

1 9
2
3
4 Experimental Philosophy on
5 Free Will: An Error Theory for
6 Incompatibilist Intuitions
7
8

9
10 *Eddy Nahmias and Dylan Murray*
11
12
13

14 **9.1 Introduction**
15

16 It's called "the problem of free will and determinism," but much depends on
17 what determinism is taken to mean and entail. *Incompatibilists* claim that it is
18 impossible for people to have free will and moral responsibility if determinism
19 is true, and they often suggest that this is the natural position to take,
20 supported by our pre-theoretical intuitions. Robert Kane, for instance, states
21 that "ordinary persons start out as natural incompatibilists" (1999, 217), and
22 Galen Strawson claims that "it is in our nature to take determinism to pose
23 a serious problem for our notions of responsibility and freedom" (1986, 89).
24 Sometimes people take "determinism" to mean "the opposite of free will,"
25 in which case incompatibilism is indeed intuitive, but at the cost of being
26 an empty tautology. In philosophical debates, *determinism* has a technical
27 meaning: a complete description of the state of the universe at one time and
28 of the laws of nature logically entails a complete description of the state of
29 the universe at any later time.¹ However, it is not obvious why determinism,
30 defined in this way, is supposed to be incompatible with free will;
31 rather, a further explanation of just why determinism precludes some ability
32 associated with free will seems required. The explanations generally offered
33 by incompatibilists are that determinism precludes either (i) the ability to
34 choose among *alternative possibilities* for action, while holding fixed the
35 actual past and the laws of nature (AP), or (ii) the ability to be the *ultimate*
36 *source* of one's actions, such that one is ultimately responsible for some aspect
37 of the conditions that led up to one's actions (US). To say that incompatibilism
38 is intuitive, then, is presumably to claim that it is natural to find one or
39 both of these conditions necessary for free will and to understand the condition
40 in such a way that determinism precludes it.

41 *Compatibilists*, who believe that determinism does *not* preclude free will
42 and moral responsibility, often develop arguments to show why AP and
43 US, as defined by incompatibilists, are *not* in fact required for free will and
44 responsibility, sometimes offering analyses of abilities meant to capture the

1 attractive features of AP and US but in ways consistent with determinism.
2 They often argue that certain premises and principles used in incompatibilist arguments are mistaken. Compatibilists have also attempted to *explain away* the intuitions incompatibilists appeal to, sometimes suggesting that these intuitions are based on mistaken interpretations of the implications of determinism. For example, determinism may *appear* to threaten free will because it is conflated with types of coercion or manipulation, which, the compatibilist argues, are importantly different from determinism.

9 These conflicting views about what determinism entails and which abilities are required for free will typically lead to stalemates, often bottoming out in disagreements about which view best captures our ordinary intuitions and conceptual usage. For instance, incompatibilists claim that it is widely accepted that free will requires the ability to do otherwise. Compatibilists respond that it is not obvious that this ability must be “unconditional” (as suggested by AP); rather, free action requires a “conditional” ability to do otherwise *if* relevant earlier conditions had been different, an ability that is consistent with determinism. Given such stalemates, it would help to gain a better understanding of people’s pre-philosophical intuitions about free will, moral responsibility, and determinism, as well as the sources of these intuitions. This could help to elucidate which position in fact accords best with ordinary thinking about these issues, or whether some of the intuitions supporting one position are produced in systematically unreliable ways. Though such information certainly won’t *resolve* the debate, it can suggest that one side needs to answer certain questions, motivate its views in new ways, or take on the argumentative burden of proof.

26 One might attempt to uncover such information about folk intuitions and their underlying psychological processes through armchair analysis, but empirical methods will often be required to supplement such analysis, especially when philosophers on opposing sides offer conflicting claims about what is intuitive. The recent movement of “experimental philosophy” does just this, drawing on the empirical methods of psychology to systematically examine people’s intuitions about philosophical issues, and then carefully considering whether and how these results impact the philosophical debates. Below, we offer a brief history of experimental philosophy on free will before presenting results from our recent study. But first we jump ahead to the conclusion we take our results to support.

37 Our goal is to develop an “error theory” for incompatibilist intuitions—to show that, when ordinary people take determinism to preclude free will and moral responsibility, they usually do so because they *misinterpret* what determinism involves. In other words, we aim to explain why people *appear* to have incompatibilist intuitions, when in fact they do not. Whereas incompatibilists have suggested that “ordinary persons have to be talked out of [their] natural incompatibilism by the clever arguments of philosophers” (Kane, 1999, 217) and that “beginning students typically recoil at the compatibilist

1 response to the problem of moral responsibility" (Pereboom, 2001, xvi), we
2 believe that ordinary persons typically need help seeing the allure of incompatibilism. As suggested above, the proper conception of determinism needs
3 to be given to them—without being presented in a misleading way—and
4 then some explanation needs to be given for why determinism, so defined,
5 is incompatible with free will and moral responsibility, perhaps by motivating
6 the idea that they require AP and US. We suggest that in this process
7 many ordinary persons (for example, beginning students) come to interpret
8 determinism as entailing threats to free will that it does *not* in fact entail.
9 We predict that laypersons often mistakenly take determinism to mean that
10 everything that happens is inevitable—it will happen *no matter what*—or that
11 agents' decisions, desires, or beliefs make no difference to what they end up
12 doing, and that such mistakes then generate people's intuitions about agents'
13 lacking free will and moral responsibility. Indeed, people may take determinism
14 to preclude the sorts of abilities *compatibilists* associate with free will, such
15 as the abilities to consciously deliberate about what to do and to control one's
16 behavior in light of one's reasons. But if people's purportedly *incompatibilist*
17 intuitions result primarily from mistakenly interpreting determinism to preclude
18 what *compatibilists* require for free will, then these intuitions do not
19 support incompatibilism.
20

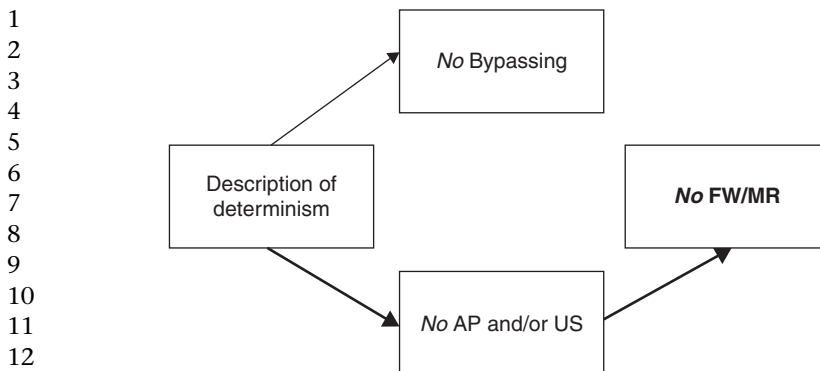
21 Suppose laypersons are presented with scenarios that describe a deterministic
22 universe, and suppose that some respond that agents in that universe
23 do *not* have free will (FW) and are *not* morally responsible (MR) for their
24 actions—they express "incompatibilist intuitions"—while others respond
25 that agents in these deterministic universes can have FW and MR—they
26 express "compatibilist intuitions." One explanation for such mixed results
27 (see below for examples) is that different people simply have different intuitions
28 about the relationship between determinism and FW or MR, perhaps
29 because they have different conceptions of "free will" or attribute moral
30 responsibility in varying ways (see Knobe & Doris, forthcoming). We think
31 that this interpretation may explain *some* of the variations in people's intuitions
32 and may even help to explain the intractability of the philosophical
33 debates. It may also be that some people who express compatibilist intuitions
34 do not understand the deterministic nature of the scenario or are not drawing
35 the intuitive connections between it and factors like AP and US. Perhaps
36 people fail to draw these purported implications of determinism due to an
37 emotional bias (Nichols & Knobe, 2007; see below). This would suggest an
38 error theory for compatibilist intuitions—that is, it would suggest that these
39 people have only *apparent*, but not *genuine* compatibilist intuitions.

40 However, the conflicting results might also be explained with an error
41 theory for the *incompatibilist* intuitions people seem to have. Our hypothesis
42 is that many people who appear to have incompatibilist intuitions are
43 interpreting determinism to entail what we will call "bypassing," and that
44 they take *bypassing* to preclude FW and MR. While *bypassing* does preclude

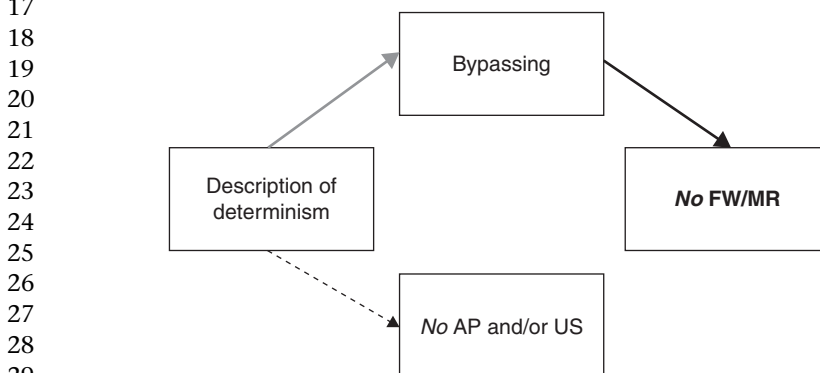
1 FW and MR, determinism *does not* entail bypassing. So, if the reason people
2 seem to express incompatibilist intuitions is that they mistakenly take deter-
3 minism to entail bypassing, then those intuitions are not *genuine* incompati-
4 bilist intuitions, and do not in fact support the conclusion that determinism,
5 properly understood, is incompatible with free will.

