

605 **Appendices**

606 **A**

607 In this appendix we present the general version of Definition 3 allowing harm and benefit to be  
608 measured along specific causal paths.

609 The path-specific counterfactual harm measures the harm caused by an action  $A = a$  compared to a  
610 default action  $A = \bar{a}$  when, rather than generating the counterfactual outcome by including all causal  
611 paths from  $A = \bar{a}$  to outcome variables  $Y$ , we consider only the effect along certain paths  $g$ . This is  
612 somewhat analogous to the path specific causal effect 5, as we are using the  $g$ -specific intervention  
613  $A = \bar{a}$  on  $Y$  in the counterfactual world relative to reference  $A = a$  (the factual action).

614 **Definition 9** (Path-specific counterfactual harm & benefit). *Let  $G$  be the DAG associated with model*  
615  *$\mathcal{M}$  and  $g$  be the edge sub-graph of  $G$  containing the paths we include in the harm analysis. The path*  
616 *specific harm caused by action  $A = a$  compared to default action  $A = \bar{a}$  is given by*

$$h_g(a, x, y; \mathcal{M}) = \int_{y^*} P(Y_{\bar{a}, \mathcal{M}_g} = y^* | a, x, y; \mathcal{M}) \max\{0, U(\bar{a}, x, y^*) - U(a, x, y)\} \quad (12)$$

$$= \int_{y^*, e} P(Y_{\bar{a}} = y^* | e; \mathcal{M}_g) P(e | a, x, y; \mathcal{M}) \max\{0, U(\bar{a}, x, y^*) - U(a, x, y)\} \quad (13)$$

617 Where  $Y_{\bar{a}, \mathcal{M}_g}$  is the counterfactual outcome  $Y$  under intervention  $do(A = \bar{a})$  in model  $\mathcal{M}_g$   
618 where  $\mathcal{M}_g$  is formed from  $\mathcal{M}$  by replacing the causal mechanisms for each variable  $f^i(pa^i, e) \rightarrow$   
619  $f_g^i(pa^i(g)^*, e) = f^i(pa^i(g)^*, pa^i(\bar{g}), e)$ , where  $Pa^i(\bar{g})$  is the set of parents of  $V^{(i)}$  that are not linked  
620 to  $V^{(i)}$  in  $g$  and  $pa^i(\bar{g})$  is the factual state of those variables.  $E = e$  is the joint state of the exogenous  
621 noise variables in  $\mathcal{M}$ . Likewise, the expected benefit is

$$b_g(a, x, y; \mathcal{M}) = \int_{y^*} P(Y_{\bar{a}, \mathcal{M}_g} = y^* | a, x, y; \mathcal{M}) \max\{0, U(a, x, y) - U(\bar{a}, x, y^*)\} \quad (14)$$

622 Note that if we following the construction of  $\mathcal{M}_g$  in 5 we get that  $\mathcal{M}_g$  is formed from  $\mathcal{M}$  by i)  
623 partitioning the parent set for each variable  $V^{(i)}$  in  $\mathcal{M}$  into  $Pa^i = \{Pa^i(g), Pa^i(\bar{g})\}$  where  $Pa^i(g)$   
624 are the parents that are linked to  $V^{(i)}$  in  $g$  and  $Pa^i(\bar{g})$  is the complimentary set, ii) replacing the  
625 mechanisms for each variable with  $f^i(pa^i, e^i) \rightarrow f_g^i(pa^i, e^i) = f^i(pa^i(g)^*, pa^i(\bar{g}), e^i)$  where  $pa^i(\bar{g})$   
626 takes the value of  $PA^i(\bar{g})_z$  in  $\mathcal{M}$  where  $A = z$  is the reference action. However, in (12) and (14)  
627 we condition on the state of all factual variables and assume no unobserved confounders, and the  
628 reference action is the factual action state. Therefore the state of  $PA^i(\bar{g})_a$  in  $\mathcal{M}$  is equal to the factual  
629 state of these variables, giving our simplified construction for  $\mathcal{M}_g$ .

630 We give examples of computing the path-specific harm in Appendices B,C.

631 **B**

632 In this appendix we discuss the omission problem and pre-emption problem 13, and the preventing  
633 worse problem 15, and show how these can be resolved using our definition of counterfactual harm  
634 (Definition 3 and its path-specific variant Definition 9).

635 **Omission Problem:** Alice decides not to give Bob a set of golf clubs. Bob would be happy if Alice  
636 had given him the golf clubs. Therefore, according to the CCA, Alice's decision not to give Bob the  
637 clubs causes Bob harm. However, intuitively Alice has not harmed Bob, but merely failed to benefit  
638 him 13.

