluated the accuracy of the derived analytical expressions for the decoding performance (DP) of the different classifiers. In order to do so, we plotted the analytical decoding performances  $DP_{th}$  against the decoding performances of the fitted classifier  $DP_{cv}$  (trained and tested) and fitted the plot with a line. If our analytical expression is a good approximation for the amount of encoded information, the percentage of variance depicted by a cross-validated linear classifier that could be explained by our analytical expressions should be large. The metric that was used to evaluate the goodness-of-fit of  $DP_{th}$  in front  $DP_{cv}$  was the percentage of variance captured by  $DP_{th}$  on the total amount of variance depicted by  $DP_{cv}$ . In other words, we found the two main directions of the  $DP_{th}$ - $DP_{cv}$  cloud (PCA) and evaluated  $\lambda_1/(\lambda_1 + \lambda_2) \times 100$ , where  $\lambda_1$  and  $\lambda_2$  corresponded to the variance captured by the first and second principal components respectively. At this point is important to remark that in those datasets involving different difficulties (monkeys 1 and 4 and surrogate data) both the  $DP_{th}$  and the  $DP_{cv}$  have been calculated by conditioning on stimulus intensity (difficulty), as evaluating the  $DP_{th}$  involves assessing the covariance matrix of the population. Not controlling for signal strength could led to important mistakes when evaluating the common trial-by-trial modulation of the network units.

We first compared the analytical expression for the decoding performance of the optimal linear classifier (Eq. (2.14)) with the decoding performance of a linear discriminant analysis (LDA) (Fig. 2.3 and 2.4) and of a logistic regression (LR) (Fig. 2.4) on the test and on the train set (5-fold cross-validation). We identified  $DP_{cv}$  with the decoding performance of the LDA because it depicted the largest performance (Fig. 2.6). This comparison was made on the original datasets of all the experiments and surrogate data (see below; see Figs. 2.3, 2.13b and 2.4). We also compared the  $DP_{th}$  of the suboptimal classifiers with the  $DP_{cv}$  of the LDA. As expected, we found that  $DP_{cv}$  was better explained by the optimal classifier's  $DP_{th}$  than by the suboptimal expressions for  $DP_{th}$  (Fig. 2.2a). Similarly, we performed the same analysis after removing pairwise corre-

lations (shuffling procedure) and found that the correlation-blind classifier had the largest explanatory power for  $DP_{cv}$  in this case (Fig. 2.5).

## 2.4.3 Bootstrap analysis

### 2.4.3.1 General Method

We performed an analysis based on bootstrap to test the effect of fluctuations of SL and PP and other statistical quantities of neuronal responses, on DP and behavior. The bootstrap analysis is based on generating distributions of a set of different statistical quantities simultaneously, which allows fixing some of those variables and test the effect of others on DP and behavior. For each recording session or dataset, a bootstrap sample was built by randomly selecting with replacement  $M$  trials (being  $M$  the total number of trials in the session or dataset) and calculating from that sample all quantities of interest (typically SL, PP, mean pairwise correlations (MPC), global activity (GA), DP and behavior). The statistical quantities of the neuronal responses and behavior relevant for our study are

- Behavioral performance of the animal, denoted B (see below for each task's definition of performance).
- Analytical decoding performance  $DP_{th}$  (Eq. (2.14)).
- Cross-validated decoding performance of a trained and tested linear classifier (LDA),  $DP_{cv}$ .
- Selectivity length SL and projected precision PP (Eq. (2.14)).
- Mean pairwise correlations (MPC), defined as the average across the off-diagonal terms of the neuronal ensemble correlation matrix for a fixed stimulus condition.
- Global activity of the network (GA), defined as the mean neuronal activity across all neurons and trials.

For each bootstrap iteration  $i$  and quantity  $j$  we will obtain the value  $x_{ij} = \tilde{x}_j + \delta x_{ij}$ , where  $\tilde{x}_j$  is the median of the distribution over all bootstrap iterations and  $\delta x_{ij}$  is the fluctuation (perturbation) around the median for that particular bootstrap iteration  $i$  (Fig. 2.7). It is important to remark that in general the median  $\tilde{x}_j$  does not need to match on the original dataset's value if  $x_j$  is a nonlinear function of neuronal responses. The number of bootstrap iterations was  $10^4$ , except for the quantity  $DP_{cv}$ , where we used  $10^3$  iterations as computing  $DP_{cv}$  requires using optimization algorithms that are in general computationally more costly.

We evaluated whether in a recording session those subsets of trials (bootstrap iterations) that showed a high amount of neuronal information also showed a higher behavioral performance. To quantify this relationship, we evaluated the co-dependence between the vector of bootstrapped performances and the cross-validated decoding performances of a linear classifier  $p(\delta B, \delta DP_{cv})$  by means of the Pearson correlation coefficient (Figs. 2.8 and 2.14b).

### 2.4.3.2 Conditioning Method

A direct evaluation of the correlation between two quantities  $\delta x_i$  and  $\delta x_j$  ( $p(\delta x_i, \delta x_j)$ ) can produce misleading results because of the dependency of either  $\delta x_i$  or  $\delta x_j$  with a third quantity  $\delta x_k$ . In order to evaluate correlations  $p(\delta x_i, \delta x_j)$  accurately in this study, we developed a novel method based on systematically conditioning to the other quantities we thought could be affecting the dependency property. Therefore, we selectively chose a particular set of bootstrap iterations such that  $\delta x_k \simeq 0$  for all those quantities we were not interested in. For instance, there is a strong inverse dependency between  $\delta MPC$  and  $\delta PP$ , mainly because the projected precision decreases when the ensemble's noise (both individual and pairwise) is large (Eq. (2.14)). Therefore, if we calculate  $p(\delta DP_{cv}, \delta MPC)$  directly using all bootstrap iterations we could find spurious results due to the dependency of MPC with PP. The following conditionings were used in each case:

- $p(\delta DP_{cv}, \delta SL | \delta PP \simeq 0, \delta MPC \simeq 0)$  and

$$p(\delta B, \delta SL | \delta PP \simeq 0, \delta MPC \simeq 0)$$

- $p(\delta DP_{cv}, \delta PP | \delta SL \simeq 0, \delta MPC \simeq 0)$  and  $p(\delta B, \delta PP | \delta SL \simeq 0, \delta MPC \simeq 0)$
- $p(\delta DP_{cv}, \delta MPC | \delta SL \simeq 0, \delta PP \simeq 0)$  and  $p(\delta B, \delta MPC | \delta SL \simeq 0, \delta PP \simeq 0)$
- $p(\delta DP_{cv}, \delta GA | \delta SL \simeq 0, \delta PP \simeq 0)$  and  $p(\delta B, \delta GA | \delta SL \simeq 0, \delta PP \simeq 0)$

To evaluate  $p(\delta DP_{cv}, \delta x_i | \delta x_j \simeq 0, \delta x_k \simeq 0)$  (Figs. 2.8 and 2.13c) we used  $10^3$  bootstrap iterations. The conditionings  $\delta x_j$  and  $\delta x_k$  were obtained by taking all those bootstrapped iterations that generated  $x_j$  and  $x_k$  to be within the  $\pm 10^{th}$  percentile (with respect to the median;  $40^{th} - 60^{th}$  percentile) of their bootstrapped distributions. Because we used the central  $\pm 10\%$  percentile of the iterations for each quantity, conditioning to both  $x_j$  and  $x_k$  left us with  $\simeq 4\%$  of the total number of iterations, 40 iterations out of 1000, that nevertheless was sufficient to evaluate statistical significance of our results. We calculated  $p(\delta DP_{cv}, \delta x_i | \delta x_j \simeq 0, \delta x_k \simeq 0)$  using the percentage change in  $DP_{cv}$  as

$$\% \text{change} DP_{cv} = \frac{\langle DP_{cv}(x_i^{high}) \rangle - \langle DP_{cv}(x_i^{low}) \rangle}{\langle DP_{cv}(x_i^{low}) \rangle}, \quad (2.29)$$

where  $\langle DP_{cv}(x_i^{high}) \rangle$  was the mean  $DP_{cv}$  for all the bootstrap iterations that produced the set  $\{x_i^{high}\}$  and where  $\langle DP_{cv}(x_i^{low}) \rangle$  was the mean  $DP_{cv}$  for all the bootstrap iterations that produced the set  $\{x_i^{low}\}$ . The sets  $\{x_i^{high}\}$  and  $\{x_i^{low}\}$  corresponded to those bootstrap iterations that were above and below the median of the  $x_i$  distribution. To evaluate the dependency with the behavior  $p(\delta B, \delta x_i | \delta x_j \simeq 0, \delta x_k \simeq 0)$  (Figs. 2.10 and 2.13d) we used  $10^4$  bootstrap iterations. The conditionings  $\delta x_j \simeq 0$  and  $\delta x_k \simeq 0$  were obtained again by taking all those bootstrapped iterations that generated  $x_j$  and  $x_k$  to be within the  $\pm 10^{th}$

percentile of their bootstrapped distributions. Because we used the central  $\pm 10\%$  percentile of the iterations for each quantity (with respect to the median), conditioning to both  $x_j$  and  $x_k$  left us with  $\simeq 4\%$  of the total number of iterations, 400 iterations out of 10000. We calculated  $p(\delta B, \delta x_i | \delta x_j \simeq 0, \delta x_k \simeq 0)$  using the percentage change in  $B$  as

$$\%change DP_{cv} = \frac{\langle B(x_i^{high}) \rangle - \langle B(x_i^{low}) \rangle}{\langle B(x_i^{low}) \rangle}. \quad (2.30)$$

The conditioning method ensured that the obtained dependency between the behavioral performance or the decoding performance of a linear classifier and the different quantities  $x_i$  was not due to an alternative pathway through a third quantity  $x_k$ . In order to control for robustness we also characterized the relationship between perturbations on behavior and features of the neural code  $p(\delta B, \delta x_i | \delta x_j \simeq 0, \delta x_k \simeq 0)$  by the Pearson correlation and by a narrower conditioning ( $\pm 5^{th}$  percentile with respect to the median; see Fig. 2.12).

## 2.4.4 Network Model

In order to validate our set of hypothesis, we built a neuronal model where we simulated the activity of an interconnected population of neurons. In each trial a particular stimulus was presented to the network, which elicited a population pattern of activity (see below for the exact generative model). We also generated surrogate behavior by reading out optimally the population activity in a trial-by-trial basis. As for monkey 1 (Fig. 2.2a) and monkey 4 (Fig. 2.2c), the performance of the virtual agent was defined by the percentage of correctly classified responses. Three different stimulus intensities were presented to the network: easy, middle and difficult trials. The amount of information encoded by the network was defined by the percentage of correct classifications (DP) performed by the optimal classifier.

### 2.4.4.1 Generative model

We generated the surrogate population activity in four main steps (see Figs. 2.13a and 2.14a). First we evaluated the generative tuning curves and correlation and covariance matrices for the population of  $N$  neurons. Then for each trial we generated a preliminary population activity pattern by drawing a  $N$ -dimensional sample from a multivariate Gaussian distribution. The mean and covariance matrix for this multivariate Gaussian corresponded to the generative mean and covariance matrix defined below. We added an additional term that accounted for the trial-by-trial noise in the presented stimulus (differential correlations). This preliminary activity pattern was transformed by an homogeneous modulation and finally the population activity pattern for a given trial was obtained by applying a Poisson step. In the following we are describing in detail each of these steps.

Each neuron's firing rate was modulated as a function of the stimulus presented to the network. We considered two stimuli  $s_1 = 1$  and  $s_2 = -1$  that were presented at a given contrast value  $\Delta$ . The contrast controls the difficulty of the trial, with larger contrast leading to easier trials (in analogy with the coherence and the motion direction parameter for monkeys 1 and 4 respectively). The mean firing rate (tuning curve) for each neuron was

$$\mu_i(s, \Delta) = m_{1i}s\Delta + m_{0i}. \quad (2.31)$$

Here,  $m_{1i}$  is the intrinsic selectivity for each neuron, and it was drawn from a normal distribution centered at the origin with a standard deviation of  $\sigma_{m_1} = 1.3$ ;  $m_{0i}$  is the baseline firing rate for each neuron in the absence of a stimulus, and it was drawn from a gamma distribution with shape and scale parameter 20 and 1 respectively. Note that from now on we will refer to the tuning curve of neuron  $i$  as  $\mu_i(s, \Delta)$  and  $f_i(s, \Delta)$  indistinctly. The contrast parameter consisted on a set of three values  $\Delta = \{0.16, 0.32, 0.48\}$ . This way of defining the tuning curve is convenient because the ensemble's total information is going to depend on basically two independent sources, the strength of the stimulus presented

to the network  $\Delta$ , and the intrinsic selectivity of each neuron in the ensemble  $m_{1i}$  (see Fig. 2.14b for the tuning curve of an example simulated neuron). As we simulated activity for 1 second, from now on I will refer to spike count, firing rate and neuronal activity indistinctly.

Neuronal activity typically shows a shared component of variance when presented with identical stimuli, the so called ‘noise correlations’. In our model, we considered limited-range correlations [46, 51, 55–57, 113], also known as exponential correlations. The generative exponential correlations were defined as

$$\rho_{ij}^{gen} = A \exp\left(-\frac{|m_{1i} - m_{1j}|}{\tau}\right) + C, \quad (2.32)$$

for  $i \neq j$ , and where we used  $A = 0.1$ ,  $C = 0$  and  $\tau = 1$ . The diagonal terms were set to  $\rho_{ii}^{gen} = 1$ . *In vivo* it is found that similarly tuned neurons show larger pairwise correlations than dissimilarly tuned neurons, and the exponential model in Eq. (2.32) naturally account for this observed feature.

Once the matrix of pairwise correlations is characterized, the generative covariance matrix was defined as

$$\Sigma_{ij}^{gen}(s, \Delta) = \sigma_i(s, \Delta)\sigma_j(s, \Delta)\rho_{ij}, \quad (2.33)$$

where  $\sigma_i^2(s, \Delta)$  is the variability of each neuron’s activity across trials. Here we used  $\sigma_i^2(s, \Delta) = \eta\mu_i(s, \Delta)$  ( $\eta = 0.5$ ) so that we obtained biologically realistic Fano Factors for the final surrogate activity. As our generative model includes differential correlations, global multiplicative and additive modulation and a Poisson step (see below), it is important to remark that  $\Sigma_{ij}^{gen}(s, \Delta)$  is not equivalent to the full model’s covariance matrix  $\Sigma_{ij}(s, \Delta)$ .

After  $\Sigma_{ij}^{gen}(s, \Delta)$  and  $\boldsymbol{\mu}(s, \Delta)$  were determined, we generated for trial  $j$  a preliminary population activity pattern by drawing a sample from a multivariate Gaussian distribution with mean and covariance matrix  $\Sigma_{ij}^{gen}(s, \Delta)$  and  $\boldsymbol{\mu}(s, \Delta)$  respectively. On each trial a particular  $s_j$  and  $\Delta_j$  were chosen pseudorandomly so that the set of presented stimuli and trial difficulties were counterbalanced. For each dataset (10 surrogate dataset)

we produce 100 trials per stimulus and difficulty ( $M = 600$  per dataset). We also included an additional term that accounted for the trial-by-trial noise on the presented stimulus (differential correlations). On each trial the preliminary activity pattern was

$$\mathbf{r}'_j = \boldsymbol{\mu}(s_j, \Delta_j) + M(s_j, \Delta_j)\mathbf{z}_j + \boldsymbol{\mu}'\delta s_j. \quad (2.34)$$

Here  $\Sigma_{ij}^{gen} = MM^T$ ,  $\mathbf{z}$  is an  $N$ -dimensional vector where each entry  $i$  has been drawn from a zero-mean and unit-variance Gaussian distribution,  $\boldsymbol{\mu}'$  is the derivative of the population's tuning curve ( $f'$ ) and  $\delta s$  is the sensory noise term, which will be drawn from a Gaussian distribution with zero mean and  $\sigma_{\delta s}^2 = \epsilon$ . The last term will produce differential correlations on the activity of the population. It is important to remark that in our particular model  $\boldsymbol{\mu}' = \mathbf{m}_1$ . Once the preliminary activity pattern was calculated, we applied a global modulation and a Poisson step to obtain the spike counts for the population

$$\mathbf{n}_j \sim Poisson(g_j\mathbf{r}'_j + g_j), \quad (2.35)$$

where  $g_j$  was drawn from a Gamma distribution with scale and shape parameter  $\theta_g = 10^{-5}$  and  $k_g = 10^5$  respectively ( $\langle g_j \rangle = 1$  and  $\sigma_g^2 = 10^{-5}$ ). We included this global modulation on our model because it has also been reported previously that networks in vivo undergo global fluctuations on their activity on a short and long time scale [42–44, 100]. The surrogate datasets depicted biologically-realistic distributions of Fano factors and pairwise correlations (Figs. 2.14c,d).

We also modeled the behavior of a virtual decision-making agent in a trial-by-trial basis by reading-out optimally (for each stimulus intensity  $\Delta$ ) the entire population of neurons ( $N = 1000$ ). The optimal classifier was calculated analytically in order not to overfit our data due to the limit size effect (200 trials per difficulty). By introducing Eq. (2.31) on the expression for the optimal classifier

$$\boldsymbol{\omega}_{opt} = 2\Sigma^{-1}\Delta\mathbf{m}_1, \quad (2.36)$$

where the analytical expression for the covariance matrix is

$$\begin{aligned} \Sigma_{ij} = & (\sigma_g^2 + 1)(\Sigma_{ij}^{gen} + \epsilon\mu'_i\mu'_j) + \\ & + \sigma_g^2 ((\mu'_i + 1)(\mu'_j + 1)) + \left( \sqrt{\mu'_i + 1} \sqrt{\mu'_j + 1} \right) \delta_{ij}. \end{aligned} \quad (2.37)$$

On each trial a decision variable will be generated

$$d_j = \boldsymbol{\omega}_{opt}^T \mathbf{n}_j + \omega_0, \quad (2.38)$$

where  $\omega_0 = -\Delta \mathbf{m}_1^T \Sigma^{-1} \mathbf{m}_0$ . The choice of the virtual agent on each trial will be the sign of the decision variable  $c_j = \text{sign}(d_j)$ . As for monkeys 1 and 4, the behavioral performance of the surrogate agent will be characterized by the number of correct classifications over the total number of trials. Under this model, a weak correlation can be created between the perturbations on the amount of encoded information and the behavioral performance of the virtual agent (Fig. 2.14f).

#### 2.4.4.2 Information encoded by the model

In this section we will present some of the analytical expressions for the decoding performance (DP) for the surrogate model. In our model  $\mathbf{f}' = \frac{\Delta \mathbf{f}}{\Delta s}$  and therefore Eq. (2.27) can be expressed as

$$d' = \sqrt{\frac{4\Delta^2 \left( \sum_{i=1}^N m_{1i}^2 \right) PP_0^2}{1 + \epsilon \left( \sum_{i=1}^N m_{1i}^2 \right) PP_0^2}}, \quad (2.39)$$

where  $PP_0^2$  is the projected precision without differential correlations. If  $N$  is large, the above equation can be rewritten as

$$d' = \sqrt{\frac{4\Delta^2 N \sigma_{m_1}^2}{PP_0^2 + \epsilon N \sigma_{m_1}^2}}. \quad (2.40)$$

From this expression it is easy to show that when  $N \rightarrow \infty$

$$d' = \frac{2\Delta}{\sqrt{\epsilon}}, \quad (2.41)$$

or in other words

$$DP = \Phi \left( \frac{\Delta}{\sqrt{\epsilon}} \right). \quad (2.42)$$

When the input signal is noisy on a trial-by-trial basis, the amount of information that can be extracted from very large ensemble sizes does not saturate to  $DP = 1.0$ , but it reaches an asymptote. So to speak, the asymptote will be fully determined by the signal to noise ratio of the input signal (see Fig. 2.14e).

### 2.4.4.3 Analysis of surrogate data

A total of 10 recording sessions of surrogate data were generated. Each dataset consisted of 200 trials per stimulus intensity (100 each stimulus) and three stimulus intensity. The network size was set to 1000 neurons.

To test how accurate our analytical expression was, we compared the  $DP_{th}$  with the cross-validated  $DP_{cv}$  of a linear classifier (LDA) in the same way as we did for the *in vivo* recordings. We randomly selected sub-ensembles of a particular ensemble size, signal intensity and dataset and evaluated  $DP_{th}$  and  $DP_{cv}$ . For each ensemble size and stimulus intensity we repeated this process 20 times and fitted the relationship between  $DP_{th}$  and  $DP_{cv}$  as for the *in vivo* datasets (Fig. 2.13b).

The perturbation analysis based on bootstrap was performed as for the *in vivo* datasets. For each bootstrap iteration we selected a particular sub-ensemble of a given size (2, 4, 6, 8 and 10 units) and calculated the following quantities: analytical decoding performance ( $DP_{th}$ ), cross-validated decoding performance of a trained and tested linear classifier ( $DP_{cv}$ ), selectivity length (SL), projected precision (PP), mean pairwise correlations (MPC), global activity (GA) and behavioral performance (B). Both the theoretical and the trained-and-tested decoding performance were calculated on inferring what was the stimulus presented to the network (either  $s_1$  or  $s_2$ ) from the activity pattern of the population. This process was repeated 20 times for each dataset and stimulus intensity. The reported dependencies (% change and Pearson correlation) between the different

features of the neural code and the encoded information and the behavior were the mean across subsampling repetitions, stimulus intensity and dataset. Significance, as for the in vivo recordings, was calculated with a two-sided Wilcoxon signed-rank test, where we tested if the set of independently obtained values was significantly above or below zero. As the subsampling procedure generated non-independent samples, 30 independent samples (10 recording sessions  $\times$  3 stimulus intensities) were used to test significance.

## **2.4.5 Experimental Methods**

Our theory was tested on three different datasets (four monkeys) with simultaneously recorded units (from 2 to  $\sim$ 50 neurons): middle temporal (MT) neurons in a monkey performing a coarse random dot motion task [35] (monkey 1); lateral prefrontal cortex (LPFC, area 8A) neurons recorded in two monkeys performing an attentional task [105] (monkeys 2 and 3); and MT neurons in a monkey performing a fine motion discrimination task (monkey 4). The details for each experiment are explained below.

### **2.4.5.1 Coarse random dot motion task recorded in MT**

This dataset has been previously described in [35]. It is freely available at the Neural Signal Archive (<http://www.neuralsignal.org>, nsa2004.2). Here we provide a brief summary.

#### **2.4.5.1.1 Subjects and recordings**

One adult macaque monkey (*Macacca mulatta*) was used in this study. Following several months of training on a coarse motion discrimination task a steel cylinder was surgically implanted over the occipital cortex. MT was identified based on its location within the temporal sulcus, directionally responsive neurons and its characteristic topography. Pairs of single neurons were recorded in MT during this experiment. Training was accomplished by operant conditioning techniques using fluids as a

positive reward. The animals were maintained in accordance with guidelines set by the U.S. Department of Health and Human Services (NIH) Guide for the Care and Use of Laboratory Animals. Electrophysiological recordings were made with tungsten microelectrodes (Micro Probe, Potomac, MD). Action potentials from multiple single neurons were discriminated using an on-line spike sorting system. In total 82 neurons were recorded (41 pairs).

#### **2.4.5.1.2 Experimental task**

The experiment consisted in a monkey performing a coarse discrimination task, a two-alternative forced-choice discrimination of a noisy motion direction signal. On each trial a random stimulus was presented for 2 seconds covering the receptive fields of both neurons. Two motion directions were defined, preferred and null. Preferred direction was set as the direction of motion to which neurons were more selective to, and null direction was defined as the opposite to the preferred direction. Because neurons in general were not tuned to the same motion direction, preferred direction in each recording session was defined as a compromise between the two directions that elicited a largest firing rate for each neuron. The motion signal consisted in a set of dynamic random patterns in which unidirectional motion signal was interspersed among random motion noise. When all dots moved randomly without a defined motion direction the stimulus was pure noise. The coherence parameter in this case was 0%. When all dots moved coherently in the same direction (either preferred or null direction) the stimulus was fully informative and the coherence parameter was then 100%. Several intermediate cases between these two extremes were used by generating patterns with a particular percentage of dots moving coherently in one direction. In each trial the monkey was presented randomly with either preferred or null direction of motion of a particular strength or coherence. The monkey's task was to discriminate correctly the direction of motion of the dots. The more dots were moving coherently, the easier was the task for the animal (Fig. 2.9).

The trial flow was the following: each trial started with the appear-

ance of a fixation point. After the monkey held its gaze for 300 ms, the stimulus was presented. It was required to hold fixation during stimulus presentation so that stimulus was presented on the receptive field of the recorded neurons. After 2 seconds of stimulus presentation the random dot pattern and the fixation point disappeared and two visual targets appeared on the screen, each one corresponding to one of the possible directions of motions (preferred or null) (see Fig. 2.2a). The monkey made a saccadic eye movement to one of the targets to report its perceived direction of motion. On 0% coherence trials the monkey was rewarded 50% of the times.

### 2.4.5.1.3 Neuronal data analysis

This experiment [35] consisted in 41 recording sessions where pairs of single-units were simultaneously recorded in the MT/V5 region of a monkey. It is important to note that only sessions with variable seed were used in this study.

In each dataset, the stimulus strength was controlled by the coherence parameter, the percentage of coherent dots moving in the same direction. Therefore, in order to control for stimulus condition, we divided each recording session according to the coherence strength in each trial. In all recording sessions, trials belonging to the 0% coherence condition were discarded because behavioral performance was not defined for this particular set of trials. We also discarded all coherences that elicited a mean behavioral performance larger than 98% to avoid ceiling effects (25.6%, 51.2% and 99.9% coherence). This is because in very easy trials the bootstrap distributions for behavior ( $B$ ) have a very low variability and therefore calculating  $p(\delta B, \delta x_i | \delta x_j \simeq 0, \delta x_k \simeq 0)$  is highly biased or even undefined. The criterion for discarding easy coherences (98% percentage correct) was grounded on the fact that when fitting a sigmoid function (psychometric curve) with a cumulative Gaussian, 0.98 is achieved for two standard deviations  $\Phi(2\sigma) = 0.98$ . After splitting recording sessions by coherence and discarding 0% and very easy coherences we obtained 187 independent sub-datasets. The mean number of trials per dataset was

75, ranging from 30 to 231 trials. For the analysis, we used a trial-by-trial population activity vector whose entries corresponded to the spike count of the whole trial, (2 seconds), for each neuron.

