

Two Wrongs Don't Make a Right: Responsibility and Overdetermination*

Carolina Sartorio

University of Arizona

(sartorio@arizona.edu)

I. Introduction

In his impressive book, *Causation and Responsibility*, Michael Moore examines the relation between metaphysics, morality, and the law (Moore (2009)). He advocates a picture according to which facts about legal responsibility are grounded in facts about moral responsibility, which in turn are grounded in the relevant metaphysical facts, such as, notably, facts about causation. In this paper I will focus on the relation that Moore sees between the metaphysical facts and the relevant moral facts in a type of scenario usually referred to as “overdetermination.” Overdetermination is an intriguing phenomenon that has puzzled legal theorists and philosophers alike. I will examine two forms of overdetermination discussed by Moore:

overdetermination by the actions of agents or by positive occurrences of some kind and overdetermination by the omissions of agents or by absences of some kind.

Moore argues for the existence of a moral asymmetry between these two forms of overdetermination: briefly, agents can be morally responsible for the outcome that ensues in scenarios of the first kind but not in scenarios of the second kind. This moral asymmetry, he argues, is rooted in a causal asymmetry: overdetermining actions are causes of the ensuing outcomes, but overdetermining omissions are not.

* For helpful comments thanks to Sara Bernstein, Randolph Clarke, Juan Comesaña, and Marc Johansen.

In this paper I will challenge the moral asymmetry defended by Moore as well as his claim that the moral asymmetry is grounded in a causal asymmetry. The question about how to assign moral responsibility to agents in scenarios of overdetermination and how to ground those attributions of responsibility in concepts like causation or other metaphysical concepts raises some interesting challenges, as we will see. I will put forth some general principles about responsibility that I think are operative in overdetermination cases and that will hopefully help to shed some light on this interesting phenomenon.

II. Moore on Overdetermination

As I noted, Moore believes that certain facts about moral responsibility are grounded in the relevant metaphysical facts, in particular, facts about causation. Does this mean that one should hope to settle certain moral questions by first settling the relevant metaphysical questions? As we will see, Moore argues for the soundness of certain “metaphysical-to-moral inferences” that apply to scenarios of overdetermination. I will examine the legitimacy of those inferences in due course. But first it will be helpful to start by laying out some more general metaphysical-to-moral inferences embraced by Moore.

In chapter 2 of his book, Moore argues that causation is relevant to moral responsibility in that its presence can increase the blameworthiness of an already blameworthy agent. If, for example, a reckless driver kills a pedestrian as a result of his reckless driving, the fact that his behavior causally resulted in harm renders him more blameworthy than he would have been if he hadn't caused any harm. In other

words, Moore believes in the existence of resultant moral luck: something that is typically beyond the agent's control—namely, whether a certain behavior by an agent causally results in harm—can determine the degree of the agent's blameworthiness.

Now, I am interested in the claim that the agent's responsibility *extends* to an outcome when the agent's behavior actually causes the outcome. (By "extends" I mean that the agent is responsible for the outcome that causally resulted from his behavior in addition to being responsible for the behavior itself, and, moreover, he is responsible for the outcome in virtue of being responsible for the behavior that resulted in it.) Arguably, there is a substantial step from this claim to the claim that this results in a genuine increase in the agent's blameworthiness. Some philosophers have argued that, although the reckless driver who kills someone as a result of his reckless driving is to blame *for* more things than he would have been responsible for otherwise, he is not more to blame than he would have been otherwise. In other words, what the agent is responsible for, or the scope of his blameworthiness, is different, but how responsible he is, or the degree of his blameworthiness, is the same (see, e.g., Thomson (1989) and Zimmerman (2002)). In this way one could grant that causation is relevant to an agent's blameworthiness without committing oneself to the possibility of resultant moral luck concerning the degree of the agent's responsibility. Here I will be interested in the claim that causation can result in an extension of an agent's responsibility to outcomes, regardless of whether or not this results in an increase in the agent's blameworthiness. I take it that this is a relatively uncontroversial claim (it is

certainly less controversial than the claim about the corresponding increase in the agent's blameworthiness). This claim can be expressed in the form of a pair of metaphysical-to-moral inferences:

M-M 1:

- (1.1) If an agent's action for which the agent is responsible caused a given outcome (and other relevant conditions obtained), then the agent's responsibility extends to the outcome.
- (1.2) If an agent's action for which the agent is responsible didn't cause a given outcome, then the agent's responsibility doesn't extend to the outcome.

The "other relevant conditions" referred to in (1.1) are the conditions that also have to obtain (in addition to the causal condition) for the agent to be morally responsible for the outcome. Presumably, one such condition is that the agent was able to foresee (or should have been able to foresee) that his action would likely cause the outcome. But there could be others.

In Moore's view, the M-M 1 inferences don't apply to agents' omissions. This is because Moore believes that omissions cannot be causes. (He believes this because he believes that omissions are mere absences and he believes, in turn, that absences cannot be causes. See, especially, chapter 18.) Still, Moore thinks that some omissions can make us morally responsible for outcomes such as harms. For example, a lifeguard can be responsible for a swimmer's death in virtue of having

omitted to jump into the water to try to save him. Hence the absence of causation doesn't entail the absence of responsibility in the case of agents' omissions.

