# The Descent of Preferences

David Spurrett (UKZN) — spurrett@ukzn.ac.za — 2018

## Abstract

More attention has been devoted to providing evolutionary scenarios accounting for the development of beliefs, or belief-like states, than for desires or preferences. Here I articulate and defend an evolutionary rationale for the development of psychologically real preference states. Preferences token or represent the expected values of discriminated states, available actions, or action-state pairings. The argument is an application the 'environmental complexity thesis' found in Godfrey-Smith and Sterelny, although my conclusions differ from Sterelny's. I argue that tokening expected utilities can, under specified general conditions, be a powerful design solution to the problem of allocating the capacities of an agent in an efficient way. Preferences are for efficient action selection, and are a 'fuel for success' in the sense urged by Godfrey-Smith for true beliefs. They will tend to be favoured by selection when environments are complex in ways that matter to an organism, and when organisms have rich behavioural repertoires with heterogenous returns and costs.

The rationale suggested here is conditional, especially on contingencies in what design options are available to selection and on trade-offs associated with the costs of generating and processing representations of value. The unqualified efficiency rationale for preferences suggests that organisms should represent expected utilities in a comprehensive and consistent way, but none of them do. In the final stages of the paper I consider some of the ways in which design trade-offs compromise the implementation of preferences in organisms that have them.

# 1. Introduction

What are preferences for? Or, a little more narrowly, why should we expect some natural, biologically evolved, agents to have preferences? The heart of the answer that I'll defend is that preferences enable efficient action selection,[1] which is to say the deployment of the relatively flexible capacities of the agent in fitness-favouring ways. I use the term 'preference' to refer to a psychological or cognitive state. To *have* preferences is not merely to *exhibit* them in behaviour that more or less consistently maximises something, as in strictly behaviourist 'revealed preference theory' (Samuelson 1938).[2] Preference-revealing behaviour is relevant here, but my primary focus is the function and development of a psychological capacity. An agent that has preferences produces and processes states that somehow represent the expected values of local world-states or outcomes that it can detect or anticipate, or actions that it can perform.[3] A few preliminary clarifications about what I mean by having preferences may be helpful.

First, I'm more interested in the usefulness of having a set or a system of preferences than with the functionality of *individual* preferences. For an agent that has *some* system of preferences, it is easy enough to pick one and consider whether or when it is advantageous. My concern is with the question why agents might be expected to have the general capacity to represent some outcomes or actions as more or less valuable than others, and to focus on that I'll mostly assume that the specific rankings are advantageous.

Second, a natural agent can have preferences more or less *completely*, in the sense of covering all or only some of the outcomes it can distinguish or actions it can perform. (For example, a creature might have richly detailed preferences over what it eats, but switch between foraging and other activities in ways not mediated by preferences.) It can also have them more or less *consistently* in the sense that the ordering may or may not respect transitivity, be stable over time, or otherwise respect requirements typically found in normative theories of preference. For much of what follows I'll write as though preferences are fairly well-behaved but this is for expository simplicity. My initial aim (§3 and §4) is to identify an ideal towards which selection could sometimes be expected to move were other constraints not relevant. Other constraints, though, are *always* relevant and after articulating the primary argument I spend some time considering them (see §5).

Third, and finally, my target here, having preferences, is different from — and more modest than — having *desires*. Although there are various theories of desire to choose between, most

---

[1] I understand 'action' here in an inclusive way largely interchangeable with 'behaviour', and covering any functional allocation of the relatively short-term capacities of an organism. (See §3.)

[2] Samuelson (1938) sought to "develop the theory of consumer's behaviour freed from any vestigial traces of the utility concept" and although economists ended up retaining the term 'utility' it was generally understood to be cognitively noncommittal.

[3] I'll assume, but not defend, Shea's 'varitel semantics' account of sub-personal representation (Shea 2018). Key advantages of his view include relaxing the teleosemantic requirement of determinate consumers (at best an

agree that desires are *personal* states, often conscious, that can — in partnership with beliefs — feature in rationalisation of intentional action. Preferences, on the other hand, are sub-personal. They're not paradigmatically conscious (although they might be consciously accessible to some agents) and to do their work they need not feature in deliberation, or reason-giving explanation. This means that my project of naturalistic explanation is different from that of Sterelny (2003) to the extent that he attempts to provide plausible evolutionary rationales for the development of states approximately corresponding to the folk-psychological categories of belief and desire.[4] In the sense at issue here, many non-human animals, including some insects and invertebrates, plausibly have preferences (see §6).

There are is another set of questions about preference and evolution that I'm not directly concerned with here. It concerns *how* preferences developed in the actual history of life and cognition, including whether they were 'invented' once (like hearts) or many times (like eyes). I'll say a little about these matters, mostly in (§6), but my primary aims here are to characterise the capacity I'm calling having preferences, and to specify fairly general conditions under which that capacity would be beneficial to an organism. I take it for granted for now that there is a genuine explanandum here: some real organisms do represent the values of some of the options open to them, i.e. that they have preferences. There's room to argue over how many types of organism do this, as well as over what else some of them might have. Although some of what follows below may be relevant to human cognition and choice, I'm not directly concerned with humans, which raise additional complications. It's probably better to have rats and insects in mind than people.

## 2. Evolution, Cognition and Complexity

The argument of the following sections is an application of a traditional approach, more recently articulated by Peter Godfrey-Smith (1996, 2002) under the label of the 'Environmental Complexity Thesis' or ECT. This is a view about the function of cognition in general, and so should be expected to apply to the restricted cognitive capacity of having preferences. In its shortest statement, the ECT says that "…the function of cognition is to enable the agent to deal with environmental complexity" (2002, p225). The ECT proposes, that is, that organisms capable of cognition can respond more effectively to significant heterogeneity ('complexity that matters') in their environments, and that at least sometimes investment in these capacities can pay their way.

Consider, for example, the predicament of a nesting reed-warbler, sometimes exploited by cuckoo brood-parasitism. The difference between an egg laid by the reed-warbler and one laid by a cuckoo is highly significant from an evolutionary point of view. This is a vivid example of

---

idealisation when applied to real neural processes), and openness to a variety of etiological content stabilising processes.

[4] Matters are complicated by the fact that Sterelny (2003, and 2001) sometimes uses 'desire' and 'preferences' interchangeably. But it's clear that he's considering states that represent the *goals* of action, which is more than I require preferences to do (see Spurrett 2015).

environmental heterogeneity. Despite its importance, the difference isn't a trivial one to detect because cuckoo eggs can closely resemble those of the birds they exploit. Cognitive mechanisms can — nobody says that they *must* — play a role in responding to this, for example by making a reed-warbler more likely to reject an egg in its nest if it has not yet laid any in that season itself, or if it has recently seen cuckoos nearby.[5] Cognitive mechanisms can, that is, operate to make the contingencies in what behaviours are produced (and when) more appropriately sensitive to such heterogeneity in the environment as an organism can detect. 'Dealing with' environmental complexity, then, is largely a matter of doing the right thing, from a repertoire which itself might change, at the right time, or doing so more often than otherwise. That is to say that Godfrey-Smith's view shares with others, including Dennett, the feature that cognition functions to help answer the question implicitly faced by any organism able to produce behaviour or action, 'Now what do I do?' (Dennett 1991, p177).[6]

In some expositions of the ECT, the design and control architecture of the organism's body is taken largely for granted, and the contribution of cognition described in terms of determining what to do with that body given the (changing) state of the external world. In contrast some, including Fred Keijzer (e.g. Keijzer 2015; Keijzer, van Duijn & Lyon 2013) have argued that the earliest job of neurons, and at least some of what deserves to be called cognition, is in co-ordinating the physical capacities of an organism independent of the state of the external world, and even in the absence of any sensory transducers. According to the 'Skin Brain Thesis' (SBT) contractile tissues pose a control problem, and almost without exception organisms with muscles also have neurons. The precursors of brains, Keijzer argues, arose to deal with making use of bodies with contractile tissues *before* external sense and environmental complexity were significant factors.

