

548 A Appendix

549 Here we provide some additional details on three topics. First, we illustrate the role played by H3 in
550 our definition of harm. Second, we discuss in more detail how our approach differs from that of RBT.
551 Third, we present four more examples that illustrate how our definition handles issues which have
552 been discussed in the harm literature.

553 A.1 Discussion of H3

554 As we mentioned, H3 is intended to capture the intuition that preventing a worse outcome is not
555 harmful. For example, following the reasoning of the car manufacturer in Example 5, the system’s
556 decision to drive into the fence rather than doing nothing is not harmful because Bob would have
557 suffered even worse injuries had the system done nothing. Since H1 and H2 are satisfied for this
558 particular contrastive event, our definition would reach the opposite verdict if it weren’t for H3.
559 Note that the counterfactual comparative account (Definition 3) also says that there is no harm: the
560 alternative event under consideration would have given a worse outcome, so that C3 is not satisfied,
561 and therefore there is no harm. Considering H3 gives more insight into the differences between the
562 counterfactual account and ours.

563 Suppose that we consider some contrastive event $\vec{X} = \vec{x}'$ such that $(M, \vec{u}) \models \vec{X} = \vec{x} \wedge O = o$ and
564 $(M, \vec{u}) \models [\vec{X} \leftarrow \vec{x}'](O = o'')$, so C1 and C2 hold, and the first half of H2 holds if $o'' \neq o$: $\vec{X} = \vec{x}$
565 rather than $\vec{X} = \vec{x}'$ causes $O = o$ rather than $O = o''$. H3 plays no role if H1 is not satisfied, so for
566 simplicity, suppose that H1 also holds. Then it is easy to see that whenever $\mathbf{u}(O = o) \neq \mathbf{u}(O = o'')$,
567 our definition gives the same verdict as the counterfactual comparative definition for this particular
568 contrast (i.e., for this choice of \vec{x}'): if $\mathbf{u}(O = o) < \mathbf{u}(O = o'')$, then $o'' \neq o$, so H2 holds,
569 as do H3 and C3; it follows that both definitions declare $\vec{X} = \vec{x}$ a harm. On the other hand, if
570 $\mathbf{u}(O = o) > \mathbf{u}(O = o'')$, then neither C3 nor H3 hold (for this choice of \vec{x}').

571 What happens if $\mathbf{u}(O = o) = \mathbf{u}(O = o'')$? This can happen for two reasons:

- 572 1. there is no but-for causation, that is, $o = o''$;
- 573 2. there is but-for causation but the counterfactual outcome $O = o''$ happens to have utility
574 identical to the actual outcome.

575 Thus, roughly speaking (and ignoring the key role played by the default utility), our definition differs
576 from the counterfactual comparative account only if $\vec{X} = \vec{x}$ rather than $\vec{X} = \vec{x}''$ is not a but-for cause
577 of the actual utility: changing \vec{x} into \vec{x}'' does not change the agent’s utility.

578 Examples in which the first reason is relevant are widespread and crucial to our analysis, for those
579 are precisely the examples in which actual causation (Definition 1) and but-for causation come apart.
580 Our Late Preemption example (Example 2) offers one illustration, the literature on actual causation
581 contains many more. An example where the second reason is relevant involves a more subtle way in
582 which but-for causation comes apart from actual causation. Consider a “Sophie’s choice” like setting:
583 An agent must choose whether $X = 1$ or $X = 2$. There are two children, who will either live or die
584 depending on the choice: if $X = i$ is chosen, then child i lives ($L_i = 1$) and child $3-i$ dies ($L_{3-i} = 0$).
585 The possible outcomes are that both children live (o_{11}), just child 1 lives (o_{10}), just child 2 lives (o_{01}),
586 and neither child lives (o_{00}), where $d = \mathbf{u}(O = o_{11}) > \mathbf{u}(O = o_{10}) = \mathbf{u}(O = o_{01}) > \mathbf{u}(O = o_{00})$.
587 In fact, $X = 1$ is chosen, so we get but-for causality, but switching from $X = 1$ to $X = 2$ gives an
588 outcome of equal utility. However, if we hold $L_1 = 1$ fixed (which we can do in our framework to
589 show causality) and switch to $X = 2$, then we get the outcome $O = o_{11}$. Thus, in our framework
590 $X = 1$ harms the agent; in the causal counterfactual framework, it does not.

