

# The Causal Decision Theorist's Guide to Managing the News

*J. Dmitri Gallow* <sup>†</sup>

[Addendum (20/3/2021): if you'd like to apply the theory proposed here to other decisions, I've created [this Jupyter notebook](#) to help.]

Causal decision theory (CDT) says that you should do your best to improve the world in which you find yourself. An act is choiceworthy to the extent that you expect it to promote valuable ends. In contrast, evidential decision theory (EDT) says that you should give yourself good news about the world. An act is choiceworthy to the extent that it indicates valuable ends. Causal decision theorists complain that EDT prescribes an irrational policy of 'managing the news'—favoring acts that provide good news about the world, even when they make the world worse than the alternatives would. CDT's solution is to not manage any news. When it evaluates an act for choiceworthiness, it ignores what the performance of that act would tell you about what the world is like.

Which act you choose to perform can give you news about which goods the world has provided for you, and this kind of news is rightly disregarded when deciding how to act. But which act you choose may *also* give you news about which goods you are in a position to *bring about*, and this kind of news ought not be disregarded. In ignoring all the news that your acts may carry, CDT ignores important correlations between your choice and the goods you are in a position to bring about. For this reason, I've come to think that CDT is in need of revision. According to the revision I favor, EDT does not err in managing the news. Instead, its error lies in managing *the wrong kind* of news—it does not discriminate between good news about the provisions of nature and good news about the extent to which you will *improve upon* the provisions of nature. In rough outline, the revision of CDT I propose counsels you to prefer acts which carry the news that they would make things better, disprefer acts which

Final Draft. Published in [the \*Journal of Philosophy\*, 117 \(3\): 117–149.](#)

<sup>†</sup> Thanks to Martín Abreu Zavaleta, David J. Barnett, R.A. Briggs, Michael Caie, Cian Dorr, Kevin Dorst, Daniel Drucker, Adam Elga, Melissa Fusco, Daniel Hoëk, James Joyce, Jason Konek, Stephen Mackereth, Japa Pallikkathayil, Bernhard Salow, James Shaw, Reuben Stern, and Erica Shumener for helpful conversations and feedback on this material. Thanks also to the University of Pittsburgh's 2019 Summer Work-in-Progress group and two anonymous reviewers.

carry the news that they would make things worse, and ignore any news about how good things are before you act.

To briefly highlight some properties of this theory: in contrast to orthodox CDT, its recommendations do not depend upon how likely you think you are to select any particular act. Unlike orthodox CDT, it will not change its verdicts if we introduce new options which are indistinguishable from existing options in all respects you care about. Additionally, if it says that an act is to be preferred to every other alternative in a pairwise choice between the two, then it will say that that act is to be preferred to every other alternative on the full menu of options. Interestingly, this isn't so for orthodox CDT, which will sometimes say that an act would be permissible, given a pairwise choice between it and every other alternative, but is impermissible to select from the full menu of options.

This theory will sometimes tell you to disprefer a causally ratifiable act (an act you will expect to do the most good possible, if you choose it) to a causally unratable act (an act you will expect to do less good than an alternative would, if you choose it). It will also, in esoteric cases, tell you to most prefer a causally dominated act—an act which would do less good than some alternative, no matter what the world is like. These two properties can seem unattractive when considered in the abstract. Nonetheless, there are particular decisions in which causally unratable acts appear more choiceworthy than their causally ratifiable alternatives; and there are particular decisions in which causally dominated acts appear rational. I'll suggest that it is a virtue of the theory that it agrees with, and helps to explain, the intuitive verdicts about these cases.

## 1 Orthodox Causal Decision Theory

### 1.1 Desirability

Let's assume that, for each epistemically possible world  $w$ , we have a measure of how strongly you desire  $w$  to be actual—we'll write that ' $\mathcal{D}(w)$ '.<sup>1</sup> If we assume the number of worlds to be finite, then with the values  $\mathcal{D}(w)$  and your subjective probability function,  $\text{Pr}$ , defined over all propositions (sets of these worlds), we may define the desirability of a *proposition*,  $\phi$ , as follows:<sup>2</sup>

$$\mathcal{D}(\phi) \stackrel{\text{def}}{=} \sum_w \text{Pr}(w \mid \phi) \cdot \mathcal{D}(w)$$

This definition says: to calculate the desirability of a proposition  $\phi$ , ask yourself how well satisfied you would expect your desires to be, were you to learn only that  $\phi$  is true. Propositions with higher  $\mathcal{D}$ -values give evidence that your

1. I assume that  $\mathcal{D}$  measures desires on an interval scale—it is unique up to positive affine transformation.
2. I abuse notation, writing ' $\text{Pr}(w \mid \phi)$ ' instead of ' $\text{Pr}(\{w\} \mid \phi)$ '.

## §1.2 Newcomb's Problem

desires are satisfied to a greater extent, and propositions with lower  $\mathcal{D}$ -values give evidence that your desires are satisfied to a lesser extent.

Suppose you are choosing whether to perform an act  $A$ , and  $S_1, S_2, \dots, S_N$  are mutually exclusive and jointly exhaustive states of the world. Then, it follows from the definition of  $\mathcal{D}$  that the desirability of your performing  $A$  is given by:<sup>3</sup>

$$\mathcal{D}(A) = \sum_S \Pr(S | A) \cdot \mathcal{D}(SA)$$

By the way, throughout, I'll use non-italic letters like 'A' and 'S' to stand for acts and states, and I'll use italic letters like '*A*' and '*S*' to stand for the proposition that you perform the act  $A$  and that the state  $S$  obtains, respectively.

### 1.2 Newcomb's Problem

Evidential decision theory (EDT) says that this quantity,  $\mathcal{D}(A)$ , measures the choiceworthiness of an act  $A$ . It says that you should prefer  $A$  to  $B$ ,  $A > B$ , iff  $\mathcal{D}(A) > \mathcal{D}(B)$ . In a slogan, it tells you to most prefer the acts that you'd be most glad to learn you had performed. As the evidential decision theorist Richard Jeffrey puts it: "there is no effective difference between asking whether you prefer  $A$  to  $B$  as a news item or as an act, for you make the news".<sup>4</sup> In most ordinary cases, giving yourself the best news coincides with doing your best to improve the world. However, there is a class of decision problems in which giving yourself good news can make matters worse.<sup>5</sup> Consider NEWCOMB'S PROBLEM.<sup>6</sup>

#### NEWCOMB'S PROBLEM

Behind door #1 is either \$1,000,000 or nothing at all. Behind door #2, there is a guaranteed \$10,000. Normally, contestants have to choose between taking a chance on winning a million dollars with door #1 or walking away with the guaranteed \$10,000 behind door #2. But this is the celebrity version of the game, and you are playing for charity, so they've made the game a bit easier: if you want, you are free to open *both* doors and take whatever money you find. Incidentally, before the show was taped, the producers analyzed

3. See §5.4 of Richard Jeffrey, *The Logic of Decision* (New York: McGraw-Hill, 1965).

4. Jeffrey, *ibid.*, p. 84.

5. There is some controversy about whether these are genuine decision problems at all. For an argument that they are not, see Richard Jeffrey, "Causality in the Logic of Decision," *Philosophical Topics*, XXI, 1 (1993): 139–51, and Richard Jeffrey, *Subjective Probability: The Real Thing* (Cambridge: Cambridge University Press, 2004). For a reply, see James M. Joyce, "Are Newcomb Problems Really Decisions?," *Synthese*, CLVI, 3 (2007): 537–62.

6. See Robert Nozick, "Newcomb's Problem and Two Principles of Choice," in Nicholas Rescher, ed., *Essays in Honor of Carl G. Hempel* (Dordrecht: D. Reidel, 1969), pp. 114–46.

your social media accounts with AI bots in an effort to predict how you would behave. If the bots predicted that you would open only door #1, then the producers put \$1,000,000 behind door #1. If, however, they predicted that you would open both doors, then they put nothing behind door #1. The predictions of these bots are 51% reliable.<sup>7</sup> You are told all of this once filming begins.

The decision you face in NEWCOMB'S PROBLEM is an easy one. You get to open both doors and take both prizes. You may find yourself in a good world in which \$1,000,000 awaits behind door #1, or you may find yourself in an unfortunate world in which no money awaits behind door #1. If the world is good, taking both prizes does the most to improve it. If the world is unfortunate, even so, taking both prizes does the most to improve it. So, in either case, taking both prizes does the most good. So you should take both prizes.

EDT, however, advises you to leave a prize behind. The reason is that, if you were to learn that you had taken both prizes, you would expect to walk away with less money; and, if you were to learn that you had left a prize behind, then you would expect to walk away with more. Let 'O' (for one) be the act of opening only door #1, and let 'B' (for both) be the act of opening both doors. Let 'M' (for million) be the proposition that there are a million dollars behind door #1. Then, the desirability of the 'news item' O is<sup>8</sup>

$$\begin{aligned} D(O) &= \Pr(M | O) \cdot D(MO) + \Pr(\neg M | O) \cdot D(\neg MO) \\ &= 51\% \cdot 1,000,000 + 49\% \cdot 0 \\ &= 510,000 \end{aligned}$$

whereas the desirability of the 'news item' B is:

$$\begin{aligned} D(B) &= \Pr(M | B) \cdot D(MB) + \Pr(\neg M | B) \cdot D(\neg MB) \\ &= 49\% \cdot 1,010,000 + 51\% \cdot 10,000 \\ &= 500,000 \end{aligned}$$

So the proposition O is more desirable than the proposition B—you should be glad to learn that you chose O, and sad to learn that you chose B. Nonetheless, B is more choiceworthy than O. Choosing B gains you \$10,000, no matter what. And choosing O loses you \$10,000, no matter what.

So NEWCOMB'S PROBLEM shows us that Jeffrey was incorrect: there *is* an effective difference between asking whether you should prefer A to B as a news

7. That is, the probability that the bots predicted you would open both doors, given that you do, is 51%. And the probability that the bots predicted you would only open door #1, given that you only open door #1, is 51%. This is how I'll continue to understand 'reliable' throughout.

8. I suppose that your desires are linear in dollars.

item, and asking whether A is more choiceworthy than B. Sometimes, choosing the act which will accomplish the most good gives you reason to think that the world is bad. Nonetheless, it is rational to choose the act which will accomplish the most good. So sometimes, you should be distressed to learn that you are choosing rationally. At least, I take these to be the lessons of NEWCOMB'S PROBLEM. There are many others who disagree,<sup>9</sup> but in what follows, I will take these lessons to heart.

### 1.3 Tickers

I say that EDT recommends leaving a prize behind, but not all defenders of EDT agree. Some say that, in my analysis of NEWCOMB'S PROBLEM, I neglected one crucial piece of evidence at your disposal. Prior to deliberation, you will know whether you are leaning towards selecting both doors or whether you are leaning towards selecting only the one—that is, you will have information about your pre-deliberation inclinations.<sup>10</sup> Let's call this additional information a 'tickle.' This tickle will tell you something about how you'll end up choosing, and this will, by itself, tell you something about whether the world is fortunate or unfortunate. It could be that, once you've taken the tickle into account by conditioning your probability function on it, your eventual choice will not give you any *additional* news about whether the world is fortunate or not. Then, you will no longer be sad to learn that you've chosen B, and EDT will recommend taking both prizes. More generally, since EDT says that your pre-deliberation probabilities and desires should determine your choice, these probabilities and desires should completely screen off your choice from the bots' prediction.<sup>11</sup> Since you have introspective access to these initial probabilities and desires, then, your choice will be independent of the bots' prediction, and EDT will recommend taking both prizes.