The ECT and SBT are not, I think, most usefully thought of as mutually exclusive competitors, but rather as differences of emphasis that can be seen as special cases of a more general view. In recent work Godfrey-Smith and Keijzer have collaborated to sketch an 'option space' for early neural evolution in just this way (Jékely, Keijzer & Godfrey-Smith 2015). I don't, therefore, need to take sides. In what follows the demands of controlling the body *and* the demands of dealing with the external environment, both of which can be complex, will feature in explaining the function and evolution of preferences. To make this explicit, I propose to work with the following gloss of the general function of cognition:

> "The function of cognition is to enable the agent to co-ordinate its (possibly complex) capacities, which can include co-ordinating those capacities with environmental complexity."

---

[5] This example is used in Sterelny (2003, Chapter 2). Darwin discussed cuckoo brood parasitism in Chapter 8 of the *Origin of Species* to illustrate how natural selection might impact on behaviour. Over half of the world's bird obligate brood parasites are species of cuckoo (a little under half of the species of cuckoo are brood parasites). Cuckoo brood parasitism evolved independently three times (Payne 2005) and there is significant variation in the details of the various arms races between host and parasite species.

[6] See also Dennett (1984) for an earlier discussion of the importance of 'producing future'.

I think of this as a friendly amendment of the ECT, standard treatments of which might have said relatively little about the demands of controlling the body (e.g. Godfrey-Smith 1996, Sterelny 2003), focusing rather on the organisation of behaviour in relation to external contingencies, but which nonetheless don't *deny* that these demands are important. The function of cognition is still dealing with complexity. I merely emphasise that some of the complexity might be in the agent itself.

Although I argue that it makes sense that preferences would have evolved (under conditions to be outlined below), the case is largely agnostic between some accounts of evolutionary advantage. In particular I think it is a matter of taste for the reader whether she interprets this in terms of fitness advantage, or as a type of 'viability explanation' (Wouters 1995).[7]

# 3. The Basic Efficiency Rationale for Preferences (ERP)

The core of my proposal, as noted, is that preferences are for enabling efficient action selection. The answer to the question why preferences evolved is:

(ERP1) Preferences enable efficient action selection.

Before defending this claim, let me clarify how it is to be understood. 'Action' is any *functional* activity that the agent produces, that is any deployment of its relatively transient 'degrees of freedom', whether muscles, glands or other kind of effector, individually or in concert.[8] This is a deliberately broad and inclusive use of 'action', more in keeping with practice in artificial intelligence than in some philosophical traditions where action is reserved for a specific class of intentional activity. The term behaviour is used in varied ways by psychologists, ethologists and others. If we restrict behaviour to cases of *functional* activity (Millikan 1993a, 1993b), action 'action' here is interchangeable with 'behaviour'. I'll use 'activity' to include deployments of degrees of freedom irrespective of whether they're functional. Not all activity is action.

There are few hard boundaries to be had in the ragged and gradual world of living things. Just as growing and reproducing may not always be strictly distinguishable, for example in the case of plants that send out runners, so it is for behaving and developing.[9] While some plant activities, such as the closing of a Venus fly trap, satisfy fairly straightforward criteria for 'behaviour', some recent discoveries in plant cognition focus on instances of what could be called discriminating development. Findings, such as that growth can be sensitive to histories

---

[7] Wouters argues persuasively that traits can be explained by reference to their contribution to viability, rather than comparative advantage. Okasha (2018) argues that it isn't generally or necessarily true that natural selection will tend to optimise fitness, and concludes that justifications for adaptationism must be empirical.

[8] At least sometimes, this includes *not* activating some capacity. (Keeping still is a perfectly good instance of behaviour.) It's also worth noting that the capacities activated need not be within the organism's own body. (Phenotypes are sometimes extended.) I'll focus mostly on within body capacities here.

[9] Godfrey-Smith (2002) contains a useful discussion of some of the difficulties here.

of exposure to conditions that are not themselves intrinsically beneficial or harmful, but were correlated with ones that were, are sometimes described in terms of 'learning' or 'conditioning' (Gagliano et al. 2016).[10] I acknowledge these difficulties, but will largely be able to ignore them here because the argument about the function of preferences concerns the paradigmatic animal actions of behavioural ecology: moving about, foraging, fighting, nesting.

Given this, *action selection* is a placeholder for *whatever* it is that determines which degrees of freedom are activated or not, and to what extent to produce behaviour. This determination needn't always be cognitive, let alone richly so. Some capacities can be quite directly triggered by local — external or internal — conditions, independent of what is happening with the rest of the organism. Saying that there is such a thing as action selection is not yet to take any view about how, or by what mixture of means, it is achieved. (Action selection is distinct from, although perhaps not independent from, action production, where production is making action from capacities.)

Finally, action selection is *efficient* in proportion to goal satisfaction achieved and in inverse proportion to costs expended. For my purposes it doesn't matter much how you conceive of the goals of an organism, as long as we're agreed that the overall goal or goals can be pursued directly or via intermediate or contributing goals. So, whether you're inclined to suppose that an organism has the single goal of reproductive success, or some mixture of survival and reproduction, or the possibly conflicted net effect of the reproductive goals of its genes, effective pursuit of *that* will involve pursuit, over smaller time-scales, of some mixture of calorie intake, hydration, rest, mating, avoiding being eaten by a tiger, and so forth. These goals needn't, and generally won't, be represented as goals *for* or *in* the agent, even though some agents might sometimes represent some of them.

Greater efficiency in selecting actions given some goals is an obvious kind of advantage over anything with the same goals. Given similar goals and capacities, an organism that allocates the capacities to more efficiently achieve the goals is doing better. The presumption that meaningful standards of efficiency can be determined, and used to interpret observed behaviour, is essential to most empirical behavioural ecology. Not only that, in some areas — for instance foraging — living organisms have been shown to be extraordinarily efficient in their patterns of behaviour. To make the basic argument for (ERP1) I defend two claims. The first is that achieving efficient action selection can often be difficult. The second is that having preferences is a plausible way of dealing with these difficulties.

The main reason that achieving efficient action selection can be difficult is that actions generally have varying (and multi-modal) costs, and varying (and multi-modal) returns. The costs include direct expenditure of energy, the depletion of specific 'fuels' or resources such as water and salt, as well as time spent, exposure to various risks such as predation in the course

---

[10] Gagliano and colleagues set up a Y-maze task, and found that a neutral cue predicting of light was associated with discriminating growth.

of the action, and the opportunity cost of forgone actions available at the same time. The returns can be as varied as the needs of the organism and include hydration, nutrition (itself potentially further varied, and sometimes including feeding dependent young), rest, access to mating opportunities, acquisition of nesting materials or control of a nesting site, and so forth.