Our omissions can make us morally responsible for harms, in Moore's view, because there can be *counterfactual dependence* between our omissions and those harms (see, again, chapter 18). Counterfactual dependence is the relation that obtains between X and Y whenever, had X not occurred, Y wouldn't have occurred either. Thus an outcome counterfactually depends on what an agent omits to do (say, the agent's omitting to perform an act A) just in case, had the agent failed to omit to A (i.e. had he performed act A), the outcome wouldn't have occurred. For example, a swimmer's death counterfactually depends on the lifeguard's failure to attempt a rescue when, had the lifeguard attempted a rescue, he would have saved the swimmer. The lifeguard can be responsible for the swimmer's death to the extent that the death counterfactually depended on his omission to attempt a rescue.

Causation and counterfactual dependence are, in Moore's view, the only two metaphysical bases for (non-inchoate)¹ responsibility, that is to say, they are the only two metaphysical relations that can ground an agent's moral responsibility for a given outcome in the world. Hence, in the case of omissions, Moore would endorse the following pair of metaphysical-to-moral inferences:

M-M 2:

¹ Non-inchoate responsibility is opposed to responsibility for mere attempts.

- (2.1) If a given outcome counterfactually depended on an agent's omission for which the agent is responsible (and other relevant conditions obtained), then the agent's responsibility extends to the outcome.
- (2.2) If a given outcome didn't counterfactually depend on an agent's omission for which the agent is responsible, then the agent's responsibility doesn't extend to the outcome.

The second inference in M-M 2, claim (2.2), has important implications in cases of overdetermination involving omissions. Overdetermination is typically defined in the following way: an outcome is overdetermined by two events if the outcome wouldn't have occurred in the absence of both events but each event was independently sufficient for the outcome's occurrence, i.e. the outcome would still have occurred if either of those events had occurred without the other (Lewis (1986), Postscript E). More particularly, our focus here is a specific kind of overdetermination sometimes called *symmetrical* overdetermination. There is symmetrical overdetermination when the overdeterminers' contributions are on a par.² A classical example of symmetrical overdetermination is two events of shooting at a victim through the heart at exactly the same time. Symmetrical overdetermination is contrasted with another kind of overdetermination called "asymmetrical" or "preemptive" overdetermination, in which there is an asymmetry

² Alternatively, some people say that there is symmetrical overdetermination when the overdeterminers are both *causes*. Lewis doesn't go this route because he believes that overdeterminers don't cause the subsequent effects (taken individually, that is; they do "collectively", since the mereological sum of the overdeterminers causes the effect). Here I will work with Lewis' formulation, for obvious reasons (Moore believes that omissions are never causes, so omissions could never give rise to overdetermination on this alternative view).

between the contributions of the overdeterminers: one of them, the preempting event, is a cause and the other one, the preempted event, isn't. Following common usage, here I will use just "overdetermination" to mean symmetrical overdetermination.

As we will see in section IV, this definition of overdetermination is not fully accurate, in particular, the counterfactual account of the relevant notion of sufficiency is flawed. But it is at least a rough approximation to the phenomenon of overdetermination and it will have to do for the time being. Note that the definition focuses on overdetermination by two events. Overdetermination by two absences, say, two omissions, would be defined in a similar way: an outcome is overdetermined by two omissions when it wouldn't have occurred in the absence of both omissions and each omission was independently sufficient for its occurrence. (Perhaps one should add a symmetry condition here too stating that the contributions of the two omissions are on a par. Note, however, that in this case one couldn't understand this condition as making reference to the *causal* contributions of the omissions, since on a view like Moore's omissions are never causally efficacious. Hence one would have to find an alternative way of drawing the distinction between the symmetrical cases and the asymmetrical cases. Since here I will only be concerned with the symmetrical cases, I will basically sidestep this delicate issue in what follows.)³

An example of overdetermination by two omissions is the following scenario. Two buttons have to be depressed at a given time in order to prevent an explosion.

³ I will revisit this issue briefly at the end of the paper. See section V below.

There is one person in charge of each of the buttons. However, when the time comes, both of the persons in charge fail to depress their buttons and the explosion ensues as a result. In this case the explosion is overdetermined by the two agents' omissions. I will call this scenario "Two Buttons" (the example is taken from Sartorio (2004)).

Note that a main feature of overdetermination scenarios is the absence of counterfactual dependence between each of the individual overdeterminers and the outcome. If two acts of shooting a victim through the heart overdetermine the victim's death, the death doesn't counterfactually depend on either of the individual overdeterminers (it would still have occurred if one of the shootings had occurred without the other). Similarly, in Two Buttons the explosion doesn't counterfactually depend on the agents' individual failures to depress their buttons (it would still have occurred if one failure had occurred without the other). Hence, even if the other conditions for responsibility are met (say, both agents had reason to believe that they could prevent the explosion by depressing their buttons, they didn't do it because they wanted the explosion to happen, etc.), it follows from claim (2.2) above that neither agent is morally responsible for the explosion in Two Buttons.

