

## **Does disrupting the Orbitofrontal Cortex alter sensitivity to punishment? A potential mechanism of compulsivity**

Karly M. Turner<sup>1</sup>, Bernard W. Balleine<sup>1</sup>, & Laura A. Bradfield<sup>2-3</sup>

1. School of Psychology, University of New South Wales, NSW, 2052, Australia
2. Centre for Neuroscience and Regenerative Medicine, University of Technology Sydney, NSW, 2007, Australia
3. St. Vincent's Centre for Applied Medical Research, St. Vincent's Hospital Sydney Limited, Sydney, NSW, 2010, Australia

Corresponding author:

**Dr. Laura Bradfield**

**[laura.bradfield@uts.edu.au](mailto:laura.bradfield@uts.edu.au)**

Centre for Neuroscience and Regenerative Medicine,  
School of Life Sciences,  
University of Technology Sydney (St. Vincent's Campus)  
405 Liverpool St  
Darlinghurst, NSW 2010  
PH: 02 83824950

**Abstract**

Abnormal orbitofrontal cortex (OFC) activity is one of the most common findings from neuroimaging studies of individuals with compulsive disorders such as substance use disorder and obsessive-compulsive disorder. The nature of this abnormality is complex however, with some studies reporting the OFC to be over-active in compulsive individuals relative to controls, whereas other studies report it being under-active, and a further set of studies reporting OFC abnormality in both directions within the same individuals. The OFC has been implicated in a broad range of cognitive processes such as decision-making and goal-directed action. OFC dysfunction could impair these processes leading to the kinds of cognitive/behavioural deficits observed in individuals with compulsive disorders. One such deficit that could arise as a result of OFC dysfunction is an altered sensitivity to punishment, which is one of the core characteristics displayed by individuals across multiple types of compulsive disorders. It is, therefore, the aim of the current review to assess the evidence implicating the OFC in adaptation to punishment and to attempt to identify the critical factors determining this relationship. We distil from this analysis some guidelines for future studies attempting to determine the precise role of the OFC in punishment.

**Keywords**

Medial orbitofrontal cortex, lateral orbitofrontal cortex, positive punishment, negative punishment, compulsion.

## Introduction

Compulsive actions are carried out repeatedly despite not necessarily achieving their intended goal. Indeed, compulsive actions will often persist despite leading to aversive outcomes or the loss of rewards that would typically be punishing enough to cause individuals to cease responding. A person with alcohol use disorder, for example, might continue to seek and consume alcohol despite suffering the breakdown of their relationships, ill health, and severe financial consequences. An individual with obsessive compulsive disorder (OCD) might continue to wash their hands despite developing painful sores or continue to spend great amounts of time checking that a door is locked despite missing out on social activities and work. On the other hand, some compulsive actions could be interpreted as an over-sensitivity to adverse outcomes. A person who obsessively washes their hands to avoid contamination, for instance, could be seen as overly-sensitive to that particular aversive outcome (germs) whilst displaying insensitivity to other more distal forms of punishment (e.g. sores on hands/social isolation; see Figure 1). Thus, the relationship between compulsivity and punishment sensitivity is clearly complex.

The relationship between orbitofrontal cortex (OFC) activity and sensitivity to punishment is also complex. Dysregulation of OFC activity among individuals with compulsive disorders is one of the most consistent findings from neuroimaging studies (Maia, Cooney, & Peterson, 2008; Moorman, 2018), although the nature of this dysregulation is multifaceted. Both hyper- and hypo-activity in OFC have been identified in neuroimaging studies of compulsive individuals (Goldstein & Volkow, 2011; Maia et al., 2008; Moorman, 2018), and even complex mixtures of both patterns have been observed within the same individuals (Goldstein & Volkow, 2011). The only thing that is clear from

these studies is that a relationship between OFC dysfunction and compulsion exists. Although the OFC has been implicated in driving behaviour and cognition across many varied behavioural paradigms over the last several decades, there is currently no clear consensus as to its primary function (or, indeed, whether there is one ‘primary’ OFC function). Nevertheless, OFC dysregulation certainly produces a broad range of deficits in paradigms that require various forms of goal-directed action and decision-making. As reviewed below, one facet of this role for OFC is that its dysfunction often causes alterations in sensitivity to punishment. Thus, it is the aim of the current review to shed light on how exactly OFC dysregulation might lead to punishment insensitivity, with a view to providing insight into how this might contribute to compulsive disorders and a focus on studies conducted within the last 10 years (see Table 1).

**Table 1.** Studies investigating the role of the orbitofrontal cortex in punishment.

Punishment sensitivity	Reference	Species	OFC subregion	Type of Manipulation	Type of Punishment	Task
Increased	O’Doherty et al. 2001	Human	IOFC	Increased BOLD signal	Reward loss	Probabilistic reversal task
	Clarke et al. 2015	NHP	antOFC	Inactivation (effect on subsequent session)	Aversive noise	Punishment task
	Mobini et al. 2002	Rat	vIOFC, some mOFC	Lesion	Reward delay	Delayed discounting
	Mar et al. 2011	Rat	OFC	Lesion	Reward delay	Delayed discounting
			IOFC	Lesion	Reward delay	Delayed discounting
	Rudebeck et al. 2006	Rat	vIOFC	Lesion	Reward delay	T-maze
	Orsini et al. 2015	Rat	IOFC	Lesion	Foot shock	Punishment task
	Ishikawa et al. 2020	Rat	IOFC	Inactivation	Foot shock	Punishment task

Decreased	O'Doherty et al. 2001	Human	mOFC	Decreased BOLD signal	Reward loss	Gambling task
	Mar et al. 2011	Rat	mOFC	Lesion	Reward delay	Delayed discounting
	Bechara et al. 1994	Human	vmOFC	Lesion	Reward loss	Iowa Gambling Task
	Winstanley et al. 2004	Rat	OFC	Lesion	Reward delay	Delayed discounting
	Pais-Vieira et al. 2007	Rat	IOFC	Lesion	Reward omission	Probabilistic discounting
	Stopper et al. 2014	Rat	mOFC	Inactivation	Reward omission	Probabilistic discounting
	Verharen et al. 2020	Rat	mOFC	Inactivation	Reward omission	Probabilistic reversal
			IOFC	Inactivation	Reward omission	Probabilistic reversal
	Jean-Richard-dit-Bressel & McNally 2016	Rat	IOFC	Inactivation	Foot shock	Punishment task
	Verharen et al. 2019	Rat	mOFC	Inactivation	Foot shock	Punishment task
	Ma et al. 2020	Rat	mOFC	Lesion	Foot shock	Punishment task
Unaffected	Fellows 2003	Human	OFC	Lesion	Reward loss	Gambling task
	Manes et al. 2002	Human	OFC	Lesion	Reward loss	Gambling task
	Rogers 1999	Human	OFC	Lesion	Reward loss	Gambling task
	Verharen et al. 2019	Rat	IOFC	Inactivation	Foot shock	Punishment task
	Pelloux et al. 2013	Rat	IOFC	Lesion	Foot shock	Punishment task
	Ostlund & Balleine 2007	Rat	IOFC	Lesion	Reward devaluation	Instrumental lever task
	Bradfield et al. 2015	Rat	mOFC	Lesion	Reward devaluation	Instrumental lever task
	Panayi & Killcross 2018	Rat	IOFC	Lesion	Reward devaluation	Instrumental lever task

