

I Introduction

WHAT IS NEUROETHICS?

Neuroethics is a new field. The term itself is commonly, though erroneously, believed to have been coined by William Safire (2002), writing in *The New York Times*. In fact, as Safire himself acknowledges, the term predates his usage.¹ The very fact that it is so widely believed that the term dates from 2002 is itself significant: it indicates the recency not of the term itself, but of widespread concern with the kinds of issues it embraces. Before 2002 most people saw no need for any such field, but so rapid have been the advances in the sciences of mind since, and so pressing have the ethical issues surrounding them become, that we cannot any longer dispense with the term or the field it names.

Neuroethics has two main branches; the *ethics of neuroscience* and the *neuroscience of ethics* (Roskies 2002). The *ethics of neuroscience* refers to the branch of neuroethics that seeks to develop an ethical framework for regulating the conduct of neuroscientific enquiry and the application of neuroscientific knowledge to human beings; *the neuroscience of ethics* refers to the impact of neuroscientific knowledge upon our understanding of ethics itself.

One branch of the “ethics of neuroscience” concerns the conduct of neuroscience itself; research protocols for neuroscientists, the ethics of withholding incidental findings, and so on. In this book I shall have little to say about this set of questions, at least directly (though much of what I shall say about other issues has implications for the conduct of neuroscience). Instead, I shall focus on questions to do with the *application* of our growing knowledge about the mind and the brain to people. Neuroscience and allied fields give us an

2 INTRODUCTION

apparently unprecedented, and rapidly growing, power to intervene in the brains of subjects – to alter personality traits, to enhance cognitive capacities, to reinforce or to weaken memories, perhaps, one day, to insert beliefs. Are these applications of neuroscience ethical? Under what conditions? Do they threaten important elements of human agency, of our self-understanding? Will neuroscientists soon be able to “read” our minds? Chapters 2 through 5 will focus on these and closely related questions.

The *neuroscience of ethics* embraces our growing knowledge about the neural bases of moral agency. Neuroscience seems to promise to illuminate, and perhaps to threaten, central elements of this agency: our freedom of the will, our ability to know our own minds, perhaps the very substance of morality itself. Its findings provide us with an opportunity to reassess what it means to be a responsible human being, apparently making free choices from among alternatives. It casts light on our ability to control our desires and our actions, and upon how and why we lose control. It offers advertisers and governments possible ways to channel our behavior; it may also offer us ways to fight back against these forces.

If the neuroscience of ethics produces significant results, that is, if it alters our understanding of moral agency, then neuroethics is importantly different from other branches of applied ethics. Unlike, say, bioethics or business ethics, neuroethics reacts back upon itself. The neuroscience of ethics will help us to forge the very tools we shall need to make progress on the ethics of neuroscience. Neuroethics is therefore not just one more branch of applied ethics. It occupies a pivotal position, casting light upon human agency, freedom and choice, and upon rationality. It will help us to reflect on what we are, and offer us guidance as we attempt to shape a future in which we can flourish. We might not have needed the term before 2002; today the issues it embraces are rightly seen as central to our political, moral and social aspirations.

NEUROETHICS: SOME CASE STUDIES

Neuroethics is not only important; it is also fascinating. The kinds of cases that fall within its purview include some of the most controversial and strange ethical issues confronting us today. In this section, I shall briefly review two such cases.

Body integrity identity disorder

Body integrity identity disorder (BIID) is a controversial new psychiatric diagnosis, the principal symptom of which is a persisting desire to have some part of the body – usually a limb – removed (First 2005). A few sufferers have been able to convince surgeons to accede to their requests (Scott 2000). However, following press coverage of the operations and a public outcry, no reputable surgeon offers the operation today. In the absence of access to such surgery, sufferers quite often go to extreme lengths to have their desire satisfied. For instance, they deliberately injure the affected limb, using dry ice, tourniquets or even chainsaws. Their aim is to remove the limb, or to damage it so badly that surgeons have no choice but to remove it (Elliott 2000).

