

Abstract: A series of recent studies have shown that people's moral judgments can affect their intuitions as to whether or not a behavior was performed intentionally. Prior attempts to explain this effect can be divided into two broad families. Some researchers suggest that the effect is due to some peculiar feature of the concept of intentional action in particular, while others suggest that the effect is a reflection of a more general tendency whereby moral judgments exert a pervasive influence on folk psychology. The present paper argues in favor of the latter hypothesis by showing that the very same effect that has been observed for *intentionally* also arises for *deciding, in favor of, opposed to,* and *advocating*.

People ordinarily distinguish between behaviors that are performed 'intentionally' and those that are performed 'unintentionally.' At least to a first glance, this distinction seems to be a perfectly ordinary part of our usual approach to understanding the mind, right alongside the concepts of belief and desire. In other words, the concept of intentional action appears to be one aspect of *folk psychology*.

Yet recent experimental work has revealed a surprising fact about the way in which people ordinarily apply this concept. It seems that people's ordinary intuitions about intentional action can actually be affected by their *moral* judgments. In particular, there seem to be cases in which people's intuitions about whether a behavior was performed intentionally depend in some way on their moral appraisal of the behavior itself. What we have here, then, is a case in which people's moral judgments appear to be influencing their folk-psychological intuitions.

A question now arises as to whether this effect is telling us anything of general significance about the relationship between folk psychology and moral judgment. Is the effect just due to some quirk in the process by which people attribute intentional action, or is it a manifestation of some more general mechanism whereby moral judgments can

have an impact on folk psychology? Here, one finds a striking divergence of views – with researchers dividing off into two basic camps.

On one side are researchers who suggest that the effect can be understood entirely in terms of certain special features of the attribution of intentional action in particular (e.g., Hindriks forthcoming; Machery forthcoming; Nichols and Ulatowski 2007; Turner 2004). These researchers propose to explain the effect by positing a process that would apply only to attributions of intentional action and would not be expected to arise for any other aspect of folk psychology.

On the other side are researchers who think that the effect can be explained in terms of some very general fact about the relationship between folk psychology and moral judgment (e.g., Alicke forthcoming; Knobe 2006; Nadelhoffer 2006; Nado forthcoming). These researchers then proceed by constructing general theories about the ways in which moral judgments impact folk psychology. The guiding hope is that, if one can arrive at the correct general theory, the specific facts about intentional action will be seen to be just one aspect of a far broader pattern.

Our aim here is to provide experimental and theoretical support for this second view. On the theory we develop here, the surprising results obtained for intuitions about intentional action so far do not really have anything to do with intentional action in particular. Rather there is a perfectly general process whereby moral judgments serve as input to folk psychology, and the effects observed for intentional action should be understood as just one manifestation of this broader phenomenon. If we are right about this, the impact of moral judgments is not merely a peculiarity of the concept of

intentional action, but instead is a pervasive feature of the theory of mind.

Background

Consider a paradigmatic case of intentional action. The agent wants to bring about an outcome, she performs a behavior specifically for that purpose, and everything proceeds exactly as planned. In a case like this one, people's intuitions will be more or less independent of moral considerations. Regardless of whether the behavior is morally good or morally bad, almost everyone will say that the agent brought about the outcome intentionally.

Now consider a behavior that is paradigmatically unintentional. The agent has no interest in bringing about the outcome, she doesn't even know that her behavior might bring it about, and she only ends up acting as a result of some sort of muscle spasm. Here again, moral considerations will have little impact on people's intuitions. No matter what moral status the behavior has, almost everyone will say that the agent brings about the outcome unintentionally.

Things get interesting, however, when we consider intermediate cases – i.e., cases that fall somewhere between the paradigmatically intentional and the paradigmatically unintentional. Thus, suppose that the agent knows that she will be bringing about a particular outcome through her behavior but that she does not care about this outcome in any way. (She has chosen to perform the behavior for some other reason entirely.) In such a case, we might say that the outcome is a 'side-effect' of her behavior. Will people say that she brought about this side-effect intentionally? It turns out that their intuitions

in cases of this type can actually be influenced by their judgments about whether the side-effect itself is morally good or morally bad.¹

The usual way of demonstrating this influence of moral judgment on attributions of intentional action is to present experimental subjects with cases in which an agent brings about a side-effect that is either morally good or morally bad. Here, for example, is a case that we will call the *harm vignette*:

The vice-president of a company went to the chairman of the board and said, 'We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment.'

The chairman of the board answered, 'I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program.'

They started the new program. Sure enough, the environment was harmed.

After reading this vignette, subjects can be asked whether they agree or disagree with the statement: 'The chairman of the board intentionally harmed the environment.'

But now suppose we construct a case that is almost exactly the same as this first one, except that the side effect is actually morally good. We then arrive at what we will call the *help vignette*:

The vice-president of a company went to the chairman of the board and said, 'We are thinking of starting a new program. It will help us increase profits, and it will also help the environment.'

The chairman of the board answered, 'I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program.'

They started the new program. Sure enough, the environment was helped.

