# The Ethics of Creating Artificial Consciousness

John Basl

Northeastern University

## 1    Introduction

The purpose of this essay is to raise the prospect that engaging in artificial consciousness research, research that aims to create artifactual entities with conscious states of certain kinds, might be unethical on grounds that it wrongs or will very likely wrong the subjects of such research. I say *might be unethical* because, in the end, it will depend on how those entities are created and how they are likely to be treated. This essay is meant to be a starting point in thinking about the ethics of artificial consciousness research ethics, not, by any means, the final word on such matters.

While the ethics of the creation and proliferation of artificial intelligences and artificial consciousnesses (see, for example, (Chalmers 2010) has often been explored both in academic settings and in popular media and literature, those discussions tend to focus on the consequences for humans or, at most, the potential rights of machines that are very much like us. However, the subjects of artificial consciousness research, at least those subjects that end up being conscious in particular ways, are research subjects in the way that sentient non-human animals or human subjects are research subjects and so should be afforded appropriate protections. Therefore, it is important to ask not only whether artificial consciousnesses that are integrated into our society should be afforded moral and legal protections and whether they are a risk to our safety or existence, but whether the predecessors to such consciousnesses are wronged in their creation or in the research involving them.

In section 2, I discuss what it means for a being to have moral status and make the case that artificial consciousnesses of various kinds will have moral status if they come to exist. I then take up the issue of whether it is thereby wrong to create such entities (section 3). It might seem obvious that the answer is "no", or at least it is no more impermissible than the creation and use of non-human research subjects. However, I argue that there should be a presumption against the creation of artificial consciousnesses.

## 2    Moral Status and Artificial Consciousness

In order to determine whether it is possible to wrong artificial consciousnesses by creating them or conducting research on them, we must first determine whether such entities have moral status and what the nature of that status is.

## 2.1   What is moral status?

The term 'moral status' is used in various ways in the ethics and applied ethics literature. Other terms, such as 'inherent worth', 'inherent value', 'moral considerability' etc., are sometimes used as synonyms and sometimes to pick out species of moral status.[1] In the broadest sense of the term, to have moral status is just to have any kind of moral significance; that is, having moral status means that in at least some contexts moral agents must be responsive to or regard the thing that has moral status.

It would be every easy to argue that artificial consciousnesses have moral status in the broad sense just described sense. After all, even a rock, if owned by someone or part of a piece of art, for example, has moral status in this sense. Instead, I will employ the term 'moral patient' to pick out a particular form of moral status. The definition of 'moral patient' as used in this paper is:

> **Moral Patient**df: X is a moral patient iff agent's like us are required to take X's interests into account in our moral deliberations for X's sake when X's interests are at stake.

This definition has the following features:

1. A being is a moral patient only if it has *interests* that are to be taken into account in moral deliberations.
2. A being's being a moral patient entitles it have its interests taken into account in moral deliberations for *its own sake*.
3. Moral patiency is a property had by an entity *relative to agents like us*.

Each of these features will be discussed in detail below, but first, it is important to discuss the relationship between moral patiency and normative theory. Some view the question of whether a being is a moral patient as dependent on which normative theory is true.[2] That is, in order to determine which beings are patients, we must first figure out whether we should be, for example, Utilitarians or Kantians, Virtue Theorists or Contractualists. If this thesis, call it the Dependency Thesis, about the relationship between moral status and normative theories is correct, we can't answer the question of whether artificial consciousnesses are moral patients without first answering the question of which normative theory is correct.

---

[1] See for example, (O'Neill 2003; Cahen 2002; Sandler and Simons 2012). For dissent on the usefulness of moral status talk see (Sachs 2011)

[2] Buchanan (2011, chap. 7), for example, discusses the differences between moral status on a Contractualist framework and moral status on a Utilitarian framework. See also (Sober 1986).