The causal decision theorist needn't, and shouldn't, deny that the tickles of our initial inclinations, probabilities, or desires can teach us something about the likely consequences of our actions. What they should deny is that tickles like these will *always* do so—for three reasons. Firstly, you may not have introspective access to your own probabilities and desires.<sup>12</sup> Secondly, though

9. See, for instance, Jeffrey, *The Logic of Decision*, *op. cit.*, Arif Ahmed, *Evidence, Decision, and Causality* (Cambridge: Cambridge University Press, 2014), Christopher Meek and Clark Glymour, "Conditioning and Intervening," *The British Journal for the Philosophy of Science*, XLV, 4 (1994): 1001–21, Christopher Hitchcock, "Conditioning, Intervening, and Decision," *Synthese*, CXCIV, 4 (2016): 1157–76, and Reuben Stern, "Interventionist Decision Theory," *Synthese*, CXCIV, 10 (2017): 4133–53.
10. See Ellery Eells, "Causality, Utility, and Decision," *Synthese*, xxxviii, 2 (1981): 295–329, and Ellery Eells, *Rational Decision and Causality* (Cambridge: Cambridge University Press, 1982).
11. To say that the tickle 'screens off' your decision from the bots' prediction is just to say that, conditional on the tickle, your act and the bots' prediction are probabilistically independent.
12. Perhaps a suitably idealized agent *would* have introspective access to their own probabilities

your probabilities and desires may determine your choice, you need not know precisely *how* they will do so. In that case, they need not screen off your choice from the bots' predictions. Thirdly, and relatedly, you need not be following the advice of EDT in order for EDT to *give* that advice. Suppose that you make this decision unreflectively, without bothering to think about it at all. In that case, even though you do not follow EDT, EDT may still be used to evaluate your act as rational or irrational. Because your deliberation provided you with no information to screen off your act from the bots' predictions, EDT will say that leaving a prize behind was the rational choice.<sup>13</sup>

#### 1.4 Utility

When EDT is evaluating the choiceworthiness of the the act  $A$ , it uses a probability function conditioned on  $A$ . But conditioning on  $A$  can provide you with two, importantly different, kinds of information. Firstly, it can provide information about states which are causally downstream of your act. In this way, conditioning on  $A$  provides you with information about the good your act stands to causally promote. Secondly, it can provide information about states which are correlated with, though not causally downstream from, your act. In this way, conditioning on  $A$  provides you with information, not about the good you stand to promote, but rather about the good the world has provided for you.

We may separate the states of the world,  $S$ , into the factors which are not causally downstream of your act—call these 'K'—and the factors which are causally downstream of your act—call those 'C'. Then,

$$\begin{aligned} \mathcal{D}(A) &= \sum_S \Pr(S | A) \cdot \mathcal{D}(SA) \\ &= \sum_K \sum_C \Pr(KC | A) \cdot \mathcal{D}(KCA) \\ &= \sum_K \underbrace{\Pr(K | A)}_{\text{evidential}} \cdot \sum_C \underbrace{\Pr(C | KA)}_{\text{causal}} \cdot \mathcal{D}(KCA) \end{aligned}$$

In the difference it makes to the terms  $\Pr(C | KA)$ , conditioning on the proposition  $A$  provides *causal* information about which states your act *promotes*. In

and desires. Even so, you may not. See David Lewis, "Causal Decision Theory," *Australasian Journal of Philosophy*, LIX, 1 (1981): 5–30.

13. See Frank Jackson and Robert Pargetter, "Where the Tickle Defense Goes Wrong," *Australasian Journal of Philosophy*, LXI, 3 (1983): 295–99. Here, I have presupposed that EDT is at least partly an *evaluative* theory which says whether a given act is rational or irrational for a given agent in a given decision scenario, even if that agent isn't being *guided by* that theory. There are other ways of understanding EDT; and these alternative understandings may escape the NEWCOMB objection.

## §1.5 Utility Can Vary With Your Act Probabilities

the difference it makes to the terms  $\Pr(K | A)$ , conditioning on the proposition  $A$  provides *evidential* information about which states your act merely *indicates*, but does not cause.

For this reason, EDT encourages you to engage in an irrational policy of ‘managing of the news’. It tells you to select only door #1, leaving a prize behind, because this gives you good news about how much money has been provided to you. However, choosing just door #1 does not make the world any better—it is, in fact, guaranteed to make matters worse. The causal decision theorist therefore suggests removing  $A$ ’s evidential influence on  $\mathcal{D}(A)$  by replacing the terms  $\Pr(K | A)$  with  $\Pr(K)$ . Thereby, CDT does not consider the merely evidential value of the act  $A$ . Rather, it considers only its *causal* value. Let’s call the resulting quantity the *utility* of an act,  $\mathcal{U}(A)$ ,<sup>14</sup>

$$\begin{aligned}\mathcal{U}(A) &\stackrel{\text{def}}{=} \sum_K \Pr(K) \cdot \sum_C \Pr(C | KA) \cdot \mathcal{D}(KCA) \\ &= \sum_K \Pr(K) \cdot \mathcal{D}(KA)\end{aligned}$$

While EDT measures the choiceworthiness of an act  $A$  with its ‘news value’,  $\mathcal{D}(A)$ , CDT measures the choiceworthiness of  $A$  with its utility,  $\mathcal{U}(A)$ .

### 1.5 Utility Can Vary With Your Act Probabilities

One important feature of the measure  $\mathcal{U}$  is that its values can depend upon how confident you are that you will end up selecting each available act (call these your *act probabilities*). For instance, in NEWCOMB’S PROBLEM, your probability that there is a million dollars behind door #1 depends upon how likely you think you are to take both prizes. If  $b$  is your probability that you’ll take both prizes, then

$$\begin{aligned}\Pr(M) &= \Pr(M | B) \cdot b + \Pr(M | O) \cdot (1 - b) \\ &= 0.51 - 0.02b \\ \text{and} \quad \Pr(\neg M) &= \Pr(\neg M | B) \cdot b + \Pr(\neg M | O) \cdot (1 - b) \\ &= 0.49 + 0.02b\end{aligned}$$

14. This is Brian Skyrms’s formulation of CDT. See his “Causal Decision Theory,” this JOURNAL, LXXIX, 11 (1982): 695–711. For alternatives, see Allan Gibbard and William Harper, “Counterfactuals and Two Kinds of Expected Utility”, in A. Hooker, J. J. Leach, and E. F. McClennan, eds., *Foundations and Applications of Decision Theory* (Dordrecht: D. Reidel, 1978): pp. 125–62, David Lewis, “Causal Decision Theory”, *op. cit.*, Jordan Howard Sobel, *Taking Chances: Essays on Rational Choice* (Cambridge: Cambridge University Press, 1994), and James M. Joyce, *The Foundations of Causal Decision Theory* (Cambridge: Cambridge University Press, 1999).

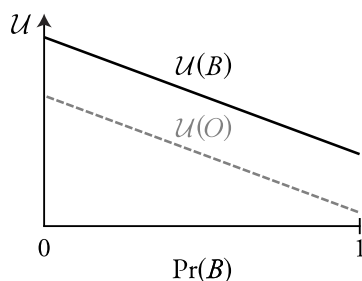


FIGURE 1: The line in solid black shows the utility of  $B$ , as a function of  $\text{Pr}(B)$ . The line in dotted grey shows the utility of  $O$ , as a function of  $\text{Pr}(B)$ .

For this reason, the utility of  $O$  and  $B$  will similarly depend upon how likely you are to take both prizes:

$$\begin{aligned} \mathcal{U}(O) &= \text{Pr}(M) \cdot 1,000,000 + \text{Pr}(\neg M) \cdot 0 \\ &= 510,000 - 20,000b \\ \text{and } \mathcal{U}(B) &= \text{Pr}(M) \cdot 1,010,000 + \text{Pr}(\neg M) \cdot 10,000 \\ &= 520,000 - 20,000b \end{aligned}$$

So as your probability for taking both prizes goes up, the utility of taking both prizes goes down. And as your probability for leaving a prize behind goes up, so too does the utility of leaving a prize behind. (See figure 1.)

Let's write ' $\mathcal{U}_B(A)$ ' for the utility you would assign to  $A$ 's performance, were you to learn that you had performed  $B$ ,

$$\mathcal{U}_B(A) \stackrel{\text{def}}{=} \sum_K \text{Pr}(K | B) \cdot \mathcal{D}(KA)$$

Notice that the desirability of  $A$ ,  $\mathcal{D}(A)$ , is just  $\mathcal{U}_A(A)$ . So EDT says that you should 'manage the news' about utility by preferring acts which give better news about their own utility. For instance, in *NEWCOMB'S PROBLEM*, it compares  $O$  with  $B$  by comparing the value of  $\mathcal{U}(O)$  on the left-hand-side of the graph from figure 1 (where  $\text{Pr}(O) = 1$ ) with the value of  $\mathcal{U}(B)$  on the right-hand-side of the graph (where  $\text{Pr}(B) = 1$ ). The value of  $\mathcal{U}(O)$  on the left-hand-side of the graph is just  $\mathcal{U}_O(O) = \mathcal{D}(O)$ , while the value of  $\mathcal{U}(B)$  on the right-hand-side of the graph is  $\mathcal{U}_B(B) = \mathcal{D}(B)$ . Since  $\mathcal{U}_O(O)$  is greater than  $\mathcal{U}_B(B)$ , EDT says that you should prefer  $O$  to  $B$ . In contrast, since  $\mathcal{U}(B) > \mathcal{U}(O)$ , no matter your probability for  $B$ , CDT says you should prefer  $B$  to  $O$ .

I agree with orthodox CDT that EDT's policy for managing the news is irrational. However, I've come to think that there are some choices in which, by ignoring the news your act will provide you, CDT ends up giving bad advice. I'll contend that, in these kinds of cases, you should manage the news—not by



## §2. Managing the Improvement News with Two Acts

$\mathcal{D}(\text{Row Col})$	$K_A$	$K_D$	$\Pr(\text{Row}   \text{Col})$	$A$	$D$
$A$	$\left[ \begin{array}{cc} 0 & 10 \\ 10 & 0 \end{array} \right]$		$K_A$	$\left[ \begin{array}{cc} 60\% & 10\% \\ 40\% & 90\% \end{array} \right]$	
$D$			$K_D$		

TABLE 1: Desirabilities and Probabilities for DEATH IN DAMASCUS. The matrix on the left-hand-side shows the desirability of the row act in the column state. The matrix on the right shows the probability of the row state, given that you select the column act.

giving yourself good news about how the world is overall, but rather by giving yourself good news about the degree to which your choice will make the world better. I will begin, in §2, with choices between two acts. As I’ll explain in §3, choices between three or more acts present special difficulties.

## 2 Managing the Improvement News with Two Acts

### 2.1 Death in Damascus

Consider the following choice.

#### DEATH IN DAMASCUS

You must choose to travel to either Aleppo or Damascus. Death has no way of learning where you go; but he has made a prediction, and he’s currently waiting in the city he predicted. If you go to the city where Death awaits, you will die; if you go the other city, you will live. Death is good at making these predictions, but he’s not perfect. Moreover, he has a tendency to guess Damascus. The probability that Death is in Damascus, given that you go to Damascus, is 90%. Whereas the probability that Death is in Aleppo, given that you go to Aleppo, is only 60%.

Whether you live or die is the only factor relevant to your decision, and you prefer living to dying. If ‘D’ is the act of going to Damascus, ‘A’ the act of going to Aleppo, and we use ‘ $K_D$ ’ and ‘ $K_A$ ’ for Death’s being in Damascus and Aleppo, respectively, then the relevant desirabilities and probabilities are shown in table 1.<sup>15</sup>

In this choice, as in NEWCOMB’S PROBLEM, the utility of the available acts depends upon your act probabilities. However, in NEWCOMB’S PROBLEM, which act has the highest utility does not vary with your act probabilities—the utility of taking both prizes exceeds the utility of taking only one, no matter how confident you are that you’ll end up taking both. In DEATH IN DAMASCUS, on the

15. The case is a modification of one from Gibbard and Harper, “Counterfactuals and Two Kinds of Expected Utility”, *op. cit.*. Similar decisions are discussed in Andy Egan, “Some Counterexamples to Causal Decision Theory”, *Philosophical Review*, CXVI, 1 (2007): 93–114.