In each sub-dataset, we calculated the bootstrap distribution and the original observed value of the following quantities: analytical decoding performance ( $DP_{th}$ ), cross-validated decoding performance of a trained and tested linear classifier ( $DP_{cv}$ ), selectivity length (SL), projected precision (PP), mean pairwise correlations (MPC), global activity state (GA) and behavioral performance ( $B$ ). The behavioral performance was defined as the fraction of correct choices of the monkey over the whole set of trials. Both the theoretical and the trained-and-tested decoding performance were calculated on inferring what was the motion direction (preferred or null) presented to the monkey on a trial-by-trial basis from the simultaneously recorded activity of the neuronal pair.

Significance for a particular dependency  $p(\delta x_i, \delta x_j | \delta x_k \simeq 0, \dots, \delta x_r \simeq 0)$  on the whole experiment was calculated with a two-sided Wilcoxon signed-rank test where we tested if the set of 187 independently obtained values (187 independent sub-datasets) was significantly above or below zero.

#### **2.4.5.2 Attentional task recorded in LPFC 8a**

This dataset is described in detail in [105]. Here we provide a brief summary.

##### **2.4.5.2.1 Subjects and recordings**

Two male monkeys (*Macaca fascicularis*) both aged 6 years old (Monkey “F”, 5.8 Kg; Monkey “JL”, 7.5 Kg) were used in this study. In each monkey a 96-channel “Utah” multielectrode arrays (Blackrock Microsystems, Utah, USA) was chronically implanted in the left caudal lateral prefrontal cortex. The multielectrode array was inserted on the prearcuate convexity posterior to the caudal end of the principal sulcus and anterior to the arcuate sulcus, a region cytoarchitectonically known as area 8A. The extracted spikes and associated waveforms were sorted offline using both manual

and semi-automatic techniques (Offline sorter, Plexon Inc., TX, USA). All procedures were in accordance with the Canadian Council of Animal Care guidelines and were preapproved by the McGill University Animal Care Committee. None of the animals were sacrificed for the purpose of this study. Neuronal recordings included both single and multiunits. The mean number of simultaneously recorded units across recording sessions was 56 for monkey “JL” and 52 for monkey “F”.

#### **2.4.5.2.2 Experimental task**

The monkeys were trained to covertly sustain attention to one of four Gabor stimuli (target) presented on a screen while ignoring the other three Gabor stimuli (distractors) (Fig. 2.2b). At the beginning of the trial a cue indicated which of the four Gabor stimuli was the target (cue period, 363 ms). Three different trial types during the attentional period (attentional period, 585 - 1755 ms) were used in each session, which were randomly interleaved from one trial to the next. In “Target” trials the target changed orientation (90° rotation) after a variable time delay interval, indicating the monkey to make a saccade towards the target within a 250 milliseconds (msec) time window to get a reward (fruit juice). In “Distractor” trials the orientation change occurred in the Gabor opposite to the cued location. Monkeys had to inhibit saccading to the distractor and maintain fixation in order to be considered as a correct trial. In “Target + Distractor” trials two simultaneous orientation changes co-occurred in the target and in the distractor opposite to the target. In this case the monkey had to make a saccade towards the target and not to the distractor to get a reward. On average, the monkeys completed ca. 1000 trials per session and performed well above chance (Fig. 2.9). Only correct “Target” trials were used in the analysis. The inclusion of the other types of trials was important for a correct and unbiased performance of the animals.

The full dataset consisted on 7 sessions per monkey. In all of them we orally administered MPH or placebo to the monkeys 30 minutes before the beginning of a session. In drug sessions we diluted 5 mg of MPH Hydrochloride (Ritalin©, Novartis, Switzerland) into 5 mL of concentrated

fruit juice vehicle and gave the juice to the monkey. A block of three drug sessions (three consecutive days) with the same dose was preceded and followed by a block of two consecutive placebo sessions. In placebo sessions, the experimental procedures were identical but the concentrated fruit juice administered before the beginning of the session did not contain the drug (4 placebo and 3 MPH sessions per monkey). Even though drug administration on the subjects was not strictly necessary for the hypothesis of the current study, it did not have an influence in the herein presented results.

#### **2.4.5.2.3 Neuronal data analysis**

The experiment was performed on two adult monkeys (“F” and “JL”) and 7 recordings sessions per monkey were produced. In order to make the task a binary classification task we considered two main approaches for this experiment: decoding the monkey’s allocation of attention on the horizontal axis and on the vertical axis. Only correct “Target” trials were used in this analysis. The mean number of trials per recording session was 209, ranging from 172 the least populated to 224 the most populated recording session for monkey “JL” and mean number of 221 trials, ranging from 198 the least populated to 246 trials the most populated recording session for monkey “F”. The mean number of simultaneously recorded units was 56, ranging from 53 to 61 for monkey “JL” and 52, ranging from 43 to 66 for monkey “F”. Neuronal recordings included both single and multiunits. Because very low firing rate units precluded any reliable statistical analysis, we excluded units firing below  $1Hz$  for all the subsequent analysis. After filtering out low firing rate units the mean number of simultaneously recorded units was 52 (ranging from 50 to 54) for monkey “JL” and 48 (ranging from 40 to 59) for monkey “F”.

The attentional period started after the cue presentation period, when all four Gabor stimuli appeared in the screen, and it ended after a random amount of time, when the target stimulus changed its orientation by  $90^\circ$ . The shortest attentional period was set to 585 ms and therefore we defined a fixed attentional time window of 585 starting right after the end

of the cue period. In this way, all trials and all units firing rate was calculated using the same time window, a crucial quantity for the amount of information a single neuron can encode.

To maximize the statistical power of our analysis, we created a larger number of independent datasets as follows. In each recording session, we selected a particular subensemble of units and calculated the following quantities: analytical decoding performance ( $DP_{th}$ ), cross-validated decoding performance of a trained and tested linear classifier ( $DP_{cv}$ ), selectivity length (SL), projected precision (PP), mean pairwise correlations (MPC), global activity (GA) and behavioral performance (B). Only “Target” correct trials were used in this experiment and thus we take behavioral performance (B) as the mean reaction time of a particular set of trials (either the original dataset or a bootstrap iteration) (see Figs. 2.9b,c for monkeys 2 and 3 respectively). Both the theoretical and the trained-and-tested decoding performance were calculated on inferring what was the monkey’s location of attention a trial-by-trial basis from the simultaneously recorded activity pattern. It is important to note that to make it a binary task we considered two different analysis: decoding top vs bottom allocation of attention (vertical axis) and left vs right allocation of attention (horizontal axis). The reported  $DP_{cv}$  corresponded to the mean across the vertical and horizontal axis in all cases and subjects.

We randomly selected ensembles of 2, 4, 6, 8 and 10 units. For each ensemble size, we selected  $N$  non-overlapping groups of units, so that we could increase the number of independent sub-datasets when assessing significance of our results. This process was repeated 5 times. For each ensemble size, we sub-selected 14, 7, 4, 3 and 2 non-overlapping subensembles of size 2, 4, 6, 8, 10 respectively. For a particular randomly constructed subensemble, the reported dependency relationship  $p(\delta x_i, \delta x_j | \delta x_k \simeq 0, \dots, \delta x_r \simeq 0)$  was evaluated as the mean across the 5 iterations and axis of classification (vertical and horizontal). For each ensemble size and monkey, the reported dependency  $p(\delta x_i, \delta x_j | \delta x_k \simeq 0, \dots, \delta x_r \simeq 0)$  was the mean across recording sessions and non-overlapping subensembles of units. We used 98 (14 non-overlapping subensembles  $\times$  7 recording sessions), 49, 28, 21, and 14 independent values to as-

sess the mean and test significance for ensemble sizes 2, 4, 6, 8 and 10 units respectively. Significance was calculated with a two-sided Wilcoxon signed-rank test where we tested if the set of independently obtained values was significantly above or below zero.

### **2.4.5.3 Fine motion discrimination task recorded in MT**

This dataset was recorded for this experiment.

#### **2.4.5.3.1 Subjects and recordings**

One adult macaque monkey (*Macaca mulatta*) was used in this experiment. Following several months of training on a fine motion discrimination task, it was surgically implanted with a head holding device, scleral coils for measuring eye movements, and a recording chamber/grid. The animal was trained using operant conditioning techniques with water or juice as a positive reward. Eye movements were measured and controlled at all times. Neuronal activity in MT/V5 was recorded while the monkey performed the task with a linear electrode array (V-probes, Plexon Inc). Spike waveforms were acquired by a BlackRock Cerebus system.

We started the recording session by inserting V-probes into MT/V5, and allow the probes to settle for  $\sim 30$  minutes. We then performed a standard battery of tests to map the receptive field and to measure the direction, speed, size, and depth tuning of the recorded neurons. These measurements are used to determine the location of the moving object such that it will adequately stimulate the receptive fields of the neurons under study. Note, however, that will not tailor the object motion stimulus to the preferred directions and speeds of the recorded neurons. This will facilitate recording from multiple neurons (with overlapping receptive fields) simultaneously. A total of 75 units were registered out of 3 sessions. The mean number of simultaneously registered neurons was 25, ranging from 24 to 26. Neuronal recordings included both single and multiunits.

### **2.4.5.3.2 Experimental task**

The monkey was trained to perform a fine motion discrimination task. In this task, an object (composed of random dots at 100% coherence; see previous section) moves upward in the visual field with either a rightward or leftward component, and the animal's task is to report by making a saccade to one of two targets whether the perceived motion was upward-rightwards or upward-leftwards. Object motion is presented in one visual hemi-field and the display follows a Gaussian velocity profile with a duration of 2015 ms and a standard deviation of 167 ms. Once the monkey gazes the fixation point the target object appears and starts the movement in the display. The moving object is a set of stereoscopic dynamic dots 100% coherent that move upward-leftwards or upward-rightwards with a particular angle. The experimental protocol involves 7 directions of object motion:  $-12^\circ$ ,  $-6^\circ$ ,  $-3^\circ$ ,  $0^\circ$ ,  $3^\circ$ ,  $6^\circ$  and  $12^\circ$ , where negative is leftwards and positive is rightwards with respect to the vertical. After 2015 ms of stimulus presentation two dots appear in the horizontal axis and the monkey has to report the perceived direction of motion by making a saccade to either the left (perceived leftwards) or to the right (perceived rightwards) (Fig. 2.2c). It is important to remark that because stimulus presentations were made at 100% coherence of the random dot pattern, the difficulty of the task was controlled by the angle of motion with respect to the vertical line (see Fig. 2.9d). The mean number of trials per session was 742, ranging from 735 to 756. In each recording session the same amount of trials were presented for each direction of motion.

### **2.4.5.3.3 Neuronal data analysis**

In each recording session the direction of motion controlled the difficulty of the trial and therefore this parameter was considered the strength of the stimulus signal. In order to control for stimulus condition, we divided each recording session according to the motion direction each trial depicted. The task of our classifier is to correctly classify a trial as belonging to rightwards or leftwards movement and therefore each recording session was split in 4 independent datasets ( $12^\circ$ ,  $6^\circ$ ,  $3^\circ$  and  $0^\circ$  stimulus

strength). As in the coarse discrimination task, we discarded all the stimulus strengths that were either ambiguous ( $0^\circ$  motion direction; behavioral performance is not defined) or too easy ( $12^\circ$ ; mean behavioral performance  $\geq 0.98$  and therefore  $p(\delta B, \delta x_i | \delta x_j \simeq 0, \delta x_k \simeq 0)$  highly biased or undefined). The analysis was performed on 6 independent datasets (3 recording sessions  $\times$  2 absolute motion directions). For the neuronal activity we used the firing rate evaluated in a time window that expanded 2 standard deviations forward and backwards from the peak of the Gaussian velocity profile of the stimulus (668 ms). Choosing 4 standard deviations centered in the presented stimulus peak ensured us we were using more than the 95% of the area of the Gaussian velocity profile.

In order to create a larger number of independent sub-datasets we proceeded as follows. In each dataset we selected a particular subensemble and calculated the following quantities: analytical decoding performance ( $DP_{th}$ ), cross-validated decoding performance of a trained and tested linear classifier ( $DP_{cv}$ ), selectivity length (SL), projected precision (PP), mean pairwise correlations (MPC), global activity (GA) and behavioral performance (B). The behavioral performance was defined as the fraction of correct choices of the monkey over the whole set of trials for that particular dataset. Both the theoretical and the trained-and-tested decoding performance were calculated on inferring what was the motion direction (leftward or rightward) presented to the monkey on a trial-by-trial basis from the activity pattern. We randomly selected ensembles of size 2, 4, 6, 8 and 10 units. For each ensemble size we selected  $N$  non-overlapping groups of units, so that we could increase the number of independent datasets when assessing significance of our results. This process was repeated 20 times. For each ensemble size we subselected 10, 5, 3, 2 and 2 non-overlapping subensembles of size 2, 4, 6, 8, 10 respectively (note that the larger the subensemble size, the lower the number of non-overlapping subensembles can be chosen). For a particular randomly constructed subensemble the reported dependency relationship  $p(\delta x_i, \delta x_j | \delta x_k \simeq 0, \dots, \delta x_r \simeq 0)$  was evaluated as the mean across the 20 iterations. For each ensemble size the reported dependency  $p(\delta x_i, \delta x_j | \delta x_k \simeq 0, \dots, \delta x_r \simeq 0)$  was the mean across recording sessions,

stimulus strength and non-overlapping subensembles of units. We used 60 (10 non-overlapping subensembles  $\times$  3 independent datasets  $\times$  2 stimulus strength), 30, 18, 12, and 12 independent values to assess the mean and test significance for ensemble sizes 2, 4, 6, 8 and 10 units respectively. Significance was calculated with a two-sided Wilcoxon signed-rank test where we tested if the set of independently obtained values was significantly above or below zero.



## Chapter 3

# INTEGRATION OF PRIOR WITH SENSORY INFORMATION

*The study presented in this chapter corresponds to the published article Nogueira, R. et al. Lateral orbitofrontal cortex anticipates choices and integrates prior with current information. Nature Communications 8, 14823 (2017) (see [15]).*

**Adaptive behavior requires integrating prior with current information to anticipate upcoming events. Brain structures related to this computation should bring relevant signals from the recent past into the present. Here we report that rats can integrate the most recent prior information with sensory information, thereby improving behavior on a perceptual decision-making task with outcome-dependent past trial history. We find that anticipatory signals in the orbitofrontal cortex about upcoming choice increase over time and are even present before stimulus onset. These neuronal signals also represent the stimulus and relevant second-order combinations**

**of past state variables. The encoding of choice, stimulus and second-order past state variables resides, up to movement onset, in overlapping populations. The neuronal representation of choice before stimulus onset and its buildup once the stimulus is presented suggest that orbitofrontal cortex plays a role in transforming immediate prior and stimulus information into choices using a compact state-space representation.**

### **3.1 Introduction**

Making a decision in real life requires the integration of preceding and current information to adaptively guide behavior [74, 114]. Previous work has investigated the neuronal regions responsible for achieving this goal by using experimental paradigms where the sequence of external events, or history, flows independently of the choices of the actor [27, 67, 114]. In many cases, however, choices of an actor can influence future external events, and so to speak, change the course of history. Relatively less work has been devoted to the study of tasks in which recent past information matters for the current choice and immediately previous choices affect the upcoming states of the world [74, 75, 115–118].

The orbitofrontal cortex (OFC), like other regions in the prefrontal cortex, is thought to play an important role in adaptive and goal-directed behavior [78, 79, 82, 85, 97, 118, 119]. Previous single-neuron accounts have demonstrated that OFC encodes a myriad of variables that are relevant for behavior in decision-making [82], such as primary rewards and secondary cues that predict them [87, 120], values of offered and chosen goods [80, 81, 121], choices and responses [81, 91–94], expected outcomes [86] and stimulus type [122], while human brain imaging studies have corroborated and largely extended these results [78, 88–90, 123]. However, in contrast to other prefrontal and parietal brain areas [67, 76, 77], the OFC displays relatively weak choice-related signals [80, 81, 92, 93]. Further, neuronal signals anticipating upcoming choices before stimulus onset have not been described, except in a single report in monkeys

[94]. This has led to the predominant view that OFC is not responsible for action initiation and selection [91, 118, 121]. Here, in contrast, we hypothesize that OFC plays a central role in decision-making, first, by representing the central latent variables of the task (state-space) and, second, by combining the most recent past with current stimulus information. We hypothesize also that this combination of information happens through a compact representation of the task's state-space, that is, by representing predominantly the variables of the immediate past that are critical to perform the task. We support this hypothesis through our findings that OFC (1) represents choice initiation and choice selection even before sensory evidence is available, (2) encodes the state-space determined by just the previous trial (here called immediate prior or immediate past information), (3) integrates the immediate prior information with current sensory evidence and (4) promotes filtering out behaviorally irrelevant variables. In this study we use an outcome-coupled perceptual decision-making task that requires integrating prior information from the previous trial with an ambiguous stimulus. This task is designed to maximize the chances of revealing choice initiation and choice selection signals that integrate both immediate prior and current information. Rats efficiently solve this task by using the relevant second-order combination of previous choice and reward and combining this most recent prior information with currently available information of a perceptually challenging stimulus. On the basis of single-neurons and simultaneously recorded neuronal ensembles in the lateral OFC (lOFC), we find a buildup of choice-related signals across time; critically, upcoming choice can be traced back to a period of time before stimulus onset. Overlapping neuronal populations encode choice, immediate prior and stimulus information stably over time up to movement onset. These neuronal populations represent behaviorally relevant variables in a task-structure dependent way. For example, information about the immediate past cease to be represented once such variables become behaviorally irrelevant due to a change in the task structure. Similarly, in the main task, the coexistence of choice-related and latent variables within the same neuronal circuits enables lOFC to play an important role in integrating prior with stimulus information to aid

choice formation using a compact state-space representation. Our results are consistent with the hypotheses that OFC plays a role in the temporal credit-assignment problem, the problem of correctly associating an action with a reward delayed in time [78, 118] and in representing latent states [97]. Furthermore, our work adds the view that IOFC might play a central role in decision-making by integrating immediate prior information with current information through a refined encoding of the state-space in the task.

## 3.2 Results

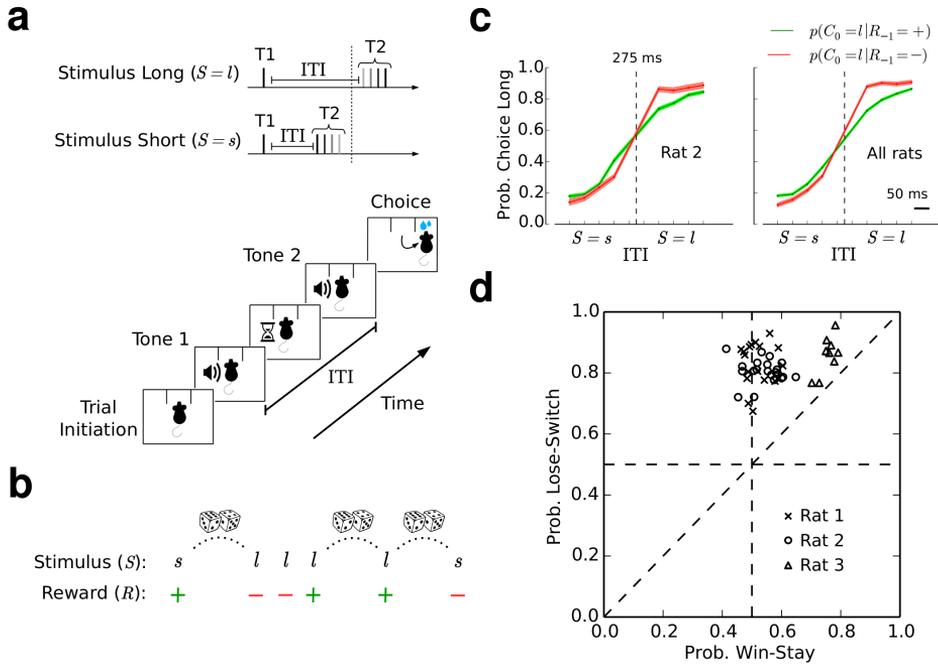
### 3.2.1 Animals use task-contingencies to improve performance

Rats performed a perceptual decision-making task (Fig. 3.1a), which in each trial consisted in classifying an inter-tone time interval (ITI), as short ( $S = s$ ) or long ( $S = l$ ). The rats self-initiated the trial with a nose poke in the central socket, after which they had to hold the position until the ITI had completely elapsed. A correct response was defined as poking into the left socket if the stimulus was short, and into the right socket if the stimulus was long, after which the rat was rewarded with water. A stimulus was considered difficult if the inter-tone interval was close to the category boundary, and easy otherwise (Fig. 3.1a). Importantly, in our task the choices of the animal influenced the history of future events. Specifically, in the trial following a correct response ( $R = +$ ), the ITI was drawn uniformly at random from eight possible values, while in trials following an incorrect response ( $R = -$ ), the stimulus was repeated (Fig. 3.1b). This sequence created a rich environment, whereby in many trials the ITIs were not drawn randomly. Rather, the environment was formally described as an outcome-coupled hidden Markov chain, that is, a Markov chain in which the sequence of trials is coupled with the outcomes of the animals' choices. The Markov chain was hidden because of two reasons (Fig. 3.2): first, due to potential limits in memory and atten-

tion, we did not consider previous trials as fully known; and second, the stimulus was not fully visible at any trial, especially so in the most difficult trials (Fig. 3.1a). The combination of independent trials after correct responses and fully dependent trials after incorrect responses allowed us to distinguish signals from the past from those that anticipated upcoming events, as discussed in the next section.

From an ideal observer's perspective, there is critical information that the animal should monitor to perform the task efficiently. The outcome in the previous trial,  $R_{-1}$ , determines whether the stimulus in the next trial will be repeated or drawn randomly: if the previous outcome was incorrect ( $R_{-1} = -1$ ), then the stimulus will be repeated in the next trial, while if the previous trial was correct ( $R_{-1} = +1$ ), then the next stimulus will be randomly drawn. Therefore, if the animal tracks the outcome  $R_{-1}$ , its behavior will improve because it could often anticipate the stimulus. In fact, the three rats learnt this task contingency by using the previous outcome to improve their behavior (Fig. 3.1c; individual rats and fits shown in Fig. 3.3). First, all animals featured a psychometric curve (computed after correct trials) with a larger fraction of correct responses for easy than for difficult trials (rat 1: difference = 9.8 pp (percentage points), non-parametric one-tailed bootstrap,  $P < 10^{-4}$ ; rat 2: difference = 10.0 pp,  $P < 10^{-4}$ ; rat 3: difference = 8.0 pp,  $P < 10^{-4}$ ; see section 3.4). Importantly, when the psychometric curve was computed after incorrect trials, the slope of this curve increased significantly for all rats (rat 1: percentage change 42%, non-parametric one-tailed bootstrap,  $P = 4.4 \times 10^{-3}$ ; rat 2: percentage change 81%,  $P < 10^{-4}$ ; rat 3: percentage change 110%,  $P = 5 \times 10^{-4}$ ). The improvement was substantial, with an average relative increase of 9 pp in performance in difficult trials after incorrect responses compared to after correct responses (non-parametric one-tailed bootstrap,  $P < 10^{-4}$ ).

Consistent with the observation that the animals use the structure of the outcome-coupled hidden Markov chain to improve their behavior, we also found that on a session by session basis animals predominantly followed the lose-switch part of a win-stay-lose-switch strategy with a substantially weaker win-stay part (Fig. 3.1d; all rats: difference



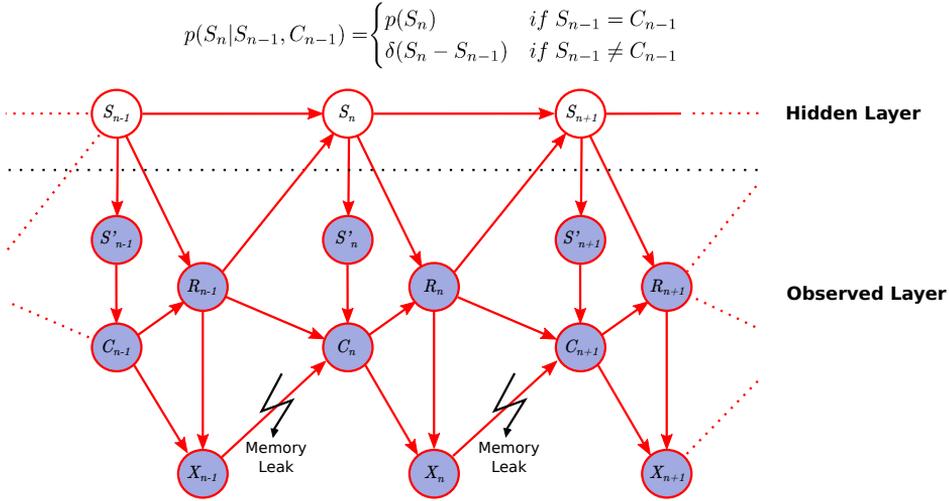
**Figure 3.1: Rats use the trial-by-trial-dependent contingencies of the task to improve their performance.** (a) Schematic of the task (see section 3.4 for details). Two identical, consecutive tones (T1 and T2) are presented to the rats (top panel). Inter-tone intervals (ITIs) can belong to two stimulus categories: short,  $S = s$ , or long,  $S = l$ . Each category has four possible ITIs (short: 50, 100, 150 and 200 ms; long: 350, 400, 450, 500 ms). The vertical dotted line represents the decision boundary at 275 ms. Difficult ITIs, depicted in gray, lie close to the decision boundary. Sequence of events within a trial (bottom panel): from trial initiation to choice. Rats self-initiate the trial and sample the stimuli in the central socket. They are rewarded with water if they poke the right socket when the stimulus is long, and the left socket when the stimulus is short (for rat 3 the contingency was the opposite). (b) The sequence of trials follows an outcome-coupled hidden Markov chain (see also Fig. 3.2): a new random stimulus condition is presented after a correct response (+), while the same stimulus condition is presented after an incorrect response (-).

