

For Kane, ed. *The Free Will Handbook*, O.U.P.  
January 9, 2010

**Who's *Still* Afraid of Determinism?  
Rethinking Causes and Possibilities<sup>1</sup>**

Christopher Taylor and Daniel Dennett  
Center for Cognitive Studies  
Tufts University

Incompatibilism, the view that free will and determinism are incompatible, subsists on two widely accepted, but deeply confused, theses concerning possibility and causation: (1) in a deterministic universe, one can never truthfully utter the sentence “I could have done otherwise,” and (2) in such universes, one can never really receive credit or blame for having caused an event, since in fact all events have been predetermined by conditions during the universe’s birth. Throughout the free will literature one finds variations on these two themes, often intermixed in various ways. When Robert Nozick<sup>2</sup> describes our longing for “originative value” he apparently has thesis (2) in mind, and thesis (1) may underlie his assertion that “we want it to be true that in that very same situation we could have done (significantly) otherwise.” John Austin, in a famous footnote, flirts with thesis (1):

Consider the case where I miss a very short putt and kick myself because I could have holed it. It is not that I should have holed it if I had tried: I did try, and missed. It is not that I should have holed it if conditions had been different: that might of course be so, but I am talking about conditions as they precisely were, and asserting that I could have holed it. There is the rub. Nor does ‘I can hole it this time’ mean that I shall hole it this time if I try or if anything else; for I may try and miss, and yet not be convinced that I could not

have done it; indeed, further experiments may confirm my belief that I could have done it that time, although I did not.<sup>3</sup>

(In later sections we discuss at length the ways in which this particular quote can lead readers astray.) Meanwhile, Robert Kane, in *The Significance of Free Will*, eloquently proclaims the importance of our presumed ability truly to cause events, the ability that thesis (2) addresses:

Why do we want free will? We want it because we want ultimate responsibility. And why do we want that? For the reason that children and adults take delight in their accomplishment from the earliest moments of their awakening as persons, whether these accomplishments are making a fist or walking upright or composing a symphony.<sup>4</sup>

Elsewhere in the free will debate one often finds authors advancing definitions that confirm the relevance of possibilities and causes. Kane describes free will itself, for instance, as “the power of agents to be the ultimate creators ... and sustainers of their own ends and purposes.”<sup>5</sup> The key words here are “power” and “creator.” Intuition suggests that the term “power” is intertwined with “possibility” roughly as follows: agent A has the power to do X if and only if it is possible that A does (will do) X. And certainly to be a “creator” one has to be the *cause* of changes in the world; one has to “make a difference” in how the world runs. .

Carl Ginet in a similar vein proposes:

Two or more alternatives are *open to me* at a given moment if which of them I do next is entirely up to my choice at that moment: Nothing that exists up to that moment stands in the way of my doing next any one of the alternatives.<sup>6</sup>

In short, the acceptance of theses (1) and (2) lies at the heart of incompatibilism. To avoid the sort of impasse that incompatibilists have apparently reached, we propose to reexamine the foundations of possibilities and causes, to understand why theses (1) and (2) look so compelling. We will tackle first possible worlds, then counterfactuals, and then causation, building each concept on its predecessors. We will then compare this account with the recent important work of Judea Pearl on causality, before turning to the implications of our approach for the two incompatibilist theses. We will discover that the desires incompatibilists describe, to have powers and to effect changes, can be satisfied without any recondite appeals to quantum indeterminacy. The suspicions to the contrary lose their force once we begin to untangle, with the aid of a little formalism, the complexities of the underlying concepts.

### Possible Worlds

While a complete account of possible worlds would require many extra pages, the following paragraphs outline an approach, compatible with modern scientific methods, that avoids various modal pitfalls identified by Quine (such as talk of “propositions,” “analyticity,” “essences,” etc.).<sup>7</sup> Ideally, science strives for a description of the universe that is as thorough and comprehensive as possible, composed in an orderly mathematical idiom. A simple example

of such ideal state-descriptions are the “Democritean” universes introduced by Quine.<sup>8</sup> A Democritean universe is completely specified using a function  $f$  that assigns to each quadruple  $(x,y,z,t)$  a value of either 0 or 1. If  $f(x,y,z,t) = 1$ , then at time  $t$  matter occupies location  $(x,y,z)$ ; otherwise point  $(x,y,z)$  is devoid of matter at  $t$ . Needless to say, modern physics has long since supplanted the tidy Democritean conception of reality, but even today the basic project of describing the world with (monstrously complex) functions remains intact. So despite its scientific shortcomings, the following definition provides a useful starting point as we struggle to discipline unruly pretheoretical intuitions:

*A possible world* is simply any function of the form described above (in mathematical notation, any function of the form  $f: \mathbb{R}^4 \rightarrow \{0,1\}$  ).<sup>9</sup>

The set of all possible worlds we will denote by  $\Omega$ ; a particularly noteworthy subset of  $\Omega$  is  $\Phi$ , which contains just the *physically or nomologically possible* worlds, in which no physical laws are violated.<sup>10</sup>

Given a possible world  $f$  we of course have many ways to describe and make assertions about it. Often it will be natural to postulate *entities* within  $f$ : connected hypersolids in  $\mathbb{R}^4$  that yield coherent life-histories for objects like stars, planets, living creatures, and everyday paraphernalia. One will also want to set up a system of *informal predicates* that apply to these entities, such as “has a length of 1 meter,” “is red”, “is human.”<sup>11</sup> We may then form sentences like

$\exists x (x \text{ is human})$

and determine whether they apply in various different possible worlds (while recognizing that

often enough one will encounter borderline worlds where incontestible verdicts prove elusive).

Worthy of special note are *identification predicates* of the form “is Socrates.” “Is Socrates,” we shall suppose, applies to any entity in any possible world that shares so many features with the well-known denizen of the actual world that we are willing to consider it “the same person.” In the actual world, of course, “is Socrates” applies to exactly one entity; in others, there may reside no such being, or one, or conceivably several to whom the predicate applies equally well. Like other informal predicates, identification predicates suffer from vagueness and subjectivity, but they do not cause unusual problems.

With this machinery in place we can now explicate such sentences as:

Necessarily, Socrates is mortal. (1)

We would propose the translation:

In every (physically?) possible world  $f$ , the sentence “ $\forall x (x \text{ is Socrates} \Rightarrow x \text{ is mortal})$ ” obtains. (2)

Here “is Socrates” and “is mortal” are informal predicates of the sort just introduced. Paraphrase (2) strikes us as both plausible and free of the logical confusions Quine decries. Of course, deciding whether (2) is true does present considerable challenges, stemming largely from the unavoidable blurriness of the predicates. Moreover, we are not specifying the set of possible worlds over which one should allow  $f$  to range; perhaps some readers will advocate set  $\Omega$  (all worlds), others  $\Phi$  (the physically possible worlds), and yet others a still more restricted set  $X$ . Logic alone cannot resolve this issue, but logical language does help us to pinpoint such

questions and recognize the sorts of vagueness we face. However we choose, we can employ the notation

$$\Box_X \varphi$$

to indicate that sentence  $\varphi$  obtains for every world in set  $X$ .

