



# Frankfurt cases and the Newcomb Problem

Arif Ahmed<sup>1</sup>

© The Author(s) 2019

**Abstract** A standard argument for one-boxing in Newcomb's Problem is 'Why Ain'cha Rich?', which emphasizes that one-boxers typically make a million dollars compared to the thousand dollars that two-boxers can expect. A standard reply is the 'opportunity defence': the two-boxers who made a thousand never had an opportunity to make more. The paper argues that the opportunity defence is unavailable to anyone who grants that in another case—a Frankfurt case—the agent is *deprived* of opportunities in the way that advocates of Frankfurt cases typically claim.

**Keywords** Decision Theory · Newcomb's Problem · Frankfurt cases · Free will

Causal Decision Theory (CDT) and Evidential Decision Theory (EDT) are two leading theories of rational choice. Briefly and informally, CDT says that a rational agent does whatever in her view most effectively *brings about* her ends. EDT says that a rational agent does whatever in her view is the best *evidence* of what she wants.<sup>1</sup>

CDT and EDT are both intuitive but not both true. They conflict over cases like:

---

<sup>1</sup> More formally (though still simplifying): let the probability function  $Cr$  and the news-value function  $V$  respectively represent the agent's credence (degree of belief) in, and desirability for, an arbitrary proposition, propositions being subsets of the set  $W$  of possible worlds. Let  $O = \{o_1 \dots o_n\}$  partition  $W$  into propositions describing the agent's options, and let  $Z = \{z_1 \dots z_m\}$  partition  $W$  into propositions that each specify the outcome in as much detail as concerns the agent. Let some counterfactual-like operator ' $>$ ' reflect causal dependence:  $o_i > z_j$  says that if the agent were to realize  $o_i$  then the outcome would be  $z_j$ . CDT says that the rational agent maximizes  $U$ , and EDT says that the rational agent maximizes  $V$ .

---

✉ Arif Ahmed  
ama24@cam.ac.uk

<sup>1</sup> Faculty of Philosophy, Cambridge University, Cambridge, UK

**Table 1** Payoffs in Newcomb's Problem

	$S_1$ : predicted $o_1$	$S_2$ : predicted $o_2$
$o_1$ : take only opaque box	$m$	0
$o_2$ : take both boxes	$m + k$	$k$

*Newcomb's Problem.* You must choose between ( $o_1$ ) taking only an opaque box ('one-boxing') and ( $o_2$ ) taking both it and a transparent one in which \$1000 ( $k$ ) is visible ('two-boxing'). You get to keep what you take. The opaque box already contains \$1 M ( $m$ ) if and only if ( $S_1$ ) a highly reliable predictor foresaw your taking only it. If ( $S_2$ ) she foresaw your taking both boxes then the opaque box is empty. So the payoff schedule (in dollars, for which we assume linear increasing utility) is as follows (Table 1).

You can't *affect* what is in the opaque box, this having been settled at the time of the prediction. So taking both boxes causes you to be richer by  $k$  than if you had taken only the opaque box. CDT therefore counts  $o_2$  as uniquely rational. But taking only the opaque box is strong evidence that you get  $m$ , and taking both boxes is strong evidence that you get  $k \ll m$  (since the predictor is very reliable). EDT therefore counts  $o_1$  as uniquely rational. CDT and EDT disagree over Newcomb's Problem.

Probably the most popular argument for one-boxing, and against CDT, is that one-boxing maximizes expected return. Since the predictor reliably foresees both one- and two-boxing, 'one-boxers' make on average just under  $m$  per trial, whereas 'two-boxers' make on average more than  $k$  but considerably less than  $m$  per trial. So followers of EDT can ask followers of CDT: 'If you're so smart, why ain'cha rich?'

To this argument, known as 'Why ain'cha rich?' or WAR, defenders of CDT have a simple reply. They say that its lower expected return has no bearing on the *rationality* of two-boxing, because two-boxers typically lack money-making *opportunities* that one-boxers typically enjoy. Most two-boxers face an opaque box that never contained any money, and most one-boxers face an opaque box that always contained \$1 M. Call this the *opportunity defence* of CDT.<sup>2</sup>

This paper defends the following material conditional: *if* agents lack opportunities in *another* case, a 'Frankfurt case', *then* the opportunity defence fails for a version of Newcomb's Problem that Sections 1–4 construct. For those who know the Frankfurt cases, the argument is that Black's role in the relevant Frankfurt cases involves prediction in a way that makes those cases interpretable as Newcomb problems. (For those who do not know the Frankfurt cases, Section 2 outlines the relevant points.)

---

Footnote 1 continued

where the  $U$  - and  $V$  -scores of an option  $o \in O$  are  $U(o) = \sum_{z \in Z} V(z)Cr(o > z)$  and  $V(o) = \sum_{z \in Z} V(z)Cr(z|o)$  respectively. For more sophisticated accounts see e.g. (Lewis 1981a; Joyce 1999 ch. 4, 5).

<sup>2</sup> E.g. (Joyce 1999: 153; Wells 2017: 5).

The conditional will prompt different people to draw different further conclusions. Anyone who maintains its antecedent will agree that the opportunity defence cannot blunt WAR. Anyone who maintains the opportunity defence will have to say that certain opportunities *are* available in Frankfurt cases. But even someone who is or becomes agnostic on both questions might find interest in the connection that the conditional establishes between the debate over practical rationality on which Newcomb's problem bears, and the debate over free will on which Frankfurt cases bear.

1. The construction starts with a base case, then modifies it in three steps. Before describing the base case I'll sketch the idea behind *Frankfurt cases*.

There is a debate over whether determinism is compatible with any kind of freedom that entails moral responsibility. One argument that they are incompatible turns on what Frankfurt called the:

*Principle of Alternate Possibilities (PAP)*. A person is responsible for doing something only if he could have done otherwise.<sup>3</sup>

If PAP is true, and if determinism entails that anyone who has done something could not have done otherwise, then determinism is incompatible with anyone's being morally responsible for what they did.

Frankfurt's case against PAP involves examples where a person is morally responsible for something that she could not have avoided doing. Such *Frankfurt cases* originally took the following form.

Suppose someone – Black, let us say – wants Jones to perform a certain action. Black is prepared to go to considerable lengths to get his way, but he prefers to avoid showing his hand unnecessarily. So he waits until Jones makes up his mind what to do, and he does nothing unless it is clear to him (Black is an excellent judge of such things) that Jones is going to decide to do something other than what [Black] wants him to do. If it becomes clear that Jones is going to decide something else, Black takes effective steps to ensure that Jones decides to do, and that he does do, what [Black] wants him to do. Whatever Jones's initial preferences and inclinations, then, Black will have his way... Now suppose that Black never has to show his hand because Jones, for reasons of his own, decides to perform and does perform the very action Black wants him to perform... In this example there are sufficient conditions for Jones's performing the action in question. What action he performs is not up to him... He has no alternative but to do what Black wants him to do. If he does it on his own, however, his moral responsibility for doing it is not affected by the fact that Black was lurking in the background with sinister intent, since this intent never comes into play.<sup>4</sup>

Frankfurt holds that Jones was responsible for what he did even though he could not have done otherwise. If so, PAP fails.