A variety of explanations of the desire for amputation of a limb have been offered by psychiatrists and psychologists. It has been suggested that the desire is the product of a *paraphilia* – a psychosexual disorder. On this interpretation, the desire is explained by the sexual excitement that sufferers (supposedly) feel at the prospect of becoming an amputee (Money *et al.* 1977). Another possibility is that the desire is the product of body dysmorphic disorder (Phillips 1996), a disorder in which sufferers irrationally perceive a part of their body as ugly or diseased. The limited evidence available today, however, suggests that the desire has a quite different aetiology. BIID stems from a mismatch between the agent's body and their *experience* of their body, what we might call their *subjective body* (Bayne and Levy 2005). In this interpretation, BIID is analogous to what is now known as gender identity disorder, the disorder in which sufferers feel as though they have been born into a body of the wrong gender.

4 INTRODUCTION

Whichever interpretation of the aetiology of the disorder is correct, however, BIID falls within the purview of neuroethics. BIID is a neuroethical issue because it raises *ethical* questions, and because answering those questions requires us to engage with the sciences of the mind. The major ethical issue raised by BIID focuses on the question of the permissibility of amputation as a means of treating the disorder. Now, while this question cannot be answered by the sciences of the mind alone, we cannot hope to assess it adequately unless we understand the disorder, and understanding it properly requires us to engage in the relevant sciences. Neuroscience, psychiatry and psychology all have their part to play in helping us to assess the ethical question. It might be, for instance, that BIID can be illuminated by neuroscientific work on phantom limbs. The experience of a phantom limb appears to be a near mirror image of BIID; whereas in the latter, subjects experience a desire for removal of a limb that is functioning normally, the experience of a phantom limb is the experience of the continued presence of a limb that has been amputated (or, occasionally, that is congenitally absent).

The experience of the phantom limb suggests that the experience of our bodies is mediated by a neural representation of a body schema, a schema that is modifiable by experience, but which resists modification (Ramachandran and Hirstein 1998). Phantom limbs are sometimes experienced as the site of excruciating pain; unfortunately, this pain is often resistant to all treatments. If BIID is explained by a similar mismatch between an unconscious body schema and the objective body, then there is every chance that it too will prove very resistant to treatment. If that's the case, then the *prima facie* case for the permissibility of surgery is quite strong: if BIID sufferers experience significant distress, and if the only way to relieve that distress is by way of surgery, the surgery is permissible (Bayne and Levy 2005).

On the other hand, if BIID has an origin that is very dissimilar to the origin of the phantom limb phenomenon, treatments less radical than surgery might be preferable. Surgery is a drastic course of

action: it is irreversible, and it leaves the patient disabled. If BIID can be effectively treated by psychological means – psychotherapy, medication or a combination of the two – then surgery is impermissible. If BIID arises from a mismatch between cortical representations of the body and the objective body, then – at least given the present state of neuroscientific knowledge – there is little hope that psychological treatments will be successful. But if BIID has its origin in something we can address psychologically – a fall in certain types of neurotransmitters, in anxiety or in depression, for instance – then we can hope to treat it with means much less dramatic than surgery. BIID is therefore at once a question for the sciences of the mind and for ethics; it is a neuroethical question.

Automatism

Sometimes agents perform a complex series of actions in a state closely resembling unconsciousness. They sleepwalk, for instance: arising from sleep without, apparently, fully awaking, they may dress and leave the house. Or they may enter a closely analogous state, not by first falling asleep, but by way of an epileptic fit, a blow on the head, or (very rarely) psychosis. Usually, the kinds of actions that agents perform in this state are routine or stereotyped. Someone who enters the state of automatism while playing the piano may continue playing if they know the piece well; similarly, someone who enters into it while driving home may continue following the familiar route, safely driving into their own drive and then simply sitting in the car until they come to themselves (Searle 1994).

Occasionally, however, an agent will engage in morally significant actions while in this state. Consider the case of Ken Parks (Broughton, *et al.* 1994). In 1987, Parks drove the twenty-three kilometres to the home of his parents-in-law, where he stabbed them both. He then drove to the police station, where he told police that he thought he had killed someone. Only then, apparently, did he notice that his hands had been badly injured. Parks was taken to hospital where the severed tendons in both his arms were repaired. He was

6 INTRODUCTION

charged with the murder of his mother-in-law, and the attempted murder of his father-in-law. Parks did not deny the offences, but claimed that he had been sleepwalking at the time, and that therefore he was not responsible for them.