As the dual of necessity, possibility yields to a similar analysis. Hence

Possibly, Socrates might have had red hair. (3)

means

There exists (within some set  $X$ ) a possible world  $f$  in which the sentence “ $\exists x (x$  is Socrates  $\wedge x$  has red hair)” obtains. (4)

Analogous to the notation “ $\Box_X \varphi$ ” we introduce

$$\Diamond_X \varphi ,$$

meaning that  $\varphi$  holds for some world within  $X$ . The familiar sentence:

Austin could have holed the putt (5)

now becomes

$\Diamond_X \exists x (x \text{ is Austin} \wedge x \text{ holes the putt}).$  (6)

Notice that in this case we need to restrict  $X$  to a narrow range of worlds, all quite similar to actuality, if we are to do justice to Austin’s meaning. For suppose that Austin is an utterly incompetent golfer, and that impartial observers are inclined to deny (5). If we let  $X$  range too widely, we may include worlds in which Austin, thanks to years of expensive lessons, winds up a championship player who holes the putt easily, thus validating (6) but distorting the presumed sense of (5). At the same time, as we shall see, there is no good reason to make  $X$  so small that

only worlds *identical* to reality in the moments before the putt get included.

### Counterfactuals

Using possible worlds, one can also profitably interpret sentences of the form

If you had tripped Arthur, he would have fallen, (7)

as David Lewis has shown.<sup>12</sup> Roughly, (7) obtains if and only if in every world approximately similar to our own where the antecedent holds, so does the consequent. In other words:

$\Box_X \varphi \Rightarrow \psi$ , (8)

where  $\varphi$  stands for “you tripped Arthur,”  $\psi$  stands for “Arthur fell,” and X is a set of worlds similar to our own. As an alternative notation, let us also write:

$\varphi \Box_{\rightarrow X} \psi$ . (9)

Choosing an optimal value for X in (8) and (9) is not always easy, but we suggest the following loose Guidelines:

In sentences like (8) and (9), X ought to:

- contain worlds in which  $\varphi$  holds,  $\sim\varphi$  holds,  $\psi$  holds, and  $\sim\psi$  holds
- contain the worlds that are otherwise very similar to the actual world (insofar as the preceding clause permits). (G)

So when analyzing (7), choose X to contain worlds in which you trip Arthur, worlds where you refrain from tripping him, worlds where he falls, and worlds where he remains upright. In the case of (10):

If the sun hadn't risen this morning, I would have overslept, (10)

X will look quite different, since it includes strange worlds in which the sun fails to rise.

In *Counterfactuals*, Lewis cleverly devises a single connective  $\Box\rightarrow$  appropriate for all  $\phi$  and  $\psi$ , but in this paper we settle for a family of connectives of type  $\Box\rightarrow_X$ . Doing so, we believe, forestalls various technical complications and accords equally well with intuition. Notice that for Lewis transitivity fails, and, worse, so does the equivalence

$$\phi \Box\rightarrow \psi \equiv \sim\psi \Box\rightarrow \sim\phi.$$

With each operator  $\Box\rightarrow_X$ , on the other hand, transitivity and contraposition succeed, provided we hold X fixed.<sup>13</sup>

Notice also that the Guidelines mitigate the difficulties raised by Fine (1975) in his famous review of Lewis. Fine considered the sentence: “Had Nixon pressed the button, a nuclear war would have started.” He observed that a world in which an electrical malfunction prevented the button from unleashing the warheads was intuitively much closer to actuality than one containing a nuclear holocaust—such glitches happen frequently, while nuclear wars are both rare and exceedingly disruptive. Someone who evaluates Fine’s sentence using Lewis’s generalized notion of similarity appears forced to deny the truth of the sentence, but our guidelines, by requiring the inclusion of worlds in which nuclear wars occur (in which  $\psi$  occurs), permit us to disregard the magnitude of the difference a nuclear war makes, while at the same time discouraging the introduction of convenient *dei ex machina* such as broken wires. Admittedly, our notions of similarity are vague, but we contend that this vagueness is an inherent feature of counterfactual thinking, and as we shall see, even the most rigorous treatments apparently rely on it.<sup>14</sup>

## Causation

Fundamental as it appears, the language of causation has stirred up interminable debate and has (perhaps for that reason) been avoided by scientists. Many philosophers apparently hope some day to unearth the one “true” account of causation, but given the informal, vague, often self-contradictory nature of the term, we think a more realistic goal is simply to develop a formal analogue (or analogues) that helps us think more clearly about the world. Our preexisting hunches about causation will provide some guidance, but we should mistrust any informal arguments that masquerade as “proofs” validating or debunking particular causal doctrines.<sup>15</sup>

When we make an assertion like

Betty’s tripping Arthur caused him to fall, (16)

a number of factors appear to be at work supporting the claim. In an approximate order of importance, we list the following:

- Causal necessity. At least since Hume, philosophers have suspected that counterfactuals play some role in our causal thinking, and this factor and the next fall within the same tradition. Our assent to sentence (16) depends on our conviction that in any world roughly similar to our own in which Arthur falls, Betty must have tripped him up. Using the notation of the previous section, we have  $\psi \Box_{\rightarrow X} \phi$ , where  $\phi$  stands for “Betty tripped Arthur,”  $\psi$  represents “Arthur fell,” and  $X$  is a set of worlds similar to our own containing instances where (i) Betty trips Arthur, (ii) Betty doesn’t trip him, (iii) Arthur falls, or (iv) he doesn’t fall. As observed above, the sentence  $\sim\phi \Box_{\rightarrow X} \sim\psi$  has the same

logical force; in other words, had Betty not tripped Arthur, he would not have fallen.

- Causal sufficiency. It may well be that whenever we affirm (16), we do so partly because we believe that (using the same notation as before)  $\phi \Box_{\times} \psi$ . In other words, we believe that Arthur's fall was an *inevitable outcome* of Betty's tripping: in any world where Betty places the obstruction in his path, Arthur goes toppling. (Or equivalently, if Arthur had *not* fallen, then Betty must in that case have refrained.) This second condition is logically entirely distinct from the first, and yet the two seem to get badly muddled in everyday thinking. Indeed, as we shall see, incompatibilist confusion often originates precisely here. Below we will discuss at greater length the relations between these two crucial conditions.
- Truth of  $\phi$  and  $\psi$  in the actual world. Although a relatively trivial requirement, it should be mentioned if only for completeness.
- Independence. We expect the two sentences  $\phi$  and  $\psi$  to be logically independent: there must exist worlds, however remote from reality, in which  $\phi$  obtains but not  $\psi$ , and vice versa. Hence "Mary's singing and dancing caused her to dance and sing" has a decidedly odd ring. This condition also helps rule out "1+1=2 causes 2+2=4."
- Temporal priority. A reliable way to distinguish causes from effects is to note that causes occur earlier.<sup>16</sup>