---

<sup>3</sup> Frankfurt (1969): 167.

<sup>4</sup> Frankfurt (1969): 173–4.

But arguably, the case only threatens both PAP *and* relevantly similar principles if Black's procedure meets certain constraints. After all, even if Jones can't get around *doing* what Black wants, he could still *try* or *choose* not to do it. Black would then force his hand – but still, the fact that Jones *could* have thus tried or chosen might suffice to make him responsible for his actual, wholehearted performance.<sup>5</sup> That may not save PAP itself: Jones still performs the act, is responsible for doing so and has no alternative to *performing* it. But it does immunize nearby principles, such as this variant of PAP: an agent is responsible for a performance only if she could have acted blamelessly in the circumstances.<sup>6</sup> That principle seems well-placed to do what the incompatibilist wanted PAP to do.

In response, defenders of Frankfurt's idea offer '*prior sign*' cases where Black can intervene *before* Jones chooses whether to do what Black wants, and so also before he can try not to. Black can do this because he can observe some precursor of Jones's choice. This example concerns the 2008 US Presidential election.

Because he dares to hope that the Democrats finally have a good chance of winning the White House, the benevolent but elderly neurosurgeon, Black, has come out of retirement to participate in yet another philosophical example... He has secretly inserted a chip in Jones's brain that enables Black to monitor and control Jones's activities. Black can exercise this control through a sophisticated computer that he has programmed so that, among other things, it monitors Jones's voting behaviour. If Jones were to show any inclination to vote for McCain (or, let us say, anyone other than Obama), then the computer, through the chip in Jones's brain, would intervene to assure that he actually decides to vote for Obama and does so vote. But if Jones decides on his own to vote for Obama (as Black, the old progressive would prefer), the computer does nothing but continue to monitor—without affecting—the goings-on in Jones's head. Now suppose that Jones decides to vote for Obama on his own, just as he would have if Black had not inserted the chip in his head. It seems, upon first thinking about this case, that Jones can be held morally responsible for his choice and act of voting for Obama, although he could not have chosen otherwise and he could not have done otherwise.<sup>7</sup>

Let Jones vote by pressing a button labelled 'M' (McCain) or a button labelled 'O' (Obama). Let Black's chip work by detecting a *necessary but not sufficient* precursor of Jones's being about to choose to press 'M'.<sup>8</sup> Black's readiness to act on this necessary condition is meant to ensure not only that Jones could not have voted for McCain but also that he could not have chosen to vote for McCain. If Jones is still morally responsible, that rules out not only PAP but also any variant of it that

---

<sup>5</sup> McKenna (2008): 774.

<sup>6</sup> Otsuka (1998).

<sup>7</sup> Fischer (2010): 316. For another '*prior sign*' case see e.g. (Sartorio 2016: 13).

<sup>8</sup> Black's oversight of a necessary but not sufficient condition is typical of Frankfurt cases. See (Levy 2008: 225).

ties moral responsibility for an act to the possibility of *choosing* or *trying* to act otherwise.

The base case is a ‘prior sign’ Frankfurt case. But its bearing on PAP is irrelevant. It won’t matter to my main argument whether Jones is morally responsible for what he does (although I shall revisit that point in Section 10). What matters is that Jones could not have *chosen or done* otherwise.

2. Start with a version of Fischer’s example—call it **F1**. Jones’s brain-state’s instantiating some neurological pattern  $N^*$  is a necessary but not sufficient condition of Jones’s being about to decide (say, in the next seven seconds) to press ‘M’. Black has secretly implanted into Jones’s brain a chip that detects  $N^*$ . If it does, it warns Black, who immediately intervenes through the chip to ensure that Jones chooses to vote for Obama (and does vote for Obama). As it happens, the chip never alerts Black, who never interferes, and Jones votes for Obama as he intended all along.

Grant the following:

(1) In F1 Jones *lacks* any opportunity to vote for McCain.

(1) is the antecedent of the target conditional: *If* Jones lacks opportunities to act otherwise in this case, *then* there is a Newcomb Problem for which the opportunity defence fails. So I may *suppose* (1) true without argument. Still, conditionals with wildly implausible antecedents have little interest, so I should motivate (1).

The intuition is that it suffices for foreclosure of an opportunity, that one’s choice be subject to the wrong kind of external control. It is hard to say precisely what counts as the ‘wrong kind’. Intuitively, *persuading* you not to buy something doesn’t count as denying you the opportunity to buy it. Nor does modestly raising its price. But neurosurgical intervention counts, as does a steep enough increase in price—like pointing a gun at your head. Still, it seems plausibly sufficient for *Black’s* having denied Jones the opportunity to vote for McCain that Black is intentionally and systematically ‘linking conditions so that all the conditions that are necessary for the opportunity to be acted upon are never jointly satisfied’.<sup>9</sup> That *does* look like the wrong kind of external control.

In F1 the following conditions are nomologically necessary for Jones’s voting for McCain: (i) Jones’s brain-state instantiates  $N^*$ ; (ii) Jones doesn’t choose to vote for Obama. In F1 Black is systematically linking (i) and (ii) so that they are never jointly satisfied. If the chip detects that (i) is true then Black makes (ii) false. Black has therefore denied Jones the opportunity to vote for McCain in F1. So (1) is true.

Next consider case **F2**: the chip only *monitors* Jones. It alerts Black if Jones’s brain-state instantiates  $N^*$ . But Black can’t influence Jones via the chip or otherwise. I claim:

(2) In F2 Jones *has* an opportunity to vote for McCain.

Mere *recording* of Jones’s brain-state, which is all that either does *or could* happen here, leaves his opportunities just as they would have been if neither Black nor the chip had been there. But if neither Black nor the chip had been there then Jones

<sup>9</sup> Dennett (2015): 129.

would have had an opportunity to vote for McCain, if anyone ever has an opportunity to do something that she does not actually do. So (2) is true.<sup>10</sup>

3. To summarize where we are: so far, I have introduced a ‘prior sign’ Frankfurt case (F1), and a variant case (F2) where Black is observing but not controlling the subject. The next case modifies F1 by changing the locus of Black’s control: instead of being able to manipulate his choice, Black can only manipulate the *effect* of his choice. I’ll argue that a comparison with F1 reveals that Jones *lacks* certain opportunities, but a comparison with F2 reveals that he now *has* certain others.