639 **Solution:** The omission problem relies on the judgement that Alice does not have a ethical obligation  
640 to provide Bob with golf clubs, therefore her choice not to do so does not constitute harm to Bob. In  
641 our definition of harm, this judgement is encoded by Alice not giving Bob clubs by default, i.e. the  
642 desired harm query is the harm 'compared to the world where Alice does not give Bob clubs'. To  
643 compute the harm we construct the model  $\mathcal{M}$  comprising of two variables; Alice's action  $A \in \{0, 1\}$

644 where  $A = 0$  indicates ‘Bob not given clubs’ and  $A = 1$  ‘Bob given clubs’, and outcome  $Y \in \{0, 1\}$   
 645 where  $Y = 1$  indicates ‘Bob has clubs’ and  $Y = 0$  indicates ‘Bob does not have clubs’. By default,  
 646 Alice is not expected to give Bob clubs, which is encoded by choosing the default action  $A = \bar{a}$  where  
 647  $\bar{a} = 0$ . The causal mechanism for  $Y$  is  $y = a$ , i.e. Bob has clubs iff he is given them. Whatever utility  
 648 function describes Bob’s preferences, the action  $A = 0$  causes no harm in this model (Lemma 3  
 649 Appendix J) as  $P(Y_0 = y^* | A = 0, Y = y) = \delta(y^* - y)$  (factual  $a$  and counterfactual  $\bar{a}$  are identical)  
 650 and for non-zero harm we require  $y^* \neq y$ .

651 Note there are other reasonable scenarios where Alice’s actions would constitute harm. For example,  
 652 if Alice was a clerk in a golf shop and Bob had pre-paid for a set of golf clubs, we could claim that  
 653 ‘the clerk Alice harmed Bob by not giving him golf clubs’. In this case, we would expect Alice to  
 654 give Bob the clubs by default (she has a ethical obligation to do so) and the harm query we want  
 655 (implied by our ethical assumptions about clerks) is where the default action is  $\bar{a} = 1$ . By choosing  
 656 not to— $A = 0$ —Alice causes harm to Bob. For example, if Bob’s utility is  $U(y) = y$  (i.e. 1 for  
 657 clubs, 0 for no clubs), then the harm caused by Alice is  $P(Y_{A=1} = 1 | A = 0, Y = 0) = 1$ . So we  
 658 can see that the choice of default action is vital for expressing these different normative assumptions.

659 **Preemption Problem:** Alice robs Bob of his golf clubs. A moment later, Eve would have robbed  
 660 Bob of his clubs. Therefore, Alice’s action does not cause Bob to be worse off as he would have  
 661 lost his clubs regardless of her actions, and so by the CCA Alice does not harm Bob by robbing him.  
 662 However, intuitively Alice harms Bob by robbing him, regardless of what occurs later [13].

663 Let  $A = \{1, 0\}$  denote Alice {robbing, not robbing} Bob respectively, and similarly  $E = \{1, 0\}$  for  
 664 Eve.  $B = \{1, 0\}$  denotes Bob {has clubs, does not have clubs}. Assume Bob’s utility is  $U(b) = b$ .  
 665 The causal mechanisms are  $e = 1 - a$  (Eve always robs Bob if Alice doesn’t) and  $b = 1 - a \vee e$   
 666 (Bob has no clubs if either Alice or Eve robs him). See Figure 3 for the causal model depicting these  
 variables.

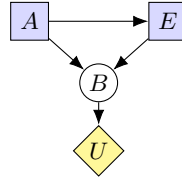


Figure 3: SCM depicting the preemption problem.

667

668 Note that while Alice’s action is an actual cause of Bob not having clubs, it is also an actual cause of  
 669 Eve not robbing Bob, which is an event equally as bad as Alice robbing Bob. Intuitively, when we  
 670 claim that Alice robbing Bob was harmful, we are making a claim about the effects of Alice’s actions  
 671 on Bob independently of their effect on Eve’s actions (independent of the effect that her action has  
 672 mediated through Eves action, preventing Eve from robbing Bob), i.e we are concerned with the  
 673 direct harm caused by Alice’s actions on Bob.

674 The relevant harm query is the path-specific harm where we compare to the default action where Alice  
 675 does not rob Bob,  $\bar{a} = 0$ . We want to determine the harm caused by Alice’s action independently of  
 676 its effect on Eve’s action, which we do by blocking the path  $\bar{g} = \{A \rightarrow E\}$ . Applying Definition 9  
 677 amounts to replacing the mechanism for  $E$  with  $f^E(a) \rightarrow f_g^E(A = 1) = 0$ , i.e.  $E$  is evaluated for  
 678 the factual value of  $A$ . We then compute the harm using the counterfactual default action  $A = 0$ ,  
 679 giving the counterfactual  $B(A = 0, E = 0) = 1$ , which gives a counterfactual utility of 1 compared  
 680 to a factual utility of 0. Therefore Alice directly harmed Bob by robbing him.

681 Note we can also choose a different model where we explicitly represent the outcomes of the two  
 682 agents decisions and the temporal order in which they occur (Figure 4). In this case the relevant harm  
 683 query is essentially the same; the path specific harm where we determine the harm caused by Alice’s  
 684 action independently of the effect it has on whether or not Eve robs Bob (i.e.  $\bar{g} = \{R_A \rightarrow R_E\}$ ).

685 **Preventing worse:** We provide two versions of the preventing worse problem [15] which have  
 686 identical causal models but intuitively different harms attributed to Alice’s action.

687 Case 1: Bob has \$2. The thief Alice is stalking Bob in the marketplace and notices that Eve (a more  
 688 effective thief) is also stalking Bob. Seeing Eve before Eve notices her, Alice decides to make her

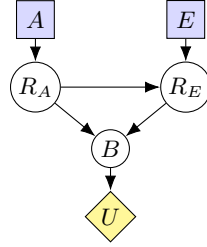


Figure 4: SCM depicting the preemption problem explicitly representing the temporal asymmetry between Alice and Eve’s actions effecting Bob.