The detectable state of the world can indicate, to varying extents, and sometimes subject to the hostile deception that Sterelny (2003) says makes some environments informationally opaque, likely costs and returns: the visible scene might include both watering-hole and predator threat, for example, or no food in the immediate vicinity, but two different patches of green in the distance. The detectable state of the organism *itself* can in turn indicate what needs are most urgent, or what costs can be most easily sustained. Sterelny (2003) argues that internal environments — being parts of an agent with united and consistent interests — will tend to be devoid of conflict, and hence that internal signals of need and capacity will tend to be highly reliable. This is, I think, over-optimistic in underestimating the importance of internal conflict, and the complexity of the task of tracking needs and resources. On the first point, it's not generally true even that the genes in a single individual have precisely the same interests (Haig 2002, Burt & Trivers 2006).[11] In addition, almost any organism large enough is itself the habitat for other organisms, with their own possibly hostile interests of their own. Internal signals, that is, can *sometimes* be subject to conflict and manipulation. That aside, though, and turning to the second point, even honest internal signalling has to reckon with the complexity of tracking needs and resources.

As Sterelny says, an environment is *informationally translucent* to organisms when states that matter to it "map in complex, one to many ways onto the cues they can detect" (Sterelny 2003, p. 21). These conditions can be satisfied in hostility-free internal environments, in several ways. One way is because of constraints on what can easily be transduced or detected. Just as in the external case, not all internal states have unique signatures that cost-effective transducers can specialise in detecting. Non-nutritive sweeteners, for example, trigger transducers whose 'proper function' is to respond to sugars that *can* be digested. The responses of salt receptors, depending on the action of ion channels, are also sensitive to the ambient sodium concentration in the organism, so the resulting neural signals can be highly ambiguous (e.g. Bertino, Beauchamp & Engelman 1982). Another way is because motivationally relevant states can depend on multiple cues. Information about temperature in humans, for example, is drawn from multiple receptors of different types that are distributed non-uniformly across the surface of the body. As Akins notes, even on the human face the ratio of cold to warm receptors varies from about 8:1 on the nose, 4:1 on the cheeks and chin, while the lips have almost no cold receptors (Akins 1996, p. 346). Any 'net' signal that might drive behaviour — to seek more or less warmth — will require these signals to be integrated in some way. More generally, bodily states can span multiple organs and tissue types, with

---

[11] Both Haig and Trivers have suggested that this intragenomic conflict predicts intrapersonal conflict, a proposal that I assess in Spurrett (2016).

varying speeds of signalling, and latencies in responding to actions that affect them. In the language of Sterelny (2003) some internal states themselves require 'robust tracking'.

Even if internal signalling was both honest *and* consistently accurate, the problem of matching the fluctuating needs and capacities of the organism to the varying profile of opportunities and risks in the detectable environment can itself be complex. Doing 'the right thing at the right time' is a trade-off problem with varying and fluctuating parameters. Behaviours, as noted, have varying returns (in calories, specific nutrients, hydration, mating opportunities, acquisition of nest-building materials, etc.) all of them uncertain to varying degrees, and their execution carries varying costs (in calories, hydration, exposure to predation, the possible returns of behaviours foregone, etc.). Some decisions involve options with the same dimensions, for example patch switching when foraging is understood as a problem of comparing actions with different expected rates of calorie intake. Others involve options or bundles with at least some dimensions that *aren't* shared, such as when choosing between pursuing hydration at a site with low predator risk and eating somewhere with risks and possible gains in social rank or mating opportunities. Making action selection efficient involves dealing with these many and varied mappings.

The discussion so far has focused on *selection* between actions. But I noted earlier than action selection and action production need not be independent. Not only that, production involves efficiency problems of its own. If behaviour is *functional* activity (or inactivity) then not all activity will amount to behaviour, because some will be noise or otherwise non-functional.[12] Making functional activity out of a collection of degrees of freedom is not generally a trivial matter. This is so *because* many behaviours require multiple degrees of freedom to be co-ordinated. They may require multiple muscles or joints to act compatibly to produce the behaviour, and also require that these ones aren't thwarted by inappropriate activation of others not directly involved. The mappings from the many degrees of freedom in a whole body to what is needed to produce even simple behaviours like pecking at a stationary key can be highly complex.

The selection of one action over another often carries an opportunity cost not merely in other allocations of the same degrees of freedom as are directly involved, but also in others whose execution have to be transiently suppressed to produce the selected one. The facts of the matter about what behaviours might be possible are not, furthermore, unchanging. The orientation of the body, and its disposition with respect to various surfaces, can determine whether jumping is an option (you can't generally do that lying down), or what combination of muscle loading across the joints of an arm is required to put something into one's mouth, or push it away.

While some actions can be achieved by modulating the activity of one or a few isolated degrees of freedom, as in the case of blinking, or releasing fluid from a tear duct, many

---

[12] Functional needn't mean 'successful' here. (On this I follow Millikan, e.g. 1984).

cannot. Consider a complex bit of animal anatomy with multiple degrees of freedom such as a primate forelimb with shoulder, elbow and a set of digits. Some combinations of allocations of the effectors, such as simultaneously flexing and relaxing the same muscle, or flexing those that would move a segment in one direction around a joint while not relaxing those that would move it in the opposite direction, are at least incompatible and could even be harmful if attempted with high intensity. For a movement comprising components around more than one joint at once, along a single limb or several — as in reaching on toe-tip for something on a high shelf — many such combinations of contraction and co-operative relaxation may be required. In many natural brains, part of this complexity is handled by means of what Sherrington called 'final common paths'. He argued that upstream of the specific effectors (in what he called the 'afferent arc') competing allocations ('reflexes') converged in a final common path, where only one reflex (as opposed to some combination or sum) would gain control of the effectors downstream (e.g. Sherrington 1906, pp. 117-118).

The general idea of a final common path is a good and useful one, but I'm not suggesting we be strictly Sherringtonian about them. We needn't suppose that all final common paths are anatomically local, or unique (both redundancy and distribution are allowed). In addition we should grant that many of the patterns of activation and inhibition required for various behaviours are learned, and likely involve what Clark (1997) called 'soft assembly': control solutions dependent on and mediated by the structure and properties of the behaving body and its environment. That said, Sherrington's insight provides a useful corrective to the tradition, going back at least to Brooks (1991) that regards almost any convergence in a control system as a symptom of allegiance to muddled models of intelligence and cognition. Not all bottlenecks are bad.

The point I'm emphasising is that the problem of making behaviours out of combinations of capacities is *itself* one that involves trade-offs, between other possible allocations of individual capacities and combinations of them, over and above whatever the metabolic and other direct costs of this or that action might be. Motor control, that is, is *not* generally independent from the problem of selection, even if many accounts elide this by taking as their starting point the already selected behaviour or movement. Besides the more obvious efficiency issues relating to getting something done without excess noise, energy expenditure, etc., is the efficiency problem of competing allocations of the same capacities.