There are important relationships between normative theory and moral status. For one thing, which normative theory is true explains the nature or source of the moral status of whichever beings have it. If contractualism is true, for example, a being's moral status is grounded in or finds its source in the consent of rational contractors; if utilitarianism is true, a being's moral status is grounded in the fact that it's being benefitted or harmed contributes to or detracts from the value of a state of affairs. Furthermore, how, in particular, moral patients are to be treated is a function of which normative theory is ultimately correct. Utilitarianism more easily licenses the killing of moral patients more easily than a Kantian ethic, for example. For this reason, the strength of the presumption against creating artificial consciousnesses defended below will depend on which normative theory is true. However, the Dependency Thesis concerns relationship between normative theory and moral patiency with respect to *which beings are moral patients*.[3]

Fortunately, the version of the Dependency Thesis that precludes us from determining whether artificial consciousnesses are moral patients independently of determining which normative theory is true is false. One point in favor of thinking that it is false is that we know that all adult humans of sound mind are moral patients, and yet we aren't sure which normative theory is true, or, at least, whether all adult humans of sound mind are moral patients is far less controversial than which normative theory is true.

One might argue that the obviousness of our patiency just serves as a condition of adequacy on normative theories and that's why we know we are patients even if we haven't settled which normative theory is true. However, it also suggests the possibility that we can make a similar case for the moral status of other beings. That is, even if some metaphysical, ontological, or supervenience version of the Dependency Thesis is true, we may have ways of specifying which things are moral patients independently of determining which normative theory is true. All that really matters for the purposes of arguing that artificial consciousnesses are or can be moral patients is that the dependency relationship between patiency and normative theory isn't epistemic, i.e. so long as we can come to know that some being is or isn't a moral patient without determining which normative theory is true.

There is good reason to think we can come to know who or is a moral patient independently. Debates about which entities have moral status and about the degree to which entities of various kinds matter happen, as it were, internal to normative theories. Utilitarians, for example, have argued about

---

[3] Another version of the Dependency Thesis might claim that the degree to which a being has moral status depends on normative theory. (Buchanan 2011) seems to suggest this as well. However, I think this version of Dependency is also false. There are ways to cash out differences in treatment owed to different kinds of beings without understanding them as having different degrees of moral status. In other words, 'degrees of moral status' can be gotten rid of without losing the ability to make the normative distinctions that talk is intended to capture. This translatability is not central to what I'll say here and so I leave it unargued for.

whether non-human animals and human infants are moral patients on part with us.[4] There are some Kantians that argue that many non-human animals should be accorded many rights in the same way that we ought.[5] So long as the intra-normative debates are coherent we can be sure, at least, that normative theories aren't fully determinate of which beings have moral status.

Furthermore, the kinds of arguments made that this or that entity is a moral patient do not typically appeal to which normative theory is true.[6] Consider, for example, a standard argument from marginal cases that non-human animals have moral status. Such arguments take for granted that so-called "marginal cases", such as infants and the severely mentally handicapped, have moral status. Then an argument is made that there is no morally relevant difference between marginal cases and certain non-human animals, for example chimps. From this it is concluded that chimps are moral patients in the same way that we are. This argument doesn't make explicit mention of normative theory, nor do the arguments typically given for the premise that there is no morally relevant difference between chimps and marginal cases.

I'm not here endorsing any particular argument from marginal cases or assessing its merits. The point is that the kinds of arguments that a Utilitarian might use to convince another Utilitarian that chimps matter are the same kinds of reasons that should convince a Contractualist or Kantian to accept that chimps are moral patients. Similarly, if Kantians could make a case that, for example, only the interests of very cognitively advanced beings are relevant to moral deliberations, that advanced cognitive capacities are a morally

---

[4] Consider for example the difference between Singer's view about the moral status of humans and Frey's view of same. Both are committed Utilitarians and yet Singer (2002) things that all sentient beings are equal, that is have equal moral status (though Singer acknowledges that typically, a human's life should often be preferred over an animals in a conflict because humans can suffer and enjoy in more ways than most non-human animals) while Frey (1983; 1988) thinks that human adults of sound mind are distinct from non-human animals, that their lives are of more value because of their capacity for certain kinds of experiments. It is worth noting that both come to similar conclusions about the ethics of animal experimentation and the differences between their views are subtle, but the fact that Frey thinks humans have additional value in virtue of having a capacity or capacities that non-human animals do not, is sufficient to demonstrate the kind of inter-normative differences in conceptions of moral status that are relevant here.