**Figure 3.1:** (cont.) (c) Psychometric curves (probability choice long versus ITI) after correct responses (green line) and after incorrect responses (red line) for an example rat (left panel) and for all rats (right). The slope of the psychometric curves after incorrect responses substantially and significantly increases relative to the slope of the curve after a correct response. Error bars (shaded) are estimated by bootstrap (one s.d.). (c) Probability of lose-switch versus probability of win-stay. Each point corresponds to a different session. Rats predominantly follow a lose-switch over a win-stay strategy. No strategy being followed corresponds to the point (0.5, 0.5) in the plot.

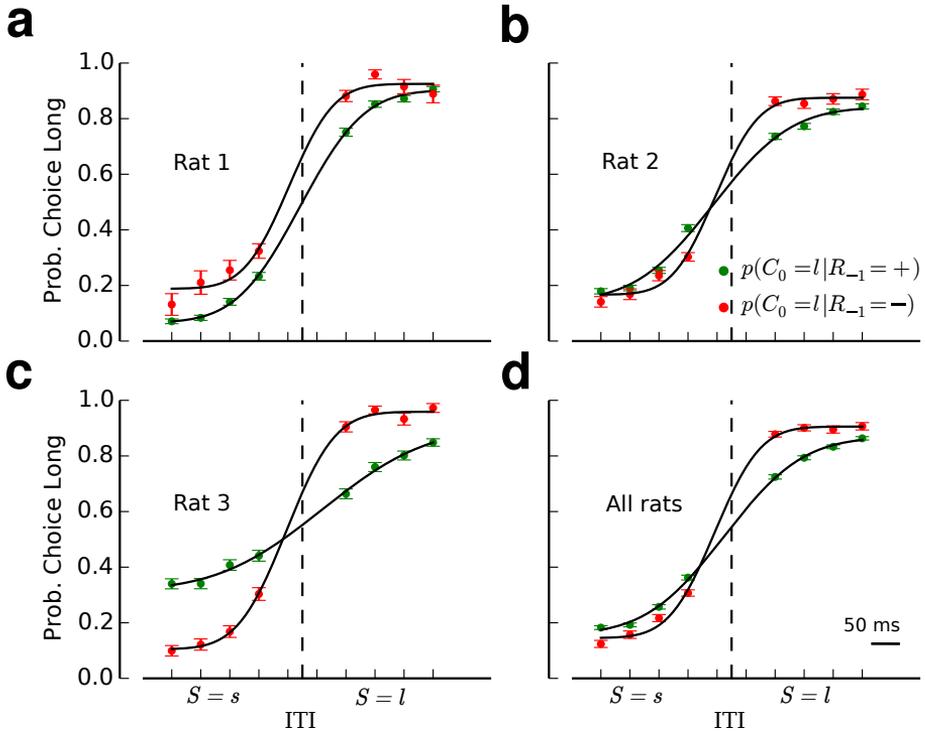
lose-switch—win-stay probabilities = 0.24 pp; non-parametric one-tailed bootstrap,  $P < 10^{-4}$ ; see section 3.4). Following a lose-switch strategy with no win-stay bias would lead to optimal behavior in our task if, ideally, the Markov chain were fully visible (not hidden). However, the actual ITI category in each trial is unobserved (because some trials are difficult) and the past might not be fully known due to memory leak. Consistent with this, the rats displayed some departures from the optimal strategy, in particular featuring a significant win-stay component in their behavior (rat 1: mean = 0.51, non-parametric one-tailed bootstrap,  $P = 1.0 \times 10^{-3}$ ; rat 2: mean = 0.54,  $P < 10^{-4}$ ; rat 3: mean = 0.75,  $P < 10^{-4}$ ).

The observed changes in the psychometric curve suggest that animals track a variable that jointly monitors previous choice  $C_{-1}$  ( $C_{-1} = -1$  if the choice was long, or  $C_{-1} = +1$  if it was short) and previous outcome  $R_{-1}$ . This second-order prior variable informs the rat about what choice it should make after an incorrect response, and mathematically is expressed as  $X_{-1} = C_{-1} \times R_{-1}$  (section 3.4). The state-space in our task consists both of the previous outcome and second-order prior, because these two variables fully define all that needs to be known by the rat to behave efficiently in this task. These two variables also fully define the prior information that is task-relevant, called immediate prior information. To confirm the prediction that the rats keep track of the second-order prior,  $X_{-1}$ , we asked how well past events are able to predict the upcoming choice  $C_0$ . Among the large number of behavioral variables that could

influence upcoming choices, we found that the second-order prior  $X_{-1}$  was the most predictive quantity, only surpassed by the stimulus itself,  $S_0$  and followed by the previous outcome  $R_{-1}$  (Fig. 3.4; section 3.4).



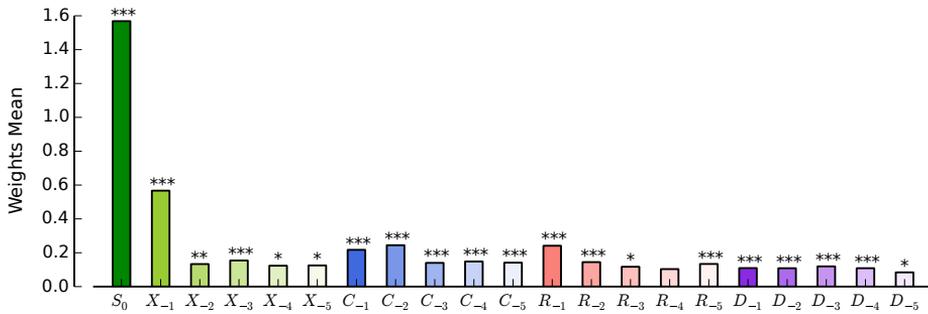
**Figure 3.2:** The sequence of trials follows an outcome-coupled hidden Markov chain: after a correct response a new random stimulus condition is drawn in the next trial, while after an incorrect response the same stimulus is repeated in the next trial (top equation). The hidden layer is composed by the actual stimulus sequence that is presented to the rat ( $\{S_n\}$ ), while the observed layer comprises the set of variables that are accessible to the rat: observed stimulus ( $S'_n$ ), choice ( $C_n$ ) reward ( $R_n$ ) and the second-order interaction between choice and reward ( $X_n$ ), (labels follow the same convention as in section 3.2 and section 3.4). Memory leak is incorporated into the model by allowing the choice to be based on a corrupted version of the previous second-order prior variable. The mathematical formulation of the outcome-coupled hidden Markov chain is displayed at the top.



**Figure 3.3:** Psychometric curves (probability choice long vs. ITI) after correct response (green dots) and after incorrect responses (red dots) for each rat (three first panels) and for all rats (fourth panel). Error bars are as in Fig. 3.1 (one standard deviation over bootstrap iterations). Fits correspond to a lapse-corrected cumulative Gaussian (see section 3.4).

### 3.2.2 Single-cells encode upcoming choice and second-order prior

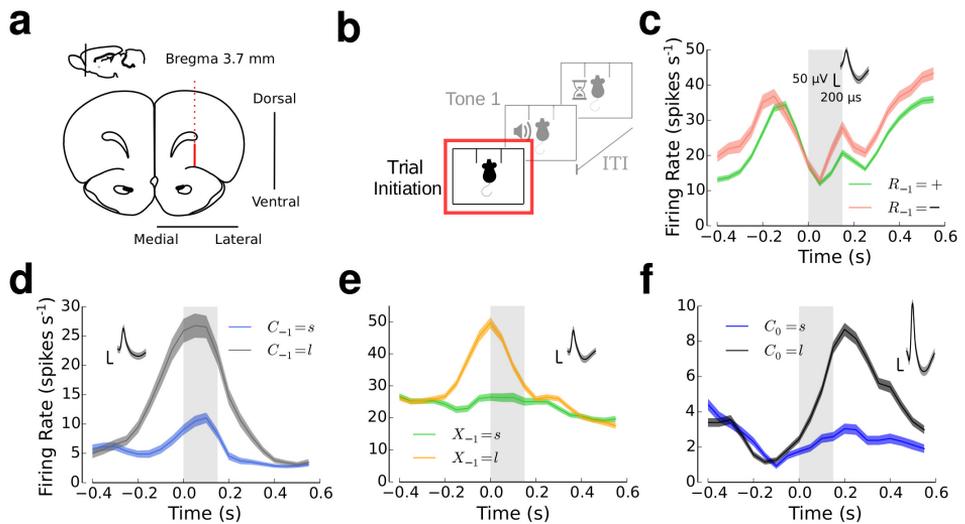
We looked for neural coding of immediate prior information and upcoming choices throughout the trial. Tetrodes were inserted in the right hemisphere of the rat IOFC (Fig. 3.5a). Small ensembles of well-isolated single units were simultaneously recorded (mean size =  $2.9 \pm 1.6$  neurons). Our dataset consisted of a total of 137 single-neurons with an average of 684 behavioral trials, eliciting a median of 9000 spikes per neuron, before



**Figure 3.4:** Logistic regression analysis predicting upcoming choice  $C_0$  based on a linear combination of binary regressors (displayed along the horizontal axis) shows that stimulus and second-order prior variables are the strongest predictors for behavior. Variables from two and more trials into the past have a weaker effect on the upcoming choice. Due to the large number of trials in our datasets and because of the hidden nature of the stimulus Markov chain, their effect is in most cases still significant (evaluated using a permutation test, see section 3.4),  $* = P < 0.05$ ,  $** = P < 0.01$ ,  $*** = P < 0.001$ .

excluding neurons with mean firing rates below 1 Hz (including all cells did not qualitatively influence the results; for a detailed description of the total number of cells for each analysis and an additional power analysis for the number of cells and rats see section 3.4). Recordings started after rats had reached a performance of at least 75%.

Our behavioral results suggest that the animals closely monitor second-order prior,  $X_{-1}$ , and other variables that correlate with it, such as previous choice  $C_{-1}$  and resulting outcome  $R_{-1}$ . We reasoned that if OFC participates in the decision-making process, then OFC neurons should encode these variables as well as reveal signals that anticipate upcoming choices. To test this prediction, we initially focused on the trial initiation period, where the stimulus has not yet been presented. We first aligned the neuronal responses to the initiation of the trial (Fig. 3.5b). Before performing pooled population-level analyses, we will first focus on the tuning of some example neurons. We found neurons whose trial-averaged activity illustrated a diversity of behaviors associated with both backward and also forward events. In Fig. 3.5 we show some individual examples.



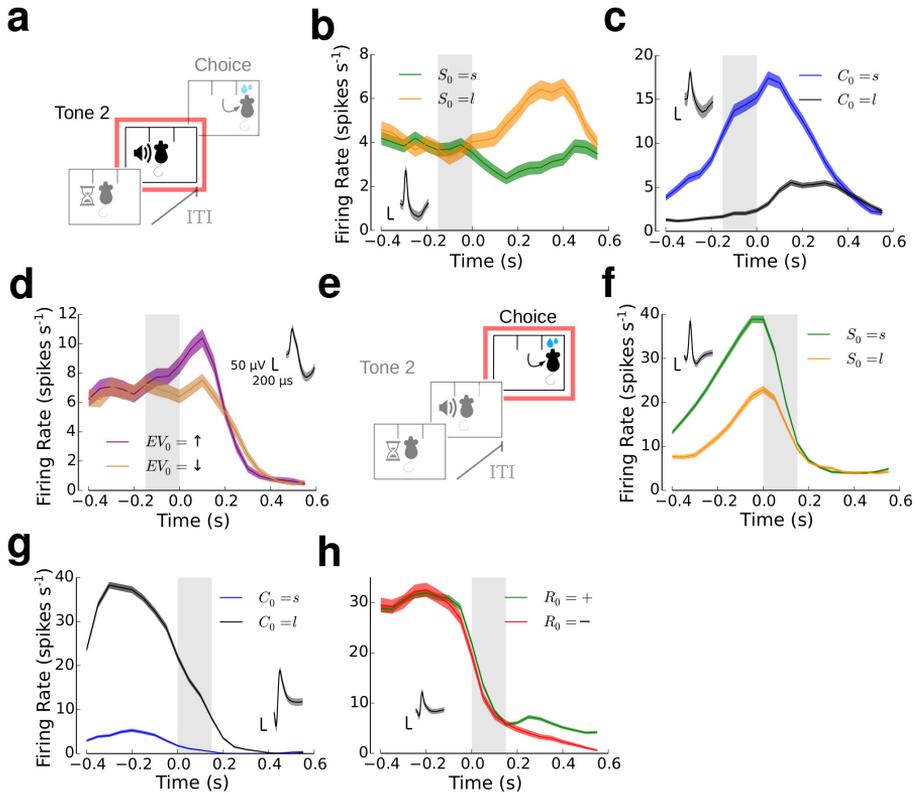
**Figure 3.5: OFC neurons encode relevant past information and anticipate upcoming choices even before stimulus onset.** (a) Electrode's path (dashed red line) and recording sites (solid red line) in rat IOFC depicted in a coronal section representation at 3.7 mm AP, 2.5 mm ML and 1.6 mm DV from Bregma. (b) Neuronal responses were aligned to trial initiation, defined as the time at which the rat starts the trial by poking the central socket. (c) Example neuron encoding reward in the previous trial  $R_{-1}$  (either + or -). This particular neuron fires more strongly for non-rewarded previous trials. (d) Example neuron representing choice in the previous trial  $C_{-1}$  (either  $s$  or  $l$ ). (e) Example neuron tracking second-order prior  $X_1 = C_{-1} \times R_{-1}$  (either  $s$  or  $l$ ). (f) Example neuron encoding upcoming choice  $C_0$  as a function of time. This neuron conveys information about rat's upcoming choice before stimulus onset (stimulus onset always happens to the right of the shaded area). (c-f): Time zero corresponds to trial initiation. The period of time between trial initiation and the shorter stimulus onset (150 ms) is indicated with shaded areas. Curves correspond to trial-averaged firing rates smoothed with a causal sliding rectangular window (size of 100 ms and step of 50 ms), and shaded areas around them correspond to s.e.m. Insets represent spike waveform for each neuron (black line, mean; shaded area, s.d.).

We identified neurons that showed conspicuous modulations as a function of the previous outcome (Fig. 3.5c), previous choice (Fig. 3.5d), second-order prior (Fig. 3.5e) and interestingly, also about upcoming choice (Fig. 3.5f). The neuron shown in Fig. 3.5f could predict upcoming choice with an accuracy of 71% (AUC, see section 3.4).

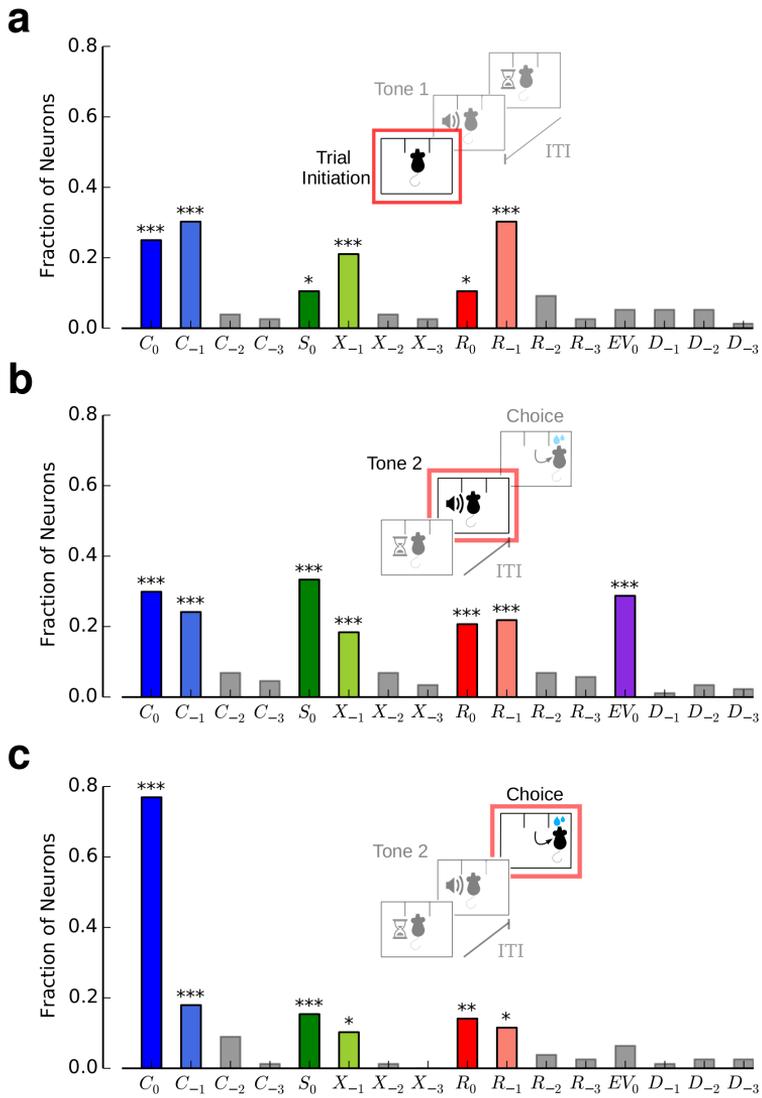
These quantities were also encoded throughout the trial (Fig. 3.6). Just before stimulus offset (Fig. 3.6a–d), when the animal is still poking into the central port, stimulus information is strongly represented in some neurons in IOFC (Fig. 3.6b). Signals about the upcoming choice were also clearly visible in this pre-movement period (Fig. 3.6c). This neuron predicted upcoming choice with 84% accuracy (AUC). Finally, the firing rate of some cells was modulated by the expected value of the outcome,  $EV_0$  (Fig. 3.6d; section 3.4). When we analyzed single-neuron responses at lateral nose poking onset (Fig. 3.6e), we found neurons whose rate was largely modulated by stimulus (Fig. 3.6f). Signals about the current choice were also strongly present, as shown by the example neuron in Fig. 3.6g. This neuron predicted the performed choice with 87% accuracy (AUC). We also observed outcome-modulated neurons in this period (Fig. 3.6h). Thus, even single-neuron activity by itself already provided strong indication that IOFC was representing the task-relevant variables.

### 3.2.3 OFC encodes immediate prior and anticipates future choices

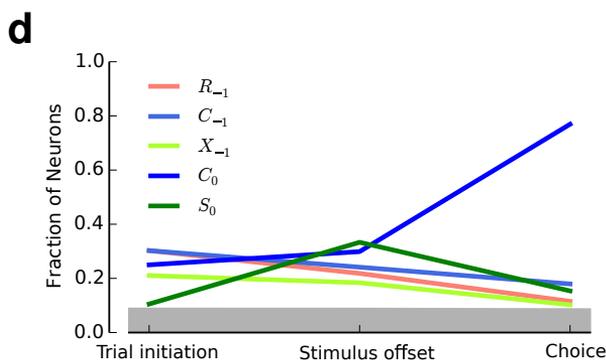
We confirmed the single-neuron observations at the population level with a Generalized Linear Model (GLM) analysis of the spike count responses of single-neurons. To do so, we regressed the spike count of each single-neuron simultaneously against a large set of variables, including the stimulus, reward, choice, difficulty and second-order prior of the current trial, the previous trial and up to three trials in back (section 3.4). This approach was preferred over a receiver operating characteristic (ROC) approach because the latter might find significant AUC values even in the absence of veridical encoding of the variable, simply due to correlations with other encoded variables (see section 3.4).



**Figure 3.6: OFC neurons encode essential quantities throughout the trial.** (a–d) Neuronal responses were aligned to stimulus offset (a). The firing rate of OFC neurons was modulated in a time period before stimulus offset (150 ms, shaded areas in b–d) by the stimulus (b), the upcoming choice (c) and the expected value of the outcome (d). (e–h) Neuronal responses were aligned to lateral nose poking onset, choice period (150 ms) (e). The firing rate of OFC neurons represented the stimulus (f), the current choice (g) and the outcome (h). In the two periods, signals about upcoming and current choice were very conspicuous. Time zero corresponds to stimulus offset (b–d) and lateral nose poking onset (f–h). Curves correspond to trial-averaged firing rates smoothed with a causal sliding rectangular window (size of 100 ms and step of 50 ms), and shaded areas around them correspond to s.e.m. Insets represent spike waveform for each neuron (black line, mean; shaded area, s.d.).



**Figure 3.7: Neurons in IOFC integrate prior with current sensory information and encode upcoming choice.**

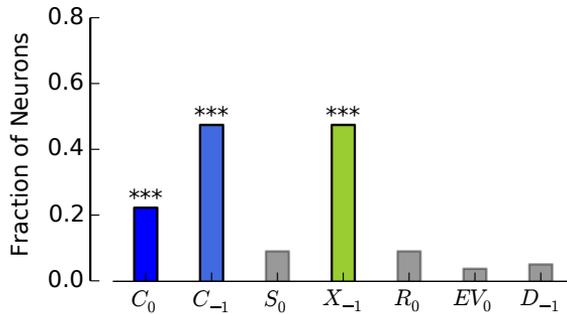


**Figure 3.7:** (cont.) **(a)** Fraction of neurons with significant regressors (see section 3.4) for each of the variables listed in the horizontal axis. Upcoming choice  $C_0$ , previous choice  $C_{-1}$ , upcoming stimulus  $S_0$ , second-order prior  $X_{-1}$ , upcoming outcome  $R_0$  and previous outcome  $R_{-1}$  are significantly encoded in the population. Upcoming choice  $C_0$  is encoded by IOFC neurons even before stimulus is presented. Variables that extend further back into the past are not significantly encoded in IOFC. **(b–c)** Fractions of neurons with significant regressors during a time period before stimulus offset **(b)**, and during a time period after lateral nose poking **(c)**. **(d)** Fractions of neurons encoding upcoming choice  $C_0$ , stimulus  $S_0$ , second-order prior  $X_{-1}$ , previous choice  $C_{-1}$  and previous reward  $R_{-1}$  at trial initiation (pre-stimulus), stimulus offset and choice periods. Note that larger fractions of neurons have choice-related signals as time progresses through the trial. **(a–c)** One-tailed binomial test,  $*$  =  $P < 0.05$ ,  $**$  =  $P < 0.01$ ,  $***$  =  $P < 0.001$ . Shaded rectangle corresponds to non-significant fraction of neurons ( $P > 0.05$ ).

Before stimulus onset, we found that a significant fraction of neurons (25%, one-tailed binomial test,  $n = 76$ ,  $P = 4.6 \times 10^{-9}$ ) predicted the upcoming choice,  $C_0$  (Fig. 3.7a). Significant fractions of cells also encoded the second-order prior  $X_{-1}$ , previous choice  $C_{-1}$ , and the previous outcome  $R_{-1}$ . Thus, the neurons shown in Fig. 3.5 represent just examples of potentially overlapping large neuronal populations that encode these variables. Interestingly, we did not find a substantial fraction of cells encoding information from two or more trials into the past, suggesting that information older than arising from the preceding trial is not present in

IOFC.

We found that cells encoded both current stimulus  $S_0$  and current outcome  $R_0$  ( $S_0$  and  $R_0$ , 11% each, one-tailed binomial test,  $n = 76$ ,  $P = 0.036$ ) even before stimulus onset. Although at first glance surprising, this result arises from the outcome-coupled hidden Markov chain structure of the environment. In fact, when we repeated our GLM analysis using only trials after correct responses—where the upcoming stimulus cannot be predicted from the stimulus used in the previous trial—we found that neither stimulus  $S_0$  (9%, one-tailed binomial test,  $n = 76$ ,  $P = 0.085$ ) nor reward  $R_0$  (9%, one-tailed binomial test,  $n = 76$ ,  $P = 0.085$ ) information was present before the onset of the stimulus (Fig. 3.8). Focusing instead only on trials after incorrect responses, we again found that a substantial fraction of cells (14%, one-tailed binomial test,  $n = 76$ ,  $P = 1.3 \times 10^{-3}$ ) can predict the stimulus.



**Figure 3.8:** Fraction of neurons with significant regressors for each of the variables listed in the horizontal axis when the linear model was fitted exclusively using trials after correct responses. Note that in this case previous choice and second-order prior are identical variables, but they are displayed separately for better comparison with Fig. 3.7. One-tailed binomial test,  $n = 76$ ,  $*$  =  $P < 0.05$ ,  $**$  =  $P < 0.01$ ,  $***$  =  $P < 0.001$ .

Altogether, our results show that, before stimulus onset, IOFC tracks the second-order prior  $X_{-1}$ , and anticipates the upcoming choice,  $C_0$ . Thus, rat IOFC carries sufficient information to play an important role in integrating immediate prior information with sensory information.

### 3.2.4 Build-up of choice-related neuronal signals

If OFC represents the integration of immediate prior with current information, then information about upcoming choices should increase as further evidence is integrated into the system. For instance, just before stimulus offset, information about the stimulus is readily available, and should be combined with prior information to inform decisions. In fact, a substantial fraction of cells encoded the upcoming choice  $C_0$  just before stimulus offset (Fig. 3.7b). This fraction was large (30%, one-tailed binomial test,  $n = 87$ ,  $P = 7.6 \times 10^{-14}$ ), and larger than during the pre-stimulus period, though not significantly (see Fig. 3.7a,b; difference = 5 pp, one-tailed non-parametric difference binomial test,  $P = 0.25$ ; see section 3.4). Integration of information at the population level could be accomplished within the same circuit, as a large fraction of cells also encoded the stimulus  $S_0$  in the current trial (33%, one-tailed binomial test,  $n = 87$ ,  $P = 1.1 \times 10^{-16}$ ). Interestingly, in the choice period, 77% of all cells (60/78 neurons) encoded choice (Fig. 3.7c) –significantly more than in the pre-stimulus periods (Fig. 3.7a,b; difference 52 pp, one-tailed non-parametric difference binomial test,  $P < 10^{-4}$ ). Thus, there is a build-up of choice-related signals in IOFC, as illustrated when plotted as a function of the analysis time period (Fig. 3.7d).

