

# Why Take Both Boxes?

Jack Spencer and Ian Wells

Forthcoming in *Philosophy and Phenomenological Research*

## Abstract

The causal dominance principle that is the crucial premise of the standard argument for two-boxing in Newcomb's problem is false. We present some counterexamples to the principle. We then offer a metaethical explanation for why the counterexamples arise. Our explanation reveals a new and superior argument for two-boxing, one that eschews the causal dominance principle in favor of a principle linking rational choice to guidance and actual value maximization.

In the classic Newcomb problem, there is an agent, a transparent box, an opaque box, and a predictor, known by the agent to be uncannily good:<sup>1</sup>

*Classic Newcomb.* The agent has two options: she can take either only the opaque box or both boxes. The transparent box contains \$1,000. The opaque box contains either \$1,000,000 or nothing, depending on a prediction made yesterday by the predictor. If the predictor predicted that the agent would take both boxes, the opaque box contains nothing. If the predictor predicted that the agent would take only the opaque box, the opaque box contains \$1,000,000. The agent knows all of this.

*One-boxing* is the claim that the agent facing *Classic Newcomb* is rationally required to take only the opaque box. *Two-boxing* is the claim that the agent is rationally required to take both boxes. In this paper, we fortify the case for two-boxing.

---

<sup>1</sup>First discussed by Nozick (1969).

Fortification is needed because the standard argument for two-boxing—a causal dominance argument—fails. The crucial premise of the standard argument is a causal dominance principle that prohibits choosing causally dominated options. The argument fails because the principle is false. As the examples that we present below establish, it is sometimes rationally permissible to choose a causally dominated option.

After presenting counterexamples to the causal dominance principle, we offer a metaethical explanation for why the counterexamples arise. The explanation reveals a new and superior argument for two-boxing, one that eschews the causal dominance principle in favor of a principle linking rational choice to guidance and actual value maximization.

## 1

The actual value of an option (sometimes called the value, utility, or actual utility of the option) is the value of the outcome that would result if the agent were to choose the option. For example, imagine that there are several boxes, each containing a sum of money. The agent is to choose one of the boxes. The outcome that would result if the agent were to choose a particular box is the agent receiving the sum of money contained therein. If money is all that matters, and more money is linearly better, then the actual value of choosing the box can be identified with the number of dollars contained therein.

The main task of decision theory is to identify the options, among those available to the agent, that the agent is rationally permitted to choose. The task is easy when the agent knows the actual values of her options, for then an option is rationally permissible to choose if and only if the option maximizes actual value.<sup>2</sup> The task is more interesting and more difficult when the agent does not know the actual values of her options.

---

<sup>2</sup>Cf. Ramsey (1990 [1926], p. 70): “Let us begin by supposing that our subject has no doubts about anything, but certain opinions about all propositions. Then we can say that he will always choose the course of action which will lead in his opinion to the greatest sum of good.”

Following Savage and Jeffrey,<sup>3</sup> many decision theorists believe that an option is rationally permissible for an agent to choose if and only if the option maximizes *expected* value, where the expected value of an option is the agent's expectation of the actual value of the option.<sup>4</sup> There are many well-defined expected value quantities and there is considerable disagreement about which of them, if any, is tied to rational choice. We will focus on two: causal expected value (hereafter *c-expected value*) and evidential expected value (hereafter *e-expected value*). Both can be defined in a common conceptual framework, which centers on the concept of a decision problem.

A decision problem is characterized by a set of options, a set of possible outcomes, and a decision-making agent. The options  $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$  are the things the agent is choosing between. We take options to be propositions that the agent can make true by choosing.<sup>5</sup> We assume that options are finite in number, mutually exclusive, and jointly exhaustive. The possible outcomes  $\mathcal{O} = \{O_1, O_2, \dots, O_m\}$  are the objects of non-instrumental desire. We take outcomes to be propositions that fully specify the desirable and undesirable consequences that might result from the choice. Like options, outcomes are assumed to be finite in number, mutually exclusive, and jointly exhaustive. We associate the agent both with a credence function  $C$  and with a valuation function  $V$ . The credence function, a probabilistically coherent function that maps propositions to the unit interval, represents the agent's beliefs. The agent's

---

<sup>3</sup>Savage (1954), Jeffrey (1965).

<sup>4</sup>See note 9.

<sup>5</sup>We follow Jeffrey (1965) and take options to be among the propositions to which the agent assigns credences. Some philosophers are skeptical of assigning credences to options, since they think that deliberation crowds out prediction. See, among others, Spohn (1977) and Levi (1997). We are convinced by the arguments in Joyce (2002), Rabinowicz (2002) and Hájek (2016) that deliberation does not crowd out prediction.

For reasons discussed in, among other places, Hedden (2015) and Pollock (2002), an agent must be certain that she will choose an option if she decides to do so. We therefore identify options and decisions. Each option, besides the null decision, is a proposition of the form:  $S$  decides to  $\phi$ .

credence in  $P$ ,  $C(P)$ , is the degree to which the agent believes that  $P$ . The valuation function, which maps outcomes to real numbers, represents the agent's desires.<sup>6</sup> The value of outcome  $O$ ,  $V(O)$ , is the degree to which the agent finds  $O$  non-instrumentally desirable.

Given this conception of a decision problem, the e-expected value of an option  $A \in \mathcal{A}$  can be written as a credence-weighted sum, wherein the relevant credences are conditional on the option in question:

$$eev(A) = \sum_O C(O | A)V(O).$$

The *rule of e-expected value* states that agents are always rationally required to choose so as to maximize e-expected value.

Let ' $\square \rightarrow$ ' be the non-backtracking counterfactual conditional. If we assume that, for each option, there is a fact of the matter about which outcome would result if the agent were to choose the option,<sup>7</sup> then the c-expected value of an option can be written as a credence-weighted sum, wherein the relevant credences are unconditional credences in counterfactual conditionals:<sup>8</sup>

$$cev(A) = \sum_O C(A \square \rightarrow O)V(O).$$

The *rule of c-expected value* states that agents are always rationally required to choose so as to maximize c-expected value.<sup>9</sup>

---

<sup>6</sup>The valuation function is unique up to positive affine transformation.

<sup>7</sup>This assumption is tantamount to counterfactual excluded middle. For a discussion of causal decision theory without counterfactual excluded middle, see, for example, Lewis (1981), Sobel (1994, p. 141-73) and Joyce (1999).

<sup>8</sup>See, for example, Gibbard and Harper (1978) and Stalnaker (1981). Some philosophers are suspicious of assigning probabilities to counterfactuals. (Thanks to an anonymous referee here.) For the purposes of this paper, alternative characterizations of c-expected work equally. We could define c-expected value using expected chance, as Skyrms (1984) does, or using dependency hypotheses, as Lewis (1981) does, or using the epistemic (as opposed to stochastic) probabilities of Kyburg (1980).

<sup>9</sup>So long as every option has an actual value (see note 15), both e-expected value and c-expected value can be defined as expectations of actual value. An  $av(A)$ -level proposition has the form  $[av(A) = v]$ , and is true just if

It is generally agreed that, in *Classic Newcomb*, the rule of e-expected value entails one-boxing,<sup>10</sup> and the rule of c-expected value entails two-boxing.<sup>11</sup> But appealing to the rule of e-expected value or the rule of c-expected value cannot settle the debate between one-boxers and two-boxers, for, as you might suspect, one-boxers typically reject the rule of c-expected value, and two-boxers typically reject the rule of e-expected value.<sup>12</sup> To move the debate forward, we need an independent argument, one that nowhere appeals to an expected value quantity.

Many two-boxers believe that they have an independent argument: namely, a causal

---

$v$  is the actual value of  $A$ . The e-expected value of an option is the agent's conditional expectation of the actual value of the option and can be written  $\sum_v vC([av(A) = v] | A)$ . The c-expected value of an option is the agent's unconditional expectation of the actual value of the option and can be written  $\sum_v vC([av(A) = v])$ .

<sup>10</sup>Let  $A_{1B}$  be the option of taking only the opaque box. Let  $A_{2B}$  be the option of taking both boxes. Conditional on  $A_{1B}$ , the agent is highly confident that the opaque box contains \$1,000,000, so, equating dollars and units of value,  $eev(A_{1B}) \approx 1,000,000$ . Conditional on  $A_{2B}$ , the agent is highly confident that the opaque box contains \$0, so  $eev(A_{2B}) \approx 1,000$ . Since  $eev(A_{1B}) > eev(A_{2B})$ , the rule of e-expected value entails one-boxing.