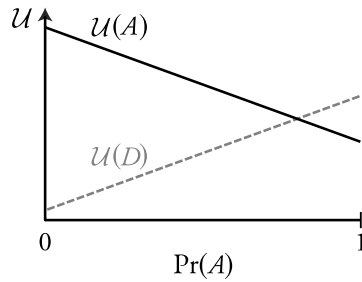


FIGURE 2: DEATH IN DAMASCUS. In solid black, the utility of A as a function of  $\text{Pr}(A)$ . In dashed grey, the utility of D, as a function of  $\text{Pr}(A)$ .

other hand, which act maximizes utility *does* vary with your act probabilities. Let  $a$  be the probability that you go to Aleppo (so that  $1 - a$  is the probability that you go to Damascus). Then,<sup>16</sup>

$$\begin{aligned} \mathcal{U}(A) &= 9 - 5a \\ \mathcal{U}(D) &= 5a + 1 \end{aligned}$$

Therefore, if  $a$  is greater than  $4/5$ ths,  $\mathcal{U}(D) > \mathcal{U}(A)$ . And, if  $a$  is less than  $4/5$ ths,  $\mathcal{U}(A) > \mathcal{U}(D)$ . When  $a$  is exactly  $4/5$ ths,  $\mathcal{U}(A) = \mathcal{U}(D)$ . (See figure 2.)

On the assumption that you've gone to Aleppo, going to Damascus has a higher utility. And on the assumption that you've gone to Damascus, going to Aleppo has a higher utility. So, as soon as you start to follow CDT's advice, it will issue new recommendations. As soon as you find yourself going to Aleppo, and you therefore raise your probability in the proposition A, CDT will tell you to go to Damascus instead. And as soon as you find yourself going to Damascus, and therefore raise your probability in the proposition D, CDT will tell you to go to Aleppo. In these kinds of cases, the verdicts of CDT are unstable.

What I find most disturbing about this instability is that it is entirely predictable. Once it advises you to go to Aleppo and you take its advice, CDT treats the information that you've decided to go to Aleppo, and that Death likely awaits there, as a reason to reconsider its initial recommendation. This treats the information as though it were *surprising*. But there's nothing surprising about this information. You were in a position to know that, by going to Aleppo, you would give yourself evidence that Death is likely in Aleppo. Information you are in a position to know you will gain, if you decide to A, is information which can be taken into account before deciding whether to A. CDT's

16. For those who wish to check the math, some advice: multiply the matrix  $\mathcal{D}(\text{Row}|\text{Col})$  (of the desirability of the row act in the column state) by the matrix  $\text{Pr}(\text{Row}|\text{Col})$  (of the probability of the row state, given the column act), to get the matrix  $\mathcal{U}_{\text{Col}}(\text{Row})$  (of the utility of the row act, given that you've selected the column act). Then: you may use the identity  $\mathcal{U}(\text{Row}) = \sum_{\text{Col}} \mathcal{U}_{\text{Col}}(\text{Row}) \cdot \text{Pr}(\text{Col})$  to derive the unconditional utilities.

## §2.2 *Cake in Damascus*

instability in cases like DEATH IN DAMASCUS is a sign that it doesn't take all of the relevant information into account before issuing its recommendations, even when it is in a position to do so. Since a theory of rational decision should take into account all the relevant information it can before issuing its recommendations, this is a reason to worry about CDT.<sup>17</sup>

### 2.2 *Cake in Damascus*

Next, consider:

#### CAKE IN DAMASCUS

You must choose to travel to either Aleppo or Damascus. You know that your fairy godmother has left you cake in one of these cities. The cake awaits wherever she predicted you would go. She is good at making these predictions, but she's not perfect. Moreover, she has a tendency to guess Damascus. The probability that cake is in Damascus, given that you go to Damascus, is 90%. Whereas the probability that cake is in Aleppo, given that you go to Aleppo, is only 60%.

Whether you have cake or not is the only factor relevant to your decision, and you'd rather have cake. Again, use 'A' for the act of going to Aleppo, 'D' for the act of going to Damascus, and ' $K_A$ ' and ' $K_D$ ' for the cake being in Aleppo and Damascus, respectively. Then, your desirabilities and probabilities are as shown in table 2.<sup>18</sup>

Again, let  $a$  be the probability that you go to Aleppo. Then,

$$U(A) = 5a + 1$$

$$U(D) = 9 - 5a$$

17. For independent arguments against this feature of CDT (which I do not endorse), see Caspar Hare and Brian Hedden, "Self-Reinforcing and Self-Frustrating Decisions," *Noûs*, 1, 3 (2016): 604–28. There are *deliberational* versions of CDT which solve this issue with instability by advising you to get yourself into a position in which you have probability 4/5ths that you will go to Aleppo and probability 1/5th that you will go to Damascus, and, from this deliberative perspective, going to Aleppo and Damascus have equal utility. From this deliberative standpoint, you should either perform the mixed act of going to Aleppo with probability 4/5ths and going to Damascus with probability 1/5th, or else, perhaps, just *pick* either Aleppo or Damascus. See Brian Skyrms, *The Dynamics of Rational Deliberation* (Cambridge: Harvard University Press, 1990), Frank Arntzenius, "No Regrets, or: Edith Piaf Revamps Decision Theory," *Erkenntnis*, LXVIII, 2 (2008): 277–97, James M. Joyce, "Regret and Instability in Causal Decision Theory," *Synthese*, CLXXXVII, 1 (2012): 123–45, and James M. Joyce, "Deliberation and Stability in Newcomb Problems and Pseudo-Newcomb Problems," in Arif Ahmed, ed., *Newcomb's Problem* (Oxford: Oxford University Press, 2018).

18. Similar decisions are discussed in Daniel Hunter and Reed Richter, "Counterfactuals and Newcomb's Paradox," *Synthese*, XXXIX, 2 (1978): 249–61.

The Causal Decision Theorist’s Guide to Managing the News

$\mathcal{D}(\text{Row Col})$	$K_A$	$K_D$	$\text{Pr}(\text{Row}   \text{Col})$	A	D
A	10	0	$K_A$	60%	10%
D	0	10	$K_D$	40%	90%

TABLE 2: Desirabilities and Probabilities for CAKE IN DAMASCUS.

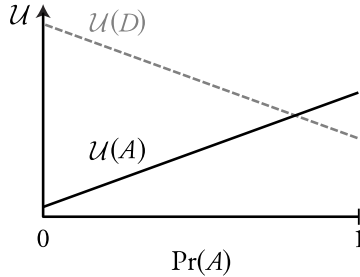


FIGURE 3: CAKE IN DAMASCUS. In solid black, the utility of A, as a function of  $\text{Pr}(A)$ . In dashed grey, the utility of D, as a function of  $\text{Pr}(A)$ .

Therefore, if  $a$  is greater than  $4/5$ ths,  $\mathcal{U}(A) > \mathcal{U}(D)$ . And, if  $a$  is less than  $4/5$ ths,  $\mathcal{U}(D) > \mathcal{U}(A)$ . When  $a$  is exactly  $4/5$ ths,  $\mathcal{U}(A) = \mathcal{U}(D)$ . (See figure 3.)

Conditional on you going to Aleppo, going to Aleppo will have a higher utility than going to Damascus. And, conditional on you going to Damascus, going to Damascus will have a higher utility than going to Aleppo. Which act CDT says is permissible depends upon how confident you are, at the beginning of deliberation, that you will end up selecting that act at deliberation’s end. If you begin deliberation thinking that you’re very likely to go to Aleppo, then CDT recommends that you go to Aleppo. On the other hand, if you begin deliberation thinking that you’re just as likely to go to Aleppo as Damascus, then CDT recommends going to Damascus.

Perhaps this is as it should be. Perhaps your initial probability that you’ll go to Aleppo gives you information about what your fairy godmother predicted. And information about what was predicted is information about where the cake awaits. If you start out thinking that you’re likely to go to Aleppo, then it’s likely that you were predicted to go to Aleppo, and so it’s likely that cake awaits in Aleppo. Perhaps. For perhaps your initial probabilities are ‘tickles’ which provide information about which prediction was made (*cf.* §1.3).

We needn’t, and shouldn’t, deny that the tickles of our initial probabilities can teach us something about the likely consequences of our actions. But, just as in our discussion of NEWCOMB’S PROBLEM, we should deny that they will *always* do so—and for the same reasons. Firstly, you may not have introspective access to your own probabilities.<sup>19</sup> Secondly, you need not be *following* the

19. Perhaps a suitably idealized agent *would* have introspective access to their own probabilities.

### §2.3 Expected Improvement

advice of CDT in order for CDT to *give* that advice. Suppose that you follow EDT. In that case, you will definitely go to Damascus, since  $U_D(D) > U_A(A)$ . If you follow EDT, then your *initial* act probabilities are not relevant to how you end up choosing; for they correlate not at all with the choice you eventually make. However, if your initial act probability for *A* was greater than 4/5ths, then CDT will say that you've chosen irrationally. It will say that you should have seen your high initial act probability for *A* as a reason to go to Aleppo. It will say this, even though your *initial* act probability for *A* doesn't provide any information about what was predicted, or where the cake awaits. (Had you been a causal decision theorist, it may have; but, since you're an evidential decision theorist, it doesn't.)

Suppose you begin deliberation more than 80% confident that you will go to Aleppo. If you then decide to go to Damascus, CDT will call your decision irrational. It will say that you're rationally obligated to go to Aleppo. However, once you learn that you've so decided, your probability for *A* will drop below 80%, and CDT will change its mind, saying that it is (now) rationally obligatory for you to go to Damascus. The advice of CDT therefore conflicts with the following plausible principle governing rational choice: if it's now in your power to do *A*, *A* is permissible, and after you do *A*, it will be permissible for you choose *B*, then, it is now permissible for you to do *A* and then choose *B*. To see the conflict, note that it is now in your power to raise your probability that you'll go to Damascus. Moreover, it is rationally permissible for you to do so. The only rational norms which govern your probabilities are *epistemic* norms. But you are now in a position to manufacture the evidence that you will go to Damascus by simply forming the intention to do so; these intentions will rationalize a higher probability for *D*.<sup>20</sup> So it is permissible to raise your probability that you will go to Damascus. Once you do so, CDT will say that it is rational for you to go to Damascus. So our principle says that it is permissible for you to raise your probability that you will go to Damascus and, then, to go to Damascus. But, so long as your probability for *A* remains above 80%, CDT will disagree. So CDT must deny the principle.

### 2.3 Expected Improvement

As I said above, I believe that what leads EDT into error is not that it manages the news; rather, it is that it manages *the wrong kind of news*. It tells you

Even so, you may not. See Lewis, "Causal Decision Theory", *op. cit.*

20. As arch-causal decision theorist James M. Joyce puts it: "...a rational agent, *while in the midst of her deliberations*, is in a position to legitimately ignore any evidence she might possess about what she is likely to do. She can readjust her probabilities for her currently available acts at will...A deliberating agent who regards herself as free need not proportion her beliefs about her own acts to the antecedent evidence that she has for thinking that she will perform them." (Joyce, "Are Newcomb Problems Really Decisions?", *op. cit.*, at p. 557, italics in original.)

to give yourself good news about how the world is overall—including factors which are outside of your control. In NEWCOMB'S PROBLEM, it says to give yourself propitious news about the provisions of nature. And this is irrational. The provisions of nature are outside of your control; they do not speak in favor of acting one way or another. In their attempt to avoid EDT's irrational managing of the news, causal decision theorists have said that you should ignore any news the performance of an act will carry with it. But there are two, importantly different, kinds of news that the performance of an act can give you. In the first place, the performance of an act can give you evidence about the provisions of nature. This kind of evidence is rightly disregarded. But so too can the performance of an act give evidence about the degree to which it will *improve upon* the provisions of nature. And this kind of evidence is not rightly disregarded. Causalists should be opposed to managing the news about *utility*. For utility encodes your (current) estimation of how desirable the world is as a whole. It conflates goods the world has provided with goods you are in a position to bring about. Causalists should instead manage the news about what your choice will do to make things better or worse.