6 What is “bypassing”? The basic idea is that one’s actions are caused by
7 forces that bypass one’s conscious self, or what one identifies as one’s “self.”
8 More specifically, it is the thesis that one’s actions are produced in a way that
9 bypasses the abilities compatibilists typically identify with free will, such as
10 rational deliberation, conscious consideration of beliefs and desires, forma-
11 tion of higher-order volitions, planning, self-control, and the like.² As such,
12 bypassing might take the form of *epiphenomenalism* about the relevant mental
13 states (that is, that deliberations, beliefs, and desires are causally irrelevant to
14 action), or it might take the form of *fatalism*—the belief that certain things
15 will happen no matter what one decides or tries to do, or that one’s actions
16 *have to happen even if* the past had been different. Bypassing suggests that
17 conscious agents have no control over their actions because they play no role
18 in the causal chain that leads to their actions, and for our study discussed
19 below, we “operationalized” bypassing in specific ways that we take to capture
20 this intuitive idea.

21 The crucial point is that determinism, as defined by philosophers debat-
22 ing free will, simply does *not* entail bypassing (certainly not in the way we
23 operationalize it below). The history of compatibilism might be caricatured
24 as an attempt to drive home this point. Compatibilists have emphasized
25 that determinism does not mean or entail that all events are inevitable, in
26 the sense that they will happen no matter what we decide or try to do. They
27 point out that determinism does not render our beliefs, desires, deliberations,
28 or decisions causally impotent. Quite the contrary. So long as our mental
29 states are part of the deterministic sequence of events, they play a crucial
30 role in determining what will happen. Of course, incompatibilists generally
31 agree with all this, but claim their arguments are not based on such mistakes.
32 Nonetheless, the pre-philosophical *intuitive* appeal of incompatibilism may
33 rest largely on such mistakes, and to the extent that it does, incompatibil-
34 ists either need to abandon the appeal to wide-scale intuitive support as a
35 motivation or a basis for their position, or they need to demonstrate that
36 incompatibilism remains intuitive *even when* people properly recognize that
37 determinism does not entail bypassing. Put another way, since incompatibil-
38 ists generally allow that determinism is compatible with the abilities compati-
39 bilists associate with free will, “*genuine* incompatibilist intuitions” are
40 those that do not involve misinterpreting determinism to involve bypassing
41 of these compatibilist abilities (see Figure 9.1). If most people who take deter-
42 minism to preclude FW and MR do so on the basis of such a mistake, then,
43 this should at least shift the burden of proof onto the incompatibilist to
44 demonstrate that people nonetheless have *genuine* incompatibilist intuitions.



14 Figure 9.1 Genuine incompatibilist intuitions



30 Figure 9.2 Apparent incompatibilist intuitions

34 Our goal here is to offer evidence that most laypersons who respond that
 35 agents do *not* have FW and MR in deterministic universes are in fact express-
 36 ing only “*apparent* incompatibilist intuitions” because they misunderstand
 37 determinism to involve bypassing (see Figure 9.2).

38
39 **9.2 Experimental philosophy on free will**

40
41 As we have seen, philosophers often appeal to ordinary intuitions and
 42 common sense about free will and moral responsibility. We think such
 43 appeals have a legitimate place in the philosophical debate. This debate,
 44 unlike others about more technical concepts, is about concepts that are

1 intimately connected to ordinary people's beliefs about and practices con-
2 cerning morality, agency, praise, blame, punishment, reward, and so on.
3 The claim is not that ordinary intuitions or conceptual usage should *exhaust*
4 the philosophical analysis of free will, much less that they will inform us
5 about any extra-semantic facts about the nature of human decision-making.
6 Rather, the claim is that folk intuitions provide important information about
7 *which* extra-semantic facts we should be looking for when we want to know
8 whether humans have free will and are morally responsible for their actions.
9 Philosophical theories should systematize such intuitions as much as possi-
10 ble, revise them when they are inconsistent or when competing theoretical
11 advantages (such as consistency with scientific facts) call for it, or explain the
12 intuitions away—that is, offer an error theory for why they *appear* to support
13 a particular position when in fact they do not. If, instead, philosophers end
14 up mired in disputes about the proper analysis of a technical concept of “free
15 will” that no longer connects with ordinary concepts and practices, then
16 these debates risk being irrelevant.

17 We take it to be particularly important for incompatibilists to establish
18 the intuitive plausibility of their position, primarily because incompatibilist
19 theories of free will are generally more metaphysically demanding than com-
20 patibilist alternatives. Incompatibilist theories require indeterminism in the
21 agent at the right time and place, and often, additionally, agent causal powers.
22 These conditions are typically required *in addition to*, rather than *instead of*,
23 compatibilist conditions. Other things being equal, incompatibilists should
24 motivate the need for these extra metaphysical conditions. Many incompati-
25 bilists have motivated their more metaphysically demanding theories, at least
26 in part, by claiming that other things are *not* equal, because our ordinary
27 intuitions, as well as our phenomenology of decision-making, support incom-
28 patibilist views. It is certainly unclear why, *without* wide-scale intuitive sup-
29 port for incompatibilism, the burden of proof would be on compatibilists.³

30 Motivated by these considerations and the lack of any empirical data on
31 what intuitions laypersons actually have, Eddy Nahmias, Stephen Morris,
32 Thomas Nadelhoffer, and Jason Turner (2005, 2006) developed the initial
33 experimental philosophy studies on folk intuitions about FW, MR, and deter-
34 minism. Using three different descriptions of determinism, they found that
35 a significant majority of participants (typically 65–85 percent) judged that
36 agents in a deterministic scenario act of their own free will and are morally
37 responsible for their actions. One of the descriptions of determinism was the
38 following “re-creating universe” scenario:
39

40 Imagine there is a universe (Universe C) that is re-created over and over
41 again, starting from the exact same initial conditions and with all the
42 same laws of nature. In this universe the same initial conditions and the
43 same laws of nature cause the exact same events for the entire history of
44 the universe, so that every single time the universe is re-created, everything

1 must happen the exact same way. For instance, in this universe a person
2 named Jill decides to steal a necklace at a particular time and then steals it,
3 and *every* time the universe is re-created, Jill decides to steal the necklace at
4 that time and then steals it.⁴
5

6 After reading the scenario, participants were asked to judge whether Jill
7 decided to steal the necklace of her own free will and whether “it would
8 be fair to hold her morally responsible (that is, blame her) for her decision
9 to steal the necklace.” 66 percent of subjects judged that Jill acted of her
10 own free will, and 77 percent judged her to be morally responsible. Similar
11 results were found using two other scenarios describing determinism in dif-
12 ferent ways, which also included variations with agents’ performing positive
13 actions (for example, saving a child) or neutral actions (for example, going
14 jogging).⁵ These results offer evidence that a significant majority of layper-
15 sons are in fact “natural *compatibilists*” and thus call for an explanation as to
16 why so many philosophers have assumed that most ordinary people begin
17 with the intuition that determinism is incompatible with free will and moral
18 responsibility.

19 In response to these results from Nahmias, Morris, Nadelhoffer, and Turner
20 (NMNT), one might argue that people only *appear* to have compatibilist intu-
21 tions, when in fact they do not. Such judgments might be unreliable, or not
22 reflect people’s *considered* beliefs or folk theories about FW and MR. In part
23 to provide such an error theory for people’s compatibilist judgments, Shaun
24 Nichols and Joshua Knobe (2007) developed experiments aimed at exploring
25 the psychological mechanisms that generate intuitions about moral respon-
26 sibility. In their studies, participants were randomly assigned to one of two
27 groups, one of which was presented with a scenario in the “abstract” condi-
28 tion, and the other in the “concrete” condition. The scenario in the *abstract*
29 condition read:

30
31 Imagine a universe (Universe A) in which everything that happens is
32 completely caused by whatever happened before it. This is true from the
33 very beginning of the universe, so what happened in the beginning of the
34 universe caused what happened next, and so on right up until the present.
35 For example one day John decided to have French Fries at lunch. Like
36 everything else, this decision was completely caused by what happened
37 before it. So, if everything in this universe was exactly the same up until
38 John made his decision, then it *had to happen* that John would decide to
39 have French Fries.

40 Now imagine a universe (Universe B) in which *almost* everything that
41 happens is completely caused by whatever happened before it. The one
42 exception is human decision making. For example, one day Mary decided
43 to have French Fries at lunch. Since a person’s decision in this universe
44 is not completely caused by what happened before it, even if everything

1 in the universe was exactly the same up until Mary made her decision, it
 2 *did not have to happen* that Mary would decide to have French Fries. She
 3 could have decided to have something different.

4 The key difference, then, is that in Universe A every decision is com-
 5 pletely caused by what happened before the decision – given the past, each
 6 decision *has to happen* the way that it does. By contrast, in Universe B,
 7 decisions are not completely caused by the past, and each human deci-
 8 sion *does not have to happen* the way that it does.

9
 10 In the *concrete* condition, this scenario was followed by another paragraph:

11
 12 In Universe A, a man named Bill has become attracted to his secretary,
 13 and he decided that the only way to be with her is to kill his wife and
 14 3 children. He knows that it is impossible to escape from his house in the
 15 event of a fire. Before he leaves on a business trip, he sets up a device in
 16 his basement that burns down the house and kills his family.

17
 18 In the abstract condition, participants were asked:

19
 20 In Universe A, is it possible for a person to be fully morally responsible
 21 for their actions?

22 Yes No

23
 24 and in the concrete condition:

25
 26 Is Bill fully morally responsible for killing his wife and children?

27 Yes No

28
 29 In the concrete condition, 72 percent of subjects gave the compatibilist
 30 response that Bill *is* fully morally responsible. In the abstract condition, how-
 31 ever, 84 percent gave the purportedly *incompatibilist* response that it is *not*
 32 possible for a person to be fully morally responsible in Universe A.