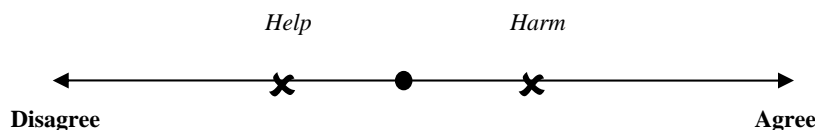
¹ Here we are simplifying for the sake of expository convenience. It is widely agreed that some sort of normative judgment is affecting people's intuitions, but no one now thinks that the relevant judgment is just a judgment about whether the side-effect is morally good or morally bad. At this point, the two major views are (a) that the effect is the result of a complex interplay among a number of different normative judgments (Machery forthcoming; Phelan and Sarkissian forthcoming; Wright and Bengson 2007) and (b) that none of our consciously-held normative judgments are playing a role and that the effect should be understood instead in terms of a fast, automatic, entirely non-conscious kind of normative appraisal (Alicke forthcoming; Knobe 2007; Nadelhoffer 2004). The details of this debate will not prove relevant to any of the questions we discuss here, and we therefore put the issue to one side.

After reading this second vignette, subjects can be asked whether they agree or disagree with the statement: ‘The chairman of the board intentionally helped the environment.’

Experimental studies concerning intuitions about cases like these consistently show a striking asymmetry (Feltz and Cokely 2007; Knobe 2003; Mallon forthcoming; Nichols and Ulatowski 2007; Phelan and Sarkissian forthcoming). Subjects who receive the harm vignette typically say that the agent harmed the environment intentionally, whereas subjects who receive the help vignette typically say that the agent helped the environment unintentionally. Yet it seems that the agent’s mental states do not differ between the two cases. The main difference lies instead in the moral status of the side-effect itself. Hence, most researchers have concluded that people’s moral judgments are somehow influencing their intuitions as to whether or not an agent acts intentionally (Knobe 2006; Malle 2006; Nadelhoffer 2006; Nado forthcoming).

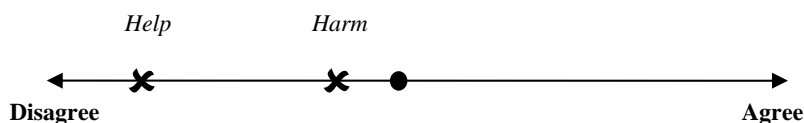
The key question now is whether this effect has something to do with the concept of intentional action in particular or whether it is simply one manifestation of a pervasive influence of moral judgment on folk psychology. In the experiment we have been discussing thus far, subjects were presented with the help and harm vignettes and asked in each case whether the agent acted *intentionally*, but what would have happened if they had instead been asked a question using some other folk-psychological concept? Suppose they had been asked whether the agent had a *desire* to help or harm the environment. Or suppose they had been asked whether the agent was *in favor* of helping or harming the environment. Would the effect then have disappeared? Or would we have found the very same asymmetry using those concepts as well?

Before we address these issues in earnest, a word is in order about the precise way in which we will be measuring levels of agreement and disagreement. In each of the experiments we report here, subjects are presented with a sentence and then asked to rate that sentence on a scale from ‘disagree’ to ‘agree.’ When subjects were given the sentences about helping and harming and asked whether the agent acted intentionally, their responses came out roughly as follows (Knobe 2005):



Note that the ratings for ‘help’ and ‘harm’ are actually on opposite sides of the midpoint, with the rating for ‘help’ on the side of disagreement and the rating for ‘harm’ on the side of agreement. It might be tempting, then, to suppose that the key result of the study is that subjects concluded on the whole that the agent acted intentionally in the harm vignette but unintentionally in the help vignette.

We think that this temptation should be resisted. For present purposes, the important thing is not whether people’s responses fall on one or the other side of the midpoint but rather whether there is a *difference* between responses to the morally bad case and responses to the morally good case. After all, suppose that some other factor – a factor that had nothing at all to do with the issues under discussion in this paper – shifted both types of responses a little bit to the left. The two responses might then come out as follows:



In this latter case, responses in both conditions are to the left of the midpoint, and a certain approach to thinking about people's intuitions might leave us with the idea that the key result was simply that, on the whole, subjects were disagreeing in both conditions. But this sort of approach serves only to obscure what is most relevant here. The important point is that moral judgments are influencing intuitions in this case in precisely the same way, and perhaps for the very same reasons, that moral judgments influence intuitions in the case described above.

Our approach will therefore be to focus not so much on the absolute levels of agreement in each case as on the differences between levels of agreement in morally good and morally bad cases. Using this basic approach, we can then examine the impact of moral judgment on people's application of an array of different concepts.

Evidence of Pervasiveness

When one pursues this research program, one quickly runs up against a surprising result. Not only does the impact of moral judgment extend beyond the concept of intentional action, moral judgments appear to be having some impact on just about every concept that involves holding or displaying a positive attitude toward an outcome. We will present data on six different concepts in this section, then turn to another two cases shortly thereafter.

1. 'Intention' and 'Intend'

One striking finding from recent work on the concept of intentional action is the surprising difference between people's use of the adverb 'intentionally' and their use of

the verb ‘intend’ and the noun ‘intention.’ Perhaps the strongest evidence here comes from a study by McCann (2005) in which subjects were given the harm vignette and asked:

- Did the chairman intentionally harm the environment?
- Did the chairman intend to harm the environment?
- Was it the chairman’s intention to harm the environment?