Assessing Parks' responsibility for his actions is a complex and difficult question, a question which falls squarely within the purview of neuroethics. Answering it requires both sophisticated philosophical analysis and neuroscientific expertise. Philosophically, it requires that we analyze the notions of "responsibility" and "voluntariness." Under what conditions are ordinary agents responsible for their actions? What does it mean to act voluntarily? We might hope to answer both questions by highlighting the role of *conscious intentions* in action; that is, we might say that agents are responsible for their actions only when, prior to acting, they form a conscious intention of acting. However, this response seems very implausible, once we realize how rarely we form a conscious intention. Many of our actions, including some of our praise- and blame-worthy actions, are performed too quickly for us to deliberate beforehand: a child runs in front of our car and we slam on the brakes; someone insults us and we take a swing at them; we see the flames and run into the burning building, heedless of our safety. The lack of a conscious intention does not seem to differentiate between these, apparently responsible, actions, and Parks' behavior.

Perhaps, then, there is no genuine difference between Parks' behavior and ours in those circumstances; perhaps once we have sufficient awareness of our environment to be able to navigate it (as Parks did, in driving his car), we are acting responsibly. Against this hypothesis we have the evidence that Parks was a gentle man, who had always got on well with his parents-in-law. The fact that the crime was out of character and apparently motiveless counts against the hypothesis that it should be considered an ordinary action.

If we are to understand when and why normal agents are responsible for their actions, we need to engage with the relevant sciences of the mind. These sciences supply us with essential data for

consideration: data about the range of normal cases, and about various pathologies of agency. Investigating the mind of the acting subject teaches us important lessons. We learn, first, that our conscious access to our reasons for actions can be patchy and unreliable (Wegner 2002): ordinary subjects sometimes fail to recognize their own reasons for action, or even *that* they are acting. We learn how little *conscious* control we have over many, probably the majority, of our actions (Bargh and Chartrand 1999). But we also learn how these actions can nevertheless be intelligent and rational responses to our environment, responses that reflect our values (Dijksterhuis *et al.* 2006). The mere lack of conscious deliberation, we learn, cannot differentiate responsible actions from non-responsible ones, because it does not mark the division between the voluntary and the non-voluntary.

On the other hand, the sciences of the mind also provide us with good evidence that *some* kinds of automatic actions fail to reflect our values. Some brain-damaged subjects can no longer inhibit their automatic responses to stimuli. They compulsively engage in *utilization behavior*, in which they respond automatically to objects in the environment around them (Lhermitte *et al.* 1986). Under some conditions, entirely normal subjects find themselves prey to stereotyped responses that fail to reflect their consciously endorsed values. Fervent feminists may find themselves behaving in ways that apparently reflect a higher valuation of men than of women, for instance (Dasgupta 2004). Lack of opportunity to bring one's behavior under the control of one's values *can* excuse. Outlining the precise circumstances under which this is the case is a problem for neuroethics: for philosophical reflection informed by the sciences of the mind.

Parks was eventually acquitted by the Supreme Court of Canada. I shall not attempt, here, to assess whether the court was right in its finding (we shall return to related questions in Chapter 7). My purpose, in outlining his case, and the case of the sufferer from BIID, is instead to give the reader some sense of how fascinating, and how strange, the neuroethical landscape is, and how significant its

8 INTRODUCTION

findings can be. Doing neuroethics seriously is difficult: it requires a serious engagement in the sciences of the mind and in several branches of philosophy (philosophy of mind, applied ethics, moral psychology and meta-ethics). But the rewards for the hard work are considerable. We can only understand ourselves, the endlessly fascinating, endlessly strange, world of the human being, by understanding the ways in which our minds function and how they become dysfunctional.

THE MIND AND THE BRAIN

This is a book about the mind, and about the implications for our ethical thought of the increasing number of practical applications stemming from our growing knowledge of how it works. To begin our exploration of these ethical questions, it is important to have some basic grasp of what the mind is and how it is realized by the brain. If we are to evaluate interventions into the mind, if we are to understand how our brains make us the kinds of creatures we are, with our values and our goals, then we need to understand what exactly we are talking about when we talk about the mind and the brain. Fortunately, for our purposes, we do not need a very detailed understanding of the way in which the brain works. We shall not be exploring the world of neurons, with their dendrites and axons, nor the neuroanatomy of the brain, with its division into hemispheres and cortices (except in passing, as and when it becomes relevant).