Stimulus also seemed to be encoded in OFC in a sensible way, with information peaking before stimulus offset. We found that the fraction of neurons encoding stimulus  $S_0$  increases significantly from trial initiation to the stimulus offset period (Fig. 3.7d; difference = 23 pp, one-tailed non-parametric difference binomial test,  $P = 2.2 \times 10^{-4}$ ) and decreases significantly thereafter (difference = 18 pp, one-tailed non-parametric difference binomial test,  $P = 3.9 \times 10^{-3}$ ). Encoding of past task events, such as second-order prior, previous choice and previous reward, declined as time progressed over the trial (Fig. 3.7d;  $C_{-1}$ : difference = 12 pp, one-tailed non-parametric difference binomial test,  $P = 0.036$ ;  $X_{-1}$ : difference = 11 pp,  $P = 0.033$ ;  $R_{-1}$ : difference = 19 pp,  $P = 2.2 \times 10^{-3}$ ; differences computed between pre-stimulus and choice periods). Altogether, these time profiles suggest that information about stimulus and

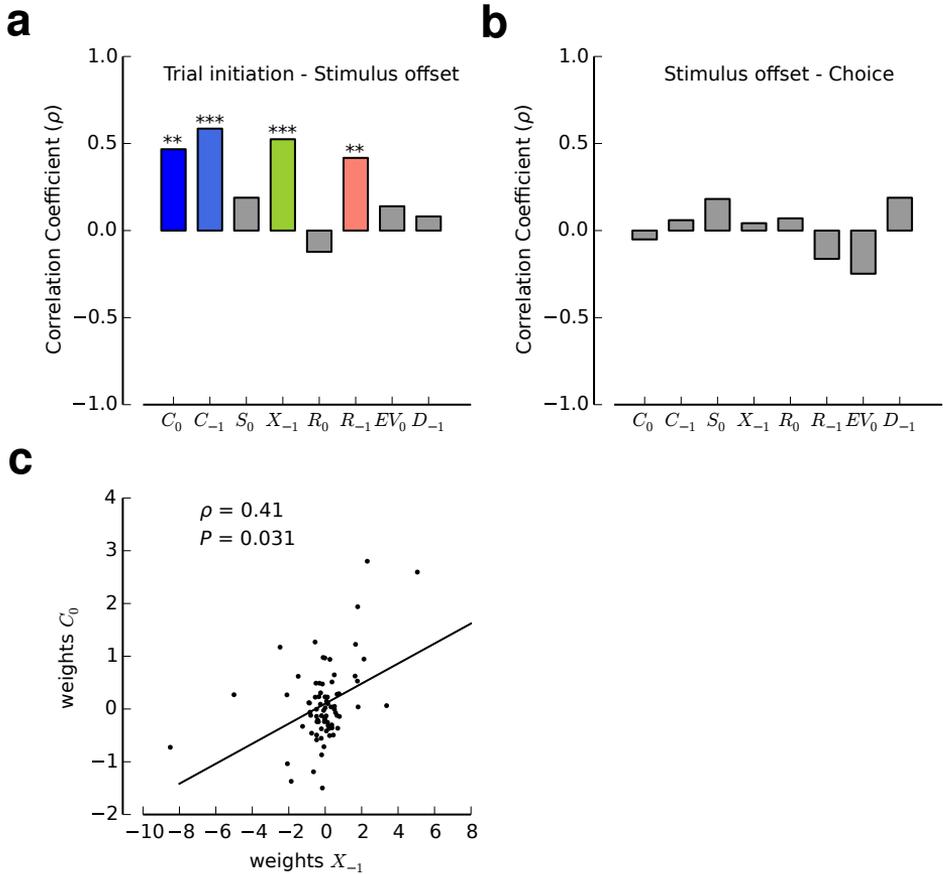
second-order prior is incorporated into choice-related signals to mediate the integration of information.

We found a correlation between the encoding weights for upcoming choice computed at the pre-stimulus and stimulus offset periods (Fig. 3.9a; section 3.4). The same was observed for the weights computed for second-order prior, previous choice and previous reward. This suggests that the encoding of these variables is partially sub-served by stable populations during the periods of time in which prior information needs to be integrated with sensory information. However, their encoding differed in the choice period, precisely when sensory information does not need to be integrated any more, as not such correlation was observed (Fig. 3.9b). In particular, the increase of choice-encoding neurons over time reported in Fig. 3.7d suggests that the lack of correlation between encoding during stimulus offset and choice periods might arise from a recruitment of additional choice-related cells, potentially motor-related. We also found that the encoding weights of second-order prior and upcoming choice were positively correlated during the pre-stimulus period (Fig. 3.9c; section 3.4), suggesting that populations of neurons encoding the previous trial's state and upcoming choice partially overlap before stimulus presentation.

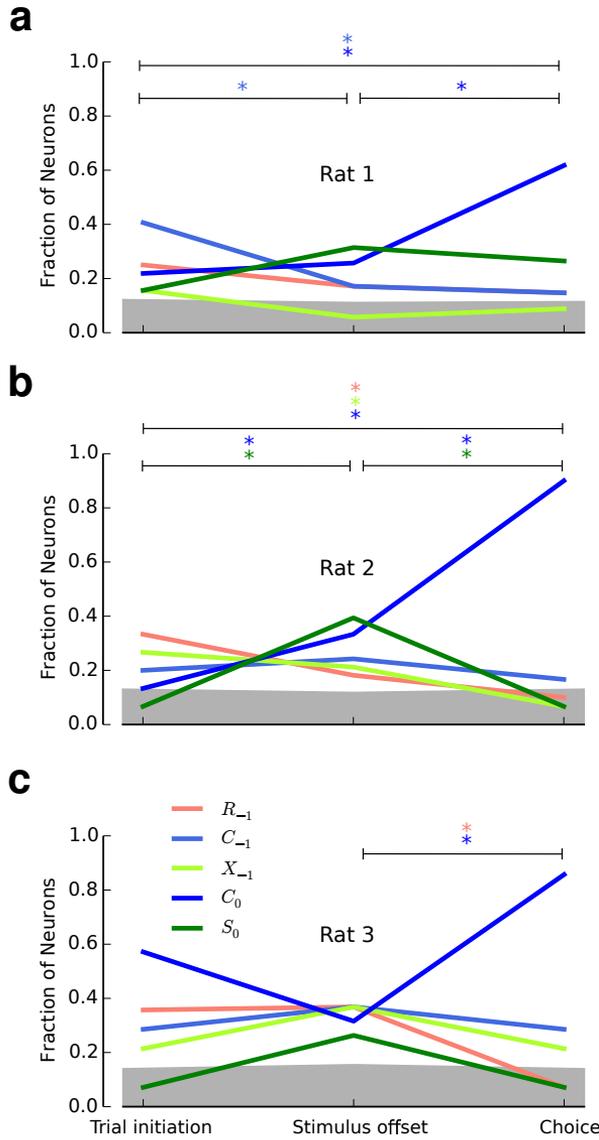
Some differences in behavior across animals were clear (see Fig. 3.1d and Fig. 3.3), with rat 3, for instance, displaying a higher lose-switch probability than the other rats. We first confirmed in a separate analysis that none of the qualitative results described above changed when neurons recorded from this rat were excluded from the analysis. We also confirmed that rat-by-rat analysis of neuronal populations delivered the same trends as reported above, generally including encoding of upcoming choice before stimulus onset and the ramping of choice-related information across time periods (Fig. 3.10).

### **3.2.5 Expected value and outcome representations**

After stimulus presentation, at stimulus offset, the animal might have a sense of how difficult the trial was. This informs about the subjective probability (confidence) of getting a reward, as easy trials should promise



**Figure 3.9: Encoding of essential variables for the task is stable before motor execution of the choice.** (a) Correlation coefficient between weights estimated at trial initiation and before stimulus offset for several variables (see section 3.4). Upcoming choice  $C_0$ , second-order prior  $X_{-1}$ , previous choice  $C_{-1}$  and previous reward  $R_{-1}$  are stably encoded in IOFC. (b) None of the correlation coefficients between weights computed just before stimulus offset and during the choice period were significantly different from zero. (c) There is a positive correlation between the encoding weights associated with second-order prior and upcoming choices at the pre-stimulus period. Two-tailed permutation test (see section 3.4), \* =  $P < 0.05$ , \*\* =  $P < 0.01$ , \*\*\* =  $P < 0.001$ .



**Figure 3.10:** Temporal evolution of the fractions of neurons with significant regressors for each rat (a-c).

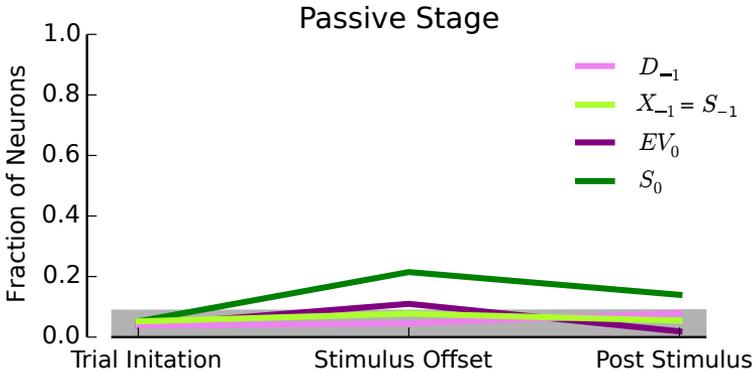
**Figure 3.10:** (cont.) Significance of the fraction of neurons corresponds to values above the gray rectangle ( $P < 0.05$ ,  $n = 76$  (trial-initiation),  $n = 87$  (stimulus-offset),  $n = 78$  (choice)). Thus, for the three rats the upcoming choice, previous reward and second-order prior are significantly encoded at trial initiation, before stimulus onset (except for rat 2 for which upcoming choice is not significant at the borderline level 0.06). Consistent for the three rats, current stimulus is strongly encoded during the stimulus offset period. Significance of the differences in fractions of neurons are indicated with elongated bars at the top and a start with the same color as the color variable ( $* = P < 0.05$ ; see section 3.4). Only significant changes are indicated (for rat 3, there are not significant changes from trial initiation to stimulus offset, likely because of noisier estimates of the fractions due to the lower number of neurons recorded in this animal). Overall, for the three rats there is a significant increase of choice-related signals across major task periods.

a more secure reward than difficult trials. Since in our experimental setup we do not vary the reward amount, encoding the subjective probability of a positive outcome amounts to the expected value in the current trial, which in turn is inversely related to the difficulty of the trial (see section 3.4). In this time epoch, the expected value was encoded in a large fraction of cells (Fig. 3.7b; 29%, one-tailed binomial test,  $n = 87$ ,  $P = 6.1 \times 10^{-13}$ ). Previous work has also found that signals about decision confidence are encoded in the activity of single-cells in rat OFC [124], and in monkey parietal cortex [125]. We also found in this period of time a large fraction of cells that encode outcome in a predictive way, as this variable can be partially inferred based on the difficulty of the trial. Outcome was also encoded at the choice period (Fig. 3.7c), consistent with the role of this area in encoding reward and outcomes [87, 120].

### 3.2.6 Behaviorally irrelevant prior is not represented in OFC

The previous results demonstrate that OFC represents state-space when rats are in an environment where it is behaviorally advantageous to keep

track of this information. We tested the encoding of immediate prior information when this information was irrelevant by placing the same rats in an environment where they were passively exposed to the same set of stimuli but rewards were not delivered. Rats were exposed to two passive stages, before and after the decision-making stage (see section 3.4). We found that OFC did no longer keep track of the immediate prior information (defined as previous stimulus  $S_{-1}$  in the passive environment, equivalent to  $X_{-1}$  in the decision-making stage; see section 3.4) at any time during the trial (Fig. 3.11). Encoding of current stimulus and difficulty at the stimulus-offset period weakly persisted in this environment, suggesting that task-irrelevant variables observable at the current trial are not completely filtered out in OFC. These results suggest that OFC does not monitor state-space from the immediate past when this information is task-irrelevant.

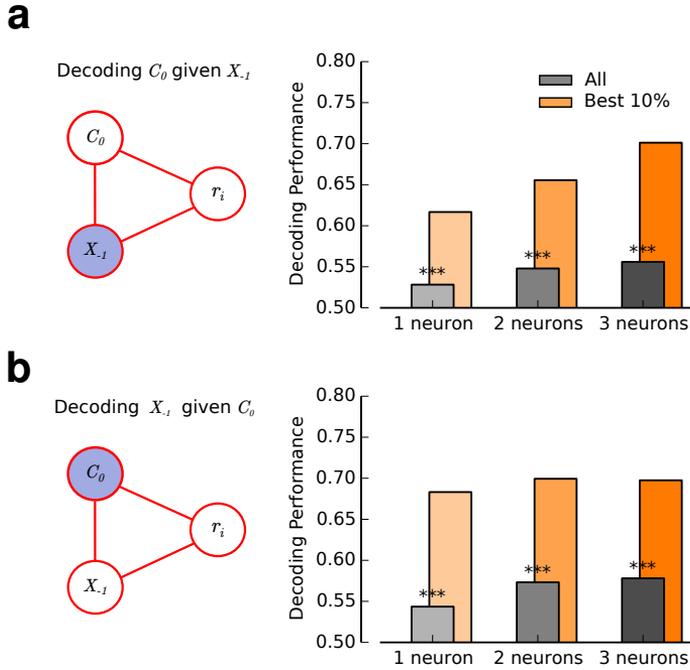


**Figure 3.11:** Temporal evolution of the fraction of neurons with significant regressors for the stage where rats are passively exposed to the same set of stimuli as in the decision-making task. The depicted fraction of neurons corresponds to the mean of the two passive stages. Neurons exclusively represent current stimulus and current expected value (difficulty) at stimulus-offset period. Thus, during the passive stage second-order prior information is no longer encoded in OFC neurons. Shaded rectangle corresponds to non-significant fraction of neurons ( $P > 0.05$ ,  $n = 76$  (trial-initiation),  $n = 87$  (stimulus-offset),  $n = 78$  (choice))

### 3.2.7 Population decoding reveals a hierarchy of variables in OFC

Our previous analysis has revealed that, following correct choices, only two variables are significantly encoded in the pre-stimulus period in single OFC neurons, namely, second-order prior and upcoming choice (Fig. 3.8). We confirmed that this result holds using a much more stringent test that does not assume that both variables are encoded linearly, as we did before. To do so, we used decoding techniques that predict one quantity at a time from the population activity of a simultaneously recorded neuronal ensemble [51], while keeping the other quantity constant (Fig. 3.12; section 3.4). We found that a classifier trained on the pre-stimulus activity of a neuronal ensemble at fixed second-order prior  $X_{-1}$  conveyed substantial information about upcoming choice (Fig. 3.12a). Similarly, when conditioning the activity to upcoming choice  $C_0$ , we found that small neuronal populations conveyed substantial information about second-order prior (Fig. 3.12b). These results hold both across all neuronal ensembles in the dataset and when selecting only the 10% most informative ensembles. Decoding performance increased monotonically with the number of neurons in the ensemble (Fig. 3.12a,b) [126, 127]. Because this conditioning-based decoding analysis does not assume that these two variables are both encoded linearly, in contrast to our previous analysis (Fig. 3.7), these results add strong support to the conclusion that both immediate prior information (that is, second-order prior) and upcoming choice are encoded in lOFC.

Which variables are most readily decoded at the population level? The analysis from the previous sections would suggest upcoming choice and prior information as strong contenders. However, this analysis was based on single neurons and ignored correlations that might be present in neuronal populations and might influence the representation of those variables. To more directly address this question, we trained a classifier as in the previous paragraph to decode per trial individual variables from the activity of small neuronal ensembles (section 3.4). Using this approach, we found that, consistent with the previous linear encoding analysis (Fig.



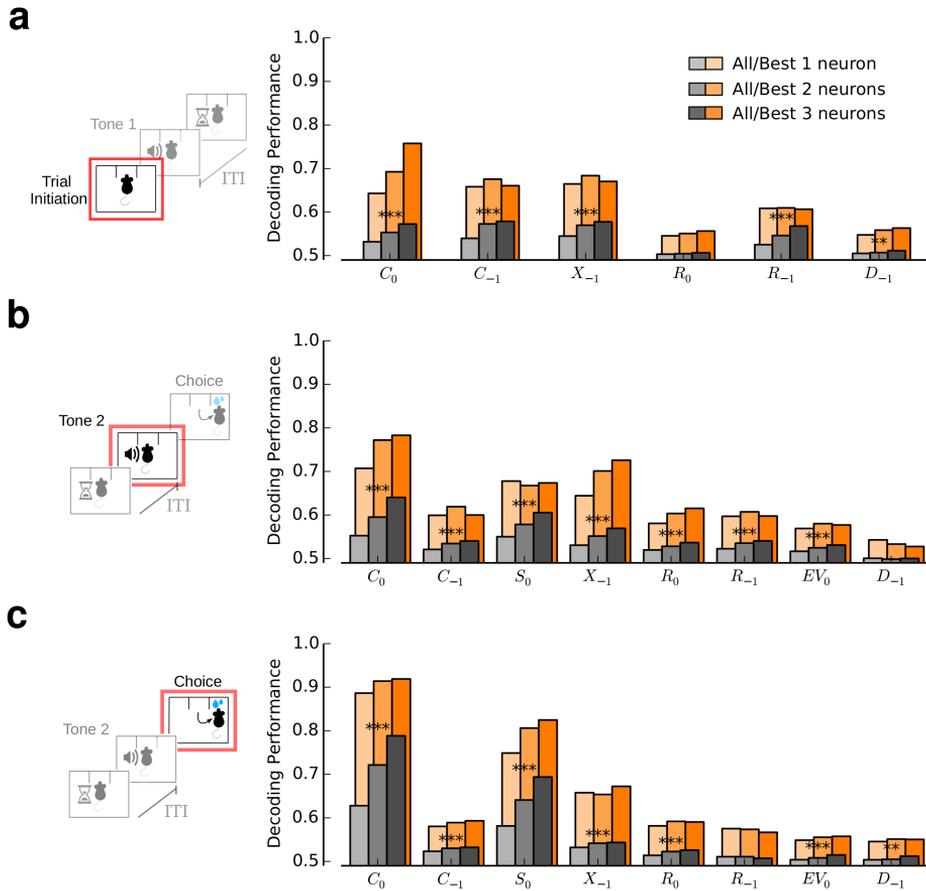
**Figure 3.12: Population decoding reveals pre-stimulus neural representations of second-order prior and upcoming choice.** (a) Decoding performance for upcoming choice at fixed second-order prior increases with the number of neurons in the ensemble (one to three) across all ensembles (gray), and does so more strongly for the 10% most informative ensembles (orange). Only trials after correct responses are used for the analysis. Left panel: schematic showing that the pre-stimulus firing rate of neuron  $i$  in the ensemble can possibly depend at most on second-order prior  $X_{-1}$  and upcoming choice  $C_0$ , as previously revealed by a linear analysis (Fig. 3.8). To show that upcoming choice truly modulates neural activity, we performed a conditioned analysis by which the value of the second-order prior is fixed (gray-blue) while a linear classifier is trained to predict upcoming choice from the activity patterns in OFC (see section 3.4). (b) Decoding performance for second-order prior at fixed upcoming choice. Colour code and analysis are as in the previous panel. One-tailed permutation test,  $*$  =  $P < 0.05$ ,  $**$  =  $P < 0.01$ ,  $***$  =  $P < 0.001$ .

3.7d), the 10% most informative neuronal ensembles had larger amounts of information about upcoming choice than about any other variable (Fig. 3.13) from the pre-stimulus to the choice periods. Information about the upcoming choice  $C_0$  was so strongly present in IOFC that it could be predicted from holdout data not used to train the classifier with an accuracy of 57% for all ensembles and 76% for the top 10% ensembles in the pre-stimulus period, 64 and 78% at the stimulus offset period, and 79 and 92% at the choice period (Fig. 3.13a–c), respectively. The population decoding analysis also again revealed second-order prior as one of the most prominently encoded variables (Fig. 3.13a–c). Other variables were also decodable from the IOFC, but less accurately. Therefore, the population decoding analysis confirms that IOFC tracks prior information on a trial by trial basis and predicts upcoming choice.

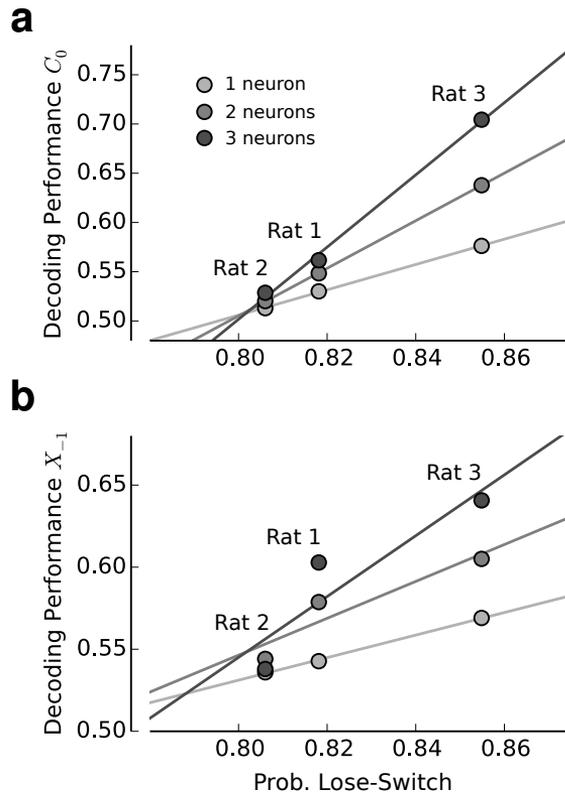
Finally, in view of the individual behavioral differences across animals, we sought to determine whether they were correlated with neuronal differences. We found a positive correlation between lose-switch probability and neuronal information about both upcoming choice and second-order prior, although this correlation did not reach significance (Fig. 3.14; permutation test,  $n = 3$ ,  $P = 0.16$ , section 3.4). Thus, animals that were more likely to switch after an incorrect response tended to provide a better information-readout in OFC ensembles about variables that are strongly linked to that switching behavior.

### 3.3 Discussion

OFC is thought to play an important role in adaptive and goal-directed behavior [78, 79, 82, 85, 97, 118, 119]. However, as OFC has been shown to encode a myriad of variables, including outcomes, expected rewards and values [80–82, 86, 87, 91–94, 120, 121], a coherent picture of its function is still missing. Previous work on reversal learning [128–130] and Pavlovian-instrumental transfer [131] has revealed that OFC function reflects crucial aspects of learning, particularly by developing novel representations of associations between cues and their predicted rewards [130,



**Figure 3.13: Population decoding analysis reveals a hierarchy of encoded variables.** (a) Decoding performance for each quantity at the pre-stimulus period as a function of the number of neurons in the ensemble (one to three) across all ensembles (gray) and for the 10% most informative ensembles (orange). All trials are used for the analysis. (b,c) Same as in the previous panel for stimulus offset and choice periods. This analysis reveals a hierarchy of encoding, with upcoming choice  $C_0$  and second-order prior  $X_{-1}$  being two of the most strongly encoded variables. One-tailed permutation test,  $*$  =  $P < 0.05$ ,  $**$  =  $P < 0.01$ ,  $***$  =  $P < 0.001$ .



**Figure 3.14:** Animals that have a higher probability of switching choice after an incorrect response also tend to have more information in OFC about both upcoming choice and second-order prior **(a)** Correlation between neuronal information about upcoming choice  $C_0$  and lose-switch probability across rats. Correlation for all ensembles sizes is strong but not significant (mean Pearson correlation = 0.999, two-tailed permutation test,  $P = 0.167$ ,  $n = 3$ , see section 3.4). **(b)** Correlation between neuronal information about second-order prior and lose-switch probability across rats. Correlation for all ensembles sizes is strong but not significant (mean Pearson correlation = 0.947, two-tailed permutation test,  $P = 0.167$ ,  $n = 3$ , see section 3.4).

132, 133], and by tracking the history of previous outcomes and choices during reward-guided decisions [134, 135]. These results show that OFC is important to process prior information that builds over an extended sequence of previous trials to guide behavior. However, it is not well known whether this goal is accomplished through a compact representation of the task's state-space, or by representing all sorts of task-relevant and task-irrelevant variables. Further, whether state variables can be represented exclusively from the previous trial at a high temporal resolution is not known.

We specifically tackled these questions by using a novel perceptual decision-making task endowed with an outcome-coupled hidden Markov chain. By introducing outcome-dependent correlations between consecutive stimuli, we ensured that the animal needed to track on a trial by trial basis the most recent past information to solve the task efficiently. This experimental design maximized the chances of finding state variables that need to be represented at high temporal resolution. It also maximized the chances of identifying interactions of these variables with choice-related signals during the decision-making process. In addition, by inserting random trials after correct responses, an analysis based on systematically conditioning on different task variables allowed us to distinguish neuronal signals that were purely associated with either the immediate past (for example, second-order prior) or future (upcoming choice) events. Thus, this task constitutes an important contribution to the classical perceptual decision-making literature by adding the necessity of considering immediate prior information. Indeed, except for some notable exceptions [74, 75, 115, 136], the study of perceptual decision-making has been dominated by paradigms where sensory information, presented in a random sequence of trials, suffices to inform a correct choice such that prior information from the previous trial can and should be ignored altogether [114, 137]. In this line, many studies have emphasized continuous integration of information over time within a trial [67, 114, 138]. As a consequence, relatively less work has focused on the discrete-like process required to integrate proximal prior events with sensory information [139].