In that study, most subjects (64%) said that the agent acted ‘intentionally,’ but less than half (42%) said that he did ‘intend’ and relatively few (27%) said that he had an ‘intention.’

At this point, one might conclude that morality does not have the same sort of effect on ‘intend’ and ‘intention’ that it does on ‘intentionally.’ (After all, the majority of subjects in the study are disagreeing with the claim that the agent ‘intended’ or had the ‘intention.’) But appearances here are misleading. While only a minority of subjects are applying these terms in the harm case, one can still see evidence of a moral asymmetry.

Thus, in one recent study (Knobe 2004), subjects were randomly assigned to receive either the help vignette or the harm vignette and then asked:

- Was it the chairman’s intention to harm [help] the environment?

Although relatively few (29%) subjects said that the agent had an intention to harm, absolutely none (0%) said he had an intention to help. So people tended not to ascribe intention in either of these cases, but they were more likely to ascribe intention in the case where the behavior was morally bad.

Similar effects have been observed for the verb ‘intend.’ Cushman (2007) developed 21 different scenarios about agents who brought about side-effects. Each

scenario was constructed with two versions – one in which the action is morally good, another in which the action is morally bad. In all 21 scenarios, subjects showed higher levels of agreement with the statement that the agent ‘intended’ to bring about the side-effect in the morally bad version than in the morally good version.

2. *‘Desire’*

Here one might suspect that the words ‘intentionally,’ ‘intend’ and ‘intention’ all express more or less the same concept and that the effect might disappear as soon as one turns to words that express other folk-psychological concepts. That, however, appears not to be the case. In fact, the effect also emerges when one looks at applications of ‘desire.’

Tannenbaum, Ditto and Pizarro (2007) conducted a study in which subjects were presented with the help and harm vignettes and then asked:

- Did the chairman have a desire to help [harm] the environment?

Subjects marked their answers to this question on a scale from 1 to 7. The mean for the help vignette was 1.6; the mean for the harm vignette was 3.4. Here again, although subjects in both conditions leaned toward a negative answer to the question, subjects assigned significantly higher ratings in the morally bad case than in the morally good case.

3. *‘Decided’*

In light of these earlier results, we suspected that the effect would also arise for ‘decided.’

We therefore conducted an additional experiment.

Subjects were 37 undergraduate students taking philosophy classes at UNC-Chapel Hill. Each subject was randomly assigned to receive either the help vignette or the harm vignette. Subjects were then asked whether they agreed or disagreed with the statement:

- The chairman decided to help [harm] the environment.

Ratings were recorded on a scale from 1 ('disagree') to 7 ('agree'). The mean rating for the help condition was 2.7; the mean for the harm condition was 4.6. This difference is statistically significant, $t(35) = 2.4, p < .05$.

4. *'Advocated' and 'In Favor Of'*

Given that the effect had emerged for so many other folk-psychological concepts, we predicted that we would be able to find it even if we simply selected arbitrary expressions that in some way indicated that an agent was holding or displaying a positive attitude toward a given outcome. We chose the expressions 'advocated' and 'in favor of.'

Subjects were 62 students taking undergraduate philosophy classes at UNC-Chapel Hill. The experiment used a 2x2 design, with each subject randomly assigned to receive a story with a particular moral status (harm or help) and also randomly assigned to a particular question type ('advocated' or 'in favor of').

Subjects in the harm condition received the following vignette:

The management of a popular coffee franchise held a meeting to discuss a new procedure for preparing and serving coffee.

The assistant manager spoke forcefully in favor of adopting the new procedure, saying:

I know that this new procedure will mean more work for the employees, which will make them very unhappy. But that is not what

we should be concerned about. The new procedure will increase profits, and that should be our goal.

Subjects in the help condition received a vignette that was almost exactly the same, except that the assistant manager argued for a policy that would mean *less* work for the employees:

I know that this new procedure will mean less work for the employees, which will make them very happy. But that is not what we should be concerned about. The new procedure will increase profits, and that should be our goal.

Subjects were then asked whether they agreed or disagreed with a particular statement about the vignette. Each subject was randomly assigned to receive either a statement claiming that the agent ‘advocated’ bringing about an effect or that the agent was ‘in favor of’ bringing about an effect. Hence, the possible statements were:

- The assistant manager advocated [was in favor of] making the employees do more work.
- The assistant manager advocated [was in favor of] making the employees do less work.

Subjects rated each statement on a scale from 1 (‘disagree’) to 7 (‘agree’). The results are displayed in Table 1.

	Harm	Help
Advocated	4.1	2.8
In Favor Of	3.8	2.6

Overall, there was a significant main effect such that subjects were more inclined to agree in the harm condition than in the help condition, $F(1, 58) = 4.6$, $p < .05$. There was no

significant difference between the two question types ('advocated' vs. 'in favor of'), nor was there any significant interaction between moral status and question type.