591 This emphasizes the point we (and RBT) made that one set of problems that occur in defining harm is
592 identical to the type of problems that occur in defining causation, and can be solved in the same way.

593 A.2 Comparison to RBT

594 In this section, we do a more careful comparison of our approach and RBT’s approach. RBT restrict
595 their analysis to choices made by agents, where different choices can be taken to have different
596 normative content (i.e., some choices are more normatively appropriate than others, although people

597 might disagree as to which is the more appropriate choice, as in our euthanasia example). This
598 assumption is critical for them, since it plays a key role in how they determine both the default
599 action and the contingency to hold fixed when checking condition AC2 in the definition of causation
600 (Definition 1). There are several problems with this approach.

601 First, as has often been pointed out in the harm literature (and is critical to the insurance industry!)
602 harm can be caused by events other than agent’s actions. Indeed we already came across such an
603 example: Batman getting a heart attack in Example 2 causes Batman harm. We would certainly like
604 an account of harm that applies to such “natural events”.

605 Second, by construction of their definition, whenever an agent performs the default action, there is
606 no harm according to RBT. Yet there are many instances in which doing what is morally preferable
607 causes harm, albeit accidentally. Simply imagine a doctor prescribing medication to a patient, and
608 the patient unfortunately suffering a very rare allergic reaction to the medication, where the reaction
609 is far worse than the initial condition that the patient had. Then clearly the doctor harmed the patient.
610 The most obvious choice of default action here is the actual action (and, in fact, RBT themselves
611 mention following “clinical guidelines” as an example of a default action in Appendix B). But this
612 means that according to RBT’s definition there would not be harm here.

613 Third, even if we focus on choices made by agents and assume that there’s a sensible default action,
614 there is a key difference between our definition of causality, which, as we said, is based on that
615 of Halpern [11, 12], and that of RBT (given in their Appendix A as Definition 9). Whereas in
616 AC2 we existentially quantify over the set \vec{W} , RBT give a definition of harm relative to a fixed
617 set \vec{W} , and assume that \vec{W} should be determined by normative considerations (as they say at the
618 end of their Appendix B, “when establishing harm the conditional contingency [i.e., choice of \vec{W}]
619 corresponds to a single contingency that is determined a priori based on our normative assumptions,
620 and taking the wrong contingency (or allowing for any contingency) will result in harm or benefit
621 being misattributed”). Nonetheless, RBT claim that their approach is equivalent to that of Halpern,
622 which is clearly not the case.

623 On a conceptual level, the same points arise for a normatively determined contingency as the ones we
624 brought up for the default action: we would also like to apply the notion of harm to natural events
625 (and thus to cases in which there does not seem to be any contingency that is morally preferable over
626 others), and there can be situations in which doing the right thing causes harm. Perhaps RBT could
627 try and resolve this by allowing there to be multiple contingencies that can be used when applying
628 Definition 9, but then they would have to somehow aggregate the different harms that we get for each
629 specific contingency; it is not at all clear how this would be done.

630 To defend their use of default actions and the idea of having a normatively determined contingency,
631 RBT consider two examples. In the first, Bob expects a government check for \$100, but does not get
632 one because, instead, Alice puts \$100 into his bank account, which disqualifies him from government
633 assistance. In the second, there are two do-gooders, Alice and Eve, who conspire to lift Bob out of
634 poverty. Alice gets there first, giving Bob \$100, but if she had not done so, Eve would have. We
635 agree with RBT that, in both examples, Alice is the cause of Bob getting \$100 rather than 0 (and
636 this follows easily from our definitions). We also agree that in the first case, Alice’s action does not
637 benefit Bob, while in the second it does. Although we do not give a definition of benefit in our paper,
638 taking the obvious analogue of our definition of harm, we would get this result by simply taking
639 different defaults in the two examples: the default in the first is that Bob gets \$100 (because that
640 is the societal expectation, given the government program) while in the second it is that Bob gets
641 \$0. We still get the arguably “right” answer although we existentially quantify in AC2. Using the
642 contingency only to establish causality as we do (rather than as a way to establish the amount of harm,
643 as RBT seem to do), we can still deal with all the examples, while also being able to deal with cases
644 where there are no obvious normative considerations that determine the appropriate contingency.