- Miscellaneous further criteria. Although less critical than the preceding points, a number of other conditions may increase our confidence when we make causal judgements. For instance, in textbook examples of causation,  $\phi$  often describes the actions of an agent, and  $\psi$  represents a change in the state of a passive object (as in “Mary causes the house to burn down”). Further, we often expect the two participants to come into physical contact during their transaction.<sup>17</sup>

In order to understand these conditions better, and in particular the distinction between necessity and sufficiency, let us try them out on a few test cases (some of which derive from Lewis).<sup>18</sup> First consider the sharpshooter aiming at a distant victim. Scrutiny of the sharpshooter’s past record shows that the probability of a successful hit in this case is 0.1; if it makes any difference, we might imagine that irreducibly random quantum events in the sharpshooter’s brain help determine the outcome. Let us suppose that in the current case the bullet actually hits and kills the victim. We unhesitatingly agree then that the sharpshooter’s actions caused the victim’s death, *despite their causal insufficiency*. Accordingly, it appears that in cases like these, people rank necessity above sufficiency when making judgments about causes.

Still, sufficiency does retain some relevance. Suppose that the king and the mayor both have an interest in the fate of some young dissident; as it happens, both issue orders to exile him, so exiled he is. This is a classic case of *over-determination*. Let  $\phi_1$  stand for “the king issues an

exile order,”  $\varphi_2$  stand for “the mayor issues an exile order,” and  $\psi$ , “the dissident goes into exile.” In the current scenario, neither  $\varphi_1$  nor  $\varphi_2$  alone is necessary for  $\psi$ : for instance, had the king failed to issue any order, the dissident would still have been exiled thanks to the mayor, and vice versa. In fact  $\varphi_1 \vee \varphi_2$  satisfies the necessity requirement, but we are (perhaps unreasonably) reluctant to posit a disjunction as a cause.<sup>19</sup> Instead, sufficiency comes to the rescue and permits a choice between the two. After all,  $\varphi_2$  fails this test: it is easy to imagine a universe where the mayor issues his decree, yet the dissident gets off (just change the king’s order into a pardon). The king’s order, on the other hand, is truly *effective*; whatever small changes we make to the universe (including changes in the mayor’s orders), the dissident’s exile follows from the king’s command. Accordingly we may dub  $\varphi_1$  the “real cause” (if we feel the need to satisfy that yearning).<sup>20</sup>

Finally, consider the tale of Billy and Susie. Both children are throwing rocks at a glass bottle, and as it happens Susie’s rock, traveling slightly faster, reaches the bottle first and shatters it. Billy’s rock arrives a moment later at exactly the spot where the bottle used to stand, but of course encounters nothing but flying shards. When choosing between  $\varphi_1$  (“Susie throws rock S”) and  $\varphi_2$  (“Billy throws rock B”), we vote for  $\varphi_1$  as the cause of  $\psi$  (“The bottle shatters”), despite the fact that neither sentence is necessary (had Susie not thrown her rock, the bottle would still have shattered thanks to Billy, and vice versa) and both are sufficient (Billy’s throw suffices to produce a broken bottle, whatever his playmate does, and likewise with Susie’s). Why? The general notion of temporal priority (introduced above in connection with distinguishing cause from effect) strikes us as one critical consideration. As with priority

disputes in science, art, and sports, we seem to put a premium on being the *first* with an innovation, and since rock S arrived in the vicinity of the bottle earlier than rock B, we give credit to Susie. Further, it is clear that, although the bottle would still have shattered without Susie's throw, the shattering event would have been significantly different, occurring at a later time with a different rock sending fragments off in different directions. We can choose set X to reflect this fact (in keeping with guidelines (G)): let it contain worlds in which either (1) the bottle doesn't shatter at all, or (2) it shatters in a way very similar to the way it shatters in reality. Then for every world in X,

$$\psi \Rightarrow \varphi_1$$

obtains; wherever in X the bottle shatters, we find Susie throwing her rock first. On the other hand,

$$\psi \Rightarrow \varphi_2$$

may well fail in X; X can certainly contain worlds where the bottle shatters but Billy refrains. In short,  $\varphi_1$  is "more necessary" than  $\varphi_2$ , provided that we choose X right. The vagueness of X, though sometimes irksome, can also break deadlocks. Not that deadlocks must always be breakable. We ought to look with equanimity on the prospect that sometimes circumstances will fail to pinpoint a single "real cause" of an event, no matter how hard we seek.

The way in which we draw the distinction between sufficiency and necessity parallels in many regards the approach in (Pearl 2000). Like us, Pearl thinks in terms of a continuum in the relative weight given to necessity and sufficiency in different causal claims, and seeks to find the right mix of the two for various circumstances. He also (Pearl 2000: 309-11) notes the

distinction between general or type-level causal claims (“car accidents cause deaths”) and singular or token causal claims (“a car accident caused Joe’s death”) and observes that the sufficiency condition tends to be paramount in the former case while necessity predominates in the latter. As a general rule, it seems that type-level claims (emphasizing sufficiency) are of particular use to epidemiologists, economists, and policy-makers in general, while token-level claims are of interest to historians, crime investigators, and other creators of non-fictional narratives.

### Pearl’s *Causality*

In the decade since the first version of this paper, Judea Pearl’s *Causality* (2000) has emerged as a particularly significant contribution to this subject. Pearl’s arsenal for attacking the topic contains three principal weapons:

(1) Bayesian probability, in which “probabilities encode degrees of belief about events in the world and data are used to strengthen, update, or weaken those degrees of belief” (Pearl 2000: 2).

(2) Causal Bayesian networks: oversimplifying somewhat, these consist of directed acyclic graphs (DAG’s) whose nodes can be used to represent variable conditions in the real world and whose links represent probabilistic and causal dependencies between them. Upon such graphs Pearl sets up a formal operator *do* ( $\cdot$ ) in which a particular node gets fixed at a specific value, producing a modified probability distribution and an altered

network of links; this operator represents *interventions*, manipulations of the real world that might be performed by real agents, at least in our imagination.<sup>21</sup>

(3) Functional causal models, which resemble causal Bayesian networks in various ways: they also begin with networks of variables  $X$ , which are then augmented by a set of *unknown* variables  $U$  and a set of deterministic functions that connect each variable  $x_i$  with other “ancestral” variables in the network along with an accompanying unknown,  $u_i$ . An analogue of the *do* ( $\cdot$ ) operator can also be applied to these networks.

In general weapon (3) can accomplish almost everything weapon (2) can (a mathematically provable result), and so Pearl employs it most frequently in the later parts of the book.