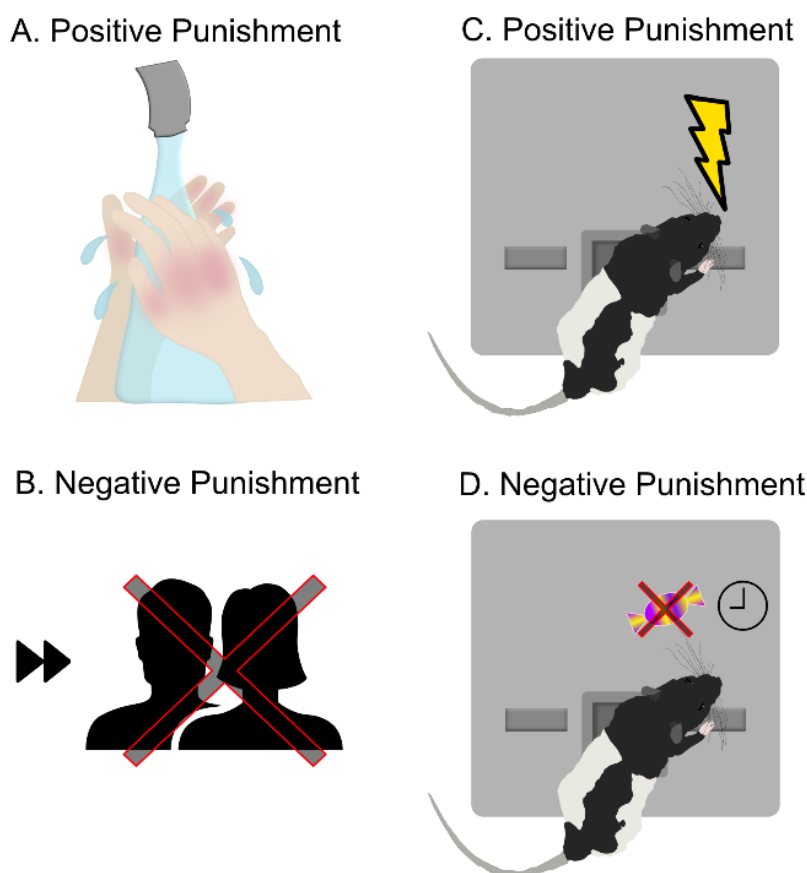
Note: OFC, orbitofrontal cortex; mOFC, medial orbitofrontal cortex; IOFC, lateral orbitofrontal cortex; vIOFC, ventrolateral orbitofrontal cortex; antOFC, anterior orbitofrontal cortex.

## **Punishment in the laboratory**

Just as individuals with compulsive disorders display an array of responses to different types of punishment, OFC manipulations in the laboratory have produced varying effects on punishment behavior. In order to try and identify common threads between studies that may enable us to make inferences about the exact nature of OFC dysregulation and punishment sensitivity, here we will group studies based on their results. First, we will review studies that demonstrated an increase in sensitivity to punishment as a result of decreased OFC function, then those that showed reduced sensitivity to punishment and finally those that manipulated OFC but found no change in punishment sensitivity. We will then attempt to identify key commonalities or differences between each of these studies with regards to the type of punishment used, the tasks employed, the OFC subregion targeted, and the species in which the study was conducted.

In order to cast a relatively broad net, we have defined 'punishment' as anything that the animal might perceive to be aversive which, following Catania (Catania, 1968) can be divided into positive punishment (defined as the presence of some aversive such as footshock or loud noise) and negative punishment (defined as the absence of something appetitive such as delayed or omitted rewards: Figure 1). For the sake of brevity, we have not considered studies of extinction or reversal learning alone as 'punishing' although they do technically involve the omission of an expected reward. Rather, the studies featured here involve the ongoing presentation of reward in such a way that its omission continues to be surprising (and therefore punishing) in, say, a probabilistic reward schedule, or after the introduction of a delay. Moreover, we have only considered tasks where the punishment was the result of a previously rewarded action and therefore have not included purely

Pavlovian studies in this review which, while incorporating aversive stimuli, are not examples of punishment per se.



**Figure 1: Positive and negative punishment in humans and rodents.** (A) In OCD pain due to excessive hand washing represents a positive punishment. (B) However, the gradual loss of rewarding social relationships is a negative punishment. In rats, we typically model these forms of punishment by (C) adding an aversive consequence, like foot shock or (D) removing or delaying an appetitive reward.

### Studies showing increased sensitivity to punishment

To the best of our knowledge, there are no laboratory-based studies conducted in humans in which OFC damage or dysregulation was directly found to increase sensitivity to punishment (outside of studies of individuals with compulsive disorders). However, an fMRI study by O'Doherty, Kringelbach, Rolls, Hornak, and Andrews (2001) conducted in healthy individuals did report an increased BOLD signal in the lateral OFC in response to monetary

loss punishments in a gambling task, and a decreased BOLD response to monetary rewards. Interestingly, they found the opposite pattern in the medial OFC, with an increased signal in response to rewards and decreased signal to punishment. Although correlational, these findings suggest that an impaired or dysregulated OFC response to gains and losses could increase sensitivity to punishment (or, indeed, reduce sensitivity to it depending on the nature of the dysregulation).