Discussion

In light of these results, we are inclined to think that the impact of moral judgment is pervasive, playing a role in the application of *every* concept that involves holding or displaying a positive attitude toward an outcome. That is, for all concepts of this basic type, we suspect that there is a psychological process that makes people more willing to apply the concept in cases of morally bad side-effects and less willing to apply the concept in cases of morally good side-effects. Of course, this effect might be difficult to detect when we turn to concepts that almost all subjects would be unwilling to apply in any such case, but our hypothesis is that the same basic psychological process is still operating, however undetectably, even on those concepts.

On this view, there is nothing special about the concept of intentional action in particular that makes moral judgments relevant to it. Rather, the significance of the concept of intentional action is simply that people fall somewhere near the midpoint in their willingness to apply this concept in cases of morally neutral side-effects (Mele and Cushman 2007). Hence, when people become less willing in the morally good case and more willing in the morally bad case, their answers end up falling on opposite sides of the midpoint, and the asymmetry is therefore especially easy to detect.

A Tentative Hypothesis

Thus far, we have been providing evidence for the view that moral considerations affect the application of a wide array of different concepts. The question now is why so many different concepts should be subject to this same basic effect.

In addressing this question, we will be adopting a somewhat unusual approach. We will not offer anything like a full picture of any of the concepts under discussion here. Instead, we will focus only on the common element that we believe they all share. The hypothesis we present here describes just this one element and does not make any claims about the other aspects of the relevant concepts.

Although this approach may be disappointing to some readers, our decision to adopt it reflects a more general view about how best to make progress in this domain. Looking on the surface, one finds that people have intuitions involving various different concepts – intuitions about intentional action, intuitions about desire, and so forth – and it is therefore natural to suppose that there must be some kind of underlying process corresponding to each of these types of intuitions. It seems to us, however, that this view is not quite right. We prefer a picture according to which the relationship between types of intuition and underlying processes is many-to-many – each type of intuition is subserved by a number of distinct underlying processes and each underlying process subserves a number of different types of intuition. Hence, we suspect that it will not be helpful just to group things together by the way they appear on the surface and go looking for something like a ‘theory of intentional action intuitions’ or a ‘theory of desire intuitions.’ Instead, the aim should be to develop theories about the various psychological processes that underlie those intuitions. Predictions about the intuitions one finds on the

surface will simply fall out of our understanding of these various processes and their interaction.

Let us begin, then, by asking what might be similar about the various concepts we have been investigating thus far. It seems to us that the common element that all of these concepts share is that they all involve the idea of a certain kind of *positivity* about an outcome – the idea of supporting or approaching or favoring an outcome (rather than, say, opposing or denouncing it). Obviously, there are numerous differences between the concept *in favor of* and the concept *advocating*, but we are suggesting that both of these concepts somehow involve the idea of positivity, and it is this shared element that we wish to investigate.

The key question then becomes how people represent the idea of positivity. We propose that the underlying representation here is not a dichotomous on/off judgment but rather a matter of *degree*.² In other words, people can represent an agent's attitude toward an outcome as being anywhere on a scale from completely positive to completely negative. At one end of the scale would be the state of an agent who has an overwhelmingly favorable attitude toward the outcome. At the other end would be the state of an agent who has an overwhelmingly unfavorable attitude toward the outcome. Intermediate cases could then be represented by points toward the middle of the scale.

It may be helpful here to represent points along this scale using numbers. The number 0 can be used to represent the neutral 'default' state. One can then use positive numbers for positive attitudes and negative numbers for negative attitudes. Thus, the number +1 could represent a slightly positive attitude, the number +50 an

² Here we abandon certain aspects of the theory in Knobe (2006) in light of the powerful objections leveled against that theory by Malle (2007).

overwhelmingly positive attitude, the number -10 a moderately negative attitude, and so forth.

We can then introduce the hypothesis that people represent different folk-psychological concepts partially in terms of different positions along this same scale. People might represent the concept *wanting* as requiring a position above +10, the concept *desperate longing* as requiring a position above +100, the concept *mild aversion* as requiring a position between -1 and -10. Of course, each of these concepts might also have many other aspects. The suggestion is simply that, whatever else might be involved in the representation of these concepts, one aspect of that representation is a position along a scale from negative attitudes to positive attitudes.

With this framework in place, we can easily make sense of the idea that people's responses to sentences applying folk-psychological concepts are not simply 'yes' or 'no' but instead exhibit various different levels of agreement. The basic idea is that people show ever more disagreement as the state of the agent on the scale grows ever farther away from the state required by the concept. Thus, suppose that the concept *intending* requires a level of at least +10. If a given agent is in a state that is represented as +8, people will recognize that she at least comes fairly close to fulfilling the requirement. Although they may not agree with the claim that she 'intends,' they will at least show a fairly low level of disagreement.

Moreover, we now arrive at a different way of understanding the relations among distinct folk-psychological concepts. Consider the concepts *intentionally*, *reluctantly*, *desire*, *trying*, *regretting*, and so forth. It seems that these various concepts are related in some way, but it has proven difficult to say precisely what the relation is. One approach

would be to suppose that some of these concepts can be ‘analyzed’ into combinations of the others. (For example, one might think that the concept *intentionally* could be analyzed in terms of the concept of *desire*.) The hypothesis under discussion here involves a rejection of this basic picture. Instead of taking certain folk-psychological concepts as primitive and then analyzing some folk-psychological concepts in terms of others, we are assuming that all of these concepts can be represented (at least in part) in the same basic way – namely, as positions along the same basic dimension. Thus, we might conclude that there is a single underlying scale on which the concept *desire* requires a value of at least +30, *intentionally* requires a value of at least +5, and so forth. It then becomes possible to explain how these various concepts might show systematic relations to each other, even if none of them can be defined in terms of any other.