Pearl’s formal results are impressive. They promise to bring clarity and algorithmic feasibility to many knotty problems in the analysis of complex phenomena studied by economists, epidemiologists, biologists, and others. Nonetheless, as a starting point for philosophical discussions of free will, we have certain reservations about his system. The first weapon, probability, has particularly wobbly credentials, as Kyburg notes in a laudatory review (2005), and as Pearl himself evidently concedes as *Causality* progresses. At the outset Pearl seems perfectly happy with a rather vague, epistemic view, where the probability of an event simply reflects a subjective measure of “our” confidence in its truth; his reply to Kyburg (2005) maintains the same stance, with the additional assertion that such a view can be translated with minimal philosophical effort into a more objective view of probability in terms of event frequencies. And yet in due course (Pearl 2000: 104) he admits to dissatisfaction with this

attitude, when he speaks of functional causal models, with their deterministic functions, providing an exciting alternative to “those slippery epistemic probabilities...with which we had been working so long.”<sup>22</sup>

The difficulty of assigning probabilities to individual events (or to events in the context of other bits of information) grows worse when we start trying to fashion fancier causal networks or functional causal models that represent a particular domain of study. Setting aside any questions about how to specify functions in a functional model and how to implement the *do* ( $\cdot$ ) operator, we see trouble both in the selection of nodes and in the connecting of them by links. First off, Pearl has no ambitions to build an (impractically) massive network that depicts all of creation; on the contrary, he posits explicitly that the average scientist must carve “a piece from the universe and proclaim that piece *in*... The rest of the universe is then considered *out* or *background*... This choice of *ins* and *outs* creates asymmetry in the way we look at things, and it is this asymmetry that permits us to talk about ‘outside intervention’ and hence about causality...” (Pearl 2000: 350). Fundamental though this carving process appears to be, Pearl provides few details on how it ought to work. Even if we overlook the difficulties in selecting nodes for our causal network, we immediately encounter the next problem of which nodes need to be connected by links. At first glance it may seem “obvious” that a node representing the price of rice in India and a node for the current length of the grass on the White House lawn would never require any linking; but to justify this feeling would seem to demand an appeal to causal intuitions which are after all the target of the current investigation. Pearl opens himself to charges of circularity when he writes: “thus [links] should be assumed to exist, by default, between two nodes in the diagram. They should be deleted only by well-motivated justifications,

such as the unlikely existence of a common cause for the two variables... Although we can never be cognizant of all the factors that may affect our variables, substantive knowledge sometimes permits us to state that the influence of a possible common factor is not likely to be significant” (Pearl 2000: 163) — vague criteria indeed.

Of course, vagueness comes with the territory in discussions of counterfactuals and causes, as our own account frankly admits. But we do feel that such vagueness must be addressed forthrightly, and that thinking in terms of varying families of similar possible worlds represents some non-circular progress. By contrast, when discussing Lewis’s work, Pearl first reiterates the challenges facing similarity metrics pointed out by Fine, then simply asserts that “such difficulties do not enter the structural account... [its] counterfactuals are not based on an abstract notion of similarity... instead they rest directly on the mechanisms... that produce those worlds” (Pearl 2000: 239; 2010). Pearl’s confidence that appealing to “mechanisms” solves causality’s problems seems quite unwarranted. The dilemmas of the previous paragraph, encountered whenever we choose models to represent pieces of the universe, are simply being brushed under the rug.

We believe that Pearl’s mostly unacknowledged troubles in choosing his causal models are in fact isomorphic to the issues similarity metrics face. Where Lewis dismisses a proposed alteration to the universe as too far-out to be relevant when evaluating a particular counterfactual or causal claim, Pearl simply omits the corresponding variable from his model. To take a concrete example, when justifying the assertion that “Had Nixon pressed the button, there would have been a nuclear holocaust,” a similarity-based account must rule out alterations to the world that involve electrical malfunctions in the button; analogously, if Pearl wants to produce a

determinate answer to the question “did a holocaust occur?”, given a modification of Nixon’s button-pushing behavior, he will have to exclude models where the electrical malfunction becomes an active variable in the network.<sup>23</sup>

In closing we mention one additional feature of Pearl’s exposition that may prove confusing to newcomers, concerning the order in which he develops his concepts. Where this paper tackles counterfactuals before causation, Pearl delays his main discussion of counterfactuals until the seventh chapter of *Causality*, creating an impression that causes have logical precedence. And yet in a sense Pearl does recognize counterfactuals’ priority, since he introduces the *do* ( $\cdot$ ) operator almost immediately and thereby implicitly defines a concept of what *would* happen, *if* we were to change the universe’s actual conditions.<sup>24</sup> In a similar way, Pearl postpones his main treatment of causal necessity and sufficiency until chapter 9, where they suddenly receive thorough and insightful treatment. A reader who stopped at chapter 8 might easily assume that only sufficiency matters to Pearl in causal analysis — that *X causes Y* simply means that, whenever X occurs, Y follows — perpetuating the confusions that this paper tries to unravel. Such quirks of expository strategy do not invalidate Pearl’s methods, but they may create philosophical misunderstandings.

### Determinism and Possibility (Thesis 1)

Now that we have some formal machinery in place, we can reconsider the spuriously “obvious” fear that determinism reduces our possibilities. We can see why the claim *seems* to have merit: let  $\varphi$  be the sentence “Austin holes the putt”, let X be the set of physically possible worlds that are *identical* to the actual world at some time  $t_0$  prior to the putt, and assume both

that Austin misses and that determinism holds. Then in fact  $\varphi$  does not hold for any world in  $X$  ( $\sim\Diamond_X \varphi$ ), because  $X$  contains only one world: the actual one. Of course, this method of choosing  $X$  (call it the *narrow method*) is only one among many. We should note that the moment we admit into  $X$  worlds that differ in a few imperceptibly microscopic ways from actuality at  $t_0$ , we may well find that  $\Diamond_X \varphi$ , even when determinism obtains. (This is, after all, what recent work on chaos has shown: many phenomena of interest to us can change radically if one minutely alters the initial conditions.) So the question is: when people contend that events are possible, are they really thinking in terms of the narrow method?

Notice that Austin evidently endorses the narrow method of choosing  $X$  when he states that he is “talking about conditions as they precisely were” whenever he asserts he could have holed the putt. Yet in the next sentence he seemingly rescinds this endorsement, observing that “further experiments may confirm my belief that I could have done it that time, although I did not.” What “further experiments” might indeed confirm Austin’s belief that he could have done it? Experiments on the putting green? Would his belief be shored up by his setting up and sinking near-duplicates of that short putt ten times in a row? If so, then he is not as interested as he claims he is in conditions as they precisely were. He is content to consider “Austin holes the putt” possible if, in situations very similar to the actual occasion in question, he holes the putt.<sup>25</sup>

We contend, then, that Austin equivocates when he discusses possibilities, and that in truth the narrow method of choosing  $X$  does not have the significance he imagines. From this it follows that the truth or falsity of determinism should not affect our belief that certain unrealized

events were nevertheless “possible,” in an important everyday sense of the word. We can bolster this last claim by paying a visit to a restricted domain in which we know with certainty that determinism reigns: the realm of chess-playing computer programs.