In non-human primates, Clarke, Horst, and Roberts (2015) trained marmosets to touch two visual stimuli on a touchscreen to gain access to banana juice, but during one session a week, touching one of these stimuli would also produce a mildly aversive loud noise. During this punishment session, this noise was not sufficiently aversive to stop touching the stimulus with which it was associated, and this was true whether animals received either saline infusions or inactivation of the anterior OFC using local infusions of GABA agonists (baclofen and muscimol). The next day, however, despite the removal of the noise and no infusions, responding was biased away from the punished side in animals that had received OFC inactivation during the punishment session. This result suggests that the anterior OFC is required to consolidate learning about punishment. Similar infusions didn't affect responding when only rewards were offered in the absence of the noise.

In rats, Mobini et al., (2002), gave animals a choice between a lever earning a small but immediate reward and another earning a larger but increasingly delayed reward. They found that animals with large excitotoxic OFC lesions, spanning most of the ventrolateral OFC and some of the medial OFC, were faster than controls in switching their preference to the smaller, immediate reward. This suggests that OFC-lesioned animals were more sensitive to the punishing delay. Mar, Walker, Theobald, Eagle, and Robbins (2011)



replicated this result using lesions confined to lateral OFC, whereas Rudebeck, Walton, Smyth, Bannerman, and Rushworth (2006) reported similar findings in rats trained to perform a T-maze task: the introduction of a delay in presenting the large reward caused animals with ventrolateral OFC lesions to switch to the arm that earned the small/immediate reward more quickly than controls. Finally, two studies in rats found that animals with lateral OFC lesions (Orsini, Trotta, Bizon, & Setlow, 2015) or inactivations (Ishikawa, Sakurai, Ishikawa, & Mitsushima, 2020) were more sensitive to footshock punishments.

The finding that reducing OFC activity enhances sensitivity to punishment has, therefore, been reported across different species (e.g. marmosets and rats) and in tasks that have employed a range of punishing outcomes, including both positive (footshock, loud noise) and negative punishers (delayed food reward, monetary loss). Nevertheless, one feature all of these studies have in common is that they each targeted the lateral OFC or a larger region that included lateral OFC. On this basis, therefore, one might be tempted to conclude (albeit very tentatively) that enhanced sensitivity to punishment is likely to result from reduced activity in the lateral rather than medial OFC. It is equally possible, given O'Doherty et al.'s (2001) findings, that excitation in the medial OFC might have a facilitating effect on punishment learning, but this remains to be tested.

### **Studies showing decreased sensitivity to punishment**

In one of the seminal findings linking OFC dysregulation to compulsivity, risk-taking, and disinhibition, Bechara, Damasio, Damasio, and Anderson (1994) tested patients that had sustained damage to their ventromedial prefrontal cortex due to meningioma resection or stroke, largely using the Iowa gambling task. For this task, individuals select from the most

advantageous decks of cards to achieve the optimal balance of rewards and punishments for long-term gain. The majority of healthy individuals learn to select the decks that are most advantageous in the long run, whereas OFC patients appeared to select more from disadvantageous decks despite losing money. This suggests that OFC patients were less sensitive to the punishment of losing money. This finding does come with some caveats, however, namely the fact that when the order in which the decks lead to punishment/reward was changed (Fellows, 2003), or when individuals received visual instructions about which choices were advantageous (Manes et al., 2002; Rogers, 1999), OFC patients performed similarly to controls (as reviewed by Floresco, Onge, Ghods-Sharifi, & Winstanley, 2008). The finding that switching the presentation order restored performance in OFC patients suggests these patients were less likely than controls to alter their initially learned responses. Interestingly, the finding that visual instructions allowed participants with OFC damage to overcome their impairment suggests that they did so when task requirements were made more observable, echoing findings from rodents showing that the OFC is important for inferring unobservable task information and is less important when information is observable (Bradfield, Dezfouli, Van Holstein, Chieng, & Balleine, 2015; Malvaez, Shieh, Murphy, Greenfield, & Wassum, 2019; Wilson, Takahashi, Schoenbaum, & Niv, 2014).

There are several animal studies in which OFC manipulations have also reduced sensitivity to punishment. For instance, Winstanley et al. (2004) reported that rats with whole OFC lesions increased their preference for an action that earned a larger, but increasingly delayed, reward relative to controls in direct contrast to the studies by Mobini et al., (2002) and Rudebeck et al., (2006) described above. Likewise, Pais-Vieira, Lima, and Galhardo (2007) found that lateral OFC lesions increased rats' preference for a lever that

earned a larger reward but was also more likely to result in the punishing omission of reward. With regards to the medial OFC, in the same study described above in which they found lateral OFC lesions reduced preference for a large, delayed reward, Mar et al., (2011) reported that medial OFC lesions increased preference for the delayed reward. Similarly, Stopper, Green, and Floresco, (2014) found that inactivating the medial OFC (via muscimol/baclofen infusions) increased rats' choices of the option that earned a larger reward but was also more likely to result in reward omission. And finally a study by Verharen, den Ouden, Adan, and Vanderschuren (2020) found that inactivation of both the medial and lateral OFC led to increased nose-poke responding that was punished with a time-out instead of the expected reward.

There are also several studies in which OFC inactivation has decreased sensitivity to footshock punishment. For example, Jean-Richard-dit-Bressel & McNally (2016) inactivated the lateral OFC using local infusions of muscimol and found that this increased responding on the punished lever during well-learned instrumental responding. Verharen, van den Heuvel, Luijendijk, Vanderschuren, and Adan (2019) found that medial OFC inactivation increased the number of footshocks received by rats who entered a food port 'early' (before the end of a cue) to retrieve a sucrose pellet, whereas inactivation of the lateral OFC did not affect performance in this task (although it did appear to reduce task engagement). Most recently, we (Ma et al., 2020) found that lesions of the medial OFC prevented animals from learning to avoid a lever that earned footshock. We further demonstrated that this response was specific to the instrumental punishment contingency because Pavlovian fear of a stimulus that predicted footshock (as measured by conditioned suppression of lever pressing) was intact in the same animals.

Similarly to studies in which OFC manipulations increase sensitivity to punishment, therefore, studies in which OFC inactivation produce reduced sensitivity to punishment also feature a variety of punishers (monetary loss, reward omission or delay, footshock, etc), species (humans and rats) and subregions (lateral and/or medial). As a result, it is difficult to discern a common thread. One commonality worth noting, however, is that although inactivation of the lateral OFC appears to lead variously to increased and decreased sensitivity to punishment, inactivation of the medial OFC has never been linked to increased punishment sensitivity but rather only to impaired (or unchanged, see below) punishment. There appears, therefore, to be something of a medial/lateral dichotomy emerging that, given more careful anatomical and task manipulations, future studies could further unpick.