This result doesn't extend to scenarios of overdetermination involving actions. For Moore believes that any overdetermining positive occurrences are causes of the ensuing outcome. Hence M-M 1 (in particular, claim (1.1)) entails that agents are morally responsible for the outcome in virtue of their overdetermining actions, in the right circumstances. For example, if depressing one button is enough to send an electrical shock that will kill a person and two agents depress their

buttons simultaneously, (1.1) entails that, in the right circumstances, both agents are responsible for the victim's death. This results in an asymmetry between scenarios of overdetermination involving actions and those involving omissions, an asymmetry that concerns the moral responsibility of the agents involved in those scenarios, and one that is rooted in an asymmetry between the causal powers of actions and omissions.⁴

III. Two Wrongs Don't Make a Right—Part I

I will argue that we should reject the moral asymmetry advocated by Moore between the two forms of overdetermination. In particular, we should reject Moore's claim that agents in overdetermination scenarios of omission are not responsible for the overdetermined outcome. I will proceed as follows. First, in this section, I will argue that the claim that agents are not responsible for the outcome in overdetermination scenarios involving omissions is extremely counterintuitive and I will pinpoint, in rough terms, the source of the counterintuitiveness. Then, in the next section, I will attempt to make more precise the nature of that source.

Imagine that they tell you that depressing the button located right in front of you will prevent an explosion, but you fail to depress the button because you want the explosion to happen. Imagine that depressing the button would indeed have prevented the explosion. Clearly, you are responsible for the explosion in this case.

⁴ Moore discusses scenarios of symmetrically culpable co-omitters in chapters 5, 6, and 18 (see especially p. 450). Moore also discusses scenarios of asymmetrical culpability involving omissions and other absences (in particular, see his discussion of the famous "thirsty traveler" case on pp. 466-7). Although I believe that the main reason to reject Moore's view in overdetermination scenarios also applies to the asymmetrical cases, for ease of exposition I will only be concerned with cases of symmetrical culpability. Asymmetrical cases give rise to a range of puzzling new questions, which I prefer to bypass in this paper. (See n. 20 below.)

Now imagine that, unbeknownst to you, there was a second agent who was told the same thing as you (he was told that he could prevent the explosion by depressing the button located right in front of him) and imagine that he also failed to depress his button because he wanted the explosion to happen. Imagine, as in the Two Buttons scenario described above, that both buttons had to be depressed in order for the explosion to be prevented, so the explosion was overdetermined by the two omissions. Could the existence of this other agent and this other button relieve you of responsibility for the explosion? Intuitively, no. For, if it did, then *your* existence and the existence of your button would also relieve *him* of responsibility and then nobody would be responsible for the explosion in that case. However, how could the fact that more people failed to do what they were supposed to do result in no one's being responsible for the explosion? Two wrongs don't make a right.⁵

Here is a real-life example with a similar structure. Every time there is an election in the United States, many people who could have voted in the election fail to vote. The outcome of the election is always overdetermined, and almost always overdetermined by a large margin. Still, oftentimes it is the case that, if several of the people who failed to vote had voted for a certain candidate (a losing candidate), that candidate would have won the election. Just like we want to hold people who voted in an election responsible for the outcome, we want to hold people who didn't vote responsible for the outcome. But, given that the outcome is (in all realistic cases) overdetermined, it follows from M-M 2 (in particular, claim (2.2)) that none of the

⁵ Moore recognizes that many people find counterintuitive to say that agents are not responsible for the outcome in cases like this, but he thinks that such an intuition can be explained away as an illicit kind of "moral clumping" phenomenon (see p. 450). I will argue that, on the contrary, the intuition can be adequately justified.

people who failed to vote are responsible for the outcome of the election (although many of the people who did vote are presumably responsible for the outcome).

The election case is particularly interesting because, given that the outcome is usually overdetermined by such a large margin, the agents involved don't even have good reason to believe that they can make a difference to the outcome. (When the advertisements for a certain candidate try to convince you that you can "make a difference" by voting for a certain candidate, they are usually simply false or only metaphorically true. In other cases they are true but irrelevant, as when they read: "*Together we can make a difference!*") Still, we think that people can be responsible for the outcome of an election, by voting or by failing to vote. In the Two Buttons case I stipulated that both agents believed (falsely) that they could make a difference, since it is particularly clear that we would want to blame them for the outcome then. But, as the election case suggests, believing that you can make a difference might not actually be required for responsibility in every case of this kind.⁶

At any rate, it seems clear that the claim that agents are not responsible in overdetermination scenarios of omission is very counterintuitive, and that its counterintuitiveness can be captured by means of the slogan "two wrongs don't make a right" (hereafter, "The Slogan"). By "wrongs" I of course mean wrongs that are blameworthy. The slogan claims that two blameworthy wrongs cannot together constitute a blameless right (the mere fact that someone else acted wrongly and in a

⁶ The election case gives rise to a puzzle because it is unclear how we can assign responsibility to the agents if they have every reason to think that their individual vote cannot make a difference. But it seems clear that we want to hold agents responsible in scenarios of that kind at least in some cases; the problem is how.

way that deserves blame doesn't get a blameworthy wrongdoer off the hook).⁷ But this is obviously very rough; in particular, as we will see next, there are some ways in which two wrongs *can* make a right. So in what follows I proceed to clarify the sense in which The Slogan states something true and I will explain how The Slogan, understood in that way, applies to overdetermination cases, including cases of overdetermination involving omissions. This will take us through a brief detour into the nature of overdetermination.

IV. Two Wrongs Don't Make a Right—Part II

Imagine two shooters who simultaneously try to kill a victim but whose bullets collide in mid-flight and are deflected away from the victim. Call this case "Colliding Bullets." Clearly, the two shooters are not responsible for the victim's death in this case. The simplest scenario of this kind is one where the victim doesn't even die. But note that the two shooters would still escape responsibility for the death even if the death were to happen in some other way—say, if a third person had shot him a few seconds later. The reason the two shooters escape responsibility for the death in these cases is that their contributions interfere with each other in such a way that they are no longer sufficient for the death. If the death still happens, it is because a different causal process resulted in the death. Plainly, two wrongs *can* make a right in this way.