689 move first. She steals \$1 from Bob. Eve was going to steal \$2 from Bob, but is incapable of doing so  
 690 if someone else robs him first (e.g. Bob realizes he’s been robbed and call for the police, making  
 691 further robbery impossible). Seeing that Bob was robbed by Alice she decides not to rob him.

692 Case 2: Eve has captured Bob and intends to torture him to death. Alice sees this, and is too far away  
 693 to prevent Eve from doing so. She has a line of sight to Bob (but not Eve) and can shoot him before  
 694 eve has a chance to torture him to death, resulting in a painless death.

695 The causal model describing both of these cases is depicted in Figure 5. Let  $E = \{1, 0\}$  denote  
 696 if Eve is present or not,  $A \in \{1, 0\}$  be Alice’s action (rob, shoot) or not,  $AB \in \{1, 0\}$  denote the  
 697 outcome following Alice’s action (Bob is robbed of \$1 / bob is shot, or not) and let  $EB \in \{1, 0\}$   
 698 denote Eves action on Bob (Bob is robbed of \$2 / Bob is tortured, or not). Let  $Y \in \{0, 1, 2\}$  denote  
 699 Bob’s outcome, with 2 being the best (Bob has \$2 in Case 1, Bob survives in Case 2), 1 being the  
 700 second worst (Bob has \$1 in Case 1, is killed painlessly in Case 2), and 0 the worst (Bob has \$0 in  
 701 Case 1, died painfully in Case 2). The causal mechanisms are  $a = e$  (e.g. Alice shoots/robs if Eve is  
 702 present),  $ab = a$  (Alice’s bullet hits with certainty / successfully robs with certainty),  $eb = e(1 - ab)$   
 703 (Eve tortures Bob if she is present and he is not shot / eve robs Bob if she is present and hasn’t been  
 704 robbed already), and  $y = ab + 2(1 - ab)(1 - eb)$  (Case 1: if Bob is shot he dies quickly, else if Eve  
 705 tortures him he dies slowly, else he lives, Case 2: Bob has \$2 if not robbed, \$1 if robbed by Alice, \$0  
 706 if not robbed by Alice and robbed by Eve).

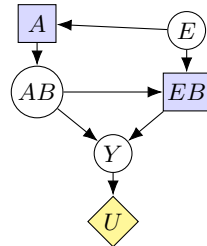


Figure 5: SCM depicting the preventing worse problem.

707 In Case 1, Alice intuitively harms Bob by robbing him. The argument supporting this is that Alice’s  
 708 robbery caused Bob to lose \$1, regardless of the fact that Alice’s action prevented a worse robbery  
 709 by Eve. However, for Case 2 it is argued in [15] that Alice intuitively didn’t harm Bob. While Bob  
 710 died due to Alice shooting him, this action was intended to prevent a worse outcome from occurring  
 711 (Bob being tortured to death), which would have happened with certainty had Alice not shot him.  
 712 However, these two scenarios are described by equivalent causal models—only the variables have  
 713 been re-labeled. However, the ethical assumptions differ between Case 1 and 2.

714 From this we conclude that to satisfactorily describe these two situations we need two different harm  
 715 queries. In either case, one of these harm queries is the morally relevant one and the other is not,  
 716 and to do this we use the path-independent and path-specific harms. Note, this is no different than  
 717 in causal analysis where in certain problems the casual effect is the desired query and in others the  
 718 path-specific effect is the desired query [5]. For Case 1 we use the path-specific harm (Definition  
 719 9) to determine the harm caused by Alice robbing Bob independently of what effect it had on Eve’s  
 720 action. We block the path  $\bar{g} = \{AB \rightarrow EB\}$  and use the default action  $\bar{a} = 0$ . In the counterfactual  
 721 world, this gives  $AB = 0$  and  $EB = f^{EB}(A = 1, E = 1) = 0$ , and therefore  $Y = 2$ , and so the

722 direct harm of Alice robbing Bob is  $2 - 1 = 1$  compared to not robbing him. For Case 2, we note that  
723 while Alice shooting Bob is arguably intrinsically harmful (as is captured by the direct harm of 1  
724 caused by  $A = 1$  if we calculate the path-specific harm as in Case 1), this is not the morally relevant  
725 harm that we are referring to when we say that intuitively Alice did not harm Bob by shooting him.  
726 The reason Alice fired the shot was precisely because of its mediating effect on  $Y$  through Eve’s  
727 actions (preventing her from torturing him to death). From this we infer that the morally relevant  
728 harm in this case is the path-independent harm. This we calculate using Definition 3 and the default  
729 action  $\bar{a} = 0$ , which in the counterfactual world gives  $AB = 0$ ,  $EB = 1$ ,  $Y = 0$  and hence  $U = 0$ ,  
730 compared to the factual utility  $U = 1$ , giving the desired result that Alice did not harm Bob compared  
731 to not shooting him. Note that if we favoured the path-independent or path-dependent harm a priori  
732 this would either fail to detect harm in Case 1 or incorrectly attribute harm to Alice in Case 2.