#### **Studies in which punishment sensitivity was unaffected by OFC manipulations**

We have touched on some of the studies that have found OFC manipulations to be without effect on punishment sensitivity above. For instance, under certain conditions (i.e. altered order of presentation/explicit instructions) OFC patients were able to perform normally on gambling tasks that involved monetary losses (Fellows, 2003; Manes et al., 2002; Rogers, 1999). We also described a study conducted by Verharen et al., (2019) in which medial OFC inactivation reduced sensitivity to a footshock punishment in rats, whereas inactivation of the lateral OFC had no effect (2019). A further study by Pelloux, Murray, and Everitt, (2013) also found no effect of lateral OFC lesions on rats' sensitivity to footshock punishment delivered while lever-pressing for cocaine.

Another form of punishment that appears to be intact in animals that have received OFC inactivation is lever pressing for a devalued food outcome (reviewed in Balleine, 2019). It has been demonstrated that the sensitivity of lever pressing to outcome devaluation (via

specific satiety) is unaffected by lesions of the lateral orbitofrontal cortex whether this is conducted in extinction or in a positive punishment situation (Balleine, Leung, & Ostlund, 2011; Ostlund & Balleine, 2007). Although lesioning or chemogenetically inactivating the medial OFC affects sensitivity to outcome devaluation in an extinction test, it appears that this effect reflects the role of this structure in inferring the outcome when it is absent rather than its role in punishment. This is because when the devalued outcome was delivered contingent on lever pressing, medial OFC lesions did not affect performance in this more direct positive punishment situation (Bradfield et al., 2015). This result suggests that animals with medial OFC lesions were just as sensitive as controls to any punishing effects delivery of the devalued outcome had on responding.

Panayi and Killcross (2018) showed that lesions of the lateral OFC also left performance intact in a punished devaluation test. Specifically, they found that instrumental outcome devaluation was intact when tested in extinction, however this 'devaluation' effect (Valued > Devalued) appeared to be larger in a separate test in which the food pellets that had previously been paired with injections of lithium chloride (LiCl) to induce taste aversion were actually delivered as a result of lever pressing. This again suggests that animals with lateral OFC lesions were still sensitive to any punishing effects of delivery of the devalued pellet outcome. This is further supported by findings from Schoenbaum, Chiba, and Gallagher (1998) who showed that although a population of lateral OFC neurons appear to encode the value of an aversive-tasting quinine solution, this occurred regardless of whether or not the subsequent response changed as a result. Although far from comprehensive, together these studies suggest that neither the lateral nor the medial OFC regulate punishment sensitivity when the punisher is a devalued or aversive-tasting outcome.

Therefore, it is also difficult to identify similarities between studies in which OFC manipulations left punishment sensitivity intact with regards to species, OFC subregions, and types of punishment. One potentially interesting commonality, however, is that (with the exception of the study reported by Pelloux et al. (2013)) all of the 'punishments' in these studies involved a change in the value of the punisher at some point during the experiment. For instance, just as a devalued outcome is initially appetitive but becomes aversive, the absence of a food or monetary reward is only aversive if the participant or animal subject has been trained to expect it. It is possible then, that whereas positive punishers with stable values are likely to always engage the OFC in order to learn how to avoid them, in some circumstances animals or humans may be able to control their actions without engaging the OFC when the punishing event is initially appetitive but changes in value. It is of interest here to note that the OFC has been implicated in the identification of changes in the incentive value of appetitive rewards (Baltz, Yalcinbas, Renteria, & Gremel, 2018; Malvaez et al., 2019), perhaps making it even more curious that these studies found no role for OFC in these punishment studies. One possibility is that if alterations in incentive value take place independently of OFC (because it has been inactivated, for example), then learning to avoid that altered outcome subsequently is also OFC-independent. A simpler explanation could be that the OFC's involvement in such changes is determined the salience of the change in value, such that increases in salience make it more (or less) likely to engage OFC.

### **Alternate task variables that could determine the role of orbitofrontal cortex in punishment**

Overall, the studies examining OFC's regulation of punishment have produced highly variable results. Multiple studies have reported increased, decreased, and/or no change in

sensitivity to punishment as a result of OFC inactivation. Thus, after reviewing the bulk of the literature, we can only make a few, tentative conclusions. For instance, it appears that inactivation of the medial OFC has never been linked to *enhanced* punishment sensitivity whereas lateral OFC inactivation has. This implies that *reducing* the activity of medial OFC will more uniformly lead to a reduction in punishment sensitivity whereas reducing lateral OFC activity will have more variable effects. Another tentative pattern emerging from these results is that responding to positive punishers is more likely to engage the OFC than negative punishers. These conclusions are far from definitive, however. In studies that failed to find an effect of lesions or inactivations on punishment, the possibility remains that the failure to find any effect was due to neural reorganisation or compensation that occurred as a result of the initial damage. Another possibility is that animals employed compensatory decision-making strategies to overcome whatever psychological deficits they suffered. On the other hand, it is equally possible that there are many more studies that have failed to find an effect of OFC manipulations on punishment but remain unreported, as studies that fail to reject the null hypotheses have traditionally been less likely to be published. For these reasons, it is possible that the studies reviewed here provide a skewed vision of the OFC's role in punishment.

Despite these possibilities, the (tentatively) variable role of lateral OFC in encoding punishing contingencies is supported by an electrophysiological recording study conducted by Hosokawa et al., (Hosokawa, Kato, Inoue, & Mikami, 2007) in which they recorded from the caudolateral OFC region during a task in which animals had to respond to cues that predict juice, water, or electric shocks. Of the 65 'cue-responsive' neurons they recorded from (i.e. neurons that responded more in the cue period than in the period just prior to the cue), around 14% seemed to be reward-preferring because they showed the greatest

responses to the most-preferred cue (juice cue), moderate responses to the water cue, and the slightest responses to the least-preferred cue (shock), whereas 23% of neurons appeared to be aversive-preferring, because they responded most to the shock cue, moderately to the water cue, and least to the juice cue. The remaining neurons showing variable patterns of activation. Although correlational, these results point to the possibility that the lesions and manipulations employed in the various studies above may have preferentially targeted the reward-preferring or the aversive-preferring neuronal subpopulation of OFC neurons, producing opposite effects on punishment responding.