[5] See, for example, (Regan 1983) who argues, using arguments very similar to those employed by Singer, to extend a Kantian conception of rights to non-human animals that are minimally conscious. See also (Rollin 2006) for a discussion that includes Contractualist discussions of animal moral status. For an excellent discussion of how a more traditional Kantian might approach the issue of animal rights see (Korsgaard 2004).

[6] Of course, which properties are taken to be morally significant are often influenced by which normative theory one takes to be true. A Kantian is more likely to think that "being an end in oneself" is a morally significant property than a Utilitarian. But, that is a sociological fact. The Kantian still owes the Utilitarian an argument as to why that property is morally significant. If the argument is sounds, the Utilitarian might agree that it is only the benefits and harms that accrue to ends in themselves that influence the value of states of affairs, just as many Utilitarians are keen to think that it is only the benefitting and harming of humans make a difference to the value of states of affairs.

relevant properties, they won't do so by appealing to the structure of Kantian normative theory, but to reasons that a Utilitarian could accept; at least they will do so if they hope to convince other Kantians that don't share their view about the relevance of advanced cognitive capacities.[7]

The above considerations provide an abbreviated, but I hope sufficient, case for the idea that we can identify moral patients without first discovering which normative theory is true.

### 2.1.1 Interests

According to the definition of moral patiency, if a being is a moral patient we must take that beings interests into account for the sake of that being. To say that a being has interests is to say that it has a welfare, that it can be benefitted or harmed.[8] Whether a being is potentially a moral patient depends, therefore, on whether it has a welfare, and that depends on which theory of welfare is true.

There are various families of views about welfare and some are more stringent about the features a being must have to have a welfare.[9] I don't intend here to settle the issue of which theory of welfare is true. Instead, below I will focus on a type of artificial consciousness that will have a welfare independently of which of many plausible theories of welfare is true.

A being's welfare can be significant in moral deliberations for a variety of reasons. For example, if I hire a dog walker, they have an obligation to me to take my dog's interests into account. However, they also, I contend, have an obligation to take my dog's welfare into account for her sake; even if I didn't own my dog, even if no one does, it would be wrong for the dog walker to kick my dog for no reason.[10]

Some being's welfare may only matter derivatively (see, for example, Feinberg (1963) on plants), but a moral patient's welfare matters for its own sake. The interests of a patient figure into our deliberations independently of their relationship to the welfare of others.[11]

---

[7] We could of course understand a normative theory to include facts about whom or what has moral status. I'm using normative theory, as is typical, to pick out a theory of right action (and, if you like, an account of the source of normativity).

[8] For a more detailed explanation see (Basl Forthcoming).

[9] For an overview of these families see (Griffin 1988; Streiffer and Basl 2011).

[10] I'm not here committing to the view that my dog's welfare matters for its own sake simply because she has a welfare. It might be that her welfare matters because she is an end in herself, or because reasonable would agree that an animal's welfare is morally significant. Again, I'm not committing to any particular normative theory or any particular source of normativity. Whichever theory is true, I explain below, my dog's welfare is relevant to moral deliberations for her own sake.

[11] This isn't to say how their welfare affects our own or others isn't also relevant to deliberations. In thinking about what to do, we must think about these conflicts of interests. That is consistent with thinking that a being's interests should be taken into account for the sake of the being under consideration.

### 2.1.2   The Agent Relativity of Moral Patiency

It might seem odd, even contradictory, to claim that a moral patient's welfare matter's in moral deliberations for their own sake while at the same time also relativizing moral patiency to a set of agents like us. However, rather than being contradictory this reflects the fact that agent's that are radically different from us might exist in an entirely different ethical world, so to speak.

Let's imagine, for example, that there is a type of being that is completely immaterial. Admittedly, I don't know how to understand how such beings interact in any sense, but I do know that whatever such beings do, they cannot have any effect on being's like us and so they are not required to take our welfare into account in whatever moral deliberations they have.

Or, assume that Lewis (2001) was right and that all possible worlds really exist in the normal everyday sense of exists. There are worlds very much like ours that are, in principle, causally cut off from us. The moral agent's in those possible worlds are under no obligation to take our welfare into account because they can't affect us in any way.