I've argued that making action selection efficient is a trade-off problem, involving co-ordinating the varied and multi-modal needs and capacities of the body, some of which constrain one another, given the also multi-modal risks and opportunities in a changing environment. Suppose that this is right. Then the argument that having preferences is a way of dealing with these difficulties can be made fairly simply. Preferences are states that attach values to possible actions or detectable world states. If those values are appropriately responsive to whatever the overall goals of the organism are, as well as its changing needs and opportunities, then they could be exploited to make action selection more efficient. If you think organisms are in the fitness business, preferences can represent expected returns in

fitness, or predictors of it. Then the members of some set of available actions would be associated with states that reflect their relative contribution to overall goals, and the processes of action selection could use these states to choose better actions.[13] Preferences just are representations of the values of actions given states (or world or organism).

Preferences, that is, are situation-specific rankings of available actions. So, for example, in a dehydrated organism actions that likely lead to water consumption should be valued more. In an animal that has just detected a predator, actions from the its defensive repertoire should come to be valued more, and so forth. (They could also be updated in the light of experience of the gains and costs of actions given detectable cues, although preferences are not necessarily associated with learning.) And the advantage comes from having a system of them that is appropriately responsive to changing conditions. Making action selection efficient is a real problem, and having preferences is a possible solution.

This shouldn't be surprising, because venerable tradition has it that effective agents will act *as if* they assign subjective probabilities to current and future states of the world (including the states that may follow available actions), and select actions that maximise expected utility (Ramsey 1931, Savage 1954). A more contemporary version of this idea suggests that effective agents should implement (and evolve to implement) some form of Bayesian rationality, including conditionalization as a mechanism of belief updating (Okasha 2018, Chapter 6).

That preferences 'could' enable efficient action selection is all I am after. Some argue that efficiency 'must' have recourse to preferences. Shizgal and Conover, for example, reporting on rat subjects that traded off mutually exclusive rewards in the form of trains of brain stimulation reward (BSR) and sugar solution, put it like this:

> "In natural settings, the goals competing for behavior are complex, multidimensional objects and outcomes. Yet, for orderly choice to be possible, the utility of all competing resources must be represented on a single, common dimension." (Shizgal & Conover 1996)

Shizgal and Conover's rats showed flexible sensitivity to the opportunity cost of foregone rewards, including in cases where one reward was a bundle including both BSR and sugar solution. Their hunch is that pulling off such tracking *without* preferences would be mysterious, or maybe — if we take the 'must' seriously — impossible. This is a stronger claim than I want to defend for now. I just need 'having preferences' to be a viable way of enabling more efficient action selection.

---

[13] It doesn't matter here whether the better of two, or best of more, action is selected always, or merely more often. Simply doing a better thing more often than otherwise is sufficient as a notion of advantage. (A variety of final outcome determination processes might respond to preferences in different ways. See §6.)

# 4. The Refined ERP

As defended above, the Evolutionary Rationale for Preferences is vulnerable to some standard objections and worries. Here I outline and briefly address the main ones. Some of them concern scenarios in which an organism might do *without* preferences, and involve differing views about how effective or successful such an organism might be. Others focus on how *expensive* building and operating a system of preferences might be in comparison to other options. It certainly isn't the case that having preferences is the *only* way to deal with action selection. Also, no reasonable approach to efficiency can ignore the costs of the means used to make other processes efficient. The upshot of thinking through the objections is a qualified restatement of the ERP.

Preferences in the sense at issue here are usefully understood as a kind of sub-personal *representation*. They are produced or modulated by systems tracking the changing needs of the organism, and by systems tracking the local opportunity space, in some cases also updated in the light of the history of consequences of actions performed in various conditions. And they have effects in systems arbitrating between mutually exclusive allocations of the scarce means, primarily bodily degrees of freedom, at the disposal of the organism. I'm not going to *defend* the claim that they're representational here — there are competing accounts of what it is to represent, and conflict tangential to my purposes over whether representations are good or bad things at all.[14] I don't think it matters much if my representation talk is replaced with reference to 'tracking' or some other notion, as long as the tracking function is discharged.

One likely objection, nonetheless, is provided by the strands of anti-representationalist thinking found in artificial intelligence and robotics. On those views, representations are unnecessary, because the world can serve as 'its own best representation' (Brooks 1991), and inefficient, because the demands of processing them will create a congested 'bottleneck' leading to paralysis or unacceptable delays (Clark 1997)[15]. The denial of preference representations here is a by-product of any *general* rejection of representations, and these often occur — as in Brooks — without explicit reference to preferences or utilities. Neither alleged problem is decisive here.

It is undoubtedly correct that degrees of freedom can be yoked to discriminators or transducers fairly directly, and hence be controlled by processes that are representational only according to notions of representation so inclusive as to be suspect. An automobile air-bag triggered by an accelerometer need not have any connections with other control systems. (Even if we do interpret the signal from accelerometer to air-bag as a representation, it

---

[14] I'm persuaded by Nick Shea's argument that reward prediction errors are meta-representational (Shea 2014). And what they meta-represent, reward predictions, are one version of what I'm calling preferences. See §6. (Not all preferences need be meta-represented by prediction errors.)

[15] Andy Clark wrote of Brooks that it is 'conceivable that much of human intelligence is based on similar environment-specific tricks and strategies' (1997, p31). Clark's (1997) discussion of Brooks is so influential that it's led many to write as though Brooks (1991) includes the expression "representational bottleneck", which it doesn't.

doesn't encode a value, and isn't compared with anything else, so isn't a *preference*.) Direct linking of transduction and activation of this kind can be combined in various ways, including some hierarchies, leading to whole agents that have been variously described as 'subsumption architectures' (Brooks 1991), and 'detection agents' (Sterelny 2003). We can think of the individual linkings as the parts of 'layers', or as 'detection systems', or 'reflexes', etc. Let's allow that such agents are both possible and actual. The issue is how efficient they are.

While it might be true that the occurrent, local environment can guide behaviour without having to be *duplicated* in a model, what preferences are supposed to represent aren't facts contained in the external environment at all, but rather facts about the needs and priorities of the agent itself, and facts about the (net) values of available actions conditional on those needs. Even if the world and body can 'represent themselves', they won't automatically represent the returns (in what matters to the organism) on actions.

Similarly, the complaints about representational bottlenecks are at their most effective in original application, which was cases where courses of action were selected by consulting richly detailed world-models, constructed 'after' perception, drawing on a stored model, and 'before' action was executed. But the value representations that preferences consist of just aren't 'world models' in that sense at all. They can be very sparse representationally speaking - either consisting of an additional layer or component of the existing motor or sensorimotor resources used to control an action, or monitor a metabolic need or state. Here is an empirical example: neuroeconomic research detected neural correlates of relative expected return from choices expressed through eye movements in the motor circuits producing the eye movements themselves (Platt & Glimcher 1999, Dorris & Glimcher 2004, Glimcher 2011). Eye movements were selected for this research partly *because* the relationship between neural activity and produced movement is comparatively simple and tractable (when the head is held stationary). It had been known since Shadlen & Newsome (1996) that activity levels in LIP neurons predict impending saccades. What this research showed was that the differing levels of neural activity for two mutually exclusive actions, prior to one being produced, stood in a regular relationship to the varying returns on each action. The only 'bottleneck' here is the Sherringtonian final common path blocking mechanically inconsistent eye movements, and the processing of preferences actually *exploits* it. Preferences, then, needn't generally involve *additional* bottlenecks, or impose an excessive representational burden. While the world might 'represent' the saccadic targets well enough, it won't itself represent the varying consequences of selection relative to the appetite of the animal.