⁷ Here I will not be concerned with other forms of responsibility, such as praiseworthiness, since they are not the focus of Moore's work. It seems to me clear, however, that similar considerations apply to praiseworthiness. Although it's harder to formulate the Slogan for praiseworthy acts ("two rights don't make a wrong" plainly doesn't capture it), a similar thought holds: your act cannot cease to be praiseworthy just because others act in similarly praiseworthy ways.

Now, note that our target scenarios, overdetermination scenarios, are importantly different. Overdeterminers don't interfere with each other or cancel each other out in the same way they do in Colliding Bullets. No additional causal process is needed to make the outcome happen, since each overdeterminer is still sufficient for the outcome's occurrence, even in the presence of the other overdeterminer. On reflection, it seems unproblematic to suggest that two wrongs can make a right if, by virtue of the two wrongs happening together, they are not actually sufficient for the harm. However, it does seem problematic to suggest that two wrongs make a right when each of the two wrongs is actually sufficient for the outcome, even in the company of the other.

The claim that I want to defend, then, is not the claim that two wrongs can never make a right, but the claim that two wrongs cannot make a right unless certain special circumstances obtain. These special circumstances are ones in which two conditions that *would* have been sufficient for the occurrence of an outcome if they had occurred in isolation from each other are not *in fact* sufficient for the outcome when they occur simultaneously. Note, also, that my main reason for wanting to say that two wrongs cannot make a right (unless special circumstances obtain) is *not* that an agent's responsibility cannot depend on what others do, or on factors that are beyond the agent's control. In accepting that two wrongs can make a right in some (special) circumstances, I am granting that sometimes an agent's responsibility *can* be partly a matter of luck (for whether your bullet collides with

another person's bullet is not fully within your control).⁸ The claim I am making is that two wrongs cannot make a right unless luck intervenes in such a way that the individual contributions are not actually sufficient, in the relevant way, for the harm to occur.

Recall that, when I gave the standard definition of overdetermination in section II, I said that it could only be taken as a first approximation to the true nature of the phenomenon. We can now see why. Schematically, the definition was the following:

(OD) Two actual events C and D overdetermine another actual event E iff:

- (i) If neither C nor D had occurred, then E wouldn't have occurred.
- (ii) If C had occurred without D, or if D had occurred without C, then E would have occurred.
- (iii) C and D were on a par (their contributions to E were symmetrical).

Scenarios like Colliding Bullets can be used (with a minor twist) to show that OD is flawed. Imagine, again, that the victim's death still occurs after the two shooters' bullets collide, due to the existence of a third shooter. Imagine also (this is the minor twist on the example) that the third shooter wouldn't have been motivated to shoot the victim in any other case, in particular, he wouldn't have shot the victim if neither shooter had shot him. Call this case "Colliding Bullets 1." OD entails that the death is overdetermined by the two shooters in this case. For (given

⁸ Although, recall that (as I pointed out in section II) there might be a further step from there to the claim that this is genuine moral luck about the *degree* of the agent's blameworthiness.

the minor twist) the death wouldn't have occurred if the two shootings had not taken place, it would have occurred if only one of the shootings had occurred, and there is no asymmetry between the contributions of the two shootings. But, clearly, the two shootings don't overdetermine the death in this case.⁹

Consider, also, this other scenario. Imagine that this time there is no third shooter; the two shooters' bullets reach the victim after they collide with each other and kill him. However, imagine that the two bullets are slowed down considerably by the collision and, as a result, by the time they reach the victim both bullets are needed to kill him. Call this scenario "Colliding Bullets 2." Again, this is not a case of overdetermination. However, OD entails that the death is overdetermined by the two shooters in this case too. For, again, the death wouldn't have occurred in the absence of both shootings, it would have occurred if only one of the shootings had taken place, and there is no asymmetry between the contributions of the two shootings. The two colliding bullets scenarios, then, show that OD fails to offer sufficient conditions for the existence of overdetermination.¹⁰

What has gone wrong? What these examples suggest is that the key notion of sufficiency involved in overdetermination cases cannot be spelled out in simple counterfactual terms, as the claim that the occurrence of one of the overdeterminers

⁹ Note that, given that in this case the third shooter only shoots because the two shooters shot and missed, we might want to say that the two shootings still contributed to the death (they are contributing causes, even though they are not overdeterminers). If so, Colliding Bullets 1 isn't a scenario where two wrongs make a right, but one where two wrongs make something like a "lesser wrong." The same goes for Colliding Bullets 2, the next case discussed in the text.

¹⁰ Arguably, OD also fails to offer necessary conditions. Imagine a "shy" bullet that can only move if it senses other bullets moving around it. It seems that an ordinary bullet and a shy bullet can overdetermine a death when shot simultaneously although, had the shy bullet been shot without the ordinary bullet, the death wouldn't have occurred (condition (ii) fails). Thanks to Sara Bernstein for the example.

in the absence of the other would have resulted in the outcome's occurrence. The colliding bullets are sufficient for the victim's death in this sense, but they are not sufficient *in the sense that matters for overdetermination*.