If you face a choice between two acts, A and B, and you are in the state K, then  $\mathcal{D}(AK) - \mathcal{D}(BK)$  tells us how much more A would improve the world than B would. You don't know which state you're in, but the expectation

$$\begin{aligned} \mathcal{I}(A, B) &\stackrel{\text{def}}{=} \sum_K \Pr(K)[\mathcal{D}(AK) - \mathcal{D}(BK)] \\ &= \mathcal{U}(A) - \mathcal{U}(B) \end{aligned}$$

says how much more than B you *expect* A to do to make the world better. (Notice that  $\mathcal{I}(B, A) = -\mathcal{I}(A, B)$ .) In a choice between two acts, A and B,  $\mathcal{I}(A, B)$  will be greater than  $\mathcal{I}(B, A)$  if and only if  $\mathcal{U}(A)$  is greater than  $\mathcal{U}(B)$ .<sup>21</sup> So, at least in a choice between two acts, saying that you should prefer A to B iff A does more to improve things than B would is equivalent to saying that you should prefer A to B iff A has a higher utility than B does.

#### 2.4 Improvement News

The reason for reformulating causal decision theory in this way is that it allows us to manage the news about what your acts will do to *improve* the world while ignoring news about the provisions of nature. In choosing A, you may give yourself news about the quantity  $\mathcal{I}(A, B)$ . Conditional on your choosing A, the degree to which you'll expect A to do more to improve things than B would is

21. To see this, note that  $\mathcal{I}(A, B) = \mathcal{U}(A) - \mathcal{U}(B)$  and  $\mathcal{I}(B, A) = \mathcal{U}(B) - \mathcal{U}(A)$ . So  $\mathcal{I}(A, B) > \mathcal{I}(B, A)$  iff  $\mathcal{U}(A) - \mathcal{U}(B) > \mathcal{U}(B) - \mathcal{U}(A)$  iff  $\mathcal{U}(A) > \mathcal{U}(B)$ .

## §2.4 Improvement News

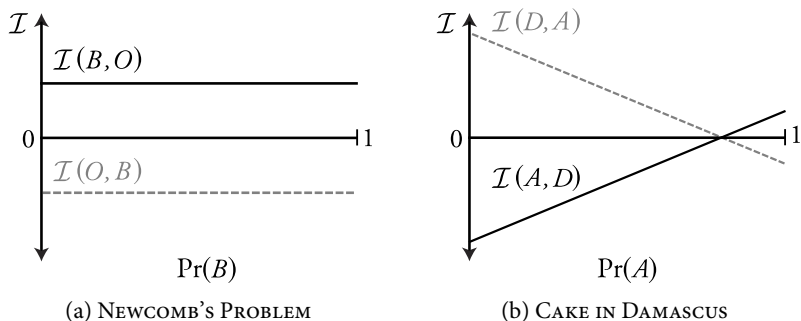


FIGURE 4: In figure 4a, the solid black line is the degree to which B is expected to do more to improve the world than O, as a function of  $\Pr(B)$ . The dashed grey line is the degree to which O is expected to improve the world more than B, as a function of  $\Pr(B)$ . In figure 4b, the solid black line is the degree to which A is expected to improve the world more than D, as a function of  $\Pr(A)$ . And the dashed grey line is the degree to which D is expected to improve the world more than A, as a function of  $\Pr(A)$ .

given by:

$$\begin{aligned} \mathcal{I}_A(A, B) &\stackrel{\text{def}}{=} \sum_K \Pr(K | A) [\mathcal{D}(AK) - \mathcal{D}(BK)] \\ &= \mathcal{U}_A(A) - \mathcal{U}_A(B) \end{aligned}$$

In a choice between two acts, A and B, EDT manages the *utility* news by comparing  $\mathcal{U}_A(A)$  to  $\mathcal{U}_B(B)$ , and preferring the act which gives the best news about utility. I suggest that the causalist manage the *improvement* news by comparing  $\mathcal{I}_A(A, B)$  to  $\mathcal{I}_B(B, A)$ , and preferring the act which gives the best news about how much more it would do to improve things than the alternative.

Let's not just ask about which of A and B gives better improvement news. Let's additionally ask about *how much* better or worse A's improvement news is than B's, in a pairwise comparison. Call that quantity ' $\mathcal{N}(A, B)$ '.

$$\mathcal{N}(A, B) \stackrel{\text{def}}{=} \mathcal{I}_A(A, B) - \mathcal{I}_B(B, A)$$

Then, if  $\mathcal{N}(A, B) > 0$ , A gives better improvement news than B does, and I say that you have reason to prefer A to B. If  $\mathcal{N}(A, B) < 0$ , then B gives better improvement news than A, and I say you have reason to prefer B to A. If  $\mathcal{N}(A, B) = 0$ , then A and B give equally good improvement news, and I say you have reason to be indifferent between A and B. (Notice that  $\mathcal{N}(B, A) = -\mathcal{N}(A, B)$ .)

In cases like NEWCOMB'S PROBLEM, there is no interesting improvement news to manage. B will get you an additional \$10,000, no matter what. (See figure 4a.) In CAKE IN DAMASCUS and DEATH IN DAMASCUS, in contrast,

there is interesting improvement news to manage. Take **CAKE IN DAMASCUS** first. While both A and D give good news about what they are doing to make the world better than the alternative would, D gives better news than A does.  $\mathcal{I}_A(A, D) = 2$  and  $\mathcal{I}_D(D, A) = 8$ , so  $\mathcal{N}(D, A) = 6$ . So I say that you have reason to prefer D to A. (See figure 4b.) In **DEATH IN DAMASCUS**, both A and D give the bad news that they are doing less to improve things than the alternative would. Even so, the bad news that D gives is worse than the bad news that A gives.  $\mathcal{I}_A(A, D) = -2$  and  $\mathcal{I}_D(D, A) = -8$ , so  $\mathcal{N}(A, D) = 6$ . So I say you have reason to prefer A to D.

### 3 Managing the Improvement News with More than Two Acts

David J. Barnett agrees with me about how to choose between two acts.<sup>22</sup> Moreover, for Barnett, this is the entire story. According to him, you should prefer A to B iff you would prefer A to B, were you forced to choose between the two. Thus: you should prefer A to B iff  $\mathcal{N}(A, B) > 0$ , and you should be indifferent between A and B iff  $\mathcal{N}(A, B) = 0$ . As Barnett recognizes, this theory leads to cyclic preferences in some cases.

#### 3.1 Improvement Cycle

Consider **IMPROVEMENT CYCLE**:

##### IMPROVEMENT CYCLE

Before you are three boxes, labeled 'P', 'Q', and 'R'. You must choose one and only one of the boxes. If it was predicted that you would choose P, then \$50 was left in R and the other boxes were left empty. If it was predicted that you would choose R, then \$100 was left in Q and the other boxes left empty. If it was predicted that you would choose Q, then \$150 was left in P and nothing was left in the other boxes. These predictions are 80% reliable.

Desirabilities and probabilities for **IMPROVEMENT CYCLE** are shown in table 3.<sup>23</sup>

**IMPROVEMENT CYCLE** is so-called because, when compared pairwise, P gives better improvement news than Q does, Q gives better improvement news

22. See David James Barnett, "Graded Ratifiability" (m.s.). If we choose an appropriate 'benchmark', then so too does Ralph Wedgewood—in his "Gandalf's Solution to the Newcomb Problem," *Synthese*, CXC, 14 (2013):2643–75. (See Barnett's "Graded Ratifiability," *op. cit.*, which shows that, in the two-act case, Barnett's account is equivalent to Wedgewood's when the 'benchmark' is given by averaging.)
23. Compare **IMPROVEMENT CYCLE** with the decision from Arif Ahmed's "Push the Button," *Philosophy of Science*, LXXIX, 3 (2012): 386–95, and the 'three crates' case from Hare and Hedden, "Self-Reinforcing and Self-Frustrating Decisions," *op. cit.*. Both of these decisions will lead Barnett's theory to advise you to adopt cyclic preferences.



### §3.1 Improvement Cycle

$D(\text{Row Col})$	$K_P$	$K_Q$	$K_R$	$\text{Pr}(\text{Row} \text{Col})$	$P$	$Q$	$R$
$P$	$\begin{bmatrix} 0 & 150 & 0 \\ 0 & 0 & 100 \\ 50 & 0 & 0 \end{bmatrix}$			$K_P$	$\begin{bmatrix} 80\% & 10\% & 10\% \\ 10\% & 80\% & 10\% \\ 10\% & 10\% & 80\% \end{bmatrix}$		
$Q$				$K_Q$			
$R$				$K_R$			

TABLE 3: Desirabilities and probabilities for IMPROVEMENT CYCLE.

than  $R$  does, and  $R$  gives better improvement news than  $P$  does. If you are given a choice between just  $P$  and  $Q$ , your choice is like *NEWCOMB'S PROBLEM*.<sup>24</sup> No matter how likely you are to choose  $P$  or  $Q$ ,  $P$  will be expected to make things better than  $Q$  would. If you take  $P$ , you'll expect  $P$  to improve the world by 5 more dollars than  $Q$  would,  $\mathcal{I}_P(P, Q) = 5$ , and if you take  $Q$ , you'll expect  $P$  to improve the world by 110 more dollars than  $Q$  would,  $\mathcal{I}_Q(Q, P) = -110$ . So  $\mathcal{N}(P, Q) = 115$ , and Barnett concludes that you should prefer  $P$  to  $Q$ . Similarly, if you are given a choice between just  $Q$  and  $R$ , you will expect that  $Q$  would make things better than  $R$  would, no matter your act probabilities.  $\mathcal{I}_Q(Q, R) = 5$  and  $\mathcal{I}_R(R, Q) = -75$ . So  $\mathcal{N}(Q, R) = 80$ , and Barnett concludes that you should prefer  $Q$  to  $R$ . Finally, if you are given a choice between just  $R$  and  $P$ , then your choice is like *DEATH IN DAMASCUS*. Both options will give you the bad news that you could likely make things better by choosing the other. However,  $P$ 's bad news is worse than  $R$ 's is.  $\mathcal{I}_P(P, R) = -25$  and  $\mathcal{I}_R(R, P) = -10$ . So  $\mathcal{N}(R, P) = 15$ , and Barnett concludes that you should prefer  $R$  to  $P$ . If you follow Barnett's advice, then you'll prefer  $P$  to  $Q$  to  $R$  to  $P$ . So your preferences will be cyclic.