This won't satisfy everyone. Producing and processing preferences still isn't *free*, and might be prohibitively expensive, at least compared to some alternative. A leading general idea, embracing several variants, about what to do instead, is for behaviours, activities or goals to be arranged in some kind of hierarchy of trumping relationships that don't depend on encoding values. What 'subsumption architecture' originally *meant*, recall, was that under certain conditions the activity of one 'layer' of a control system would simply over-ride ('subsume') that of another. A natural example is provided by sea-slugs, that are indiscriminate

omnivores and would eat the eggs that they had laid themselves if it wasn't that during egg-laying they released a hormone that inhibited eating behaviour (Davis, *et al* 1977). The value of eating their own eggs isn't *represented* as low or negative. The one behaviour simply suppresses the other.

Sterelny (2003) has argued that in most creatures there's no need for the development of preference states because internal environments will tend to be 'transparent' in the sense of being characterised by highly reliable signals, unlike external environments which are prone to pollution from the conflicting interests of other organisms. In these transparent internal environments, he suggests, some kind of relatively fixed motivational hierarchy, perhaps of 'drives', can be a sufficiently effective mechanism of action selection. As noted in (§1) Sterelny is partly talking about a somewhat different topic to me, even though he sometimes uses the term 'preferences', since the project in his (2003) is to develop a plausible natural history of cognitive states approximately corresponding to folk-psychological beliefs and desires. The goal-representing proto-desire states he's concerned with are fit partners for proto-beliefs, perhaps enjoying some degree of accessibility to the agent that has them. Preferences in the sense at issue here are more modest states. That said, let me address the two lines of argument in Sterelny (2003).

First, as noted in (§3), while it may be approximately true that internal environments contain less hostility than external ones, it doesn't follow that tracking the needs of an organism is a trivial or simple matter, let alone that the state of the organism can be treated as 'its own best representation'. The needs of an organism vary in many dimensions, including over various time-scales, may involve quantities that aren't or can't be tracked in any way that is both cheap and reliable, and can vary in their urgency, either because of their own values or in interaction with other factors.[16] The needs of an organism can be many-dimensional, independently varying, and opaque.

Second, it is possible that the view defended here and Sterelny's is more compatible (in some respects) that may superficially appear. If 'drives' represent the strength of needs, or the strength of motivation for certain behaviours, then they're doing the same general thing that I'm saying preferences do, and we're just using different terms for (at least roughly) the same target. I'd like to think that is at least partly true, and that the appearance of disagreement is because Sterelny is using 'preference' some of the time as a stylistic variant for desire-like state, and aiming to provide a natural history of cognitive capacities approximately like folk-psychological beliefs and desires, where a key part of what desires do is represent *goals*. Sub-personally representing values of action-state pairings, though, needn't involve representing goals.

---

[16] McFarland & Sibly (1975) has an instructive discussion of the challenges of constructing an idealised 'command space' representing the behavioural relevance of the changing needs of an organism (see also Houston & McFarland 1981).

Matters are complicated because drive talk is far from univocal, the term having been used to name various theoretical constructs in motivational psychology, most of which have been superseded and displaced by reward-based concepts more compatible with preferences (Berridge 2004). Not only that, behavioural ecologists sometimes use drive talk in ways that are agnostic about cognitive mechanism, referring to variation in the tendency to behave in some way. Sterelny doesn't explicitly *endorse* any of the models of drives that were rejected along the way. Still, he's undoubtedly suggesting that a lot could be achieved in behaviour selection by means of motivational states more modest than desires, which includes that they don't represent goals, but may represent urgency or strength of motivation in some way integrating bodily needs and detected opportunities. Subject to keeping a close eye on what words we use when, I don't think we're substantially at odds here.

This leaves one outstanding matter: Sterelny's explicit endorsement of an idea of 'fixed motivational hierarchies', which I take to mean that some needs or priorities can — under specific conditions — 'trump' others in ways that circumvent any need for representations of their relative values to be processed. The proposal here is a version of the family of views — already noted — including subsumption architectures and other non-representational trumping hierarchies. And it is undoubtedly true that complicated organisational matters of various kinds can be made simpler and more tractable by switching between distinct 'modes' of activity within which a sub-set of the total behavioural options are in play, and the mix of 'internal' and 'external' information gathering is tailored to the current mode. Not only that, many animals arguably *do* something like this - switching between territory-defending mode, mate competition mode, etc.

This doesn't dissolve the argument for preferences for several reasons. First, preferences have a role to play *within* modes. Consider an animal in foraging mode. Its options have varying consequences, even if its sole aim is chasing some superficially simple target such as net rate of calorie intake. Herring gulls exhibit preferences over which eggs to brood, and are sensitive to dimensions including size, shape, position, colour, visual texture (Baerends & Kruijt 1971, McFarland & Sibly 1975). Second, having a set of modes means having the problem of when to be in which. The question of whether or not to switch mode — say from foraging to dealing with a threat — is itself an economic problem, involving costs and benefits. Preferences could have a role mediating *between* modes. Finally, fixed motivational hierarchies have the property that they 'leave money on the table'. That is, at least if 'fixed' is interpreted to mean exhibiting — which may not involve representing — a lexical preference.[17] If any amount of some good, no matter how small, is worth more to you than any amount of another, no matter how large, you're unable to make trade-offs, or secure available marginal returns. Lexical ordering might make some hard choices more tractable, but it carries a cost.

---

[17] In Rawls (1971) a lexical ordering is a strict ranking, where one criterion (in his case a principle of justice) has to be completely satisfied before another can be applied at all. (

I've been saying that preferences represent *values*. The right way to think of this, is that they represent a kind of subjective utility, which is to say a value that is abstract compared to any of the specific modalities of both cost and return which matter to the agent. It might seem as though there's one more way of resisting my position, or diluting it, which is to say that something less abstract could serve as the 'common currency' here, for example calories. This line of thinking might be encouraged by the influence of models of foraging behaviour, in which net rate of calorie intake feature prominently. But the point doesn't generalise. The reason for this is that there isn't one concrete quantity that is generally the salient limiting resource that determines what allocation of means is efficient. Sometimes it is calories, sometimes it is hydration or some other specific nutrient, sometimes it might be time, or something else entirely. The limiting resource when people who aren't poor buy groceries is more likely to be time than money. (See McFarland & Bösser 1993: pp48-49) More generally, the opportunity cost, composed of the varied costs and returns of actions foregone, provides the right standard of comparison. So to do their work well, preferences have to be 'in' a quantity more abstract than any of the dimensions in which choices can vary, some of which don't feature in all options, i.e. utility.

The upshot of the preceding discussion, I argue, is that no decisive objection to the ERP has been found, although various good reasons to qualify or refine it have. The following formulation is the result:

(ERP2a) Preferences enable efficient action selection for degrees of freedom with alternative uses in non-transparent environments.

(ERP2b) Preferences will be selected when the gains outweigh the costs, and design options are accessible to selection.

The stipulation about degrees of freedom with alternative uses exempts truly single-use capacities the control of which can be entirely independent of any other control system (the biological equivalent of an air-bag). The restriction to non-transparent environments (in the sense used in Sterelny 2003) is to allow that in an environment where behaviour could be efficient while being entirely cue-bound, perhaps preferences would have no work to do (especially if the demands of internal environments are set aside).