645 A.3 More Examples

646 We present four further examples that illustrate how our approach deals with the difficulties of
647 defining harm that have been highlighted in the literature.

648 The cases in Sections 4.1 and 4.2 all involved a binary outcome; there were only two relevant events
649 that could occur. Carlson et al. [4] discuss cases that involve more than two possible events in order to
650 argue against existing causal accounts. The following example forms one instance of their argument.

651 **Example 6 (Tear Gas)** The Joker sprays tear gas in exactly one of Batman’s eyes. If he had not
652 done that, he would have sprayed tear gas in both of Batman’s eyes, which would have made Batman
653 even worse off. One of the alternatives available to the Joker, however, was to simply leave Batman
654 alone.

655 Intuitively here Joker harms Batman when he sprays him. To argue that the “incorrect” answer is
656 obtained by the definition of harm they focus on, Carlson et al. consider a specific alternative event,
657 namely, that Joker sprays tear gas in both of Batman’s eyes, while observing that other alternatives
658 (like leaving Batman alone) are also available. Rather than existentially quantifying over \vec{x}' , as
659 we have done, (both in Definition 2 and the gloss of the counterfactual harm definition given in
660 Definition 3), they take a version of counterfactual harm where $\vec{X} = \vec{x}'$ is taken to be the closest
661 alternative to $\vec{X} = \vec{x}$ (according to some implicit, but unspecified, notion of closeness). Both our
662 definition of harm and our gloss of the counterfactual definition (with the obvious assumptions about
663 utility, and taking the default utility to be that of Batman being unharmed for our definition) agree
664 that Joker did harm Batman in this case, as we would expect.

665 In this example, there are three events of interest (Joker sprays tear gas in one eye; Joker sprays tear
666 gas in both eyes; Joker doesn’t spray tear gas at all). We can model this using a variable TG that
667 takes on three possible values (say, 0, 1, and 2). According to Definition 3, as long as one of them
668 leads to a better utility than what actually happened, there was harm. But as the golf clubs example
669 shows, this conclusion is not always justified; in general, we need to take defaults into account. \square

670 Now we present an example, due to Shiffrin [25], that illustrates the role of both the choice of the
671 range of variables in the causal model and the choice of default.

672 **Example 7** Betty is drowning in a fast-moving river. Veronica rescues her by grabbing her arm and
673 pulling her out, accidentally fracturing Betty’s humerus.

674 Did Veronica’s rescue harm Betty? Shiffrin claims it does because Veronica could have pulled her
675 out without breaking her arm. Indeed, Klocksiam [18], in his analysis, points out that “it seems
676 possible to rescue someone from drowning without breaking her arm”. The first step in our analysis
677 is to decide whether we should allow this possibility. That is, suppose that we have a variable P that
678 describes how and whether Veronica pulls out Betty. We can take $P = 0$ if Veronica does not pull out
679 Betty, $P = 1$ if she pulls her out by grabbing (and breaking) her arm. The modeler must then decide
680 whether to allow P to take a value, say 2, where $P = 2$ if Veronica rescues Betty in such a way that
681 Betty’s arm is not broken. Reasonable people might disagree whether such an event is possible. First
682 suppose we decide that P can take only values 0 and 1. Then the possible outcomes are that Betty
683 drowns ($O = 0$) or Betty is saved ($O = 1$). In this model, any utility function that makes the utility
684 of drowning worse than that of being saved would result in Veronica’s rescue not harming Betty.

685 Now suppose that we allow $P = 2$. Then we would take $O = 1$ to represent Betty being saved but
686 her arm being broken, and $O = 2$ to represent Betty being saved without her arm being broken. In
687 that case, whether Veronica harms Betty depends on the default. If we take the default utility to be
688 $u(O = 2)$ then Veronica does cause Betty harm, while if we take it to be $u(O = 0)$, she does not.
689 Note that the latter choice is quite defensible. Given Betty’s situation, making it out alive in whatever
690 way possible would presumably be all that matters to her. \square

691 This example clearly shows that to apply our framework in practice, it is important to have some
692 guidelines on what count as a reasonable choice, both in the choice of variables and values and the
693 choice of default value. As we mentioned in the introduction, Halpern and Hitchcock [13] discuss this
694 issue in the context of causal models; to the best of our knowledge, this issue has not been discussed
695 in the context of default values. While this issue is beyond the scope of the current paper, we should
696 make clear that we would not, in general, expect there to be a unique “correct” model. As we have
697 said repeatedly, reasonable people can disagree about these choices.