One important feature of our task is that relevant prior information

was exclusively present in the previous trial. This immediate prior information was encapsulated in the second-order prior variable  $X_{-1}$ , the interaction between previous trial choice and reward. The second-order prior along with the previous outcome fully defined the state-space in our task. Our results show that IOFC represents the structure of the task in a compact way, as we found that second-order prior was among the most strongly encoded variables in IOFC. Our results are in line with a theoretical proposal [97] recently supported by human functional magnetic resonance imaging (fMRI) and rat inactivation studies [132, 136] that OFC represents the state-space, and hence add electrophysiological single-cell and neuronal population evidence for such theoretical scenario. In contrast, previous work has shown that in other brain areas, like the dorsolateral prefrontal cortex in monkeys, both task-relevant and task-irrelevant information is encoded in value-based decision-making [137, 140]. In addition, we also embedded animals in an environment in which they had to ignore prior information. In this environment, immediate prior information seemed to be abolished in OFC, suggesting that OFC differentially represents state variables that are relevant for the task.

Another important question is the degree of involvement of OFC in the decision-making process. We found a definite encoding of choice-related variables throughout the decision process, appearing even before stimulus onset. This result is consistent with recent work where monkey OFC population activity has been postulated to represent an internal deliberation mediating the choice between two options [141]. It is also in line with a large body of work showing that OFC plays an important role in goal-directed behavior and thus in action initiation and selection (for example, see refs [136, 142–146]). Previous work has also found evidence that a multitude of areas are involved in action initiation and selection, such as parietal and prefrontal areas [27, 67, 76, 81, 91, 93]. However, our results constitute the first report of the existence of neurons in the rodent OFC that have predictive power about upcoming choices before stimulus onset. Interestingly, some of these neurons were found to anticipate upcoming choice with a success probability of 69% (out of 750 test trials not used for training, 520 were correctly predicted using

logistic regression), thus demonstrating the presence of strong choice-encoding neurons in OFC even during the pre-stimulus period. At the population level the fraction of neurons encoding for upcoming choice before stimulus onset was strong and highly significant.

Finally, we found evidence that the observed compact representation of state-space in OFC can play a role in integrating immediate prior with current information. First, we found a strong representation of current stimulus information that declined after stimulus offset, an effect that was accompanied by a large increase of choice-related signals representing the integration of stimulus with prior information. This result suggests that the neuronal representation of the state-space interacts in the OFC with the decision-making process, potentially by facilitating the combination of prior with current information. This result is consistent with a recent human fMRI study suggesting that OFC represents posterior probability distributions by integrating extended prior experience with current information [147].

All in all, our results provide an integrative view of the rodent IOFC by showing that it predominately represents state-space (in particular, second-order prior), the integration of immediate past with current information, and the initiation and selection of choices. Our results, finally, open an interesting door to study the link between individual differences in behavior and detailed OFC electrophysiological encoding, by suggesting that animals that lose-switch more also have a stronger neuronal representation of past behaviorally relevant variables, and support the notion that across-subjects OFC differences modulate overall behavior, such as risk-seeking [148] and drug-seeking [149] behaviors.

## **3.4 Methods**

### **3.4.1 Behavioral task**

Three Wistar rats were trained to perform an auditory time-interval categorization task. Trials were self-initiated by the animals by nose poking, which elicited a pure tone of 50 ms duration after a random delay drawn

from a uniform distribution with values 50, 100, 150, 200, 250 and 300 ms. A second tone, identical in duration and frequency to the first one, was presented after a time interval, called ITI. The task is to categorize the ITI, as short ( $S = s$ ) or long ( $S = l$ ). ITIs are drawn randomly (see below for incorrect trials) from a uniform discrete distribution with values 50, 100, 150 or 200 ms for short intervals ( $S = s$ ) and 350, 400, 450 or 500 ms for long intervals ( $S = l$ ). Reward is provided in trials in which the animal sampled the full stimulus and poked to the left (right) socket, when the stimulus was short (resp. long). False alarms (poking in the opposite side) or early withdrawals (withdrawal before stimulus termination) were punished with a 3-s time out and a white noise (WAV-file, 0.5 s, 80-dB sound pressure level). After an incorrect trial, the ITI of the previous trial was repeated. This experimental design created correlations across trials based on the behavior of the animal. The mean fraction of false alarms was 0.08, 0.11 and 0.15 and the mean fraction of early withdrawals was 0.37, 0.31 and 0.14 for rat 1–3, respectively. All trials during task performance were self-initiated.

The animals went through two additional passive stages before and after the decision-making stage described above. During the passive stages rats were presented with the same set of stimuli as in the decision-making stage while they could freely move around the environment. Rewards were not provided at any time during the passive stages. Passive stage A occurred before the decision-making stage and it lasted a fixed set of stimulus presentations (rat 1: 400 trials; rat 2: 600 trials and rat 3: 600 trials). Passive stage B occurred after the decision-making stage, and it lasted the same number of stimulus presentations as in passive stage A. The experiment was approved by the animal Ethics Committee of the University of Barcelona. Rats were cared for and treated in accordance with the Spanish regulatory laws (BOE 256; 25-10-1990), which comply with the European Union guidelines on protection of vertebrates used for experimentation (EUVD 86/609/EEC).

### 3.4.2 Presentation of acoustic stimuli

The protocols of stimulation were controlled through MATLAB, a National Instrument card (BNC-2110), and a breakout box (FS 300 kHz). Sound triggers had microsecond precision. Sound tones were delivered through earphones (ER.6i Isolator, Etymotic Research), which were screwed in each recording session to the earphone holders, chronically attached to the animal skull with dental cement. The earphones were adjusted inside the ear with silicone tips with a separating distance of 1 mm from the ear canal. Similarly, sound calibration was performed inside the acoustic isolation box with a microphone (MM1, Beyerdynamic) placed 1 mm away from the earphone and using a preamplifier (USB Dual Pre, Applied Research and Technology). The sound tones had a duration of 50 ms, with an intensity of 80-dB SPL pure tones of 5322 Hz, and 6-ms rise/fall cosine ramps.

### 3.4.3 Logistic regression of behavior

Rats' choices were classified on a trial-by-trial basis by a logistic regression (linear classifier). Classification was based on a decision variable  $DV$ : when  $DV > 0$  the trial was classified as belonging to class 1 ( $C_0 =$  short choice), when  $DV < 0$  the trial was classified as belonging to class 2 ( $C_0 =$  long choice). The decision variable  $DV$  was a weighted sum of the all task variables we thought might be influencing rat's behavior:  $DV = \sum_{i=1}^M \omega_i x_i + \omega_0$ , where  $\omega_i$  and  $x_i$  were each task variable's contribution to the decision and its particular value ( $x_i = \pm 1$ ) respectively,  $\omega_0$  was the offset term and  $M$  was the total number of variables used for predicting rat's behavior (21 regressors in total). Logistic regression assumes that the probability of  $C_0 = 1$  (short choice) to be the correct class given the task variables is given by

$$p(C_0 = 1 | \{x_i\}) = \sigma \left( \sum_{i=1}^M \omega_i x_i + \omega_0 \right), \quad (3.1)$$

where  $\sigma(\cdot)$  is the logistic function. The set of task variables  $\{x_i\}$  is  $S_0$ ,  $R_{-n}$ ,  $D_{-n}$ ,  $C_{-n}$  and  $X_{-n}$  where  $n$  is the number of trials back in time, which ranged from 1 to 5. Here  $R_{-n}$  is the reward given to the rat  $n$  trials back in time, that is, the correctness of the response (+1 correct, rewarded, -1 incorrect, non-rewarded);  $D_{-n}$  is the trial difficulty defined on the basis of the distance between the presented ITI and the category boundary (50, 100, 450 and 500 ms, easy trial,  $D_{-n} = +1$ ; 150, 200, 350 and 400 ms, difficult trial,  $D_{-n} = -1$ );  $C_{-n}$  is rat's choice (+1 short choice, -1 long choice) and  $X_{-n}$  ( $n$ -back second-order prior) is the interaction term between reward and choice,  $X_{-n} = R_{-n} \times C_{-n}$ . Thus, the variable  $X_{-n}$  is also binary and it takes the value  $X_{-n} = +1$ , when  $R_{-n}$  was correct (incorrect) and  $C_{-n}$  was short (long) and the value  $X_{-n} = -1$  when  $R_{-n}$  was incorrect (correct) and  $C_{-n}$  was short (long). The variable  $S_0$  is the stimulus category (short or long) presented to the rat on the current trial.

For most sessions, the number of trials belonging to class 1 did not match the number of trials belonging to class 2, in other words, conditions were unbalanced. We addressed this problem by subsampling [150, 151], which consists in balancing the number of trials for the two classes by randomly excluding trials from the most populated class. A large imbalance can be problematic when comparing classifier's performance among data sets: if class 1 and class 2 are unbalanced, then Decoding Performance (DP) can be larger than chance ( $DP > 0.5$ ) even when there is no information in any of the regressors. Subsampling was repeated 20 times. Each time the model was trained and tested by 5-fold cross validation. The reported decoding performance (DP; fraction of correct classifications) corresponds to the mean DP over all recording sessions, subsampling and cross-validation iterations. To test statistical significance of each task-variable's weight we used a permutation test that sampled the null hypothesis. For each subsampling iteration (20 iterations) we shuffled choice labels ('short' or 'long' choice) and computed the task-variable weights by 5-fold cross-validation. This procedure was repeated 1,000 times. The null hypothesis distribution for each task variable's weight  $\omega_i$  was the mean absolute value across recording sessions,

subsampling and cross-validation iterations. We defined the probability that a particular variable did not influence the rat's behavior by the fraction of samples that fell above the estimated weight's absolute value. The reported one-tailed  $P$ -values were equal to that fraction. We preferred employing a permutation to test for significance in the regressors against more traditional methods that are based on the assumption that the residuals are Gaussian [121, 137, 140] because the residuals we observed in our data were strongly non-Gaussian. Furthermore, permutation tests are in general more conservative (lower probability of type I errors). Finally, permutation tests sample the null hypothesis while taking into account correlations in the regressors.

### 3.4.4 Psychometric curve analysis

Each rat's psychometric curve was defined as the fraction of long choices over all completed trials (correct trials and false alarms), as a function of the ITI after merging all the sessions for that animal. The all-rats psychometric curve was computed by merging all sessions from all rats. We compared the percentage of correct answers (performance) when trials were easy (ITI = 50, 100, 450 and 500 ms; far from category boundary) against the percentage of correct answers when trials were difficult (ITI = 150, 200, 350, 400 ms; close to category boundary). Significance testing of the difference of animals' performance between easy and difficult trials was based on the non-parametric bootstrap, as follows. We randomly selected with replacement  $k$  trials (where  $k$  is the total number of trials after merging all sessions for a particular animal or all sessions from all animals for the all-rats case) from the set of trials and assessed each rat and all-rats performances on easy and on difficult trials. We repeated this procedure 10,000 times and compared the difference of the resulting two distributions to a reference value, in this particular case zero. We defined the probability that performance on easy trials was equal to performance on difficult trials by the fraction of samples that fell above zero. The reported one-tailed  $P$  values were equal to that fraction. Psychometric curves from trials after correct (error) responses were computed by con-

sidering only those trials that followed a correct (incorrect) response. For each rat and all-rats we compared the psychometric curve after correct trials with the psychometric curve after incorrect trials. Each curve was fitted with the following function [152]

$$P_l(ITI|\mu, \sigma, \gamma, \lambda) = \gamma + (1 - \gamma - \lambda) \left( \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{ITI} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \right), \quad (3.2)$$

where  $P_l(ITI)$  is the probability of long choice as a function of the time difference between tones. The fitted parameters  $\gamma$ ,  $1 - \lambda$ , correspond to the lapse rates for short ITI and long ITI respectively, whereas the parameters  $\mu$  and  $\sigma$  correspond to the centre and the inverse slope of the sigmoid function, respectively. We included lapse rates to avoid biased slope and centre parameter estimates [152]. The parameter estimates corresponded to the maximum likelihood solution of a binomial process with an expected value as a function of ITI defined by equation (3.2). We compared the steepness of the psychometric curve after correct and incorrect responses by means of the difference in inverse slope parameters  $\sigma$  for the two conditions divided by the slope after correct trials (percentage change). Statistical significance was assessed by a non-parametric one-tailed bootstrap (10,000 repetitions), where we assigned uncertainty intervals to the estimated parameters and compared their difference to the reference value zero, as above. To test for significance of performance increase of the psychometric curves computed after incorrect and correct trials we used non-parametric one-tailed bootstrap as described above. The same test was used to test significance for the win-stay and lose-switch probabilities, as well as for testing if they differed.

### 3.4.5 Surgical procedure

Recordings were obtained from three Wistar rats that were chronically implanted with tetrodes in the lateral orbital frontal cortex (IOFC) (see Fig. 3.5a). Animals were trained for 21 days. After 1 week of water and food ad libitum, a microdrive holding the tetrodes was implanted.

To perform the surgery, anesthesia was induced using intraperitoneal injections of ketamine (60 mg/kg) and medetomidine (0.5 mg/kg). The animals were then mounted in a stereotaxic frame, and their skulls exposed. A 3-mm-diameter craniotomy was made, with its center at 1600 microns dorso-ventral, 3.7 mm anterior-posterior and 2.5 mm medium-lateral from bregma [153]. Body temperature was monitored through a rectal thermometer and maintained (36–38 °C) using an electric blanket. Heart rate and blood oxygen levels were monitored. Reflexes were regularly checked during surgery to assure deep anesthesia. Other drugs were given during surgery and recovery period to prevent infection, inflammation, and as analgesia: antibiotics (enrofloxacin; 10 mg/kg sc) and topical application of neomycin and bacitracin in powder (Cicatrín), analgesic (buprenorphine; 0.05 mg/kg sc), anti-inflammatory (methylprednisolone; 10 mg/kg ip), and atropine (0.05 mg/kg sc) to prevent secretions during surgery. Once the animals went through all experimental sessions, they were sacrificed by means of an overdose of pentobarbital (0.8 ml).

### **3.4.6 Tetrodes and microdrives**

Each tetrode was made from four twisted strands of HM-L-coated 90% platinum-10% iridium wire of 17  $\mu\text{m}$  diameter (California Fine Wire, Grover Beach, CA). Gold plating decreased their impedance to ca. 300 – 500 k $\Omega$ . Four tetrodes were held by a cannula attached to a microdrive supplied by Axona (St. Albans, UK). This microdrive allowed for dorsal to ventral tetrode movement to search for new units. Microdrives were attached to the skull with dental cement and seven stainless steel screws. The OFC was reached by vertical descent, and the tetrodes were lowered 1600 microns during the surgery. Vertical descent performed after surgery was of 50 microns/day until the OFC was reached [153]. This depth estimation was verified on by histological reconstruction of the electrode's tracks (see Fig. 3.5a).

### **3.4.7 Electrophysiological recordings from awake, freely moving rats**

During the training period, animals lived in large cages of  $28 \times 42 \times 30$  cm (Charles River) in a rich environment, under a 12:12-h light-dark cycle, and with food ad libitum and water restriction. Before training and after 1 week of postoperative recovery period, the animals were accustomed to the recording chamber. The electrode wires were AC-coupled to unity-gain buffer amplifiers. Lightweight hearing aid wires (2–3 m) connected these to a preamplifier (gain of 1,000), and to the filters and amplifiers of the recording system (Axona, St. Albans, UK). Signals were amplified ( $\times 15,000$ – $40,000$ ), high-pass filtered (360 Hz), and acquired using software from Axona (St. Albans, UK). Each channel was continuously monitored at a sampling rate of 48 kHz. Action potentials were stored as 50 points per channel (1 ms; 200  $\mu$ s prethreshold and 800  $\mu$ s postthreshold) whenever the signal from any of the pre-specified recording channels exceeded a threshold set by the experimenter for subsequent offline spike sorting analysis. Data were excluded if any drift was detected. Before each experimental session, tetrodes were screened for neuronal activity. Once spikes could be well isolated from background noise, the experimental protocol started.

### **3.4.8 Experimental setup**

The recordings were performed inside a black acrylic box of dimensions  $22 \times 25.5 \times 35$  cm. This box was placed inside two wooden boxes placed one inside the other. Between each box, two isolating foam rubbers (4 and 2 cm thick) were placed to soundproof for low and high frequencies. A wooden cover and soundproof foams closed the entire recording chamber, with only a hole to allow the entry of a recording wire (2 mm thick) connected to the preamplifier. Water valves were placed outside the recording chamber. The animals poked their noses into three different sockets (2 cm wide and separated by 3 cm each, and with no cover in the top part to avoid being hit by the microdrive). Recordings were obtained in

darkness, and the experiment was filmed with an infrared camera placed above the recording chamber.

### 3.4.9 Neural data

Recordings were obtained from three Wistar rats that were chronically implanted with tetrodes in their lateral orbital frontal cortex (IOFC) (Fig. 3.5a). We used the pre-stimulus (or trial-initiation), stimulus offset and choice periods for neuronal data analysis. The trial-initiation period starts with the rat nose-poking into the central socket and lasts for 150 ms. The stimulus offset period starts 100 ms before the second tone onset and it lasts until tone offset (150 ms in total). The choice period corresponds to a 150 ms time window that starts with nose-poking into one of the two lateral sockets.

A total of 137 single units were recorded from three rats (53, 62 and 22 from rats 1–3, respectively). On average  $2.9 \pm 1.6$  neurons (max 8) across all rats and sessions were recorded simultaneously. We excluded all neurons firing at  $< 1$  Hz from further analysis, because their low firing rate precluded any reliable statistical analysis. All results remained qualitatively similar when including these cells. For the pre-stimulus, stimulus offset and choice periods, 76 (rat 1: 32; rat 2: 30; rat 3: 14), 87 (rat 1: 35; rat 2: 33; rat 3: 19) and 78 (rat 1: 34; rat 2: 30; rat 3: 14) single-units fulfilled the criterion, respectively (firing above 1 Hz). After filtering out low-activity units, the mean number of simultaneously recorded neurons across all rats and all sessions was  $2.0 \pm 1.0$ . Figures 3.5 and 3.6 were generated using a 100 ms causal rectangular window, sliding in steps of 50 ms. The total mean number of trials across sessions was 684, with an average number of 538 correct and 145 error trials. This led to a median of 9,000 spikes per neuron, before neuron exclusion, and a high-signal to noise ratio quality for hypothesis testing (see section 3.2).

### 3.4.10 ROC analysis

For each neuron we computed the area under the curve (AUC) for a particular task variable as the probability of sampling a larger spike rate  $r$  from  $p(r|z = 1)$  than from  $p(r|z = -1)$ , where  $z$  refers to any of the binary task variables [49, 154]. For AUC values below one half we reversed the populations, to ensure AUCs of at least one half.

### 3.4.11 Generalized linear model for neuronal activity

For the GLM analysis, for each neuron we fitted the spike count in one of the three periods defined previously by

$$n_j \sim \text{Poisson} \left( f^{-1} \left( \sum_{i=0}^k \omega_i x_i \right) \right), \quad (3.3)$$

where the link function  $f(\cdot)$  was taken to be the natural logarithm. The argument of the link function is a weighted sum over an exhaustive family of  $k$  binary regressors

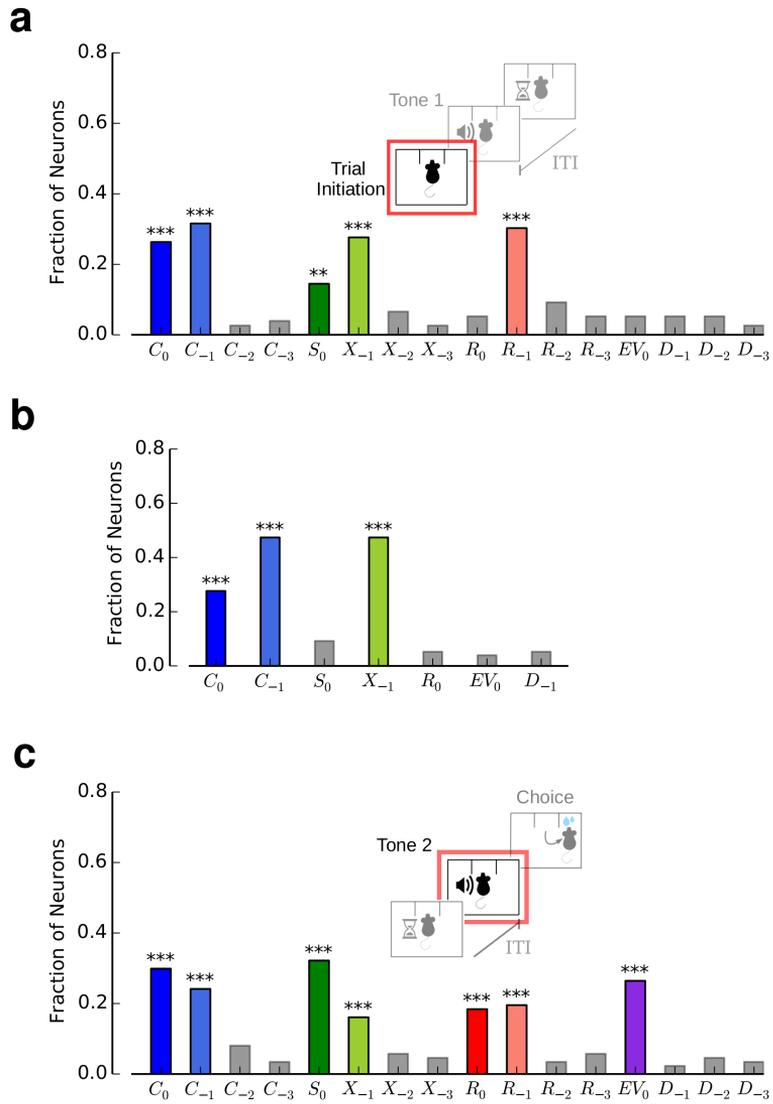
$$\begin{aligned} \sum_{i=0}^k \omega_i x_i = & \omega_0 + \omega_1 R_{-3} + \omega_2 D_{-3} + \omega_3 C_{-3} + \omega_4 X_{-3} + \\ & + \omega_5 R_{-2} + \omega_6 D_{-2} + \omega_7 C_{-2} + \omega_8 X_{-2} + \\ & + \omega_9 R_{-1} + \omega_{10} D_{-1} + \omega_{11} C_{-1} + \omega_{12} X_{-1} + \\ & + \omega_{13} R_0 + \omega_{14} EV_0 + \omega_{15} C_0 + \omega_{16} S_0. \end{aligned} \quad (3.4)$$

Here  $R_{-n}$  is the reward given to the rat  $n$  trials back in time, that is, the correctness of the response (+1 correct, rewarded, -1 incorrect, non-rewarded);  $D_{-n}$  is the trial difficulty defined on the basis of the distance between the presented ITI and the category boundary (50, 100, 450 and 500 ms, easy trial,  $D_{-n} = +1$ ; 150, 200, 350 and 400 ms, difficult trial,  $D_{-n} = -1$ );  $C_{-n}$  is rat's choice (+1 short choice, -1 long choice) and  $X_{-n}$  ( $n$ -back second-order prior) is the interaction term between reward and choice,  $X_{-n} = R_{-n} \times C_{-n}$ . Thus, the variable  $X_{-n}$  is also binary

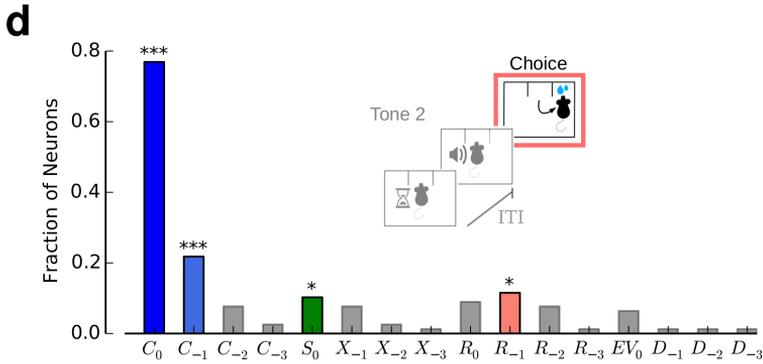
and it takes the value  $X_{-n} = +1$ , when  $R_{-n}$  was correct (incorrect) and  $C_{-n}$  was short (long) and the value  $X_{-n} = -1$  when  $R_{-n}$  was incorrect (correct) and  $C_{-n}$  was short (long). For the current trial ( $n = 0$ ), we renamed difficulty  $D_0$  by  $EV_0$ , and referred to it as expected value, because it is of more conventional use. As  $S_{-n}$  and  $X_{-n}$  are the same variable, we excluded in equation (3.4) the former for past trials and the latter for the current trial.

The GLM fit was applied to different subsets of the data: (i) including all trials (Fig. 3.7) or (ii) including only trials after a correct response (Fig. 3.8) and also to the datasets corresponding to the two passive stages, where the animals were presented the same set of stimuli in a passive manner (Fig. 3.11). In analysis (i), the GLM included all regressors as specified in equation (3.4). For each regressor and neuron, statistical significance was assessed using a permutation test that sampled the null hypothesis. We shuffled each neuron’s spike count across trials and fitted the model on each of 10,000 random shuffles. We defined the probability that a particular regressor was not modulating neuron’s spike count by the fraction of samples that fell above or below the real regressor value for  $\omega_i > 0$  or  $\omega_i < 0$ , respectively. Two-tailed  $P$  values for each regressor and neuron were twice that fraction. The reported fraction of neurons (Figs. 3.7, 3.8 and 3.15) was the number of neurons that had the firing rate significantly modulated by each task-variable over the total number of neurons used in the analysis. We preferred employing a permutation to test for significance in the regressors against more traditional methods that assume that the residuals are Gaussian [121, 137, 140], because the residuals that we observed in our data were strongly non-Gaussian. Furthermore, permutation tests are in general more conservative (lower probability of type I errors). Finally, permutation tests sample the null hypothesis while taking into account correlations in the regressors. Note that it is not necessary to apply Bonferroni correction in our case as we always included all variables of interest in the GLM simultaneously rather than running individual tests for each variable separately.