Thus far I've not discussed preferences in relation to learning. This is deliberate. The benefits of preferences — in efficient allocation of behavioural capacities — could in principle be realised in an agent incapable of learning. This would be analogous to a reinforcement learning system tuned by 'evolutionary methods' (Sutton & Barto 1998). Although such a system might solve a 'reinforcement learning problem' the individual system itself does no learning, but its properties are found by comparing the success levels of many variant systems. Such a system would still, in the language of reinforcement learning, have a 'policy' that matched states to actions, and ranked multiple available actions, which is close enough to having preferences. I don't know if any *real* biological agents have this property - of having some kind of preference implementation while themselves lacking the capacity to learn. On

the other hand, at least some important kinds of learning, in particular operant or reward-based learning, are impossible *without* preferences. Reward-based learning involves tuning the agent's tendencies to perform various actions on the basis of the history of the consequences of performing them in various conditions. This is very explicit in the computational literature: the 'policy' is associated with expected rewards, and the behaviours associated with actual rewards, which are compared in various ways, so that — when there is a mis-match — the policy can be updated. And the rewards are quantitative signals. To the extent that living systems approximate this kind of functionality, signals encoding expected rewards, and actual rewards, and capable of comparing them (preferences) must be implemented somehow.

The ERP as defended here seems to imply that selection will drive some organisms (where the gains outweigh the cost, and where the design of the agent includes multi-use degrees of freedom, etc.) towards consistently tracking the costs and returns of actions and environmental states in a highly efficient way. But we know very well that living agents mostly *don't* do this, and behavioural economics and behavioural psychology have produced a long catalogue of systematic deviations from strict economic efficiency, as well as a few candidate general theories under which most of the deviations turn out to be special cases. (e.g. Kahneman & Tvesky 1979, Ainslie 2017). Whatever general theory, if any, is correct, a part of the explanatory story about the deviations will likely concern design and computational constraints. In the following section I briefly outline a few that seem especially important, and that can be expected to result in real preference implementations being compromised in various ways.

## 5. Design Constraints and Trade-offs

One response to the ways in which living agents apparently deviate from the norms of efficient agency (sometimes called 'rationality') is to argue either that there's something wrong with the data, perhaps because it was acquired in contexts that were ecologically bizarre, or that we're mistaken about what the rights norms of efficient agency are. Okasha, for example, critically discusses several ingenious responses to the phenomenon of intransitive choice, involving scenarios in which it is arguably optimal to violate transitivity including by exhibiting inter-temporal inconsistency (Okasha 2018, pp185f). Okasha's approach focuses on comparing the predictions and consequences of theories of fitness optimising behaviour and rational choice, largely independent of considerations of cognitive implementation. That isn't the approach I take here. I've defended a prediction about cognitive design innovation in the history of living organisms. Consequently specifically *cognitive* design factors, insofar as they bear on the implementation of preferences, are directly relevant. I emphasise three broad groups of factors.

The first is *efficient coding*. As noted in passing above in connection with the responses of salt receptors (§3), psychophysical processes aren't consistently perfect. More generally the neural encoding of physical facts doesn't 'attempt' to represent objective magnitudes, but rather compresses transduced variation into a baseline-dependent encoding, where the baseline itself is variable (Barlow, 1961). Even before the encoding, the processes of transduction themselves abandon information about objective magnitudes. Retinal cells, for example, respond relative to a transient average intensity (Burns & Baylor 2001). Measures like this are a pretty good way of making effective use of communication channels with fixed and often low bandwidth, and are found in psychophysical processes quite generally. If the most raw data about magnitudes in the world has this property, then — barring the action of some specialised system of correction and recalibration[18] — any represented magnitude that combines, integrates or compares those will inherit that property, and perhaps repeat the 'efficient coding' leading to further deviations from tracking of objective magnitudes. There is no reason to suppose that representations of value would be immune to this process, and some evidence that it exhibits it (e.g. Tremblay & Schultz 1999; Tobler, Fiorillo & Schultz 2005). Not only that, if we suppose that they exhibit it, then we find plausible explanations for several well-documented deviations from norms of efficiency. Among these are the specific inter-temporal inconsistency observed in many living organisms, sometimes referred to as their 'hyperbolic discounting' of delayed rewards (e.g. Gibbon 1977). If delays are 'efficiently' encoded in the psychophysical sense, then the ranking of options at different delays may fail to remain stable in the absence of new information. Another is the effect of irrelevant alternatives, since efficient coding including additional low-value options reduces the represented difference between others (Glimcher 2011, p242f).

The second is also related to efficiency, and concerns *chunking and hierarchies*. Bandwidth and processing constraints often favour 'chunking' and simplifying heuristics including decisions by 'trumping'. Chunking refers to grouping several things together to reduce the number of options that have to be dealt with (Miller 1956). In a decision case this might involve selecting first between, say, foraging and resting, and only after that selecting between different foraging actions. Switching between categories needn't involve comparison of any, let alone all, of the sub-elements of each. So, for example, switching from foraging to flight (because of predation risk) might be handled by a process that simply invokes flight behaviour when the detected

---

[18] Trained experts in highly scaffolded environments (including systems of financial mathematics and computers) can, for example, far more closely approximate economic rationality than lone living agents. (Conversely some environment, like shopping malls, are densely scaffolded to exploit unwary agents.)

risk passes some threshold, and remaining in flight mode until some later discrimination of sufficient safety. As argued above, these transitions needn't be devoid of sensitivity to economic factors: very hungry animals often tolerate foraging risks that more sated individuals would avoid. But this sensitivity doesn't establish that a single system of preferences is being applied to every alternative behaviour that is is in play. A patchwork of incomplete and variously chunked comparisons might be the most that time and specifically cognitive capacities allow. In addition, some behaviours might be sufficiently urgent or important that it makes evolutionary sense to make their execution under certain conditions effectively unconditional, or reflex. The basic initial actions required of dependant infants to secure nutrition might fall into this category. Doing anything else would be a terrible way of getting the export-explore trade-off wrong! If they are, then those allocations may bypass evaluation against available alternatives entirely.

Third, and finally, is *distributed and embodied cognition*. In (§4) I had reason to reject some of the more radical claims made on behalf of embodied or situated cognition. But this rejection was quite narrow in scope, limited to defending some kinds of bottlenecks, and one specific kind of representation. Not only that, the argument in favour of preferences made above is *prima facie* consistent with the idea that the implementation or representation of preferences — like all cognition — is a distinct stage in between sensation and action. That is to say, the argument for preferences might seem encourage a centralised and disembodied conception of cognition. I count myself among those who are pretty sure that this is *not* the right general picture of natural cognition, or a particularly helpful model for much artificial cognition. Real brains exhibit considerable parallelism in their architecture, fail strictly to segregate sensory from motor processes, and operate with 'incomplete' encodings that take for granted — and cannot function without — aspects of the structure of the bodies they control or sense with, or the environments in which they live. In addition many processes operate in quick and dirty ways, often not engaging in much communication with other systems. This picture is not at all to congenial to the idea of a single, central register of preferences against which available courses of action and their possible consequences (whether or not 'chunked' in line with the preceding paragraph) are judged. It is more plausible that preference implementation is widely distributed, with cues and candidate actions having preference relevant processing early, and with evaluation distributed and duplicated in various ways.[19] One neurally specific model of