What is the right account of overdetermination, then? It is hard to say, and this is not the place to try to give a better account.¹¹ Instead of attempting to do this, I will simply point at one possible direction in which one could go. It seems to me that one could try a more sophisticated counterfactual account that appealed to strategies sometimes used to give an account of the concept of causation itself. In particular, one could appeal to the strategy of uncovering certain dependencies that exist while holding certain actual facts fixed. This is a strategy employed by, e.g., Yablo (2002) in giving his counterfactual account of causation. Yablo argues that, although sometimes effects don't counterfactually depend on their causes (as in cases of preemption), in those cases there is still "de facto" dependence, or dependence *in the actual circumstances*: the effect depends on the cause holding fixed certain facts about potential routes to the effect (facts in virtue of which those other routes were closed or causally inefficacious). Similarly, one could argue that the relevant notion of sufficiency behind overdetermination is a kind of *de facto sufficiency*. Take the two bullets in Colliding Bullets 2. If either one of those two bullets hadn't been shot, the victim would still have died. However, if either one of

¹¹ Hall (2004) offers a different account that seems to avoid the problems raised by the colliding bullets scenario. On Hall's view, an overdetermination scenario is, roughly, a scenario that can be broken up into two or more intrinsic causal structures culminating in the effect. But it is unclear that this account will capture every case of overdetermination. For instance, we could imagine a colliding bullets scenario where each bullet contributes to the intrinsic properties of the causal process involving the other bullet (say, it helps to determine the other bullet's momentum after the two bullets collide) but where the resulting momentum of each bullet is still enough to kill the victim. This is a case of overdetermination but it would be hard to capture it by appeal to Hall's account.

those two bullets hadn't been shot *but somehow the other bullet had retained its actual momentum*, then the victim would not have died. In other words, neither bullet is sufficient for the victim's death, when one holds fixed its actual momentum. Similarly for Colliding Bullets 1: if either bullet hadn't been shot, then, holding fixed the fact that the other bullet was deflected away from its path, the death wouldn't have occurred. By contrast, in a genuine scenario of overdetermination, such as one where the two bullets pierce the victim's heart at exactly the same time without colliding with each other, the shooting of each bullet is sufficient for the death in the relevant sense (if either bullet hadn't been shot, then, holding fixed the actual momentum of the other bullet, the death would still have occurred).¹²

Note that an account of overdetermination along these lines (a "de facto" counterfactual account of overdetermination) is still in some important ways a counterfactual account, but it is a counterfactual account where the actual facts play an important role too (and not just by virtue of fixing the truth of the relevant counterfactuals). Setting aside the question of whether an account of this type could succeed, it seems to me that this is a main virtue of such an account. For an account of this type strives to capture an important insight about overdetermination brought out by the colliding bullets scenarios: the fact that overdetermination has to do in some important sense with the *actual* contributions of the overdeterminers. We feel that two colliding bullets cannot be overdeterminers if the *actual* momentum of each bullet is not sufficient to bring about the death by itself, even if

¹² Note that an account along these lines would also have the right result in the shy bullet case described in n. 10 (the scenario that seemed to show that OD fails to offer necessary conditions for overdetermination).

each one of them *would* have been sufficient to bring about the death by itself.

“Actual sufficiency”, not “counterfactual sufficiency”, is what seems to matter for overdetermination.

Now we may finally return to The Slogan (“two wrongs don’t make a right”). My suggestion is that we should read it as the claim that two wrongs don’t make a right when each of the two wrongs is sufficient for the occurrence of the ensuing harm, in the sense of sufficiency that matters. That is to say, two wrongs don’t make a right when each of the two wrongs is sufficient for the occurrence of the harm, not in the sense of counterfactual sufficiency, but in the sense of *actual* sufficiency.

What follows from this? One thing that follows is that, in the case of two blameworthy acts, where the actual contribution of each one of the two acts is sufficient for the occurrence of a harm, the two agents are responsible for the harm in virtue of their blameworthy acts. The different shooting scenarios we have examined in this section involved acts by two agents. But it is possible to apply The Slogan thus understood to omissions as well. In the case of omissions it is a bit harder to get a grip on the notion of actual sufficiency and the difference between actual sufficiency and counterfactual sufficiency. But the difference still seems to exist. Consider the following variant on the Two Buttons scenario discussed before. Again, if two buttons were to be depressed simultaneously, an explosion would be prevented. Again, each of the agents in charge intentionally fails to depress his button. However, imagine that, unbeknownst to them, an emergency mechanism has been set up to go off automatically if (and only if) *both* buttons fail to be depressed. So, when the two agents fail to depress their buttons, the emergency mechanism

goes off and prevents the explosion. This is another scenario where two wrongs make a right. And we can see the distinction between counterfactual and actual sufficiency here too. Although, in this case, each omission is counterfactually sufficient for the explosion (if one agent had omitted to depress his button but the other agent hadn't, then the explosion wouldn't have been prevented), neither omission is *actually* sufficient for the explosion. By contrast, in the original Two Buttons scenario the two omissions *are* actually sufficient for the explosion. Hence, whereas two wrongs do make a right in the scenario where the emergency mechanism goes off and prevents the explosion, they don't in Two Buttons.