698 There is one final issue we would like to address: why we consider a contrastive definition rather
699 than just giving a definition in the spirit of the causal-counterfactual account. Definition 2 explicitly
700 invokes a contrastive outcome o' whose utility is better than that of the actual outcome. We could

701 have instead just defined harm as the result of causing an outcome whose utility is worse than the
702 default.

703 One reason why we did not do so is that the default utility is not always achievable, and it would be
704 counterintuitive to say that the agent was harmed if the outcome has a utility lower than the default,
705 even though it is the best possible outcome. For example, there are diseases for which a surgery can
706 only provide a temporary cure; in this case, a successful surgery gives the patient a temporary relief,
707 and an unsuccessful surgery results in the patient's death. While the default utility for the patient, as
708 for all people, is to be alive and healthy, saying that a successful surgery harmed the patient seems
709 wrong. In fact, defining harm as the result of causing an outcome with the utility worse than the
710 default provides counter-intuitive results even when the default utility is achievable, as the following
711 example demonstrates.

712 **Example 8 (Pills)** Consider the following vignette, again taken from [4] (where it is presented as a
713 problem for both the causal-counterfactual and contrastive causal-counterfactual accounts):

714 Barney suffers from a painful disease. On Monday, he can either take Pill A or
715 not. On Tuesday, he will have another choice, between taking Pill B or not. Barney
716 believes that he will be completely cured just in case he takes only Pill A, and
717 partially cured just in case he takes both pills. Accordingly, he takes Pill A on
718 Monday and does not take Pill B on Tuesday . . . He is, however, misinformed about
719 the effects of the pills. Taking only Pill A causes his disease to be merely partially
720 cured. If he had taken both pills, he would have been completely cured. Had he
721 not taken Pill A on Monday, on the other hand, nothing he could have done later
722 would have produced even a partial cure.

723 To capture this in our framework, let O be a three-valued variable that captures Barney's health:
724 $O = 2$ if he is fully cured, $O = 1$ if he is partially cured, and $O = 0$ if he is not cured at all. A and
725 B capture whether or not Barney takes pills A and B respectively. The equation for O is then such
726 that $O = 2$ if $A = B = 1$, $O = 1$ if $A = 1$ and $B = 0$, and $O = 0$ otherwise. As Barney considers
727 taking pill B only if he fails to take pill A, the equation for B is $B = \neg A$. The context is such that
728 $A = 1$; therefore, $B = 0$ and $O = 1$.

729 Carlson et al. claim that taking the pill does not harm Barney; we agree. Yet it easy to see that
730 $A = 1$ does cause $O = 1$. Indeed it is a but-for cause: had Barney not taken the pill, O would have
731 been 0. It is easy to see why this is a problem for the causal-counterfactual account: Barney would
732 have been better off if $O = 1$ had not obtained; specifically, he would be better off if O had been 2
733 (although this is not the outcome that results when changing A to 0 and therefore is not a problem for
734 the counterfactual comparative account). Carlson et al. also view it as a problem for the contrastive
735 causal-counterfactual account, because in applying it, they compare $O = 1$ to the outcome $O = 2$,
736 (which, again, is not the outcome that obtains by switching A to 0), since they take the closest world
737 to the one where Barney takes just one pill to be the world where he takes both pills. Our definition
738 avoids this problem. We do not consider the "closest" state of affairs. Rather, we compare $O = 1$ to
739 the outcome $O = 0$ caused by switching to $A = 0$. $O = 0$ has utility worse than that of the outcome
740 obtained from $A = 1$, so it is not a harm according to our definition, for what we view as the "right"
741 reasons. Assuming that the default utility is $u(O = 2)$, $A = 1$ does cause an outcome whose utility
742 is worse than the default and therefore a non-contrastive version of our definition would not have
743 given the desired result. □