33 Nichols and Knobe (N&K) claim that this disparity between participants'
 34 responses to the abstract and concrete cases is due to the psychological mech-
 35 anisms driving people's intuitions. They suggest that the immoral action in
 36 the concrete case engages people's emotions in a way that leads them to
 37 offer compatibilist judgments, but that these judgments are the result of a
 38 performance error of one's normal capacity to make correct attributions of
 39 moral responsibility. According to this *affective performance error model*, the
 40 emotions induced by high-affect scenarios, such as an agent's murdering
 41 his spouse and children, skew people's attributions of MR. Compatibilist
 42 responses are "performance errors brought about by affective reactions. In
 43 the abstract condition, people's underlying theory is revealed for what it
 44 is—incompatibilist" (2007, 672). The performance error model thus presents

1 an error theory for compatibilist intuitions. Because they are the result of
 2 affect-produced error, these intuitions should not be assigned any real sig-
 3 nificance in a theory of moral responsibility. As N&K say, “if we could elimi-
 4 nate the performance errors, the compatibilist intuitions should disappear”
 5 (2007, 678). N&K are rightly tentative about the affective performance error
 6 model and discuss the possibility that an *affective competence model*, a *concrete*
 7 *competence model*, or some hybrid model might instead be correct. N&K note
 8 that “we don’t yet have the data we need to decide between these competing
 9 models,” but they claim that “the philosophical implications of the perform-
 10 ance error model have a special significance because the experimental evi-
 11 dence gathered thus far [including NMNT’s] seems to suggest that the basic
 12 idea behind this model is actually true” (2007, 678).⁶

13 We believe that there are several problems with Nichols and Knobe’s (2007)
 14 study, as well as with their favored performance error model. For instance,
 15 N&K do not ask any questions about free will because, they say, “the expres-
 16 sion ‘free will’ has become a term of philosophical art, and it’s unclear how
 17 to interpret lay responses concerning such technical terms” (682, note 3).
 18 We do not think “free will” should be treated as a technical term nor that
 19 ordinary intuitions about free will are irrelevant to philosophical debates, so
 20 we continue to ask questions about it in our studies here. N&K instead use
 21 just one experimental question, asking whether it is possible for a person to
 22 be “fully morally responsible” for their actions. This phrase (itself somewhat
 23 technical sounding) is ambiguous and likely to be understood differently by
 24 different people. First, the “fully” may be contrasted either with being *partially*
 25 responsible or with being responsible in some *lesser* sense not involving moral
 26 desert. Indeed, attributions of MR are notoriously ambiguous between hold-
 27 ing people responsible for forward-looking (for example, deterrent) reasons
 28 and holding them responsible for backward-looking (for example, retributive)
 29 reasons. Thus, we think it is difficult to move directly from people’s responses
 30 to questions about being “fully morally responsible” to their intuitions about
 31 the compatibility of determinism and the sort of MR (or FW) involved in
 32 philosophical debates. Though we’re unsure how best to address these issues,
 33 we think questions about whether agents *deserve* blame (or praise) are more
 34 likely to elicit the relevant notion of moral responsibility, and we use this
 35 language in our studies.⁷

36 We are more concerned with N&K’s description of determinism. For
 37 instance, they write that in Universe A, “given the past, each decision *has*
 38 *to happen* the way that it does.” This wording leaves the scope of the modal
 39 operator ambiguous. It may be interpreted as: “Given past events, it is neces-
 40 sary (or inevitable) that later events (for example, decisions) happen the way
 41 they do” rather than: “It is necessary that, given past events, later events (for
 42 example, decisions) occur.”⁸ The latter, correct reading allows that later events
 43 (effects) *could* be otherwise as long as earlier events (causes) were otherwise.
 44 The former reading, however, mistakenly conflates determinism with fatalism

1 (that all actual events are necessary or inevitable), and it negates a compati-
2 bilist, conditional understanding of the ability to do otherwise (see above),
3 because one's actions *have to happen even if* the past (for example, one's rea-
4 sons) had been different (see Nahmias 2006 and Turner and Nahmias 2006).
5 Furthermore, the concluding sentence of N&K's abstract scenario reads,
6 "By contrast, in Universe B, decisions are not completely caused by the past,
7 and each human decision *does not have to happen* the way that it does." Some
8 participants may read this to mean that (by contrast), in Universe A, each
9 human decision *does have to happen* the way it does (full stop). We suspect
10 that this reading, perhaps along with the other issues we've raised, may lead
11 people to interpret N&K's description of determinism to suggest one or more
12 of the following: that agents' actions could not happen otherwise *even if* the
13 past had been different; that agents' decisions, beliefs, and desires are not
14 playing a role in influencing their actions; or that agents have no control over
15 what they do. That is, we predict that N&K's scenario will lead many people
16 to interpret determinism to involve bypassing.

17 N&K claim that "one cannot plausibly dismiss the high rate of incom-
18 patibilist responses in the abstract condition as a product of some subtle
19 bias in our description of determinism. After all, the concrete condition
20 used precisely the same description, and yet subjects in that condition were
21 significantly more likely to give compatibilist responses" (670–71). This
22 response, however, neglects the possibility that the description of determin-
23 ism has potentially misleading features that imply bypassing in the abstract
24 case but *not* in the concrete case. We do not dispute the idea that the high
25 negative affect likely induced in N&K's concrete case—with Bill's selfish,
26 premeditated, and wanton murder of his wife and three children—may bias
27 many participants to judge Bill to be morally responsible, but it may also lead
28 participants to neglect features of the scenario that might otherwise mitigate
29 their responsibility attributions. Indeed, it may be that the high negative
30 affect causes participants to neglect the *bypassing* features of the scenario. In
31 other words, N&K's description of determinism may lead people to make a
32 mistake, which is then "cancelled out" in the concrete case—but not in the
33 abstract case—by high negative affect. Hence, we predict that most people
34 will *not* read N&K's *concrete* scenario to involve bypassing, which may help to
35 explain why they are generally willing to attribute MR to Bill.

36 While we agree with N&K that very high negative affect likely biases
37 people to neglect potential responsibility-mitigating factors, we do not
38 agree with the more general assumption that people are more competent
39 in making judgments about FW, MR, and determinism when they consider
40 *abstract* cases than when they consider *concrete* cases including specific agents
41 performing specific actions. On the contrary, assuming all else is equal
42 (such as degree of affect), we believe that concrete conditions likely *facilitate*
43 participants' comprehension and capacity to make accurate attributions of
44 responsibility (in N&K's terms, we advocate a type of "concrete competence

1 model"). Specifically, we believe that judgments about responsibility—
2 including whether agents deserve credit or blame for their actions—will be
3 more reliable if they engage our capacities to think about the beliefs, desires,
4 and intentions of agents (for example, our “theory of mind” capacities), and
5 these are presumably more likely to be engaged when we consider specific
6 agents in specific circumstances. More generally, it may be that people’s
7 intuitions are more reliable when they have more details about a scenario,
8 which is likely part of the reason why philosophers construct thought experi-
9 ments with specific details to probe (or prime) our intuitions.⁹ Hence, while
10 we agree with N&K that concrete cases that *also* involve high affect may lead
11 to errors, we do not believe this is a product of concreteness *per se*. Rather,
12 we believe that, in general, concrete cases are more likely to reveal reliable
13 intuitions about MR and FW than abstract cases. For instance, we believe that
14 N&K’s description of determinism is more likely to lead to interpretations of
15 bypassing in the abstract case.¹⁰

16 Finally, some other explanation is required for why so many more partici-
17 pants express incompatibilist intuitions in N&K’s abstract scenario than in
18 NMNT’s cases, since the performance error model simply cannot account for
19 this difference. The majority of participants in NMNT’s (2006) studies gave
20 compatibilist responses, even for those scenarios that *did not involve high*
21 *negative affect*—that is, those that involved positive actions, such as saving a
22 child from a burning building or returning money one finds in a lost wallet, as
23 well as those that involved neutral actions, such as going jogging. Moreover,
24 no significant differences in responses were found between these cases and
25 those that did involve negative actions—robbing a bank, stealing a necklace,
26 and keeping the money one finds in a lost wallet (see note 5 above).¹¹ Thus,
27 the performance error model does not provide an explanation for these previ-
28 ous results. Some other explanation for the difference in responses to N&K’s
29 and NMNT’s cases is required. One possibility is that all of NMNT’s scenarios
30 describe *concrete* agents and actions and ask about those agents’ FW and
31 MR, whereas N&K’s abstract scenario does not include or ask about specific
32 agents or actions, but again, we believe there is no good reason to think con-
33 creteness alone leads to performance errors. Another (non-exclusive) possi-
34 bility is that N&K’s description of determinism primes bypassing judgments
35 significantly more than NMNT’s descriptions. Our new study explores these
36 possibilities.

37 An initial attempt to explore the issue of bypassing was developed in
38 Nahmias, Justin Coates, and Trevor Kvaran (2007). They found that, across
39 several different scenarios, most people responded that MR and FW were
40 possible in a deterministic universe *if* the scenario described the decisions
41 of agents in that universe as being “completely caused by the specific
42 thoughts, desires, and plans occurring in our minds.” In contrast, most people
43 responded that FW and MR were *not* possible in a deterministic universe *if*
44 the scenario described agents’ decisions as “completely caused by the specific

1 chemical reactions and neural processes occurring in our brains.”¹² The latter,
 2 reductionistic description seems to prime people to think that agents’ mental
 3 states are not playing the proper role in their actions—that their conscious self
 4 is bypassed. Thus, even though determinism in the technical sense is equally
 5 present in both scenarios, people tend to think determinism is compatible
 6 with FW and MR unless they take determinism to involve bypassing.