<sup>11</sup>Let  $O_0$ ,  $O_T$ ,  $O_M$ , and  $O_{M+T}$  be the outcomes of receiving \$0, \$1,000, \$1,000,000, and \$1,001,000, respectively. The agent knows that either  $O_0$  or  $O_M$  will result if she takes only the opaque box and that either  $O_T$  or  $O_{M+T}$  will result if she takes both boxes. Moreover, she knows that her choice has no causal bearing on what sum of money is contained in the opaque box, so  $C([A_{1B} \square \rightarrow O_0]) = C([A_{2B} \square \rightarrow O_T])$  and  $C([A_{1B} \square \rightarrow O_M]) = C([A_{2B} \square \rightarrow O_{M+T}])$ . Hence, no matter what credence function she has, the c-expected value of taking both boxes is exactly 1,000 greater than the c-expected value of taking only the opaque box:

$$\begin{aligned}
cev(A_{2B}) &= \sum_O C([A_{2B} \square \rightarrow O])V(O) \\
&= C([A_{2B} \square \rightarrow O_T])V(O_T) + C([A_{2B} \square \rightarrow O_{M+T}])V(O_{M+T}) \\
&= (1 - C([A_{2B} \square \rightarrow O_{M+T}]))(1,000) + C([A_{2B} \square \rightarrow O_{M+T}])V(O_{M+T}) \\
&= 1,000 + C([A_{2B} \square \rightarrow O_{M+T}])V(O_{M+T}) \\
&= 1,000 + C([A_{1B} \square \rightarrow O_0])V(O_0) + C([A_{1B} \square \rightarrow O_M])V(O_M) \\
&= 1,000 + \sum_O C([A_{1B} \square \rightarrow O])V(O) \\
&= 1,000 + cev(A_{1B}).
\end{aligned}$$

<sup>12</sup>Heterodoxically, Eells (1982) argues that the rule of e-expected value entails two-boxing, and Spohn (2012) argues that the rule of c-expected value entails one-boxing.

dominance argument.<sup>13</sup>

### 3

A natural way to argue for two-boxing is by disjunctive syllogism. We can imagine running through the argument from the agent's point of view:

The opaque box contains either \$1,000,000 or nothing. If it contains \$1,000,000, then both boxes together contain \$1,001,000, and hence I would make more money if I took both boxes. If it contains nothing, then both boxes together contain \$1,000, and hence I would make more money if I took both boxes. Either way, I would make more money if I took both boxes. So I should take both boxes.

This argument, although unregimented, seems compelling and nowhere invokes an expected value quantity.

A preliminary attempt to regiment the argument appeals to states and (strict) dominance. A set of propositions  $\mathcal{S} = \{S_1, S_2, \dots, S_n\}$  is a set of *states* if its members are mutually exclusive and jointly exhaustive, and each  $S \in \mathcal{S}$  is compossible with each  $A \in \mathcal{A}$ . Let  $AS$  be the conjunction of option  $A$  and state  $S$ . If the options and states are sufficiently fine-grained (and let us choose them so that they are), then every  $AS$  necessitates a unique outcome. If  $AS$  necessitates  $O$ , we set  $V(AS)$  equal to  $V(O)$ . Option  $A_i$  *dominates* option  $A_j$ , then, if and only if there is a set of states  $\mathcal{S}$  such that, for every  $S \in \mathcal{S}$ ,  $V(A_i S)$  exceeds  $V(A_j S)$ .

One might allege the following connection between dominance and rational choice:

**Dominance:** If option  $A_i$  dominates option  $A_j$ , then it is not rationally permissible for the agent to choose  $A_j$ .

But it is common ground between one-boxers and two-boxers that Dominance is false. It is

---

<sup>13</sup>See, for example, Joyce (1999, p. 152-54), Lewis (1981, p. 309-12), Skyrms (1984, p. 67) and Sobel (1994).

sometimes rationally permissible to choose dominated options, as cases like the following make vivid:<sup>14</sup>

*The Extortionist.* A moviegoer parks her car in the lot. An extortionist, who the moviegoer has excellent reason to trust, says to the moviegoer, “If you pay me \$10, I’ll ensure that your windshield is unbroken when you return. But I’ll smash your windshield if you don’t pay me.”

Let the set of states be  $\{S_B, S_{\bar{B}}\}$ , where  $S_B$  is the proposition that the windshield is broken when the moviegoer returns and  $S_{\bar{B}}$  is the proposition that the windshield is not broken when the moviegoer returns. Let  $A_P$  be the option of paying the extortionist and let  $A_{\bar{P}}$  be the option of not paying.  $V(A_{\bar{P}}S_{\bar{B}}) > V(A_P S_{\bar{B}})$ , since it would be better by the moviegoer’s lights not to pay the extortionist and return to an unbroken windshield than to pay the extortionist and return to an unbroken windshield.  $V(A_{\bar{P}}S_B) > V(A_P S_B)$ , since it would be better by the moviegoer’s lights not to pay the extortionist and return to a broken windshield than to pay the extortionist and return to a broken windshield. Dominance therefore entails that the agent is rationally required to not pay the extortionist—which is absurd. The moviegoer is rationally required to pay the extortionist: paying \$10 is much better than paying \$1,000 for a new windshield.

A bit of reflection reveals why Dominance fails. Reasoning by Dominance is supposed to put the agent in a position to conclude a fact about the ordinal ranking of options vis-à-vis actual value. The disjunctive syllogism mentioned at the outset of this section, for example, is supposed to put the agent in a position to conclude that the actual value of taking both boxes exceeds the actual value of taking only the opaque box. But actual value does not respect dominance: the fact that  $A_i$  dominates  $A_j$  does not entail that the actual value of  $A_i$  exceeds the actual value of  $A_j$ . Since we are assuming that, for each option, there is a fact of the matter about which outcome would result if the agent were to choose the option, we can characterize actual value as a sum.

---

<sup>14</sup>An adaptation of an example from Joyce (1999, p. 114-19). Jeffrey (1965, p. 9-10) uses the example of nuclear disarmament.

Where  $T$  is an indicator function that assigns one to truths and zero to falsehoods, the actual value of an option,  $av(A)$ , can be written:

$$av(A) = \sum_O T(A \Box \rightarrow O)V(O).$$

Given our assumptions, there is exactly one  $O \in \mathcal{O}$  for which  $[A \Box \rightarrow O]$  is true. If  $[A \Box \rightarrow O]$  is true, then  $O$  is the outcome that would result if the agent were to choose  $A$ , and hence  $av(A)$  equals  $V(O)$ .<sup>15</sup> Let  $S_{@}$  be the state that actually obtains. The fact that  $A_i$  dominates  $A_j$  entails that  $V(A_i S_{@})$  exceeds  $V(A_j S_{@})$ . It might be tempting to identify the actual values of  $A_i$  and  $A_j$  with  $V(A_i S_{@})$  and  $V(A_j S_{@})$ , respectively. But that temptation must be resisted. The actual value of  $A$  is equal to  $V(AS)$  only if  $S$  would have obtained had the agent chosen  $A$ . If the agent chooses  $A_i$ , then the actual value of  $A_i$  is equal to  $V(A_i S_{@})$ . But the actual value of an unchosen option  $A_j$  need not be equal to  $V(A_j S_{@})$ .

To illustrate, return to *The Extortionist*, and suppose that the extortionist is trustworthy. The moviegoer irrationally chooses to not pay the extortionist and returns to a broken windshield.  $S_B$  is true, and the actual value of not paying is equal to  $V(A_{\bar{P}} S_B)$ .  $V(A_{\bar{P}} S_B)$  exceeds  $V(A_P S_B)$ , of course, since not paying dominates paying. But the actual value of paying is not  $V(A_P S_B)$ ; rather, the actual value is  $V(A_P S_{\bar{B}})$ , since the outcome that would result if the agent were to pay the extortionist is that she would have \$10 fewer and an unbroken windshield. Moreover,  $V(A_P S_{\bar{B}})$  exceeds  $V(A_{\bar{P}} S_B)$ .

But, importantly, while actual value does not respect dominance, it does respect causal dominance. A state is *causally act-independent* for an agent if and only if the agent knows that she has no causal influence over whether the state obtains. (More formally,  $S$  is causally act-independent for an agent if and only if the agent knows, for each  $A \in \mathcal{A}$ ,  $S \leftrightarrow [A \Box \rightarrow S]$ .)

---

<sup>15</sup>If counterfactual excluded middle fails, unchosen options might fail to have actual values. When  $[A \Box \rightarrow O]$  is true, the chance of  $O$  conditional on  $A$ , i.e.  $Ch(O | A)$ , is one, so we could broaden the notion of actual value by setting  $av(A)$  equal to  $\sum_O Ch(O | A)V(O)$ . But it is unclear whether the broadened notion of actual value can do the metaethical work done by the narrower notion.



