



# The Right to Lie: Kant on Dealing with Evil

## Citation

Korsgaard, Christine. 1986. The right to lie: Kant on dealing with evil. *Philosophy and Public Affairs* 15, no. 4: 325-349.

## Published Version

<http://www.wiley.com/bw/journal.asp?ref=0048-3915>

## Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:3200670>

## Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

## **The Right to Lie: Kant on Dealing with Evil<sup>1</sup>**

**Christine M. Korsgaard**

One of the great difficulties with Kant's moral philosophy is that it seems to imply that our moral obligations leave us powerless in the face of evil. Kant's theory sets a high ideal of conduct and tells us to live up to that ideal regardless of what other persons are doing. The results may be very bad. But Kant says that the law "remains in full force, because it commands categorically." (G 438-439/57)<sup>2</sup> The most well-known example of this "rigorism", as it is sometimes called, concerns Kant's views on our duty to tell the truth.

In two passages in his ethical writings, Kant seems to endorse the following pair of claims about this duty: First, one must never under any circumstances or for any purpose tell a lie. Second, if one does tell a lie one is responsible for all of the consequences that ensue, even if they were completely unforeseeable.

One of the two passages occurs in the *Metaphysical Principles of Virtue*. There Kant classifies lying as a violation of a perfect duty to oneself. In one of the casuistical questions, a servant, under instructions, tells a visitor the lie that his master is not at home. His master, meanwhile, sneaks off and commits a crime, which would have been prevented by the watchman sent to arrest him. Kant says:

Upon whom ... does the blame fall? To be sure, also upon the servant, who here violated a duty to himself by lying, the consequence of which will now be imputed to him by his own conscience. (MMV 431/93)

The other passage is the infamous one about the murderer at the door from the essay, "On A Supposed Right to Lie From Altruistic Motives." Here Kant's claims are more extreme, for he says that the liar may be held legally as well as ethically responsible for the consequences, and the series of coincidences he imagines is even more fantastic:

After you have honestly answered the murderer's question as to whether his intended victim is at home, it may be that he has slipped out so that he does not come in the way of the murderer, and thus that the murder may not be committed. But if you had lied and said he was not at home when he had really gone out without your knowing it, and if the murderer had then met him as he went away and murdered him, you might justly be accused as the cause of his death. For if you had told the truth as far as you knew it, perhaps the murderer might have been apprehended by the neighbors while he searched the house and thus the deed might have been prevented. (SRL 427/348)

Kant's readers differ about whether Kant's moral philosophy commits him to the claims he makes in these passages. Unsympathetic readers are inclined to take them as evidence of the horrifying conclusions to which Kant was led by his notion that the necessity in duty is rational necessity - as if Kant were clinging to a logical point in the teeth of moral decency. Such readers take these conclusions as a defeat for Kant's ethics, or for ethical rationalism generally; or they take Kant to have confused principles which are merely general in their application and *prima facie* in their truth with absolute and universal laws. Sympathetic readers are likely to argue that Kant here mistook the implications of his own theory, and to try to show that, by careful

construction and accurate testing of the maxim on which this liar acts, Kant's conclusions can be blocked by his own procedures.

Sympathetic and unsympathetic readers alike have focused their attention on the implications of the first formulation of the categorical imperative, the Formula of Universal Law. The *Foundations of the Metaphysics of Morals* contains two other sets of terms in which the categorical imperative is formulated: the treatment of humanity as an end in itself, and autonomy, or legislative membership in a Kingdom of Ends. My treatment of the issue falls into three parts. First, I want to argue that Kant's defenders are right in thinking that, when the case is treated under the Formula of Universal Law, this particular lie can be shown to be permissible. Second, I want to argue that when the case is treated from the perspective provided by the Formulas of Humanity and the Kingdom of Ends, it becomes clear why Kant is committed to the view that lying is wrong in every case. But from this perspective we see that Kant's rigorism about lying is not the result of a misplaced love of consistency or legalistic thinking. Instead, it comes from an attractive ideal of human relations which is the basis of his ethical system. If Kant is wrong in his conclusion about lying to the murderer at the door, it is for the interesting and important reason that morality itself sometimes allows or even requires us to do something that from an ideal perspective is wrong. The case does not impugn Kant's ethics as an *ideal* system. Instead, it shows that we need special principles for dealing with evil. My third aim is to discuss the structure that an ethical system must have in order to accommodate such special principles.

### Universal Law

The Formula of Universal Law tells us never to act on a maxim that we could not at the same time will to be a universal law. A maxim which cannot even be conceived as a universal law without contradiction is in violation of a strict and perfect duty, one which assigns us a particular action or omission. A maxim which cannot be willed as

universal law without contradicting the will is in violation of a broad and imperfect duty, one which assigns us an end, but does not tell us what or how much we should do towards it. Maxims of lying are violations of perfect duty, and so are supposed to be the kind that cannot be conceived without contradiction when universalized.

The sense in which the universalization of an immoral maxim is supposed to "contradict" itself is a matter of controversy. On my reading, which I will not defend here<sup>3</sup>, the contradiction in question is a "practical" one: the universalized maxim contradicts itself when the efficacy of the action as a method of achieving its purpose would be undermined by its universal practice. So, to use Kant's example, the point against false promising as a method of getting ready cash is that if everyone attempted to use false promising as a method of getting ready cash, false promising would no longer *work* as a method of getting ready cash, since, as Kant says, "no one would believe what was promised to him but would only laugh at any such assertion as vain pretense." (G 422/40)

Thus the test question will be: could this action be the universal method of achieving this purpose? Now when we consider lying in general, it looks as if it could not be the universal method of doing anything. For lies are usually efficacious in achieving their purposes because they deceive, but if they were universally practiced they would not deceive. We believe what is said to us in a given context because most of the time people in that context say what they really think or intend. In contexts in which people usually say false things - e.g., when telling stories that are jokes - we are not deceived. If a story that is a joke and is false counts as a lie, we can say that a lie in this case is not wrong, because the universal practice of lying in the context of jokes does not interfere with the *purpose* of jokes, which is to amuse and does not depend on deception. But in most cases lying falls squarely into the category of the sort of action Kant considers wrong: actions whose efficacy depends upon the fact that most people

do not engage in them, and which therefore can only be performed by someone who makes an exception of himself. (G 424/42)

When we try to apply this test to the case of the murderer at the door, however, we run into a difficulty. The difficulty derives from the fact that there is probably already deception in the case. If murderers standardly came to the door and said: "I wish to murder your friend - is he here in your house?" then perhaps the universal practice of lying in order to keep a murderer from his victim would not work. If everyone lied in these circumstances the murderer would be aware of that fact and would not be deceived by your answer. But the murderer is not likely to do this, or, in any event, this is not how I shall imagine the case. A murderer who expects to conduct his business by asking questions must suppose that you do not know who he is and what he has in mind.<sup>4</sup> If these are the circumstances, and we try to ascertain whether there could be a universal practice of lying in these circumstances, the answer appears to be yes. The lie will be efficacious even if universally practiced. But the reason it will be efficacious is rather odd: it is because the murderer supposes you do not know what circumstances you are in - that is, that you do not know you are addressing a murderer - and so does not conclude from the fact that people in those circumstances always lie that *you* will lie.