In analysis (ii) only regressors from the previous and the current trials were included, except for  $R_{-1}$  which, by construction, was constant for



**Figure 3.15:** Using a linear regression instead of a GLM on the data presented in Fig. 3.7 and Fig 3.8.



**Figure 3.15:** (cont.) (a-d) Linear regression-based analysis gives virtually identical results to those in Fig. 3.7 and Fig 3.8 obtained from a GLM analysis. Panels a, c, and d correspond to Fig. 3.7a-c. Panel b corresponds to Fig. 3.8. \* =  $P < 0.05$ , \*\* =  $P < 0.01$ , \*\*\* =  $P < 0.001$ . ( $n = 76$  for panel a and b,  $n = 87$  for panel c and  $n = 78$  for panel d).

this particular set of trials. Regressors from previous trials were not included to avoid overfitting due to the reduced set of trials for this analysis. After correct trials, regressors  $C_{-1}$  and  $X_{-1}$  were equivalent and the pair was treated as a single-variable. In Fig. 3.8 fractions of neurons encoding  $C_{-1}$  and  $X_{-1}$  were reported separately only to allow a better comparison with Fig. 3.7. Significance of each regressors was tested using a permutation test. We also fitted the GLM using only trials after an incorrect response. The procedure was identical to (ii) but in this case, because of the experimental protocol,  $-C_{-1}$  and  $X_{-1}$  and  $S_0$  were identical and  $EV_0$  and  $D_{-1}$  were identical as well.

For the passive stages all trials were used. The set of regressors in this particular case comprised current stimulus  $S_0$ , current expected value or difficulty  $EV_0$ , second-order prior  $X_{-1}$  (from  $-1$  to  $-3$  trials in back) and previous difficulty  $D_{-1}$  (from  $-1$  to  $-3$  trials in back as well). It is important to note that because in the passive stage rewards are not delivered, the second-order prior variable  $X_{-1}$  is undefined. However, in the decision-making stage the second-order prior variable is equivalent to the previous stimulus for all trials, that is,  $X_{-1} = S_{-1}$ . Thus, we

take  $S_{-1}$  in the passive stages as the analogous to the state-space in the decision-making task. The reported fraction of neurons (Fig. 3.11) was the number of neurons that had the firing rate significantly modulated by each task-variable over the total number of neurons used in the analysis. Significance for each regressor was calculated as described above.

For each regressor a binomial test was used to assess if the fraction of neurons that had their firing rates modulated by that particular regressor was significantly greater than chance [121, 140] (5%; one-tailed). Statistical significance for the difference in fractions between two conditions was tested by a non-parametric difference binomial test that sampled the null hypothesis as follows. Independent samples from two identical binomial distributions were drawn 10,000 times and the null hypothesis was built as the difference of these binomial processes. The expected values of the two identical binomial processes were the weighted mean of the two fractions to be compared. We defined the probability that the two fractions were instances of the same underlying binomial process by the proportion of samples that fell above the observed fraction difference. The reported one-tailed  $P$  values corresponded to that proportion. One-tailed  $P$  values were used instead of two-tailed  $P$  values because the study's hypothesis was to test whether previous trial regressors (such as previous choice  $C_{-1}$  or previous second-order prior  $X_{-1}$ ) were decreasing over the course of the trial, and whether upcoming choice  $C_0$  was increasing as rats went through trial's stages. For the case of upcoming stimulus  $S_0$  and upcoming expected value  $EV_0$  our hypothesis was that they had to peak during the stimulus presentation period.

It is important to note that it is not possible to directly compare the fractions of neurons with significant regressors after correct, incorrect or all trials, because of the large difference on the correlation structure among regressors across conditions. First, several task variables that are different on after-correct trials become the same variable for after-incorrect trials, and vice versa. For instance,  $X_{-1}$  and  $C_{-1}$  are the same variable after correct trials, while after incorrect trials  $X_{-1}$ ,  $-C_{-1}$  and  $S_0$  are all three the same variable, and  $EV_0$  and  $D_{-1}$  are again the same. In addition, as depicted in Fig. 3.1d, rats after an incorrect response tend to

switch choice more often than repeat the same choice after a correct response. Therefore, the regressor  $C_{-1}$  is more strongly correlated with  $C_0$  after an incorrect response than after a correct response. The differential increase of correlations between regressors, when conditioned after correct or incorrect trials and the resulting differential biases obtained from fitting a model precluded a direct comparison of the reported fractions of significant neurons across conditions.

### 3.4.12 Correlation of regression weights

We tested the stability of the neuronal representations over time by correlating the fitted values of weights in the GLM across different time periods. Correlations among weights could simply arise because of different responsiveness of the neurons, such that for instance when a neuron that is more responsive in the pre-stimulus period might also be more responsive in the offset stimulus period. To avoid creating correlations due to differences in overall firing rate across neurons in the population, we first normalized each firing rate by subtracting and dividing it by its mean and s.d. respectively ( $z$ -score) for a particular time window. This normalization can result in negative normalized rates, violating the assumptions of the previously used GLM model since a natural logarithmic function was used (equation (3.3)). To overcome this problem, we instead fitted the data by linear regression (see section 3.4.11). Fig 3.15 shows that using linear regression instead of a GLM (Fig. 3.7) does not qualitatively change the results. Subsequent analysis for correlated weights was performed on the linear regression coefficients, using the same set of regressors, equation (3.4), as for the GLM.

Stability of the neuronal representation for each variable (for example, the upcoming choice  $C_0$ ) across the trial was assessed by using the correlation coefficient (Pearson correlation) between two vectors, each with the  $i$ -th entry being the regression coefficient for that variable (for example, upcoming choice  $C_0$ ) of neuron  $i$ , computed at two different periods, namely pre-stimulus and stimulus offset periods (Fig. 3.9a) or stimulus offset and choice periods (Fig. 3.9b). Statistical significance of the

correlation coefficient was assessed by a permutation test that sampled the null hypothesis. For each regressor (for example, upcoming choice  $C_0$ ) the null-hypothesis distribution was built from the set of correlation coefficients obtained after shuffling the relationship between each neuron's  $z$ -scored firing rate and the regressor, and computing their respective Pearson correlation coefficient as before. This process was repeated 10,000 times. We defined the probability that a particular regressor was not stable across time by the fraction of samples that fell above the real correlation coefficient value (if  $\rho > 0$ ) or below the real correlation coefficient value (if  $\rho < 0$ ). The reported two-tailed  $P$  values for each regressor were twice that fraction.

We tested whether the second-order prior and upcoming choice at trial initiation are encoded by the same neurons. Unfortunately, we cannot use the same approach as just described, as computing the vectors of the regressors across neurons for both  $X_{-1}$  the  $C_0$ , and then computing the correlation coefficient between them will lead to biases due to using two regressors from the same model in the same dataset [140]. We avoided this problem by instead computing regression weights for each variable while fixing the value of the other variable, as follows. We first restricted our analysis to trials that followed a correct response and focused on the pre-stimulus period, where only information about two variables is found,  $C_0$  and  $X_{-1}$  (see Figs 3.8 and 3.15b shows how the linear regression model gives qualitatively similar results as the GLM model when focusing on trials that followed correct responses). The weights for  $C_0$  were therefore computed by fitting the model on the subset of trials where the variable  $X_{-1}$  was constant ( $C_{-1} = X_{-1}$  for this particular set of trials; see section 3.4.11). This conditioning procedure ensured that the estimated weight for  $C_0$  was not affected by its intrinsic correlation with  $X_{-1}$ . Because  $X_{-1}$  is a binary variable, the reported weight for  $C_0$  was the mean between the weight estimated for set of trials where  $X_{-1} = +1$  and where  $X_{-1} = -1$ . The same procedure was applied for the weight associated to  $X_{-1}$ , where again the final weight for this variable was the mean between the weight fitted on the subset of trials where  $C_0 = +1$  and  $C_0 = -1$ . The reported correlation coefficient was computed from two vectors, one composed of

the mean weight for  $C_0$  (mean across conditionings and  $X_{-1} = +1$  and where  $X_{-1} = -1$ ) of each neuron  $i$ , and the other composed of the mean weight for  $X_{-1}$  (mean across conditionings  $C_0 = +1$  and  $C_0 = -1$ ) of each neuron  $i$ .

Statistical significance of the correlation coefficient was again assessed by a permutation test that sampled the null hypothesis. The null-hypothesis distribution was built from the set of correlation coefficients obtained after shuffling the relationship between each neuron's  $z$ -scored firing rate and the regressors, and yielded one correlation coefficient sample by following the same computations as described in the previous paragraph. This process was repeated 10,000 times. We defined the probability that neurons encoding  $C_0$  do not tend to encode  $X_{-1}$  by the fraction of samples that fell above the real correlation coefficient value (if  $\rho > 0$ ) or below the real correlation coefficient value (if  $\rho < 0$ ). The reported two-tailed  $P$  values for each regressor were twice that fraction.

### 3.4.13 Population decoding

Small populations (two or three neurons) of simultaneously recorded single-neurons were used to classify a set of trials as belonging to either class 1 or class 2 (for example, class 1 and class 2 can correspond to short and long choices for the variable  $C_0$ , or to correct and incorrect responses for the variable  $R_{-1}$ ). Classification is based on a decision variable  $DV$ : when  $DV > 0$  the trial is classified as class 1, and when  $DV < 0$  the trial is classified as belonging to class 2. The decision variable  $DV$  is a weighted sum of the population activity  $DV = \sum_{i=1}^{i=N} \omega_i r_i + \omega_0$ , where  $\omega_i$  and  $r_i$  are each neuron's contribution to the decision variable and spike rate respectively,  $\omega_0$  is the offset term, and  $N$  is the total number of neurons used in the classifier. Logistic regression assumes that the probability of class 1 to be the correct class given the activity pattern of the population is given by  $p(\text{class1}|\{r_i\}) = \sigma\left(\sum_{i=1}^{i=N} \omega_i r_i + \omega_0\right)$ , where  $\sigma(\cdot)$  is the logistic function. The model was trained and tested using five-fold cross validation.

For most sessions, the number of trials belonging to class 1 did not

match the number of trials belonging to class 2, in other words, conditions were unbalanced. We addressed this problem by subsampling [150, 151], which consists in balancing the number of trials for the two classes by randomly excluding trials from the most populated class. A large imbalance can be problematic when comparing classifier's performance among data sets: if class 1 and class 2 are unbalanced, then Decoding Performance (DP) can be larger than chance ( $DP > 0.5$ ) even when there is no information in any of the regressors. Subsampling was repeated 20 times. Each time the model was trained and tested by 5-fold cross validation. The reported decoding performance (DP; fraction of correct classifications) corresponds to the mean DP over all recording sessions, subsampling and cross-validation iterations.

Statistical significance of DP was tested using a permutation test that sampled the null hypothesis. For the set of trials (the whole recording session when class 1 and class 2 were balanced and the particular subsampling iteration when class 1 and class 2 were unbalanced) we shuffled each trial's class label and estimated DP through the five-fold cross-validation method (20 repetitions for the subsamplings). This procedure was repeated 1,000 times. Each of the samples of the null hypothesis distribution was computed as the mean across recording sessions, subsampling and cross-validation for a particular shuffling iteration. We defined the probability that the neuronal ensemble had no information about that particular task variable by the fraction of samples that fell above the real DP. The reported one-tailed  $P$  values were that fraction.

#### **3.4.14 Conditioned population decoding**

As many of the variables are partially correlated (for example, choice with stimulus), being able to decode one of them necessarily means that we can decode the others. To test if we can read out both of a pair of partially correlated variables independently, we performed a conditioning decoding analysis in which we tested for information of one variable while keeping the values of the other variable fixed (Fig. 3.12). We restricted our analysis to trials after correct responses. As shown in Fig. 3.8, the

GLM analysis revealed that single-neurons seemed to encode only two variables: upcoming choice  $C_0$  and second order prior  $X_{-1}$ . We therefore decoded upcoming choice  $C_0$  by fitting a classifier on the subset of trials where  $X_{-1} = +1$  and  $X_{-1} = -1$  independently (subsampling method and five-fold cross validation, see section 3.4.13). The reported DP when classifying upcoming choice given second-order prior was the mean between the two conditioned DP. To decode  $X_{-1}$  the same procedure was applied but conditioning on each of the two possible values of  $C_0$  instead. The reported DP when classifying second order prior given upcoming choice was the mean between the two conditioned DP. In this way, even though decoded quantities might be correlated, reported population information content about  $C_0$  and  $X_{-1}$  could not be explained simply by a correlation to other variables (Fig. 3.12).  $P$  values were computed using a permutation test, as described in section 3.4.13.

### **3.4.15 Information ranking**

We used decoding performance (DP) for each variable that was deemed significant by the GLM analysis as a proxy for the amount of information that the neuronal population contained about that variable (Fig. 3.7). DP is computed as described above. Our analysis provides the intuitive result that decoding performance increases with the number of neurons in the ensemble (Figs 3.12 and 3.13). Some previous population analysis violated this due to misusing linear classifiers [155].

### **3.4.16 Correlation between behavior and neuronal activity across rats**

We correlated rats' probability of switching choice after an incorrect response (lose-switch probability) with the Decoding Performance of both upcoming choice  $C_0$  and second-order prior  $X_{-1}$  calculated during the trial-initiation period across rats. Three different sizes of neuronal ensembles were used for this analysis: single units, pairs and triplets of neurons. The reported DPs for each rat were the mean DP across all the

groups of simultaneously recorded neurons of a given size for that particular animal. Significance was tested using a permutation test that sampled the null hypothesis. For each ensemble size we shuffled the probability lose-switch vector and the Decoding Performance vector. This process gave us a total of thirty-six possible combinations ( $3! \times 3!$ ). We defined the probability that there was no correlation between the two vectors by the fraction of samples that fell above the real correlation value. None was above the real correlation value but six samples were equal. Therefore, we considered that out of these six, three were above and three were below the real correlation value. The reported two-tailed  $P$ -values were twice that fraction.

### 3.4.17 Power Analysis

To address whether the sample size used in this study was sufficiently large, we performed a power analysis. The power ( $\pi$ ) in hypothesis testing is defined as the probability of correctly rejecting the null hypothesis (when it is false). It is related to the  $\beta$ , the probability of a type II error, through the equality  $\pi = 1 - \beta$ . We first calculated the statistical power of the fraction of neurons that had their firing rate significantly modulated by each one of the reported regressors. Because significance for each neuron and regressor was calculated using a Binomial test (one-tailed), the power analysis was calculated using a Binomial distribution as well. First, we calculated the threshold on the fraction of neurons that would make us reject the null hypothesis when it is true. This value was the smallest fraction of neurons that fulfilled  $P < 0.05$  for a one-tailed Binomial test. As this quantity only depends on the total number of neurons used in the Binomial test, it is the same across all regressors. We calculated the statistical power of the fraction of neurons encoding for a particular regressor as the probability of sampling a largest or equal value than the threshold from a binomial distribution, with probability of success defined as the actual fraction of neurons reported in the study. For a fixed probability of type I error ( $\alpha = 0.05$ ) the statistical power increases with the strength of the result itself (reported fraction of neurons) and the number of samples

used in the Binomial test (number of neurons). In other words, the probability of making a type II error when inferring the significance of that result decreases if these quantities increase. The same procedure was applied to the statistical power analysis associated with the number of rats. Every rat that showed results in favor of the hypothesis was considered a success, and otherwise it was a failure. Because all rats in the study belonged to the success class (our main results are consistent across rats) we calculated the power if an additional rat without significant results (belonging to failure class) was added to the study. We considered a result had sufficient statistical power when it exceeded the standard threshold  $\pi = 0.80$ . For the number of single units ( $n = 76$ ), the statistical power of the fraction of neurons encoding  $C_0$  before stimulus presentation was  $\pi = 0.9996$ , and  $\pi = 0.9947$  for  $X_{-1}$ . Thus, the two central results of our study are statistically rather solid. Regarding the number of rats, we found significant results consistently across all rats. If we included a fourth rat with negative results the power of our conclusions would be  $\pi = 0.95$ . Thus, despite the small sample size, the strength of the results supported sufficient statistical power.

# Chapter 4

## DISCUSSION

The main goal of theoretical neuroscience is to determine the fundamental principles of information processing in the brain and ultimately characterize the link between neuronal activity and behavior. Even though this dissertation is very far away from providing general principles and closed solutions, trying to provide an answer to this question has been the driving-force of the studies presented in this thesis. In particular I have tried to provide mathematical and experimental evidence on how information is processed in the brain of rats and monkeys and how the representation of this information correlates with behavior. In the following sections I will first discuss and further contextualize the experimental findings reported in chapters 2 and 3. Then I will present in brief what future experiments and lines of research the experimental findings herein presented suggest to perform and follow. Finally I will comment and give my personal thoughts and predictions on a historical and philosophical contextualization of science, or in other words, on the methodology for acquiring valid and true knowledge of the hidden rules governing our world.

## 4.1 Encoding of Information

Identifying the features of the neural code that affect the amount of information encoded by a neural network is an important step on the quest for fully characterizing the link between neuronal activity and behavior. In this chapter we identified the selectivity length (SL) and the projected precision (PP) as the only features affecting information by deriving an analytic expression on the decoding performance (DP) of a binary linear classifier. We tested our prediction on four datasets encompassing two different brain areas (middle temporal, MT; and lateral prefrontal cortex 8a, LPFC 8a) and three different tasks (a coarse and a fine visual discrimination task and an attentional task) using a novel method based on perturbing the original dataset. The perturbation method consisted on randomly selecting with replacement trials from a particular dataset (bootstrap) which could be seen as a procedure to create virtual instances of the same experiment. By selectively conditioning to those perturbations that fixed the features of the neural code we were not interested in, we further confirmed that SL and PP were the only features affecting the amount of information encoded by the neural network. Moreover, SL and PP were found to be the features affecting the most the behavioral performance of the monkeys. Finally, we showed that a simple model of the cerebral cortex was able to reproduce our findings.

Although many different features of the neural code could be defined besides SL and PP, this chapter focused principally on discrediting mean pairwise correlations as affecting the encoding of information and behavior. This was motivated, partially, because of the classical idea that mean pairwise correlations play an important role on information and behavior [34, 35, 58, 60]. The idea herein proposed is not entirely novel and, even though at first glance it might look a bit redundant with respect to some existing literature [47, 48, 113], it constitutes an important contribution to the field mainly because it is based on experimentally-realistic information measures and ensemble sizes. Let's have a closer look to what I mean by experimentally-realistic information measures and how does it relate to previous studies. As stated in chapters 2 and 1, most of the previous

theoretical studies are based on using (linear) Fisher information (FI) as the proxy for the amount of information encoded by a network. Yet this is the information metric that has to be used when dealing with continuous parameter estimates, it is not very well suited for estimations that have to be performed on a discrete set of values, which characterizes the majority of the experimental protocols used in the field. Not knowing the exact shape of the tuning curve  $f(s)$  or its derivative with respect to the stimulus  $f'(s)$  can produce wrong estimates of FI, as  $I_F(s) = \mathbf{f}'^T \Sigma^{-1} \mathbf{f}'$ . The natural extension of linear FI when dealing with discrete parameter estimates is the DP of a linear classifier. Interestingly, for the discrete binary case, it is easy to show that the variability on the estimate becomes  $\text{Var}(\hat{s}) = \frac{N_{Error}}{N} = 1 - \text{DP}$ , where  $\hat{s} = \{0, 1\}$  is the inferred parameter, and  $N_{Error}$  is the total number of mistakes performed by the classifier out of  $N$  trials. As FI's definition is based on the variability on the estimate of an unbiased classifier, FI and DP are highly related information measurements. Although important theoretical contributions have been made in order to correct for possible biases when assessing FI on electrophysiological recordings [156], using the DP of a linear classifier is a more natural way of evaluating encoded information on experiments involving neuronal recordings and a discrete set of experimental conditions.

The other contribution of this study with respect to previous theoretical works is that it analyzes the role of pairwise correlations for fixed and experimentally-realistic ensemble sizes. Previous theoretical studies were in part motivated by the experimental finding that the amount of information encoded on single neurons did not largely outperformed the performance of behaving animals [49] (see chapter 1). Roughly speaking, the solution they proposed was based on introducing pairwise correlations in order to limit the amount of information encoded in a neural network when  $N \rightarrow \infty$  [35, 47]. Even though these studies represent in my opinion outstanding contributions to the field of theoretical neuroscience, they did not explicitly explore what was the role of pairwise correlations for fixed and middle-sized populations of neurons. In chapter 2 we provided an answer to that question by showing through an analytic expression for the DP of a linear classifier and a non-parametric perturbation method

based on bootstrapping and conditioning, that only SL and PP have an effect on the amount of encoded information.

In this part of the thesis we also showed that global activity played no role on the amount of information encoded by a neural network. However, it is worth mentioning that there is an important difference on how previous works have studied mean pairwise correlations and global activity. While mean pairwise correlations have been historically approached by their impact on the encoding of information, global activity has been generally used to fit the statistics of the neuronal code regardless of its functional role [42–44]. By showing that mean pairwise correlations had no influence on the decoding performance of a linear classifier we complemented some seminal theoretical studies [35, 47, 48] and contradicted some experimental ones [58, 60]. However, by showing the equivalent result for global activity, we provided further experimental evidence supporting one of the main claims of the study [51].

One of the greatest differences between the results presented in chapter 2 and previous studies regarding global activity [42–44, 51] lies on the fact that in our approach we present results of how global activity can affect both encoding of information and behavior. Even though [51] shows how this feature of the neural code affects the DP when using ensemble sizes in the order of  $\sim 100$  neurons, this study is novel on analyzing how global activity affects the performance of behaving monkeys. Our results indeed showed a modulation on performance with respect to global activity on monkey 2, but it was not consistent across the rest of datasets. On the contrary, we found that selectivity length and projected precision were generally the most influential factors on behavior consistently across the four datasets. This result provides further evidence, although not in a conclusively manner, that information in the cerebral cortex is read-out optimally by downstream units.

As stated earlier in this thesis, previous literature has indeed explored how mean pairwise correlations affected the performance of behaving animals [58–60]. Roughly speaking these studies reported that performance increased with attention, which in turn modulates behavioral performance. Therefore they claimed for a functional role of mean pairwise

correlations on the amount of information encoded by a neural network in an indirect manner. Although this claim seems to be in odds with the results reported in chapter 2, we think that they find this indirect correlation between mean pairwise correlations and performance due to the strong link that exists between projected precision and mean pairwise correlations. As our perturbation method is able to condition on the different features of the neural code, it would be enlightening to test on their datasets whether projected precision or mean pairwise correlations affect the most behavioral performance when selectively conditioning to the undesired quantity. Indeed, in [59] it is reported that similarly tuned neurons decrease their correlation with attention whereas dissimilarly tuned neurons increase mean pairwise correlations with attention. Attention here is directly linked with behavioral performance as attending to a stimulus increases the performance on a perceptual decision-making task. This is exactly what one would expect if the performance is directly linked with projected precision, being the correlation with mean pairwise correlation just an illusion arising from the relationship between these two features of the neural code.

Another important aspect to be mentioned in this discussion is how our approach is related with the well-known method for studying how mean pairwise correlations affect the encoding of information: destroying pairwise correlations by shuffling the activity of each neuron across trials for a fixed experimental condition. If the amount of encoded information is characterized before ( $I_{unsh}$ ) and after the shuffling procedure ( $I_{sh}$ ), the sign of  $\Delta I \equiv I_{unsh} - I_{sh}$  will determine whether correlations are helpful or harmful for encoding [61, 105]. In my opinion this is a very powerful non-parametric method to address the problem. The perturbation approach taken in chapter 2 complements the shuffling procedure because it is able to change (perturb) as well in an offline manner a feature of the neural code. The two approaches have to be thought as complementary rather than antagonistic. Additionally, by using the analytic expression on the DP of a linear classifier (encoded information) we would be able to assess the effect of removing correlations on the amount of encoded information by evaluating the projected precision and its equivalent after

destroying pairwise correlations (see chapter 2). This could be accomplished just by removing the off-diagonal terms in the precision matrix (inverse of the covariance matrix) and evaluating its projection into the stimulus axis, which is independent from the trial-by-trial shared and individual noise of the population of neurons.

By building a very simple model of the cortex we were able to qualitatively reproduce the experimental findings reported in this part of the thesis. In general terms, the model consisted on a population of heterogeneous Poisson neurons that encoded on their firing rate the value of a hidden parameter, the stimulus presented to the network on a trial-by-trial basis. Additionally, the stimulus presented was added with a term of noise that produced differential correlations [47] and a common additive and multiplicative gain modulated the activity of the neuronal population on each trial [51]. The introduction of differential correlations was very convenient because it made information to saturate for large ensemble sizes ( $N \rightarrow \infty$ ) and therefore we could reconcile the well-known experimental finding that behavioral performance is not outstandingly superior to the amount of information encoded by single-neurons [49]. We also simulated the behavior of a virtual agent by reading out optimally this pool of neurons on each trial. By analyzing this surrogate dataset using the same set of techniques we were able to report qualitatively equivalent results as those reported for the experimental recordings. This was a very relevant finding as it proved that the set of results reported in chapter 2 can be obtained when considering a very simple encoding-decoding scheme of the cerebral cortex, where encoding neurons are read-out by downstream units to produce behavior.

## 4.2 Prior Information

Integrating sensory with prior information is a fundamental feature of adaptive behavior. Although it is not clear yet whether this integration is performed by an optimal combination of these two sources of information or by simple phenomenological rules that work relatively well in a

particular environment (see chapter 1), experimental evidence indicates indeed that animals combine them to guide their behavior [72, 73].