To my mind, cyclic preferences like these are irrational. Strict preference should be transitive and it should be irreflexive, so it should be acyclic.<sup>25,26</sup> Fortunately, managing the improvement news need not lead us into cyclic preferences like these. In *IMPROVEMENT CYCLE*, there are three pairwise compar-

24. When I say that the choice is 'like *NEWCOMB'S PROBLEM*', I simply mean that one act is expected to make things better than the alternative, no matter your act probabilities. I don't mean to say that one option causally dominates the other.
25. From transitivity and irreflexivity, acyclicity follows. For suppose there is a cycle leading from  $A_1$  to itself. By transitivity,  $A_1 > A_1$ , contradicting irreflexivity. So, if  $>$  is transitive and irreflexive, there can be no cycles.
26. In fact, I'm inclined to say something stronger: cyclic preferences aren't just irrational—they're impossible. It's not possible to prefer an option to itself, and preferring  $A$  to  $B$  and  $B$  to  $C$  is incompatible with preferring  $C$  to  $A$ . (To prefer  $A$  to  $B$  and  $B$  to  $C$  requires you to be robustly disposed to trade  $C$  for  $B$ , and to trade  $B$  for  $A$ ; but since trading  $C$  for  $B$  for  $A$  just is a way of trading  $C$  for  $A$ , and you know this, you'll not have a robust disposition to trade  $A$  for  $C$ , so you won't count as preferring  $A$  to  $C$ .) However, you needn't accept this stronger claim to accept everything else I have to say here.

isons, and three corresponding pieces of improvement news to manage:

$$\mathcal{N}(P, Q) = 115 \quad \mathcal{N}(Q, R) = 80 \quad \mathcal{N}(R, P) = 15$$

The first two pieces of improvement news are more noteworthy than the third. They give you strong *pro tanto* reason to prefer P to Q and to prefer Q to R. The final piece of improvement news gives you a weaker *pro tanto* reason to prefer R to P. If you are to avoid cycles, you cannot respond to all of these reasons.<sup>27</sup> But you can respond to the strongest reasons you have by preferring P to Q to R. And this is what I recommend you do.

In general: I recommend you draw up a list of all your reasons, and form preferences which respond to as many of them as you can, while giving priority to the stronger reasons. So, take an arbitrary pair of options, A, B (with  $A \neq B$ ). If  $\mathcal{N}(A, B) > 0$ , then you have *pro tanto* reason to prefer A to B. So include this preference, 'A > B', on your list. If  $\mathcal{N}(A, B) = 0$ , then you have *pro tanto* reason to be indifferent between A and B. So include this indifference, 'A ~ B', on your list. Do this for every pair of distinct options. Now: *order* the preferences on your list by the strength of the reasons speaking in their favor. That is: if  $\mathcal{N}(A, B) > \mathcal{N}(C, D) \geq 0$ , then A > B should come before C > D (or C ~ D) on your list. Start at the top, and work your way down. As soon as you encounter an entry such that it, together with the previous entries, would lead to a cycle, strike it from the list.<sup>28</sup> You have *pro tanto* reason to hold that preference, but that reason has been overridden by the stronger reasons above it. Continue down the list, striking any preference you encounter which, together with the unstruck preferences above it, would lead to a cycle. Once you've reached the bottom of your list, transitively close the strict preferences and indifferences which remain,<sup>29</sup> and extend strict preferences along indifferences—that is, if either  $A \sim B > C$  or  $A > B \sim C$ , then include the strict preference  $A > C$ .<sup>30</sup>

That's almost my entire theory of practical rationality—but for a complication with cycles generated by equally weighty reasons (§3.2). The impatient

27. If you agree with me that cyclic preferences are impossible, then the antecedent will be unnecessary, and we can simply say: you cannot respond to all of these reasons.
28. Terminology: a *cycle* is a sequence of options  $A_1, A_2, \dots, A_N$  such that: (a) for each  $i \in \{1, 2, \dots, N - 1\}$ , either  $A_i > A_{i+1}$  or  $A_i \sim A_{i+1}$ ; (b) either  $A_N > A_1$  or  $A_N \sim A_1$ ; and (c) either  $A_i > A_{i+1}$  for some  $i \in \{1, 2, \dots, N - 1\}$  or else  $A_N > A_1$ .
29. That is: if  $A > B$  and  $B > C$ , then include the strict preference  $A > C$ . And, if  $A \sim B$  and  $B \sim C$ , then include the indifference  $A \sim C$ .
30. This theory of rational preference is formally identical to the voting theory of Tideman, with the strength of improvement news  $\mathcal{N}(A, B)$  swapped out for the size of candidate A's majority over candidate B. This isn't an accident—the theory was inspired by Tideman's. See Thorwald Nicolaus Tideman, "Independence of Clones as a Criterion for Voting Rules," *Social Choice and Welfare*, IV, 3 (1987): 185–206. For more on the connections between voting and decision theory, see R.A. Briggs, "Decision-Theoretic Paradoxes as Voting Paradoxes," *Philosophical Review*, CXIX, 1 (2010): 1–30.

### §3.2 Symmetric Improvement Cycle

reader may feel free to skip ahead to §4.

#### 3.2 Symmetric Improvement Cycle

Above, I imposed only one requirement on your list: that the preferences appearing on it be ordered by the strength of the reasons speaking in their favor. That is, if  $\mathcal{N}(A, B) > \mathcal{N}(C, D) \geq 0$ , then  $A > B$  should come before  $C > D$  (or  $C \sim D$ ). But what if  $\mathcal{N}(A, B) = \mathcal{N}(C, D)$ ? Then, either preference could come first. You face an arbitrary choice of which of these two, equally strong, reasons to respond to first. This arbitrary choice can end up making a difference to the preferences you form at the end of the day.

For instance, consider:

#### SYMMETRIC IMPROVEMENT CYCLE

Everything is as in IMPROVEMENT CYCLE, except: if it was predicted that you would take P, then \$100 was left in R; if it was predicted that you would take R, then \$100 was left in Q; and, if it was predicted that you would take Q, then \$100 was left in P.

If you take P, then you'll expect that R would make things better. If you take Q, then you'll expect that P would make things better. And, if you take R, then you'll expect that Q would make things better.

$$\begin{array}{lll} \mathcal{I}_P(P, Q) = 0 & \mathcal{I}_Q(Q, P) = -70 & \mathcal{I}_R(R, P) = 0 \\ \mathcal{I}_P(P, R) = -70 & \mathcal{I}_Q(Q, R) = 0 & \mathcal{I}_R(R, Q) = -70 \end{array}$$

So you have three, equally good, pieces of improvement news to manage:

$$\mathcal{N}(P, Q) = 70 \quad \mathcal{N}(Q, R) = 70 \quad \mathcal{N}(R, P) = 70$$

So: you have *pro tanto* reason to prefer P to Q, *pro tanto* reason to prefer Q to R, *pro tanto* reason to prefer R to P, and each of these reasons is as strong as the others. If you are to avoid cyclic preferences, you cannot respond to all of these reasons. By the symmetry of the case, it seems clear that you should end up indifferent between P, Q, and R. However, following my advice from §3.1, you could draw up any of the following lists:

- |            |            |            |
|------------|------------|------------|
| 1. $P > Q$ | 1. $P > Q$ | 1. $Q > R$ |
| 2. $Q > R$ | 2. $R > P$ | 2. $P > Q$ |
| 3. $R > P$ | 3. $Q > R$ | 3. $R > P$ |
| 1. $Q > R$ | 1. $R > P$ | 1. $R > P$ |
| 2. $R > P$ | 2. $P > Q$ | 2. $Q > R$ |
| 3. $P > Q$ | 3. $Q > R$ | 3. $P > Q$ |

For each list, I advised you to keep the first two preferences and strike the third. You could then end up with any of the following three preference orderings:

$$P > Q > R$$

$$Q > R > P$$

$$R > P > Q$$

Arbitrary choices about which of three equally strong reasons you choose to respond to last should not end up making a difference to the preferences you hold at the end of the day. So we should not say that any of these preferences are rationally permissible. Instead, we should find some way of undoing the effects of this arbitrary choice.

In cases like SYMMETRIC IMPROVEMENT CYCLE, your reasons for holding some preference are *tied* with your reasons for holding another preference. When you draw up your list, then, you will end up breaking this tie by placing one of these preferences below the other. If, together, these tied preferences would complete a cycle, then your arbitrary choice about which preference to place beneath the other is a choice about which preference in the potential cycle to strike from your list. Thereby, it is a choice about which option in the potential cycle to most prefer and which to least prefer. For example, in SYMMETRIC IMPROVEMENT CYCLE, whenever the preference  $P > Q$  is placed at the bottom of your list,  $Q$  ends up being the most preferred option and  $P$  ends up being the least preferred option.

In general, making a choice to place a preference  $A > B$  *lower* on your list is a way of arbitrarily favoring  $B$  and disfavoring  $A$ . (And, therefore, placing  $A > B$  *higher* is arbitrarily favoring  $A$  and disfavoring  $B$ .) For placing  $A > B$  lower on your list means that it is more likely to be struck. And if  $A > B$  is struck, then it would have completed a cycle, so then there is some sequence of options  $C_1, C_2, \dots, C_N$ , such that the preferences  $B \geq C_1 \geq C_2 \geq \dots \geq C_N \geq A$  (with at least one of these preferences strict) remain unstruck. So, if  $A > B$  is struck, then  $A$  will be (at least weakly) dispreferred to each of  $C_1, C_2, \dots, C_N$ , and  $B$  will be (at least weakly) preferred to each of  $C_1, C_2, \dots, C_N$ .

To undo the effects of these arbitrary tie breaks, I suggest that we consider every arbitrary way of favoring some options over others that would lead you to break these ties in one way or another. And if an option ends up at the top of your preference ordering on at least one of these ways of breaking ties, then I will say that it is permissible to choose that option.<sup>31</sup> More carefully, I'll say that a way of breaking ties is *rationalizable* iff it is determined by some (arbitrary) strict order of your options,  $\triangleright$ , in a way that meets the following

31. An option  $A$  is at the *top* of your preference ordering iff there is no option  $B$  such that  $B > A$ . Note that there could be more than one option at the top of an ordering.

### §3.2 Symmetric Improvement Cycle

three constraints. Firstly, if  $\mathcal{N}(A, C) = \mathcal{N}(B, C) > 0$ , then the preference  $A > C$  comes before  $B > C$ , favoring  $A$  and disfavoring  $B$ , iff  $A \triangleright B$ . Secondly, if  $\mathcal{N}(C, A) = \mathcal{N}(C, B) > 0$ , then  $C > B$  comes before  $C > A$ , favoring  $A$  and disfavoring  $B$ , iff  $A \triangleright B$ . And finally, if  $\mathcal{N}(A, C) = \mathcal{N}(B, C) = \mathcal{N}(A, D) = \mathcal{N}(B, D) = 0$ , then  $A \sim C$  comes before  $C \sim B$  iff  $A \sim D$  comes before  $D \sim B$ .

This third constraint requires some explanation. When it comes to indifference relations like  $A \sim B$ , placing them higher will sometimes favor  $A$  and sometimes favor  $B$ . Due to the symmetry of the indifference relation, we can't say in general how breaking these ties will affect the final preferences. However, if we have a *chain* of indifference running from  $A$  to  $B$ ,  $A \sim C \sim B$ , and placing  $A \sim C$  before  $C \sim B$  favors  $A$  over  $B$  (or *vice versa*) then, if we're given *another* chain of indifference,  $A \sim D \sim B$ , placing  $A \sim D$  before  $D \sim B$  will *also* favor  $A$  over  $B$  (or *vice versa*). For, whatever potential cycles running from  $A$  to  $B$  might have motivated placing  $A \sim C$  above  $C \sim B$  in the first case are *also* potential cycles which will motivate placing  $A \sim D$  above  $D \sim B$  in the second case. So, while we can't say in general how someone who favored  $A$  over  $B$  would be motivated to break ties, we can note a consistency constraint on those choices: for any  $A, B, C$ , and  $D$ : if  $A \sim C$  is placed above  $C \sim B$ , then  $A \sim D$  should also be placed above  $D \sim B$ .