7 We designed our current study in order to further explore the possible
 8 effects of bypassing on people’s judgments of FW and MR, and to test
 9 our error theory for incompatibilist intuitions. We presented participants
 10 with different descriptions of determinism (N&K’s scenario versus NMNT’s
 11 “re-creating universe” scenario, with abstract and concrete versions of each),
 12 and then asked participants not only about FW and MR but also about
 13 bypassing. We predicted that

- 14
- 15 1. In general, participants’ judgments about bypassing would correlate sig-
 16 nificantly with their judgments about MR and FW. That is, when making
 17 judgments about agents in a deterministic universe, (a) most participants
 18 who respond that the agents do *not* have MR and FW would also respond
 19 that the agents’ decisions, beliefs, and desires do *not* affect what happens—
 20 that is, such participants would interpret the deterministic nature of the
 21 scenario to involve bypassing—whereas (b) most participants rejecting the
 22 bypassing claims would respond that the agents *do* have MR and FW. That
 23 is, bypassing judgments would explain away most *apparent* incompatibilist
 24 intuitions, whereas most people who do *not* misunderstand determinism
 25 to involve bypassing would express *prima facie* compatibilist intuitions.¹³
- 26 2. Judgments of FW and MR would be *lower*, while judgments of bypassing
 27 would be *higher*, in N&K’s abstract scenario compared to NMNT’s abstract
 28 scenario. That is, N&K’s description of determinism would, in the abstract
 29 case, lead more people to misunderstand determinism.
- 30 3. Judgments of FW and MR would be *lower*, while judgments of bypassing
 31 would be *higher*, in the abstract scenarios compared to the concrete sce-
 32 narios, with this difference especially pronounced in N&K’s high-affect
 33 scenario.

34 35 9.3 Methods

36
 37 Participants included in the analysis were 249 undergraduate students at
 38 Georgia State University (Atlanta, GA) who were randomly assigned to com-
 39 plete one of four versions of the experimental task.¹⁴ We used software from
 40 QuestionPro to develop and administer these surveys online. Using a 2x2
 41 between-subjects design, four scenarios were generated by systematically
 42 varying (1) whether the deterministic scenario was N&K’s or NMNT’s, and
 43 (2) whether the scenario was *abstract* or *concrete*. AQ1

1 Participants began by reading a general description of the task, providing
 2 informed consent, and then reading one of the four scenarios. After reading the
 3 scenario, participants answered a series of experimental questions designed to
 4 probe their intuitions about FW and MR, as well as whether they interpreted
 5 the scenario to involve bypassing. N&K's abstract and concrete scenarios
 6 read exactly as they are presented above, as did NMNT's concrete scenario.
 7 NMNT's abstract scenario replaces the last sentence of the concrete version
 8 with:

9
 10 For instance, in this universe whenever a person decides to do something,
 11 every time the universe is re-created, that person decides to do the same
 12 thing at that time and then does it.

13
 14 In order to replicate N&K's study, participants given those surveys were first
 15 asked:

16
 17 Which of these universes do you think is most like ours? Universe A
 18 Universe B

19
 20 Participants who were given the NMNT surveys were first asked:

21
 22 Is it possible that our universe could be like Universe C, in that the same
 23 initial conditions and the same laws of nature cause the exact same
 24 events for the entire history of the universe? Yes No

25
 26 Participants were next asked to indicate their level of agreement with each of
 27 a series of statements using a 6-point rating scale (strongly disagree, disagree,
 28 somewhat disagree, somewhat agree, agree, strongly agree). The first statement
 29 was always the *moral responsibility* (MR) question (replicating N&K's format).
 30 The remaining statements in each survey were randomized to decrease the
 31 likelihood of order effects. The most important experimental questions we
 32 asked read as follows (variations between scenarios are in brackets: *N&K*
 33 *abstract* scenario asks about Universe A; *NMNT abstract* asks about Universe C;
 34 *N&K concrete* asks about Bill; *NMNT concrete* asks about Jill):

35 36 9.3.1 The MR/FW questions

37
 38 **MR:** In Universe [A/C], it is possible for a person to be fully morally
 39 responsible for their actions.

40 [Bill/Jill] is fully morally responsible for [killing his wife and children/
 41 stealing the necklace].

42 **FW:** In Universe [A/C], it is possible for a person to have free will.

43 It is possible for [Bill/Jill] to have free will.
 44

1 **Blame:** In Universe [A/C], a person deserves to be blamed for the bad
 2 things they do.
 3 [Bill/Jill] deserves to be blamed for [killing his wife and children/stealing
 4 the necklace.]
 5

6 **9.3.2 The Bypassing questions**

7 (These questions represent our way of operationalizing “bypassing”; we take
 8 it that philosophers on all sides of the free will debate should agree that
 9 if one *properly* understands determinism, one should *not* agree with these
 10 statements.)
 11

12 **Decisions:** In Universe [A/C], a person’s decisions have no effect on what
 13 they end up being caused to do.
 14 [Bill’s/Jill’s] decision to [kill his wife and children/steal the necklace] has
 15 no effect on what [he/she] ends up being caused to do.

16 **Wants:** In Universe [A/C], what a person wants has no effect on what
 17 they end up being caused to do.
 18 What [Bill/Jill] wants has no effect on what [he/she] ends up being caused
 19 to do.

20 **Believes:** In Universe [A/C], what a person believes has no effect on what
 21 they end up being caused to do.

22 What [Bill/Jill] believes has no effect on what [he/she] ends up being
 23 caused to do.

24 **No Control:** In Universe [A/C], a person has no control over what
 25 they do.

26 [Bill/Jill] has no control over what [he/she] does.

27 **Past Different:** In Universe A, everything that happens *has to* happen,
 28 even if what happened in the past had been different.

29 Bill *has* to kill his wife and children, even if what happened in the past
 30 had been different.¹⁵
 31

32 After providing responses to these questions, participants then answered two
 33 comprehension questions to ensure that they understood the scenario and
 34 several demographic questions (for example, gender, age, religious affiliation,
 35 and so on).
 36

37 **9.4 Main results**

38
 39 In order to examine the relationship between participants’ judgments about
 40 bypassing and their judgments about MR and FW, we created two com-
 41 posite scores which we used for the analyses below: an *MR/FW composite*
 42 *score*, which was obtained by computing the average of each participant’s
 43 responses to the MR, FW, and Blame questions, and a *Bypassing composite*
 44 *score*, obtained by computing the average of each participant’s responses to

1 the Decisions, Wants, Believes, and No Control questions.¹⁶ Initial visual
 2 inspection of the data suggested that, as we predicted, (i) MR/FW scores
 3 were lower and *Bypassing* scores higher in response to the abstract scenarios
 4 compared to the concrete scenarios, and (ii) MR/FW scores were lower,
 5 while *Bypassing* scores were higher, in response to N&K's *abstract* scenario
 6 compared to NMNT's *abstract* scenario (see Figure 9.3). Moreover, across
 7 conditions, (i) the majority of participants who gave *apparent* incompati-
 8 bilist responses (MR/FW scores less than the 3.5 midpoint) also gave *high*
 9 bypassing responses (Bypassing scores > 3.5), whereas (ii) most participants
 10 who gave *prima facie* compatibilist responses (MR/FW scores > 3.5) also gave
 11 *low* (< 3.5) bypassing responses (see Figure 9.4).¹⁷ Given these findings, we
 12 employed a series of analyses in order to determine whether these results
 13 were statistically significant.¹⁸

14 To determine whether MR/FW scores were significantly lower in N&K's
 15 surveys than in NMNT's surveys and lower in the abstract conditions than
 16 in the concrete conditions, we ran a 2 (survey: N&K, NMNT) x 2 (condition:
 17 abstract, concrete) Analysis of Variance (ANOVA) on the mean MR/FW com-
 18 posite scores (see Figure 9.5). The ANOVA showed a significant main effect for
 19 survey: $F(1, 245) = 5.396, p = .021$, a significant main effect for condi-
 20 tion: $F(1, 245) = 61.058, p < .001$, and a marginally significant interaction

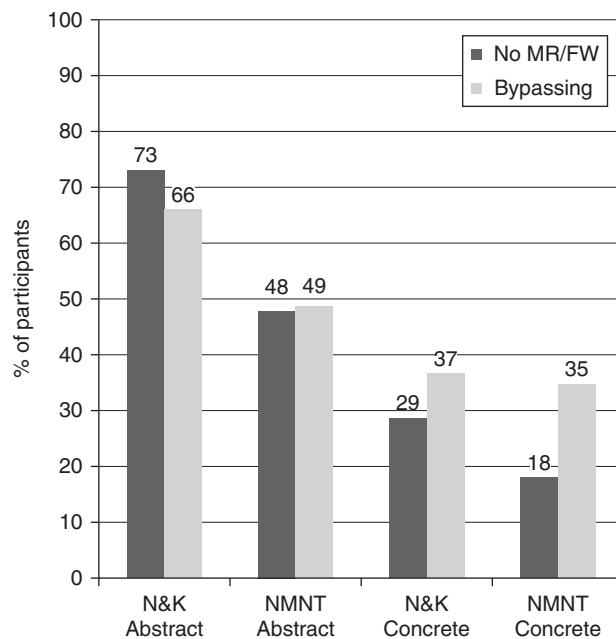
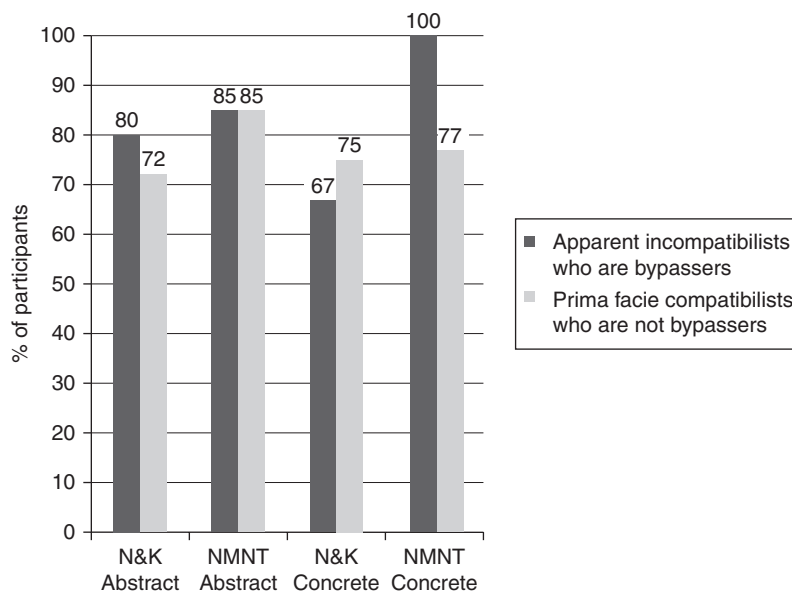