Computers are marvels of determinism. Even their so-called random number generators only execute pseudo-random functions, which produce *exactly* the same sequence of “random” digits each time the computer reboots. That means that computer programs that avail themselves of randomness at various “choice” points will nevertheless spin out exactly the same sequence of states if run over and over again from a cold start.<sup>26</sup> Suppose, for instance, you install two different chess-playing programs on your computer, and yoke them together with a little supervisory program that pits them against each other, game after game, in a potentially endless series. Will they play the same game, over and over, until you turn off the computer? Perhaps; but if either chess program consults the random number generator during its calculations (if, for instance, it periodically “flips a coin” to escape from Buridan’s ass difficulties in the course of its heuristic search), then in the following game the state of the random number generator will have changed. Accordingly different alternatives will be “chosen” and a variant game will blossom, resulting in a series in which the games, like snowflakes, are no two alike.<sup>27</sup> Nevertheless, if you turned off the computer, and then restarted it running the same program, exactly the same variegated series of games would spin out.

This gives us a toy model of a deterministic Democritean universe, in which kazillions of bits get flipped in sequence, governed by a fixed physics. Rewinding and replaying the tape of

life is really possible in such a toy world. Suppose we create such a chess universe involving two programs, A and B, and study the results of a lengthy run. We will find lots of highly reliable patterns. Suppose we find that A (almost) always beats B. That is a pattern that we will want to explain, and saying “Since the program is deterministic, A was *caused* always to beat B” would fail to address that curiosity. We will want to know what it is about the structure, methods, and dispositions of A that accounts for its superiority at chess. A has a competence or power that B lacks, and we need to isolate this interesting factor.<sup>28</sup> When we set about exploring the issue, availing ourselves of the high level perspective from which the visible “macroscopic” objects include representations of chess pieces and board positions, evaluations of possible moves, decisions about courses to pursue, and so forth, we will uncover a host of further patterns: some of them endemic to chess wherever it is played (*e.g.*, the near certainty of B’s loss in any game where B falls a rook behind) and some of them peculiar to A and B as particular chess players (*e.g.*, B’s penchant for getting its queen out early).<sup>29</sup> In short, we will find a cornucopia of *explanatory* regularities, some exceptionless (in our voluminous run) and others statistical.

These macroscopic patterns are salient moments in the unfolding of a deterministic pageant that, looked at from the perspective of micro-causation, is pretty much all the same. What from one vantage point appear to us to be two chess programs in suspenseful combat, can be seen through the “microscope” (as we watch instructions and data streaming through the CPU) to be a single deterministic automaton unfolding in the only way it can, its jumps already predictable by examining the precise state of the pseudo-random number generator. There are no “real” forks or branches in its future; all the “choices” made by A and B are already determined.

Nothing, it seems, is really *possible* in this world other than what actually happens. Suppose, for instance, that an ominous mating-net looms over B at time *t* but collapses when A runs out of time and terminates its search for the key move one pulse too soon; that mating net *was never going to happen*.<sup>30</sup> (This is something we could prove, if we doubted it, by running the same tournament another day. At exactly the same moment in the series, A would run out of time again and terminate its search at exactly the same point.)

So what are we to say? Is our toy world really a world without prevention, without offense and defense, without lost opportunities, without the thrust and parry of genuine agency, without genuine possibilities? Admittedly, our chess programs, like insects or fish, are much too simple agents to be plausible candidates for morally significant free will, but we contend that the determinism of their world does not rob them of their different powers, their different abilities to avail themselves of the opportunities presented. If we want to understand what is happening in that world, we may, indeed must, talk about how their choices cause their circumstances to change, and about what they *can* and *cannot* do.

Suppose we find two games in the series in which the first twelve moves are the same, but with A playing White in the first game and B playing White in the second. At move 13 in the first game, B “blunders” and it’s all downhill from there. At move 13 in the second game, A, in contrast, finds the saving move, castling, and goes on to win. “B *could have castled* at that point in the first game,” says an onlooker, echoing Austin. True or false? The move, castling, was just as legal the first time, so in *that* sense, it was among the “options” available to B. Suppose we

find, moreover, that castling was not only one of the represented candidate moves for B, but that B in fact undertook a perfunctory exploration of the consequences of castling, abandoned, alas, before its virtues were revealed. Could B have castled, then? Looking at *precisely* the same case, again and again, is utterly uninformative, but looking at *similar* cases is in fact diagnostic. If we find that in many similar circumstances in other games, B *does* pursue the evaluation slightly farther, discovering the virtues of such moves and making them—if we find, in the minimal case, that flipping a single bit in the random number generator would result in B’s castling—then we support (“with further experiments”) the observer’s conviction that B could have castled then. We would say, in fact, that B’s failure to castle was a fluke, bad luck with the random number generator. If, on the contrary, we find that discovering the reasons for castling requires far too much analysis for B to execute in the time available (although A, being a stronger player, is up to the task), then we will have grounds for concluding that no, B, unlike A, could not have castled. To imagine B castling would require too many alterations of reality; we would be committing an error alluded to earlier, making X too large.

In sum, using the narrow method to choose X is useless if we want to explain the patterns that are manifest in the unfolding data. It is only if we “wobble the events” (as David Lewis has said), looking *not* at “conditions as they precisely were” but at nearby neighboring worlds, that we achieve any understanding at all.<sup>31</sup> Once we expand X a little, we discover that B has additional options, in a sense both informative and morally relevant (when we address worlds beyond the chessboard). The burden rests with incompatibilists to explain why “real” possibility demands a narrow choice of X—or why we should be interested in such a concept of possibility,

regardless of its “reality.”

As we have seen, possibilities of the broader, more interesting variety can exist quite comfortably in deterministic worlds. Indeed, introducing indeterminism adds nothing in the way of worthwhile possibilities, opportunities, or competences to a universe. If in our sample deterministic world program A always beats program B, then replacing the pseudo-random number generator with a genuinely indeterministic device will not help B at all: A will *still* win every time. Though pseudo-random generators may not produce genuinely random output, they come so close that no ordinary mortal can tell the difference. A superior algorithm like A’s will hardly stumble when faced with so inconsequential a change. And analogous conclusions could well apply in meatier universes like ours. To put it graphically, the universe could be deterministic on even days of the month and indeterministic on odd days, and we’d never notice a difference in human opportunities or powers; there would be just as many triumphs—and just as many lamentable lapses—on October 4 as on October 3 or October 5. (If your horoscope advised you to postpone any morally serious decision to an odd numbered day, you would have no more reason to follow this advice than if it told you to wait for a waning moon.)