733 We argue from these two examples that there is no single causal formula for harm that is correct in all  
734 scenarios—in some the morally relevant measure of harm is path-specific (e.g. the direct harm), in  
735 others it is the path-independent harm. This is in contrast to other approaches to define harm with  
736 a single causal formula that applies to all scenarios, namely 8, and we discuss this approach and  
737 provide counterexamples to it in Appendix D

## 738 C

739 In this Appendix we discuss selecting and interpreting default actions, harmful events, and various  
740 edge cases not covered in the main body of our paper such as harmful default actions. Note that while  
741 the CCA (Definition 2) states ‘[the action] had not been performed’, this should not be interpreted as  
742 ‘do nothing’, as doing nothing is often a valid action choice and should be included as an element of  
743  $A$ . Instead, we argue that statements about harm often implicitly assume some default action, often  
744 following from ethical or normative assumptions (although this is not always the case). Indeed, in  
745 Appendix D we show in Example 3 that being able to enforce a unique default action is vital in some  
746 scenarios to give intuitive results.

747 Our definition of harm treats the default action as an integral part of the harm query, just as a reference  
748 treatment is necessary when defining treatment effects 79. These default-dependent measures of  
749 harm can be converted to default-independent measures if desired, e.g. by taking the max over all  
750 default actions, but in all of the examples we explore this is not desirable. We also note that while  
751 the examples outlined in the main text assume deterministic default actions, it is trivial to extend our  
752 definitions to non-deterministic default actions by replacing  $\text{do}(A = \bar{a})$  in Definition 3 with a soft  
753 intervention (e.g. 17). For examples of how the default action resolves the omission problem, and  
754 when path-specific and path-independent harm should be used, see Appendix B

755 **Default actions:** In some cases harm is attributed to an agent by comparing to normative actions  
756 or policies, and so the default action is often implied by the situation or determined by normative  
757 assumptions (e.g. Example 1 below). For example, in a case of negligence a doctor’s actions may  
758 be compared to clinical guidelines, or in a randomized control trial the harm caused by a drug is  
759 typically determined by comparing to the outcomes that would have occurred if the trial participants  
760 had instead been given a placebo. This is not always the case however (Example 2). The relevant  
761 harm query can also compare to actions that the agent could never take (Example 3). While some  
762 have argued against comparative accounts on the grounds that it is not always clear which comparison  
763 is needed 35, this problem arises due to the ambiguity of statements about harm rather than due to a  
764 problem with its formal definition (note, we do not consider scenarios where the agent’s action alters  
765 the user’s utility function). Clearly, there is not a single universal comparison or default action that is  
766 suitable for all situations (this assumption leads to the omission problem, described in Appendix B),  
767 and the ability to explicitly choose the comparison is a feature rather than a fault with the CCA.

768 **Example 1:** The claim ‘the doctor harmed the patient by not treating them’ and ‘the bystander with no  
769 medical training failed to benefit the patient by not treating them’ both tacitly assume different default  
770 actions. In the first, the doctor has an ethical obligation to treat the patient (e.g. the Hippocratic oath),  
771 and likewise the patient can expect to be treated by the doctor. Hence if they are not treated, harm  
772 can occur. In the second, the bystander may have no ethical obligation to help the patient (depending  
773 on our ethical assumptions) and so the intuitive choice of default action is to not treat the patient. In  
774 both of these examples, the ‘correct’ default action depends on the situation and in these examples is  
775 informed by our assumptions as to the ethical obligations of the agent.

776 **Example 2:** Consider a drug that a doctor is expected to provide to a patient which rarely causes  
 777 severe side effects. For a given patient, those side effects occur, and clearly the drug has harmed  
 778 the patient. Perhaps the most obvious harm measure to capture this would be the total harm caused  
 779 by the treatment compared to the default action where the doctor did not treat the patient at all, or  
 780 provided them with a different treatment. Each of these is a different but valid harm query, and the  
 781 correct one will depend on the situation. For example if we are measuring the harm due to the doctors  
 782 negligence, we should compare to the normative default action alone (and should find zero harm  
 783 due to negligence as the doctor followed the correct protocol), whereas if we are trying to establish  
 784 harm caused by the drug to this patient due to the side effects it caused, we should use the default ‘no  
 785 treatment’.

786 **Example 3:** How can we deal with cases where every action available to the agent is harmful? In  
 787 this case harm is still measured compared to some default action, even if the action is idealized and  
 788 not actually available to the agent. For example, if a doctor is forced at gun point to choose between  
 789 administering two poisons that will harm the patient, we can still measure this harm compared to the  
 790 counterfactual action where the doctor does not treat the patient, even if this action is not available to  
 791 the doctor.

792 **Harmful events:** Finally, we note that while we focus on harmful actions due to our focus on training  
 793 ethical artificial agents, our results extend trivially to harmful events as actions are formally equivalent  
 794 to events in the causal models we consider, and instead of default actions we can use default events.

## 795 D Comment on Beckers et. al.

796 In this appendix we discuss an alternative proposal for qualitatively defining harm [8], which was  
 797 developed following the presentation of our preliminary results. We describe this definition (which  
 798 we refer to as BCH) and three examples where BCH leads to counter-intuitive results (intuitively  
 799 harmful actions being identified as not harmful or vice versa). First we present a simplified version of  
 800 the BCH definition of harm where we restrict our attention to attributing harm to single actions.