Figure 9.3 Judgments about MR, FW, and bypassing

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20



21 Figure 9.4 Apparent incompatibilists who are bypassers and prima facie compatibilists
22 who are not

23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

effect: $F(1, 245) = 3.297, p < .071$. We ran two additional pre-planned t -tests specifically comparing the mean MR/FW composite responses to *N&K abstract* vs. *NMNT abstract*, and to *N&K concrete* vs. *NMNT concrete*. We found that the mean MR/FW score was significantly lower in *N&K abstract* than in *NMNT abstract*: $t(1, 131) = -2.973, p = .004$, but was *not* significantly different between the two concrete conditions: $t(1, 114) = 2.346, p = .128$. Thus, the abstract scenarios received significantly lower MR/FW ratings than the concrete scenarios, independent of whether the survey was N&K's or NMNT's, and N&K's surveys received lower MR/FW scores than NMNT's in both the concrete and abstract conditions, but only significantly lower scores in the concrete condition. We believe that MR/FW scores are not significantly lower in *NMNT concrete* because of the very high affect involved in *N&K concrete*, which we suspect drives down participants' *Bypassing* scores, and thereby drives up their MR/FW scores (see note 20).

To determine whether *Bypassing* scores were significantly higher in N&K's surveys than in NMNT's surveys and higher in the abstract conditions than in the concrete conditions, we ran a 2 (survey: N&K, NMNT) x 2 (condition: abstract, concrete) ANOVA on the mean *Bypassing* composite scores (see Figure 9.5). The ANOVA showed a significant main effect for condition: $F(1, 245) = 20.665, p < .001$, but only a near-significant effect for survey: $F(1, 245) = 3.463, p = .064$ (see note 20). There was no significant interaction

Survey	Condition	N	MR/FW		Bypassing	
			Mean	Std. Dev.	Mean	Std. Dev.
N&K	Abstract	77	2.818	1.204	3.958	1.205
	Concrete	56	4.363	1.298	3.018	1.216
NMNT	Abstract	56	3.482	1.360	3.442	1.346
	Concrete	60	4.444	1.174	2.946	1.183

Figure 9.5 MR/FW and bypassing descriptive statistics

effect. Thus, the abstract scenarios received significantly higher *Bypassing* scores than the concrete scenarios, independent of whether the survey was N&K's or NMNT's, and N&K's surveys received marginally higher *Bypassing* scores than NMNT's, independent of whether the condition was abstract or concrete. We ran an additional pre-planned *t*-test specifically comparing the mean *Bypassing* responses to *N&K abstract* vs. *NMNT abstract*. As hypothesized, we found that the mean *Bypassing* score was significantly higher in *N&K abstract* than in *NMNT abstract*: $t(1, 131) = 2.319, p = .022$.¹⁹

In order to statistically assess the relationship between these two variables of interest, we computed Pearson correlation coefficients between the *MR/FW* and *Bypassing* composite scores for each scenario. Consistent with our hypothesis, but even more dramatically than we expected, we found a strong inverse correlation between *Bypassing* and *MR/FW* scores—that is, the higher a participant's *Bypassing* score, the lower his or her *MR/FW* score, and vice versa—in each of the four scenarios (*N&K abstract*: $r(75) = -0.695, p < .001$; *N&K concrete*: $r(54) = -0.569, p < .001$; *NMNT abstract*: $r(54) = -0.803, p < .001$; *NMNT concrete*: $r(58) = -0.708, p < .001$). Collapsing across all four surveys, the correlation coefficient between *Bypassing* and *MR/FW* scores was strikingly high: $r(247) = -0.734, p < .001$.

Consistent with our hypothesis, then, average scores in response to the *Bypassing* questions were significantly *higher* in *N&K's abstract* scenario than in *NMNT's abstract* scenario, average scores in response to the *MR/FW* questions were significantly *lower* in *N&K abstract* than in *NMNT abstract*, and responses to the *Bypassing* and *MR/FW* questions were strongly inversely correlated across scenarios. Given these results, we further suspected that responses to the *MR/FW* questions were lower for *N&K abstract* compared to *NMNT abstract* precisely *because* participants interpreted N&K's abstract scenario to involve a higher degree of bypassing. We hypothesized that the degree to which one interpreted a scenario to involve bypassing would *mediate* the relationship between survey and *MR/FW* responses. That is, we hypothesized that the difference in *MR/FW* responses between the two *abstract* conditions of the surveys was *caused* largely by people's bypassing judgments. In order to test this causal hypothesis more directly, we used a mediation analysis.²⁰

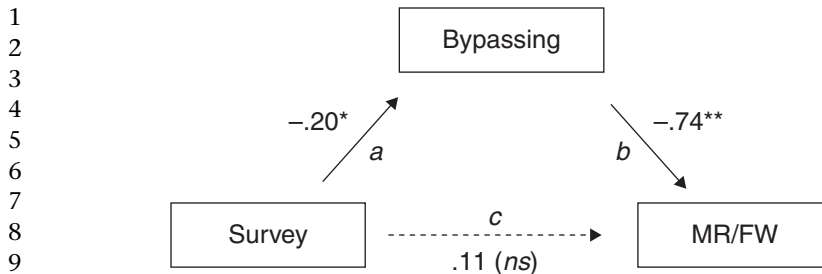


Figure 9.6 Mediation analysis

We conducted three regression analyses to test for mediation (see Figure 9.6), as outlined by Baron and Kenny (1986) and more recently by MacKinnon et al. (2002).²¹ The first regression equation used survey type (*N&K abstract*, *NMNT abstract*) to predict MR/FW score (path c), and yielded a significant effect: $t(132) = 2.973, p = .004$, corroborating the above results showing that *N&K abstract* prompts participants to respond that agents do not have MR and FW more than does *NMNT abstract*. The second regression equation estimated changes in *Bypassing* score using survey type (path a) and also yielded a significant effect: $t(132) = -2.319, p = .022$, corroborating the above results showing that participants interpreted *N&K abstract* to involve bypassing more than they did *NMNT abstract*. The third equation estimated MR/FW score using both survey type and *Bypassing* score. The link between *Bypassing* and MR/FW scores (path b) was highly significant: $t(132) = -12.799, p < .001$, and the relation (path c) between survey and MR/FW score was reduced to *non-significance* once *Bypassing* was included in the model: $t(132) = 1.821, ns$; Sobel test = $2.286, p < .022$. Thus, all the conditions of mediation were met: survey type was a significant predictor of MR/FW and of *Bypassing* scores, and *Bypassing* was a significant predictor of MR/FW scores while controlling for survey.²²

Thus, as we hypothesized, the effect that survey type (that is, the description of determinism in the abstract scenarios) had on a participant's MR/FW responses was mediated by whether the participant interpreted the description to involve bypassing. That is, these results suggest that which survey participants read (*N&K* versus *NMNT abstract*) had *no significant causal effect* on their MR/FW responses *over and above* the effect it had in virtue of causing different interpretations of whether the scenario involved bypassing.

9.5 Discussion

We predicted that participants across surveys and conditions would be *more* likely to judge determinism to threaten free will and moral responsibility

1 when they interpreted determinism to involve bypassing. We also predicted
2 that participants would be *more* likely to interpret determinism to involve
3 bypassing in N&K's abstract scenario than in NMNT's abstract scenario,
4 and that this would help to explain why the description of determinism
5 in N&K's scenario leads people to be *less* likely to attribute FW and MR to
6 agents. Finally, we predicted that participants would be *less* likely to interpret
7 determinism to involve bypassing in concrete scenarios, especially in
8 N&K's high-affect case, and hence *more* likely to attribute FW and MR to the
9 agents in those scenarios. Our results strongly support each of these predictions.
10 Indeed, not only do our results show that (i) the vast majority of
11 participants who express *apparent* incompatibilist intuitions interpret determinism
12 to involve bypassing, while those who express *prima facie* compatibilist
13 intuitions tend *not* to misinterpret determinism in this way, and that
14 (ii) there is a dramatic correlation between the degree to which participants
15 take a scenario to involve bypassing and the degree to which they attribute
16 MR and FW to agents in that scenario, but the results also suggest that
17 (iii) the difference in attributions of MR and FW between the abstract conditions
18 is *caused by* people's bypassing interpretations.

19 We think it is safe to conclude from our results that there is an important
20 connection between (a) whether people take a description of determinism
21 to entail bypassing and (b) whether people take the scenario so described to
22 preclude free will and moral responsibility. The most plausible interpretation
23 of this connection is that when a person takes determinism to entail
24 bypassing it generally *causes* him or her to judge that determinism precludes
25 MR and FW, and hence to offer *apparent* incompatibilist intuitions.
26 Conversely, when a person does *not* take determinism to entail bypassing,
27 they are likely to offer compatibilist intuitions—they do not see any conflict
28 between determinism and FW or MR (below we will further consider whether
29 such intuitions should count as supporting compatibilism). This causal
30 conclusion also draws support from Nahmias et al.'s (2007) study,
31 which manipulated a type of bypassing directly—rather than measuring
32 responses to it—and found that it significantly influenced participants'
33 attributions of MR and FW. Thus, previous evidence supports the conclusion
34 that bypassing judgments mediate attributions of MR and FW, rather than
35 the other way around.