The same point can be made readily using Kant's publicity criterion. (PP 381-383/129-131) Can we announce in advance our intention of lying to murderers without, as Kant says, vitiating our own purposes by publishing our maxims? (PP 383/131) Again the answer is yes. It does not matter if you say publicly that you will lie in such a situation, for the murderer supposes that you do not know you are in that situation.<sup>5</sup>

These reflections might lead us to believe, then, that Kant was wrong in thinking that it is never all right to lie. It is permissible to lie to deceivers in order to counteract the intended results of their deceptions, for the maxim of lying to a deceiver is

universalizable. The deceiver has, so to speak, placed himself in a morally unprotected position by his own deception. He has created a situation which universalization cannot reach.

### Humanity

When we apply the Formula of Humanity, however, the argument against lying that results applies to any lie whatever. The formula runs:

Act so that you treat humanity, whether in your own person or in that of another, always as an end and never as a means only. (G 429/47)

In order to use this formula for casuistical purposes, we need to specify what counts as treating humanity as an end. "Humanity" is used by Kant specifically to refer to the capacity to determine ends through rational choice. (G 437/56; MMV 392/50) Imperfect duties arise from the obligation to make the exercise, preservation, and development of this capacity itself an end. The perfect duties - that is, the duties of justice, and, in the realm of ethics, the duties of respect - arise from the obligation to make each human being's capacity for autonomous choice the condition of the value of every other end.

In his treatment of the lying promise case under the Formula of Humanity, Kant makes the following comments:

For he whom I want to use for my own purposes by means of such a promise cannot possibly assent to my mode of acting against him and cannot contain the end of this action in himself. ... he who transgresses the rights of men intends to make use of the persons of others merely as means, without considering that as rational beings, they must always be esteemed at the same time as ends, i.e.

only as beings who must be able to contain in themselves  
the end of the very same action.

(G 429-430/48)

In these passages, Kant uses two expressions that are the key to understanding the derivation of perfect duties to others from the Formula of Humanity. One is that the other person "cannot possibly assent to my mode of acting toward him" and the second is that the other person cannot "contain the end of this action in himself." These phrases provide us with a test for perfect duties to others: an action is contrary to perfect duty if it is not possible for the other to assent to it or to hold its end.

It is important to see that these phrases do not mean simply that the other person *does not* or *would not* assent to the transaction or that she does not happen to have the same end I do, but strictly that she *cannot* do so: that something makes it impossible. If what we cannot assent to means merely what we are likely to be annoyed by, the test will be subjective and the claim that the person does not assent to being used as a means will sometimes be false. The object you steal from me may be the gift I intended for you, and we may both have been motivated by the desire that you should have it. And I may care about you too much or too little to be annoyed by the theft. For all that this must be a clear case of your using me as a mere means.<sup>6</sup>

So it must not be merely that your victim will not like the way that you propose to act, that this is psychologically unlikely, but that something makes it impossible for her to assent to it. Similarly, it must be argued that something makes it impossible for her to hold the end of the very same action. Kant never spells out why it is impossible, but it is not difficult to see what he has in mind.

People cannot *assent* to a way of acting when they are given no chance to do so. The most obvious instance of this is when coercion is used. But it is also true of deception: the victim of the false promise cannot assent to it because he doesn't know it is what he is being offered. But even when the victim of such conduct does happen to



know what is going on there is a sense in which he cannot assent to it. Suppose, for example, that you come to me and ask to borrow some money, falsely promising to pay it back next week, and suppose that by some chance I know perfectly well that your promise is a lie. Suppose also that I have the same end you do, in the sense that I want you to have the money, so that I turn the money over to you anyway. Now here I have the same end that you do, and I tolerate your attempts to deceive me to the extent that they do not prevent my giving you the money. Even in this case I cannot really assent to the transaction *you* propose. We can imagine the case in a number of different ways. If I call your bluff openly and say "never mind that nonsense, just take this money" then what I am doing is not accepting a false promise, but giving you a handout, and scorning your promise. The nature of the transaction is changed: now it is not a promise but a handout. If I don't call you on it, but keep my own counsel, it is still the same. I am not accepting a false promise. In this case what I am doing is *pretending* to accept your false promise. But there is all the difference in the world between actually doing something and pretending to do it. In neither of these cases can I be described as accepting a false promise, for in both cases I fix it so that it is something else that is happening. My knowledge of what is going on makes it *impossible* for me to accept the deceitful promise in the ordinary way.

The question whether another can assent to your way of acting can serve as a criterion for judging whether you are treating her as a mere means. We will say that knowledge of what is going on and some power over the proceedings are the conditions of possible assent; without these, the concept of assent does not apply. This gives us another way to formulate the test for treating someone as a mere means: Suppose it is the case that if the other person knows what you are trying to do and has the power to stop you, then what you are trying to do cannot be what is really happening. If this is the case, the action is one that by its very nature is impossible for the other to assent to. You cannot wrest from me what I freely give to you; and if I have the power to stop

you from wresting something from me and do not use it, I am in a sense freely giving it to you. This is of course not intended as a legal point: the point is that any action which depends for its nature and efficacy on the other's ignorance or powerlessness fails this test. Lying clearly falls into this category of action: it only deceives when the other does not know that it is a lie.<sup>7</sup>

A similar analysis can be given of the possibility of holding the end of the very same action. In cases of violation of perfect duty, lying included, the other person is unable to hold the end of the very same action because the way that you act prevents her from *choosing* whether to contribute to the realization of that end or not. Again, this is obviously true when someone is forced to contribute to an end, but it is also true in cases of deception. If you give a lying promise to get some money, the other person is invited to think that the end she is contributing to is your temporary possession of the money: in fact, it is your permanent possession of it. It doesn't matter whether that would be all right with her if she knew about it. What matters is that she never gets a chance to choose the end, not knowing that it is to be the consequence of her action.<sup>8</sup>

According to the Formula of Humanity, coercion and deception are the most fundamental forms of wrongdoing to others - the roots of all evil. Coercion and deception violate the conditions of possible assent, and all actions which depend for their nature and efficacy on their coercive or deceptive character are ones that others cannot assent to. Coercion and deception also make it impossible for others to choose to contribute to our ends. This in turn makes it impossible, according to Kant's value theory, for the ends of such actions to be good. For on Kant's view "what we call good must be, in the judgement of every reasonable man, an object of the faculty of desire." (C2 60/62-63) If your end is one that others cannot choose - not because of what they want, but because they are not in a position to choose - it cannot, as the end of that action, be good. This means that in any cooperative project - whenever you need the

decisions and actions of others in order to bring about your end - everyone who is to contribute must be in a position to *choose* to contribute to the end.