There are also likely to be a number of other task variables that we have not considered in the forgoing that could be critical in determining whether or not the OFC is engaged in punishment. For instance, all punishment studies necessarily involve some kind of interaction between appetitive and aversive learning, and it is possible that OFC regulates such interactions generally, not just in relation to punishment. Indeed, there are several Pavlovian studies examining such interactions that do not involve the punishment of an instrumental response, but do engage OFC (Bissonette, Gentry, Padmala, Pessoa, & Roesch, 2014; Morrison & Salzman, 2011). Again, however, the conclusion that OFC only regulates punishment learning due to regulating such interactions must again be approached with caution due to both the possible existence of unpublished studies that support this conclusion and the existence of studies that have explicitly separated the OFC's role in instrumental from Pavlovian conditioning (Balleine et al., 2011; Ma et al., 2020; Ostlund & Balleine, 2007).

Another variable that may determine OFC engagement is whether or not the punishment could be considered to be cued or uncued. As mentioned above, the OFC has



been particularly implicated in the inference of unobservable outcomes and appears to play little role in decision-making or reward learning when outcomes are fully observable. That the OFC might therefore regulate sensitivity to so-called ‘passive’ punishment – for which there is no cue indicating that an action must be performed to avoid a punisher – fits neatly with this observability hypothesis. This is because, in the absence of a cue signalling that punishment is impending, animals must mentally infer the punisher in order to avoid it; a process that likely engages OFC. The vast majority of the studies reviewed here fall into the passive punishment category because although there may be a cue that accompanies the punisher itself (e.g., a light preceding footshock, or the noise of a devalued pellet being dispensed), this cue is present only *after* the performance of an action and thus cannot inform the action itself. As noted above, when task instructions were made explicit to OFC-lesioned patients taking the Iowa gambling task *before* they made their action selection, intact performance for these patients was restored. This appears, therefore, to mimic more closely an ‘active’ avoidance-type task in which a cue is presented before the action is taken to avoid punishment. Interestingly, as we have noted previously (Manning, Bradfield, & Iordanova, 2020), there have been very few published studies investigating the role of OFC in active avoidance. As we would expect this procedure to be largely independent of OFC given its reliance on observable cues, it is possible that the lack of studies is due to OFC manipulations producing null results that were less likely to be published.

Another potentially key determinant of OFC’s engagement during punishment could be whether or not the punished animal is acting under goal-directed or habitual control. OFC has typically been thought to regulate goal-directed actions rather than habits (Bradfield et al., 2015; Jones et al., 2012), and there is evidence that punishment learning, like reward-related instrumental learning, is initially goal-directed (Bolles, Holtz, Dunn, &

Hill, 1980). As actions become well-practised, however, there is evidence that they become insensitive to punishment, at least after drug treatments of various kinds (Furlong, Corbit, Brown, & Balleine, 2018; Furlong, Supit, Corbit, Killcross, & Balleine, 2017). Nevertheless, in the majority of the studies reviewed here, it is not possible to say whether the animal was responding in a goal-directed or habitual fashion as the punishment was delivered 'online', making it impossible to separate the learning in each system as a result of observing the animals' performance. Thus, it is possible that the OFC was only recruited in the punishment studies reviewed above when behaviour was under goal-directed and not habitual control, and future studies will be required to determine this.

Finally, it is worth noting that the role of the OFC in flexible and adaptive responding under changing contingencies (e.g. reversal learning) is well established and could be a key contributor to these mixed findings. Loss of OFC function restricts the ability of an individual to adapt when contingencies change (Parkes et al., 2018; Rudebeck & Murray, 2011; Schoenbaum, Roesch, Stalnaker, & Takahashi, 2011; Schoenbaum, Nugent, Saddoris, & Setlow, 2002), therefore future studies will need to test shifts in both directions – i.e. away from increasing punishment and towards reduced punishment. This will help to rule out the effect of reduced flexibility in responding to reversed or changed contingencies and provide a clearer picture regarding punishment sensitivity.

### **Future Directions**

Clearly, there is still much work to be done in order to unravel how the OFC and its dysregulation might affect punishment responding. Here we make some specific suggestions for interested researchers – and indeed ourselves – to follow should they wish to answer some of these questions with the necessary specificity to avoid confusion.

First, in order to disentangle the relative roles of different OFC subregions on punishment learning, it will be necessary to implement a consistent punishment task and compare the effects of manipulating the different subregions on identical versions of this task. Ideally, researchers will investigate the effects of both inactivating and exciting neural activity in these populations, possibly with greater anatomical specificity than has been achieved previously. If manipulating activity in these opposing directions were to have opposing effects on punishment, this would inspire confidence regarding the role of these subregions in punishment sensitivity. Moreover, bidirectional manipulation of OFC activity mimics the dysregulation that is observed endogenously in the brains of compulsive individuals. Although it is difficult if not impossible to replicate these endogenous conditions, such artificial manipulations will still provide a strong basis on which to make causal conclusions about the effect of OFC dysregulation on punishment sensitivity.

An alternative to this approach is to employ a consistent OFC subregion manipulation and assess the behaviour of animals on a punishment task that differs on only one variable, such as in the identity of the punisher. This would allow for the direct comparison of the role of lateral OFC, say, on punishment learning after receiving footshock versus receiving a food pellet that has been paired with LiCl to render it aversive. The types of consequences proposed when defining human behaviour as compulsive are often only indirectly related to the punished action or are delayed in their effects. For example, losing ones' job, financial loss, or the breakdown of personal relationships often involve the loss of reward (negative punishment) rather than experiencing the addition of an aversive consequence (positive punishment). The use of immediate, highly salient, positive punishment in experimental settings (e.g., footshock) may not reflect the majority of punishment experiences for those suffering from compulsive disorders. It would be

interesting to explore whether sensitivity to negative and positive punishment correlates within individuals so as to unpack the role of the OFC in these different processes. In addition, the loss of rewarding events will require careful examination of the relative balance between their rewarding and punishing properties. For example, the loss of a job or a relationship may be less consequential if it is already less rewarding to the individual. Understanding how these processes interact may help to explain the persistence of compulsive actions in the context of real-life consequences.

It is also important to consider that most of the rodent studies to date have used either lesion or inactivation techniques to explore the role of the OFC through loss of function. Some compensation and broader loss of the many cognitive functions assigned to the OFC is likely to have impacted the results of these studies, particularly for lesion studies. Future studies with increased temporal specificity to capture and separate task components, as well as maintaining greater functional integrity, will hopefully provide a more nuanced understanding of which processes are controlled by OFC circuits.