Prior information is a term encompassing a wide range of possible sources of information. For instance, it can consist on providing the set of presented stimuli with some statistical regularities [72] or on coupling the environment with the behaving agent's preceding choices and outcomes [15], among others. In this thesis our goal was to study prior information when the stream of upcoming stimuli is coupled with the most recent set of choices and outcomes. In particular we defined a perceptual decision-making task where the same stimulus was presented after an incorrect response and a random stimulus was presented after a correct response. Coupling the features of the environment (presented stimulus) with the most recent set of choices and outcomes corresponds in my opinion to a more naturalistic experimental protocol as real-life decisions and their associated outcomes have an impact on the future state of the environment.

Under this environment we found that rats' behavior was consistent with an integration of sensory with prior information as an important factor guiding their choices. We used three different approaches to test this hypothesis. First we compared the psychometric curve after a correct and an incorrect response and found that all rats consistently underwent an increase in sensitivity after a mistake. If after an incorrect response the previous choice and its associated outcome are predictive for the upcoming stimulus, the obtained results are consistent with the fact that rats are indeed integrating prior information with sensory information while performing the task. The next analysis we performed on behavior was based on evaluating the probability of making the same choice after a correct response (win-stay; WS) and the probability of changing response after a mistake (lose-switch; LS). Based on the definition of our experimental protocol the optimal behavior in our environment was characterized by  $P(WS) = 0.5$  and  $P(LS) = 1.0$ . All rats showed a very strong component of the lose-switch strategy and a weak, but significantly different from 0.5, component of win-stay. Finally, we fitted the choice of each rat on a trial-by-trial basis using a logistic regression, where we used as regressors the set of variables we thought could be affecting their response.

As expected, we found that presented stimulus was the variable with the largest explanatory power on the choice of all rats. Interestingly, we also found previous choice, previous reward and a second order combination of these two variables as regressors with an important predictive power. It has been formerly reported that previous choices and their associated outcomes can shape the response of behaving animals [74, 75, 115], however the novelty of this study relies on the fact that previous choices and outcomes were incorporated in the form of sequential prior information in a perceptual decision-making task.

Besides studying behavior in a protocol where prior information is defined as the recent choice and its associated outcome, we also aimed to characterize what brain region is the responsible for combining sensory with prior information. As stated in chapter 1, the prefrontal cortex (PFC) is a region of the brain that plays a fundamental role in cognition in general and decision-making in particular. Within this region, located in the ventral surface of the frontal lobe (in primates, see chapter 1), it can be found the orbitofrontal cortex (OFC). The OFC has been shown to encode a large number of variables related with decision-making: expected outcomes [86], cues associated with particular rewards [87] and the values of presented and chosen goods [80]. Altogether these set of results seem to indicate that the OFC plays an important role in economic-decision making [82, 85]. However, recent experimental findings have found the OFC to encode choice-related variables in both rats [91, 93] and monkeys [94]. This opens the door to a reevaluation of the role of the OFC in decision-making that is able to account for the classical and the recent findings.

A recent study proposes that the role of the OFC is to represent the current state within a cognitive-map of the task [97]. This concept is tightly linked to the concept of state representation in the Reinforcement Learning (RL) literature. Even though this theory is not very well reconciled with the classical role of OFC in economic decision-making, it can explain a large set of classical findings like reversal learning [128, 130, 157, 158], delayed alternation [159], extinction [160] and devaluation [161]. In all of these cases the OFC would be critical for efficiently

representing the state space, in particular the OFC would be necessary for disambiguating states that are perceptually equivalent but correspond to different positions of the state space. For example, in some particular instances of the reversal learning task [158], the state space would be fully determined by two states:  $s = 1$  represents the situation where option A is rewarded and option B is not rewarded; and  $s = 0$  represents the opposite contingency. Alternating between the two states is relatively easy and fast in a RL paradigm with the states defined as above, however it becomes much more difficult when just a single state is defined and the outcomes associated with each action have to be learned and unlearned in each alternation of the action-reward contingency.

This theory is in high agreement with the electrophysiological results reported in our study. We found that the rats' OFC was principally encoding previous choice ( $C_{-1}$ ), previous reward ( $R_{-1}$ ), previous second order history ( $X_{-1} = C_{-1} \times R_{-1}$ ), presented stimulus ( $S_0$ ) and upcoming choice ( $C_0$ ). The encoding of previous trial variables decreased during the course of the trial whereas information about presented stimulus peaked during the presentation of the acoustic information. Surprisingly, we found information about upcoming choice before the stimulus onset. The build-up of the choice-related signal suggests that OFC would be responsible for integrating sensory with prior information in the form of previous choice and reward. Under the RL paradigm, these results indicate that the OFC represented the current location within the state space by disambiguating the current sensory with prior information (previous choice and reward). The existence of upcoming choice information further validates their hypothesis on our dataset because this signal can be understood as the integrated version of the state-space variable that is build up on the basis of the information available from all the different sources.

In [97] an interesting prediction is made: if the OFC is indeed representing the current location within the state space, then information encoded in this brain area should satisfy both the *representation* and the *specificity* condition. The *representation* condition states that the OFC should encode all the variables involved in generating a full representa-

tion of the state space. This condition was fulfilled on our results because of the set of variables we found encoded in the OFC (see above). The *specificity* condition states that this representation should be, so to say, compact, or in other words, the *specificity* conditions states that all the irrelevant task variables should be filtered out from being encoded in the OFC as they are not necessary for representing the state space. As our behavioral protocol represented a Markov chain, information from two trials in back or more should be discarded because they are irrelevant for characterizing the current location within the state space of the task. In our study we only found encoding of very recent (previous trial) choices and rewards. Moreover, we also analyzed what variables were represented in the OFC when rats were just passively exposed to the acoustic information. In this particular situation the state space could be fully determined by just the presented stimulus. We indeed found that neurons in the OFC only encoded for current stimulus under this protocol while all the irrelevant information (the set of previous stimuli) was filtered out. Even though the electrophysiological results reported in chapter 3 are highly aligned with the claim of study [97], I think it is important to remark here how does it relate with the classical functional theory of the OFC as being responsible for the encoding of expected value and reward. In the RL literature, the expected value associated to a state  $s$  is represented by the state value  $V(s)$ . However, the OFC to be representing  $V(s)$  is in contradiction with the statements of [97]. In this sense these two theories seem to be in a clear disagreement. Nevertheless, when analyzing what quantities were encoded in rats' OFC during the performance of the task, we found that at the time when the stimulus was presented, the OFC had information about the difficulty of the present trial. In other words, once the stimulus was presented to the rat, their OFC was encoding for expected reward as it can be defined as the total amount of reward to receive times the probability of obtaining it (difficulty of the trial). In this sense our results were also in agreement with the classical functional role associated with the OFC.

I would like to end this section with some personal comments on this topic. Brain anatomy is not homogeneous but there are large differences

across species. The OFC in rats has a very different cytoarchitecture than the OFC in primates [162, 163]. This is the reason why some researchers have questioned the idea that it exists an analogous of the primates' OFC in rats (and many other frontal regions) [162]. Interestingly, the classical functional role of the OFC relies generally on experimental evidence collected in monkeys and humans while modern hypothesis about the OFC being involved in the encoding of the current location within the state space [97], action selection and initialization [91, 93] and representation of the time delays [92] are most often performed on rats. Even though my opinion is not intrinsically novel [82], I think that most of the experimental contradictions that have been reported while trying to find a unified theory of the OFC are principally based on the differences across species, in particular between the OFC of rats and primates. Characterizing the functional role of the OFC (and many other brain areas) would perhaps be easier when hypothesis were limited to specific species or at least to very similar groups of them.

### **4.3 Future research**

It is my intention to discuss in this section the two general kinds of research I would like to perform in the near future. First I will start with a set of ideas that arise as a natural consequence of chapters 2 and 3. Then I will present in brief the research that, although it is inspired, is not directly linked to this dissertation but could be very beneficial for a deeper understanding of neuroscience.

The most natural extension of chapter 2 would involve testing our predictions in more tasks, species and brain areas. If we have to find general computational principles in the brain, our predictions should be validated in the largest amount of possible datasets. For instance, showing that selectivity length and projected precision are the most important features on information encoding and behavior across the auditory, visual, somatosensory and olfactory cortices would provide our results with much more robustness. Obviously, in order to study whether each of

these sensory channels fulfill our predictions we should have to design perceptual decision-making tasks corresponding to that particular sense and perform electrophysiological recordings on the corresponding primary sensory cortical areas.

Another interesting path of research derived from the study presented in chapter 2, would be analyzing whether the information in the brain is extracted optimally or sub-optimally from the pool of encoding neurons. Nevertheless, it is not a trivial question to be answered experimentally. One of the most relevant studies in our field aiming to provide a satisfactory answer to this question was published in 2015 by Pitqow and colleagues [111]. Even though their approach was quite convoluted and several approximations had to be done for the sake of mathematical tractability, they claimed for an optimal read-out of the sensory neurons under the presence of differential correlations [47] to account for the experimental results. They approached the problem by analyzing the relationship between the trial-by-trial noise of sensory neurons and monkeys' choice (choice correlation or choice probability). Our approach would be rather different, we would use the perturbation method (bootstrap) to study how perturbations in the amount of encoded information would correlate with perturbations on the performance of the behaving animal. We could derive then the mathematical expressions corresponding to the different read-out strategies and compare it with the experimental result. Despite the efforts in this direction, we have not obtained conclusive results yet. Consequently, it would be very interesting to analyze how the perturbation method could be applied to the study of choice correlations.

Although the relationship between encoding and behavior variability has been extensively studied in a trial-by-trial basis, [111, 154, 164], it has never been addressed by the perturbation method described in this thesis. By combining it with the conditioning method, we could dissect what are the real factors influencing the decision, in the same way as we dissected what were the most important factors on the performance of the animal. In chapter 2 we studied encoded information from a static point of view. Information time profiles are often very enlightening because they can reveal rich hidden dynamics or processes that can be very insightful. De-

spite of being computationally very expensive, making the same analysis on discrete time windows across the course of the trial could be a very interesting path to take as well.

Regarding the second scientific project of this dissertation, many interesting future extensions can be also considered. As the number of subjects and neurons is not very extense, complementing chapter 3 with a larger sample would provide further strength and validity to our claims. However, as detailed in the methodology, a statistical power analysis performed on the number of animals and neurons revealed that the sample size was sufficiently large due to the magnitude of our results.

In chapter 3 we claimed for the OFC to be the responsible of integrating prior with sensory information. Although strong correlations can sometimes indicate causality, it is very difficult to conclusively show that a particular brain region is necessary for a specific brain function. Nowadays, the best way to test for causality relies on inactivating the particular brain area and analyze the behavioral impairments. This could be accomplished by removing the hypothesized brain area or by silencing it with cryogenic, pharmacological or optogenetic techniques. This should be, in my opinion, one of the most important steps to take in the future for this line of research. By studying the behavioral effects of shutting down the OFC we could conclusively determine whether the reported findings are just phenomenological correlations or if they correspond to a truly causal connection between the OFC and the integration of sensory with prior information.

Interestingly, and in this line of thinking, a recent study found that by optogenetically inactivating the posterior parietal cortex (PPC) of rats performing an auditory parametric working memory task (PWM) [73], the animals' performance experienced a substantial increase. Their analysis showed that, even though the task was not designed to consider previously presented stimuli, rats still combined past information together with current sensory information to perform it. When shutting down the PPC, rats' performance was enhanced because they started performing the task just relying on the acoustic sensory information. Finally, I think that, due to the fact that the reported results were much in agreement with

the hypothesis stated by Wilson and colleagues [97], it would be very interesting to further analyze the rats' behavior under a RL scheme. Once characterized, we could seek for distinctive signatures of RL in the rats' brain activity.

Besides proposing possible extensions to the research reported in this doctoral thesis I also found convenient to briefly state here some of the main lines of research I want to explore in the near future. Identifying efficient ensembles of neurons for the encoding of information has been a topic of discussion for quite some time in theoretical neuroscience [165, 166]. The analytic formula for the amount of encoded information derived in chapter 2 could be used to find an optimal algorithm for building efficient ensembles of neurons under some network constraints like the tuning curves or the distribution of correlations. This question is generally tackled by methods that tend to be computationally expensive [165]. With the analytic expression we could provide a grounded solution for this problem and understand for instance whether non-tuned neurons help the encoding of information or not.

The nature of differential correlations [47] and its relationship with information is, in my opinion, a very elegant solution to the well-known problem of behavioral performances not being outstandingly larger than the amount of information found in single-neurons. Nevertheless, I think that it would be very interesting to study a natural extension of this idea by analyzing how differential correlations propagate along the stages of a feedforward network. As differential correlations arise simply by the noise that can be found in the input layer, this rationale could be further extended to an arbitrarily deep network, where the input to each layer corresponds to the (noisy) output of the preceding one. It is important to remark that a similar idea was the proposed in [167]. The main question I would like to answer is: is the noise going to be amplified, decrease or stay constant across layers? As deep feedforward networks are nowadays playing an important role in the field of AI [168, 169], characterizing how noise propagates in such a system could potentially have many useful applications.

Finally, I would like to present an idea I have been thinking about re-

cently which has been mostly inspired by the recent successes obtained by deep convolutional networks (deep feedforward networks) on pattern classification and in particular in image recognition. Previous to this improvement, most of the human efforts in this field were directed towards identifying the most efficient features of a set of images for classification. By including large datasets, very flexible models and efficient optimization algorithms, researchers realized that these algorithms were outstandingly better than humans on this task. When extracting the highly convoluted statistical structure of real-world images, optimization algorithms were way more efficient than humans. I think this same idea could be applied for performing science. As the causal graph defining the set of relationships in a system can be very complicated, I think algorithms, provided with enough data, would be able to extract these regularity patterns in a much more efficient way than humans. This is not an entirely novel idea, and indeed many important steps have been done already in this direction [170–172]. In the future I would like to explore this idea in detail and start applying it in systems biology and systems neuroscience.

## 4.4 General perspective

As stated in the introduction of this thesis, living beings, and in particular humans, can be thought as devices that exploit the statistical regularities that can be found in the particular environment they are embedded in so that they can maximize their chances for survival. So to say, we are inferential machines. Inference can be done at very different levels. For example, perception per se is an inferential process where information gathered through the different sensory channels is combined with prior information in order to make sense of the current state of the environment. Perception is performed by most of the living beings, from amoebas to humans. As the complexity of our universe is potentially infinite across the temporal and the spatial domain, it exists a very large set of scales at which living beings can try to find statistical regularities. In other words, there are different levels of complexity at which living beings could po-

tentially try to infer the hidden underlying variables or regularities. The more complex is the system used by the living beings to infer these hidden variables or rules, the more abstract will be the domain of the inference performed by them. Under this framework, humans can be thought of as the living beings that make inference in the more abstract domain, that is, making inference on the general hidden rules and relationships of how nature works as a whole.

Humans have tried to infer the general rules of the world for thousands of years. Paleolithic societies already tried to make some sense of the origins and future of their environment, their society and themselves. The oldest evidence for burial rituals can be dated back to 400,000 years ago. Later, when agriculture played a fundamental role on the development of a society, finding regularities on the seasonal cycles and on the movement and location of the celestial bodies could be the key difference between eating or starving to death. Since then, human societies have proposed thousands and thousands of hypothesis regarding the underlying truth ruling all domains of nature. Being able to make accurate predictions on all of these domains has acted many times as the most important factor on determining whether a society succeeded or failed. Making such inferences is not trivial at all and this is the reason why the method itself has also evolved across time.

Many different methodologies have emerged throughout history but here I will focus on two of them: rationalism and empiricism. Roughly speaking, rationalism is a school of thought that states that truth will be mainly achieved by reason. Just by the intellect and deduction we will be able to reach the hidden structure of the world. Rene Descartes and Immanuel Kant are very important figures of rationalism. Mathematics and metaphysics, among others, correspond to knowledge that can be acquired only through reason. Empiricism, on the contrary, is the epistemological school of thought that states that real knowledge can only be acquired through sensory experience. John Locke and David Hume are two important figures of this philosophical movement. During many years Rationalism was the dominant school of thought in western society and many important achievements were accomplished by it. The development

of modern mathematics would be the most important one in my opinion. It was also responsible for many important mistakes, being the most notorious the importance of religions in all domains of life and knowledge for more than 2,000 years.

Modern science emerged from a combination between rationalism and empiricism. The corner stone of modern science is the scientific method. Hypothesis about the underlying structure or hidden rules of a natural phenomena are generated through the intellect and then they are tested experimentally by comparing the predictions derived from the hypothesis with the outcomes of the experiment. This is a really powerful method that has been proven to be extremely successful on extracting general principles and rules from nature. Science can be therefore defined as the process of acquiring knowledge through the scientific method. Thanks to science we have been able to understand the rules that govern the motion and interaction of physical bodies, the behavior and evolution of the electromagnetic fields, the stability and potential of the fundamental constituents of matter, the relationship between time, space, mass and velocity, and the evolution of all forms of life, among others.

The amount of knowledge we have acquired thanks to the scientific method is in my opinion one of the most important achievement humankind has ever made. Nevertheless, it is astonishing how, after all these amazing achievements, there are still communities that resist to accepted how powerful is this method for acquiring knowledge. It is always important to remember that the difference between scientific knowledge and non-scientific "knowledge" is that the former assumes in a humble way that the human inferential process for general rules of nature can not be entirely performed by our intellect. It understands our limited computational capabilities and "asks" nature whether the proposed hidden rule is true or not. In contrast, the non-scientific "knowledge" is not interested in asking nature whether their hypothesis are correct or not, that is, they are not interested in experimentally testing their ideas about the world. It does not mean that non-scientific hypothesis are automatically wrong, but they are likely to be wrong.

I think that we are nowadays experiencing an epistemological change

of paradigm. In the same way as the emergence of modern science constituted an epistemological change of framework by combining rationalism (generation of hypothesis) with empiricism (testing hypothesis empirically), we are historically located in the dawn of a new epoch, where we no longer can generate hypothesis that account for the complexity of the natural phenomena we are trying to understand. I would call this new epoch as *empirical artificialism*. Understanding how the physical world works is relatively "easy". It took us around 300 hundred years to have a relatively good and profound understanding of how the physical world at our middle scale works. Mechanics, electromagnetism, general and special relativity and the quantum theory are robust theories that account with relatively high precision for most of the physical phenomena that expand from the Angström ( $10^{-10}$  m) to the Mega-parsec ( $10^{22}$  m). In my opinion, for systems biology we are no longer able to generate hypothesis that account for the all the complexity involved on it.

Our computational capabilities are limited and so it is our capability of understanding. In the same way as the combination of very flexible models with very large datasets and efficient optimization rules has lead to a substantial improvement on artificial image recognition because of the task complexity, systems biology (and science in general) would experience a substantial improvement if we apply the same rationale for the generation of hypothesis. These artificially generated hypothesis should be tested on experiments (*empirical artificialism*), as the ultimate razor for determining what is the truth will always be nature itself.

I would like to call attention on what does it mean to understand something. For that I would like to start with the most common reply I am told when I expose this set of ideas: "science is about understanding nature through mathematical principles that explain a particular process. By first deriving the mathematics and logical connections, and then by testing them on experiments, we can understand such processes". I totally agree with this statement, however the causal relationships involved in such complex systems can not be represented in many cases by the set of computational resources humans have for abstraction and inference. I think that we have to re-define the concept of understanding something

under this new framework. Understanding a natural process or general law is being able to predict the future state of the system. When we say that the parabolic shot is understood, we are implicitly saying that, given the initial position and velocity of the object, we could predict with high accuracy what will be the future positions and velocities of the object for a particular time window. The further in time we can characterize the state of the system, the better will be our model or theory.

This rationale could be applied to all the fields of physics mentioned above. Quantum physics needs further clarification. The definition of states in quantum physics is what makes it different from most of the rest of disciplines. In there, it makes no longer sense to define the state of a system by its position, velocity, energy, time, momentum, etc, because surprisingly, at the atomic level, these different properties are sometimes mutually exclusive. In quantum physics, a state is fully determined by its wave function (quantum state). Under this framework, we can say that we understand quantum mechanics because, given the initial properties of the system, we can potentially characterize the time evolution of its wave function. This is, in my opinion, what has to be the final aim of science, predicting the future state of a system given the initial conditions. If the hidden rule or relationship is so complex that humans can not understand it, it is not because we are not doing science, it is just because our computational capabilities are limited when compared to the complexity of the world we are embedded in. At this point I would like to mention that, although I thought about this idea some time before reading it, an equivalent idea was published by Ted Chiang back in the year 2000 [173]. Ted Chiang calls this type of science *metascience*, however, I would rather call it *computational epistemology*.

As a final comment, I would like to state that, even though the scientific method has been proven to be the most efficient epistemological methodology, and also I think that it needs to evolve to the *empirical artificialism* for further understanding nature, I am also concerned that eventually these epistemological methodologies will become outdated. The theory of evolution is able to explain not only the evolution of living beings but many more phenomena, namely the evolution of general ideas

(memes) in a society [6]. I think that the methodology we apply for understanding nature can be also understood as a meme and therefore it is also ruled by the theory of evolution. Many different epistemological strategies have been generated throughout the course of history. When a society developed an epistemological methodology that was able to predict accurately the future state of world (understand its hidden rules), these societies tended to persist throughout time alongside with this methodology. Science has been the most successful epistemological meme for the last three centuries, but small variants will start arising soon and the best adapted ones (best predictive power) will take over. My prediction is that *empirical artificialism* is going to be the upcoming best-adapted epistemological meme. Who knows what will be the next step.

# Bibliography

1. Darwin, C. *On the origin of species* (John murray, london, 1859).
2. Mach, E. Contributions to the analysis of the sensations (CM Williams, Trans.) *La Salle, IL: Open Court.*(Original work published 1890) (1980).
3. Von Helmholtz, H. *Treatise on Physiological Optics: Translated from the 3rd German Ed* (Optical Society of America, 1925).
4. Kersten, D., Mamassian, P. & Yuille, A. Object perception as Bayesian inference. *Annu. Rev. Psychol.* **55**, 271–304 (2004).
5. Tenenbaum, J. B. & Griffiths, T. L. Generalization, similarity, and Bayesian inference. *Behavioral and brain sciences* **24**, 629–640 (2001).
6. Dawkins, R. *The selfish gene* (Oxford university press, 2016).
7. Kandel, E. R., Schwartz, J. H., Jessell, T. M., Siegelbaum, S. A., Hudspeth, A. J., *et al.* *Principles of neural science* (McGraw-hill New York, 2000).
8. Knill, D. C. & Richards, W. *Perception as Bayesian inference* (Cambridge University Press, 1996).
9. Rao, R. P. An optimal estimation approach to visual perception and learning. *Vision research* **39**, 1963–1989 (1999).
10. Knill, D. C. & Pouget, A. The Bayesian brain: the role of uncertainty in neural coding and computation. *TRENDS in Neurosciences* **27**, 712–719 (2004).

11. Wolpert, D. M., Ghahramani, Z. & Jordan, M. I. An internal model for sensorimotor integration. *Science*, 1880–1882 (1995).
12. Ernst, M. O. & Banks, M. S. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature* **415**, 429–433 (2002).
13. Mamassian, P. & Landy, M. S. Interaction of visual prior constraints. *Vision research* **41**, 2653–2668 (2001).
14. Van Beers, R. J., Sittig, A. C. & van Der Gon, J. J. D. Integration of proprioceptive and visual position-information: An experimentally supported model. *Journal of neurophysiology* **81**, 1355–1364 (1999).
15. Nogueira, R. *et al.* Lateral orbitofrontal cortex anticipates choices and integrates prior with current information. *Nature Communications* **8**, 14823 (2017).
16. Mamassian, P. & Goutcher, R. Prior knowledge on the illumination position. *Cognition* **81**, B1–B9 (2001).
17. Berkes, P., Orbán, G., Lengyel, M. & Fiser, J. Spontaneous cortical activity reveals hallmarks of an optimal internal model of the environment. *Science* **331**, 83–87 (2011).
18. Körding, K. P. & Wolpert, D. M. Bayesian decision theory in sensorimotor control. *Trends in cognitive sciences* **10**, 319–326 (2006).
19. Bishop, C. M. *Pattern recognition and machine learning* (Springer, 2006).
20. Rosenblatt, F. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological review* **65**, 386 (1958).
21. Fisher, R. A. The use of multiple measurements in taxonomic problems. *Annals of human genetics* **7**, 179–188 (1936).
22. Duda, R. O., Hart, P. E. & Stork, D. G. *Pattern classification* (Wiley, New York, 1973).

23. Acerbi, L., Vijayakumar, S. & Wolpert, D. M. On the origins of suboptimality in human probabilistic inference. *PLoS computational biology* **10**, e1003661 (2014).
24. Hubel, D. H. & Wiesel, T. N. Receptive fields of single neurones in the cat's striate cortex. *The Journal of physiology* **148**, 574–591 (1959).
25. Brunton, B. W., Botvinick, M. M. & Brody, C. D. Rats and humans can optimally accumulate evidence for decision-making. *Science* **340**, 95–98 (2013).
26. Abolafia, J. M., Martinez-Garcia, M., Deco, G. & Sanchez-Vives, M. V. Variability and information content in auditory cortex spike trains during an interval-discrimination task. *Journal of neurophysiology* **110**, 2163–2174 (2013).
27. Romo, R & Salinas, E. Flutter discrimination: neural codes, perception, memory and decision making. *Nat. Rev. Neurosci.* **4**, 203–218 (2003).
28. Romo, R., Hernández, A., Zainos, A., Lemus, L. & Brody, C. D. Neuronal correlates of decision-making in secondary somatosensory cortex. *Nature neuroscience* **5**, 1217 (2002).
29. Chen, A., DeAngelis, G. C. & Angelaki, D. E. Macaque parieto-insular vestibular cortex: responses to self-motion and optic flow. *Journal of Neuroscience* **30**, 3022–3042 (2010).
30. Chen, A., DeAngelis, G. C. & Angelaki, D. E. Convergence of vestibular and visual self-motion signals in an area of the posterior sylvian fissure. *Journal of Neuroscience* **31**, 11617–11627 (2011).
31. Uchida, N. & Mainen, Z. F. Speed and accuracy of olfactory discrimination in the rat. *Nature neuroscience* **6**, 1224 (2003).
32. Wilson, R. I. & Mainen, Z. F. Early events in olfactory processing. *Annu. Rev. Neurosci.* **29**, 163–201 (2006).