The sense in which a good end is an object for everyone is that a good end is in effect one that everyone, in principle, and especially everyone who contributes to it, gets to cast a vote on. This voting, or legislation, is the prerogative of rational beings; and the ideal of a world in which this prerogative is realized is the Kingdom of Ends.

### The Kingdom of Ends

The Kingdom of Ends is represented by the kingdom of nature; we determine moral laws by considering their viability as natural laws. On Kant's view, the will is a kind of causality. (G 446/64) A person, an end in itself, is a free cause, which is to say a first cause. By contrast a thing, a means, is a merely mediate cause, a link in the chain. A first cause is, obviously, the initiator of a causal chain, hence a real determiner of what will happen. The idea of deciding for yourself whether you will contribute to a given end can be represented as a decision whether to initiate that causal chain which constitutes your contribution. Any action which prevents or diverts you from making this initiating decision is one that treats you as a mediate rather than a first cause; hence as a mere means, a thing, a tool. Coercion and deception both do this. And deception treats you as a mediate cause in a specific way: it treats your reason as a mediate cause. The false promiser thinks: if I tell her I will pay her back next week, then she will choose to give me the money. Your reason is worked, like a machine: the deceiver tries to determine what levers to pull to get the desired results from you. Physical coercion treats someone's person as a tool; lying treats someone's *reason* as a tool. This is why Kant finds it so horrifying; it is a direct violation of autonomy.

We may say that a tool has two essential characteristics: It is there to be used, and it does not control itself: its nature is to be directed by something else. To treat

someone as a mere means is to treat her as if these things were true of her. Kant's treatment of our duties to others in the *Metaphysical Principles of Virtue* is sensitive to *both* characteristics. We are not only forbidden to use another as mere means to our private purposes. We are also forbidden to take attitudes towards her which involve regarding her as not in control of herself, which is to say, as not using her reason.

This latter is the basis of the duties of respect. Respect is violated by the vices of calumny and mockery (MMV 466-468/131-133): we owe to others not only a practical generosity toward their plans and projects - a duty of aid - but also a generosity of attitude toward their thoughts and motives. To treat another with respect is to treat him as if he were using his reason and as far as possible as if he were using it well. Even in a case where someone evidently *is* wrong or mistaken, we ought to suppose he must have what he takes to be good reasons for what he believes or what he does. This is not because, as a matter of fact, he probably does have good reasons. Rather, this attitude is something that we *owe* to him, something that is his right. And he cannot forfeit it. Kant is explicit about this:

Hereupon is founded a duty to respect man even in the logical use of his reason: not to censure someone's errors under the name of absurdity, inept judgement, and the like, but rather to suppose that in such an inept judgment there must be something true, and to seek it out. ... Thus it is also with the reproach of vice, which must never burst out in complete contempt or deny the wrongdoer all moral worth, because on that hypothesis he could never be improved either -- and this latter is incompatible with the idea of man, who as such (as a moral being) can never lose all predisposition to good. (MMV 463-464/128-129)

To treat others as ends in themselves is always to address and deal with them as rational beings. Every rational being gets to reason out, for herself, what she is to think or to choose or to do. So if you need someone's contribution to your end, you must put the facts before her and ask for her contribution. If you think she is doing something wrong, you may try to convince her by argument but you may not resort to tricks or force. The Kingdom of Ends is a democratic ideal, and poor judgment does not disqualify anyone for citizenship. In the *Critique of Pure Reason*, Kant says:

Reason depends on this freedom for its very existence.  
For reason has no dictatorial authority; its verdict is  
always simply the agreement of free citizens, of whom  
each one must be permitted to express, without let or  
hindrance, his objections or even his veto.<sup>9</sup>

This means that there cannot be a good reason for taking a decision out of someone else's hands. It is a rational being's prerogative, as a first cause, to have a share in determining the destiny of things.

This shows us in another way why lying is for Kant a paradigm case of treating someone as a mere means. Any attempt to control the actions and reactions of another by any means except an appeal to reason treats her as a mere means, because it attempts to reduce her to a mediate cause. This includes much more than the utterance of falsehoods. In the *Lectures on Ethics*, Kant says "whatever militates against frankness lowers the dignity of man." (LE 231)<sup>10</sup> It is an everyday temptation, even (or perhaps especially) in our dealings with those close to us, to withhold something, or to tidy up an anecdote, or to embellish a story, or even just to place a certain emphasis, in order to be sure of getting the reaction we want.<sup>11</sup> Kant holds the Socratic view that any sort of persuasion that is aimed at distracting its listener's attention from either the reasons that she ought to use or the reasons the speaker thinks she will use is wrong.<sup>12</sup>

In light of this account it is possible to explain why Kant says what he does about the liar's responsibility. In a Kantian theory our responsibility has definite boundaries: each person as a first cause exerts some influence on what happens, and it is your part that is up to you. If you make a straightforward appeal to the reason of another person, your responsibility ends there and the other's responsibility begins. But the liar tries to take the consequences out of the hands of others; he, and not they, will determine what form their contribution to destiny will take. By refusing to share with others the determination of events, the liar takes the world into his own hands, and makes the events his own. The results, good or bad, are imputable to him, at least in his own conscience. It does not follow from *this*, of course, that this is a risk one will never want to take.

#### Humanity and Universal Law

If the foregoing casuistical analyses are correct, then applying the Formula of Universal Law and applying the Formula of Humanity lead to rather different answers in the case of lying to the murderer at the door. The former seems to say that this lie is permissible, but the latter says that coercion and deception are the most fundamental forms of wrongdoing. In a Kingdom of Ends coercive and deceptive methods can never be used.