Another potentially crucial task variable that could be easily manipulated involves the use of cues. As touched upon previously, the OFC has been particularly implicated in governing responding when it relies on the inference of unobservable information (Bradfield et al., 2015; Malvaez et al., 2019; Wilson et al., 2014). Thus, it is possible that studies in which the punishing outcome is explicitly cued prior to any action being selected by, say a tone, or in which discrete trials are signalled by a light turning on, will differentially engage the OFC relative to a task that is similar but in which the outcome or trial is not explicitly cued. This should be straightforward to test; however, there are some additional considerations that need to be borne in mind in these specific cases. Whereas medial OFC

appears to be important to inferring unobservable instrumental outcomes (Bradfield et al., 2015), ventral and lateral OFC have been posited to infer outcomes in the Pavlovian conditioning situation (Balleine et al., 2011; Takahashi et al., 2013). As such, whereas one might expect a primary role of medial OFC in punishment, any involvement of the lateral OFC is likely to be indirect, mediated by changes in the ability of the subject to predict the punishing event.

Finally, and perhaps most importantly, continuing to unravel the exact homology of the OFC and other prefrontal cortical regions between rodents, primates, and humans is going to be vital in comparing OFC functions in various tasks across species. There has been great progress with respect to this question in recent years, with Paxinos & Watson (Paxinos & Watson, 2014) revising their famous rodent brain atlases to reflect prefrontal homologies more accurately, and Ongur and Price (2000; Öngür, Ferry, & Price, 2003) offering greater specificity to the initial Walker's (1940) and Brodmann's areas offered by each author respectively. This work is ongoing, however. One avenue that could be particularly fruitful to pursue in this respect would be the use of a common task across rodents, nonhuman primates, and humans, to contribute to knowledge of functional homology, something that has been found to be of value previously (Balleine & O'Doherty, 2010; Griffiths, Morris, & Balleine, 2014).

## **Conclusion**

Despite inconsistencies in the direction of effect, it is clear that the OFC plays an important role in establishing sensitivity to punishment. Further, inaccurate assessment of punishment severity (either increased or decreased sensitivity) can be linked to compulsive

actions in disorders such as obsessive-compulsive disorder and addiction, although this is something that has received little attention in the laboratory-based studies so far.

Therefore, understanding the relationship between OFC dysregulation, punishment sensitivity and compulsive actions provides a promising framework for understanding the underlying mechanisms that support the persistence of behaviour in the face of negative consequences. This knowledge may aid the development of novel therapeutic strategies for interventions and treatment of compulsive behaviours.

**Funding:** This work was supported by grants from the National Health & Medical Research Council of Australia to L.A.B and B.W.B (project grant number 1148244) and to K.T. (Early Career Fellowship number 1122221), and grants from the Australian Research Council to L.A.B (Discovery project grant number 200102445) and to B.W.B (Discovery Project grant number 200103401).

## References

- Balleine, B. W. (2019). The Meaning of Behavior: Discriminating Reflex and Volition in the Brain. *Neuron*, *104*(1), 47–62. <https://doi.org/10.1016/j.neuron.2019.09.024>
- Balleine, B. W., Leung, B. K., & Ostlund, S. B. (2011). The orbitofrontal cortex, predicted value, and choice. *Annals of the New York Academy of Sciences*. <https://doi.org/10.1111/j.1749-6632.2011.06270.x>
- Balleine, B. W., & O'Doherty, J. P. (2010). Human and Rodent Homologies in Action Control: Corticostriatal Determinants of Goal-Directed and Habitual Action. *Neuropsychopharmacology*, *35*(1), 48–69. <https://doi.org/10.1038/npp.2009.131>
- Baltz, E. T., Yalcinbas, E. A., Renteria, R., & Gremel, C. M. (2018). Orbital frontal cortex updates state-induced value change for decision-making. *ELife*. <https://doi.org/10.7554/eLife.35988>
- Bechara, A., Damasio, A. R., Damasio, H., & Anderson, S. W. (1994). Insensitivity to future consequences following damage to human prefrontal cortex. *Cognition*, *50*(1–3), 7–15. [https://doi.org/10.1016/0010-0277\(94\)90018-3](https://doi.org/10.1016/0010-0277(94)90018-3)
- Bissonette, G. B., Gentry, R. N., Padmala, S., Pessoa, L., & Roesch, M. R. (2014). Impact of appetitive and aversive outcomes on brain responses: linking the animal and human literatures. *Frontiers in Systems Neuroscience*, *8*. <https://doi.org/10.3389/fnsys.2014.00024>
- Bolles, R. C., Holtz, R., Dunn, T., & Hill, W. (1980). Comparisons of stimulus learning and response learning in a punishment situation. *Learning and Motivation*, *11*(1), 78–96. [https://doi.org/10.1016/0023-9690\(80\)90022-3](https://doi.org/10.1016/0023-9690(80)90022-3)
- Bradfield, L. A., Dezfouli, A., Van Holstein, M., Chieng, B., & Balleine, B. W. (2015). Medial Orbitofrontal Cortex Mediates Outcome Retrieval in Partially Observable Task Situations. *Neuron*, *88*(6), 1268–1280. <https://doi.org/10.1016/j.neuron.2015.10.044>
- Catania, A. C. (1968). *Contemporary research in operant behavior*. New York: Scott Foresman & Co.
- Clarke, H. F., Horst, N. K., & Roberts, A. C. (2015). Regional inactivations of primate ventral prefrontal cortex reveal two distinct mechanisms underlying negative bias in decision making. *Proceedings of the National Academy of Sciences*, *112*(13), 4176–4181. <https://doi.org/10.1073/pnas.1422440112>
- Fellows, L. K. (2003). Ventromedial frontal cortex mediates affective shifting in humans: evidence from a reversal learning paradigm. *Brain*, *126*(8), 1830–1837. <https://doi.org/10.1093/brain/awg180>
- Floresco, S. B., Onge, J. R. St., Ghods-Sharifi, S., & Winstanley, C. A. (2008). Cortico-limbic-striatal circuits subserving different forms of cost-benefit decision making. *Cognitive*,