If there is a set of causally act-independent states  $\mathcal{S}$  such that, for every  $S \in \mathcal{S}$ ,  $V(A_i S)$  exceeds  $V(A_j S)$ , then option  $A_i$  *causally dominates* option  $A_j$ .<sup>16</sup> The alleged connection between causal dominance and rational choice is structurally identical to the alleged connection between dominance and rational choice:

**Causal Dominance:** If option  $A_i$  causally dominates option  $A_j$ , then it is not rationally permissible for the agent to choose  $A_j$ .

But Causal Dominance is more plausible than Dominance. Causal Dominance avoids the absurd recommendation, in *The Extortionist*, that the moviegoer rationally ought to not pay.<sup>17</sup>

Causal Dominance is weaker than Dominance but still strong enough to entail two-boxing. Let the set of states be  $\{S_{\$0}, S_{\$M}\}$ , where  $S_{\$0}$  is the proposition that the opaque box contains \$0 and  $S_{\$M}$  is the proposition that the opaque box contains \$1,000,000. Since

$$V(A_{2B}S_{\$0}) = 1,000 > 0 = V(A_{1B}S_{\$0}), \text{ and}$$

$$V(A_{2B}S_{\$M}) = 1,001,000 > 1,000,000 = V(A_{1B}S_{\$M}),$$

$A_{2B}$  dominates  $A_{1B}$ . Moreover, the agent knows that she has no causal influence over the amount of money in the opaque box, so  $S_{\$0}$  and  $S_{\$M}$  are causally act-independent states. Hence, taking both boxes causally dominates taking only the opaque box, a fact exploited in the *Causal Dominance Argument* for two-boxing:

(P1) If option  $A_i$  causally dominates option  $A_j$ , then it is not rationally permissible for the agent to choose  $A_j$ .

(P2) In *Classic Newcomb*, taking both boxes causally dominates taking only the opaque box.

---

<sup>16</sup>If  $A_i$  causally dominates  $A_j$ , then  $V(A_i S_{\@}) > V(A_j S_{\@})$ . Since  $S_{\@}$  is causally act-independent,  $[A_i \square \rightarrow A_i S_{\@}]$  and  $[A_j \square \rightarrow A_j S_{\@}]$  both are true, so  $av(A_i) = V(A_i S_{\@})$  and  $av(A_j) = V(A_j S_{\@})$ . Hence,  $av(A_i) > av(A_j)$ .

<sup>17</sup>The moviegoer knows that she exerts causal influence over the future state of her windshield, so neither  $S_B$  nor  $S_{\bar{B}}$  is causally act-independent.

(C) Therefore, an agent facing *Classic Newcomb* is rationally required to take both boxes.

The Causal Dominance Argument is the aforementioned standard argument for two-boxing.<sup>18</sup>

Note the intimate relation between Causal Dominance and the rule of c-expected value. Given a set of causally act-independent states  $\mathcal{S}$ , the c-expected value of an option can be characterized as a function of the agent's unconditional credences in the members of  $\mathcal{S}$ :

$$cev(A) = \sum_O C(A \square \rightarrow O)V(O) = \sum_S C(S)V(AS).$$

As the last sum in the equation makes clear, a causally dominated option cannot maximize c-expected value: the rule of c-expected value entails Causal Dominance.

## 4

We believe that Causal Dominance is false, and hence that the Causal Dominance Argument is unsound. We will offer two counterexamples to the rule of c-expected value, and then transform them into counterexamples to Causal Dominance.

The first counterexample is non-ideal. An ideal agent is both introspective—she knows all of the facts about her own beliefs and desires—and logically omniscient. A non-ideal agent is introspective but not logically omniscient. An ideal counterexample features an ideal agent, and a non-ideal counterexample, like the following, features a non-ideal agent:

*The Fire.* The fire alarm rings and the agent, a firefighter, hurries onto the truck. On the ride over she deliberates. She has three options: she can enter the building through the left door, the middle door, or the right door. Since she does not know the exact distribution of residents in the building, she does not know which option will result in the most rescues. Based on her credences about the distribution of

---

<sup>18</sup>Some prefer an informational variant; see, for example, Pollock (2010, p. 57-82). Not every two-boxer relies on dominance reasoning. See, for example, Levi (1975).

residents, she calculates the c-expected value of each option and writes the value on a notecard. After exiting the truck and attaching the water hose, she races toward the building. She reaches into her pocket, but the notecard is gone! Time is of the essence. She knows that all of the residents will die in the time it would take her to recalculate the c-expected values. Her credences about the distribution of residents are unchanged, so she knows that her current c-expected values are what they were when she calculated them. But she cannot fully remember the results of her calculations. She remembers that the c-expected value of entering through the middle door is 9. Of the other two options, she remembers that one has a c-expected value of 0 and that the other has a c-expected value of 10, but she cannot remember which c-expected value goes with which option. (In fact, entering through the right door has a c-expected value of 10, as the lost notecard attests.)<sup>19</sup>

We say that the agent facing *The Fire* is rationally required to enter through the middle door, even though it is true, by hypothesis, that the option that uniquely maximizes c-expected value is entering through the right door.<sup>20</sup>

The second counterexample is ideal.

*The Frustrater.* There is an envelope and two opaque boxes, A and B. The agent has three options: she can take box A, box B, or the envelope. (The three options may be labeled  $A_A$ ,  $A_B$ , and  $A_E$ , respectively.) The envelope contains \$40. The two boxes

---

<sup>19</sup>*The Fire* is an elaboration of a case discussed by Kagan (MS). The fact that non-ideal agents are not always able to access expected value is also discussed in, among other places, Feldman (2006) and Weirich (2004, ch. 5). Some philosophers believe that decision theory applies only to ideal agents, and hence that examples like *The Fire* cannot be relevant to decision theory. We think that decision theory should extend to non-ideal agents. But, even if decision theory applies only to ideal agents, we think that non-ideal cases, such as *The Fire*, help shed light on how an agent must be related epistemically to a value quantity in order to be rationally required to maximize that value quantity, an issue that we discuss in more detail in sections 5, 6, and 7.

<sup>20</sup>We think that the intuition in this case speaks for itself. But we offer a theoretical account of why the agent facing *The Fire* is rationally required to enter through the middle door in note 31.

together contain \$100. How the money is distributed between the boxes depends on a prediction made yesterday by the Frustrater, a reliable predictor who seeks to frustrate. If the Frustrater predicted that the agent would take box A, box B contains \$100. If the Frustrater predicted that the agent would take box B, box A contains \$100. If the Frustrater predicted that the agent would take the envelope, each box contains \$50. The agent knows all of this.<sup>21</sup>

We say that an ideal agent facing *The Frustrater* is rationally required to choose the envelope. But the options that maximize c-expected value are  $A_A$  and/or  $A_B$ , depending on the agent's credences.<sup>22</sup> (*Proof*: No matter what credence function the agent has,  $cev(A_E) = 40$  and  $cev(A_A) + cev(A_B) = 100$ . Two numbers smaller than 40 cannot sum to 100.)<sup>23</sup>

With a few alterations, both *The Fire* and *The Frustrater* can be transformed into counterexamples to Causal Dominance. Start with a variation on *The Fire*:

---

<sup>21</sup>This example is inspired by other purported counterexamples to the rule of c-expected value: Bostrom (2001), Egan (2007) and especially Ahmed (2014b). Lewis (1981) shows that there are realistic cases that have the same structure as *Classic Newcomb*. There are also realistic cases that have the same structure as *The Frustrater*. A commuter wants to get home as soon as possible. There are three routes home: highway A, highway B, and the ferry. The commuter knows that the ferry is second-fastest. The commuter does not know which highway is faster—that varies depending on the day—but the commuter knows that one of the highways is slightly faster than the ferry and that the other is much, much slower. Moreover, the commuter reasonably believes that commuters are like-minded. Conditional on taking highway A/B, the commuter is highly confident that highway B/A is the fastest route.

<sup>22</sup>We assume that the agent facing *The Frustrater* cannot play a mixed strategy. Perhaps the agent is unable to randomize her choice, or perhaps it is simply unwise to play a mixed strategy, since the Frustrater is very good at detecting whether an agent is playing a mixed strategy and punishes the agent severely for doing so.

<sup>23</sup>If we transform *The Frustrater* into a sequence of choices—first a choice between  $A_E$  and eliminating  $A_E$ , and then a choice, if  $A_E$  is eliminated, between  $A_A$  and  $A_B$ —the rule of c-expected value as applied to the sequence recommends  $A_E$ . We note three things. First, this is a different decision problem. *The Frustrater* remains a counterexample to the rule of c-expected value. Second, it may not be rationally permissible for the agent to choose between  $A_E$  and eliminating  $A_E$ —perhaps because the Frustrater punishes agents who do so. Third, not all of the counterexamples to the rule of c-expected value can be transformed into a sequence of choices, cf. Egan (2007). Thanks to Bernhard Salow and Caspar Hare for discussion on this point.