801 **Definition 10.**  $A = a$  rather than  $A = a'$  causes  $Y = y$  rather than  $Y = y'$  in the model  $M$  for  
 802 exogenous noise state  $E = e$  iff;

- 803 1.  $A(e) = a$  and  $Y(e) = y$
- 804 2. There exists a set of environment variables  $W$  with factual state  $W(e) = w$  such that  
 805  $Y_{A=a', W=w}(e) = y'$
- 806 3.  $A = a$  is minimal; There is no strict subset of the set of variables  $\tilde{A} \subset A$  such that for  
 807  $\tilde{A} = \tilde{a}$  we can satisfy conditions 1. and 2.

808 In the following we will focus on scenarios where we can consider single action variables alone and  
 809 so we can ignore condition 3 in Definition 10.

810 **Definition 11** (BCH harm).  $A = a$  harms the user in model  $\mathcal{M}$  and exogenous noise state  $E = e$ , if  
 811 there exists an outcome  $Y = y$  an action  $A = a'$  such that,

- 812 H1  $U(y) < d$  where  $d$  is the default utility
- 813 H2  $\exists Y = y'$  s.t.  $A = a$  rather than  $A = a'$  causes  $Y = y$  rather than  $Y = y'$  and  
 814  $U(y') > U(y)$ .
- 815 H3  $U(y) \leq U(y'')$  for the unique  $y''$  such that  $Y_{a'}(e) = y''$

816 If we restrict to deterministic models (i.e.  $P(E = e)$  deterministic), attempt to only determine if  
 817 harm is non-zero rather than quantify how much harm is caused (i.e. map all non-zero harm values to  
 818 1), and assume that the users utility function is independent of the agents action  $A$  and the context  $X$   
 819 given the outcome  $Y$ , then it is possible to directly compare our harm measure to that proposed in  
 820 BCH. First, we present three problematic cases where the BCH gives counter-intuitive results. We  
 821 then attempt to diagnose why our approaches give different answers in these cases.

822 **Example 1: two thieves.** (repeat of Case 1 in Appendix B). Bob has \$2. The thief Alice is stalking  
 823 Bob in the marketplace and notices that Eve (a more effective thief) is also stalking Bob. Seeing Eve

824 before Eve notices her, Alice decides to make her move first. She steals \$1 from Bob. Eve was going  
 825 to steal \$2 from Bob, but is incapable of doing so if someone else robs him first (e.g. Bob realizes  
 826 he’s been robbed and call for the police, making further robbery impossible). Seeing that Bob was  
 827 robbed by Alice she decides not to rob him. The causal model for this scenario is described below,

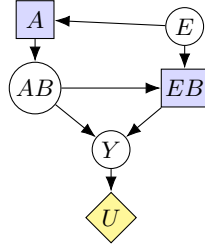


Figure 6: SCM depicting the two robbers problem.  $E \in \{1, 0\}$  denote Eve {present, not present},  $A = \{1, 0\}$  denotes Alice decides to {rob, not rob} Bob,  $AB \in \{1, 0\}$  denotes Bob is {robbed, not robbed} by Alice,  $EB = \{1, 0\}$  Eve attempts to {rob, not rob} Bob.  $Y$  denotes how much money Bob has finally. Causal mechanisms  $a = e$  (Alice robs if Eve is present),  $ab = a$  (Alice always succeeds in robbing Bob),  $eb = e(1 - ab)$  (Eve robs Bob if she is present and he hasn’t been robbed already), and  $y = ab + 2(1 - ab)(1 - eb)$  (if Bob is not robbed at all he has \$2, if Alice Robs him he has \$1, and if Eve robs him and Alice does not he has \$0).

828 Intuitively Alice harmed Bob by robbing him, but by Definition [II](#) she did not. The only available  
 829 counterfactual action for Alice is  $\bar{a} = 0$ . This counterfactual action (with no contingencies) leads  
 830 to the counterfactual outcome  $Y_{A=0}(e) = 0$ , i.e. if Alice doesn’t rob Bob then Eve will, resulting  
 831 in a lower utility  $U(Y = 1) > U(Y = 0)$ . Therefore H3 is not satisfied and Alice did not harm  
 832 Bob by robbing him. We discuss this problem further in Appendix [B](#) and argue that the morally  
 833 relevant harm query in this scenario is the direct (path-specific) harm of Alice robbing Bob compared  
 834 to not robbing him ( $\bar{a} = 0$ ), independent of the benefit caused by preventing Eve from robbing him  
 835 (blocking  $\bar{g} = \{AB \rightarrow EB\}$ ). Applying Definition [9](#) it is simple to check the path-specific harm  
 836 described is 1.