### Determinism and Causation (Thesis 2)

The hunch that determinism would eliminate some worthwhile type of causation from the universe has even less merit than the claim that it eliminates possibilities. We suspect this fear stems from the conflation of causal necessity with causal sufficiency—as we have seen, our language makes this confusion all too easy. Determinism is essentially a doctrine concerned

with sufficiency: if  $\sigma_0$  is a (mind-bogglingly complex) sentence that specifies in complete detail the state of the universe at  $t_0$  and  $\sigma_1$  similarly specifies the universe at a later time  $t_1$ , then determinism dictates that  $\sigma_0$  is sufficient for  $\sigma_1$  in all physically possible worlds. But determinism tells us nothing about what earlier conditions are *necessary* to produce  $\sigma_1$ , or any other sentence  $\psi$  for that matter. Hence, since causation generally presupposes necessity, the truth of determinism would have little bearing on the validity of our causal judgments.<sup>32</sup>

For example: according to determinism, the precise condition of the universe one second after the big bang (call the corresponding sentence  $\sigma_0$ ) causally sufficed to produce the assassination of John F. Kennedy in 1963 (sentence  $\psi$ ). Yet there is no reason at all to claim that  $\sigma_0$  caused  $\psi$ . Though sufficient,  $\sigma_0$  is hardly necessary. For all we know, Kennedy might well have been assassinated anyway, even if some different conditions had obtained back during the universe's birth.<sup>33</sup> More plausible causes of the event would include: "A bullet followed a course directed at Kennedy's body"; "Lee Harvey Oswald pulled the trigger on his gun"; perhaps "Kennedy was born"; conceivably "Oswald was born."<sup>34</sup> But conspicuously absent from this list are microscopically detailed descriptions of the universe billions of years prior to the incident. Incompatibilists who assert that under determinism  $\sigma_0$  "causes" or "explains"  $\psi$  miss the main point of causal inquiry.

In fact, determinism is perfectly compatible with the notion that some events have no cause at all. Consider the sentence "The devaluation of the rupiah caused the Dow Jones average to fall." We rightly treat such a declaration with suspicion; are we really so sure that among

nearby universes the Dow Jones fell *only* in those where the rupiah fell first? Do we even imagine that every universe where the rupiah fell experienced a stock market sell-off? Might there not have been a confluence of dozens of factors which jointly sufficed to send the market tumbling but none of which by itself was essential? On some days, perhaps, Wall Street's behavior has a ready explanation; yet at least as often we suspect that no particular cause is at work. And surely our opinions about the market's activities would remain the same, whether we happened to adopt Newton's physics or Schrödinger's. In general, instances of apparent randomness (such as stock market fluctuations, coin flips and dice throws) are described in our account as having no cause at all; even if determinism is true, they are determined by chaotic processes that cannot be reduced to tractable sets of necessary and sufficient conditions.

Of course, one might wonder why it is that causal necessity matters to us as much as it does. Let us return for a moment to chess programs A and B. Suppose our attention is drawn to a rare game in which B wins, and we want to know "the cause" of this striking victory. The trivial claim that B's win was "caused" by the initial state of the computer is totally uninformative. Of course the total state of the toy universe at prior moments was *sufficient* for the occurrence of the win; we want to know which features were *necessary*, and thereby understand what such rare events have in common. We want to discover those features, the absence of which would most directly be followed by B's loss, the default outcome. Perhaps we will find a heretofore unsuspected flaw in A's control structure, a bug that has only just now surfaced. Or perhaps the victory is a huge coincidence of conditions that require no repair, since the probability of their recurrence is effectively zero. Or we might find an idiosyncratic island of

brilliance in B's competence, which once diagnosed would enable us to say just what circumstances in the future might permit another such victory for B.

In closing, let us return to the human desire pinpointed by Kane that motivates so much of this debate: the desire to be able to take full credit as the creators and causes of change in the world. Consider for instance the wish that we (Taylor and Dennett) have to be acknowledged as the authors of this paper. Suppose that determinism turns out to be true. Would that in any way undercut our claim that our activity nevertheless played an essential role in this paper's creation? Not in the least, even after we factor in the earlier deeds of our parents and teachers. Without our efforts, it is safe to say that no paper exactly like this (or even closely similar) would have been produced.<sup>35</sup> Hence we are entitled to claim some "originative value" for our unique accomplishment. The thirst for originality and causal relevance is not to be quenched by abstruse quantum events: all that we require is the knowledge that without our presence, the universe would have turned out significantly different.

## Appendix: Van Inwagen's Consequence Argument

Peter Van Inwagen (1975) hopes to bolster the incompatibilist sense of lost causal powers with the following basic argument:

1. Let  $\varphi$  be some event that actually occurs in agent A's life (missing a putt, say). Also let  $\sigma_0$  be a comprehensive description of the universe's state at some time in the remote past, and let  $\lambda$  be a statement of the laws of nature.
2. Then, assuming determinism,  $\lambda \wedge \sigma_0 \Rightarrow \varphi$  applies in every possible world.  
Equivalently,  $\sim\varphi \Rightarrow \sim(\lambda \wedge \sigma_0)$ .
3. If A has the power to cause  $\alpha$  and  $\alpha \Rightarrow \beta$  obtains in every possible world, then A has the power to cause  $\beta$ .
4. So if A has the power to cause  $\sim\varphi$ , then A has the power to cause the falsity of either  $\lambda$  or  $\sigma_0$ , which is absurd.
5. Therefore A lacks the power to cause  $\sim\varphi$ .

This argument illustrates nicely the confusion that causal necessity and sufficiency engender. As we have argued, counterfactual necessity is the single most crucial condition for causation, and accordingly we would recommend that Van Inwagen's "power to cause  $\alpha$ " be rendered as follows:

A has the *power to cause*  $\alpha$  iff for some sentence  $\gamma$  describing an action of A and a world  $f$  close to actuality,  $\gamma \wedge \alpha$  holds in  $f$  and  $\alpha \Rightarrow \gamma$  in every world similar to  $f$ .

In other words, within some cluster of nearby worlds, there is a possible action of A (called  $\gamma$ )

that is a necessary condition for  $\alpha$  to occur. But under this definition, line 3 above has no warrant whatever. Line 3 hypothesizes that  $\alpha \Rightarrow \gamma$  in a cluster of nearby worlds, and that  $\alpha \Rightarrow \beta$  in every world; if we could deduce that  $\beta \Rightarrow \gamma$  in this cluster, we would be home free. But of course in Logic 101 we learn that  $\alpha \Rightarrow \gamma$  and  $\alpha \Rightarrow \beta$  do not entail  $\beta \Rightarrow \gamma$ , and so line 3 fails, and Van Inwagen's argument with it.<sup>36</sup>