- Affective, & Behavioral Neuroscience*, 8(4), 375–389.  
<https://doi.org/10.3758/CABN.8.4.375>
- Furlong, T. M., Corbit, L. H., Brown, R. A., & Balleine, B. W. (2018). Methamphetamine promotes habitual action and alters the density of striatal glutamate receptor and vesicular proteins in dorsal striatum. *Addiction Biology*, 23(3), 857–867.  
<https://doi.org/10.1111/adb.12534>
- Furlong, T. M., Supit, A. S. A., Corbit, L. H., Killcross, S., & Balleine, B. W. (2017). Pulling habits out of rats: adenosine 2A receptor antagonism in dorsomedial striatum rescues meth-amphetamine-induced deficits in goal-directed action. *Addiction Biology*, 22(1), 172–183. <https://doi.org/10.1111/adb.12316>
- Goldstein, R. Z., & Volkow, N. D. (2011). Dysfunction of the prefrontal cortex in addiction: neuroimaging findings and clinical implications. *Nature Reviews Neuroscience*, 12(11), 652–669. <https://doi.org/10.1038/nrn3119>
- Griffiths, K. R., Morris, R. W., & Balleine, B. W. (2014). Translational studies of goal-directed action as a framework for classifying deficits across psychiatric disorders. *Frontiers in Systems Neuroscience*. <https://doi.org/10.3389/fnsys.2014.00101>
- Harsh, J., & Badia, P. (1975). Choice for signalled over unsignalled shock as a function of shock intensity. *Journal of the Experimental Analysis of Behavior*, 23(3), 349–355.  
<https://doi.org/10.1901/jeab.1975.23-349>
- Hosokawa, T., Kato, K., Inoue, M., & Mikami, A. (2007). Neurons in the macaque orbitofrontal cortex code relative preference of both rewarding and aversive outcomes. *Neuroscience Research*. <https://doi.org/10.1016/j.neures.2006.12.003>
- Ishikawa, J., Sakurai, Y., Ishikawa, A., & Mitsushima, D. (2020). Contribution of the prefrontal cortex and basolateral amygdala to behavioral decision-making under reward/punishment conflict. *Psychopharmacology*, 237(3), 639–654.  
<https://doi.org/10.1007/s00213-019-05398-7>
- Jean-Richard-dit-Bressel, P., & McNally, G. P. (2016). Lateral, not medial, prefrontal cortex contributes to punishment and aversive instrumental learning. *Learning & Memory*, 23(11), 607–617. <https://doi.org/10.1101/lm.042820.116>
- Jones, J. L., Esber, G. R., McDannald, M. A., Gruber, A. J., Hernandez, A., Mirenzi, A., & Schoenbaum, G. (2012). Orbitofrontal Cortex Supports Behavior and Learning Using Inferred But Not Cached Values. *Science*, 338(6109), 953–956.  
<https://doi.org/10.1126/science.1227489>
- Ma, C., Jean-Richard-dit-Bressel, P., Roughley, S., Vissel, B., Balleine, B. W., Killcross, S., & Bradfield, L. A. (2020). Medial Orbitofrontal Cortex Regulates Instrumental Conditioned Punishment, but not Pavlovian Conditioned Fear. *Cerebral Cortex Communications*.  
<https://doi.org/10.1093/texcom/tgaa039>
- Maia, T. V., Cooney, R. E., & Peterson, B. S. (2008). The neural bases of obsessive–compulsive disorder in children and adults. *Development and Psychopathology*, 20(4), 1251–1283. <https://doi.org/10.1017/S0954579408000606>
- Malvaez, M., Shieh, C., Murphy, M. D., Greenfield, V. Y., & Wassum, K. M. (2019). Distinct cortical–amygdala projections drive reward value encoding and retrieval. *Nature Neuroscience*, 22(5), 762–769. <https://doi.org/10.1038/s41593-019-0374-7>
- Manes, F., Sahakian, B., Clark, L., Rogers, R., Antoun, N., Aitken, M., & Robbins, T. (2002). Decision-making processes following damage to the prefrontal cortex. *Brain*, 125(3), 624–639. <https://doi.org/10.1093/brain/awf049>
- Manning, E. E., Bradfield, L. A., & Iordanova, M. D. (2020). Adaptive behaviour under



- conflict: Deconstructing extinction, reversal, and active avoidance learning. *Neuroscience & Biobehavioral Reviews*.  
<https://doi.org/10.1016/j.neubiorev.2020.09.030>
- Mar, A. C., Walker, A. L. J., Theobald, D. E., Eagle, D. M., & Robbins, T. W. (2011). Dissociable Effects of Lesions to Orbitofrontal Cortex Subregions on Impulsive Choice in the Rat. *Journal of Neuroscience*, *31*(17), 6398–6404. <https://doi.org/10.1523/JNEUROSCI.6620-10.2011>
- Mobini, S., Body, S., Ho, M.Y., Bradshaw, C., Szabadi, E., Deakin, J., & Anderson, I. (2002). Effects of lesions of the orbitofrontal cortex on sensitivity to delayed and probabilistic reinforcement. *Psychopharmacology*, *160*(3), 290–298.  
<https://doi.org/10.1007/s00213-001-0983-0>
- Moorman, D. E. (2018). The role of the orbitofrontal cortex in alcohol use, abuse, and dependence. *Progress in Neuro-Psychopharmacology and Biological Psychiatry*, *87*, 85–107. <https://doi.org/10.1016/j.pnpbp.2018.01.010>
- Morrison, S. E., & Salzman, C. D. (2011). Representations of appetitive and aversive information in the primate orbitofrontal cortex. *Annals of the New York Academy of Sciences*, *1239*(1), 59–70. <https://doi.org/10.1111/j.1749-6632.2011.06255.x>
- O’Doherty, J., Kringelbach, M. L., Rolls, E. T., Hornak, J., & Andrews, C. (2001). Abstract reward and punishment representations in the human orbitofrontal cortex. *Nature Neuroscience*, *4*(1), 95–102. <https://doi.org/10.1038/82959>
- Ongur, D. (2000). The Organization of Networks within the Orbital and Medial Prefrontal Cortex of Rats, Monkeys and Humans. *Cerebral Cortex*, *10*(3), 206–219.  
<https://doi.org/10.1093/cercor/10.3.206>
- Öngür, D., Ferry, A. T., & Price, J. L. (2003). Architectonic subdivision of the human orbital and medial prefrontal cortex. *The Journal of Comparative Neurology*, *460*(3), 425–449.  
<https://doi.org/10.1002/cne.10609>
- Orsini, C. A., Trotta, R. T., Bizon, J. L., & Setlow, B. (2015). Dissociable Roles for the Basolateral Amygdala and Orbitofrontal Cortex in Decision-Making under Risk of Punishment. *Journal of Neuroscience*, *35*(4), 1368–1379.  
<https://doi.org/10.1523/JNEUROSCI.3586-14.2015>
- Ostlund, S. B., & Balleine, B. W. (2007). Orbitofrontal Cortex Mediates Outcome Encoding in Pavlovian But Not Instrumental Conditioning. *Journal of Neuroscience*, *27*(18), 4819–4825. <https://doi.org/10.1523/JNEUROSCI.5443-06.2007>
- Pais-Vieira, M., Lima, D., & Galhardo, V. (2007). Orbitofrontal cortex lesions disrupt risk assessment in a novel serial decision-making task for rats. *Neuroscience*, *145*(1), 225–231. <https://doi.org/10.1016/j.neuroscience.2006.11.058>
- Panayi, M. C., & Killcross, S. (2018). Functional heterogeneity within the rodent lateral orbitofrontal cortex dissociates outcome devaluation and reversal learning deficits. *eLife*, e37357. <https://doi.org/10.7554/eLife.37357.001>
- Parkes, S. L., Ravassard, P. M., Cerpa, J.-C., Wolff, M., Ferreira, G., & Coutureau, E. (2018). Insular and Ventrolateral Orbitofrontal Cortices Differentially Contribute to Goal-Directed Behavior in Rodents. *Cerebral Cortex*, *28*(7), 2313–2325.  
<https://doi.org/10.1093/cercor/bhx132>
- Paxinos, G., & Watson, C. (2014). *The Rat Brain in Stereotaxic Coordinates* (7th ed.). Retrieved from <https://www.elsevier.com/books/the-rat-brain-in-stereotaxic-coordinates/paxinos/978-0-12-391949-6>
- Pelloux, Y., Murray, J. E., & Everitt, B. J. (2013). Differential roles of the prefrontal cortical