837 **Example 2: robber & Samaritan** An intuitive property of harm is that the harms caused by one  
 838 agent’s actions should not by default be cancelled out by the another agents beneficial actions—e.g.  
 839 stabbing someone is harmful, regardless of whether or not a doctor will treat the wound in response.  
 840 This ability to disentangle agent A’s harm from agent B’s benefit is vital for determining harm  
 841 in complex scenarios involving multiple actions or events. For example: Bob has \$1 and Alice  
 842 steals it. Seeing this, Eve feels bad for Bob and later gifts Bob a dollar, restoring him to his initial  
 843 funds. Intuitively we would say Alice harmed Bob and Eve benefited him, or at least it would be  
 844 counter-intuitive to say that Alice robbing Bob was not harmful because at a later time Bob’s finances  
 845 were restored by a second agent (Eve). The causal model describing this situation is depicted below,

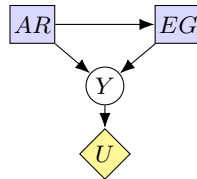


Figure 7: SCM depicting the preventing worse problem. Alice {robs, doesn’t rob} Bob (AR) is denoted  $AR = \{1, 0\}$ . Eve gives (EG) Bob money if she sees he has been robbed ( $eg = ar$ ). Bobs money is his initial money, minus any theft and adding any gifts  $y = 1 - ar + eg$ . Let  $U(y) = y$ .

846 The intuitive default utility for this scenario is  $d = 1$  (Bob expects to have \$1), but as the factual  
 847 outcome is  $Y = 1$  then H1 cannot be satisfied. To satisfy H1 we would need to choose a default  
 848 utility  $d > 1$  which amounts to Bob having by default more money than he would have regardless  
 849 of Alice robbing him (e.g. assuming Bob can expect to become richer following a robbery). We  
 850 therefore either recover a counter-intuitive answer (Alice did not harm Bob by robbing him), or have  
 851 to use a counter-intuitive default utility that is hard to justify beyond choosing whichever value gives  
 852 the desired answer.



853 Our approach is to measure the direct harm caused by  $A = 1$  (Alice robs) compared to  $\bar{a} = 0$  (Alice  
854 doesn't rob), blocking the path  $\bar{g} = \{AS \rightarrow EG\}$ . This disentangles that harm caused by Alice  
855 robbing Bob from the benefit due to this action causing Eve to help Bob. It is simple to check that  
856 this results in a harm of 1. As described in Appendix B, this is the intuitive choice of harm query  
857 as we are interested in the harm caused directly to Bob by Alice robbing him, independent of the  
858 indirect effect of causing Eve to benefit him. In other (causally equivalent) scenarios described in  
859 Appendix B, the intuitive harm query we desire if the total (path-independent) harm, and as with  
860 default actions this has to be implied from the context.

861 **Example 3: omission problem.** In this example we present an extension of the omission problem  
862 that is violated by the BCH definition of harm. The Phoenicians are a moderately wealthy people,  
863 collectively owning \$2. The Romans can decide to gift them an extra \$2 or do nothing, and they have  
864 no moral obligation to give them anything. Unbeknownst to the Romans, the Carthaginians decide  
865 that if the Romans don't give the Phoenicians anything they will attack them, stealing all of their  
866 money. But if the Romans do gift the Phoenicians \$2, the Phoenicians will become too powerful and  
867 the Carthaginians won't attack. The Phoenician's utility is equal to how much money they have.

868 The Romans decide not to gift the Phoenicians anything, and they are attacked by the Carthaginians  
869 and have all their money stolen. Intuitively, the Romans didn't harm the Phoenicians (any harm was  
870 caused by the Carthaginians)—instead the Romans failed to benefit them. However, by the BCH  
871 account the Roman's harmed the Phoenicians.

872 To see this, first note that if  $d \leq 0$  then the Carthaginians actions do not constitute harm, as the  
873 factual utility is equal to the default and H1 cannot be satisfied. This would be a counter-intuitive  
874 result, so we assume that  $d > 0$ . The Phoenicians end up with no money, so H1 is satisfied as  
875  $U(y) < d$ . H2 is also satisfied by a simple but-for counterfactual because if the Romans had given  
876 the Phoenicians money, the Carthaginians wouldn't have attacked and the Phoenicians would have \$4  
877 which is more than their factual \$0. Finally, H3 is satisfied by the same argument as H2. Therefore  
878 the Romans harmed the Phoenicians. Applying our methods, it is sufficient to note that the implied  
879 default action should be  $A = 0$  (By default we do not expect the Romans to give money to the  
880 Phoenicians, reflecting the ethical assumptions implicit in the 'failure to benefit' assertion) and this  
881 gives a counterfactual harm of zero because the factual and counterfactual actions are identical.

882 **Analysis:** Why do these issues arise? Firstly, the BCH account of harm proposes a single casual  
883 formula for harm that applies to all scenarios, allowing for any counterfactual action or contingency  
884 to establish harm much as is done in actual causality [34]. If H3 was not included, this could result in  
885 harm being attributed in cases of 'preventing worse' (as pointed out in [8] and described in Appendix  
886 B), but H3 is included to fix this by requiring that benefit does not occur in the case where no  
887 contingency is taken, which in these examples is the same as requiring that an action cannot be  
888 harmful if its total (path-independent) benefit is non-zero. But this is precisely the case in Example  
889 1, where Alice prevents a worse outcome but, intuitively, we would want to ascribe harm to her  
890 actions. By separating direct and indirect harm, we can see that her actions were indirectly beneficial  
891 (she prevented a worse robbery), but directly harmful, and in this scenario the morally relevant (i.e.  
892 'intuitive') measure of harm is the direct harm. In the equivalent Case 2 example in Appendix B, the  
893 harm query we intuitively want is the total harm rather than the path-specific harm. This points to the  
894 conclusion that a one-size-fits-all harm query is not tenable, given that the intuitive measures of harm  
895 we desire are sometimes path-dependent and sometimes path-independent.