36 If these interpretations of the data are correct, then people's interpreting
37 determinism to entail bypassing may be the best explanation for ordinary
38 people's intuitions that *appear* to support the incompatibility of determinism
39 and FW and MR. But these intuitions do *not* properly support incompatibilism
40 because determinism does *not* properly entail bypassing as we have
41 operationalized it here. Determinism, properly understood, simply does *not*
42 entail that our decisions, beliefs, and desires have no effect on what we end
43 up doing, nor that we have no control over what we do (see note 15), nor that
44 our actions *have to* happen just as they do *even if* the past had been different

1 (see below). When people do not misinterpret determinism in these ways,
2 they usually do not take it to threaten FW and MR.

3 Furthermore, our results suggest that certain descriptions of determinism
4 and certain conditions increase the degree to which people will misinter-
5 pret determinism to entail bypassing (and hence to lower their attributions
6 of MR and FW). Specifically, Nichols and Knobe's (2007) description of
7 determinism in the abstract condition has this effect, and this likely occurs
8 because of the problems we pointed out earlier with their description of
9 determinism, including the use of language that may suggest that decisions
10 and actions in Universe A *have to happen even if* the past had been different.
11 Indeed, 48 percent of participants in the *N&K abstract* condition responded
12 that, in Universe A, everything that happens has to happen the way it does
13 even if the past had been different, and these responses were significantly
14 correlated with participants' *MR/FW* scores.²³ The degree to which people
15 interpret determinism to involve bypassing is also higher in the abstract
16 conditions than in the concrete conditions, perhaps because descriptions
17 of concrete agents and actions prime people to think about the effective-
18 ness of agent's beliefs, desires, and decisions (for example, consider how
19 difficult it is to agree with the statement: "What Bill wants has no effect on
20 what he ends up being caused to do"—his desire to be with his secretary
21 was precisely what led to his murderous actions!). And when the concrete
22 conditions induce *high* negative affect, as with N&K's murderous Bill, this
23 tends to lower people's bypassing judgments while leaving their judgments
24 of FW and MR high.

25 As we said, we find it plausible that very high negative affect can induce
26 biases in MR judgments; there is, after all, evidence of such biases (see Nichols &
27 Knobe 2007, 672). But this does not suggest that, in general, people get
28 things wrong in concrete cases and right in abstract cases. Indeed, given
29 that interpreting determinism to entail bypassing is a *mistake*, it appears that
30 scenarios with abstract descriptions of unspecified agents performing unspec-
31 ified actions may bias people toward making this mistake. This is unsurpris-
32 ing, since thinking in terms of the efficacy of agents' mental states—taking
33 the intentional stance—is more likely to occur when one is thinking about
34 specific agents performing specific actions.

35 Due to these considerations, we believe that, in general, concrete scenarios
36 involving low affect are more useful for eliciting reliable intuitions about
37 FW, MR, and (in)compatibilism than either abstract or high-affect scenarios.
38 Indeed, fewer participants made the mistake of interpreting determinism to
39 involve bypassing in *NMNT concrete* than in any other scenario. Furthermore,
40 *every* one of the few participants who did respond as an apparent incom-
41 patibilist in *NMNT concrete* also interpreted the scenario to involve bypass-
42 ing, while *every* participant who did *not* interpret the scenario to involve
43 bypassing responded as a *prima facie* compatibilist. Thus, the scenario caused
44 fewer participants to mistake the scenario to involve bypassing and less

1 variability in responses among participants' who did not make this mistake.
2 Future research should explore these issues. For instance, it would be useful
3 to compare responses to bypassing and MR/FW questions in other concrete
4 cases, such as N&K's low-affect concrete case of the tax cheat and NMNT's
5 cases involving positive actions (for example, agent saving a child) and
6 neutral actions (for example, agent mowing the lawn), and perhaps also to
7 develop, if possible, high-affect *abstract* cases.

8 Hence, we conclude that the extant research in experimental philosophy
9 on free will converges on the conclusion that most laypersons do *not* have
10 *genuine* incompatibilist intuitions. They do have the intuition that bypassing
11 undermines FW and MR, and they can be primed to judge that determinism
12 entails bypassing. But the latter judgment is based on a mistaken interpreta-
13 tion of determinism. The judgment that bypassing undermines FW and MR
14 does not support incompatibilism. If anything, it supports compatibilist the-
15 ories of freedom and responsibility, since those theories emphasize that FW
16 and MR require that our conscious, rational deliberative processes play the
17 right role in producing our decisions and actions. So, folk intuitions that take
18 bypassing of these processes to preclude FW and MR support compatibilist
19 theories. People seem to be attending to compatibilist conditions for FW and
20 MR—whether agent's actions are properly caused by their decisions, beliefs,
21 desires, and so on. The more likely it is that people believe that these condi-
22 tions are met, the more likely they are to attribute FW and MR; the more
23 likely it is that people believe that these conditions are *not* met—for example,
24 because of bypassing—the less likely they are to attribute FW and MR. These
25 results support the claim that people have merely *apparent* incompatibilist
26 intuitions—most people do *not* seem to think that determinism—*without*
27 bypassing—precludes FW and MR (see Figure 9.2).

28 There are several ways in which incompatibilists might object to these
29 interpretations of our results, but space limits us to consider only some of
30 them briefly—we leave it to our critics to complete the task.

31 First, one might argue that even *after* recognizing that determinism does
32 *not* involve bypassing, people would (or should) recognize that determinism
33 threatens FW and MR, and that those who do not recognize this are likely
34 failing to understand the deterministic nature of the scenario. This response,
35 in effect, offers a debunking explanation for compatibilist intuitions by
36 arguing that participants who express such judgments do so because they
37 fail to understand determinism or its implications. Perhaps they do not
38 recognize that determinism is incompatible with AP or US, but if they *did*,
39 they would think it was incompatible with MR and FW (even while also
40 recognizing that determinism does not involve bypassing).

41 This is an interesting objection, and we would be intrigued to see experi-
42 ments to test for it. Recall that the vast majority of our participants who
43 did *not* take determinism to involve bypassing also attributed FW and MR
44 to agents in those scenarios, so this objection requires that most of these

1 participants are failing to understand the deterministic nature of the scenario
2 or failing to understand that these agents do not meet conditions *the partici-*
3 *pants themselves* take to be necessary for free will (for example, AP or US *in the*
4 *incompatibilists' sense*).²⁴ Future research should try to elicit whether people
5 understand determinism to conflict with AP and US and whether people take
6 AP and US—understood in ways that are incompatible with determinism—to
7 be necessary for free will or moral responsibility. Unfortunately, it is difficult
8 to properly describe what it means for an agent to be the ultimate source of
9 her decisions, or to have an *unconditional* ability to do otherwise, without a
10 good bit of explanation, which might verge on “intuition coaching.” While
11 we are inclined to think that this difficulty suggests that AP and US are not
12 particularly “natural” or intuitive to non-philosophers, others may argue
13 that they *are* intuitive once people properly understand the relevant ideas.
14 This suggests a second objection to our interpretation of the data.

15 The incompatibilist might argue that untutored intuitions are simply
16 *irrelevant* to the philosophical debates about free will and moral responsibil-
17 ity or, what is different, that the sort of studies carried out by experimental
18 philosophers cannot uncover information about the relevant intuitions.
19 These objections might be motivated by the belief that the issues are so
20 complex that responses from untrained individuals reveal little to nothing
21 about the truth, or by the conviction that philosophers can discern the rel-
22 evant intuitions by considering their own intuitions or those adduced from
23 their students and other folk, or by the belief that the methods employed by
24 experimental philosophers simply cannot do the job they aim to do (see, for
25 example, Kauppinen, 2007).

26 Extensive responses to such objections applied to experimental philoso-
27 phy in general have been offered elsewhere (for example, Nadelhoffer and
28 Nahmias, 2007; Nahmias et al., 2006; Knobe and Nichols, 2008; Weinberg,
29 2007). We reiterate that philosophical debates about free will and moral
30 responsibility require an understanding of the way non-philosophers think
31 about these issues in order to develop a theory that accords with folk intui-
32 tions, where possible, and to know what it is that we’re revising (and why)
33 where revision is advocated. We agree with David Lewis when he says that
34 in developing philosophical theories “we are trying to improve *that* theory,
35 that is to leave it recognizably the same theory we had before” (1986, 134).
36 If the studies we have presented here or previous studies have design flaws,
37 then attempts should be made to improve them, rather than abandoning
38 the very idea of understanding folk intuitions in an empirically informed
39 way. Part of this process might include making sure that participants under-
40 stand the concepts involved as clearly as possible, but one of the motiva-
41 tions for surveying people untrained in the philosophical debates is the
42 worry that philosophical training may end up shaping intuitions toward
43 a certain theory. For instance, it is not uncommon for philosophy teach-
44 ers to initially present determinism using metaphors that suggest fatalism

1 or epiphenomenalism (or, on the other side, to present indeterminism as
 2 involving entirely random uncaused events). Finally, since professional
 3 philosophers writing about free will tend to have theoretical commitments
 4 and hence *post*-theoretical “intuitions,” and since these “intuitions” (as well
 5 as reports about folk, for example, students’ intuitions) tend to conflict with
 6 each other, it is appropriate to attempt, as best we can, to uncover informa-
 7 tion about *pre*-theoretical intuitions, their sources, and their reliability.