Most importantly, we now have the resources available to explain how moral judgments could have a pervasive impact on people’s application of folk-psychological concepts. The key notion here is that a wide variety of different folk-psychological concepts can be represented at least in part in terms of positions along the very same kind of continuous scale. Hence, there is no need to suppose that moral features actually figure in the representations associated with each of these individual concepts. *Moral judgments* only need to directly affect one particular type of representation – the underlying representation of an agent’s positive attitude. It then follows automatically that these judgments will come to have an impact on people’s application of a huge variety of different concepts.

If one accepts this basic framework, a question immediately arises as to *why* people’s moral judgments would affect their underlying representations of an agent’s

attitude. We will propose one possible answer to this question, but we first want to emphasize that the basic framework we have been developing thus far does not depend in any way on the assumption that this answer is correct. That is to say, even if the specific hypothesis we offer in the paragraphs below turns out to be entirely mistaken, there might be no need to reject the basic idea behind the approach offered here – that the pervasive impact of moral judgment arises because people’s moral judgments are affecting a kind of underlying representation of positivity.

That being said, let us now introduce the hypothesis. We propose that moral judgments affect the underlying representation through a kind of *calibration*. To see what we mean here, consider the steps one might take in making a thermometer. One might begin by constructing a device such that the level of mercury in a tube is appropriately correlated with the outside temperature, but the process would not simply end there. One would also have to find some way of drawing lines on the tube to signify which level counts as 0° C, which as 10° C, etc. This latter step is the process of ‘calibrating’ the thermometer. Now, what we want to suggest is that there is a similar process of calibration at work in people’s representation of positivity. People might be able to construct some sort of representation of the degree to which an agent’s attitude is positive, but there is also a second and very important step in the process. People need some way of determining which level will count as 0, which as +10, etc. It is here, we propose, that moral judgment plays a role.

We said above that the 0 point would be assigned to the ‘default’ attitude. But which exact attitude qualifies as the *default*? Before offering any general theory here, we can simply look to a few intuitions about cases. First consider the case of attitudes

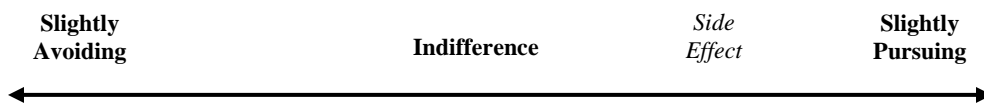
toward environmental harm. Independent of any theory, it just strikes us as intuitive that the default attitude toward harming the environment is one that involves being at least slightly motivated to prevent this outcome. (We would then say that an agent who has even less motivation to prevent the outcome is ‘especially unconcerned’ and that an agent could only be ‘especially concerned’ if she showed a rather strong interest in preventing it.) When it comes to helping the environment, however, our intuitions look very different. There, we take it that the default level is to be at least slightly motivated to pursue the outcome and that the various other levels of motivation can then be defined in terms of their differences from this default.

If these intuitions are on the right track, it seems that the level of motivation regarded as the default depends in some way on moral considerations. Perhaps people assign the default level to a kind of *moral mean* – not a statistically average level of motivation, but a normative level of motivation that we morally expect of others. It matters little, however, whether this specific account proves correct. The key point for present purposes is just that the level of motivation people regard as the ‘default’ depends in some way on moral considerations.

What we have been constructing so far is a general framework for thinking about the attribution of positive and negative attitudes. We can now ask what implications this framework would have for the attribution of positive attitudes in cases of side-effects.

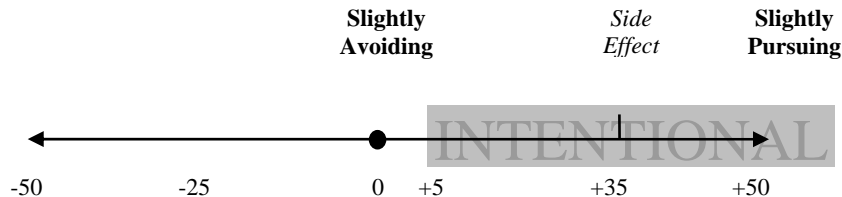
If our approach here is correct, the attitude of the agent toward the side-effect should be represented as a particular position along a continuous scale. It would be represented as slightly less positive than the attitude of an agent who truly did feel positively about the outcome in and of itself, but it would nonetheless be represented as

somewhat more positive than the attitude of an agent who did not hold any positive attitude toward it at all. In other words, the attitude of an agent who welcomes something as a side-effect would be represented as having a kind of intermediate status. We can depict this intermediate status along the continuum as follows:



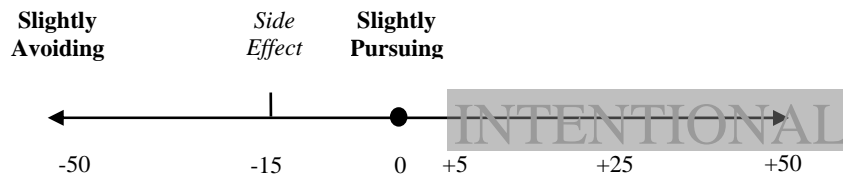
Suppose now that people only regard a behavior as intentional when its position falls somewhere above +5. The question now is what sort of attitude the agent needs to have to fulfill this requirement.