This result impugns Kant's belief that the formulas are equivalent. But it is not necessary to conclude that the formulas flatly say different things, and are unrelated except for a wide range of coincidence in their results. For one thing, lying to the murderer at the door was not shown to be permissible in a straightforward manner: the maxim did not so much pass as evade universalization. For another, the two formulas can be shown to be expressions of the same basic theory of justification. Suppose that your maxim is in violation of the Formula of Universal Law. You are making an exception of yourself, doing something that everyone in your circumstances could not

do. What this means is that you are treating the reason *you* have for the action as if it were stronger, had more justifying force, than anyone else's exactly similar reason. You are then acting as if the fact that it was in particular *your* reason, and not just the reason of a human being, gave it special weight and force. This is an obvious violation of the idea that it is your humanity - your power of rational choice - which is the condition of all value and so which gives your needs and desires the justifying force of *reasons*. Thus, any violation of the Formula of Universal Law is also a violation of the Formula of Humanity. This argument, of course, only goes in one direction: it does not show that the two formulas are equivalent. The Formula of Humanity is more strict than the Formula of Universal Law - but both are expressions of the same basic theory of value: that your rational nature is the source of justifying power of your reasons, and so of the goodness of your ends.

And although the Formula of Humanity gives us reason to think that all lies are wrong, we can still give an account in the terms it provides of what vindicates lying to a liar. The liar tries to use your reason as a means - your honesty as a tool. You do not have to passively submit to being used as a means. In the *Lectures on Ethics*, this is the line that Kant in fact takes. He says:

if we were to be at all times punctiliously truthful we might often become victims of the wickedness of others who were ready to abuse our truthfulness. If all men were well-intentioned it would not only be a duty not to lie, but no one would do so because there would be no point in it. But as men are malicious, it cannot be denied that to be punctiliously truthful is often dangerous... if I cannot save myself by maintaining silence, then my lie is a weapon of defense. (LE 228)

The common thought that lying to a liar is a form of self-defense, that you can resist lies with lies as you can resist force with force, is according to this analysis correct.<sup>13</sup> This should not be surprising, for we have seen that deception and coercion are parallel. Lying and the use of force are attempts to undercut the two conditions of possible assent to actions and of autonomous choice of ends, namely, knowledge and power. So, although the Formula of Universal Law and the Formula of Humanity give us different results, this does not show that they simply express different moral outlooks. The relation between them is more complex than that.

### Two Casuistical Problems

Before I discuss this relation, however, I must take up two casuistical problems arising from the view I have presented so far. First, I have argued that we *may* lie to the murderer at the door. But most people think something stronger, that we ought to lie to the murderer - that we will have done something wrong if we do not. Second, I have argued that it is permissible to lie to a deceiver in order to counter the deception. But what if someone lies to you for a good end, and, as it happens, you know about it? The fact that the murderer's *end* is evil has played no direct role in the arguments I have given so far. We have a right to resist liars and those who try to use force because of their methods, not because of their purposes. In one respect this is a virtue of my argument. It does not license us to lie to or to use violence against persons *just* because we think their purposes are bad. But it looks as if it may license us to lie to liars whose purposes are good. Here is a case<sup>14</sup>: suppose someone comes to your door and pretends to be taking a survey of some sort. In fact, this person is a philanthropist who wants to give his money to people who meet certain criteria, and this is his way of discovering appropriate objects for his beneficence. As it happens, you know what is up. By lying, you could get some money, although you do not in fact meet his criteria. The argument that I derived from the Formula of Universal Law about lying to the



murderer applies here. Universalizing the lie to the philanthropist will not destroy its efficacy. Even if it is a universal law that everyone will lie in these circumstances, the philanthropist thinks you do not know you are in these circumstances. By my argument, it is permissible to lie in this case. The philanthropist, like the murderer, has placed himself in a morally unprotected position by his own deception.

Start with the first casuistical problem. There are two reasons to lie to the murderer at the door. First, we have a duty of mutual aid. This is an imperfect duty of virtue, since the law does not say exactly what or how much we must do along these lines. This duty gives us *a* reason to tell the lie. Whether it makes the lie imperative depends on how one understands the duty of mutual aid, on how one understands the "wideness" of imperfect duties.<sup>15</sup> It may be that on such an urgent occasion, the lie is imperative. Notice that if the lie were impermissible, this duty would have no force. Imperfect duties are always secondary to perfect ones. But if the lie is permissible, this duty will provide a reason, whether or not an imperative one, to tell the lie.

The second reason is one of self-respect. The murderer wants to make you a tool of evil; he regards your integrity as a useful sort of predictability. He is trying to use you, and your good will, as a means to an evil end. You owe it to humanity in your own person not to allow your honesty to be used as a resource for evil. I think this would be a perfect duty of virtue; Kant does not say this specifically but in his discussion of servility (the avoidance of which is a perfect duty of virtue) he says "Do not suffer your rights to be trampled underfoot by others with impunity." (MMV 436/99)

Both of these reasons spring from duties of virtue. A person with a good character will tell the lie. Not to tell it is morally bad. But there is no duty of justice to tell the lie. If we do not tell it, we cannot be punished, or, say, treated as an accessory to the murder. Kant would insist that even if the lie ought to be told this does not mean that the punctiliously truthful person who does not tell it is somehow implicated in the murder. It is the murderer, not the truthful person, who commits this crime.

Telling the truth cannot be part of the crime. On Kant's view, persons are not supposed to be responsible for managing each other's conduct. If the lie were a duty of justice, we would be responsible for that.

These reflections will help us to think about the second casuistical problem, the lie to the philanthropist. I think it does follow from the line of argument I have taken that the lie cannot be shown to be impermissible. Although the philanthropist can hardly be called evil, he is doing something tricky and underhanded, which Kant's view disapproves. He should not use this method of getting the information he wants. This is especially true if the reason he does not use a more straightforward method is that he assumes that if he does people will lie to him. We are not supposed to base our actions on the assumption that other people will behave badly. Assuming this does not occur in an institutional context, and you have not sworn that your remarks were true<sup>16</sup>, the philanthropist will have no recourse to justice if you lie to him. But the reasons that favor telling the lie that exist in the first case do not exist here. According to Kant, you do not have a duty to promote your own happiness. Nor would anyone perform such an action out of self-respect. This is, in a very trivial way, a case of dealing with evil. But you can best deal with it by telling the philanthropist that you know what he is up to, perhaps even that you find it sneaky. This is *because* the ideal that makes his action a bad one is an ideal of straightforwardness in human relations. This would also be the best way to deal with the murderer, if it *were* a way to deal with a murderer. But of course it is not.