9 9.6 Conclusion

10 Incompatibilists suggest that free will and moral responsibility require condi-
 11 tions that are incompatible with determinism—conditions that are generally
 12 more demanding than those required by compatibilists. One way to moti-
 13 vate the claim that these conditions are indeed necessary is to argue that
 14 incompatibilism is intuitive, and that compatibilism is thus a “quagmire of
 15 evasion,” a revision of the way ordinary people think about these issues—to
 16 suggest that “ordinary persons have to be talked out of this natural incom-
 17 patibilism by the clever arguments of philosophers” (Kane, 1999, 217).
 18 Our evidence here suggests that people may instead need to be talked *into*
 19 incompatibilism by the clever arguments, or subtle thought experiments, of
 20 philosophers. We suggest that incompatibilism only *appears* to be intuitive,
 21 largely because determinism is misinterpreted. Indeed, it is misinterpreted
 22 such that it precludes the very conditions compatibilists identify with free
 23 and responsible agency. It may be that incompatibilism is intuitive even
 24 *after* this mistake is corrected—that people find determinism threatening
 25 even if they understand that it does *not* involve bypassing (for example,
 26 fatalism or epiphenomenalism). We await the evidence. And obviously, *some*
 27 people—for example, some philosophers—do have genuine incompatibilist
 28 intuitions. But if most people think that free will and moral responsibility
 29 can exist even if determinism (properly construed) is true, the argumenta-
 30 tive burden shifts to these philosophers to explain why people’s intuitions
 31 need to be revised so that they accept a more demanding theory of free will.
 32 We await the argument.²⁵

35 Notes

- 36
- 37 1. This definition is drawn from van Inwagen (1983). Two less technical, though not
 38 quite equivalent, ways of stating determinism are: (1) In a deterministic universe,
 39 necessarily, re-creating identical initial conditions and laws of nature produces
 40 identical later events; (2) Determinism is the thesis that every event is *completely*
 41 *caused* by earlier events, such that, necessarily, *given* the earlier events and the laws
 42 of nature, the later events occur. These two descriptions are more similar to the
 43 ones used in the studies described below.
 - 44 2. For compatibilist accounts of these abilities, see, for example, Fischer and Ravizza
 (1998), Frankfurt (1971), Watson (1976), and Wolf (1990). Even incompatibilists

- 1 generally take these sorts of capacities to be *necessary* for free and responsible
 2 agency—see, for example, O'Connor (2005). Related use of “bypassing” language
 3 is introduced in Blumenfeld (1988) and Mele (1995).
- 4 3. For further development of the points raised in this paragraph see Nahmias *et al.*
 5 2006, 30–33. Moreover, if *revisionism* is called for (Vargas, 2005), it's unclear why
 6 philosophers should revise the concept of free will to be more metaphysically
 7 demanding than required by our ordinary intuitions.
- 8 4. The wording of the scenario as presented here is slightly altered from that used
 9 by Nahmias *et al.* (2006) in order to reflect the exact wording we used in our new
 10 study presented below.
- 11 5. In one scenario (Jeremy), affirmative responses for FW were 76 percent (negative
 12 action), 68 percent (positive), and 79 percent (neutral), and for MR they were
 13 83 percent (negative) and 88 percent (positive). For the other scenario (Fred &
 14 Barney), affirmative responses for FW were 76 percent (negative action) and
 15 76 percent (positive), and for MR they were 60 percent (negative) and 64 percent
 16 (positive) (see Nahmias *et al.*, 2006, 39). All results were significantly different
 17 from chance, as determined by χ^2 goodness-of-fit tests.
- 18 6. Interestingly, Paul Edwards anticipated this model 50 years earlier: “The very
 19 same persons, whether educated or uneducated, use it [MR] in certain contexts in
 20 the one sense and in other contexts in the other. Practically all human beings ...
 21 use what [C.A.] Campbell calls the unreflective conception when they are domi-
 22 nated by violent emotions like anger, indignation, or hate, and especially when
 23 the conduct they are judging has been personally injurious to them. On the other
 24 hand, a great many people, whether they are educated or not, will employ what
 25 Campbell calls the reflective conception when they are not consumed with hate
 26 or anger—when they are judging a situation calmly and reflectively and when
 27 the fact that the agent did not ultimately shape his own character has been viv-
 28 idly brought to their attention” (1958, 111). Edwards goes on to suggest that the
 29 “reflective conception” is the right one because the “unreflective conception” is
 30 driven by emotional bias.
- 31 7. Moreover, some participants may interpret the MR question in the concrete
 32 case as: “Should Bill be punished for his action?” and even if they do *not* think
 33 he has free will or “full moral responsibility,” they may think he needs to be
 34 punished for his multiple homicides. This interpretation might not be primed
 35 in the abstract case since there is no specific human action to be (potentially)
 36 punished. Finally, notice a subtle but important difference: Universe A does not
 37 mention *humans*, whereas Universe B explicitly mentions that the “one exception
 38 is human decision making,” and it concludes: “each human decision *does not*
 39 *have to happen* the way that it does.” This may prime some readers to think that
 40 Universe B is more like our universe, especially in the abstract condition, which
 41 is *not* then followed by the description of Bill, whose behavior suggests that he
 42 is human. In that case, some of the differences in results between N&K's abstract
 43 and concrete cases might also be explained by a difference in intuitions people
 44 have about moral responsibility when asked about an ‘alternate universe’ (A) vs.
 a ‘real-world universe’ (B), differences that have been demonstrated in Nahmias,
 Coates, and Kvaran (2007) and Nichols and Roskies (2008).
8. That is, the description suggests that determinism entails $[(Po \& L) \supset P]$, rather
 than the proper $[(Po \& L) \supset P]$.
9. Consider, by analogy, linguistic surveys intended to elicit people's intuitions concern-
 ing grammaticality. These surveys generally ask people to consider specific

- 1 sentences rather than asking them to consider abstract questions about whether
2 various constructions of sentences could be grammatical.
- 3 10. For more discussion of differences in judgments about FW and MR based on
4 abstract/concrete differences, as well as real world/alternate world differences, see
5 Nahmias, Coates, and Kvaran (2007) and Nichols and Roskies (2008).
- 6 11. Even NMNT's scenarios that do involve negative actions are not very "high-
7 affect." Stealing a necklace, robbing a bank, and keeping \$1000 found in a wallet
8 seem more similar to N&K's (2007) *low-affect* condition, in which Bill cheats on
9 his taxes, than N&K's *high-affect* condition, in which Bill murders his family.
- 10 12. For instance, in one version (the real world cases), when the agents were
11 described with the psychological predicates, 89 percent of participants said that
12 agents should be held morally responsible and 83 percent said they had free will,
13 whereas when the agents were described with the "neuro-reductionistic" predi-
14 cates, only 40 percent said they had MR and 38 percent said they had FW.
- 15 13. We call them "*prima facie* compatibilist intuitions" in part because we are not
16 committed to the idea that ordinary people have the (positive) intuition that
17 determinism is compatible with FW and MR—their intuitions may not be so
18 theoretically rich. But we think that lacking intuitions that (genuinely) support
19 incompatibilism is sufficient to say that people are "natural compatibilists." We
20 also accept that there may be alternative explanations that suggest people are
21 expressing only *apparent* compatibilist intuitions (see section V).
- 22 14. Participants were 436 undergraduate students in critical thinking or psychology
23 courses at Georgia State University who completed the entire survey. We excluded
24 187 participants prior to analysis who (a) responded incorrectly to either of two
25 comprehension questions or (b) completed the survey too quickly (less than one
26 half of one standard deviation from the mean time for completion), leaving 249
27 (42 percent male, 58 percent female) participants whose data we analyzed. Studies
28 were carried out under previous approval of the University's Institutional Review
29 Board.
- 30 15. This question was not asked in the NMNT surveys, though the following similar
31 statements were used: "If Universe C were re-created with *different* initial condi-
32 tions or *different* laws of nature, it is possible Jill would *not* [mow her lawn/steal
33 the necklace] at that time." Because of these differences, however, the Past
34 Different question was not used in the composite scores described below.
- 35 16. These composite scores provide a more robust measure of people's intuitions
36 concerning MR, FW, and bypassing. Lest one worry about averaging these scores,
37 and in so doing losing information about responses to each particular question,
38 responses to all questions factored into each composite score were, with very few
39 exceptions, highly positively intracorrelated with one another in all four survey
40 conditions. Across conditions, reliability analyses produced a Cronbach's alpha
41 of .807 among the questions used to compute the *MR/FW* composite score, and
42 a Cronbach's alpha of .823 among the questions used to compute the *Bypassing*
43 composite score, indicating that each composite score was strongly internally
44 consistent. Some might worry that the "no control" question is an inappropriate
measure to assess bypassing because they think that determinism *does* entail that
one has *no* control over what one does. We believe that this is mistaken on philo-
sophical grounds, but removing the "no control" question from the *Bypassing*
composite score also lowers the Cronbach's alpha among the questions used to
compute it from .823 to .797, which suggests that the "no control" question is
accessing the same folk concept as the other questions about bypassing.