The answer, of course, is that it depends on the moral status of the behavior itself. Consider first what happens when the behavior in question is *harming the environment*. Since we morally expect others to be motivated to avoid bad outcomes, the 0-point will be a state of being slightly motivated to avoid the outcome, and the continuum will end up looking like this:



Since accepting an outcome as a side-effect now falls above +5, people will regard agents with this attitude as having performed a behavior intentionally.

But now consider what occurs when the behavior is *helping the environment*. The 0 point then gets assigned to a state of being slightly motivated to pursue the outcome, and the continuum ends up looking more like this:



In this latter configuration, the state of accepting an outcome as a side-effect falls below +5, and the behavior is regarded as unintentional.

In this way, we can explain the effects observed for intuitions about intentional action without assuming anything special about the concept of intentional action in particular. The very same explanation can thus be applied to each of the other concepts discussed above.

Testing the Hypothesis

Let us now be frank. Numerous hypotheses have already been offered to explain these phenomena, and each in turn has fallen prey to experimental falsification. It therefore seems overwhelmingly likely that our own hypothesis will turn out to be false as well. Perhaps the model adopted here will help to guide some future research, but it seems to us that we have good reason to expect that, sooner or later, some clever researcher will conduct an experiment that refutes our view.

Oddly enough, we think there might actually be good grounds for an even more specific prediction. Not only can we predict that our hypothesis will be refuted, we can

predict roughly *how* it will be refuted. In each of the earlier cases, someone offered a hypothesis that predicted that the effect would only arise in a fairly circumscribed range of cases, but the hypothesis was then falsified when further research showed that the effect also arises in cases that fall outside of that range. We suspect that our own hypothesis will fall victim to the same difficulty. Sooner or later, someone is going to show that the effects we have been investigating thus far are really just one special case of a far more pervasive phenomenon.

In the meantime, though, it can be seen that the hypothesis we put forward above predicts an impact of moral judgments that goes considerably beyond the cases discussed in the previous section. In other words, the hypothesis was developed to explain the impact of moral judgments on attributions of positive attitudes in cases of side-effects, but one can see immediately that the hypothesis ends up predicting an impact of moral judgments on many other kinds of cases as well. We therefore conducted two additional experiments to examine the impact of moral considerations in cases of a rather different sort.

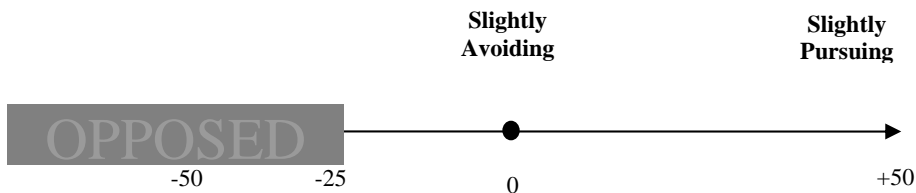
5. *'Opposed'*

Thus far, we have been concerned exclusively with the attribution of positive attitudes: 'intending,' 'desiring,' 'in favor of,' and so forth. In each of these cases, one finds an attitude whereby the agent is favorably disposed to an outcome or motivated to pursue it. But suppose we now try to extend our investigation to negative attitudes. For example, instead of simply considering intuitions about whether an agent is 'in favor' of a given

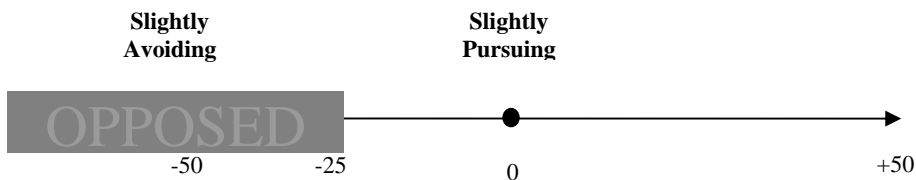
outcome, suppose we consider intuitions about whether the agent is ‘opposed’ to an outcome.

It follows from the hypothesis we advanced above that people’s moral judgments will have an impact here too – but that this time the impact will go in the opposite direction. While people were more inclined to say that an agent was ‘in favor’ of harming the environment than helping it, they should be more inclined to say that an agent is ‘opposed’ to helping the environment than to harming it.

To see why this is so, one need only suppose that the concept *opposed* requires some negative value along an underlying scale (say, a value at least below -25). Then people’s representation of an agent’s attitude towards a harmful outcome might look like this:



But the very same sort of attitude toward a helpful outcome might look like this:



The result is that the very same sort of attitude should be able to count as being ‘opposed’ in the case of a helpful outcome but not in the case of a harmful outcome.

To test this hypothesis, we conducted an additional experiment. Subjects were 56 students taking philosophy classes at UNC-Chapel Hill. Each subject was randomly assigned either to the ‘harm’ condition or the ‘help’ condition.