Finally, imagine that rocks have a welfare but that it is impossible for us to come to know about that welfare. In such a case, while we may make these beings worse off, we are either under no obligation to take their welfare into account, or if we are so required, we are excused for failing to do so because of our ignorance and so for all practical purposes rocks are not moral patients.[12]

These examples show, at least in principle, that whether a being is a moral patient is agent relative; it is relative to agent**s** sufficiently like us that engage in causal interactions with potential patients and which can come to know or have reasonable beliefs that their actions affect the welfare of potential patients.

## 2.2   Can artificial consciousnesses be moral patients?

There is not a single question of whether artificial consciousnesses could satisfy the conditions of moral patiency. There is a  technological version of the question: will we ever be in a technological position to create artificial consciousnesses that satisfy the conditions of patiency?

The answer to that question depends in part on an answer to a nomological version of the question: do the laws of our universe make it possible to create consciousness out of something other than the kind of matter of which we are composed and configured in a way that's very similar to consciousnesses we know of?

The technological and nomological questions just raised are interesting and important, especially to those who wish to create artificial consciousnesses. However, as a philosopher, I'm in no position to answer them. I'm going to

---

[12] For a discussion of the distinction between obligation and excuse see (McMahan 2009).

assume that artificial consciousnesses with a large range of cognitive capacities are creatable and instead focus on the following conceptual question: is it conceptually possible to create an artificial consciousness that is a moral patient?

I think the answer to this question is clearly "yes". To see why, just imagine that we've managed to create an artificial consciousness and embodied it, certainly a conceptual possibility. This being is, we know, mentally very much like us. It is a moral agent, it has a similar phenomenology, it goes about the world much like we do, etc.. What would we owe to this being? I think it is our moral equal and that denying that would make one, to use Singer's term, a *speciesist*. But, even if you think that such a being would not be our moral equal, it would certainly be wrong to hit such a thing in the face with a bat, or to cut off its arm because of the effect such actions would have on the welfare of such a being. That is, even if we have some special obligations to the members of our own species and some degree of partiality justified, this kind of artificial consciousness is a moral patient.

The more interesting question isn't whether an artificial consciousness very much like us is a moral patient, but what are the minimal conditions for an artificial consciousness to be a moral patient. After all, it seems plausible that merely being conscious does not make a thing a moral patient. Imagine that we can create an artificial being that has conscious experiences of color but nothing else (Basl Forthcoming). Such a being is not a moral patient because it doesn't have a welfare; which color experience it is having, by hypothesis, doesn't make its life go better or worse.

So, what are the minimal conditions for an artificial consciousness to have a welfare that we should care about for its own sake? That's a much harder question to answer. Fortunately, we can say something about the ethics of creating artificial consciousnesses without fully answering it by thinking about existing non-human moral patients.

It is, I hope, relatively uncontroversial that at least some non-human animals, mammals and birds in particular, are moral patients. It is at least pertinent to our deliberations about whether to experiment on such animals, whether it is permissible to withhold food from our pets, to encourage them to fight for our pleasure, that these activities will affect the welfare of these beings (and that welfare of such beings are pertinent for their own sake). It is more controversial whether such beings matter because they have the capacity for suffering and enjoyment, whether they have desires that can be satisfied or frustrated, whether they are sufficiently rational etc. That is, it is controversial in virtue of which properties they have a welfare, but relatively uncontroversial whether their welfare ought figure into our moral deliberation for the sake of those animals whose welfare is at stake.

For the purposes of evaluating the ethical case against the creation of artificial consciousnesses, we can restrict that evaluation to the creation of artificial consciousnesses with capacities similar to non-human animals that we

take to be moral patients. If it turns out that the conditions for a being a moral patient are less restrictive, the conclusions below will apply to artificial consciousnesses of that type.