- 1 17. Data reported in Figures 9.3 and 9.4 do not include the composite scores for
 2 20 of the 249 participants whose *Bypassing* composite scores were equal to the
 3 3.5 midpoint.
- 4 18. While we did replicate N&K's overall findings, our results were slightly different
 5 than theirs. In our study, 68 percent of participants in the abstract condition gave
 6 the apparent incompatibilist response that it is *not* possible for a person to be fully
 7 morally responsible in Universe A, compared to the 84 percent in N&K's previous
 8 study. Also, 87.5 percent of participants in our present study responded that Bill is
 9 fully morally responsible in N&K's concrete scenario, compared to 72 percent in
 10 N&K's study.
- 11 19. A *t*-test comparing the mean *Bypassing* responses in *N&K concrete* and *NMNT con-*
 12 *crete* was *not* significant, which explains why the ANOVA did not show a signifi-
 13 cant main effect for survey. We suspect that this lack of effect is due to the very
 14 high affect in the *N&K concrete* scenario driving down participants' *Bypassing*
 15 scores.
- 16 20. Because *N&K's concrete* scenario involves high affect in a way that *NMNT's concrete*
 17 scenario does not, we did not include either concrete scenario in the mediation
 18 analysis, as doing so would introduce another, potentially confounding, variable
 19 (high affect) in addition to condition (concrete vs. abstract) (see also note 20). AQ2
- 20 21. Mediation analysis involves the specification of a causal model between three
 21 variables. Suppose that an *initial variable*, *X* (in our case, survey type), is assumed
 22 to have a causal effect on an *outcome variable*, *Y* (in our case, MR/FW responses).
 23 Call *c* the direct effect of *X* on *Y*. A mediational model of the relationship between
 24 them, then, is one in which the causal effect of *X* on *Y* is *mediated* by an *interven-*
 25 *ing variable*, *M* (in our case, bypassing judgments). Call *a* the effect of the initial
 26 variable *X* on *M*, and *b* the effect of *M* on the outcome variable *Y*. *Complete*
 27 *mediation* obtains when variable *X* no longer has any direct effect on *Y* when *M* is
 28 controlled for, such that path *c* is zero. *Partial mediation* obtains when *c* is reduced
 29 if *M* is controlled for, but not to zero, because paths *a* and *b* account for some, but
 30 not all, of the overall causal effect of *X* on *Y*. Mediational models can be assessed
 31 statistically by mediation analysis, which uses multiple regression analyses to
 32 estimate the values of paths *a*, *b*, and *c*.
- 33 22. The mediating variable, *Bypassing*, explains 58 percent of the total effect of sur-
 34 vey type on MR/FW score (see Kenny, Kashy, & Bolger, 1998, 260–1).
- 35 23. Pearson correlation coefficient between Past Different ($M = 3.66$, $SD = 1.57$)
 36 and MR/FW composite scores in *N&K abstract*: $r(75) = -.277$, $p = .015$. By com-
 37 parison, in the *NMNT abstract* scenario only 7 percent disagreed with (that is,
 38 "missed") the following similar question ($M = 2.68$, $SD = 1.40$): "Suppose that
 39 in Universe C, a person named Jill decides to mow her lawn at a particular time
 40 and then does it. If Universe C were re-created with *different* initial conditions or
 41 *different* laws of nature, it is possible Jill would *not* mow her lawn at that time."
- 42 24. The only data we have that appears relevant does not support this hypothesis.
 43 In the *NMNT abstract* scenario we asked: "Suppose that in Universe C, a person
 44 named Jill decides to mow her lawn at a particular time and then does it. If
 Universe C were re-created with the *same* initial conditions and the *same* laws of
 nature, it is possible Jill would *not* mow her lawn at that time." (In *NMNT concrete*
 we asked whether it is possible Jill would *not* steal the necklace). The objection
 under consideration would predict there to be a significant correlation between
 responses to these questions and responses to the MR/FW questions—for example,
 people who say it *is* possible for Jill to do otherwise may be neglecting to see that
 determinism rules out AP and so should be more likely to attribute MR and FW

1 to her. However, there was no significant correlation between responses to this
 2 question and MR/FW composite scores in *NMNT abstract*: $r(54) = .137$, *ns*, though
 3 there was a marginally significant correlation in *NMNT concrete*: $r(58) = .242$,
 4 $p = .06$.

5 25. We would like to thank the following people for helpful comments on earlier
 6 drafts: Shaun Nichols, Al Mele, Jason Turner, Stephen Morris, Neil Levy, Tamler
 7 Sommers, George Graham, Dan Weiskopf, Joshua Knobe, Reuben Stern, Jason
 8 Shepard, Thomas Nadelhoffer, Trevor Kvaran, Fiery Cushman, David Blumenfeld,
 9 and especially Bradley Thomas. This chapter was completed in part with support
 10 from a grant (for E.N.) from the University of Chicago Arete Initiative and the
 11 John Templeton Foundation.

12 References

- 13 Baron, R. M. and Kenny, D. A. (1986), "The Moderator-Mediator Variable Distinction in
 14 Social Psychological Research: Conceptual, Strategic, and Statistical Considerations,"
 15 *Journal of Personality and Social Psychology* 51: 1173–82.
- 16 Blumenfeld, D. (1988), "Freedom and Mind Control," *American Philosophical Quarterly*
 17 25: 215–27.
- 18 Chalmers, D. (1996), *The Conscious Mind: In Search of a Fundamental Theory*. New York:
 19 Oxford University Press.
- 20 Edwards, P. (1958), "Hard and Soft Determinism," in S. Hook, ed., *Determinism and*
 21 *Freedom in the Age of Modern Science*, New York: University Press.
- 22 Fischer, J. and Ravizza, M. (1998), *Responsibility and Control: A Theory of Moral*
 23 *Responsibility*. Cambridge: Cambridge University Press.
- 24 Frankfurt, H. (1971), "Freedom of the Will and the Concept of a Person," in *The*
 25 *Importance of What We Care About* (Cambridge University Press, 1988), 11–25.
- 26 Kane, R. (1999), "Responsibility, Luck, and Chance: Reflections on Free Will and
 27 Indeterminism," *Journal of Philosophy* 96: 217–40.
- 28 Kauppinen, A. (2007), "The Rise and Fall of Experimental Philosophy," *Philosophical*
 29 *Explorations*. AQ3
- 30 Kenny, D. A. Deborah. A. K. and Bolger, N. (1998), "Data Analysis in Social
 31 Psychology," in D. Gilbert, S. Fiske, and G. Lindzey, eds, *The Handbook of Social*
 32 *Psychology: Vol. 1* (4th edn). Boston: McGraw-Hill, 233–65.
- 33 Knobe, J. and Nichols, S. (2008), "An Experimental Philosophy Manifesto," in
 34 J. Knobe and S. Nichols, eds, *Experimental Philosophy*, Oxford University Press. AQ4
- 35 Knobe, J. and Doris, J. (Forthcoming), "Strawsonian Variations: Folk Morality and
 36 the Search for a Unified Theory," in J. Doris, ed., *The Handbook of Moral Psychology*,
 37 Oxford University Press). AQ5
- 38 Lewis, D. (1986), *The Plurality of Worlds*. Oxford: Blackwell Publishers. AQ6
- 39 Lycan, W. (2003), "Free Will and the Burden of Proof," in Anthony O'Hear, ed.,
 40 *Proceedings of the Royal Institute of Philosophy for 2001–02*, Cambridge University
 41 Press, 107–22. AQ7
- 42 MacKinnon, D. P. Lockwood, C. M. Hoffman, J. M. West, S. G. and Sheets, V. (2002),
 43 "A Comparison of Methods to Test Mediation and other Intervening Variable
 44 Effects," *Psychological Methods* 7: 83–104.
- 45 Mele, A. (2005), *Autonomous Agents*. New York: Oxford University Press.
- 46 Nadelhoffer, T. and Nahmias, E. (2007), "The Past and Future of Experimental
 47 Philosophy," *Philosophical Explorations* 10.2: 123–49.
- 48 Nahmias, E. (2006), "Folk Fears about Freedom and Responsibility: Determinism vs.
 49 Reductionism," *Journal of Cognition and Culture* 6: 215–37.

- 1 Nahmias, E. Morris, S. Nadelhoffer, T. and Turner, J. (2005), "Surveying Freedom: Folk
- 2 Intuitions about Free Will and Moral Responsibility," *Philosophical Psychology* 18:
- 3 561–84.
- 4 Nahmias, E. Morris, S. Nadelhoffer, T. and Turner, J. (2006), "Is Incompatibilism
- 5 Intuitive?" *Philosophy and Phenomenological Research* 73: 28–53.
- 6 Nahmias, E. Coates, J. and Kvaran, T. (2007), "Free Will, Moral Responsibility, and
- 7 Mechanism: Experiments on Folk Intuitions," *Midwest Studies in Philosophy* 31:
- 8 214–42.
- 9 Nichols, S. and Knobe, J. (2007), "Moral Responsibility and Determinism: The
- 10 Cognitive Science of Folk Intuitions," *Nous* 41: 663–85.
- 11 Nichols, S. and Roskies, A. (Forthcoming), "Bringing Moral Responsibility Down to AQ8
- 12 Earth."
- 13 O'Connor, T. (2005), "Freedom With a Human Face," *Midwest Studies in Philosophy*,
- 14 29, 207–27.
- 15 Pereboom, D. (2001), *Living Without Free Will*. Cambridge University Press. AQ9
- 16 Strawson, G. (1986), *Freedom and Belief*. Oxford: Clarendon.
- 17 Turner, J. and Nahmias, E. (2006), "Are the Folk Agent Causationists?" *Mind and*
- 18 *Language* 21: 597–609.
- 19 van Inwagen, P. (1983), *An Essay on Free Will*. Oxford: Clarendon Press.
- 20 Vargas, M. (2005), "The Revisionist's Guide to Responsibility," *Philosophical Studies*
- 21 125: 399–429.
- 22 Watson, G. (1975), "Free Agency," in G. Watson, ed., *Free Will*, Oxford University
- 23 Press, 96–110.
- 24 Weinberg, J. (2007), "How to Challenge Intuitions Empirically Without Risking
- 25 Skepticism," *Midwest Studies in Philosophy* 31: 318–43.
- 26 Wolf, S. (1990), *Freedom within Reason*. Oxford University Press. AQ10
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44

QUERY FORM

BOOK TITLE:	AGUILAR
CHAPTER NO:	Chapter 9

Queries and / or remarks

Query No.	Query / remark	Response
AQ1	Retain italics here?	
AQ2	it is maybe unclear what you mean by see note 20 as this is note 20?	
AQ3	Please provide other publication details.	
AQ4	Place of publication, please.	
AQ5	Please update info, if available.	
AQ6	Place of publication, please.	
AQ7	Place of publication, please.	
AQ8	Please update info, if available.	
AQ9	Place of publication, please.	
AQ10	Place of publication, please.	