*The Dominating Fire.* Everything is the same as in *The Fire*, except that, unbeknownst to the agent, the option of entering through the right door causally dominates the other two options.

From the standpoint of rationality, *The Dominating Fire* is no different than *The Fire*. A non-ideal agent might not be in a position to know which options causally dominate which others. (We can imagine that the  $V(AS)$ 's are stored in the agent's brain, in the form of a payoff matrix, and that it takes the agent a non-trivial amount of time to survey the matrix.) If an agent is not in a position to know that an option is causally dominated, then the fact that the option is causally dominated is not relevant to what the agent rationally ought to choose. Therefore, as in *The Fire*, an agent facing *The Dominating Fire* is rationally required to enter through the middle door, even though entering through the middle door is causally dominated by entering through the right door.

Causal Dominance is an elimination principle, which marks options as rationally impermissible to choose. But it entails the following selection principle:

**Causal Dominance Selection:** If option  $A_i$  causally dominates all other options, then the agent is rationally required to choose  $A_i$ .

*The Dominating Fire* is a counterexample not just to Causal Dominance, but also to Causal Dominance Selection.

There are no ideal counterexamples to Causal Dominance Selection, a fact that we will return to, and explain, later. But there are ideal counterexamples to Causal Dominance:

*The Semi-Frustrater.* There are two buttons, a white button and a black button. The agent has four options: she can press either button with either hand. (The four options may be labeled  $A_{RH:W}$ ,  $A_{LH:W}$ ,  $A_{RH:B}$ , and  $A_{LH:B}$ .) The white button connects to the white box, the black button connects to the black box, and the agent will receive the contents of whatever box is connected to the button she presses.

One of the boxes contains \$0 and the other contains \$100. Which box contains which sum depends on a prediction made yesterday by the Semi-Frustrater. The Semi-Frustrater seeks to frustrate. If the Semi-Frustrater predicted that the agent would press the black button, the white box contains \$100. If the Semi-Frustrater predicted that the agent would press the white button, the black box contains \$100. There are two left-right asymmetries. First, the agent will receive an extra \$5 if she presses a button right-handedly. Second, because the Semi-Frustrater bases her prediction on a scan of merely half of the agent’s brain, the Semi-Frustrater is a 90% reliable predictor of right-handed button pressings but only a 50% reliable predictor of left-handed button pressings. The agent knows all of this.

We say that *The Semi-Frustrater*, like *The Frustrater*, is an ideal counterexample to the rule of c-expected value. In our view, an ideal agent facing *The Semi-Frustrater* is rationally required to choose  $A_{LH:W}$  or  $A_{LH:B}$ , and rationally permitted to choose either, even though the options that maximize c-expected value are, depending on the agent’s credences,  $A_{RH:W}$  and/or  $A_{RH:B}$ .<sup>24</sup> What is more surprising is that we have an ideal counterexample to Causal Dominance. The claim that an (ideal) agent is never rationally permitted to choose a (strictly) causally dominated option is a staple of game theory, where it appears in textbooks as the injunction against playing strategies that can be iteratively eliminated by (strict) causal domination,<sup>25</sup> and is regarded as sacrosanct by many expert decision theorists.<sup>26</sup> But  $A_{RH:W}$

---

<sup>24</sup>Either  $S_W$ , the white box contains \$100, or  $S_B$ , the black box contains \$100. Since the agent knows that her choice has no causal influence over the contents of the boxes,  $\{S_W, S_B\}$  is a set of causally act-independent states. Equating dollars and units of value,  $V(S_W A_{RH:W}) = 105 = 5 + V(S_W A_{LH:W})$ ;  $V(S_B A_{RH:W}) = 5 = 5 + V(S_B A_{LH:W})$ ;  $V(S_W A_{RH:B}) = 5 = 5 + V(S_W A_{LH:B})$ ; and  $V(S_B A_{RH:B}) = 105 = 5 + V(S_B A_{LH:B})$ . So  $cev(A_{RH:W})$  maximizes if  $C(S_W) \geq 0.5$ , and  $cev(A_{RH:B})$  maximizes if  $C(S_B) \geq 0.5$ .

<sup>25</sup>See, for example, Fudenberg and Tirole (1991, ch. 2) and Myerson (1991, s. 3.1).

<sup>26</sup>Briggs (2015, p. 836): “The following is an independently compelling claim about rationality: if it is knowable a priori that strategy  $a$  yields a better result than strategy  $b$ , then it is pragmatically irrational to choose strategy  $b$  when strategy  $a$  is available.” Pettigrew (2015, p. 806): “the so-called Dominance Principle, which says that an

causally dominates  $A_{LH:W}$ , and  $A_{RH:B}$  causally dominates  $A_{LH:B}$ , so an ideal agent facing *The Semi-Frustrater* is rationally required to choose a causally dominated option.

One might object that *The Semi-Frustrater* is really no different than *Classic Newcomb*. In both examples there is some intuitive pull toward choosing a causally dominated option, since in both examples consistently choosing a causally dominated option results in greater long run wealth. Consistent one-boxers end up wealthier than do consistent two-boxers. Consistent left-handers end up wealthier than do consistent right-handers. Two-boxers resist the one-boxing intuition, so, if the intuition that an agent rationally ought to press a button left-handedly in *The Semi-Frustrater* is really no different than the intuition that an agent rationally ought to take only the opaque box in *Classic Newcomb*, two-boxers should also resist the left-handed intuition.

But, in at least two respects, *The Semi-Frustrater* and *Classic Newcomb* are importantly different. The first and most important difference is metaethical—that is, it concerns how the agent is epistemically related to the relevant value quantities. For an agent facing *Classic Newcomb*, maximizing c-expected value is non-accidentally doable. If the agent seeks to maximize c-expected value, she can be confident both about which option she will choose and that she will choose an option that maximizes c-expected value. She will be confident that she will take both boxes and confident that by doing so she will choose an option that maximizes c-expected value. By contrast, for an agent facing *The Semi-Frustrater*, maximizing c-expected value is doable only accidentally. An agent facing *The Semi-Frustrater*, who seeks to maximize c-expected value, cannot be confident both about which option she will choose and that she will choose an option that maximizes c-expected value. If she is confident about which option she will choose, then she is confident that by choosing that option she will *fail* to maximize c-expected value. As we say in more detail below, in our view, the fact that the c-expected value

---

option is irrational if there is an alternative that is guaranteed to be better than it, and if there is nothing that is guaranteed to be better than that alternative [...] is an uncontroversial principle of decision theory.” Also see, for example, Buchak (2015), Briggs (2010), Gibbard and Harper (1978), Lewis (1981), Joyce (1999), Nozick (1969), Sobel (1994) and Skyrms (1984). In epistemic decision theory, too, the claim that (strictly) dominated options are *ipso facto* irrational is relied upon heavily. See, for example, Joyce (1998).

of an option exceeds the c-expected value of another option is relevant to what an agent rationally ought to choose only if the agent is appropriately related to c-expected value maximization. An agent facing *Classic Newcomb* is appropriately related to c-expected value maximization, but an agent facing *The Semi-Frustrater* is not. Ultimately it is this epistemological and metaethical difference that marks the crucial divide between *The Semi-Frustrater* and *Classic Newcomb*.

But even before we get into metaethics, there is a simple descriptive difference between *The Semi-Frustrater* and *Classic Newcomb*. In *Classic Newcomb*, there is unequal environmental fortune. Imagine that the choices are made in a room containing only the agent and the boxes. Consistent one-boxers almost always make their choices in lucrative rooms: they almost always choose between two options, each worth \$1,000,000 or more, in a room that contains more than \$1,000,000. Consistent two-boxers almost always make their choices in impoverished rooms: they almost always choose between two options, each worth no more than \$1,000, in a room that contains \$1,000. The argument for one-boxing, based on the claim that consistent one-boxers are wealthier than consistent two-boxers, is undermined by the unequal environmental fortune. What explains why consistent one-boxers are wealthier than consistent two-boxers is that one-boxers make their choices in lucrative rooms, not that one-boxers choose more wisely.<sup>27</sup> Consistently choosing unwisely in lucrative rooms leads to greater long run wealth than does consistently choosing wisely in impoverished rooms. Notice, in *The Semi-Frustrater*, however, that there is no difference in environmental fortune. Like consistent left-handers, consistent right-handers always make their choices in rooms that contain exactly \$105. What explains why consistent left-handers are wealthier than consistent right-handers is a difference of rationality, not a difference of environmental fortune. Consistent right-handers end up poorer than do consistent left-handers because they choose irrationally.

---

<sup>27</sup>For more on this point, see Wells (Forthcoming).



Although the Causal Dominance Argument is unsound, a successful, independent argument for two-boxing is in the nearby vicinity. The successful argument relies on a metaethical principle connecting guidance and actual value maximization to rational choice.