All of this is fascinating, and much of it is of philosophical, and sometimes even ethical, relevance. But it is more important, for our purposes, to get a grip on how minds are constituted at much a higher level of abstraction, in order to shake ourselves free of an ancient and persistent view of the mind, the view with which almost all of us begin when we think about the mind, and from which few of us ever manage entirely to free ourselves: dualism. Shaking ourselves free of the grip of dualism will allow us to begin to frame a more realistic image of the mind and the brain; moreover, this more realistic image, of the mind as composed of mechanisms, will itself prove to be

important when it comes time to turn to more narrowly neuroethical questions.

Dualism – or more precisely *substance* dualism (in order to distinguish it from the more respectable *property* dualism) – is the view that there are two kinds of basic and mutually irreducible substances in the universe. This is a very ancient view, one that is perhaps innate in the human mind (Bloom 2004). It is the view presupposed, or at least suggested by, all or, very nearly all, religious traditions; it was also the dominant view in philosophical thought for many centuries, at least as far back as the ancient Greeks. But it was given its most powerful and influential exposition by the seventeenth century philosopher René Descartes, as a result of which the view is often referred to as Cartesian dualism. According to Descartes, roughly, there are two fundamental kinds of substance: *matter*, out of which the entire physical world (including animals) is built, and *mind*. Human beings are composed of an amalgam of these two substances: mind (or soul) and matter.

It is fashionable, especially among cognitive scientists, to mock dualists, and to regard the view as motivated by nothing more than superstition. It is certainly true that dualism's attractions were partly due to the fact that it offered an explanation for the possibility of the immortality of the soul and therefore of resurrection and of eternal reward and punishment. If the soul is immaterial, then there is no reason to believe that it is damaged by the death and decay of the body; the soul is free, after death, to rejoin God and the heavenly hosts (themselves composed of nothing but soul-stuff). But dualism also had a more philosophical motivation. We can understand, to some extent at least, how mere matter could be cleverly arranged to create complex and apparently intelligent behavior in animals. Descartes himself used the analogy of clockwork mechanisms, which are capable of all sorts of useful and complex activities, but are built out of entirely mindless matter. Today, we are accustomed to getting responses magnitudes more complex from our machines, using electronics rather than clockwork. But even so, it remains difficult to

10 INTRODUCTION

see how mere matter could really *think*; be rational and intelligent, and not merely flexibly responsive. Equally, it remains difficult to see how matter could be *conscious*. How could a machine, no matter how complex or cleverly designed, be capable of experiencing the subtle taste of wine, the scent of roses or of garbage; how could there be something that it is *like* to be a creature built entirely out of matter? Dualism, with its postulation of a substance that is categorically different from mere matter, seems to hold out the hope of an answer.

Descartes thought that matter could never be conscious or rational, and it is easy to sympathize with him. Indeed, it is easy to agree with him (even today some philosophers embrace property dualism because, though they accept that matter could be intelligent, they argue that it could never be conscious). Matter is unconscious and irrational – or, better, *arational* – and there is no way to make it conscious or rational simply by arranging it in increasingly complex ways (or so it seems). It is therefore very tempting to think that since we are manifestly rational and conscious, we cannot be built out of matter alone. The part of us that thinks and experiences, Descartes thought, must be built from a different substance. Animals and plants, like rocks and water, are built entirely out of matter, but we humans each have a thinking part as well. It follows from this view that animals are incapable not only of thought, but also of experience; notoriously, this doctrine was sometimes invoked to justify vivisection of animals. If they cannot feel, then their cries of pain must be merely mechanical responses to damage, rather than expressions of genuine suffering (Singer 1990).

It's easy to share Descartes' puzzlement as to how mere matter can think and experience. But the centuries since Descartes have witnessed a series of scientific advances that have made dualism increasingly incredible. First, the idea that there is a *categorical* distinction to be made between human beings and other animals no longer seems very plausible in light of the overwhelming evidence that we have all evolved from a common ancestor. Human beings have not always been around on planet Earth – indeed, we are a