# 3 The Case Against Creating Artificial Consciousness

Just as there's nothing intrinsically wrong with creating biological consciousnesses through traditional means (i.e. breeding animals or having children), there's nothing intrinsically wrong with creating an artificial consciousness, at least not from the perspective of the created being.[13] If scientists were to create an artificial consciousness, like that described above, that goes about the world as we do and has a mental life like ours, and those same scientists and society generally were to treat that being in a way that was commensurate with its being a moral patient, that would be morally permissible. However, there are reasons to expect that such a being would not be treated in a way that is commensurate with its being a moral patient and in such a case, we have good reason not to allow the creation of such a consciousness, at least not without adequate protections. The case against the creation of artificial consciousnesses is thus conditional: there is an ethical reason not to create artificial consciousnesses when there is sufficient risk that such beings will be moral patients and when there is also sufficient risk that these patients will be mistreated.

To illustrate the kinds of risks to artificial consciousnesses associated with their creation it is useful to first discuss some cases having to do with traditional organisms. Consider the following case:

> **Neural Chimera**: Researchers are attempting to create human-animal neural chimeras by injecting human stem cells into the brains of guppies. In light of some recent developments in stem-cell research and in neuroscience, the scientists think that they can significantly alter the cognitive capacities of the resultant guppies by doing so. Their hope is to create guppies with brains much more like ours that they can use in Parkinson's research. Since guppies are relatively cheap to feed and reproduce quickly, they think this would be an excellent solution to the need for better animal models for Parkinson's research.

Let's assume that researchers, after conducting the research described intend to care for the resulting guppies in the same way that they care for typical guppies. If so, such research is almost certainly unethical? If the scientists are right that there is a significant chance that the resulting guppies

---

[13] See (Shiffrin 1999) for a dissenting argument that bringing a child into existence is a pro tanto wrong.

will be, mentally, a lot more like us, they would be owed much more than what is typically accorded to guppies in normal research contexts. By creating a moral patient that is much like us, the scientists obligate themselves to treat these subjects commensurately with that moral status, but the research as described fails to do so.

This kind of case seems far-fetched, but scientists are concerned with the creation of such chimeras and these sorts of ethical worries have been raised by others about this research (Streiffer 2005). And, it is not the only sort of research that raises these worries. I've argued elsewhere that testing cognitive enhancements on non-human research subjects has the potential to alter their capacities in ways that increase the risks that they will be mistreated (Basl 2013).

The case above serves to illustrate what might make creating artificial consciousness unethical. That's not to say all artificial consciousness research will be unethical. In assessing the ethics of creating artificial consciousness research programs from the perspective of the research subjects (the artificial consciousnesses that might be created) the following questions must be addressed:

1. How probable is it that a given research program will result in the creation of an artificial consciousness that is a moral patient?
2. How probable is it that such patients will fail to be treated appropriately?

These questions are difficult to answer in any precise fashion. It will vary from research program to research program and it will depend on what safeguards are put in place to protect the interests of the research subjects. Still, there are some considerations that suggest that these probabilities are sufficiently high (or at least should be judged to be so).

## 3.1   The probability of creating artificial consciousness

How probable the creation of artificial consciousness is in a given research context is extremely difficult to determine. This is in part because it depends on the nature of consciousness. For example, if consciousness requires neural correlates, and those correlates aren't realizable using the methods or materials in use in some research program, the probability of creating an artificial consciousness is low.

The problem is that the nature of consciousness is a difficult problem and so we can't be sure which research programs are most promising with respect to creating artificial consciousness. In fact, it might be that the question of consciousness goes unsettled until researchers are able to create an artificial consciousness to confirm one or other of various theories.

Given these difficulties, what should we say about the probability of creating an artificial consciousness? Are we stuck thinking there is no way to

assign a probability one way or another and so need to concern ourselves with the ethical risk to research subjects? I think not. The reason is that every attempt to create artificial consciousness is taken with the aim of success and because of the ethical risk success carries.

Attempts to create artificial consciousness will not be made at random. Researchers will attempt methods they think promising or that have a better chance of success than alternatives. This doesn't tell us that the probability that an artificial consciousness of the sort that would be a moral patient will be created is high, but I think that it provides us a reason to assume that it is high if there are ethical risks associated with success. That is, since artificial consciousness researchers are engaging in a research project with an eye towards success and since, as I argue below, success carries with it certain ethical risks that would not arise if the research were not pursued, we should, perhaps artificially, assume that the risk of creating artificial consciousnesses is relatively high until we have reason to think otherwise.[14]

## 3.2   The probability of mistreatment

Whereas it is extremely difficult to predict how probable it is that a given research program will result in an artificial consciousness that is a moral patient, it is not so difficult to see why if a program were successful there would be a substantial chance that the created consciousness would be mistreated.