There are two ‘ought’s of decision-making, an objective ‘ought’ and a rational ‘ought’. Decision theory, being consequentialist in nature, takes both to be reducible to value quantity maximization.

The objective ‘ought’ reduces to actual value maximization: agents are always objectively required to choose so as to maximize actual value.

The objective ‘ought’ is not our main concern. Our main concern is the rational ‘ought’, which can, and often does, come apart from the objective ‘ought’. For example:

*Boxes like Miners.* There are three opaque boxes: the left box, the middle box, and the right box. The agent must choose exactly one box. The agent knows that the middle box contains \$9. Of the other two boxes, the agent knows that one contains \$0 and that the other contains \$10, but does not know which box contains which sum. (In fact, the right box contains \$10.)

An agent facing *Boxes like Miners* is, though objectively required to choose the right box, rationally required to choose the middle box.

At the metaethical level, the most important difference between the objective ‘ought’ and the rational ‘ought’ is a difference of guidance. The objective ‘ought’ is not always capable of guiding the agent’s choice. Actual value is the value quantity the maximization of which makes options objectively permissible for the agent to choose, but agents are not always capable of being guided by actual value. A necessary condition on being capable of being guided by actual value is being in a position to know of some option that it maximizes actual value, and agents often are in no such position. An agent facing *Boxes like Miners*, for example, is in no such position.

The rational ‘ought’, by contrast, is always capable of guiding the agent’s choice. An agent is

always capable of being guided by the value quantity the maximization of which makes options rationally permissible for the agent to choose.

It is here that we break with the metaethical orthodoxy. The orthodoxy has it that a value quantity is choice-guiding if and only if the facts about which options maximize the value quantity supervene on the facts about the agent's beliefs and desires.<sup>28</sup> Actual value fails this supervenience test. The actual value of an option is a function of the truth-values of certain counterfactual claims, and such truth-values float free of the agent's psychology. By contrast, e-expected value and c-expected value pass the supervenience test. Both are functions of the agent's beliefs and desires.

In our view, the orthodoxy is mistaken twice over. First, it is a mistake to try to divide value quantities into those that are, and those that are not, choice-guiding. Whether a value quantity is capable of guiding an agent's choice is settled, in our view, occasion by occasion, not once and for all. Second, it is a mistake to identify choice-guidance with supervenience on the agent's beliefs and desires. On some occasions an agent is capable of being guided by a value quantity that does not supervene on her beliefs and desires, and on some occasions an agent is incapable of being guided by a value quantity that supervenes on her beliefs and desires.

We claim that a value quantity is capable of guiding an agent's choice on an occasion if and only if the agent has stable access to the value quantity on that occasion. Stable access is defined in terms of being in a position to know. An agent is *in a position to know* a proposition if and only if there is no obstacle blocking her from knowing the proposition.<sup>29</sup> An agent has *access* to a value quantity  $Q$  if and only if there is an option  $A \in \mathcal{A}$  such that the agent is in a position to know of  $A$  that it maximizes  $Q$ . An agent has *stable access* to  $Q$  if and only if there is an option  $A \in \mathcal{A}$  such that (i) the agent is in a position to know of  $A$  that it maximizes  $Q$ , and (ii) conditional on  $A$ , the agent still is in a position to know of  $A$  that it maximizes  $Q$ .<sup>30</sup> If an agent has stable access

---

<sup>28</sup>Or, more generally, supervene on the agent's internal mental states. See, for example, Conee and Feldman (2004).

<sup>29</sup>Cf. Williamson (2000, p. 95).

<sup>30</sup>By "conditional on  $A$ ," we have the following in mind. Take the agent's credence function and conditionalize

to  $Q$ , then the agent is *stably* in a position to know of some option  $A$  that it maximizes  $Q$ .

The stability is crucial. An agent who chooses option  $A$  is guided by  $Q$  only if she can know both that she will choose  $A$  and that  $A$  is  $Q$ -maximizing. It is for this reason that access alone is not sufficient for guidance. An agent who has access but lacks stable access to  $Q$  cannot know both which option she will choose and that the option she will choose is  $Q$ -maximizing. Such an agent either is surprised by which option she chooses, in which case her choice is not guided at all, or she anticipates choosing an option that is not  $Q$ -maximizing, in which case her choice is not guided by  $Q$ . Stability plugs this gap. An agent who has stable access to  $Q$  is capable of being guided by  $Q$  because she can know both that she will choose  $A$  and that  $A$  is  $Q$ -maximizing.<sup>31</sup>

## 6

If an agent is incapable of being guided by a value quantity, then the maximization of that value quantity is not what makes options rationally permissible for the agent to choose. In our view, this is the fact that explains why the rule of c-expected value admits of counterexamples. Agents are not always capable of being guided by c-expected value—that is, agents do not always have stable access to c-expected value. An agent facing *The Fire* or *The Dominating Fire* does not have access, let alone stable access, to c-expected value, since the external time constraints, together with the agent's limited powers of deduction, form an obstacle blocking her from knowing that entering through the right door maximizes c-expected value.<sup>32</sup> An agent facing *The Frustrater* or

---

it on  $A$ . Then ask whether the agent still is in a position to know that  $P$ , relative to her updated credence function. If she is, then she is stably in a position to know that  $P$ . If not, not.

<sup>31</sup>*Question:* Does choice-guidance supervene on the agent's mental states? *Answer:* If being in a position to know is a mental state, then yes. Otherwise, no.

<sup>32</sup>*Question:* What value quantity is an agent facing *The Fire* rationally required to maximize? *Answer:* A value quantity that stands to c-expected value as c-expected value stands to actual value; we might call it c-expected<sub>2</sub> value. A  $cev(A)$ -level proposition is of the form  $[cev(A) = v]$ . Just as the c-expected value of an option is a credence-weighted average of the agent's hypotheses about the actual value of the option (see note 9), the c-expected<sub>2</sub> value of an option is a credence-weighted average of the agent's hypotheses about the c-expected value of the option:

*The Semi-Frustrater* has access but lacks stable access to c-expected value, since there is no option available in either decision problem that maximizes c-expected value conditional on itself.<sup>33</sup> We claim that the rule of c-expected value admits of counterexamples only when agents lack stable access to c-expected value. In other words, we accept the *guiding rule of c-expected value*: that agents who have stable access to c-expected value are rationally required to choose so as to maximize c-expected value.

In spirit, the guiding rule of c-expected value is similar to the ratificationisms defended by Harper, Jeffrey, Sobel, and others.<sup>34</sup> But at the level of detail, the views differ in important ways, and we think that the guiding rule of c-expected value is an improvement upon the more familiar ratificationisms.

The key notion for any ratificationism is that of an option being ratifiable. An option  $A$  is *ratifiable* if and only if, conditional on  $A$ ,  $A$  maximizes c-expected value, and *nonratifiable*, otherwise. There are two sorts of ratificationisms. According to *principled* ratificationism, it is never rationally permissible to choose nonratifiable options. Principled ratificationism is implausible: *The Frustrater* and *The Semi-Frustrater* are counterexamples. According to *lexical* ratificationism, the more plausible version of the view, options are lexically ordered by

---

$cev_2(A) = \sum_v vC([cev(A) = v])$ . There is also a value quantity that stands to c-expected<sub>2</sub> value as c-expected value stands to actual value; we might call it c-expected<sub>3</sub> value. In general, for any  $n > 1$ ,

$$cev_n(A) = \sum_v vC([cev_{n-1}(A) = v]).$$

<sup>33</sup>*Question*: What value quantity is an agent facing *The Frustrater* rationally required to maximize? *Answer*: The most causally fine-grained expected value quantity to which the agent has stable access. See Spencer and Wells (MS), in which we develop a theory of rational choice in the face of decision instability. Much of decision theory rests on the assumption that there is a single value quantity that any agent facing any decision problem is rationally required to maximize. We reject this claim. We explore the prospects for a unified decision theory that rejects this assumption.

<sup>34</sup>For discussion of ratificationism, see, among others, Eells (1982), Egan (2007), Gustafsson (2011), Hare and Hedden (2016), Harper (1986), Jeffrey (1983), Joyce (2007), Rabinowicz (1988), Sobel (1994), Skyrms (1984), and Weirich (1988, 2004).

ratifiability: ratifiable options are infinitely more choiceworthy than are nonratifiable options; hence it is rationally permissible to choose nonratifiable options only if none of the available options are ratifiable. Lexical ratificationists can disagree with one another about how to choose among options at the same lexical order. Some lexical ratificationists use e-expected value to choose among options at the same lexical order.<sup>35</sup> They claim that an agent is rationally required to choose a ratifiable option the e-expected value of which is not exceeded by that of any other ratifiable option, unless there are no ratifiable options, in which case the agent is rationally required to choose a nonratifiable option that maximizes e-expected value. Other lexical ratificationists use c-expected value to choose among options at the same lexical order.<sup>36</sup> Note that both forms of lexical ratificationism entail two-boxing, since, in *Classic Newcomb*, taking both boxes is the only ratifiable option.