33. Tolhurst, D. J., Movshon, J. A. & Dean, A. F. The statistical reliability of signals in single neurons in cat and monkey visual cortex. *Vision research* **23**, 775–785 (1983).
34. Shadlen, M. N. & Newsome, W. T. The variable discharge of cortical neurons: implications for connectivity, computation, and information coding. *Journal of neuroscience* **18**, 3870–3896 (1998).
35. Zohary, E., Shadlen, M. N. & Newsome, W. T. Correlated neuronal discharge rate and its implications for psychophysical performance. *Nature* **370**, 140–143 (1994).
36. Faisal, A. A., Selen, L. P. & Wolpert, D. M. Noise in the nervous system. *Nature reviews. Neuroscience* **9**, 292 (2008).
37. Rolls, E. & Deco, G. *The noisy brain. Stochastic dynamics as a principle of brain function* (Oxford Univ. Press, UK, 2010).
38. Van Vreeswijk, C., Sompolinsky, H., *et al.* Chaos in neuronal networks with balanced excitatory and inhibitory activity. *Science* **274**, 1724–1726 (1996).
39. Brunel, N. Dynamics of sparsely connected networks of excitatory and inhibitory spiking neurons. *Journal of computational neuroscience* **8**, 183–208 (2000).
40. Renart, A. *et al.* The asynchronous state in cortical circuits. *science* **327**, 587–590 (2010).
41. Carandini, M. Amplification of trial-to-trial response variability by neurons in visual cortex. *PLoS biology* **2**, e264 (2004).
42. Lin, I.-C., Okun, M., Carandini, M. & Harris, K. D. The nature of shared cortical variability. *Neuron* **87**, 644–656 (2015).
43. Ecker, A. S. *et al.* State dependence of noise correlations in macaque primary visual cortex. *Neuron* **82**, 235–248 (2014).
44. Goris, R. L., Movshon, J. A. & Simoncelli, E. P. Partitioning neuronal variability. *Nature neuroscience* **17**, 858–865 (2014).
45. Cohen, M. R. & Kohn, A. Measuring and interpreting neuronal correlations. *Nature neuroscience* **14**, 811–819 (2011).

46. Kohn, A., Coen-Cagli, R., Kanitscheider, I. & Pouget, A. Correlations and neuronal population information. *Annual review of neuroscience* **39**, 237–256 (2016).
47. Moreno-Bote, R. *et al.* Information-limiting correlations. *Nature neuroscience* **17**, 1410–1417 (2014).
48. Abbott, L. & Dayan, P. The effect of correlated variability on the accuracy of a population code. *Neural computation* **11**, 91–101 (1999).
49. Britten, K. H., Shadlen, M. N., Newsome, W. T. & Movshon, J. A. The analysis of visual motion: a comparison of neuronal and psychophysical performance. *Journal of Neuroscience* **12**, 4745–4765 (1992).
50. Bair, W., Zohary, E. & Newsome, W. T. Correlated firing in macaque visual area MT: time scales and relationship to behavior. *Journal of Neuroscience* **21**, 1676–1697 (2001).
51. Arandia-Romero, I., Tanabe, S., Drugowitsch, J., Kohn, A. & Moreno-Bote, R. Multiplicative and additive modulation of neuronal tuning with population activity affects encoded information. *Neuron* **89**, 1305–1316 (2016).
52. Beck, J., Bejjanki, V. R. & Pouget, A. Insights from a simple expression for linear fisher information in a recurrently connected population of spiking neurons. *Neural computation* **23**, 1484–1502 (2011).
53. Salinas, E. & Abbott, L. Vector reconstruction from firing rates. *Journal of computational neuroscience* **1**, 89–107 (1994).
54. Seung, H. S. & Sompolinsky, H. Simple models for reading neuronal population codes. *Proceedings of the National Academy of Sciences* **90**, 10749–10753 (1993).
55. Kohn, A. & Smith, M. A. Stimulus dependence of neuronal correlation in primary visual cortex of the macaque. *Journal of Neuroscience* **25**, 3661–3673 (2005).

56. Smith, M. A. & Kohn, A. Spatial and temporal scales of neuronal correlation in primary visual cortex. *Journal of Neuroscience* **28**, 12591–12603 (2008).
57. Ecker, A. S. *et al.* Decorrelated neuronal firing in cortical microcircuits. *science* **327**, 584–587 (2010).
58. Cohen, M. R. & Maunsell, J. H. Attention improves performance primarily by reducing interneuronal correlations. *Nature neuroscience* **12**, 1594–1600 (2009).
59. Ruff, D. A. & Cohen, M. R. Attention can either increase or decrease spike count correlations in visual cortex. *Nature neuroscience* **17**, 1591–1597 (2014).
60. Mitchell, J. F., Sundberg, K. A. & Reynolds, J. H. Spatial attention decorrelates intrinsic activity fluctuations in macaque area V4. *Neuron* **63**, 879–888 (2009).
61. Averbeck, B. B., Latham, P. E. & Pouget, A. Neural correlations, population coding and computation. *Nature reviews. Neuroscience* **7**, 358 (2006).
62. McAdams, C. J. & Maunsell, J. H. Effects of attention on orientation-tuning functions of single neurons in macaque cortical area V4. *Journal of Neuroscience* **19**, 431–441 (1999).
63. Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic population codes. *Nature neuroscience* **9**, 1432–1438 (2006).
64. Baillargeon, R. Infants' physical world. *Current directions in psychological science* **13**, 89–94 (2004).
65. Baillargeon, R. The acquisition of physical knowledge in infancy: A summary in eight lessons. *Blackwell handbook of childhood cognitive development* **1**, 1 (2002).
66. Lake, B. M., Ullman, T. D., Tenenbaum, J. B. & Gershman, S. J. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 1–101 (2016).

67. Roitman, J. D. & Shadlen, M. N. Response of neurons in the lateral intraparietal area during a combined visual discrimination reaction time task. *J. Neurosci.* **22**, 9475–9489 (2002).
68. Katz, L. N., Yates, J. L., Pillow, J. W. & Huk, A. C. Dissociated functional significance of decision-related activity in the primate dorsal stream. *Nature* **535**, 285–288 (2016).
69. Rudolph, K. & Pasternak, T. Transient and permanent deficits in motion perception after lesions of cortical areas MT and MST in the macaque monkey. *Cerebral Cortex* **9**, 90–100 (1999).
70. Miura, K., Mainen, Z. F. & Uchida, N. Odor representations in olfactory cortex: distributed rate coding and decorrelated population activity. *Neuron* **74**, 1087–1098 (2012).
71. Rolls, E. T., Scott, T. R., Sienkiewicz, Z. J. & Yaxley, S. The responsiveness of neurones in the frontal opercular gustatory cortex of the macaque monkey is independent of hunger. *The Journal of physiology* **397**, 1–12 (1988).
72. Hanks, T. D., Mazurek, M. E., Kiani, R., Hopp, E. & Shadlen, M. N. Elapsed decision time affects the weighting of prior probability in a perceptual decision task. *The Journal of neuroscience* **31**, 6339–6352 (2011).
73. Akrami, A., Kopec, C. D., Diamond, M. E. & Brody, C. D. Posterior parietal cortex represents sensory stimulus history and is necessary for its effects on behavior. *bioRxiv*, 182246 (2017).
74. Barraclough, D. J., Conroy, M. L. & Lee, D. Prefrontal cortex and decision making in a mixed-strategy game. *Nat. Neurosci.* **7**, 404–410 (2004).
75. Seo, M., Lee, E. & Averbeck, B. B. Action selection and action value in frontal-striatal circuits. *Neuron* **74**, 947–960 (2012).
76. Leon, M. I. & Shadlen, M. N. Representation of time by neurons in the posterior parietal cortex of the macaque. *Neuron* **38**, 317–327 (2003).

77. Hanes, D. P. & Schall, J. D. Neural control of voluntary movement initiation. *Science* **274**, 427–430 (1996).
78. Rushworth, M. F., Noonan, M. P., Boorman, E. D., Walton, M. E. & Behrens, T. E. Frontal cortex and reward-guided learning and decision-making. *Neuron* **70**, 1054–1069 (2011).
79. Stalnaker, T. A. Orbitofrontal neurons infer the value and identity of predicted outcomes. *Nat. Commun.* **5**, 3926 (2015).
80. Padoa-Schioppa, C & Assad, J. A. Neurons in the orbitofrontal cortex encode economic value. *Nature* **441**, 223–226 (2006).
81. Kennerley, S. W., Dahmubed, A. F., Lara, A. H. & Wallis, J. D. Neurons in the frontal lobe encode the value of multiple decision variables. *J. Cogn. Neurosci.* **21**, 1162–1178 (2009).
82. Furuyashiki, T & Gallagher, M. Neural encoding in the orbitofrontal cortex related to goal-directed behavior. *Ann. N. Y. Acad. Sci.* **1121**, 193–215 (2007).
83. Fuster, J. M. *The Prefrontal Cortex* (Raven Press, New York, 1997).
84. Kringelbach, M. L. & Rolls, E. T. The functional neuroanatomy of the human orbitofrontal cortex: evidence from neuroimaging and neuropsychology. *Progress in neurobiology* **72**, 341–372 (2004).
85. Wallis, J. D. Orbitofrontal cortex and its contribution to decision-making. *Annu. Rev. Neurosci.* **30**, 31–56 (2007).
86. Schoenbaum, G, Chiba, A. A. & Gallagher, M. Orbitofrontal cortex and basolateral amygdala encode expected outcomes during learning. *Nat. Neurosci.* **1**, 155–159 (1998).
87. Tremblay, L & Schultz, W. Relative reward preference in primate orbitofrontal cortex. *Nature* **398**, 704–708 (1999).
88. O’Doherty, J. P. Reward representations and reward-related learning in the human brain: insights from neuroimaging. *Curr. Opin. Neurobiol.* **14**, 769–776 (2004).

89. O'Doherty, J, Kringelbach, M. L., Rolls, E. T., Hornak, J & Andrews, C. Abstract reward and punishment representations in the human orbitofrontal cortex. *Nat. Neurosci.* **4**, 95–102 (2001).
90. Hare, T. A., O'Doherty, J, Camerer, C. F., Schultz, W & Rangel, A. Dissociating the role of the orbitofrontal cortex and the striatum in the computation of goal values and prediction errors. *J. Neurosci.* **28**, 5623–5630 (2008).
91. Feierstein, C. E., Quirk, M. C., Uchida, N, Sosulski, D. L. & Mainen, Z. F. Representation of spatial goals in rat orbitofrontal cortex. *Neuron* **51**, 495–507 (2006).
92. Roesch, M. R., Taylor, A. R. & Schoenbaum, G. Encoding of time-discounted rewards in orbitofrontal cortex is independent of value representation. *Neuron* **51**, 509–520 (2006).
93. Furuyashiki, T, Holland, P. C. & Gallagher, M. Rat orbitofrontal cortex separately encodes response and outcome information during performance of goal-directed behavior. *J. Neurosci.* **28**, 5127–5138 (2008).
94. Padoa-Schioppa, C. Neuronal origins of choice variability in economic decisions. *Neuron* **80**, 1322–1336 (2013).
95. Shadlen, M. N. & Newsome, W. T. Neural basis of a perceptual decision in the parietal cortex (area LIP) of the rhesus monkey. *Journal of neurophysiology* **86**, 1916–1936 (2001).
96. Williams, Z. M., Elfar, J. C., Eskandar, E. N., Toth, L. J. & Assad, J. A. Parietal activity and the perceived direction of ambiguous apparent motion. *Nature neuroscience* **6**, 616 (2003).
97. Wilson, R. C., Takahashi, Y. K., Schoenbaum, G & Niv, Y. Orbitofrontal cortex as a cognitive map of task space. *Neuron* **81**, 267–279 (2014).
98. Sutton, R. S. & Barto, A. G. *Reinforcement learning: An introduction* **1** (MIT press Cambridge, 1998).

99. Panzeri, S., Harvey, C. D., Piasini, E., Latham, P. E. & Fellin, T. Cracking the neural code for sensory perception by combining statistics, intervention, and behavior. *Neuron* **93**, 491–507 (2017).
100. Arandia-Romero, I., Nogueira, R., Mochol, G. & Moreno-Bote, R. What can neuronal populations tell us about cognition? *Current opinion in neurobiology* **46**, 48–57 (2017).
101. Mountcastle, V. B., Talbot, W. H., Darian-Smith, I. & Kornhuber, H. H. Neural basis of the sense of flutter-vibration. *Science* **155**, 597–600 (1967).
102. Smith, M. A., Kohn, A. & Movshon, J. A. Glass pattern responses in macaque V2 neurons. *Journal of Vision* **7**, 5–5 (2007).
103. Hung, C. P., Kreiman, G., Poggio, T. & DiCarlo, J. J. Fast readout of object identity from macaque inferior temporal cortex. *Science* **310**, 863–866 (2005).
104. DeAngelis, G. C., Ohzawa, I. & Freeman, R. D. Depth is encoded in the visual cortex by a specialized receptive field structure. *Nature* **352**, 156 (1991).
105. Tremblay, S., Pieper, F., Sachs, A. & Martinez-Trujillo, J. Attentional filtering of visual information by neuronal ensembles in the primate lateral prefrontal cortex. *Neuron* **85**, 202–215 (2015).
106. Romo, R., Hernández, A., Zainos, A. & Salinas, E. Correlated neuronal discharges that increase coding efficiency during perceptual discrimination. *Neuron* **38**, 649–657 (2003).
107. Gu, Y. *et al.* Perceptual learning reduces interneuronal correlations in macaque visual cortex. *Neuron* **71**, 750–761 (2011).
108. Treue, S. & Trujillo, J. C. M. Feature-based attention influences motion processing gain in macaque visual cortex. *Nature* **399**, 575 (1999).
109. Sompolinsky, H., Yoon, H., Kang, K. & Shamir, M. Population coding in neuronal systems with correlated noise. *Physical Review E* **64**, 051904 (2001).

110. Chen, Y., Geisler, W. S. & Seidemann, E. Optimal decoding of correlated neural population responses in the primate visual cortex. *Nature neuroscience* **9**, 1412 (2006).
111. Pitkow, X., Liu, S., Angelaki, D. E., DeAngelis, G. C. & Pouget, A. How can single sensory neurons predict behavior? *Neuron* **87**, 411–423 (2015).
112. Graf, A. B., Kohn, A., Jazayeri, M. & Movshon, J. A. Decoding the activity of neuronal populations in macaque primary visual cortex. *Nature neuroscience* **14**, 239–245 (2011).
113. Kanitscheider, I., Coen-Cagli, R. & Pouget, A. Origin of information-limiting noise correlations. *Proceedings of the National Academy of Sciences* **112**, E6973–E6982 (2015).
114. Gold, J. I. & Shadlen, M. N. The neural basis of decision making. *Annu. Rev. Neurosci.* **30**, 535–574 (2007).
115. Averbeck, B. B., Sohn, J. W. & Lee, D. Activity in prefrontal cortex during dynamic selection of action sequences. *Nat. Neurosci.* **9**, 276–282 (2006).
116. Glascher, J., Daw, N., Dayan, P & O’Doherty, J. P. States versus rewards: dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron* **66**, 585–595 (2010).
117. Doll, B. B., Simon, D. A. & Daw, N. D. The ubiquity of model-based reinforcement learning. *Curr. Opin. Neurobiol.* **22**, 1075–1081 (2012).
118. Lee, D, Seo, H & Jung, M. W. Neural basis of reinforcement learning and decision making. *Annu. Rev. Neurosci.* **35**, 287–308 (2012).
119. Rudebeck, P. H. & Murray, E. A. The orbitofrontal oracle: cortical mechanisms for the prediction and evaluation of specific behavioral outcomes. *Neuron* **84**, 1143–1156 (2014).

120. Rolls, E. T., Critchley, H. D., Mason, R & Wakeman, E. A. Orbitofrontal cortex neurons: role in olfactory and visual association learning. *J. Neurophysiol.* **75**, 1970–1981 (1996).
121. Sul, J. H., Kim, H, Huh, N, Lee, D & Jung, M. W. Distinct roles of rodent orbitofrontal and medial prefrontal cortex in decision making. *Neuron* **66**, 449–460 (2010).
122. Watson, K. K. & Platt, M. L. Social signals in primate orbitofrontal cortex. *Curr. Biol.* **22**, 2268–2273 (2012).
123. Rangel, A & Hare, T. Neural computations associated with goal-directed choice. *Curr. Opin. Neurobiol.* **20**, 262–270 (2010).
124. Neural correlates, computation and behavioural impact of decision confidence. *Nature* **455**, 227–31 (2008).
125. Kiani, R. & Shadlen, M. N. Representation of confidence associated with a decision by neurons in the parietal cortex. *science* **324**, 759–764 (2009).
126. Lapish, C. C., Balaguer-Ballester, E, Seamans, J. K., Phillips, A. G. & Durstewitz, D. Amphetamine exerts dose-dependent changes in prefrontal cortex attractor dynamics during working memory. *J. Neurosci.* **35**, 10172–10187 (2015).
127. Balaguer-Ballester, E, Lapish, C. C., Seamans, J. K. & Durstewitz, D. Attracting dynamics of frontal cortex ensembles during memory-guided decision-making. *P Lo S Comput. Biol.* **7**, e1002057 (2011).
128. Izquierdo, A, Suda, R. K. & Murray, E. A. Bilateral orbital prefrontal cortex lesions in rhesus monkeys disrupt choices guided by both reward value and reward contingency. *J. Neurosci.* **24**, 7540–7548 (2004).
129. Jang, A. I. The role of frontal cortical and medial-temporal lobe brain areas in learning a Bayesian prior belief on reversals. *J. Neurosci.* **35**, 11751–11760 (2015).

130. Schoenbaum, G, Setlow, B, Saddoris, M. P. & Gallagher, M. Encoding predicted outcome and acquired value in orbitofrontal cortex during cue sampling depends upon input from basolateral amygdala. *Neuron* **39**, 855–867 (2003).
131. Ostlund, S. B. & Balleine, B. W. Orbitofrontal cortex mediates outcome encoding in Pavlovian but not instrumental conditioning. *J. Neurosci.* **27**, 4819–4825 (2007).
132. Schuck, N. W., Cai, M. B., Wilson, R. C. & Niv, Y. Human orbitofrontal cortex represents a cognitive map of state space. *Neuron* **91**, 1402–1412 (2016).
133. Schoenbaum, G, Nugent, S. L., Saddoris, M. P. & Setlow, B. Orbitofrontal lesions in rats impair reversal but not acquisition of go, no-go odor discriminations. *Neuroreport* **13**, 885–890 (2002).
134. Noonan, M. P. Separate value comparison and learning mechanisms in macaque medial and lateral orbitofrontal cortex. *Proc. Natl Acad. Sci. Usa* **107**, 20547–20552 (2010).
135. Walton, M. E., Behrens, T. E., Buckley, M. J., Rudebeck, P. H. & Rushworth, M. F. Separable learning systems in the macaque brain and the role of orbitofrontal cortex in contingent learning. *Neuron* **65**, 927–939 (2010).
136. Bradfield, L. A., Dezfouli, A, van Holstein, M, Chieng, B & Balleine, B. W. Medial orbitofrontal cortex mediates outcome retrieval in partially observable task situations. *Neuron* **88**, 1268–1280 (2015).
137. Genovesio, A, Tsujimoto, S, Navarra, G, Falcone, R & Wise, S. P. Autonomous encoding of irrelevant goals and outcomes by prefrontal cortex neurons. *J. Neurosci.* **34**, 1970–1978 (2014).
138. Drugowitsch, J., Moreno-Bote, R., Churchland, A. K., Shadlen, M. N. & Pouget, A. The Cost of Accumulating Evidence in Perceptual Decision Making. *The Journal of Neuroscience* **32**, 3612–3628 (2012).

139. Lange, F. P., Rahnev, D. A., Donner, T. H. & Lau, H. Prestimulus oscillatory activity over motor cortex reflects perceptual expectations. *J. Neurosci.* **33**, 1400–1410 (2013).
140. Donahue, C. H. & Lee, D. Dynamic routing of task-relevant signals for decision making in dorsolateral prefrontal cortex. *Nature neuroscience* **18**, 295–301 (2015).
141. Rich, E. L. & Wallis, J. D. Decoding subjective decisions from orbitofrontal cortex. *Nat. Neurosci.* **19**, 973–980 (2016).
142. Gourley, S. L. The orbitofrontal cortex regulates outcome-based decision-making via the lateral striatum. *Eur. J. Neurosci.* **38**, 2382–2388 (2013).
143. Gremel, C. M. & Costa, R. M. Orbitofrontal and striatal circuits dynamically encode the shift between goal-directed and habitual actions. *Nat. Commun.* **4**, 2264 (2013).
144. Gremel, C. M. Endocannabinoid modulation of orbitostriatal circuits gates habit formation. *Neuron* **90**, 1312–1324 (2016).
145. Rhodes, S. E. & Murray, E. A. Differential effects of amygdala, orbital prefrontal cortex, and prelimbic cortex lesions on goal-directed behavior in rhesus macaques. *J. Neurosci.* **33**, 3380–3389 (2013).
146. Sleezer, B. J., Castagno, M. D. & Hayden, B. Y. Rule encoding in orbitofrontal cortex and striatum guides selection. *J. Neurosci.* **36**, 11223–11237 (2016).
147. Chan, S. C., Niv, Y. & Norman, K. A. A probability distribution over latent causes, in the orbitofrontal cortex. *J. Neurosci.* **36**, 7817–7828 (2016).
148. Galvan, A. Earlier development of the accumbens relative to orbitofrontal cortex might underlie risk-taking behavior in adolescents. *J. Neurosci.* **26**, 6885–6892 (2006).
149. Bolla, K. I. Orbitofrontal cortex dysfunction in abstinent cocaine abusers performing a decision-making task. *Neuroimage* **19**, 1085–1094 (2003).

150. Friedman, J., Hastie, T. & Tibshirani, R. *The elements of statistical learning* (Springer series in statistics New York, 2001).
151. He, H & García, E. A. Learning from imbalanced data. *Ieee Trans. Knowl. Data Eng.* **21**, 1263–1284 (2009).
152. Wichmann, F. A. & Hill, N. J. The psychometric function: I. Fitting, sampling, and goodness of fit. *Percept Psychophys.* **63**, 1293–1313 (2001).
153. Paxinos, G. & Watson, C. *The rat brain in stereotaxic coordinates* (Academic Press, San Diego CA, 1998).
154. Britten, K. H., Newsome, W. T., Shadlen, M. N., Celebrini, S. & Movshon, J. A. A relationship between behavioral choice and the visual responses of neurons in macaque MT. *Visual neuroscience* **13**, 87–100 (1996).
155. Schoenbaum, G & Eichenbaum, H. Information coding in the rodent prefrontal cortex. I. Single-neuron activity in orbitofrontal cortex compared with that in pyriform cortex. *J. Neurophysiol.* **74**, 733–750 (1995).
156. Kanitscheider, I., Coen-Cagli, R., Kohn, A. & Pouget, A. Measuring Fisher information accurately in correlated neural populations. *PLoS computational biology* **11**, e1004218 (2015).
157. Kim, J. & Ragozzino, M. E. The involvement of the orbitofrontal cortex in learning under changing task contingencies. *Neurobiology of learning and memory* **83**, 125–133 (2005).
158. Butter, C. M. Perseveration in extinction and in discrimination reversal tasks following selective frontal ablations in *Macaca mulatta*. *Physiology & Behavior* **4**, 163–171 (1969).
159. Mishkin, M., Vest, B., Waxler, M. & Rosvold, H. E. A re-examination of the effects of frontal lesions on object alternation. *Neuropsychologia* **7**, 357–363 (1969).
160. Bouton, M. E. Context and behavioral processes in extinction. *Learning & memory* **11**, 485–494 (2004).

161. Colwill, R. M. & Rescorla, R. A. Postconditioning devaluation of a reinforcer affects instrumental responding. *Journal of experimental psychology: animal behavior processes* **11**, 120 (1985).
162. Wise, S. P. Forward frontal fields: phylogeny and fundamental function. *Trends in neurosciences* **31**, 599–608 (2008).
163. Wallis, J. D. Cross-species studies of orbitofrontal cortex and value-based decision-making. *Nature neuroscience* **15**, 13–19 (2012).
164. Romo, R., Hernández, A. & Zainos, A. Neuronal correlates of a perceptual decision in ventral premotor cortex. *Neuron* **41**, 165–173 (2004).
165. Leavitt, M. L., Pieper, F., Sachs, A. J. & Martinez-Trujillo, J. C. Correlated variability modifies working memory fidelity in primate prefrontal neuronal ensembles. *Proceedings of the National Academy of Sciences* **114**, E2494–E2503 (2017).
166. Ince, R. A., Panzeri, S. & Kayser, C. Neural codes formed by small and temporally precise populations in auditory cortex. *Journal of Neuroscience* **33**, 18277–18287 (2013).
167. Zylberberg, J., Pouget, A., Latham, P. E. & Shea-Brown, E. Robust information propagation through noisy neural circuits. *PLoS computational biology* **13**, e1005497 (2017).
168. Krizhevsky, A., Sutskever, I. & Hinton, G. E. *Imagenet classification with deep convolutional neural networks* in *Advances in neural information processing systems* (2012), 1097–1105.
169. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
170. Yamins, D. L. *et al.* Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences* **111**, 8619–8624 (2014).
171. Carleo, G. & Troyer, M. Solving the quantum many-body problem with artificial neural networks. *Science* **355**, 602–606 (2017).

172. Schütt, K. T., Arbabzadah, F., Chmiela, S., Müller, K. R. & Tkatchenko, A. Quantum-chemical insights from deep tensor neural networks. *Nature communications* **8**, 13890 (2017).
173. Chiang, T. Catching crumbs from the table. *Nature* **405**, 517–518 (2000).