### Ideal and Non-Ideal Theory

I now turn to the question of what structure an ethical theory must have in order to accommodate this way of thinking. In *A Theory of Justice*,<sup>17</sup> John Rawls proposes a division of moral philosophy into ideal and non-ideal theory. In that work, the task of ideal theory is to determine "what a perfectly just society would be like,"

while non-ideal theory deals with punishment, war, opposition to unjust regimes, and compensatory justice. (§2,p. 8-9) Since I wish to use this feature of Rawls's theory for a model, I am going to sketch his strategy for what I will call a double-level theory.

Rawls identifies two conceptions of justice, which he calls the general conception and the special conception. (§§11,26,39,46) The general conception tells us that all goods distributed by society, including liberty and opportunity, are to be distributed equally unless an unequal distribution is to the advantage of everyone, and especially those who fall on the low side of the inequality. (§13) Injustice, according to the general conception, occurs whenever there are inequalities that are not to the benefit of everyone.(§11, p. 62) The special conception in its most developed form removes liberty and opportunity from the scope of this principle and says they must be distributed equally, forbidding tradeoffs of these goods for economic gains. It also introduces a number of priority rules, for example, the priority of liberty over all other considerations, and the priority of equal opportunity over economic considerations. (§§ 11,46,82)

Ideal theory is worked out under certain assumptions. One is strict compliance: it is assumed that everyone will act justly. The other, a little harder to specify, is that historical, economic, and natural conditions are such that realization of the ideal is feasible. Our conduct towards those who do not comply, or in circumstances which make the immediate realization of a just state of affairs impossible, is governed by the principles of non-ideal theory. Certain ongoing natural conditions which may always prevent the full realization of the ideal state of affairs also belong to non-ideal theory: the problems of dealing with the seriously ill or mentally disturbed, for instance, belong in this category. For purposes of constructing ideal theory, we assume that everyone is "rational and able to manage their own affairs." (§39, p. 248) We also assume in ideal theory that there are no massive historic injustices, such as the oppression of blacks and women, to be corrected. The point is to work out our ideal view of justice on the

assumption that people, nature, and history will behave themselves so that the ideal can be realized, and then to determine - in light of that ideal - what is to be done in actual circumstances, when they do not. The special conception is not applied without regard to circumstances. Special principles will be used in non-ideal conditions.

Non-ideal conditions exist when, or to the extent that, the special conception of justice cannot be realized effectively. In these circumstances our conduct is to be determined in the following way: the special conception becomes a goal, rather than an ideal to live up to: we are to work towards the conditions in which it is feasible. For instance, suppose there is a case like this: widespread poverty or ignorance due to the level of economic development is such that the legal establishment of the equal liberties makes no real difference to lot of the disadvantaged members of society. It's an empty formality. On the other hand, some inequality, temporarily instituted, would actually tend to foster conditions in which equal liberty could become a reality for everyone. In these circumstances, Rawls's double-level theory allows for the temporary inequality. (§§ 11,39) The priority rules give us guidance as to which features of the special conception are most urgent. These are the ones that we should be striving to achieve as soon as possible. For example, if formal equal opportunity for blacks and women is ineffective, affirmative action measures may be in order. If some people claim that this causes inefficiency at first, it is neither here nor there, since equality of opportunity has priority over efficiency. The special conception may also tell us which of our non-ideal options is least bad, closest to ideal conduct. For instance, civil disobedience is better than a resort to violence not only because violence is bad in itself, but because of the way in which civil disobedience expresses the democratic principles of the just society it aspires to bring about. (§ 59) Finally, the general conception of justice commands categorically. In sufficiently bad circumstances none of the characteristic features of the special conception may be realizable. But there is no excuse, *ever*, for violation of the general conception. If inequalities are not benefiting those on the lower end of them

in some way, they are simply oppression. The general conception, then, represents the point at which justice becomes uncompromising.<sup>18</sup>

A double-level theory can be contrasted to two types of single-level theory, both of which in a sense fail to distinguish the way we should behave in ideal and in non-ideal conditions, but which are at opposite extremes. A consequentialist theory such as utilitarianism does not really distinguish ideal from non-ideal conditions. Of course, the utilitarian can see the difference between a state of affairs in which everyone can be made reasonably happy and a state of affairs in which the utilitarian choice must be for the "lesser of evils", but it is still really a matter of degree. In principle we do not know what counts as a state in which everyone is "as happy as possible" absolutely. Instead, the utilitarian wants to make everyone as happy as possible relative to the circumstances, and pursues this goal holds regardless of how friendly the circumstances are to human happiness. The difference is not between ideal and non-ideal states of affairs but simply between better and worse states of affairs.

Kant's theory as he understood it represents the other extreme of single-level theory. The standard of conduct he sets for us is designed for an ideal state of affairs: we are always to act as if we were living in a Kingdom of Ends, regardless of possible disastrous results. Kant is by no means dismissive towards the distressing problems caused by the evil conduct of other human beings and the unfriendliness of nature to human ideals, but his solution to these problems is different. He finds in them grounds for a morally motivated religious faith in God.<sup>19</sup> Our rational motive for belief in a moral author of the world derives from our rational need for grounds for hope that these problems will be resolved. Such an author would have designed the laws of nature so that, in ways that are not apparent to us, our moral actions and efforts do tend to further the realization of an actual Kingdom of Ends. With faith in God, we can trust that a Kingdom of Ends will be the consequence of our actions as well as the ideal that guides them.

In his *A Critique of Utilitarianism*<sup>20</sup>, Bernard Williams spells out some of the unfortunate consequences of what I am calling single-level theories. According to Williams, the consequentialist's commitment to doing whatever is necessary to secure the best outcome may lead to violations of what we would ordinarily think of as integrity. There is no kind of action that is so mean or so savage that it can *never* lead to a better outcome than the alternatives. A commitment to always securing the best outcome never allows you to say "bad consequences or not, this is not the sort of thing I do; I am not that sort of person." And no matter how mean or how savage the act required to secure the best outcome is, the utilitarian thinks that you will be irrational to regret that you did it, for you will have done what is in the straightforward sense the right thing.<sup>21</sup> A Kantian approach, by defining a determinate *ideal* of conduct to live up to rather than setting a *goal* of action to strive for, solves the problem about integrity, but with a high price. The advantage of the Kantian approach is the definite sphere of responsibility. Your share of the responsibility for the way the world is is well-defined and limited, and if you act as you ought, bad outcomes are not your responsibility. The trouble is that in cases such as that of the murderer at the door it seems grotesque simply to say that I have done my part by telling the truth and the bad results are not my responsibility.