---

[19] Attentional processes, for example, would be quicker if they had their own local, perhaps abbreviated or simplified, 'copies' of the preferences instead of having to consult head office. Redundancy brings risks of

this is Cisek's 'Affordance Competition Hypothesis' (Cisek 2007). Cisek rejects 'classical sandwich' models of cognition (the term is due to Hurley 1998), proposing that "the processes of action selection and specification occur simultaneously" (2007, p1586) and argues that sensory information selectively informs the generation of multiple incompletely specified behaviours, which may be released into execution prior to full specification. Reading his account of parallel, incomplete, competing behaviour specification processes with Dennett in mind, it is tempting to call it a 'multiple drafts' model of affordance competition: "From this perspective, behaviour is viewed as a constant competition between internal representations of the potential actions which Gibson (1979) termed 'affordances'" (Cisek 2007, p1586). My point here doesn't depend on endorsing Cisek's specific, though promising and interesting, proposal. What matters is that preference functionality does not depend on clear segregation of input and output, or serial ordering of decision processes. It is more likely that preference implementation is pervasive and interleaved in all processes from sensory transduction to activation of degrees of freedom

Together these suggest that in natural agents the implementation of preferences, when found, will be objectively inaccurate because of efficient coding, sometimes subject to chunking and other simplifications, sometimes bypassed entirely, and distributed around the brain in ways that introduce additional scope for volatility and even conflict. Because of the cognitive costs of implementing preferences, they may be found in ways that are limited in scope - for example an omnivorous creature with an otherwise simple lifestyle might have relatively fine-grained ranking of food types, accompanied by less flexible and cost-sensitive systems selecting between foraging and other behaviours. What real preferences have to do is represent, accurately enough, the returns on available behaviours so that (often enough) better one can be selected than would be otherwise. Like any cognitive system, they're expensive, and additional increments of accuracy drive up the expense. This is, furthermore, what we find in nature, as I argue below.

## 6. Preferences and Real Brains

Although including occasional real biological illustrations, my argument so far has been largely theoretical. One way of seeing it is as a prediction: given certain factors, evolution can

---

conflict, though, such as cases where the abstinent addict still orients towards predictors of the target of their addiction.

be expected to invent a preference-like capacity. A fair question is whether or not this prediction is correct. Here I briefly defend the position that, although the available evidence contains substantial gaps, it is correct: Many natural organisms have preferences.

This evidence is perhaps most comprehensive in the case of humans, but I'll begin with monkeys. Final common paths provide a natural architectural or functional 'place' for the operation of preferences. Circuits specialised for filtering out mechanically incompatible actions might allow competition between options that are not ruled out to be expressed. And if preferences are to influence allocation, they are the last place for them to do so. Some key early experiments in neuroeconomics depended on this very line of thinking to generate empirical hypotheses or to interpret their results (e.g. Platt & Glimcher 1999, Dorris & Glimcher 2004). As discussed in (§4), Monkey subjects were trained to express choices through saccades to different targets. Saccades are a convenient behaviour precisely because they're controlled by small networks of muscles with their own series of neural topographic maps, corresponding to a two-dimensional 'dart board' of skull-relative fixation targets. These maps are, in part, a final common path that prevents the eyes from attempting conflicted movements like simultaneously turning up and down. They're also a bottleneck for competition between fixation targets. In the interval preceding choice levels of activity in regions of the topographic maps corresponding to the saccadic targets under study varied with the expected return on that movement. Platt and Glimcher's monkeys were 'paid' in juice, while Dorris and Glimcher's subjects played an inspection game with returns in water, but in both cases local neural activity and (relative) expected subjective utility from the corresponding target were correlated. In later work Klein, Deaner and Platt (2008) found that activity in neurons specialised for saccades reflected values of both social *and* fluid rewards, so this isn't merely a fact about saccades for fluid rewards.

This isn't, furthermore, merely something about monkeys, or restricted to saccadic movements. In neuroeonomic study of human choices, correlates of utility are usually sought further 'upstream' of anatomically detailed final common paths, because electrode recordings are rarely used in human neuroeconomic experiments, and choices expressed by button presses and other hand movements don't correspond to somatotopic maps that are as conveniently tractable as those for eye movements (Glimcher 2011). Even so, the evidence is rather compelling. Levy and Glimcher (2012) survey relevant experimental work up to 2012. First they detail studies showing that activity in the ventral striatum was positively associated

with, among other things, monetary gains and losses, cumulative monetary rewards, anticipation of varying monetary rewards, expected values of uncertain monetary rewards, and discounted value of delayed monetary rewards. Second, they consider studies with at least one incentive other than money, including consumer goods, gustatory rewards (water, juice, food), physical pain, social reputation, again finding consistent correlations. Whether or not Shizgal and Conover were correct to say that orderly choice means that there "must" be value representations in a common scale, the work of some neuroeconomists suggests that there is *in fact* one for a wide range of choice types.

Comparative neuroeconomics is a small field, but approximately analogous results have been found in other vertebrates and some invertebrates with indications that the neural systems across phyla share structural similarities, and similar functional roles for dopamine or related molecules.[20] Capacity for reinforcement learning is diagnostic of having preferences. That is, although preferences need not be exploited in learning, reward learning requires both sensitivity to rewards and updating behavioural dispositions in light of reward-based consequences of earlier behaviour. Besides humans, monkeys, and other widely studied chordates (especially rodents and bird) reward sensitivity and capacity for reinforcement learning or operant conditioning has been found in insects (including bees, flies, cockroaches and ants), crustaceans (crabs, crayfish, lobsters) and cephalopod molluscs (Perry, Baron & Cheng 2013). In addition, dopamine or similar molecules play functionally similar roles in modulation of behaviour towards reward (Barron, Søvik & Cornish 2010). The fact that in some species apparently incapable of reinforcement learning, such as nematodes, dopamine plays a role modulating motor activity (Barron, Søvik & Cornish 2010) encourages the thought that the earliest steps towards implementing preferences were elaborations of motor control systems. That is, that final common paths paying their way in helping with the problem of producing behaviour out of capacities for activity provided the initial platform for the development of processes prioritising between available behaviours.

# 6. Conclusion

I have argued that psychologically real preference states can play a valuable role in prioritising the allocation of the capacities of living agents, and that this plausibly explains their evolution. The states I've focused on, preferences, are more cognitively modest than

---

[20] [I'm working on selecting appropriate additional references for this paragraph.]

desires on most accounts of desires. The argument I've offered is a version of the Environmental Complexity Thesis. Whereas standard expositions of that tend to focus on the development of world-representing states, whether understood as decoupled representations, or proto-beliefs, I've focused primarily on the role of value representations in dealing with the problem of matching the capacities of an agent to the changing contingencies of the world around it. That problem is often complex, involving mappings from varied capacities to action, given changing needs and changing external contingencies, with costs and benefits of multiple types and magnitudes. Doing *better* at dealing with this problem means deploying the available behavioural repertoire more effectively. This is an obvious kind of advantage, especially if the goals more efficiently achieved are understood in terms of fitness. In an unqualified form, the argument suggests something that we don't find. Real organisms aren't *that* efficient, and their deviations from efficiency exhibit some patterns. Many of these are the result of cognitive design limitations and trade-offs. Real preference implementations are compromises. Nonetheless, they are widely - though not universally - found in real organisms with nervous systems, and appear to have deep evolutionary roots.