I conclude that the right understanding of The Slogan undermines Moore's asymmetry between overdetermination scenarios of action and overdetermination scenarios of omission. For in scenarios of both kinds, as in all genuine cases of overdetermination, the two wrongs are each actually sufficient for the harm. And, whereas two wrongs can make a right when they are each merely counterfactually sufficient for the harm, they cannot make a right when they are each actually sufficient for the harm.¹³

¹³ There are some scenarios where it might appear, at first sight, that two *actually* sufficient wrongs make a right. But I think that these are cases where we don't have two wrongs to start with. Consider, for example, this variant on the Two Buttons case. The two agents, X and Y, find out about each other's existence. X is told by a reliable source that Y won't press his button, and Y is told by another reliable source that X won't press his button. As a matter of fact, however, both want the explosion to happen so they fail to press their buttons wishing it to happen. Arguably, neither X nor Y is responsible for the explosion in this case, but their omissions are actually sufficient for the explosion. However, I think that, to the extent that we think that the agents are not responsible for the explosion in this case, it's because we think that their omissions weren't blameworthy. Failing to do something that you have very good reason to believe won't do any good is not usually blameworthy, regardless of what your intentions regarding the outcome are. (If a child is drowning and a bystander has good reason to believe that throwing him a lifesaver won't be enough to save him, then the bystander is not blameworthy for the child's death in virtue of not throwing him the lifesaver, even if he wanted the child to die.)

V. In Search of a Responsibility Basis

Recall the metaphysical-to-moral inferences that Moore embraces in the case of omissions:

M-M 2:

- (2.1) If a given outcome counterfactually depended on an agent's omission for which the agent is responsible (and other relevant conditions obtained), then the agent's responsibility extends to the outcome.
- (2.2) If a given outcome didn't counterfactually depend on an agent's omission for which the agent is responsible, then the agent's responsibility doesn't extend to the outcome.

I have argued that we should reject M-M 2. In particular, I have argued that we should reject claim (2.2): when presented with cases of overdetermination by two omissions such as the Two Buttons case (where the outcome fails to depend counterfactually on the individual omissions), we shouldn't conclude that the agents involved are not responsible for the outcome. This means that, if our initial ideas about metaphysical bases for responsibility conflict with the claim that the agents are responsible for the outcome in these cases, we should not revise the claim about the agents' responsibility but, instead, our initial ideas about what constitutes an appropriate metaphysical basis for responsibility. In essence, this is what Moore thinks we should do in paradigm cases of responsibility by omission (where, he believes, causation fails to ground responsibility): we should take as a starting point

the assumption that the agent is responsible and then look for an appropriate relation (such as counterfactual dependence) to ground his responsibility. What I am suggesting is that here, too, we should reason in this way: we should take the assumption that the agents are responsible as a starting point and then look for an appropriate way to ground their responsibility.

Now, what could this basis for responsibility be? If one believed, unlike Moore, that causation by omission is possible, then one could of course try to argue that the relevant responsibility basis is causation. That is to say, one could try to claim that the agents in Two Buttons cause the explosion, although the explosion is not counterfactually dependent on their individual omissions. Now, I don't think that this is a promising strategy. For, even if causation by omission were possible, I think that it would be very implausible to suggest that the agents in Two Buttons cause the explosion. I have argued for this elsewhere (Sartorio (2004)), so here I will limit myself to a brief outline of the main argument. It is the following. Imagine a variant on the case, "Two Buttons-One Stuck," where one button is controlled by a person and the other button is controlled by an automated mechanism. On this particular occasion, where both buttons had to be depressed to prevent the explosion, the agent fails to depress his button and the mechanism also fails (the button controlled by the mechanism becomes stuck). Presumably, the agent's omission is not a cause of the explosion in this case. But, I argued, there is no metaphysically significant difference between Two Buttons and Two Buttons-One Stuck in virtue of which the agent could be a cause of the explosion in Two Buttons and not in Two Buttons-One Stuck. Therefore, the agent is not a cause of the

explosion in Two Buttons. In my opinion, this is what makes Two Buttons (and, in general, overdetermination cases involving omissions) particularly interesting: the fact that, in these cases, agents seem to be morally responsible for the ensuing outcomes without causing those outcomes.¹⁴

If not causation, what could the responsibility basis be, in cases like Two Buttons? In Sartorio (2004) I gave a partial answer to this question. The first thing to note is that, if we allow causation by absences, then, although neither of the individual failures is a cause of the explosion in Two Buttons, the failure of the two buttons to be simultaneously depressed presumably is. (If one doesn't believe in causation by absences, then one can simply rephrase this in terms of counterfactual dependence: although the explosion doesn't counterfactually depend on the individual failures, it does depend on the failure of the two buttons to be simultaneously depressed. One can then run the rest of the story in terms of counterfactual dependence.) Note that the fact that agent X failed to depress his button entails the fact that the two buttons weren't depressed at the relevant time, since the fact that the two buttons weren't depressed at the relevant time obtains just in case *either* X didn't depress his button *or* Y didn't depress his (or neither did). The same goes for the fact that agent Y failed to depress his button. In other words, the fact that the two buttons weren't depressed at the relevant time is a disjunctive

¹⁴ By the way, I sympathize with Moore's claim that overdetermining *actions* are causes. However, if this turned out to be wrong and overdeterminers were never causes (as, for example, Lewis (1986) believes; see n. 2 above), I would say the same thing about scenarios of overdetermination involving actions: agents are responsible for the outcome in these cases, and their responsibility is grounded in some other way.

fact, which obtains just in case at least one of the individual failures obtains. I will refer to this fact as 'D'.