In F3 then, Black *can* intervene. But if the chip detects N\* then Black interferes not with Jones but with *the voting machine*. He remotely alters its program so that pressing *either* button records a vote for Obama. As before, Jones votes for Obama without interference, because his brain-state never instantiates N\*, and Black never intervenes.

If (1) and (2) are true then so are:

(3) In F3 Jones *lacks* any opportunity to vote (i.e. to record a vote) for McCain.

(4) In F3 Jones *has* an opportunity to press ‘M’.

My case for (3) is that F3 shares with F1 those features that make (1) true: Black is controlling Jones’s environment so as to prevent joint realization of all the necessary conditions for Jones’s voting for McCain. These necessary conditions include: (i) Jones’s brain-state instantiates N\*; (ii) the machine is running normally. If the chip detects that (i) is true then Black makes (ii) false. Black is intentionally and systematically preventing it from being the case that Jones votes for McCain. So (3) is true.

The case for (4) is that F3 shares with F2 those features that make (2) true. What makes it true in F2 that Jones can vote for McCain is that no external agent is stopping him. In F2, the chip only *records* N\*. Black does nothing to influence Jones’s subsequent vote. Equally in F3, if the chip records N\*, Black does nothing to prevent Jones from pressing ‘M’. He intervenes if at all only on the causal connection between the pressing of ‘M’ and the subsequent recording of a vote. If anyone ever has an opportunity to do anything that she doesn’t actually do, then Jones has an opportunity in F3 to press ‘M’. So (4) is true.

4. So far, I have argued that a comparison with F2 and the base case shows that Jones has certain opportunities in F3 and lacks certain others. I’ll now modify F3, but not in a way that makes any structural difference to Jones’s opportunities. All that changes are the possible outcomes. The result is a Newcomb Problem.

This fourth case, F4, resembles F3 except that the machine doesn’t record votes but transfers money. Jones chooses whether to press ‘M’ or ‘O’. As things are initially set up, pressing ‘M’ has no effect. But pressing ‘O’ transfers \$1000 to Jones’s account. As before, a necessary condition of Jones’s being about to choose to press ‘M’ is that his brain-state instantiates N\*.

<sup>10</sup> Incompatibilists will say that if determinism is true then (2) fails: nobody ever has an opportunity to do otherwise. But none of the cases that I discuss require determinism: see Section 7.

**Table 2** Payoffs in F4

	$S_1$ : Black detects $N^*$	$S_2$ : Black does not detect $N^*$
$o_1$ : press 'M'	$m$	0
$o_2$ : press 'O'	$m + k$	$k$

Again, a chip in Jones's brain can reliably detect  $N^*$ . But this time Black's overriding intention is to reward restraint. More particularly, he wants to award \$1 M to anyone who presses 'M'. So if the chip detects  $N^*$ , Black remotely reprograms the machine so that Jones's pressing *either* button causes \$1 M to be wired to Jones's account, on top of the possible \$1000 that Jones stands to gain by pressing 'O'. Finally, we suppose that in F4 Jones knows about Black's powers and intentions.

As it happens, Black doesn't act. Jones presses 'O'. His brain-state never instantiates  $N^*$ . He ends up \$1000 richer.

It is as true in F4 as in F3 that Black is systematically controlling Jones's environment to prevent the joint realization of the necessary conditions for something: not for Jones's voting for McCain, but for his realizing any outcome in which he presses 'M' and does *not* get \$1 M. Jones's knowing that and how Black is controlling his situation has no tendency to diminish this control. Nor does it *give* Black control over Jones's choice between pressing 'M' and pressing 'O'.

Claim: if (3) and (4) are true then so are:

- (5) In F4 Jones *lacks* any opportunity to: press 'M' and not make \$1 M.
- (6) In F4 Jones *has* an opportunity to press 'M'.

The case for (5) is that F4 shares with F3 those features that make (3) true in F3. Two necessary conditions for the realization of this conjunction (i.e. Jones presses 'M' and does *not* make \$1 M) are: (i) Jones's brain-state at some point instantiates  $N^*$ ; (ii) the machine is running normally. If the chip detects that (i) is true, Black makes (ii) false. Jones therefore lacks any opportunity to: press 'M' and not make a million. F4 is a Frankfurt case where Black is determined that anyone who presses 'M' is rewarded: and to get his way, he is prepared to act to ensure that Jones does not both press 'M' and miss out on the big prize. So (5) is true.

The case for (6) is that F4 shares with F3 those features that make (4) true in F3. As in F3, if the chip detects  $N^*$ , Black does nothing to stop Jones *from pressing 'M'*. Jones has an opportunity to press 'M' if anyone ever has an opportunity to do anything that she doesn't do. So (6) is true.

5. F4 is a Newcomb Problem where pressing 'M' corresponds to one-boxing and pressing 'O' corresponds to two-boxing. The payoffs in F4 are as follows (Table 2):

In F4 Jones chooses the two-boxing option (i.e. he presses 'O').

In F4, pressing 'O' dominates. If Black *has* detected  $N^*$  then Jones gets at least \$1 M. But pressing 'O' returns an additional \$1000. If Black *has not* detected  $N^*$ , then pressing 'M' leaves Jones empty-handed, but pressing 'O' returns \$1000.

Moreover, whether Black has *already* detected  $N^*$  when Jones chooses is *causally* independent of Jones's choice, because it precedes Jones's choice by about seven seconds. CDT therefore prefers pressing 'O'.

But what Jones chooses to do *is evidentially* relevant to whether Black has detected  $N^*$ , given Jones's evidence when he chooses. Certainly, if Jones chooses to press 'M', Black has detected  $N^*$ . But if Jones chooses to press 'O', it is less than certain that Black has detected  $N^*$ , and may be most unlikely according to Jones's credences. If Jones thinks that Black detects  $N^*$  on fewer than 99.9% of occasions on which the agent presses 'O', EDT recommends pressing 'M'.<sup>11</sup>

This completes the construction of the advertised Newcomb Problem i.e. F4. Here as in all Newcomb Problems, it is foreseeable that the average return to 'one-boxing' (pressing 'M') exceeds the average return to 'two-boxing' (pressing 'O'). At any rate this is so if Black detects  $N^*$  on all occasions on which the agent presses 'M' and on fewer than 99.9% of occasions on which the agent presses 'O'. 'One-boxers', who all make \$1 M per trial, can therefore press the WAR objection against Jones and other 'two-boxers', who on average make less, possibly much less, than \$1 M per trial.

6. The next step argues that in *this* Newcomb Problem the 'opportunity defence' is not available. Here is the informal argument. For formalization in a weak modal logic, see the Appendix.