If an option ends up at the top of your preference ordering, given at least one rationalizable way of breaking ties, then let us say that it is permissible to choose that option. More carefully, here is the complete algorithm for managing the improvement news: first, list all of the preferences which you have *pro tanto* reason to hold, ordered by the strength of the reasons you have to hold them. If there are ties, where you have equally strong reasons to hold two or more preferences, then consider all the lists which result from breaking those ties in a rationalizable way.<sup>32</sup> Call these the *initial* lists. On each initial list, if a preference, combined with the preceding unstruck preferences, would lead to a cycle, then strike it from the list. Once this is done, generate a preference ordering from each initial list, in the way described in §3.1. Call each resulting preference ordering *preliminary*. For instance, in SYMMETRIC IMPROVEMENT CYCLE, each of  $P > Q > R$ ,  $Q > R > P$ , and  $R > P > Q$  are preliminary preference orderings. If an option is at the top of some preliminary preference ordering, then say that it is a 'tier 1' option. The tier 1 options are all and only the permissible options. In SYMMETRIC IMPROVEMENT CYCLE, for instance,  $P$ ,

32. If I did not require ties to be broken in a rationalizable way, then my theory of rational preference would not be independent of clones in the sense defined in §4. (Without a rationalizable tie-breaking method, Lemma 1 in the appendix would be false.) The example which illustrates this is too complicated to include here—it involves seven acts—but the interested reader should consult Example 1 in T. M. Zavist and T. N. Tideman, "Complete Independence of Clones in the Ranked Pairs Rule," *Social Choice and Welfare*, xi, 2 (1989): 167–73. (Note: what I am calling a 'rationalizable' tie-breaking method is *not* what Zavist and Tideman call an 'impartial' tie-breaking method.)

Q, and R are all tier 1 options, and they are all permissible.

Given this procedure for selecting permissible options, we may construct a full preference ordering. If an option would be permissible, were all the tier 1 options removed from the menu, then say that it is a 'tier 2' option. If an option would be permissible, were all the tier 1 and tier 2 options removed from the menu, then say that it is a 'tier 3' option. And so on. In general, for  $k > 1$ , if an option would be permissible, were all options of tier  $k - 1$  or lower removed from the menu, then say that it is a 'tier  $k$ ' option. If A is a tier  $n$  option and B is a tier  $k$  option, then A should be preferred to B if  $n < k$ , B should be preferred to A if  $n > k$ , and you should be indifferent between A and B if  $n = k$ . Or, more concisely, A is to be weakly preferred to B iff A's tier is no greater than B's.

#### 4 Further Discussion

This completes my theory of how to rationally manage the improvement news when you're choosing between arbitrarily many options. In this section, I'll apply the theory to some illustrative cases and note some of its properties.

##### 4.1 *The Frustrator*

Consider the following choice:

###### THE FRUSTRATOR

Before you are three boxes, labeled 'A', 'B', and 'C'. You must choose one and only one of the boxes. Yesterday, a reliable predictor, known as *the Frustrator*, made a prediction about how you would choose. If she predicted that you would choose A, then \$100 was left in B and nothing was left in A. If she predicted that you would choose B, then \$100 was left in A and nothing was left in B. If she predicted that you would choose C, then \$40 was left in both A and B. There is a guaranteed \$40 in C, no matter what was predicted. These predictions are 80% reliable.

Desirabilities and probabilities for THE FRUSTRATOR are in table 4.<sup>33</sup>

If you choose A, then you'll expect that C would make things better, and B would make things better still. If you choose B, then you'll expect that C would make things better, and A would make them better still. If you choose C, then you'll expect that either A or B would make things (slightly) better.

$$\mathcal{I}_A(A, B) = -70 \quad \mathcal{I}_B(B, A) = -70 \quad \mathcal{I}_C(C, A) = -2$$

33. Similar decisions are discussed in Arif Ahmed "Dicing with Death," *Analysis*, LXXIV, 4 (2014): 587–92, and Jack Spencer and Ian Wells, "Why Take Both Boxes?," *Philosophy and Phenomenological Research*, XCIX, 1 (2019): 27–48.

### §4.1 The Frustrator

$D(\text{Row } \text{Col})$	$K_A$	$K_B$	$K_C$	$\text{Pr}(\text{Row} \text{Col})$	A	B	C
A	0	100	40	$K_A$	80%	10%	10%
B	100	0	40	$K_B$	10%	80%	10%
C	40	40	40	$K_C$	10%	10%	80%

TABLE 4: Desirabilities and probabilities for THE FRUSTRATOR.

$$\mathcal{I}_A(A, C) = -26 \qquad \mathcal{I}_B(B, C) = -26 \qquad \mathcal{I}_C(C, B) = -2$$

So no matter which option you select, you will be giving yourself the bad news that some other option would likely make things better. In a pairwise comparison, A and B both give equally bad improvement news, so you've no more reason to prefer A to B than you do to prefer B to A. C, however, gives much better news than either A or B when compared pairwise:

$$\mathcal{N}(A, B) = 0 \qquad \mathcal{N}(C, B) = 24 \qquad \mathcal{N}(C, A) = 24$$

So I say that you have *pro tanto* reason to prefer C to A, C to B, and to be indifferent between A and B. There's no obstacle to responding to all of these reasons, so I say that you should have the preference ordering  $C > A \sim B$ .

Orthodox CDT disagrees. It says that you should prefer at least one of A and B to C. Let  $a$  be your probability that A, and let  $b$  be your probability that B. (So that  $1 - a - b$  is your probability that C.) Then,

$$\mathcal{U}(A) = 42 + 42b - 28a$$

$$\mathcal{U}(B) = 42 + 42a - 28b$$

$$\mathcal{U}(C) = 40$$

So, if  $a > b$ , then  $\mathcal{U}(B) > \mathcal{U}(C)$ . If  $b > a$ , then  $\mathcal{U}(A) > \mathcal{U}(C)$ . If  $a = b$ , then  $\mathcal{U}(A) = \mathcal{U}(B) > \mathcal{U}(C)$ . For *no* act probability does C have a higher utility than both A and B.

Suppose that you always begin deliberation thinking you're equally likely to select any of the available options. And suppose that you are given a choice between just A and C—B is taken off the menu (though there's still a 10% probability that the Frustrator falsely predicted that B). In that case, at the beginning of deliberation,  $\mathcal{U}(A) = 42 - (28/2) = 28$ , while  $\mathcal{U}(C) = 40$ . So CDT would tell you to prefer C to A. Suppose, on the other hand, that you are given a choice between just B and C—A is taken off of the menu. In that case, at the beginning of deliberation,  $\mathcal{U}(B) = 28$  and  $\mathcal{U}(C) = 40$ . So CDT would tell you to prefer C to B. C is, in other words, a *Condorcet winner*. In voting theory, a Condorcet winner, named after the Marquis de Condorcet, is a candidate who

would win in a one-on-one contest with any other candidate. Likewise, in THE FRUSTRATOR, were you given a choice between C and any other alternative, and your act probabilities are uniform, then orthodox CDT would advise you to most prefer C. Nonetheless, if you are given a choice between A, B, and C, and your act probabilities are uniform, orthodox CDT will say to disprefer C to both A and B.

In contrast, if you manage the improvement news, then, if you would prefer an option to every alternative in a pairwise choice between the two, then it will always be your most preferred option overall. Similarly: if you manage the improvement news, then, if you would *disprefer* an option to every alternative in a pairwise choice between the two, then it will always be your *least* preferred option overall. To understand why, see the proofs of **Propositions 1** and **2** in the appendix.

#### 4.2 Clone-Independence

Say that an option  $C^*$  is a *clone* of another option C iff  $C^*$  is identical to C in all respects that you care about—that is, for all states of nature K,  $\Pr(K \mid C^*) = \Pr(K \mid C)$  and  $\mathcal{D}(KC^*) = \mathcal{D}(KC)$ . And say that a theory of rational preference is *independent of clones* iff adding or removing a clone doesn't affect your preferences between the other options. This seems to me to be a property that any reasonable theory of rational choice should possess; but notice that, if you always begin deliberation thinking that you're equally likely to select any option, then orthodox CDT will not be independent of clones.

Take, for instance, CAKE IN DAMASCUS, from §2.2. If you begin deliberation thinking that you're just as likely to choose A as D, then orthodox CDT will tell you to prefer D (and it won't at any point reverse this judgment as your act probability for D rises to 100%). But suppose that we introduce 4 additional 'clones' of the Aleppo option—perhaps there are 5 different paths leading to Aleppo, all equally good, and only one path to Damascus. Then, your choice is not between A and D, but instead between  $A_1, A_2, A_3, A_4, A_5$  and D. If you start out thinking that you're equally likely to pick each each path, then orthodox CDT will tell you to prefer each path to Aleppo to the path to Damascus.<sup>34</sup>

In contrast, if you manage the improvement news in the way I've suggested, then you will prefer the path to Damascus to each path to Aleppo (and you'll be indifferent between the different paths to Aleppo). This follows from a more general principle: managing the improvement news is independent of clones,

34. A defender of orthodox CDT may wish to protest that, in order for  $C^*$  to really be a *clone* of C, your initial probability for  $C^*$  must be the same as it was for C. In that case, consider CAKE IN DAMASCUS with an initial probability for A of  $5/11$ ths. Then, orthodox CDT will say to prefer D to A. If a clone for A is then introduced, your probability for D will be  $1/11$ th, in which case, orthodox CDT will say to prefer A to D.



### §4.3 Causal Ratifiability

$\mathcal{D}(\text{Row Column})$	$K_F$	$K_G$	$K_H$
$F$	1	1	1
$G$	0	9	10
$H$	0	10	9

TABLE 5: Desirabilities for THREE SHELLS.

and always counsels indifference between clones. To understand why this is so, see the proofs of **Propositions 3** and **4** in the appendix.<sup>35</sup>

#### 4.3 Causal Ratifiability

Say that an option  $A$  is *causally ratifiable* iff  $\mathcal{I}_A(A, B) \geq 0$ , for all  $B$  (or, equivalently, iff  $\mathcal{U}_A(A) \geq \mathcal{U}_A(B)$ , for all  $B$ ). A causally ratifiable act is one that you will expect to do the most good possible, once you've chosen it. Managing the improvement news will sometimes lead you to least prefer the only causally ratifiable option. Consider the following decision:

##### THREE SHELLS

Before you are three shells, labeled 'F', 'G', and 'H'. You must choose one, and only one, of the shells. There's a guaranteed \$1 under shell F, no matter what. If it was predicted that you would choose F, then nothing was left under shells G and H. If it was predicted that you would choose G, then \$9 was left under G and \$10 was left under H. If it was predicted that you would choose H, then \$10 was left under G and \$9 was left under H.

Desirabilities for this case are shown in table 5. For simplicity, assume that these predictions are 100% reliable.<sup>36</sup>

Note that F is causally ratifiable, since, if you take F, you'll be certain that there's nothing under either G or H. And note that neither G nor H is causally ratifiable, since, if you take G, then you'll be certain that there's more money under H; and, if you take H, you'll be certain that there's more money under G.

$$\mathcal{I}_F(F, G) = 1 \qquad \mathcal{I}_G(G, F) = 8 \qquad \mathcal{I}_H(H, F) = 8$$

35. Managing the improvement news has this property for roughly the same reason that Tideman's voting method (as emended in Zavist and Tideman, "Complete Independence of Clones in the Ranked Pairs Rule", *op. cit.*) is independent of clones—though Tideman's definition of 'clone' is different from mine, and Zavist and Tideman's criterion of 'impartiality' is different from my criterion of 'rationalizability'.

36. The decision comes from Brian Skyrms, *Pragmatics and Empiricism* (New Haven: Yale University Press, 1984).

$$\mathcal{I}_F(F, H) = 1 \quad \mathcal{I}_G(G, H) = -1 \quad \mathcal{I}_H(H, G) = -1$$

So F is the only option which is causally ratifiable. Nonetheless, if you manage the improvement news in the way that I've suggested, then you'll end up preferring both G and H to F. For you'll have three pieces of improvement news to manage:

$$\mathcal{N}(G, F) = 7 \quad \mathcal{N}(G, H) = 0 \quad \mathcal{N}(H, F) = 7$$

Thus, you have reason to prefer G to F, you have reason to prefer H to F, and you've no more reason to prefer G to H than you have reason to prefer H to G. There's no obstacle to responding to all of these reasons, so you should be indifferent between G and H and prefer both to F:  $G \sim H > F$ .