Should we say, then, that each of the agents in Two Buttons is responsible for the explosion because the fact that he omitted to depress his button *entails* D and D was a cause of the explosion (or is something on which the explosion counterfactually depends)? Although it is quite natural to want to say this, again, unfortunately I don't think it works. For all of this is true of the agent in Two Buttons-One Stuck, where we (presumably) *don't* want to say that the agent is responsible for the explosion.¹⁵ What we should say, instead, is that the agents in Two Buttons are responsible for the explosion because they are *responsible* for D, which in turn caused the explosion (or is something on which the explosion depended). In contrast, the agent in Two Buttons-One Stuck is *not* responsible for D. Hence, even if the fact that he failed to depress his button entails D, and even if D was a cause of the explosion (or something on which the explosion depended), the agent doesn't come out responsible for the explosion in Two Buttons-One Stuck.

That is what I said in Sartorio (2004). Obviously, a full account of the agents' responsibility in Two Buttons would have to say something about *why* the two agents in Two Buttons are responsible for D but the one agent in Two Buttons-One Stuck is not. I didn't attempt to do that back then because my aim at the time was not to give an account of the bases for responsibility but only to explain how one can *consistently* claim that the two agents are responsible in Two Buttons but the one

¹⁵ Moore agrees with me about the agent's lack of responsibility in cases like Two Buttons-One Stuck. He then goes on to suggest that, in light of this fact, it is very hard to justify the claim that the agents are responsible in cases like Two Buttons (see, in particular, p. 451 and 467). In what follows I suggest a way in which to ground the difference in responsibility between the two scenarios.

agent is not responsible in Two Buttons-One Stuck. The thought was that, just as it is plausible to think that the agents in Two Buttons are responsible for the explosion but the agent in Two Buttons-One Stuck is not, it is also plausible to think that the agents in Two Buttons are responsible for D but the agent in Two Buttons-One Stuck is not.

Now, if one wanted to do more than this, what could one try to say? How could one try to explain the difference in the agents' responsibility between Two Buttons and Two Buttons-One Stuck with respect to D? This is what I want to explore in this final part of the paper.

As I explained, D is a disjunctive fact: it is the fact that obtains just in case agent X doesn't depress his button (at the relevant time, t) or agent Y doesn't depress his button (at t). X's failure to depress his button entails D, and so does Y's failure to depress his button. So consider the question: Under what conditions is it plausible to claim that an agent is/is not responsible for a disjunctive fact? In the literature on moral responsibility, several people have endorsed a principle along the following lines (I will call it the "Responsibility for Disjunctions" principle):

(RD) If both A and B obtain, and an agent is responsible for A but not for B, then the agent is not responsible for the disjunctive fact: A or B.¹⁶

¹⁶ See van Inwagen (1978), section III, Heinaman (1986), p. 306, Rowe (1989), p. 320, and Fischer and Ravizza (1993b), p. 346. van Inwagen would argue that the agent is not responsible for A or B because, given that he is not responsible for B, he couldn't have prevented the disjunction, and agents are not responsible for what they couldn't have prevented. Heinaman argues that an agent can be responsible for the disjunction A or B when he is responsible for A and when B is false, but agrees with van Inwagen that, if he is not responsible for B and B is *true*, then he is not responsible for A or B. Rowe argues that an agent is not responsible for A or B when B is true and it is made true by a completely independent process outside of the agent's control. And Fischer and Ravizza argue that an

Obviously, RD is not a general principle about responsibility for disjunctive facts, since it merely states a sufficient condition for the absence of responsibility for such facts. Still, it will be helpful for our purposes here.

As I will explain momentarily, I don't think that RD is universally true. However, I do think that something *in the neighborhood of* RD must be true and that this closely related principle helps explain our intuitions about certain cases. Consider, in particular, the following examples discussed by van Inwagen.¹⁷ Imagine that Gunnar shoots and kills Ridley and is responsible for doing so. Is he also responsible for the fact that Gunnar killed Ridley or $2+2=4$? Or for the fact that Gunnar killed Ridley or grass is green? It seems not. If so, a principle like RD would help explain why, since Gunnar is not responsible for the fact that $2+2=4$ or for the fact that grass is green.

Note that RD entails that the one agent in Two Buttons-One Stuck is not responsible for D. For he is not responsible for the fact that the other button wasn't depressed (this was the result of a mechanism over which he had no control). Note, however, that RD *also* entails that the two agents aren't responsible for D in Two Buttons. For neither agent is responsible for the status of the other button. This is where I think that the principle goes wrong. Claiming that the agents are not responsible for D in Two Buttons would violate The Slogan (note that the special circumstances described in section IV don't apply to the case at hand). Clearly, the

agent is not responsible for A or B when B does not result from a "responsive sequence". So all of these philosophers seem to be endorsing RD.

¹⁷ van Inwagen (1978), pp. 213-4.

agents in Two Buttons *are* responsible for the fact that the buttons weren't depressed at the relevant time. Hence, despite what philosophers seem to think, RD is not universally true. It is too strong.

However, I do think that some principle in the neighborhood of RD (a weaker principle) must be true. As van Inwagen's examples seem to suggest, agents are *typically* not responsible for disjunctive facts when they are not responsible for one of the obtaining disjuncts. As far as I can tell, agents are not responsible in circumstances of this kind *unless* absolving them of responsibility would result in a violation of The Slogan. Hence, at least as a first pass, we could try the following revision of the principle:

(RD*) If both A and B obtain, and an agent is responsible for A, and *no one (neither A nor anyone else) is responsible for B*, then the agent is not responsible for the disjunctive fact: A or B.