In arguing that research involving the use of cognitive enhancement technologies on non-human research subjects, I've raised the worry that for a variety of reasons, we might be more likely to mistreat research subjects than in traditional research contexts (Basl 2013). This is because cognitive enhancement research might alter research subjects in ways that aren't detectable and that can't be communicated by the research subjects themselves. Also, without further education about the concerns associated with cognitive enhancement, researchers might fail to take required precautions.

Similar worries arise with respect to artificial consciousnesses. Artificial consciousness research, unlike research involving non-human research subjects, is not subject to oversight designed to protect research subjects. Without oversight and researcher education, researchers are less likely to take the welfare of research subjects into account.

Furthermore, depending on which methods of creating artificial consciousness are successful, researchers may be in a poor epistemic situation with respect to determining whether they've created a moral patient at all. To see why, consider the following highly stylized case:

---

[14] Some projects that might be classified as "artificial consciousness projects" but are thought to involve only preliminary research or are being done as development steps should be excluded from the scope of this assumption.

**Selection**: Artificial consciousness researchers, informed by evolutionary biologists, have devised a series of problems that they think will encourage the evolution of consciousness. Programs are written that mutate with imperfect replication and reproduce proportionally to the efficacy with which they solve the various problems.

Let's say that the research program described in Selection is very likely to lead to programs that are conscious in ways that make those programs moral patients. It doesn't thereby follow that researchers will immediately know when such beings evolve or how to promote or avoid frustrating the interests of the created beings. Just because chimps and dogs are both moral patients, it doesn't thereby follow that treating them appropriately means treating them similarly. The same goes for artificial consciousnesses. Even if we become fairly confident that we've created an artificial consciousness, we can't be sure we know what is thereby required of us.

Of course, if it is impossible for researchers to determine what the interests of such research subjects are, they may be excused for any and all treatment that isn't commensurate with the moral patiency of these beings. But, researchers are required to try to make determinations about the interests of these beings and to try to treat them appropriately in the context of doing research. What would be problematic in the above case would be for the researchers to experiment on such intelligences without making attempts to determine what's good for them.

## 3.3 The case against creating artificial consciousnesses revisited

What does the case against creating artificial consciousnesses amount to? The above concerns provide a *pro tanto* case, or an overrideable presumption against artificial consciousness research which aims to create artificial consciousnesses that have capacities, such as sentience, self-awareness, or desires of one form or other, like mammals, birds, or humans. It is a presumption only, not an all things considered objection to such research. Furthermore, what is required to override the presumption depends, in part on which normative theory is true which will partially determine what is owed to moral patients.

What kinds of considerations will override the presumption against creating artificial consciousnesses? First, as research becomes more advanced, researchers might be able to determine that particular research program is valuable but that the risk of creating a moral patient is extremely low. In such case, the ethical risk to the research subject is low and the presumption is overridden.

Second, researchers might be engaged in research that is likely to result in the creation of moral patients, but are taking or intend to take sufficient care

to determine what constitutes appropriate treatment of the created research subjects. That is, they will not merely conduct their research, but also conduct research to determine what promotes or frustrates the interests of created beings. In doing so, the researchers lower the probability of mistreatment and the presumption is overridden.

Finally, it might be that artificial consciousness research is so valuable, that the cost of not doing it so costly, that it is worth doing no matter how poorly life goes for the entities created and where doing the research efficiently rules out taking the time to discern what's good or bad for research subjects. There may be an ethical imperative to engage in artificial consciousness research. And if so, and if doing so efficiently, requires that we ignore the interests of the created entities, then the presumption may be overridden. Just as with research involving non-human animals, the relevant question is whether the value of doing the research justifies the ethical costs accrued in harming research subjects. Sometimes, the answer is "yes".

However, whether the presumption is overridden for this reason is going to be extremely sensitive to the nature of our moral obligations, and, thereby, to which normative theory is ultimately true.[15] For a Utilitarian, it will be permissible to ignore the interests of artificial consciousnesses if doing so maximizes utility. The conditions under which a Kantian will agree to ignore or allow the interests of such patients to be overridden will be much more stringent.