The fatal flaw in lexical ratificationism is the lexical ordering of options. By treating ratifiable options as infinitely more choiceworthy than nonratifiable options, lexical ratificationists effectively claim that ratifiability is an infinite value. But ratifiability is not a value, let alone an infinite one. An option is not made more choiceworthy by being ratifiable. The counterexamples to lexical ratificationism, like the following, due to Skyrms (1984), exploit precisely this flaw:

*Three Shells.* The agent has three options: there are three shells, shell J, shell K, and shell L, and the agent can choose any one of them. How much money is contained in each shell depends on a prediction made yesterday by a reliable predictor. If the predictor predicted that the agent would choose shell J, then shell J contains \$1, shell K contains \$0, and shell L contains \$0. If the predictor predicted that the agent would choose shell K, then shell J contains \$0, shell K contains \$9, and shell L contains \$10. If the predictor predicted that the agent would choose shell L, then shell J contains \$0, shell K contains \$10, and shell L contains \$9.

---

<sup>35</sup>Cf. Jeffrey (1983).

<sup>36</sup>See Egan (2007) for discussion.

Choosing shell J is the only ratifiable option. Hence, according to any lexical ratificationism, an agent facing *Three Shells* is rationally required to choose shell J, no matter what credence function the agent has. Lexical ratificationism is thereby refuted. An agent facing *Three Shells* is rationally required to choose shell J only if, at the time of decision, she is highly confident that she will. If she is highly confident that she will choose shell J, then she is highly confident that shell J contains \$1 and that the other two shells contain nothing. But if the agent is not highly confident that she will choose shell J, then it is not even rationally permissible for her to choose shell J. In the extreme case, in which the agent is highly confident that she will not choose shell J, the agent regards shell J as the worst of her three options by far. The claim, made by lexical ratificationists, that the agent is nevertheless rationally required to choose shell J is clearly false.

The guiding rule of c-expected value avoids the counterexamples to lexical ratificationism by rectifying the fatal flaw. The guiding rule of c-expected value is not built on the notion of ratifiability. It is built on the notion of having stable access to a value quantity—specifically, having stable access to c-expected value. Whereas ratifiability operates at the ethical level, affecting the choiceworthiness of options, stable access operates at the *metaethical* level, affecting the relevance of value quantities. Stable access is a necessary condition for rational relevance: facts about which options maximize a given value quantity can be relevant to what an agent rationally ought to choose only if the agent has stable access to that value quantity.

The guiding rule of c-expected value thus can handle cases like *Three Shells* correctly. If an agent facing *Three Shells* is highly confident that she will choose shell J, then she has stable access to c-expected value, and the guiding rule of c-expected value correctly entails that she is rationally required to choose shell J. If the agent is not highly confident that she will choose shell J, then she lacks stable access to c-expected value, and the guiding rule of c-expected value is silent. When an agent, ideal or nonideal, lacks stable access to c-expected value, the facts about the c-expected values of options are not relevant to what the agent rationally ought to choose.

Since we accept the guiding rule of c-expected value, we think that there is a sound argument from c-expected value to two-boxing. A competent agent facing *Classic Newcomb* has stable

access to c-expected value because (i) she is in a position to know that  $A_{2B}$  (uniquely) maximizes c-expected value, and (ii) conditional on  $A_{2B}$ , she still is in a position to know that  $A_{2B}$  (uniquely) maximizes c-expected value. Hence, by the guiding rule of c-expected value, she is rationally required to take both boxes.

But, as noted above, arguing from c-expected value to two-boxing fails to move the debate forward. What we need is an independent argument for two-boxing.

## 7

We think that the best independent argument for two-boxing goes through the guiding rule of actual value.

According to the *rule of actual value*, agents are always rationally required to choose so as to maximize actual value. Everyone rejects the rule of actual value, and for good reason. Counterexamples abound. Rational permission and actual value maximization often come apart. But the rule of c-expected value also fails: rational permission and c-expected value maximization also come apart. An agent is rationally required to choose so as to maximize c-expected value only when she has stable access to c-expected value. We think that the same holds for actual value. We accept the *guiding rule of actual value*: that agents who have stable access to actual value are rationally required to choose so as to maximize actual value.

The guiding rule of actual value entails the uncontroversial claim that rational permission and actual value maximization can come apart when agents lack access to actual value. In *Boxes like Miners*, for example, the agent is rationally required to choose the middle box, even though choosing the right box uniquely maximizes actual value.

The guiding rule of actual value also entails that rational permission and actual value maximization can come apart when an agent has access but lacks stable access to actual value. Not much attention has been paid to the question of whether rational permission and actual value maximization can come apart in such cases, in part because it requires some fancy

footwork to devise an example. Here is one:

*Unstable Boxes like Miners.* There are four boxes, the outside-left box, the middle-left box, the middle-right box, and the outside-right box. The outside boxes are opaque and the middle boxes are transparent. The middle-left box and the middle-right box each contain \$9. One of the outside boxes contains \$0 and the other contains \$10. Which outside box contains which sum depends on a prediction made yesterday by a reliable predictor. If the predictor predicted that the agent would choose either the middle-left box or the outside-left box, the outside-right box contains \$10. If the predictor predicted that the agent would choose either the middle-right box or the outside-right box, the outside-left box contains \$10. The agent knows all this. The agent also believes that she will choose the middle-left box.

Since the agent believes that she will choose the middle-left box and believes that the predictor is extremely reliable, she believes that the outside-right box contains \$10. Moreover, let us suppose that it is true that the outside-right box contains \$10. Presumably, then, if the predictor is reliable enough, the agent *knows* that choosing the outside-right box uniquely maximizes actual value. But her epistemic position is unstable. Conditional on choosing the outside-right box, she ceases to be in a position to know that choosing the outside-right box maximizes actual value. It seems to us clear that an agent facing *Unstable Boxes like Miners* is rationally required to choose either the middle-left box or the middle-right box, and rationally permitted to choose either. Even when an agent knows which option uniquely maximizes actual value, rational permission and actual value maximization can come apart, if the agent's knowledge is unstable. This is a somewhat surprising result.

But, as concerns *Classic Newcomb*, the real substance of the guiding rule of actual value is what it says about the coincidence between rational permission and actual value maximization: namely, that rational permission and actual value maximization cannot come apart if the agent has stable access to actual value.



The simplest cases in which an agent has stable access to actual value are cases in which the agent knows the actual values of her options. (Imagine an agent choosing among transparent boxes, each containing a sum of money.) *Classic Newcomb* is interesting in part because it is a case in which the agent has stable access to actual value without being in a position to know the actual values of her options. The agent is not in a position to know whether the actual value of  $A_{2B}$  is 1,000 or 1,001,000, for example, because she does not know whether the opaque box contains \$0 or \$1,000,000. Nevertheless, she is in a position to know that taking  $A_{2B}$  (uniquely) maximizes actual value, and, conditional on  $A_{2B}$ , she still is in a position to know that  $A_{2B}$  (uniquely) maximizes actual value.

The guiding rule of actual value entails that if an agent is stably in a position to know of an option that it uniquely maximizes actual value, the agent is rationally required to choose the option. This claim is the crucial premise of the *Objective Argument* for two-boxing:

(P1) If an agent is stably in a position to know of an option that it uniquely maximizes actual value, then the agent is rationally required to choose the option.

(P2) An agent facing *Classic Newcomb* is stably in a position to know of taking both boxes that it uniquely maximizes actual value.

(C) Therefore, an agent facing *Classic Newcomb* is rationally required to take both boxes.

The Causal Dominance Argument and the Objective Argument are closely related. If  $A_i$  causally dominates  $A_j$ , then, unless the agent is otherwise epistemically disabled, the agent is stably in a position to know that the actual value of  $A_i$  exceeds the actual value of  $A_j$ . Pointing out that taking both boxes causally dominates taking only the opaque box therefore helps to justify the minor premise of the Objective Argument.

The crucial difference between the Causal Dominance Argument and the Objective Argument lies in their respective major premises.<sup>37</sup> The major premise of the Objective

---

<sup>37</sup>Ahmed (2014a, ch. 7) claims that the best argument for two-boxing goes through a principle, akin to Causal

Argument amounts to the claim that agents are rationally required to be guided by actual value when they are capable of being guided by actual value. Or to put the point in deontological terms (since agents are always objectively required to choose so as to maximize actual value): in the rare cases in which the objective ‘ought’ provides the agent with guidance, the guidance provided by the rational ‘ought’ cannot conflict with the guidance provided by the objective ‘ought’.<sup>38</sup>

The major premise of the Causal Dominance Argument—namely, Causal Dominance—is refuted by cases like *The Dominating Fire* and *The Semi-Frustrater*. But such cases pose no threat to the major premise of the Objective Argument, since they are not cases in which the agent has stable access to actual value.