The point of a double-level theory is to give us both a definite and well-defined sphere of responsibility for everyday life and some guidance, at least, about when we may or must take the responsibility of violating ideal standards. The common sense approach to this problem uses an intuitive quantitative measure: we depart from our ordinary rules and standards of conduct when the consequences of following them would be "very bad." This is unhelpful for two reasons. First, it leaves us on our own about determining *how* bad. Second, the attempt to justify it leads down a familiar consequentialist slippery slope: if very bad consequences justify a departure from ordinary norms, why do not slightly bad consequences justify such a departure? A

double-level theory substitutes something better than this rough quantitative measure. In Rawls's theory, for example, a departure from equal liberty cannot be justified by the fact that the consequences of liberty are "very bad" in terms of mere efficiency. This does not mean that an endless amount of inefficiency will be tolerated, because presumably at some point the inefficiency may interfere with the effectiveness of liberty. One might put the point this way: the measure of "very bad" is not entirely intuitive but rather, bad enough to interfere with the reality of liberty. Of course this is not an algorithmic criterion and cannot be applied without judgment, but it is not as inexact as a wholly intuitive quantitative measure, and, importantly, does not lead to a consequentialist slippery slope.

Another advantage of a double-level theory is the explanation it offers of the other phenomenon which Williams is concerned about: that of regret for doing a certain kind of action even if in the circumstances it was the "right" thing. A double-level theory offers an account of at least some of the occasions for this kind of regret. We will regret having to depart from the ideal standard of conduct, for we identify with this standard and think of our autonomy in terms of it. Regret for an action we would not do under ideal circumstances seems appropriate even if we have done what is clearly the right thing.<sup>22</sup>

### Kantian Non-Ideal Theory

Rawls's special conception of justice is a stricter version of the egalitarian idea embodied in his general conception. In the same way, it can be argued that the Formula of Universal Law and the Formula of Humanity are expressions of the same idea - that humanity is the source of value, and of the justifying force of reason. But the Formula of Humanity is stricter, and gives implausible answers when we are dealing with the misconduct of others and the recalcitrance of nature. This comparison gives rise to the idea of using the two formulas and the relation between them to construct a

Kantian double-level theory of individual morality, with the advantages of that sort of account. The Formulas of Humanity and the Kingdom of Ends will provide the ideals which govern our daily conduct. When dealing with evil circumstances we may depart from this ideal. In such cases, we can say that the Formula of Humanity is inapplicable because it is not designed for use when dealing with evil. But it can still guide our conduct. It defines the goal towards which we are working, and if we can generate priority rules we will know which features of it are most important. It gives us guidance about which of the measures we may take is the least objectionable.

Lying to deceivers is not the only case in which the Formula of Humanity seems to set us a more ideal standard than the Formula of Universal Law. The arguments made about lying can all be made about the use of coercion to deal with evil-doers. Another, very difficult case in which the two formulas give different results, as I think, is the case of suicide. Kant gives an argument against suicide under the Formula of Universal Law, but that argument does not work.<sup>23</sup> Yet under the Formula of Humanity we can give a clear and compelling argument against suicide: nothing is of any value unless the human person is so, and it is a great crime, as well as a kind of incoherence, to act in a way that denies and eradicates the source of all value. Thus it might be possible to say that suicide is wrong from an ideal point of view, though justifiable in circumstances of very great natural or moral evil.

There is also another, rather different sense of "rigorism" in which the Formula of Humanity seems to be more rigorous than that of Universal Law. It concerns the question whether Kant's theory allows for the category of merely permissible ends and actions, or whether we must always be doing something that is morally worthy: that is, whether we should *always* pursue the obligatory ends of our own perfection and the happiness of others, when no other duty is in the case.

The Formula of Universal Law clearly allows for the category of the permissible. Indeed, the first contradiction test is a test of permissibility. But in the *Metaphysical*



*Principles of Virtue*, there are passages which have sometimes been taken to imply that Kant holds the view that our conduct should always be informed by morally worthy ends. (MMV 390/48) The textual evidence is not decisive. But the tendency in Kant's thought is certainly there: for complete moral worth is only realized when our actions are not merely in accordance with duty but from duty, or, to say the same thing a different way, perfect autonomy is only realized when our actions and ends are completely determined by reason, and this seems to be the case only when our ends are chosen as instantiations of the obligatory ends.

Using the Formula of Humanity it is possible to argue for the more "rigorous" interpretation. First, the obligatory ends can be derived more straightforwardly from Humanity than from Universal Law. Kant does derive the obligatory ends from the Formula of Universal Law, but he does it by a curiously round-about procedure in which someone is imagined formulating a maxim of rejecting them and then finding it to be impermissible. This argument does not show that there would be a moral failing if the agent merely unthinkingly neglected rather than rejecting these ends. The point about the pervasiveness of these ends in the moral life is a more complicated one, one that follows from their adoption by this route: Among the obligatory ends is our own moral perfection. Pursuing ends that are determined by reason, rather than merely acceptable to it, cultivates one's moral perfection in the required way. (MMV 380-381/37-38; 444-447/108-111)

It is important to point out that even if this is the correct way to understand Kant's ideal theory, it does not imply that Kantian ethics commands a life of conventional moral "good deeds." The obligatory ends are one's own perfection and the happiness of others; to be governed by them is to choose instantiations of these larger categories as the aim of your vocation and other everyday activities. It is worth keeping in mind that natural perfection is a large category, including all the activities that cultivate body and mind. Kant's point is not to introduce a strenuous moralism but to

find a place for the values of perfectionism in his theory. But this perfectionism will be a part of ideal theory if the argument for it is based on the Formula of Humanity and cannot be derived from that of Universal Law. This seems to me to be a desirable outcome. People in stultifying economic or educational conditions cannot really be expected to devote all their spare time to the cultivation of perfectionist values. But they can be expected not to do what is impermissible, not to violate the Formula of Universal Law. Here again, the Formula of Humanity sheds light on the situation even if it is not directly applied: it tells us why it is morally as well as in other ways regrettable that people should be in such conditions.

### Conclusion

If the account that I have given is correct, the resources of a double-level theory may be available to the Kantian. The Formula of Humanity and its corollary, the vision of a Kingdom of Ends, provide an ideal to live up to in daily life as well as a long term political and moral goal for humanity. But it is not feasible always to live up to this ideal, and where the attempt to live up to it would make you a tool of evil, you should not do so. In evil circumstances, but only then, the Kingdom of Ends can become a goal to seek rather than an ideal to live up to, and this will provide us with some guidance. The Kantian priorities - of justice over the pursuit of obligatory ends, and of respect over benevolence - still help us to see what matters most. And even in the worst circumstances, there is always the Formula of Universal Law, telling us what we must in not in any case do. For whatever bad circumstances may drive us to do, we cannot possibly be justified in doing something which others in those same circumstances could not also do. The Formula of Universal Law provides the point at which morality becomes uncompromising.