896 Secondly, Example 2 suggests that approaches using default utilities are not tenable, because they  
897 preclude the possibility of the user being harmed by any action or event that occurred previously  
898 if, in the end, the user obtains the default utility. Clearly this is not the case in general—users can  
899 achieve the expected or default outcome (e.g. leaving the market with as much money as they came  
900 in with) and still have been harmed. The BCH therefore cannot robustly detect harm in cases where  
901 both benefit and harm occur unless we choose large values of  $d$  (essentially removing  $d$  from the  
902 analysis). But it is not clear how these default utility values can be justified beyond fixing a problem  
903 that they cause—seeing as these large values of  $d$  do not correspond to any utility the user can expect  
904 to have (e.g. in Example 2 it would require that the user can expect to be richer than they initially  
905 were following a robbery).

906 Thirdly, in Example 3 we see that by allowing for any default action the BCH account can end up  
907 misattributing harm in cases where the agent has no ethical duty to act (or by extension, has a duty  
908 to perform specific actions). This is because the BCH allows the counterfactual action to take any

909 value—in this case, the Romans gifting the Phoenicians money. This can be avoided by by not  
 910 attributing harm to the Romans using counterfactual actions that they could never be expected to take  
 911 from an ethical standpoint (i.e. using default actions). Just by allowing the Romans to (in theory)  
 912 give any positive amount of money, we could even make the harm they cause by not giving the  
 913 Phoenicians anything arbitrarily large. In conclusion, by evaluating over all possible actions that the  
 914 agent could take the BCH doesn't allow normative assumptions about actions (e.g. to do with the  
 915 ethical responsibility to act or not act) to be included in the harm query.

## 916 E

917 In this Appendix we prove Theorem [1](#). Noting that  $\max\{0, U(a, x, y) - U(\bar{a}, x, y)\} -$   
 918  $\max\{0, U(\bar{a}, x, y^*) - U(a, x, y)\} = U(a, x, y) - U(\bar{a}, x, y^*)$ , subtracting the expected harm from  
 919 the expected benefit (Def [3](#)) gives,

$$\mathbb{E}[b|a, x; \mathcal{M}] - \mathbb{E}[h|a, x; \mathcal{M}] \quad (15)$$

$$= \int_{y, y^*} P(y, Y_{\bar{a}} = y^* | a, x; \mathcal{M}) (U(a, x, y) - U(\bar{a}, x, y^*)) \quad (16)$$

$$= \int_y P(y|a, x) U(a, x, y) - \int_{y^*} P(Y_{\bar{a}} = y^* | x) U(\bar{a}, x, y^*) \quad (17)$$

$$= \mathbb{E}[U|a, x] - \mathbb{E}[U|\bar{a}, x] \quad (18)$$

## 920 F

921 In this Appendix we derive the SCM model for the treatment decision task in examples 1 and 2, and  
 922 calculate the average treatment effect and counterfactual harm.

923 Patients who receive the default ‘no treatment’  $T = 0$  have a 50% survival rate.  $T = 1$  has a 60%  
 924 chance of curing a patient, and a 40% chance of having no effect, with the disease progressing as if  
 925  $T = 0$ , whereas  $T = 2$  has a 80% chance of curing a patient as a 20% chance of killing them, due to  
 926 some unforeseeable allergic reaction to the treatment.

927 Next we evaluate this expression for our two treatment by constructing an SCM for the decision task.  
 928 The patient’s response to treatment is described by three independent latent factors (for example  
 929 genetic factors) that we model as exogenous variables. Firstly, half of the patients exhibit a robustness  
 930 to the disease which means they will recover if not treated, which we encode as  $E^1 \in \{0, 1\}$  where  
 931  $e^1 = 1$  implies robustness with  $P(e^1 = 1) = 0.5$ . Secondly, the patients may exhibit a resistance to  
 932 treatment 1 indicated by variable  $E^2$ , with  $e^2 = 1$  implying resistance with  $P(e^2 = 1) = 0.4$ . Finally,  
 933 the patients can be allergic to treatment 2, indicated by variable  $E^3$  with  $e^3 = 1$  and  $P(e^3 = 1) = 0.2$ .  
 934 Given knowledge of these three factors the response of any patient is fully determined, and so we  
 935 define the exogenous noise variable as  $E^Y = E^1 \times E^2 \times E^3$  with  $P(e^Y) = P(e^1)P(e^2)P(e^3)$ .

936 Next we characterise the mechanism  $y = f(t, e^Y) = f(t, e^1, e^2, e^3)$  where  $f(0, e^Y) = [e^1 = 1]$   
 937 (untreated patients recover if they are robust),  $f(1, e^Y) = [e^1 = 1] \vee [e^2 = 0]$  (patients with  $T = 1$   
 938 recover if they are robust or non-resistant) and  $f(2, e^Y) = [e^3 = 0]$  (patients with  $T = 2$  recover if  
 939 they are non-allergic), where  $[X = x]$  are Iverson brackets which return 1 if  $X = x$  and 0 otherwise,  
 940 and  $\vee$  is the Boolean OR.