(5) and (6) imply that Jones has an opportunity to press 'M' but not to press 'M' without making \$1 M. But if Jones *can* in this sense press 'M' but *can't* press 'M' without making \$1 M, then he *can* get \$1 M. Therefore, Jones *has* an opportunity to make \$1 M. He cannot invoke the opportunity defence.

Or if he does, then he must explain which of (5) or (6) he rejects, and why. (6) looks unassailable; and I argued that (5) is true if (1) is. Hence my conditional conclusion: *if* in Frankfurt cases like F1 Black has deprived Jones of any opportunity to realize what Black wants to avert, *then* F4 is a Newcomb case where the opportunity defence is unavailable.

That concludes the main argument. I turn to three objections.

7. Nothing in these examples presupposes determinism about Jones's choices, in the sense that the state of the world at any earlier time nomologically determines Jones's choice. It is not true, for instance, that  $N^*$  nomologically necessitates that Jones will choose to press 'M' in F4, nor that the absence of  $N^*$  nomologically necessitates that Jones will choose to press 'O'. The absence of  $N^*$  at any time only necessitates that Jones will not choose to press 'M' *in the next few seconds*; this is nomologically consistent with his doing so later.

But even that limited determinism, not about Jones's ultimate choice but about any choice he makes in a given time window, might seem to threaten (6). The objection is that *when* he presses 'O' Jones *lacks* any opportunity to press 'M', because his brain-state was not recently instantiating  $N^*$ . Nor did he have that

<sup>11</sup> EDT strictly prefers pressing 'M' iff  $m > k + mCr(S_1|O_2)$  i.e. iff  $Cr(S_1|O_2) < 0.999$ .

opportunity *earlier* if his brain-state never instantiated  $N^*$ . So Jones can complain (contrary to (6)) that he *never* had an opportunity to press ‘M’.<sup>12</sup>

There are three responses. The most defensive response is that the objection would only appeal to someone who thinks that *any* pre- $t$  state of the world, that nomologically determines that one will not choose an option at  $t$ , eliminates any opportunity to choose that option at  $t$ . One might instead think that the *only* determining factors that are genuinely opportunity-depriving are those involving external impediments or external interference by another. On this view, the fact that your present choice is *somehow* determined by your past (internal) brain-state cannot by itself deprive you of any present opportunity to do whatever that brain-state rules out. Thus (apparently) Dennett, who grants in support of (5) that Frankfurt-style interference deprives Jones of an opportunity,<sup>13</sup> whilst maintaining—in support of (6)—that mere determination-by-the-past does not.

[A] real opportunity is an occasion where a self-controller ‘faces’ (is informed about) a situation in which the outcome of its subsequent ‘deliberation’ will be a decisive (as we say) factor. In such a situation more than one outcome is ‘possible’ so far as the agent or self-controller is concerned; that is, the critical nexus passes through its deliberation. That is what we mean when we say that the outcome is up to the agent.<sup>14</sup>

Anyone who agrees with Dennett will accept (5) but reject this objection to (6).

Waiving that, the second point is that even if Jones lacked any opportunity to press ‘M’ when he chose to press ‘O’, it doesn’t follow that he lacked any opportunity to press ‘M’. When he pressed ‘O’ he could instead have waited. And there need have been nothing about the state of his brain or of the world at that time that was nomologically inconsistent with his pressing ‘M’ later. So even if Jones lacked any opportunity to press ‘M’ when he chose to press ‘O’, he had an opportunity to press ‘M’, and he had an opportunity to make \$1 M.<sup>15</sup>

The third, more tentative point is that it may be unnecessary to suppose that the prior appearance of  $N^*$  is a *necessary* condition for Jones to press ‘M’. Suppose instead that  $N^*$  is ‘nearly’ necessary: at the outset there is a very high *chance*  $p$  that his brain-state instantiates  $N^*$  at some later time, given that he does at some later time choose to press ‘M’; and that there is a lower chance  $q < p - \frac{k}{m}$  that his brain-state ever instantiates  $N^*$ , given that he chooses to press ‘O’. If these conditional chances control the corresponding frequencies then this indeterministic case is a Newcomb Problem where pressing ‘M’ generates a foreseeably higher average return than pressing ‘O’. If Jones makes \$1 K, can he now plead the opportunity defence?

<sup>12</sup> Cf. Widerker (1995): 251.

<sup>13</sup> Dennett (2015): 144–5. Here Dennett is discussing the original type of Frankfurt case, where Black stands ready to manipulate Jones after observing Jones’s choice. But there is reason to think he would say the same about the ‘prior sign’ cases under discussion here: see especially p. 129 as quoted in §2.

<sup>14</sup> Dennett (2015): 129.

<sup>15</sup> Cf. Hunt (2005): 135.

Clearly (6) still holds, because now Jones's choice is undetermined by *anything* in its past. Everything turns on whether (5) holds i.e. on whether, in this indeterministic case, Black has deprived Jones of an opportunity to: press 'M' and not make \$1 M. Has he? To see why he has, compare these cases.

- (a) Black shuts Jones in a room from which the only exit is a door that Black has secured with a steel padlock, and Jones doesn't have a key.
- (b) Black shuts Jones in a room from which the only exit is a door that Black has secured with a combination lock with 100 million settings, and Jones doesn't know the combination.

Could anyone say with a straight face that Jones has an opportunity to leave the room in (b) though not in (a)? It may in *some* sense be possible that Jones leaves in (b) but not in (a). For instance, it may be that laws of nature plus simple physical facts rule out Jones's getting past the padlock in (a), but not his hitting on the right combination in (b). But this sort of 'opportunity' lacks real interest.<sup>16</sup> For all practical purposes, Black *has* ensured that Jones stays put in (b). For that matter, if quantum tunnelling is physically possible then it is equally possible in that sense that Jones leaves in (a).

(6) can therefore tolerate some injection of chanciness into F4. Black can ensure that Jones does not achieve this or that outcome without making anything physically impossible. For Black to do that, it suffices that (i) it doesn't happen because (ii) Black has made it sufficiently unlikely. But then even in this *indeterministic* version of F4, Black has eliminated Jones's opportunity to get anything other than \$1 M by pressing 'M'. Since Jones *has* an opportunity to press 'M', it follows that he has an opportunity to make \$1 M, so cannot plead the opportunity defence.

Somebody might object that there is, in this indeterministic version, something that Jones knows how to do—pressing 'M'—such that, if he were to do it, then he *would* be falsifying the material conditional  $o_1 \rightarrow m$ , at least on a non-back-tracking reading of 'If he were to do it...'.<sup>17</sup> But then in (b), there is something that he knows how to do—entering the sequence 01271756—such that, if he were to do it, then he would be opening the combination lock. If it is consistent with the latter that Black has locked Jones into the room in (b) then it is consistent with the former that Black has locked Jones out of falsifying  $o_1 \rightarrow m$  in F4.