#### 4.4 Causal Dominance

Managing the improvement news will sometimes lead you to violate the principle of CAUSAL DOMINANCE:

CAUSAL DOMINANCE

If D *causally dominates* A—i.e., if  $\mathcal{D}(DK) > \mathcal{D}(AK)$ , for each state of nature K—then D is to be preferred to A.

I expect many causalists to see this as a deficit of a theory of rational choice; however, I urge caution. Note that the theory will always say that you have *pro tanto* reason to prefer a dominating option to the option it dominates. After all, if  $\mathcal{D}(DK) > \mathcal{D}(AK)$ , for every K, then it must be that  $\mathcal{N}(D, A) > 0$ . In most cases, you will be able to respond to this *pro tanto* reason without forming cyclic preferences, and the theory will say that D is to be preferred to A. Nonetheless, there are some esoteric cases in which your *pro tanto* reason to disprefer a causally dominated option can be overridden. Moreover, in these esoteric cases, I believe that managing the improvement news agrees with, and helps to explain, our intuitive verdicts.

For a case like this, consider:

THE SEMI-FRUSTRATOR

Before you are two boxes, a white box and a black box. You may point to either box, and you will be given that box's contents. Yesterday, *the Frustrator* made a prediction about which box you would point to. If she predicted that you would point to the black box, then she put \$100 in the white box and left the black box empty. If she predicted that you would point to the white box, then she put \$100 in the black box and left the white box empty. She is excellent at these predictions—but only when you point to the boxes

§4.4 Causal Dominance

$\mathcal{D}(\text{Row Column})$	$K_B$	$K_W$			
$B$	0	100	]		
$W$	100	0			
$B_\Delta$	$-\Delta$	$100 - \Delta$			
$W_\Delta$	$100 - \Delta$	$-\Delta$			
$\Pr(\text{Row} \text{Column})$	$B$	$W$	$B_\Delta$	$W_\Delta$	
$K_B$	100%	0%	50%	50%	]
$K_W$	0%	100%	50%	50%	

TABLE 6: Desirabilities and probabilities for THE SEMI-FRUSTRATOR.

with your left hand. (Perhaps because she bases the prediction on a brain scan of only your right hemisphere.) Given that you point with your left hand, you are 100% sure that she correctly predicted where you'd point. However, given that you point with your right hand, you're only 50% sure that she predicted correctly. You may point with either hand, but, if you use your right, then you must pay a pittance,  $\Delta$ .

IN THE SEMI-FRUSTRATOR, you have four available acts: point to black with your left hand,  $B$ , point to black with your right hand,  $B_\Delta$ , point to white with your left hand,  $W$ , and point to white with your right hand,  $W_\Delta$ . And there are two states of nature: either the Frustrator predicted black,  $K_B$ , or the Frustrator predicted white,  $K_W$ . Desirabilities and probabilities for the case are as shown in table 6.<sup>37</sup>

If you choose  $B$ , you'll expect that  $W$  and  $W_\Delta$  would make things much better, and that  $B_\Delta$  would make things slightly worse. Likewise, if you choose  $W$ , you'll expect that  $B$  and  $B_\Delta$  would make things much better, and that  $W_\Delta$  would make things slightly worse. If you choose either of  $B_\Delta$  or  $W_\Delta$ , then you won't expect that the other would make things any better, and you'll expect that both of  $B$  and  $W$  would make things slightly better,

$$\begin{array}{lll}
 \mathcal{I}_B(B, W) = -100 & \mathcal{I}_B(B, W_\Delta) = \Delta - 100 & \mathcal{I}_B(B, B_\Delta) = \Delta \\
 \mathcal{I}_W(W, B) = -100 & \mathcal{I}_W(W, B_\Delta) = \Delta - 100 & \mathcal{I}_W(W, W_\Delta) = \Delta \\
 \mathcal{I}_{B_\Delta}(B_\Delta, B) = -\Delta & \mathcal{I}_{B_\Delta}(B_\Delta, W) = -\Delta & \mathcal{I}_{B_\Delta}(B_\Delta, W_\Delta) = 0 \\
 \mathcal{I}_{W_\Delta}(W_\Delta, B) = -\Delta & \mathcal{I}_{W_\Delta}(W_\Delta, W) = -\Delta & \mathcal{I}_{W_\Delta}(W_\Delta, B_\Delta) = 0
 \end{array}$$

So, you have the following pieces of improvement news to manage:

$$\mathcal{N}(B_\Delta, W) = 100 - 2\Delta \quad \mathcal{N}(B, B_\Delta) = 2\Delta \quad \mathcal{N}(B, W) = 0$$

37. This decision comes from Spencer and Wells, "Why Take Both Boxes?," *op. cit.*.

$$\mathcal{N}(W_{\Delta}, B) = 100 - 2\Delta \quad \mathcal{N}(W, W_{\Delta}) = 2\Delta \quad \mathcal{N}(B_{\Delta}, W_{\Delta}) = 0$$

You have strong *pro tanto* reason to prefer  $B_{\Delta}$  to  $W$ , and strong *pro tanto* reason to prefer  $W_{\Delta}$  to  $B$ . And you have a weaker *pro tanto* reason to prefer  $B$  to  $B_{\Delta}$  and  $W$  to  $W_{\Delta}$ . You can't respond to all of these reasons, since they lead to the cycle  $B_{\Delta} > W > W_{\Delta} > B > B_{\Delta}$ . But you can respond to the strongest reasons you have by forming the preferences  $B_{\Delta} > W$  and  $W_{\Delta} > B$ . Following the advice from §3.1, you then face an arbitrary choice about whether to put the preference  $B > B_{\Delta}$  above or below the preference  $W > W_{\Delta}$  on your initial list. The former will lead you to the preference ordering  $W_{\Delta} > B > B_{\Delta} > W$ , while the latter will lead you to the preference ordering  $B_{\Delta} > W > W_{\Delta} > B$ . Undoing the effects of this arbitrary choice in the way outlined in §3.2, you'll end up with the final preferences:  $B_{\Delta} \sim W_{\Delta} > B \sim W$ .

But note that, in THE SEMI-FRUSTRATOR, pointing to either box with your right hand is *causally dominated* by pointing to that box with your left. No matter whether the Frustrator has predicted black or white, pointing to a box with your left hand would get you a pittance more than pointing to that box with your right. Nonetheless, pointing to a box with your right hand does *seem* to be more rational than pointing to that box with your left. Indeed, Spencer and Wells presented THE SEMI-FRUSTRATOR as a counterexample to CAUSAL DOMINANCE. We should be cautious about our first-pass judgments in these kinds of cases, but, on reflection, I'm inclined to endorse the judgment, and regard this as a problem for the CAUSAL DOMINANCE principle, and not a problem for managing the improvement news.

If you manage the improvement news, then you'll think that what's true in the CAUSAL DOMINANCE principle is that you will always have a distinctive and compelling reason to disprefer a causally dominated option: you are certain that, no matter what the world is like, some other option would do more good than it will. However, in cases like THE SEMI-FRUSTRATOR, something similar is true of the causally *dominating* option: given that you choose it, you'll be certain that some option would do more good than it. For instance: you have a compelling reason to disprefer  $B_{\Delta}$  to  $B$ —namely,  $B$  certainly would do a pittance more good than  $B_{\Delta}$ . But something similar is true of  $B$ : if you choose  $B$ , then you'll be certain that  $W$  would do more good than *it*.

If you manage the improvement news, you'll think that THE SEMI-FRUSTRATOR is analogous to the following counterexample to CAUSAL DOMINANCE: God asks you to name any natural number,  $n$ , whereupon he will provide you with  $n$  days in heaven. In this decision, *every* option is causally dominated—nonetheless, it does not appear that no option is permissible. One common thing to say about this case is that selecting a causally dominated option is usually irrational because it means passing up a more preferable option. However, when

you have an infinite hierarchy of ever-more-preferable options, then, for every option, there is another option which is preferable to it. In that case, you've no choice *but* to pass up a more preferable option. If something speaks against *every* option, then it speaks against none of them. Everyone's sin is no one's sin.

If you manage the improvement news, then you should say something analogous about THE SEMI-FRUSTRATOR: for every option, there is an alternative which you have *pro tanto* reason to prefer to it. (This can happen with a finite number of options because your *pro tanto* reasons for preference can form cycles.) The fact that  $B_{\Delta}$  is causally dominated by B means that, if you choose  $B_{\Delta}$ , then you're passing up a *pro tanto* more preferable option; however, in this decision, you've no choice *but* to pass up a *pro tanto* more preferable option. Everyone's sin is no one's sin; so, in cases like this, being causally dominated need not speak against an option.

## A Technicalities

To review, this is the algorithm for managing the improvement news ('MIN'): first, list the preferences you have *pro tanto* reason to hold, ordered by the strength of the reasons you have to hold those preferences. So: for all A, B, C, D, if  $\mathcal{N}(A, B) > \mathcal{N}(C, D) \geq 0$ , then  $A > B$  comes before  $C > D$  (or  $C \sim D$ ). If there are ties, where you have equally strong reasons to hold multiple preferences, then consider all the lists which result from breaking ties in a rationalizable way. (A tie-breaking procedure is rationalizable iff there is some strict total ordering of the options,  $\triangleright$ , such that, for all A, B, C, and D: (a) if  $\mathcal{N}(A, C) = \mathcal{N}(B, C) > 0$ , then  $A > C$  comes before  $B > C$  iff  $A \triangleright B$ , (b) if  $\mathcal{N}(C, A) = \mathcal{N}(C, B) > 0$ , then  $C > B$  comes before  $C > A$  iff  $A \triangleright B$ , and (c) if  $\mathcal{N}(A, C) = \mathcal{N}(B, C) = \mathcal{N}(A, D) = \mathcal{N}(B, D) = 0$ , then  $A \sim C$  comes before  $C \sim B$  iff  $A \sim D$  comes before  $D \sim B$ .) Call these the *initial* lists. For each initial list, if a preference, combined with the preceding unstruck preferences, would lead to a cycle, then strike it from the list. With the remaining preferences, transitively close strict preferences and indifferences and extend preference along indifference. Call each resulting preference ordering *preliminary*. If an option is at the top of *some* preliminary preference ordering, then it is a permissible option, and we say that it is a 'tier 1' option. For  $k > 1$ , if an option would be permissible were all options of tier  $k - 1$  or lower eliminated, then say that it is a 'tier  $k$ ' option. A is to be weakly preferred to B iff A's tier is no greater than B's.

**Definition 1.**  $\top$  is a Condorcet winner iff, for every alternative  $A \neq \top$ ,  $\mathcal{N}(\top, A) > 0$ .

**Proposition 1.** *If  $\top$  is a Condorcet winner, then MIN says that  $\top$  is to be preferred to every other alternative.*

*Proof.* If  $\top$  a Condorcet winner, then all initial lists will include  $\top > A$ , for every alternative  $A$ . No other preferences will lead you to strike  $\top > A$  from the list. For take any two alternatives,  $A, B$ . Either  $A > B$ ,  $B > A$ , or  $A \sim B$ . None of these possibilities lead to a cycle when combined with the preferences  $\top > A$  and  $\top > B$ . So every preliminary preference ordering will have  $\top$  as the unique top element. So  $\top$  will be the only tier 1 option.  $\square$

**Definition 2.**  $\perp$  is a Condorcet loser iff, for every alternative  $A \neq \perp$ ,  $\mathcal{N}(A, \perp) > 0$ .