Let me close with some reflections about the extent to which Kant himself might have agreed with this modification of his views. Throughout this paper, I have

portrayed Kant as an uncompromising idealist, and there is much to support this view. But in the historical and political writings, as well as in the *Lectures on Ethics*, we find a somewhat different attitude. This seems to me to be especially important: Kant believes that the Kingdom of Ends on earth, the highest political good, can only be realized in a condition of peace. (MMJ 354-355/127-129) But he does not think that this commits a nation to a simple pacifism that would make it the easy victim of its enemies. Instead, he draws up laws of war in which peace functions not as an uncompromising ideal to be lived up to in the present but as a long range goal which guides our conduct even when war is necessary. (PP 343-348/85-91; MMJ 343-351/114-125) If a Kantian can hold such a view for the conduct of nations, why not for that of individuals? If this is right, the task of Kantian moral philosophy is to draw up for individuals something analogous to Kant's laws of war: special principles to use when dealing with evil.

---

<sup>1</sup> This paper was delivered as the Randall Harris Lecture at Harvard in October, 1985. Versions of the paper have been presented at the University of Illinois at Urbana-Champaign, the University of Wisconsin at Milwaukee, the University of Michigan, and to the Seminar on Contemporary Social and Political Theory at Chicago. I owe a great deal to the discussions on these occasions. I want to thank the following people for their comments: Margaret Atherton, Charles Chastain, David Copp, Stephen Darwall, Michael Davis, Gerald Dworkin, Alan Gewirth, David Greenstone, John Koethe, Richard Kraut, Richard Strier, and Manley Thompson. And I owe special thanks to Peter Hylton and Andrews Reath for extensive and useful comments on the early written versions of the paper.

<sup>2</sup> Where I have cited or referred to any of Kant's works more than once in this paper I have inserted the reference into the text. The following abbreviations are used:

**G** *Foundations of the Metaphysics of Morals.*(1785) The first page number is that of the Prussian Academy Edition Volume IV; the second is that of the translation by Lewis White Beck. Indianapolis: Bobbs-Merrill Library of Liberal Arts, 1959.

**C2** *Critique of Practical Reason.* (1788) Prussian Academy Volume V; Lewis White Beck's translation. Indianapolis: Bobbs-Merrill Library of Liberal Arts, 1956.

**MMV** *The Metaphysical Principles of Virtue.* (1797) Prussian Academy Volume VI; James Ellington's translation in *Immanuel Kant: Ethical Philosophy.* Indianapolis: Hackett, 1983.

---

**MMJ** *The Metaphysical Elements of Justice.* (1797) Prussian Academy Volume VI; John Ladd's translation. Indianapolis: Bobbs-Merrill Library of Liberal Arts, 1965.

**PP** *Perpetual Peace.* (1795) Prussian Academy Volume VIII, translation by Lewis White Beck in *On History*, edited by Lewis White Beck. Indianapolis: Bobbs-Merrill Library of Liberal Arts, 1963.

**SRL** "On a Supposed Right to Lie from Altruistic Motives" (1797) Prussian Academy Volume VIII; translation by Lewis White Beck in *Immanuel Kant: Critique of Practical Reason and Other Writings in Moral Philosophy.* Chicago: University of Chicago Press, 1949; rpt: New York: Garland Publishing Company, 1976.

**LE** *Lectures on Ethics.*(1775-1780) edited by Paul Menzer from the notes of Theodor Friedrich Brauer, using the notes of Gottlieb Kutzner and Chr. Mrongovius; translated by Louis Infield. London: Methuen & Co., Ltd., 1930; rpt: New York, Harper Torchbooks, 1963; current rpt: Indianapolis, Hackett Press.

<sup>3</sup> I defend it in "Kant's Formula of Universal Law", forthcoming in *Pacific Philosophical Quarterly.*

<sup>4</sup> I am relying on an assumption here, which is that when people ask us questions they give us some account of themselves and of the context in which the questions are asked. Or, if they don't, it is because they are relying on a context that is assumed. If someone comes to your door looking for someone, you assume that there's a family emergency or some such thing. I am prepared to count such reliance as deception if the questioner knows about it and uses it, thinking that we would refuse to answer his questions if we knew the real context to be otherwise. Sometimes people ask me, "Suppose the murderer just asks whether his friend is in your house, without saying anything about why he wants to know?" I think that, in our culture anyway, people do not *just ask* questions of each other about anything except the time of day and directions for getting places. After all, the reason why refusal to answer is an

---

unsatisfactory way of dealing with this case is that it will almost inevitably give rise to suspicion of the truth, and this is because people normally answer such questions. Perhaps if we did live in a culture in which people regularly *just asked* questions in the way suggested, refusal to answer would be commonplace and would not give rise to suspicion; it would not even be considered odd or rude. Otherwise there would be no way to maintain privacy.

<sup>5</sup> In fact, it will now be the case that if the murderer supposes that you suspect him, he is not going to ask you, knowing that you will answer so as to deceive him. Since we must avoid the silly problem about the murderer being able to deduce the truth from his knowledge that you will speak falsely, what you announce is that you will say whatever is necessary in order to conceal the truth. There is no reason to suppose that you will be mechanical about this. You are not going to be a reliable source of information. The murderer will therefore seek some other way to locate his victim.

On the other hand, suppose that the murderer does, contrary to my supposition, announce his real intentions. Then the arguments that I have given do not apply. In this case, I believe, your only recourse is refusal to answer (whether or not the victim is in your house, or you know his whereabouts). If an answer is extorted from you by force you may lie, according to the argument I will give later in the paper.

<sup>6</sup> Kant himself takes notice of this sort of problem in a footnote to this passage in which he criticizes Golden-Rule type principles for, among other things, the sort of subjectivity in question: such principles cannot establish the duty of beneficence, for instance, because "many a man would gladly consent that others should not benefit him, provided only that he might be excused from showing benevolence to them." (G 430n/48n)

<sup>7</sup> Sometimes it is objected that someone could assent to being lied to in advance of the actual occasion of the lie, and that in such a case the deception might still succeed.