For that matter there is, on the non-back-tracking reading, something that Jones knows how to do in the original Frankfurt case—pressing 'M'—such that if he were to do it, he would be voting for McCain. So anyone who takes this line against (5) must also reject (1). The objection is therefore impotent against the *conditional: if*

<sup>16</sup> Hume's famous example is another case where we regard an opportunity as having been foreclosed without regarding its realization as *physically* impossible. 'A prisoner who has neither money nor interest, discovers the impossibility of his escape, as well when he considers the obstinacy of the gaoler, as the walls and bars with which he is surrounded; and, in all attempts for his freedom, chooses rather to work upon the stone and iron of the one, than upon the inflexible nature of the other' (Hume 1975: 90).

<sup>17</sup> A *non-back-tracking* conditional  $\supset$  is true iff all the closest worlds at which *a* is true (the closest '*a*-worlds') are *b*-worlds, where the closest *a*-worlds (i) match actuality over all particular matters of fact that are causally independent of whether *a* is true; and (ii) contain no violations of the actual laws except possibly for a small, local miracle that shortly precedes and also brings about the realization of *a*.

Jones lacks the opportunities that he is supposed to lack in ‘prior sign’ Frankfurt cases, *then* the opportunity defence fails for both F4 and this indeterministic version of it.

A chancy case like this may be feasible. In a recent study, fMRI subjects would ‘fixat[e] on the centre of [a] screen where a stream of letters was presented. At some point, when they felt the urge to do so, they were to freely decide between one of two buttons, operated by the left and right index fingers, and press it immediately’.<sup>18</sup> Experimenters found that:

two brain regions encoded with high accuracy whether the subject was about to choose the left or right response prior to the conscious decision... The first region was in frontopolar cortex, BA10. The predictive information in the fMRI signals from this brain region was already present 7 s. before the subject’s motor decision.<sup>19</sup>

Suppose we find that for almost all subjects who choose to press the left-hand button, an indication of this is present in the relevant region of BA10 a few seconds before they choose. Similarly, suppose that for most subjects who choose to press the right-hand button, an indication of this is present in the same region a few seconds beforehand. We could arrange things so that pressing the right-hand button caused the release of \$1000 to the subject’s bank account and pressing the left-hand button caused the termination of the experiment without this reward. Black could then further arrange to release \$1 M to the subject’s bank account as soon as fMRI detected that the subject would press the left-hand button.<sup>20</sup> In this case all subjects have, but only those who press the left-hand button realize, an opportunity to win \$1 M.

8. A second objection is that F4 is not a Newcomb case at all. A genuine Newcomb case involves exactly two options: ‘one-boxing’ and ‘two-boxing’. But F4 offers Jones *three* options at any time. He can press ‘M’ now; he can press ‘O’ now; and he can *wait* before choosing.

This makes a difference because CDT seems to recommend waiting. After all, the *effect* of waiting can only be to make N\* more likely to occur at *some* time before Jones chooses what to press. But F4 is a Newcomb case only if CDT recommends pressing ‘O’.

But even if CDT sometimes recommends waiting, this doesn’t help CDT. Whether or not F4 is strictly a Newcomb case, followers of CDT will wait either forever (and never win anything) or until the marginal costs of waiting exceed the marginal benefits of increasing the probability that N\* is at some point realized in their brains, and at that point press ‘O’ and end up with \$1000. In either case those

<sup>18</sup> Soon et al. (2008): 543.

<sup>19</sup> Soon et al. (2008): 544.

<sup>20</sup> The \$1 M reward would not be immediate: given current constraints, the time required to determine that activity in BA10 was symptomatic of impending left-button-pressing exceeds the time-lapse between that activity and the subject’s choice. But it is still a Newcomb case. Even if the reward is subsequent to the choice it is still causally independent of the choice, because it depends only on activity that strictly preceded the choice.

agents end up worse than followers of EDT, and in either case they can't complain that they lacked the opportunity to win \$1 M.

Besides, it isn't clear that CDT *does* recommend waiting. If (as we might suppose) N\* only ever occurs within a fixed time interval preceding the subject's choice, there is no reason to expect waiting to increase the probability that N\* precedes whatever choice the subject ultimately makes. Compare: my space capsule only deploys a parachute when it is a minute away from landing. Waiting before starting my descent (rather than descending immediately) doesn't make it more likely that the parachute is deployed at some time before landing.

But if CDT doesn't recommend waiting (and of course neither does EDT) then the only plausible options at any time are pressing 'M' and pressing 'O', so F4 is after all a Newcomb case.

9. The third objection involves a notion of opportunity that falsifies (5).

**C-opportunity:** At any time:

- (i) You have a *counterfactual opportunity* (a *C-opportunity*) to realize the truth of a proposition  $p$  by choosing  $o$  iff  $o$  is an option at that time such that if you *had* at that time realized  $o$ , then  $p$  *would* have been true.<sup>21</sup>
- (ii) You have a *C-opportunity* to  $\varphi$  by choosing  $o$  iff you have a C-opportunity to realize the truth of the proposition that you  $\varphi$  by choosing  $o$ .
- (iii) You have a *C-opportunity* to  $\varphi$  iff one of your options is such that you have a C-opportunity to  $\varphi$  by choosing it.

Crucially, the counterfactual in (i) is a non-back-tracking conditional  $> : a > b$  is true iff all of the closest possible worlds at which  $a$  is true (the closest ' $a$ -worlds') are  $b$ -worlds, where the closest  $a$ -worlds (i) match the actual world over all particular matters of fact that are causally independent of whether  $a$  is true; and (ii) contain no violations of the actual laws except possibly for a small, local miracle that very shortly precedes and also brings about the realization of  $a$ .<sup>22</sup>

On this definition of opportunity, (5) is false when Jones chooses to press 'O'. At that time, he has the option to press 'M'. And if—in the non-back-tracking sense—Jones *had* then pressed 'M', it would have been true that: he presses 'M' and fails to

<sup>21</sup> For present purposes I treat an *option* as a proposition that the agent can make true if she pleases (Jeffrey 1983: 84), i.e. one that would be true if the agent were to choose to make it true, where this counterfactual has a non-back-tracking sense. One problem with this proposal as a general definition is that a proposition's being an option in this sense is neither necessary nor sufficient for the agent to have the *ability* to realize it, since abilities are dispositions that can be masked or finked through external interference (Vihvelin 2004: 437 ff.). But that makes no difference to the cases under consideration, where Black does not, and is not disposed to, interfere with the causal connection between Jones's choice and any of his options.