- subregions and basolateral amygdala in compulsive cocaine seeking and relapse after voluntary abstinence in rats. *European Journal of Neuroscience*, n/a-n/a.  
<https://doi.org/10.1111/ejn.12289>
- Rogers, R. (1999). Dissociable Deficits in the Decision-Making Cognition of Chronic Amphetamine Abusers, Opiate Abusers, Patients with Focal Damage to Prefrontal Cortex, and Tryptophan-Depleted Normal Volunteers Evidence for Monoaminergic Mechanisms. *Neuropsychopharmacology*, 20(4), 322–339.  
[https://doi.org/10.1016/S0893-133X\(98\)00091-8](https://doi.org/10.1016/S0893-133X(98)00091-8)
- Rudebeck, P. H., & Murray, E. A. (2011). Balkanizing the primate orbitofrontal cortex: distinct subregions for comparing and contrasting values. *Annals of the New York Academy of Sciences*, 1239(1), 1–13. <https://doi.org/10.1111/j.1749-6632.2011.06267.x>
- Rudebeck, P. H., Walton, M. E., Smyth, A. N., Bannerman, D. M., & Rushworth, M. F. S. (2006). Separate neural pathways process different decision costs. *Nature Neuroscience*, 9(9), 1161–1168. <https://doi.org/10.1038/nn1756>
- Schoenbaum, G., Roesch, M. R., Stalnaker, T. A., & Takahashi, Y. K. (2011). Orbitofrontal Cortex and Outcome Expectancies: Optimizing Behavior and Sensory Perception. In *Neurobiology of Sensation and Reward*. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/22593899>
- Schoenbaum, G., Chiba, A. A., & Gallagher, M. (1998). Orbitofrontal cortex and basolateral amygdala encode expected outcomes during learning. *Nature Neuroscience*, 1(2), 155–159. <https://doi.org/10.1038/407>
- Schoenbaum, G., Nugent, S. L., Saddoris, M. P., & Setlow, B. (2002). Orbitofrontal lesions in rats impair reversal but not acquisition of go, no-go odor discriminations. *Neuroreport*, 13(6), 885–890. <https://doi.org/10.1097/00001756-200205070-00030>
- St. Onge, J. R., Chiu, Y. C., & Floresco, S. B. (2010). Differential effects of dopaminergic manipulations on risky choice. *Psychopharmacology*, 211(2), 209–221.  
<https://doi.org/10.1007/s00213-010-1883-y>
- Stopper, C. M., Green, E. B., & Floresco, S. B. (2014). Selective Involvement by the Medial Orbitofrontal Cortex in Biasing Risky, But Not Impulsive, Choice. *Cerebral Cortex*, 24(1), 154–162. <https://doi.org/10.1093/cercor/bhs297>
- Takahashi, Y. K., Chang, C. Y., Lucantonio, F., Haney, R. Z., Berg, B. A., Yau, H.-J., ... Schoenbaum, G. (2013). Neural Estimates of Imagined Outcomes in the Orbitofrontal Cortex Drive Behavior and Learning. *Neuron*, 80(2), 507–518.  
<https://doi.org/10.1016/j.neuron.2013.08.008>
- Verharen, J. P. H., den Ouden, H. E. M., Adan, R. A. H., & Vanderschuren, L. J. M. J. (2020). Modulation of value-based decision making behavior by subregions of the rat prefrontal cortex. *Psychopharmacology*, 237(5), 1267–1280.  
<https://doi.org/10.1007/s00213-020-05454-7>
- Verharen, J. P. H., van den Heuvel, M. W., Luijendijk, M., Vanderschuren, L. J. M. J., & Adan, R. A. H. (2019). Corticolimbic Mechanisms of Behavioral Inhibition under Threat of Punishment. *The Journal of Neuroscience*, 39(22), 4353–4364.  
<https://doi.org/10.1523/jneurosci.2814-18.2019>
- Walker, A. E. (1940). A cytoarchitectural study of the prefrontal area of the macaque monkey. *The Journal of Comparative Neurology*, 73(1), 59–86.  
<https://doi.org/10.1002/cne.900730106>
- Wilson, R. C., Takahashi, Y. K., Schoenbaum, G., & Niv, Y. (2014). Orbitofrontal Cortex as a

Cognitive Map of Task Space. *Neuron*, 81(2), 267–279.

<https://doi.org/10.1016/j.neuron.2013.11.005>

Winstanley, C. A. (2004). Contrasting Roles of Basolateral Amygdala and Orbitofrontal Cortex in Impulsive Choice. *Journal of Neuroscience*, 24(20), 4718–4722.

<https://doi.org/10.1523/JNEUROSCI.5606-03.2004>