---

Dominance, which he calls CDB: “If you know that a certain available option makes you worse off, given your situation, than you would have been on some identifiable alternative, then that first option is irrational” (p. 202). He then formulates a weaker principle, CDB-sequence: “If you know that a certain available sequence of choices makes you worse off, given your situation, than you would have been on some identifiable alternative, then that first sequence is irrational” (p. 211, italics original). He offers a counterexample to CDB-sequence and argues that “accepting CDB and not CDB-sequence looks completely unmotivated” (p. 211). As it turns out, both *Unstable Boxes like Miners* and *The Semi-Frustrater* are counterexamples to CDB. But we do not need anything nearly as strong as CDB to motivate two-boxing. Neither *Unstable Boxes like Miners* nor *The Semi-Frustrater* are counterexamples to the guiding rule of actual value. As for Ahmed’s counterexample to CDB-sequence—namely, *Newcomb Insurance*—it matters whether there is a single choice or a sequence of choices, since the value quantities to which the agent has stable access depends on it. If there is a single choice, even a single choice among sequences, we agree with Ahmed’s judgments. If there is a sequence of choices, each among non-sequential options, we agree with the recommendations of the rule of c-expected value.

<sup>38</sup>Kotzen (MS) also proposes a connection between the objective ‘ought’ and the rational ‘ought’ and suggests that the proposed connection justifies two-boxing. In broad strokes, we agree. But at the level of detail, our suggestion is importantly different from Kotzen’s. Kotzen’s proposed connection, unlike ours, does not require *stable* access. As such, it is false: *Unstable Boxes like Miners* is a counterexample, as is a variation on *The Semi-Frustrater* in which the agent knows that she will not press the black button.

The second premise of the Objective Argument is uncontroversial. An agent facing *Classic Newcomb* is stably in a position to know of taking both boxes that it uniquely maximizes actual value. If one-boxers want to resist the Objective Argument, they must reject the guiding rule of actual value.

The guiding rule of actual value is an instance of a schema that all sides should accept. The schema involves two elements: the objective value quantity, which is the value quantity the maximization of which makes options objectively permissible to choose, and the guidance relation, which an agent facing a decision problem bears to value quantities. The schema is an objective guidance constraint on the rational ‘ought’:

**Objective Guidance:** If an agent facing a decision problem bears the guidance relation to the objective value quantity, then the options that are rationally permissible for the agent to choose must maximize the objective value quantity.

To get from Objective Guidance to the guiding rule of actual value, we need two further claims: that the guidance relation is stable access, and that actual value is the objective value quantity.

One-boxers could disagree with our claim that the guidance relation is stable access, but this will not be of much help in resisting the Objective Argument. After all, the Objective Argument is not wedded to any particular conception of guidance. It can be recast using any conception of guidance on which an agent facing *Classic Newcomb* bears the guidance relation to actual value, and, so far as we can tell, *every* plausible conception of guidance is one on which an agent facing *Classic Newcomb* bears the guidance relation to actual value. To resist the Objective Argument, one-boxers therefore must deny that actual value is the objective value quantity.

In principle, there are three ways that one-boxers could deny that actual value is the objective value quantity: they could deny that there is an objective ‘ought’; they could grant that there is an objective ‘ought’, but deny that there is an objective value quantity; or they could argue that some value quantity besides actual value is the objective value quantity. Only the third way is plausible.

The first two ways effectively abandon the consequentialist common ground between one-boxers and two-boxers. One-boxers and two-boxers should agree that, when an agent faces a decision problem, there are facts about which options are objectively permissible for the agent to choose, and that what makes an option objectively permissible for the agent to choose is the maximization of some value quantity, the objective value quantity. The dispute between one-boxers and two-boxers thus reduces to a dispute about what the objective value quantity is. Two-boxing is true if actual value is the objective value quantity. For one-boxing to be true, some other value quantity would have to be the objective value quantity.

It is not clear what one-boxers could take the objective value quantity to be. Heretofore, it has been common ground between one-boxers and two-boxers that actual value is the objective value quantity. The most promising proposal we have yet to see was suggested to us by Arif Ahmed, in personal communication. Ahmed, a committed one-boxer, suggests that one-boxers take *e-actual value* to be the objective value quantity.

We can introduce e-actual value in a way that makes its similarity to actual value clear. Actual value, an explicitly causal notion, can be explicated using e-expected value and knowledge. Let  $\mathcal{S} = \{S_1, \dots, S_n\}$  be the finest partition of causally act-independent states. Let  $S_{@}$  be the member of  $\mathcal{S}$  that is true at the actual world, and let  $C^{S_{@}}$  be the agent's credence function conditionalized on  $S_{@}$ . The *actual value* of an option is then equal to the e-expected value of the option relative to  $C^{S_{@}}$ . The e-actual value of an option also can be explicated using e-expected value and knowledge. Say that a state,  $T$ , is evidentially act-independent just if  $T$  is compossible with each  $A \in \mathcal{A}$  and, for each  $A \in \mathcal{A}$ ,  $C(T) = C(T | A)$ . Let  $\mathcal{T} = \{T_1, \dots, T_m\}$  be the finest partition of evidentially act-independent states. Let  $T_{@}$  be the member of  $\mathcal{T}$  that is true at the actual world, and let  $C^{T_{@}}$  be the agent's credence function conditionalized on  $T_{@}$ . The *e-actual value* of an option is then equal to the e-expected value of the option relative to  $C^{T_{@}}$ . We can think of it this way, then: the actual value of an option is what the agent expects the value of the option to be, given full causal knowledge (that is, knowledge of which member of the finest partition of causally act-independent states is true), and the e-actual value is what

the agent expects the value of the option to be, given full evidential knowledge (that is, knowledge of which member of the finest partition of evidentially act-independent states is true).

There is nothing incoherent about the suggestion that e-actual value is the objective value quantity. But if one-boxing stands and falls with the claim that e-actual value is the objective value quantity, then one-boxers are in serious trouble; for the claim that agents are always objectively required to choose so as to maximize e-actual value is highly unintuitive. To see this, return to two examples from above.

Consider *The Frustrater*, and let us suppose that box A contains \$100, box B contains \$0, and the envelope contains \$40. It seems clear, then, that an agent facing *The Frustrater* is objectively required to take box A, the most lucrative box. But if e-actual value is the objective value quantity, then an agent facing *The Frustrater* is objectively required to take the envelope. The true member of the finest partition of evidentially act-independent states,  $T_{@}$ , does not specify how much money is in box A, since any state that specifies how much money is in the boxes is not evidentially act-independent. Hence, the option that uniquely maximizes e-actual value is taking the envelope.

In *Boxes like Miners*, actual value and e-actual value coincide. Taking the right box uniquely maximizes both. But if we imagine that the agent regards herself as slightly intuitive—that the agent’s credence that the left box contains \$10 conditional on taking the left box is slightly greater than her unconditional credence that the left box contains \$10—then, although taking the right box still uniquely maximizes actual value, and although it still seems that the agent is objectively required to take the right box, taking the middle box uniquely maximizes e-actual value.

There may be other proposals one-boxers could produce about what the objective value quantity is, and we can evaluate them individually. But we doubt that any will be plausible. To us, it seems obvious that an agent facing *Classic Newcomb* is objectively required to take both boxes, so it seems obvious that the objective value quantity, whatever it proves to be, is uniquely maximized by two-boxing in *Classic Newcomb*.

Strictly speaking, the Objective Argument does not require that actual value be the objective value quantity or that stable access be the guidance relation. All that it requires is three claims: (1) that Objective Guidance is true; (2) that two-boxing in *Classic Newcomb* uniquely maximizes the objective value quantity; and (3) that an agent facing *Classic Newcomb* bears the guidance relation to the objective value quantity.

That being said, we believe that the objective value quantity is actual value, and we believe that the guidance relation is stable access. In our view, the guiding rule of actual value is virtually undeniable.

## 9

The foregoing discussion provides us not only with a sound argument for two-boxing, but also with the resources needed to explain why Causal Dominance admits of counterexamples.