<sup>22</sup> For present purposes, this analysis of the non-backtracking counterfactual is equivalent to Lewis's (1979a, b) theory, though like Edgington's improvement upon it (Edgington 2004), it relies explicitly on the notion of causal independence. The analysis implies that we have a very wide range of counterfactual opportunities. For instance (and as Lewis notes), if determinism is true then we have abundant C-opportunities to realize the truth of the proposition that a law of nature is broken (Lewis 1981b).

make \$1 M. Jones therefore has a C-opportunity to press ‘M’ and fail to make \$1 M.

Similarly, if opportunities are just C-opportunities then Jones lacks any opportunity to make \$1 M. This is because when he chooses, Black has not detected  $N^*$  and has not transferred \$1 M to Jones’s account. These facts are causally independent of what Jones chooses. So if—in the non-back-tracking sense—Jones had chosen either option, it would *not* have been the case that he made \$1 M. Jones therefore lacks a C-opportunity to make \$1 M.<sup>23</sup>

I concede that if opportunities are C-opportunities then (5) is false. But I deny that this undermines the argument. I defended the material conditional: *if* Jones lacks an opportunity to defy Black in a Frankfurt case like F1, then he has an opportunity to make \$1 M in a Newcomb case like F4. If opportunities are C-opportunities then the consequent of this conditional is false, but so is its antecedent. The conditional is still true.

The antecedent is false in F1 because the implant has *not* detected  $N^*$  when Jones chooses to vote for Obama. So Black is *not* intervening. This fact is causally independent of Jones’s choice. So at the closest worlds at which Jones chooses to vote for McCain, Black is not intervening. At these worlds Jones presses ‘M’ and votes for McCain. So Jones has a C-opportunity to vote for McCain. Anyone who thinks that opportunities are C-opportunities should deny that Black eliminates them in ‘prior sign’ Frankfurt cases.<sup>24</sup> But this is no objection to the conditional defended here: *if* Jones lacks an opportunity to defy Black in a Frankfurt case like F1, *then* he has an opportunity to make  $m$  in a Newcomb case like F4.<sup>25</sup>

10. The next step in this dialectic – which I cannot pursue at length – turns on a question that this conditional raises. We saw two kinds of opportunity in play: C-opportunity, and the kind of opportunity that Black’s presence, intentions and

<sup>23</sup> Supposing that Jones has an option of waiting doesn’t change this: it is not true that if he had waited then his brain-state would at some later point instantiated  $N^*$  (see Section 8—although perhaps if he had waited then his brain-state *might* have instantiated  $N^*$ : that is, if  $o_3$  is the proposition that he waits at the time that he actually chooses to press ‘O’, and  $n^*$  is the proposition that his brain-state instantiates  $N^*$  at some later time, then perhaps  $\sim(o_3 > \sim n^*)$ .) So neither is it true that if he had waited then he would have had a C-opportunity to make \$1 M. Besides, if Jones *did* have a C-opportunity to make \$1 M in F4 then that only helps WAR.

<sup>24</sup> Cf. Smith (2004): 100 ff. For illuminating further discussion, see (Whittle 2010: 12 ff).

<sup>25</sup> Another reason to care about the difference, between C-opportunities and those being denied in Frankfurt cases, arises from Pettit’s distinction between Liberal and Republican concepts of negative liberty. I am free in the Liberal sense to the extent that others don’t interfere with me. I am free in the Republican sense to the extent that this non-interference is *resilient* i.e. not down to the whim of any non-interfering other—she *could* not have interfered had she so chosen (Pettit 1993: 24–8). For instance, suppose I live in a state where it is common knowledge that all crimes are detected and severely punished; but my neighbour owns a gun. This state is depriving my neighbour of a ‘Frankfurtian’ opportunity to shoot me. After all, one necessary condition of my neighbour’s shooting me is that he *chooses* to shoot me, and the state is preventing him from choosing to shoot me by maintaining and advertising its forensic and penal capacities. But it isn’t depriving him of a C-opportunity to shoot me, because *nothing* is preventing him from shooting me *if he chooses*. This state preserves my freedom in the Liberal sense but not in the Republican sense.

capacities are denying to Jones in the Frankfurt cases—call it ‘F-opportunity’.<sup>26</sup> The question is whether lack of C-opportunity to do better is by itself *sufficient* for exculpation. If so, then Jones cannot be blamed for making only \$1000, because he *lacks* any C-opportunity to make more. On the other hand, if lack of C-opportunity is insufficient, then the opportunity defence is at best incomplete, because as we saw, it appeals only to Jones’s lack of C-opportunity, not to any lack of F-opportunity.

One *prima facie* basis for thinking that lack of C-opportunity cannot by itself be exculpatory arises from reflection on something like the *original* Frankfurt case (not the ‘prior sign’ cases). Suppose that Black plans to wait until Jones has decided what to do, and then if necessary—if he has decided to vote for McCain—to coerce or force him into voting for Obama. Jones then lacks any C-opportunity to vote for McCain, but he is still responsible for voting for Obama if he does it wholeheartedly. Similarly, we might suspect that someone who autonomously chooses to exercise, and does exercise, the ‘two-boxing’ option in the Newcomb Problem is responsible for his (relative) impoverishment. (Of course this is at best a sketch of an argument).<sup>27</sup>

It might also help to examine moral Newcomb Problems: cases like F4 where the payoffs carry moral significance. The literature hasn’t covered these extensively, but one notable recent treatment is due to MacAskill, who offers the following:

Box A is opaque; Box B, transparent. If the Predictor predicts that you choose Box A only, then he puts one wish into Box A. With that wish, you would save the lives of one million terminally ill children. If he predicts that you choose both Box A and Box B, then he will put nothing into Box A. Box B—transparent to you—contains a stick of gum. You have two options only: Choose Box A, or choose both Box A and Box B.<sup>28</sup>

I am inclined to agree with MacAskill’s own intuition in favour of one-boxing here; more to the present point, it seems reasonable to blame any two-boxer for betting with other people’s lives that he could outwit the predictor. (‘But when I chose to two-box I *couldn’t* have saved them!’—‘But in the sense in which you couldn’t outwit the predictor, you *could* have saved them.’) If so, the lack of any C-opportunity to do better than one actually did cannot in general be exculpatory.

<sup>26</sup> One way to cash out the distinction is in terms of the interpretation of counterfactuals. If, in the definition of C-opportunity, we read the counterfactual in (i) as non-backtracking, we get the intended reading. But if we read it as suitably backtracking, in the manner of Horgan (1981), we may get a parallel notion of F-opportunity. For a recent discussion connecting the EDT/CDT dispute and the interpretation of counterfactuals, see (Price and Liu 2018 §2.3). (Of course the conditional that the main part of this essay defends is true on a *uniform* application of either interpretation).