# References

Ainslie, G. (2017) De Gustibus Disputare: Hyperbolic delay discounting integrates five approaches to impulsive choice. *Journal of Economic Methodology*, 24:2, 166-189.

Barlow, H. B. (1961) The coding of sensory messages. In: Thorpe and Zangwill (eds.), *Current Problems in Animal Behaviour*. New York: Cambridge University Press, pp. 330-360.

Barron, A.B., Søvik, E., & Cornish, J.L. (2010) The roles of dopamine and related compounds in reward-seeking behaviour across animal phyla. *Frontiers in Behavioural Neuroscience*, 4(163).

Brooks, R.A. (1991) Intelligence without representation, *Artificial Intelligence*, 47: 139–159.

Burns, M.E., and Baylor, D.A. (2001) Activation, deactivation, and adaptation in vertebrate photoreceptor cells. *Annual Review of Neuroscience*, 24: 779-805.

Burt A, Trivers R (2006) *Genes in Conflict: The biology of selfish genetic elements* Bellknap/Harvard University Press, Cambridge, Mass.

Clark, A. (1997) *Being There*, Cambridge, Mass.: MIT Press.

Davis, W. J., Mpitsos, G. J., Pinneo, J. M., & Ram, J. L. (1977) Modification of the behavioral hierarchy of Pleurobranchaea. I. Satiation and feeding motivation, *Journal of Computational Physiology*. 117: 99-125.

Dennett, D.C. (1991), *Consciousness Explained*, Boston, Little, Brown.

Dennett, D.C. (1996), *Kinds of minds,* New York, Basic Books.

Dennett, D.C. (2017), *From Bacteria to Bach and Back,* New York, Norton.

Dennett, D.C. (2018), 'Reflections on Peter Godfrey-Smith', in Huebner, B. (ed) *The Philosophy of Daniel Dennett,* Oxford, Oxford University Press, pp. 250-253.

Gagliano, M., Vyazovskiy, V.V., Borbély, A.A., Gromonprez, M. & Depczynski, M. (2016) Learning by Association in Plants. *Scientific Reports*, 6, 38427.

Gibbon, J. (1977) Scalar expectancy theory and Weber's law in animal timing. *Psychological Review,* 84: 279–325.

Glimcher, P. (2011) *Foundations of Neuroeconomic Analysis*, Oxford: Oxford University Press.

Godfrey-Smith, P. 1996. *Complexity and the Function of Mind in Nature*, Cambridge: Cambridge University Press.

Godfrey-Smith, P. 2002. Environmental Complexity and the Evolution of Cognition. In R. Sternberg and J. Kaufman (eds.) *The Evolution of Intelligence*. Mahwah: Lawrence Erlbaum, pp. 233-249.

Godfrey-Smith, P. (2018), 'Towers and Trees in Cognitive Evolution', in Huebner, B. (ed) *The Philosophy of Daniel Dennett,* Oxford, Oxford University Press, pp. 225-249.

Haig D (2002) *Genomic Imprinting and Kinship.* Rutgers University Press, New Brunswick.

Hurley, S. (1998) *Consciousness in Action.* London: Harvard University Press.

J_ékely, G., Keijzer, F., Godfrey-Smith, P. (2015) An option space for early neural evolution. *Philosophical Transactions of the Royal Society B,* 370(1684), 1-10.

Kahneman, D., and Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2): 263-291.

Keijzer, F. (2015) Moving and Sensing Without Input or Output: Early Nervous Systems and the Origins of the Animal Sensorimotor Organization. *Biology & Philosophy,* 30, pp311-331.

Keijzer, F., van Duijn, M. & Lyon, P. (2013) What nervous systems do: early evolution, input-output, and the skin brain thesis. *Adaptive Behavior*, 21, 67-85. (doi:10.1177/1059712312465330)

Klein, J.T., Deaner, R.O., & Platt, M.L. (2008), 'Neural correlates of social target value in macaque parietal cortex' *Current Biology*, vol. 18, pp. 419-424.

Levy, D.J. & Glimcher, P.W. (2012), 'The root of all value: a neural common currency for choice', *Current Opinion in Neurobiology*, vol. 22, pp. 1027-1038.

McFarland, D.J. and Sibly, R.M. (1975) The behavioural final common path, *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 270(907), 265-293.

McFarland, D. and Bösser, T. (1993) *Intelligent Behaviour in Animals and Robots*. Cambridge, Mass: Bradford Books.

McNamara, J.M. and Houston, A.I. (1986) The Common Currency for Behavioral Decisions, *The American Naturalist*, 127(3), pp358-378.

Millikan, R. (1984) *Language, Thought and Other Biological Categories*. Cambridge, Mass: MIT Press.

Millikan, R. (1993) What is Behavior? A Philosophical Essay on Ethology and Individualism in Psychology, Part 1. In *White Queen Psychology and Other Essays for Alice*. Cambridge, Mass: Bradford Books. (135-150)

Millikan, R.G. (2017), *Beyond Concepts*, Oxford, Oxford University Press.

Okasha, S. (2018) *Agents and Goals in Evolution*. Oxford: Oxford University Press.

Payne, R.B. (2005) *The Cuckoos*, Oxford University Press.

Perry, C.J., Barron, A. & Cheng, K. (2013), 'Invertebrate learning and cognition: relating phenomena to neural substrate', *Wiley Interdisciplinary Reviews: Cognitive Science*, vol. 4, pp. 5610582.

Platt, M.L., & Glimcher, P.W., (1999), 'Neural correlates of decision variables in parietal cortex', *Nature*, vol. 400, pp. 233-238.

Rawls, J. (1971) *A Theory of Justice*. Bellknap Press.

Samuelson, P. (1938) A Note on the Pure Theory of Consumer's Behaviour. *Economica*, 5(17), pp. 61-71.

Shadlen, M. N., and Newsome, W. T. ( 1996). Motion perception: Seeing and deciding. *Proceedings of the National Academy of Sciences USA*, 93: 628-633.

Shea, N. (2014) Reward Prediction Error Signals are Meta-Representational, *Noûs,* 48(2), pp. 314–341.

Shea, N. (2018) *Representation in Cognitive Science*. Oxford: Oxford University Press.

Spurrett, D. (2014), 'Philosophers Should Be Interested in 'Common Currency' Claims in the Cognitive and Behavioural Sciences'. *South African Journal of Philosophy,* vol. 33. pp. 211-221.

Spurrett, D. (2016) Does Intragenomic conflict predict Intrapersonal conflict? *Biology and Philosophy,* 31(3), pp313-333.

Spurrett, D. (2015) The Natural History of Desire, *South African Journal of Philosophy,* 34(3), pp304-313.

Sterelny, K. (2003). *Thought in a Hostile World*, Oxford: Blackwell.

Sutton, R.S. & Barto, A.G. (1998) Reinforcement Learning. Cambridge, Mass: MIT Press.

Tobler, P.N., Fiorillo, C.D., & Schultz, W, (2005) Adaptive coding of reward value by dopamine neurons. *Science*, 307:1642-1645.

Tremblay, L, and Schultz, W. (1999) Relative reward preference in primate orbitofrontal cortex. *Nature,* 398: 704-708.

Wouters, A. (1995) Viability Explanation. *Biology & Philosophy*, 10, pp435-457.

## Acknowledgments