## Bibliography

- Austin, John. 1961. "Ifs and Cans." In *Philosophical Papers*, ed. J. O. Urmson and G. Warnock. Oxford: Clarendon Press.
- Dennett, Daniel. 1978. *Brainstorms*. Cambridge, Mass.: MIT Press.
- Dennett, Daniel. 1984. *Elbow Room: the Varieties of Free Will Worth Wanting*. Cambridge, Mass.: MIT Press.
- Dennett, Daniel. 1988. "Coming to Terms with the Determined" (review of Honderich 1988). *The Times Literary Supplement*, November 4-10, 1988: 1219-20.
- Dennett, Daniel. 1991. "Real Patterns." *Journal of Philosophy* 88: 27-51.
- Dennett, Daniel. 1995. *Darwin's Dangerous Idea*. New York: Touchstone.
- Dennett, Daniel. 2005. "Natural Freedom." *Metaphilosophy* 36: 449-459.
- Fine, K. 1975. Review of Lewis, *Counterfactuals*, *Mind* 84, 451-8.
- Fischer, J. 2005. "Dennett on the Basic Argument." *Metaphilosophy* 36: 427-435.
- Gasking, Douglas. 1955. "Causation and Recipes." *Mind* 64: 479-487.
- Ginet, Carl. 1990. *On Action*. Cambridge: Cambridge University Press.
- Hall, Ned. 2000. "Causation and the Price of Transitivity." *Journal of Philosophy* 97: 198-222.
- Hart, H. L. A. and A. M. Honoré. 1959. *Causation in the Law*. Oxford: Clarendon Press.
- Honderich, Ted. 1988. *A Theory of Determinism: The Mind, Neuroscience, and Life-Hopes*. Oxford: Clarendon Press.
- Kane, Robert. 1998. *The Significance of Free Will*. Oxford: Oxford University Press.

- Kyburg, H. 2005. Review of Pearl (2000), *Artificial Intelligence*, **169**, 174-79.
- Lewis, David. 1973. *Counterfactuals*. Cambridge, Mass.: Harvard University Press.
- Lewis, David. 2000. "Causation as Influence." *Journal of Philosophy* 97: 182-197.
- McLaughlin, J. A. 1925. "Proximate Cause." *Harvard Law Review* 39: 149-155.
- Nozick, Robert. 1981. *Philosophical Explanations*. Cambridge: Harvard University Press.
- Paul, L. A. 2000. "Aspect Causation." *Journal of Philosophy* 97: 235-256.
- Pearl, J. 2000. *Causality: Models, Reasoning, and Inference*, Cambridge: Cambridge Univ. Press.
- Pearl, J. 2005. "Response to review by Kyburg," *Artificial Intelligence*, **169**, 180.
- Pearl, J. 2010. "If Oswald Hadn't Used Counterfactuals, My Robot Would Have,"  
 UCLA Cognitive Systems Laboratory, Technical Report (R-360), February 2010.  
 Available at [http://bayes.cs.ucla.edu/csl\\_papers.html](http://bayes.cs.ucla.edu/csl_papers.html).
- Quine, W.V.O. 1969. "Propositional Objects." In *Ontological Relativity*. New York: Columbia University Press.
- Quine, W.V.O. 1980. "Reference and Modality." In *From a Logical Point of View*. Second ed., revised. Cambridge, Mass.: Harvard University Press.
- Schaffer, Jonathan. 2000. "Trumping Preemption." *Journal of Philosophy* 97: 165-181.
- Taylor, Christopher and Daniel Dennett. 2002. "Who's Afraid of Determinism: Rethinking Causes and Possibilities." In Kane, R. Ed. *Oxford Handbook of Free Will*. Oxford: Oxford University Press.
- Tooley, Michael. 1987. *Causation: A Realist Approach*. Oxford: Oxford University Press.
- Van Inwagen, Peter. 1975. "The Incompatibility of Free Will and Determinism." *Philosophical*

*Studies 27: 185-99.*

Woodward, James. 2003. *Making Things Happen: A Theory of Causal Explanation*, Oxford:  
Oxford University Press.

## Notes

---

1. This paper incorporates almost all of “Who’s Afraid of Determinism: Rethinking Causes and Possibilities” (Taylor and Dennett, 2002) but includes significant additions and adjustments; hence it supersedes that paper.

2. Nozick 1981: 313. “We want it to be true that in that very same situation we could have done (significantly) otherwise, so that our actions will have originative value.”

3. Austin 1961: 166.

4. Kane 1998: 100.

5. Kane 1998: 4.

6. Ginet 1990: 90.

7. See Quine 1980 for a discussion of these pitfalls.

8. Quine 1969: 147-55.

9. The average educated person’s casual working assumptions about the cosmos still resemble the Democritean account, and philosophers traditionally rely on nothing more sophisticated when exploring the implications of determinism and indeterminism, causation and possibility.

Our suggestion that possible worlds simply *are* functions of the appropriate form may seem disturbingly reductive, particularly when one contemplates the particular function(s) that correspond to the actual world; accordingly David Lewis takes pains to distinguish possible worlds from their mathematical “handles.” However one wishes to address these ontological scruples, nothing in the following discussion hinges on them.

---

10. Since we are restricting ourselves to the scientifically old-fashioned Democritean worlds, we would have trouble specifying the contents of  $\Phi$  precisely—and besides, of course, we do not yet *know* all the laws of nature!—but we can pretend that we know, and hence we can pretend that in most cases one can judge whether or not a particular world  $f$  accords with natural law.

John Horton Conway's Game of Life can be viewed as a particularly simple pseudo-Democritean universe, eliminating one spatial dimension and quantizing time. (See Dennett 1991: 27-51 or Dennett 1995, for an introduction to Life.) The set of all possible sequences of bitmaps is then  $\Omega$ , and the single (deterministic) rule of Life "physics" applied to every "initial" state gives us the subset  $\Phi$  of  $\Omega$ . Every variation on Conway's "physics" generates a different subset  $\Phi$ .

11. Of course, these predicates unleash a horde of problems concerning vagueness, subjectivity, and (in such cases as "believes that snow is white") intentionality, but difficulties along these lines do not imperil the basic approach.

12. Lewis 1973, *passim*.

13. Of course, X can vary in practice, as already observed, so that we must treat assumptions of transitivity with care. But notice that the Guidelines (G) tend to yield the same set X for the two sentences  $\varphi \Rightarrow \psi$  and  $\sim\psi \Rightarrow \sim\varphi$ ; hence the rule of contraposition is in general reliable.

14. We are indebted to Gary Drescher (personal communication) for clarifying our thinking on this topic.

15. See, e.g., Tooley 1987.

16. A vast amount of ink has been spilled arguing that the direction of causation is either independent of or logically prior to the direction of time, and to address the matter here would

---

require too lengthy a digression. So we merely note the issue, and tentatively take the direction of time as a given (originating ultimately in the Second Law of Thermodynamics) from which the direction of causation derives.

Gasking (1955) raises a number of interesting cases in which cause and effect appear to be simultaneous: for instance, if a piece of iron attains a temperature of (say) 1000°C and thereupon starts to glow, we still distinguish the former as cause and the latter as effect. But this apparent exception to the rule has a ready explanation that Gasking himself hints at: when a speaker refers to the iron “reaching 1000°,” she is envisioning this event as the endpoint in a lengthy heating process. The heating process *does* precede the glowing, and so the latter is considered an effect.

Another category of “exceptions” includes diseases and their symptoms (say, a cold and sneezing), which might sometimes arise simultaneously. Yet often enough diseases *do* precede their symptoms, while symptoms (by definition) never appear before their diseases. Accordingly we grant diseases the status of “cause.”