# 4 Conclusion

The above is not meant to amount to a complete defense of the presumption against creating artificial consciousnesses, and in fact, it leaves open how strong the ethical presumption is. Instead, I hope that it is raised the ethical concerns associated with such research enough to start a conversation about the ethics of engaging in it. Ultimately, each artificial consciousness research program will have to be evaluated individually to assess the ethical risks, just as each research program involving non-human animals is evaluated. However, it is important that we recognize that there are ethical risks concerning the research subjects and not only those risks that accrue to us or that we face once we've realized artificial consciousness like our own.

Works Cited

Basl, John. Forthcoming. "Machines as Moral Patients We Shouldn't Care About (Yet)." *Philosophy & Technology*

---

[15] The same can be said of traditional research. While some Kantians, for example, endorse a complete prohibition on any animal research, most recognize that the rights of all beings are not, in fact, inviolable come what may. If circumstances are such that great harms can't be avoided, some individuals may be sacrificed for others (Regan 1983; Buchanan 2011).

———. 2013. "Sensitivity Enhancement: The Ethics of Engineering Non-Human Research Subjects." In *Designer Biology: The Ethics of Intensively Engineering Biological and Ecological Systems*, 219–232. Lexington Books.

Buchanan, Allen E. 2011. *Beyond Humanity?: The Ethics of Biomedical Enhancement*. Oxford University Press, USA.

Cahen, Harley. 2002. "Against the Moral Considerability of Ecosystems." In *Environmental Ethics: An Anthology*, edited by Andrew Light and H. Rolston III. Blackwell.

Chalmers, David. 2010. "The Singularity: A Philosophical Analysis." *Journal of Consciousness Studies* 17 (9-10): 7–65.

Feinberg, Joel. 1963. "The Rights of Animals and Future Generations." *Columbia Law Review* 63: 673.

Frey, R. G. 1983. "Vivisection, Morals and Medicine." *Journal of Medical Ethics* 9 (2): 94.

———. 1988. "Moral Standing, the Value of Lives, and Speciesism." *Between the Species: a Journal of Ethics* 4 (3): 191.

Griffin, James. 1988. *Well-Being: Its Meaning, Measurement, and Moral Importance*. Oxford University Press, USA.

Korsgaard, Christine M. 2004. "Fellow Creatures: Kantian Ethics and Our Duties to Animals." *The Tanner Lectures on Human Values* 25: 26.

Lewis, David K. 2001. *On the Plurality of Worlds*. Malden, Mass.: Blackwell Publishers.

McMahan, Jeff. 2009. *Killing in War*. 1st ed. OUP Oxford.

O'Neill, John. 2003. "The Varieties of Intrinsic Value." In *Environmental Ethics: An Anthology*, edited by Holmes III Rolston and Andrew Light.

Regan, Tom. 1983. *The Case for Animal Rights*. Berkeley: University of California Press.

Rollin, Bernard E. 2006. *Animal Rights & Human Morality*. Amherst, N.Y.: Prometheus Books.

Sachs, Benjamin. 2011. "The Status of Moral Status." *Pacific Philosophical Quarterly* 92 (1): 87–104.

Sandler, Ronald, and Luke Simons. 2012. "The Value of Artefactual Organisms." *Environmental Values* 21 (1): 43–61.

Shiffrin, S.V. 1999. "Wrongful Life, Procreative Responsibility, and the Significance of Harm." *Legal Theory* 5 (02): 117–148.

Singer, Peter. 2002. *Animal Liberation*. Vol. 1st Ecco pbk. New York: Ecco.

Sober, Elliott. 1986. "Philosophical Problems for Environmentalism." In *The Preservation of Species*, edited by Bryan Norton. Princeton University Press.

Streiffer, Robert. 2005. "At the Edge of Humanity." *Kennedy Institute of Ethics Journal* 15 (4): 347–370.

Streiffer, Robert, and John Basl. 2011. "Applications of Biotechnology to Animals in Agriculture." In *The Oxford Handbook of Animal Ethics*, edited by T Beauchamp and R Frey. Oxford.