Just like RD, RD* only offers sufficient conditions for the lack of responsibility. But, unlike RD, RD* is consistent with the truth of The Slogan. As we have seen, in Two Buttons RD entails that neither agent is responsible for D, where The Slogan entails that the agents are responsible. In contrast, RD* is silent about the agents' responsibility for D. On the other hand, RD* explains our intuitions about van Inwagen's examples and about Two Buttons-One Stuck (where claiming that the agent is not responsible clearly doesn't violate The Slogan, since there is only one wrong in those cases, not two). In sum, RD* is a plausible and theoretically fruitful

principle of responsibility about disjunctive facts, which accounts for our intuitions about the responsibility of agents (or lack thereof) in some paradigmatic cases while remaining consistent with the truth of The Slogan.

I don't think that this is the end of the story, unfortunately. There are further questions that arise for scenarios where the agents' contributions are asymmetrical in certain ways.¹⁸ The existence of these scenarios suggests that RD* might be too weak. For example, consider a variant on the Two Buttons-One Stuck scenario where the mechanism controlling one of the buttons is a hundred years old. Imagine that the person who installed the mechanism at the time failed to take some necessary precautions that would have prevented the failure of the mechanism today. In that case we would probably want to hold that person responsible for the failure of the mechanism today.¹⁹ As a result, in this case RD* doesn't entail that the agent in charge of the other button (the agent who failed to depress his button today) is off the hook (the principle is consistent with his being responsible for the disjunctive fact that resulted in the explosion). I don't think it's obvious what one should say about this case, but the fact that the mechanism was bound to fail given what happened a hundred years ago suggests that it might be unreasonable to blame the agent for the explosion today. If so, we should look for a principle that is stronger than RD* (without being as strong as RD), to account for this fact.

I am not going to try to offer such a principle here, since it seems clear that potential complications of this kind arise only in connection with scenarios

¹⁸ Thanks to Randy Clarke for discussion of this point.

¹⁹ Even if that person is already dead. Presumably, we hold Hitler responsible for some things that happen nowadays.

involving asymmetrical contributions by agents. I have restricted my focus to the much “cleaner” symmetrical overdetermination cases, so these complications needn’t concern us. For our purposes here we may simply stick with D* while restricting the principle’s intended application to such cases.²⁰

On the basis of both The Slogan and RD*, it is possible to substantiate the claim that the agents in Two Buttons are responsible for the explosion but the agent in Two Buttons-One Stuck is not. The agents in Two Buttons are responsible for the explosion by virtue of being responsible for D but the agent in Two Buttons-One Stuck is not responsible for the explosion, since he is not equally responsible for D. A more general conclusion about the metaphysical bases for responsibility seems warranted: we should conclude that causation and counterfactual dependence are not the relevant responsibility bases in the case of responsibility for outcomes (or, at least, they are not the only ones). Being responsible for a cause (or being responsible for something on which the outcome depends) can function as a responsibility basis, and it is the one that is presumably operative in the relevant overdetermination scenarios.

References

Fischer, J. and M. Ravizza (1993a) (eds.) *Perspectives on Moral Responsibility*, Ithaca: Cornell University Press.

²⁰ Scenarios of asymmetrical contributions are famously tricky (see n. 4 above). Some of the questions that arise for them are: Does the mere fact that one agent acts before the other suggest that their contributions are asymmetrical? If they are asymmetrical, in what way, exactly, are they asymmetrical? (Is the person who acts first more responsible, in virtue of having acted first? Or is the person who acts second more responsible, in virtue of having acted closer to the occurrence of the harm?) It is not at all clear how one should adjudicate these questions.

Fischer, J. and M. Ravizza (1993b) "Responsibility for Consequences," in Fischer and Ravizza (1993a), pp. 321-47.

Hall, N. (2004) "The Intrinsic Character of Causation," in D. Zimmerman (ed.), *Oxford Studies in Metaphysics 1*, New York: Oxford University Press, pp. 255-300.

Heinaman, R. (1986) "Incompatibilism without Alternative Possibilities," *Australasian Journal of Philosophy* 64: 266-76, reprinted in Fischer and Ravizza (1993a), pp. 296-309.

Lewis, D. (1986) "Causation," in *Philosophical Papers II*, New York: Oxford University Press, pp. 159-213.

Moore, M. (2009) *Causation and Responsibility*, Oxford: Oxford University Press.

Rowe, W. (1989) "Causing and Being Responsible for what is Inevitable," *American Philosophical Quarterly* 26: 153-9, reprinted in Fischer and Ravizza (1993a), pp. 310-21.

Sartorio, C. (2004) "How to be Responsible for Something without Causing It," *Philosophical Perspectives* 18, 1: 315-36.

Thomson, J. J. (1989) "Morality and Bad Luck," *Metaphilosophy* 20: 203-21, reprinted in D. Statman (ed.), *Moral Luck*, Albany: State University of New York Press, 1993, pp. 195-215.

van Inwagen, P. (1978) "Ability and Responsibility," *The Philosophical Review* 83, 2: 201-24.

Yablo, S. (2002) "De Facto Dependence," *Journal of Philosophy* 99, 3: 130-48.

Zimmerman, M. (2002) "Taking Luck Seriously," *Journal of Philosophy* 99, 11: 553-76.