<sup>27</sup> It might matter, for instance, that in the standard Frankfurt case Jones doesn’t *know* that he lacks a C-opportunity to do otherwise, whereas in the Newcomb Problem he may know that he lacks a C-opportunity to make \$1 M. But if the predictor is imperfect (but still very good) in the Newcomb Problem then he *doesn’t* know this; nor does he know it if he hasn’t yet decided what to do. Obviously the matter invites further investigation.

<sup>28</sup> MacAskill (2016): 429.

Obviously there is more to be said about the interplay of ethical, decision-theoretic and forensic reasoning that this case involves, and also about whether the conditional that this essay defends serves better as the major premise of a modus ponens or as the major premise of a modus tollens (if either). Still, the conditional itself is true: whatever kind of opportunity is being foreclosed in the Frankfurt cases is *not* being denied to two-boxers like Jones in the Newcomb case that we studied. And I hope that it gets us closer to resolving at least one of the two long-standing philosophical problems that it connects.

**Acknowledgements** I am most grateful to Yael Loewenstein, Jack Spencer and Robert Stalnaker for helpful discussion, and to two referees for very helpful comments on an earlier draft of this paper. I am also grateful to the Leverhulme Trust, which supported me via a Research Fellowship (REF-2018-231) during the writing of this paper.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

## Appendix

This appendix argues formally for the claim in Section 5 that Jones has an opportunity to make \$1 M in F4.

I assume two schematic principles of the logic of opportunity ascriptions. For any proposition  $p$ , write  $\nabla p$  to mean that Jones has an opportunity to realize an outcome in which  $p$  holds. Write  $\Delta p \stackrel{\text{def.}}{=} \sim \nabla \sim p$  for its dual:  $\Delta p$  means that Jones *lacks* any opportunity to *avoid* an outcome in which  $p$  holds. I'll assume that any instances of the following schemata are true if  $t$  is replaced by any tautology and  $p$  and  $q$  by any propositions.

- (7)  $\Delta t$   
 (8)  $\Delta(p \rightarrow q) \rightarrow (\Delta p \rightarrow \Delta q)$

These are theorems of the weakest modal logic  $K$ .<sup>29</sup> In English, (7) says that Jones cannot falsify a tautology and (8) says that if Jones cannot achieve  $p$  without  $q$ , and cannot avoid  $p$ , then he cannot avoid  $q$ .<sup>30</sup>

Now if every instance of (7) and (8) is true then so is every instance of:

<sup>29</sup> The axioms of  $K$  are (a) all tautologies (b) every instance of (8). The rules are *modus ponens* and necessitation: if  $p$  is a theorem then so is  $\Delta p$ .

<sup>30</sup> (8) resembles van Inwagen's beta, which says that from  $Np$  and  $N(p \rightarrow q)$  we may infer  $Nq$ , where  $Np$  says that  $p$  is true and nobody had any choice about it (van Inwagen 1983: 94). One objection to beta is that it entails *agglomeration*:  $(Np \& Nq) \rightarrow N(p \& q)$ , which is false. For a counterexample (McKay and Johnson 1996: 115-6), suppose this coin was not tossed, so ( $p$ ) it didn't land heads and ( $q$ ) it didn't land tails. Nobody had a choice about whether it landed heads or about whether it landed tails, so  $Np \& Nq$ , but somebody *did* have a choice about whether it was tossed, so  $\sim N(p \& q)$ .

Could one similarly object to (8)? No: although (8) implies an agglomeration principle  $(\Delta p \& \Delta q) \rightarrow \Delta(p \& q)$ , the counterexample is less telling. Consider a coin that Jones did not toss but might have tossed. It seems that Jones had an opportunity to achieve an outcome in which it landed heads,

$$(9) \Delta(p \rightarrow q) \rightarrow (\nabla p \rightarrow \nabla q)$$

This says that if Jones cannot achieve  $p$  without  $q$ , but *can* achieve  $p$ , then he can achieve  $q$ . Proof: any instance of  $\Delta((p \rightarrow q) \rightarrow (\sim q \rightarrow \sim p))$  is an instance of (7). By an instance of (8), it follows that  $\Delta(p \rightarrow q) \rightarrow \Delta(\sim q \rightarrow \sim p)$  and by another instance of (8) that  $\Delta(p \rightarrow q) \rightarrow (\Delta \sim q \rightarrow \Delta \sim p)$ . The latter is logically equivalent to  $\Delta(p \rightarrow q) \rightarrow (\sim \Delta \sim p \rightarrow \sim \Delta \sim q)$ . From this, it follows by the definition of  $\Delta$  and the interchange of logical equivalents licensed by (7) and (8) that  $\Delta(p \rightarrow q) \rightarrow (\nabla p \rightarrow \nabla q)$ .

(5) and (6) say respectively that Jones cannot press ‘M’ without getting \$1 M and that he can press ‘M’. In formal terms, and writing  $m$  for the proposition that Jones gets \$1 M:

$$(10) \Delta(o_1 \rightarrow m)$$

---

Footnote 30 continued

and so an outcome in which  $p$  fails, just as anyone who can buy a ticket in a fair lottery has an opportunity to win it. So  $(\Delta p \& \Delta q) \rightarrow \Delta(p \& q)$  is true because its antecedent is false.

More generally, it seems plausible that for *some* closeness relation  $m$  on possible worlds, we have:  $\Delta q \leftrightarrow \forall p(\nabla p \rightarrow (p >_m q))$ , where ‘ $p > q$ ’ is true if and only some  $(p \& q)$  world is  $m$ -closer to actuality than any  $(p \& \sim q)$ -world: that is,  $p$  is inevitable for Jones iff it would be true *whichever* of his opportunities he were to realize. Then whatever the  $m$  turns out to be, (8) is also true. Informal proof: for any  $m$ , if  $p \rightarrow q$  is true at all  $m$ -closest  $o$ -worlds for any opportunity  $o$ , and if  $p$  is true at all  $m$ -closest  $o$ -worlds for any opportunity  $o$ , then  $q$  is true at all  $m$ -closest  $o$ -worlds for any opportunity  $o$ . The analogous principle for  $N$  would be  $Nq \leftrightarrow \forall p(Mp \rightarrow (p >_m q))$ , writing ‘ $M$ ’ for the dual of  $N$ . This principle seems false for any reading of ‘ $>_m$ ’ on which it is false that: if Jones had tossed this coin it *would not* have landed tails.

Objection: I said in Section 7 that Jones *lacks* any opportunity to enter the right combination into a combination lock. What is the difference between (i) buying a ticket in a fair lottery and (ii) entering *some* number into a combination lock, such that my having an opportunity to do (i) gives me an opportunity to win, but Jones’s having an opportunity to do (ii) does not give *him* an opportunity to leave the room?