17. Notice that we do not in the above clauses make any provision to ensure the transitivity of causation. Lewis (2000: 191-5), among others, feels it important to guarantee transitivity by making “causation” the ancestral of “causal dependence.” But Lewis himself provides many examples of transitivity’s counter-intuitive consequences. For instance, suppose that agent A wants to travel to New York. Agent B, hoping to thwart A, lets the air out of the tires on A’s car. In consequence, A takes the train instead and reaches New York only slightly behind schedule. If causation is transitive, then B has “caused” A’s successful arrival, despite the fact that the two

---

sentences “B lets the air out of A’s tires” and “A arrives in New York” satisfy none of our more crucial conditions. Lewis finds the awkward implications of transitivity acceptable; we remain unpersuaded.

Hall (2000) goes to even greater lengths defending transitivity. His account would seemingly imply that a pebble on the train tracks south of Paris that minutely alters the course of the Orient Express is a "cause" of the train's arrival in Istanbul several days later. Paul's "Aspect Causation" (2000) suggests a possible diagnosis for Hall's willingness to countenance such bizarre conclusions, as stemming from an overeager acceptance of the premiss that causation is a relation between "events" (however this problematic term may be defined). At any rate, notice that on our account one can consistently consider false the sentence "Pebble *p*'s lying on the tracks south of Paris caused the train's arrival in Istanbul," while accepting "Pebble *p*'s lying on the tracks south of Paris caused the train's arrival in Istanbul via a minutely altered course in France."

18.Lewis 2000.

19.Obviously, a sentence like “Drugs or aliens caused Elvis’s premature demise” abbreviates the cumbersome “Drugs caused Elvis’s premature demise or aliens caused Elvis’s premature demise” — a disjunction of two separate causes, not a single disjunctive cause.

20.Invoking causal sufficiency in this way solves, to our satisfaction, all of the analogous problem cases raised by Schaffer (2000). Note that Schaffer rather misleadingly suggests that “counterfactual accounts of causation” must always be formulated solely in terms of necessity (2000: 176). We, on the contrary, consider our account essentially “counterfactual” even though

---

it allows for sufficiency along with necessity.

Lewis's formulation (2000) of "Causation as Influence" can be viewed as an indirect way of introducing sufficiency into an originally necessity-centered account. For present purposes we consider our approach more illuminating, but both strategies point in the same general direction.

21. Woodward (2003) also builds his "manipulationist" notion of causal explanation on interventions of this sort; interestingly, he is avowedly unperturbed by the threatened circularity of basing an account of causality on such an undeniably causal notion as intervention. (2003: 106)

22. Notice that when he launches his first chapter with a discussion of probability theory Pearl is tacitly using a notion of what events and combinations of events — even combinations that have never occurred in reality — are *possible*. His starting point therefore resembles ours.

23. Woodward (2003: 110) also points out that Pearl's formalism is ideal for his purposes but that it is dependent on choosing "the *correct* causal graph for the system in which the intervention occurs." See Pearl 2010 for a more explicit discussion of counterfactuals, prompted in part by our discussion in an earlier draft of this essay.

24. Pearl has acknowledged this point (personal communication).

25. When Austin speaks of further experiments, could he be referring to experiments in the high-tech labs of physicists and microbiologists, experiments that would convince him that his brain amplifies indeterministic quantum events? Given the extreme impracticality of such

---

experiments, and Austin's overall skepticism about the relevance of science in these contexts ("[A modern belief in science] is not in line with the traditional beliefs enshrined in the word *can*," Austin 1961: 166), this interpretation seems unlikely. But this is precisely the direction in which Kane and some other incompatibilists have headed. See also Dennett 1984: 133-37.

26. We are restricting our attention to programs that do not require or accept input from the external world, which could, of course, be random in any of several senses. The easiest way to ensure that there is variation in subsequent runs of a program is to have it call for inputs of these sorts: the time taken from the computer's clock, the presence or absence of a pulse from a Geiger counter, the last digit in the latest Dow Jones Industrial Average as taken off the Internet, etc.

27. All this is independent of whether or not either chess program can "learn from its experience," which is another way their internal state could change over time to guarantee that no two games were the same.

28. Another case in which we could know *all* the deterministic micro-details but be baffled about how to explain the causal regularities is Dennett's example of the two black boxes (1995: 412-22).

29. Dennett 1978: 107.

30. Cf. the comet plunging towards earth that gets intercepted at the last minute by the other comet, unnoticed till then, that had been on its collision trajectory since its birth millions of years ago (Dennett 1984: 124).

---

31.If we exclude such variation, then trivially, castling in the second game was not “open to B,” to use Ginet’s terminology. Recall that Ginet requires that “nothing that exists up to that moment stands in the way of my doing next any one of the alternatives.” The narrow method has the effect of treating the precise state of B’s contemplation of the option of castling as something *external*, as something that can itself “stand in the way” at the moment of choosing, guaranteeing that *nothing about B* could *explain* B’s choice, whatever it is. As Dennett notes, “If you make yourself really small, you can externalize virtually everything” (1984: 143).

32.See the Appendix for an additional example of the conflation of necessity and sufficiency (in Van Inwagen’s Consequence Argument).

33.Imagine that we take a snapshot of the universe at the moment of Kennedy’s assassination, then alter the picture in some trivial way (by moving Kennedy 1 mm to the left, say). Then, following the (deterministic) laws of physics in reverse, we can generate a movie running all the way back to the Big Bang, obtaining a world in which  $\sigma_0$  subtly fails.

34.Of course, the last two options fail the sufficiency test so badly that we prefer not to countenance them as causes. As explained earlier, sufficiency does have *some* relevance in assigning causes, just not the overwhelming importance that incompatibilists imply.

35.Similarly, Deep Blue, in spite of its being a deterministic automaton, authored the games of chess that vanquished Kasparov. No one *else* was their author; Murray Campbell and the IBM team that created Deep Blue can’t take credit for those games; *they* didn’t see the moves. It was the vast exploratory activity of Deep Blue itself that was the originating cause of those

---

magnificent games.

36. Fischer (2005) attempts to sidestep the argument in this Appendix by formulating a definition of causation solely in terms of sufficiency, thus revitalizing premise (3). He admits that this definition does not “capture some ordinary, commonsense idea” and that on this reading “I have it in my power to render it true that the sun continues to shine.” We concede to Fischer that this curious conception of causation would keep the Consequence Argument alive, but only barely. Fischer’s parting shot in his discussion is to formulate a new version of what he calls the Basic Argument based on “the extremely plausible and intuitively attractive Principle of the Fixity of the Past and Laws: an agent has it within his power to do A only if his doing A can be an extension of the actual past, holding the natural laws fixed.” (Fischer, 2005: 32) This principle is simply an affirmation of the narrow method of interpreting possibility, repeating Austin’s mistake. Once this mistake is recognized, much of the literature of incompatibilism loses its foundation. As Dennett (2005) observes, “We wouldn’t want to say farewell to something as much fun as the Basic Argument in an appendix would we? Well, yes.”