Subjects in the harm condition received the following vignette:

The CEO of a company was talking with his assistant. The assistant said: ‘We have conducted an in-depth study of the company’s proposed new policy. Our study shows that the new policy would decrease profits for the company and that it would also harm the environment.’

The CEO said: ‘Look, I don’t really care about what happens to the environment. What I care about is making sure that our profits don’t decrease. So, with that in mind, let’s definitely not implement that new policy.’

Subjects in the help condition received a vignette that was exactly the same, except that the word ‘harm’ was replaced with ‘help.’ Thus, the vignette in this condition told the story of a policy that would help the environment and an executive who was against that policy because he knew that it would decrease profits.

After reading their vignettes, subjects were asked whether they agreed or disagreed with the statements:

- The CEO was opposed to harming [helping] the environment.
- The CEO deserves *blame* for what he did.

All statements were rated on a scale from 1 to 7.

There was no significant difference between conditions on the statement about blameworthiness. For the statement about being ‘opposed,’ ratings for subjects in the harm condition, $M = 2.3$, were significantly lower than ratings for subjects in the help condition, $M = 3.4$, $t(54) = 2.0$, $p < .05$.

6. *Trying without Foresight*

The hypothesis we have been developing thus far predicts an impact of moral judgment on just one particular type of representation – the underlying representation of the positivity of an agent’s attitude. A question arises, however, as to whether one might find similar effects for other types of representation as well.

Suppose, for example, that we turn away from representations of positivity and look instead at representations of *credence*. We will be turning, then, to representations of the degree to which an agent regards an outcome as likely or unlikely. Here too, one might posit a continuous dimension. At one extreme would be the state of an agent who is absolutely certain that an outcome will occur. At the other extreme would be the state of an agent who is completely convinced that it will not. The various intermediate positions could represent different degrees of uncertainty.

What we want to know now is whether people’s moral judgments affect their representation of an agent along this continuum. We can proceed here using our usual method. First, we need to find a concept that can be correctly applied only to agents who show a particular level of credence. Then we can ask whether intermediate cases in the application of this concept are treated differently depending on their moral status.

Consider in this light the concept *intending*. If a person wishes that she could bring about an outcome but thinks that there is absolutely no chance she will succeed, one would not normally say that she ‘intends’ to bring about this outcome. But the situation is very different when her credences go to the opposite extreme. If she wishes to bring about an outcome and is completely convinced that she will succeed, she seems clearly to ‘intend’ to bring the outcome about. What happens then in the intermediate case?

Suppose she thinks it is unlikely that she will succeed but nonetheless believes that there is at least some chance that she will bring about the outcome. The burning question is whether or not people's moral judgments will affect their application of the concept of intending in cases in this intermediate category.

To test address this question, we ran one final experiment. Subjects were 37 undergraduates taking philosophy classes at UNC- Chapel Hill. Each subject was randomly assigned to either the 'help condition' or the 'harm condition.' Subjects in the help condition read the following vignette:

A man wants to defuse a bomb that will kill thousands of innocent tourists if it explodes. The only way to defuse the bomb is to enter the correct code on a keypad, but the man does not know the code. There is only a one in ten million chance of his guessing the code.

The man is fully aware that there is virtually no chance that he will successfully defuse the bomb, but he desperately wants to save the tourists. So, without even looking at the keypad, he just randomly presses some keys.

Subjects in the harm condition read a vignette that was exactly the same, except that the word 'defuse' was replaced with 'detonate.' Thus, the vignette in the harm condition involved a man who is trying to detonate a bomb that will kill thousands of innocent tourists but who believes that he will probably fail.

After reading their vignettes, subjects were asked whether they agreed or disagreed with the statement:

- The man *intended* to defuse [detonate] the bomb.

To help subjects understand precisely what we were getting at here, we also included a second statement, namely:

- The man *wanted* to defuse [detonate] the bomb.

The order of these two statements was counterbalanced.

There were no significant order effects and no significant differences on the question about whether the agent ‘wanted.’ Indeed, for this question, subjects showed a ceiling effect, with a mean rating of 6.5.

On the question about whether the agent ‘intended,’ the mean rating for subjects in the help condition was 3.8; the mean rating for subjects in the harm condition was 5.6. This difference is statistically significant, $t(35) = 2.5, p < .05$.

The one sad thing about this result is that it is in no way predicted by the hypothesis we presented above. That hypothesis predicts an effect in cases that are intermediate in positivity but does not also predict an effect in cases that are intermediate in credence level. (Of course, the hypothesis does not specifically predict that one will *not* find an effect in this latter type of case.)

Now, one possible reaction at this point would be to suggest that there is a process whereby moral judgments affect representations of positivity and then another, completely separate process whereby moral judgments affect representations of credence. But we hope that someone will be able to do better than that. Our hope is that someone will be able to find a single underlying process that explains all of the phenomena we have been discussing here.

Discussion

In an earlier section, we showed an influence of moral considerations in intuitions about positive attitudes toward side-effects. It now appears that the effect is not limited to that narrow range of cases. Far from it: the data presented here suggests that that effect

actually arises both in cases that do not involve positive attitudes and in cases that do not involve side effects.