---

One can therefore agree to be deceived. I think it depends what circumstances are envisioned. I can certainly agree to remain uninformed about something, but this is not the same as agreeing to be deceived. I could say to a doctor: "don't tell me if I am fatally ill, even if I ask" for instance. But if I then do ask the doctor whether I am fatally ill, I cannot be certain whether she will answer me truly. Perhaps what's being envisioned is that I simply agree to be lied to, but not about anything in particular. Will I then trust the person with whom I've made this odd agreement?

<sup>8</sup> A similar conclusion about the way in which the Formula of Humanity makes coercion and deception wrong is reached by Onora O'Neill in "Between Consenting Adults," *Philosophy and Public Affairs* Volume 14, No. 3 (Summer, 1985), pp. 252-277.

<sup>9</sup> *Immanuel Kant's Critique of Pure Reason*, translated by Norman Kemp Smith. (New York: St. Martin's Press, 1965) A738-739/B766-767, p. 593.

<sup>10</sup> It is perhaps also relevant that in Kant's discussion of perfect moral friendship the emphasis is not on good will towards one another but on complete confidence and openness. See MMV 471-472/139-139.

<sup>11</sup> Some evidence that Kant is concerned with this sort of thing may be found in the fact that he identifies two meanings of the word "prudence" (Klugheit); "The former sense means the skill of a man in having an influence on others so as to use them for his own purposes. The latter is the ability to unite all these purposes to his own lasting advantage." (G 416n/33n) A similar remark is found in *Anthropology from a Pragmatic Point of View*. (1798) See the translation by Mary J. Gregor (The Hague: Martinus Nijhoff, 1974) p. 183. Prussian Academy Edition Volume VII, p.322.

<sup>12</sup> I call this view Socratic because of Socrates's concern with the differences between reason and persuasion and, in particular, because in the *Apology*, he makes a case for the categorical duty of straightforwardness. Socrates and Plato are also concerned with a troublesome feature of this moral view that Kant neglects. An argument must come

---

packaged in some sort of presentation, and one may well object that it is impossible to make a straightforward presentation of a case to someone who is close to or admires you, without emphasis, without style, without taking some sort of advantage of whatever it is about you that has your listener's attention in the first place. So how can we avoid the non-rational influence of others? I take it that most obviously in the *Symposium*, but also in other dialogues concerned with the relation of love and teaching such as the *Phaedrus*, Plato is at work on the question whether you can use your sex appeal to draw another's attention to the reasons he has for believing or doing things, rather than as a distraction that aids your case illicitly.

<sup>13</sup> Of course you may also resist force with lies, if resisting it with force is not an option for you. This gives rise to a question about whether these options are on a footing with each other. In many cases, lying will be the better option. This is because when you use coercion you risk doing injury to the person you coerce. Injuring people unnecessarily is wrong, a wrong that should be distinguished from the use of coercion. When you lie you do not risk doing this extra wrong. But Kant thinks that lying is in itself worse than coercion, because of the peculiarly direct way in which it violates autonomy. So it should follow that if you can deal with the murderer by coercion, this is a *better* option than lying. Others seem to share this intuition. Cardinal John Henry Newman, responding to Samuel Johnson's claim that he would lie to a murderer who asked which way his victim had gone, suggests that the appropriate thing to do is "to knock the man down, and to call out for the police." (*Apologia Pro Vita Sua: Being a History of His Religious Opinions*. (London: Longmans, Green & Co., 1880) p. 361. I am quoting from Sissela Bok, *Lying*. (New York: Vintage Books, 1979) p 42.) If you can do it without seriously hurting the murderer, it is, so to speak, cleaner just to kick him off the front porch than to lie. This treats the *murderer himself* more like a human being than lying to him does.



---

<sup>14</sup> I owe this example to John Koethe.

<sup>15</sup> For a discussion of this question see Barbara Herman, "Mutual Aid and Respect for Persons" *Ethics* 94 (July 1984) pp. 577-602.

<sup>16</sup> In the *Lectures on Ethics*, Kant takes the position that you may lie to someone who lies to or bullies you as long as you don't say specifically that your words will be true. He claims this is not lying, because such a person should not expect you to tell the truth. (LE 227,229)

<sup>17</sup> John Rawls, *A Theory of Justice*. Cambridge, Massachusetts: Harvard University Press, 1971. Section and page numbers referring to this work will appear in the text.

<sup>18</sup> In a non-ideal case, one's actions may be guided by a more instrumental style of reasoning than in ideal theory. But non-ideal theory is not a form of consequentialism. There are two reasons for this. One is that the goal set by the ideal is not just one of good consequences, but of a just state of affairs. If a consequentialist view is one that defines right action entirely in terms of good consequences (which are not themselves defined in terms of considerations of rightness or justice) then non-ideal theory is not consequentialist. The second reason is that the ideal will also guide our choice among non-ideal alternatives, importing criteria for this choice other than effectiveness. I would like to thank Alan Gewirth for prompting me to clarify my thoughts on this matter, and David Greenstone for helping me to do so.

<sup>19</sup> See the "Dialectic of Pure Practical Reason" of the *Critique of Practical Reason*, and the *Critique of Teleological Judgment*, §87.

<sup>20</sup> Bernard Williams, in *Utilitarianism For and Against*, by J.J.C. Smart and Bernard Williams (Cambridge: Cambridge University Press, 1973), pp. 75-150.

<sup>21</sup> Williams also takes this issue up in "Ethical Consistency" originally published in the Supplementary Volumes to the *Proceedings of the Aristotelian Society* XXXIX, 1965, and

---

reprinted in his collection *Problems of the Self* (Cambridge: Cambridge University Press, 1973), pp. 166-186.

<sup>22</sup> It is important here to distinguish two kinds of exceptions. As Rawls points out in "Two Conceptions of Rules" (*The Philosophical Review*, Volume 64 (January 1965)), a practice such as promising may have certain exceptions built into it. Everyone who has learned the practice understands that the obligation to keep the promise is cancelled if one of these obtains. When one breaks a promise because this sort of exception obtains, regret would be inappropriate and obsessive. And these sorts of exceptions may occur even in "ideal" circumstances. The kind of exception one makes when dealing with evil should be distinguished from exceptions built into practices.

<sup>23</sup> Kant's argument depends on a teleological claim: that the instinct whose office is to impel the improvement of life cannot universally be used to destroy life without contradiction. (G 422/40) But as I understand the contradiction in conception test, teleological claims have no real place in it. What matters is not whether nature assigns a certain purpose to a certain motive or instinct, but whether everyone with the same motive or instinct could act in the way proposed and still achieve their purpose. There is simply no argument to show that everyone suffering from acute misery could not commit suicide and still achieve the purpose of ending that misery.