Reply: getting the right combination is, and buying a winning lottery ticket isn’t, *quasi-miraculous* in the sense of suggesting some conspiracy (Lewis 1986: 58 ff.). In a big lottery, the fact that my ticket is the winner has little tendency to suggest a conspiracy; whereas Jones’s getting the combination right—or a monkey’s spontaneously typing a 950-page dissertation on anti-realism—though perhaps no more unlikely, *is* evidence of collusion or similar.

If Jones had been one of millions of prisoners in the same situation, or if the monkey had been one of *very* many monkeys in the same situation, then it would have been different. If Jones is the only prisoner then the prior probability of  $(c_j)$  there being a conspiracy involving Jones is relatively high, so  $(r_j)$  Jones’s getting the right combination raises this probability to a high level. But if Jones is one of many prisoners then the prior probability of a conspiracy involving *Jones in particular* is relatively low, so the posterior probability of a conspiracy involving Jones is also low, although the Bayes factor  $Pr(r_j|c_j)/Pr(r_j|\sim c_j)$  is the same.

There are arguments that we should not put the difference between (i) and (ii) in terms of quasi-miracles but rather in terms of the mathematically more sophisticated notion of *typicality* (Williams 2008: 406 ff.; Elga 2004: 71; Gaifman and Snir 1982: 544). However, typicality supports the distinction equally well. Jones’s buying a winning ticket in a lottery, or tossing a fair coin which then lands heads, are both typical in the relevant sense, but his hitting on the right combination is *not* typical.

Note finally that the avoidance of quasi-miracles or of atypicality should also figure in whatever closeness measure  $m$  witnesses the counterfactual constraint  $\Delta q \leftrightarrow \forall p(\nabla p \rightarrow (p >_m q))$ . That means that if Jones lacks an opportunity to leave in (b), then since Jones has the opportunity to enter some combination or other, the  $m$ -closest worlds at which he enters some combination or other are worlds at which he does not enter the right combination. This is no objection to the constraint: it just reveals a conceptual connection between quasi-miracles or typicality and opportunities.

(11)  $\nabla o_1$ 

But it follows from an instance of (9) together with (10) and (11) that  $\nabla m$  i.e. Jones *does* have an opportunity to make \$1 M.

## References

- Dennett, D. (2015). *Elbow room: The varieties of free will worth wanting* (2nd ed.). Cambridge, MA: MIT Press.
- Edgington, D. (2004). Counterfactuals and the benefit of hindsight. In P. Dowe & P. Noordhof (Eds.), *Cause and chance: Causation in an indeterministic world* (pp. 12–27). London: Routledge.
- Elga, A. (2004). Infinitesimal chances and laws of nature. *Australasian Journal of Philosophy*, 82, 67–76.
- Fischer, J. M. (2010). The Frankfurt cases: The moral of the stories. *Philosophical Review*, 119, 315–336.
- Frankfurt, H. (1969). Alternate possibilities and moral responsibility. *Journal of Philosophy* 66:829–839. In G. Watson (Ed.), *Free will* (Vol. Second, pp. 167–177). Oxford: OUP 2003.
- Gaifman, H., & Snir, M. (1982). Probabilities over rich languages, testing and randomness. *Journal of Symbolic Logic*, 47, 495–548.
- Horgan, T. (1981). Counterfactuals and Newcomb's problem. *Journal of Philosophy*, 78, 331–356.
- Hume, D. (1975). *Enquiries concerning human understanding and concerning the principles of morals*, ed. L. A. Selby-Bigge. Oxford: Clarendon Press.
- Hunt, D. P. (2005). Moral responsibility and buffered alternatives. *Midwest Studies in Philosophy*, 29, 126–145.
- Jeffrey, R. C. (1983). *The logic of decision* (2nd ed.). Chicago: Chicago UP.
- Joyce, J. M. (1999). *Foundations of causal decision theory*. Cambridge: CUP.
- Levy, N. (2008). Counterfactual intervention and agents' capacities. *Journal of Philosophy*, 105, 223–239.
- Lewis, D. K. (1979a). Counterfactual dependence and time's arrow. *Noûs*, 13, 455–476.
- Lewis, D. K. (1979b). Postscript to Lewis. In his *Philosophical Papers, Vol. II*. Oxford: OUP.
- Lewis, D. K. (1981a). Causal decision theory. *Australasian Journal of Philosophy*, 59, 5–30.
- Lewis, D. K. (1981b). Are we free to break the laws? *Theoria*, 47, 113–121. Reprinted in his *Philosophical Papers, Vol. II*. Oxford: OUP 1986: 291–298.
- Lewis, D. K. (1986). Postscripts to Lewis 1979. In D. K. Lewis (Ed.), *Philosophical Papers* (Vol. II, pp. 52–66). Oxford: OUP.
- MacAskill, W. (2016). Smokers, psychos, and decision-theoretic uncertainty. *Journal of Philosophy*, 123, 425–445.
- McKay, T., & Johnson, D. (1996). A reconsideration of an argument against compatibilism. *Philosophical Topics*, 24, 113–122.
- McKenna, M. (2008). Frankfurt's argument against alternative possibilities: Looking beyond the examples. *Noûs*, 42, 770–793.
- Otsuka, M. (1998). Incompatibilism and the avoidability of blame. *Ethics*, 108, 685–701.
- Pettit, P. (1993). Negative liberty, liberal and republican. *European Journal of Philosophy*, 1, 15–38.
- Price, H., & Liu, Y. (2018). 'Click!' bait for causalists. In A. Ahmed (Ed.), *Newcomb's problem* (pp. 160–179). Cambridge: CUP.
- Sartorio, C. (2016). *Causation and free will*. Oxford: OUP.
- Smith, M. (2004). A theory of freedom and responsibility. In M. Smith (Ed.), *Ethics and the a priori: Selected essays on moral psychology and meta-ethics* (pp. 84–113). Cambridge: CUP.
- Soon, C. S., Brass, M., Heinze, H.-J., & Haynes, J. D. (2008). Unconscious determinants of free decisions in the human brain. *Nature Neuroscience*, 11, 543–545.
- van Inwagen, P. (1983). *An essay on free will*. Oxford: Clarendon Press.
- Vihvelin, K. (2004). Free will demystified: A dispositional account. *Philosophical Topics*, 32, 427–450.
- Wells, I. (2017). Equal opportunity and Newcomb's problem. *Mind*. <https://doi.org/10.1093/mind/fzx018>.
- Whittle, A. (2010). Dispositional abilities. *Philosophers' Imprint*, 10, 1–23.
- Widerker, D. (1995). Libertarianism and Frankfurt's attack on the principle of alternate possibilities. *Philosophical Review*, 104, 247–261.

Williams, J. R. G. (2008). Chances, counterfactuals and similarities. *Philosophy and Phenomenological Research*, 77, 385–420.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.