At this point, we think there is little hope of developing a theory that somehow treats each type of case individually. What is needed is a general theory about the impact of moral judgment on folk psychology. It should then be possible to show how that general theory can explain the specific types of judgments we find in these various types of cases. Given the general theory, and given some general facts about how people represent positive attitudes and side-effects, it should be possible to derive the prediction that people's attribution of positive attitudes in side-effect cases will depend in part on their moral judgments.

Although the hypothesis proposed here is probably mistaken in certain details, it at least has the right form. The theory posits a role for moral considerations in absolutely all attributions of positive or negative attitudes. It then predicts that moral considerations will not be sufficient to shift people's intuitions in many types of cases but that they will be sufficient to shift people's judgments in cases that have an 'intermediate' character. In other words, the framework introduced above makes it possible to derive the effects observed for positive attitudes toward side-effects from a perfectly general theory.

Conclusion

When experimental studies first began showing that moral considerations could influence the application of folk-psychological concepts, it might have been thought that this effect would be limited to a tightly constrained range of cases. One could have supposed, e.g., that the effect would only arise for the concept of intentional action, or that it would only

arise in cases of side-effects, or that there would be some other, fairly narrow range of circumstances in which it could be found. It could then have been supposed that there was a kind of ‘core’ of folk-psychology that was entirely free of the impact of moral judgment.

Plausible though it may have seemed, this view appears not to be correct. On the contrary, as we learn more and more about the application of various different folk-psychological concepts, we are coming to find an impact of moral considerations in more and more places. It seems to us that there is now good reason to believe there are no concepts anywhere in folk psychology that enable one to describe an agent’s positive or negative attitudes in a way that is entirely independent of moral considerations. The impact of moral judgments, we suspect, is utterly pervasive.

*Department of Philosophy
UNC Chapel Hill*

References

Alicke, M. forthcoming: Blaming badly. *Journal of Cognition and Culture*.

Cushman, F. 2007: The effect of moral judgment on causal and intentional attribution: What we say, or how we think? Unpublished manuscript. Harvard University.

Feltz, A. and Cokely, E. 2007: An anomaly in intentional action ascription: more evidence of folk diversity. *Proceedings of the Cognitive Science Society*.

Hindriks, F. A. forthcoming: Intentional action and the praise-blame asymmetry. *Philosophical Quarterly*.

Knobe, J. 2003: Intentional action and side effects in ordinary language. *Analysis*, 63, 190-193.

Knobe, J. 2004: Intention, intentional action and moral considerations. *Analysis*, 64, 181-187.

- Knobe, J. 2005: Folk psychology and folk morality: response to critics. *Journal of Theoretical and Philosophical Psychology*, 24, 252-258.
- Knobe, J. 2006: The concept of intentional action: A case study in the uses of folk psychology. *Philosophical Studies*, 130, 203-231
- Knobe, J. 2007: Reason explanation in folk psychology. *Midwest Studies in Philosophy*, 31, 90-106.
- Machery, E. forthcoming: Understanding the folk concept of intentional action: philosophical and experimental issues. *Mind and Language*.
- Malle, B. 2006: Intentionality, morality, and their relationship in human judgment. *Journal of Cognition and Culture*, 6, 87-113.
- Malle, B. F. 2007: The puzzle of intentionality and moral cognition. Paper presented at the Society of Personality and Social Psychology Annual Convention, Memphis, TN.
- Mallon, R. forthcoming: Knobe vs. Machery: testing the trade-off hypothesis. *Mind and Language*.
- McCann, H. 2005: Intentional action and intending: recent empirical studies. *Philosophical Psychology*, 18, 737-748.
- Mele, A. and Cushman, F. 2007: Intentional action, folk judgments, and stories: sorting things out. *Midwest Studies in Philosophy*, 31, 184-201.
- Nadelhoffer, T. 2004: Blame, badness, and intentional action: a reply to Knobe and Mendlow. *Journal of Theoretical and Philosophical Psychology*, 24, 259-269.
- Nadelhoffer, T. 2006: Bad acts, blameworthy agents, and intentional actions: some problems for jury impartiality. *Philosophical Explorations*, 9(2), 203-220.
- Nado, J. forthcoming: Effects of moral cognition on judgments of intentionality. *British Journal for the Philosophy of Science*.
- Nichols, S. and Ulatowski, J. 2007: Intuitions and individual differences: the Knobe effect revisited. *Mind and Language*, 22, 346-365.
- Phelan, M. and Sarkissian, H. forthcoming: The folk strike back; or, why you didn't do it intentionally, though it was bad and you knew it. *Philosophical Studies*.
- Phelan, M. and Sarkissian, H. forthcoming: Is the 'trade-off hypothesis' worth trading for? *Mind and Language*.
- Tannenbaum, D., Ditto, P.H., and Pizarro, D.A. 2007: Different moral values produce different judgments of intentional action. Unpublished manuscript. University of California-Irvine.

Turner, J. 2004: Folk intuitions, asymmetry, and intentional side effects. *Journal of Theoretical and Philosophical Psychology*, 24, 214-219.

Wright, J. and Bengson, J. 2007: Asymmetries in folk judgments of responsibility and intentional action. Unpublished manuscript. University of Wyoming.