Nothing about what an agent rationally ought to do follows from the relations of causal dominance among the options. Of course, both actual value and c-expected value respect causal dominance, so, if  $A_i$  causally dominates  $A_j$ , the actual value of  $A_i$  exceeds the actual value of  $A_j$ , and the c-expected value of  $A_i$  exceeds the c-expected value of  $A_j$ . But nothing about what an agent rationally ought to do follows from the actual values of the options, and nothing about what an agent rationally ought to do follows from the c-expected values of the options. There is no direct connection between dominance or value quantity maximization and rational choice. In order to derive conclusions about what an agent rationally ought to do, we need to know, in addition to the facts about which options maximize which value quantities, the facts about which value quantities the agent has stable access to.

Once we appreciate that stable access mediates the connection between value quantity maximization and rational choice, we can explain the pattern of counterexamples to Causal Dominance that we find. Since both actual value and c-expected value respect causal dominance, and since both the guiding rule of actual value and the guiding rule of c-expected

value are true, we should expect counterexamples to Causal Dominance to arise when, but only when, agents lack stable access both to actual value and to c-expected value. This is exactly what we find. There are non-ideal counterexamples to Causal Dominance because a non-ideal agent may lack stable access both to actual value and to c-expected value, despite the fact that one of her options causally dominates another (e.g., *The Dominating Fire*). There are ideal counterexamples to Causal Dominance because an ideal agent might lack stable access both to actual value and to c-expected value, despite the fact that one of her options causally dominates another (e.g., *The Semi-Frustrater*). There are non-ideal counterexamples to Causal Dominance Selection because a non-ideal agent might lack stable access both to actual value and to c-expected value, despite the fact that one of her options causally dominates all others (e.g., *The Dominating Fire*). There are no ideal counterexamples to Causal Dominance Selection because an ideal agent is guaranteed to have stable access to actual value if one of her options causally dominates all others.<sup>39,40</sup>

## References

Ahmed, A. 2014a. *Evidence, Decision and Causality*. Cambridge University Press.

—. 2014b. “Dicing with Death.” *Analysis* 74:587–94.

---

<sup>39</sup>An ideal agent knows that the states  $\mathcal{S} = \{S_1, S_2, \dots, S_n\}$  are causally act-independent. Hence, for each  $A \in \mathcal{A}$  and each  $S \in \mathcal{S}$ , she knows that  $(S \leftrightarrow [A \square \rightarrow S])$ . The states are causally act-independent, so there is some  $S_j \in \mathcal{S}$  such that, for any  $A \in \mathcal{A}$ ,  $av(A) = V(AS_j)$ . Hence, if option  $A_i$  dominates all other options, the ideal agent knows that  $A_i$  uniquely maximizes actual value. Moreover, conditional on  $A_i$ , the  $V(AS)$ 's remain unchanged, and the ideal agent still knows that some  $S \in \mathcal{S}$  obtains; hence the ideal agent still is in a position to know that  $A_i$  uniquely maximizes actual value.

<sup>40</sup>Thanks to Arif Ahmed, Sara Aronowitz, Dave Chalmers, Mikaël Cozic, Nilanjan Das, Kevin Dorst, Paul Egré, Branden Fitelson, Caspar Hare, Brian Hedden, Wes Holliday, Uriah Kriegel, David Nicolas, François Recanati, Bernhard Salow, Ginger Schultheis, Brad Skow, Bob Stalnaker, Quinn White, the members of the Rutgers Formal Epistemology and Decision Theory Reading Group (2015), and an anonymous referee.

- Bostrom, N. 2001. "The Meta-Newcomb Problem." *Analysis* 61:309–10.
- Briggs, R. 2010. "Decision-Theoretic Paradoxes as Voting Paradoxes." *The Philosophical Review* 119:1–30.
- . 2015. "Costs of Abandoning the Sure-Thing Principle." *Canadian Journal of Philosophy* 45:827–40.
- Buchak, L. 2015. "Revisiting Risk and Rationality: A Reply to Pettigrew and Briggs." *Canadian Journal of Philosophy* 45:841–62.
- Conee, E. and Feldman, R. 2004. *Evidentialism*. Oxford University Press.
- Eells, E. 1982. *Rational Decision and Causality*. Cambridge University Press.
- Egan, A. 2007. "Some Counterexamples to Causal Decision Theory." *Philosophical Review* 116:94–114.
- Feldman, F. 2006. "Actual Utility, the Objection from Impracticality, and the Move to Expected Utility." *Philosophical Studies* 129:49–79.
- Fudenberg, D. and Tirole, J. 1991. *Game Theory*. MIT Press.
- Gibbard, A. and Harper, W. 1978. "Counterfactuals and Two Kinds of Expected Utility." In Leach J. Hooker, A. and E. McClennen (eds.), *Foundations and Applications of Decision Theory*, 125–62. Reidel.
- Gustafsson, J. 2011. "A Note in Defense of Ratificationism." *Erkenntnis* 75:147–150.
- Hájek, A. 2016. "Deliberation Welcomes Prediction." *Episteme* 507–28.
- Hare, C. and Hedden, B. 2016. "Self-Reinforcing and Self-Frustrating Decisions." *Noûs* 50:604–28.
- Harper, W. 1986. "Mixed Strategies and Ratifiability in Causal Decision Theory." *Erkenntnis* 24:25–36.



- Hedden, B. 2015. "Options and Diachronic Tragedy." *Philosophy and Phenomenological Research* 90:423–45.
- Jeffrey, R. 1965. *The Logic of Decision*. University of Chicago Press.
- . 1983. *The Logic of Decision*. 2nd ed. University of Chicago Press.
- Joyce, J. 1998. "A Nonpragmatic Vindication of Probablism." *Philosophy of Science* 65:575–603.
- . 1999. *The Foundations of Causal Decision Theory*. Cambridge University Press.
- . 2002. "Levi on Causal Decision Theory and the Possibility of Predicting One's Own Actions." *Philosophical Studies* 110:69–102.
- . 2007. "Are Newcomb Problems Really Decisions?" *Synthese* 156:537–62.
- Kagan, S. MS. "The Paradox of Methods." Unpublished. Accessed Spring 2012.
- Kotzen, M. MS. "Three Principles of Inference and Deliberation."
- Kyburg, H. 1980. "Acts and Conditional Probabilities." *Theory and Decision* 12:149–71.
- Levi, I. 1975. "Newcomb's Many Problems." *Theory and Decision* 6:161–75.
- . 1997. *The Covenant of Reason: Rationality and the Commitments of Thought*. Cambridge University Press.
- Lewis, D. 1981. "Causal Decision Theory." *Australasian Journal of Philosophy* 59:5–30.
- Myerson, R. 1991. *Game Theory: Analysis of Conflict*. Harvard University Press.
- Nozick, R. 1969. "Newcomb's Problem and Two Principles of Choice." In N. Rescher (ed.), *Essays in Honor of Carl G. Hempel*, 114–46. Reidel.
- Pettigrew, R. 2015. "Risk, Rationality, and Expected Utility Theory." *Canadian Journal of Philosophy* 45:798–826.

- Pollock, J. 2002. "Rational Choice and Action Omnipotence." *Philosophical Review* 111:1–23.
- . 2010. "A Resource-Bounded Agent Addresses the Newcomb Problem." *Synthese* 176:57–82.
- Rabinowicz, W. 1988. "Ratifiability and Stability." In P. Gärdenfors and N. Sahlin (eds.), *Decision, Probability, and Utility*, 406–25. Cambridge University Press.
- . 2002. "Does Practical Deliberation Crowd Out Self-Prediction?" *Erkenntnis* 57:91–122.
- Ramsey, F. 1990 [1926]. "Truth and Probability." In D. H. Mellor (ed.), *Philosophical Papers*. Cambridge University Press.
- Savage, L. 1954. *The Foundations of Statistics*. Wiley Publications in Statistics.
- Skyrms, B. 1984. *Pragmatics and Empiricism*. Yale University Press.
- Sobel, J. H. 1994. *Taking Chances: Essays on Rational Choice*. Cambridge University Press.
- Spencer, J. and Wells, I. MS. "The Metaethical Foundations of Decision Theory."
- Spohn, W. 1977. "Where Luce and Krantz Do Really Generalize Savage's Decision Model." *Erkenntnis* 11:113–34.
- . 2012. "Reversing 30 Years of Discussion: Why Causal Decision Theorists Should One-Box." *Synthese* 187:95–122.
- Stalnaker, R. 1981. "Letter to David Lewis." In Stalnaker R. Harper, W. L. and G. Pearce (eds.), *Ifs: Conditionals, Belief, Decision, Chance and Time*, 151–53. Reidel.
- Weirich, P. 1988. "Hierarchical Maximization of Two Kinds of Expected Utility." *Philosophy of Science* 55:560–82.
- . 2004. *Realistic Decision Theory: Rules for Nonideal Agents in Nonideal Circumstances*. Oxford University Press.
- Wells, I. Forthcoming. "Equal Opportunity and Newcomb's Problem." *Mind* .

Williamson, T. 2000. *Knowledge and its Limits*. Oxford University Press.