**Proposition 2.** *If  $\perp$  is a Condorcet loser, then MIN says that  $\perp$  is to be dispreferred to every alternative.*

*Proof.* If  $\perp$  is a Condorcet loser, then all initial lists will include  $A > \perp$ , for every alternative  $A$ . No other preference will lead you to strike  $A > \perp$  from the list. For take any two alternatives,  $A, B$ . Either  $A > B$ ,  $B > A$ , or  $A \sim B$ . None of these possibilities lead to a cycle when combined with the preferences  $A > \perp$  and  $B > \perp$ . So every preliminary preference ordering on any menu will have  $\perp$  as the unique bottom element. So every alternative will have a lower tier than  $\perp$ . So  $\perp$  will be dispreferred to every alternative.  $\square$

Some notation: for distinct options  $X$  and  $Y$ , use ' $[X, Y]$ ' for the preference  $X > Y$  if  $\mathcal{N}(X, Y) > 0$ ,  $Y > X$  if  $\mathcal{N}(X, Y) < 0$ , and  $X \sim Y$  if  $\mathcal{N}(X, Y) = 0$ . And use ' $\mathcal{N}[X, Y]$ ' for  $\max\{\mathcal{N}(X, Y), \mathcal{N}(Y, X)\}$ .

**Definition 3.**  $C$  and  $C^*$  are clones iff, for each state of nature  $K$ ,  $\Pr(K | C) = \Pr(K | C^*)$  and  $\mathcal{D}(KC) = \mathcal{D}(KC^*)$ .

It follows immediately from this definition that, if  $C$  and  $C^*$  are clones, then, for any option  $X \neq C, C^*$ ,  $\mathcal{N}[X, C] = \mathcal{N}[X, C^*]$ .

**Lemma 1.** *If  $C$  and  $C^*$  are clones, then, for all options  $A \neq C, C^*$ ,  $[A, C]$  will be struck from an initial list iff  $[A, C^*]$  is also struck.*

*Proof.* Define a relation of  $C$ -precedence over the options besides  $C$  and  $C^*$  as follows: if  $[A, C]$  comes before  $[B, C]$  on the initial list, then  $A$   $C$ -precedes  $B$ . And similarly define a relation of  $C^*$ -precedence: if  $[A, C^*]$  comes before  $[B, C^*]$  on the initial list, then  $A$   $C^*$ -precedes  $B$ . If ties are broken in a rationalizable way, it follows that  $A$   $C$ -precedes  $B$  iff  $A$   $C^*$ -precedes  $B$ ,<sup>38</sup> and we can simply

38. Suppose  $A$   $C$ -precedes  $B$ . Then, either  $\mathcal{N}[A, C] > \mathcal{N}[B, C]$  or  $\mathcal{N}[A, C] = \mathcal{N}[B, C]$ . In the first case,  $\mathcal{N}[A, C^*] > \mathcal{N}[B, C^*]$ , so  $A$   $C^*$ -precedes  $B$ . If  $\mathcal{N}[A, C] = \mathcal{N}[B, C]$ , then: either (i)  $\mathcal{N}[A, C] > 0$  or (ii)  $\mathcal{N}[A, C] = 0$ . In case (i), either (i)  $[A, C] = A > C$  or (ii)  $[A, C] = C > A$ .

say that  $A$  precedes  $B$ . Precedence provides a strict total ordering of the options besides  $C$  and  $C^*$ , which we may then use to induce an enumeration of those options,  $A_1, A_2, \dots, A_N$ , where  $A_i$  precedes  $A_{i+1}$ , for each  $i \geq 1$ .

Suppose now that there is some option  $A$  such that either  $[A, C]$  is struck and  $[A, C^*]$  is not, or else  $[A, C^*]$  is struck and  $[A, C]$  is not. Let  $A_f$  be the first such option, given our enumeration, and wLOG, suppose that  $[A_f, C]$  is struck and  $[A_f, C^*]$  is not. Since the preference between  $A_f$  and  $C$  is struck, it must form a cycle with some earlier unstruck preferences. At least one of these earlier preferences must involve  $C$ . Choose one and call it  $[A_e, C]$ . Then, we have a potential cycle of unstruck preferences: either  $C \geq A_e \geq \dots \geq A_f$  or  $C \leq A_e \leq \dots \leq A_f$  (where at least one of these preferences is strict). Since  $[A_e, C]$  is above  $[A_f, C]$  on the initial list,  $A_e$  precedes  $A_f$ . So the preference  $[A_e, C^*]$  cannot be struck (since  $[A_e, C]$  is unstruck, and by hypothesis  $A_f$  is the first option  $A$  in our enumeration such that exactly one of  $[A, C]$  and  $[A, C^*]$  is struck). So, when we arrive at the preference  $[A_f, C^*]$ , we will have a potential cycle: either  $A_f \geq C^* \geq A_e \geq \dots \geq A_f$  or  $C^* \leq A_e \leq \dots \leq A_f$  (with at least one of these preferences strict). So we will strike  $[A_f, C^*]$ . This contradicts our assumption that there is an option  $A$  such that one of  $[A, C]$  and  $[A, C^*]$  was struck and the other was not. So there can be no such option.  $\square$

**Proposition 3.** *If  $C$  and  $C^*$  are clones, then, on every preliminary preference ordering,  $C \sim C^*$ .*

*Proof.* Suppose, for *reductio*, that  $C \not\sim C^*$ . Since  $C$  and  $C^*$  are clones,  $\mathcal{N}(C, C^*) = 0$ , and the initial list includes  $C \sim C^*$ . Therefore, if  $C > C^*$  in the preliminary ordering, there must be some options  $A_1, \dots, A_N$  such that the sequence of preferences  $C \geq A_1 \geq \dots \geq A_N \geq C^*$  (with at least one of these preferences strict) all come before  $C \sim C^*$  on the list, and all remain unstruck. Since  $C$  and  $C^*$  are clones,  $C^* \geq A_1$  must also appear on the list. Since this is inconsistent with the preliminary preference ordering, the preference  $C^* \geq A_1$  must be struck. So  $C^* \geq A_1$  is struck and  $C \geq A_1$  is not. But this is impossible, by **Lemma 1**. So  $C \sim C^*$ .  $\square$

Notation: given a menu of option  $\mathbf{M}$ , let  $\mathcal{P}(\mathbf{M})$  be the set of permissible (or ‘tier 1’) options on  $\mathbf{M}$ .

In case (i),  $A > C$  comes before  $B > C$ , so there’s some strict total order  $\triangleright$  which rationalizes the initial list such that  $A \triangleright B$ . But then, since  $\mathcal{N}(A, C^*) = \mathcal{N}(B, C^*)$  and  $A \triangleright B$ , it must be that  $A > C^*$  comes before  $B > C^*$ , wherefore  $A$   $C^*$ -precedes  $B$ . In case (ii),  $C > A$  comes before  $C > B$ , so there’s some strict total order  $\triangleright$  which rationalizes the initial list such that  $B \triangleright A$ . But then, since  $\mathcal{N}(C^*, A) = \mathcal{N}(C^*, B)$  and  $B \triangleright A$ , it must be that  $C^* > A$  comes before  $C^* > B$ , wherefore  $A$   $C^*$ -precedes  $B$ . In case (2), since  $A \sim C$  comes before  $C \sim B$ ,  $A \sim C^*$  must come before  $C^* \sim B$  (since the initial list is rationalizable). So  $A$  will  $C^*$ -precede  $B$ . (The opposite direction is exactly the same, with ‘ $C$ ’ swapped out for ‘ $C^*$ ’ and ‘ $C^*$ ’ swapped out for ‘ $C$ ’.)

**Definition 4.** A theory of rational choice is independent of clones iff: if  $C$  and  $C^*$  are clones,  $\mathbf{M}^*$  is a menu of options which includes  $C$  and  $C^*$ ,  $\mathbf{M} \stackrel{\text{def}}{=} \mathbf{M}^* \setminus \{C^*\}$ , and  $A \in \mathbf{M}$ , then it is permissible to select  $A$  from the menu  $\mathbf{M}$  iff it is permissible to select  $A$  from the menu  $\mathbf{M}^*$ :  $A \in \mathcal{P}(\mathbf{M}^*)$  iff  $A \in \mathcal{P}(\mathbf{M})$ .

**Proposition 4.** MIN is independent of clones.

*Proof.* Suppose that  $A \in \mathcal{P}(\mathbf{M})$ . Then, there is some preliminary preference  $\succeq_{\mathcal{L}}$  over the options in  $\mathbf{M}$ , constructed from an initial list  $\mathcal{L}$ , such that  $A$  is at the top of  $\succeq$ . Turning to the larger menu  $\mathbf{M}^*$ , consider the list  $\mathcal{L}^*$  which is identical to  $\mathcal{L}$ , except that, for all  $A \neq C, C^*$ , the preference  $[A, C^*]$  is added immediately after  $[A, C]$ . Let  $\succeq^*$  be the preliminary preference constructed from this list. There is no preference which is struck from  $\mathcal{L}$  but not from  $\mathcal{L}^*$ , since there are fewer preferences above each entry on  $\mathcal{L}$ —if a preference doesn't lead to a cycle with the preferences above it on  $\mathcal{L}^*$ , it can't lead to a cycle with a subset of those preferences. So, in particular, there is no preference of the form  $B \succ A$  which is struck from  $\mathcal{L}$  but not from  $\mathcal{L}^*$ . Since all such preferences must have been struck from  $\mathcal{L}$  for  $A$  to be a top element, all such preferences will be struck from  $\mathcal{L}^*$  as well, and  $A$  will be a top element of  $\succeq^*$ . So  $A \in \mathcal{P}(\mathbf{M}^*)$ .

Going in the other direction, suppose that  $A \in \mathcal{P}(\mathbf{M}^*)$ . Then, there is a preliminary preference  $\succeq^*$  over the options in  $\mathbf{M}^*$ , constructed from an initial list  $\mathcal{L}^*$ , such that  $A$  is at the top of  $\succeq^*$ . Let  $\mathcal{L}$  be just like  $\mathcal{L}^*$ , except that all preferences involving  $C^*$  have been removed. Suppose that there is some preference  $[X, Y]$  which is struck from  $\mathcal{L}^*$  but which remains unstruck on  $\mathcal{L}$ . Since  $[X, Y]$  is struck from  $\mathcal{L}^*$ , it must enter into a cycle with some of the unstruck preferences above it.  $C^*$  must be in this cycle (else, the same cycle would be created in  $\mathcal{L}$ ). So  $[X, Y]$  either creates a cycle of the form  $X \succeq^* Y \succeq^* \dots \succeq^* C^* \succeq^* \dots X$  or a cycle of the form  $Y \succeq^* X \succeq^* \dots \succeq^* C^* \succeq^* \dots Y$  (with at least one of these preferences strict). By Lemma 1, if preferences involving  $C^*$  remain unstruck, then so too do the corresponding preferences involving  $C$ . So, on  $\mathcal{L}^*$ ,  $[X, Y]$  either leads to a cycle of the form  $X \succeq^* Y \succeq^* \dots \succeq^* C \succeq^* \dots X$  or a cycle of the form  $Y \succeq^* X \succeq^* \dots \succeq^* C \succeq^* \dots Y$ . All of these preferences will remain unstruck on  $\mathcal{L}$ , since there are fewer preferences above them on  $\mathcal{L}$  than there are on  $\mathcal{L}^*$ . So  $[X, Y]$  will also be struck from  $\mathcal{L}$ . So there are no preferences struck from  $\mathcal{L}^*$  which are not also struck from  $\mathcal{L}$ . And, in particular, there is no preference of the form  $B \succ A$  which is struck from  $\mathcal{L}^*$  but not from  $\mathcal{L}$ . Since all such preferences are struck from  $\mathcal{L}^*$ , all such preferences are also struck from  $\mathcal{L}$ . So  $A$  will be at the top of  $\succeq$ . So  $A \in \mathcal{P}(\mathbf{M})$ .  $\square$