941 The recovery rate for  $T = 1$  and  $T = 2$  can be calculated with [1](#) to give  $P(Y_1 = 1) = P(e^1 =$   
 942  $1 \vee e^2 = 0) = 1 - P(e^1 = 0)P(e^2 = 1) = 0.8$ , and likewise  $P(Y_2 = 1) = P(e^3 = 0) = 0.8$ . Hence  
 943 the two treatments have identical outcome statistics (recovery/mortality rates), and all observational  
 944 and interventional statistical measures are identical, such as risk, expected utility and the effect of  
 945 treatment on the treated. Note as there are no unobserved confounders the recovery rate for action  
 946  $A = a$  is equal to  $\mathbb{E}[Y_a]$ .

947 We compute the counterfactual expected harm by evaluating [4](#), noting that  $Y_0^*(e) = 1$  if  $e^1 = 1$ ,  
 948  $Y_1^*(e) = 0$  if  $e^1 = 0$  and  $e^2 = 1$ , and  $Y_2^*(e) = 0$  if  $e^3 = 1$ . This gives  $P(Y_1 = 0, Y_0^* = 1) = 0$ ,  
 949 i.e. there are no values of  $e^Y$  that satisfy both  $Y_1(e) = 0$  and  $Y_0(e) = 1$ , and therefore  $\text{do}(T_1 = 1)$   
 950 causes zero harm. However,  $P(Y_2 = 0, Y_0^* = 1) = P(e^1 = 1)P(e^3 = 1) = 0.1$ , and so  $\text{do}(T_2 = 2)$



951 causes non-zero harm. This is due to the existence of allergic patients who are also robust, and will  
 952 die if treated with  $T = 2$  but would have lived had  $T = 0$ .  
 953

## 954 G

955 In this Appendix we derive the policies of agents 1-3 in Example 3. We note that outcome  $Y$  is  
 956 described by a heteroskedastic additive noise model with the default action  $\bar{a}$  (no action) corresponding  
 957 to  $A = 1$ ,  $K = 1$ . The expected harm is given by Theorem 5 with  $\sigma(\bar{a}) = 100$ ,  $\sigma(A = 2) = 100$ ,  
 958  $\sigma(A = 3) = 0$  and  $\sigma(A = 1, K) = 100K$ .  $\mathbb{E}[U|\bar{a}] = 100$   $\mathbb{E}[U|A = 2] = 110$ ,  $\mathbb{E}[U|A = 3] = 80$   
 959 and  $\mathbb{E}[U|A = 1] = 100K$ , where we have used  $\text{Var}(KY) = K^2\text{Var}(Y)$  and  $\text{Var}(Y + 10) = \text{Var}(Y)$

960 Agent 1 takes action 1 and the maximum value  $K = 20$  as this extremizes  $\mathbb{E}[U|a]$ .

961 Agent 2 chooses  $a = \arg \max_a \{\mathbb{E}[Y|a] - \text{Var}(Y|a)\}$  which for each action is given by,

$$E[Y|A = 1] - \lambda \text{Var}(Y|A = 1) = 100K - 100^2 K^2 \lambda \quad (19)$$

$$E[Y|A = 2] - \lambda \text{Var}(Y|A = 2) = 110 - 100^2 \lambda \quad (20)$$

$$E[Y|A = 3] - \lambda \text{Var}(Y|A = 3) = 80 \quad (21)$$

962 For action 1 the optimal  $K = 1/200\lambda$ , which gives  $E[Y|A = 1] - \text{Var}(Y|A = 1) = 1/4\lambda$ . Note  
 963 that  $1/4\lambda > 110 - 100^2 \lambda$  for  $\lambda < 0.0032$ , which  $80 > 1/4\lambda$  for  $\lambda > 0.003125$ . Therefore there is  
 964 no value of  $\lambda$  for which agent 2 selects action 2, choosing action 1 for  $\lambda < 0.003125$  and action 3  
 965 otherwise.

966 For agent 3 applying Theorem 5 gives,

$$\mathbb{E}[Y|A = 1] - \lambda \mathbb{E}[h|A = 1, K] = 100K - \lambda \left[ \frac{|100(K - 1)|}{\sqrt{2\pi}} e^{-\frac{1}{2}} + \frac{100(K - 1)}{2} \left( \text{erf} \left( \frac{\text{sign}(K - 1)}{\sqrt{2}} \right) - 1 \right) \right] \quad (22)$$

$$= \begin{cases} 100K - 8.332(K - 1)\lambda, & K \geq 1 \\ 100K - 59.937(1 - K)\lambda, & K < 1 \end{cases} \quad (23)$$

$$E[Y|A = 2] - \lambda \mathbb{E}[h|A = 2] = 110 \quad (24)$$

$$E[Y|A = 3] - \lambda \mathbb{E}[h|A = 3] = 80 - \lambda \left[ \frac{100}{\sqrt{2\pi}} e^{-\frac{20^2}{2 \times 100^2}} + \frac{20}{2} \left( \text{erf} \left( \frac{20}{\sqrt{2} \times 100} \right) - 1 \right) \right] \quad (25)$$

967 Clearly, the agent will never take action 3 as its expected HPU is smaller than that for action 2 for  
 968 all  $\lambda$ . For action 1, for  $K < 1$  the expected HPU is also smaller than that for action 2, for all  $\lambda$ . For  
 969 action 1 with  $K > 1$ , if  $\lambda < 12.002$  the optimal  $K = 20$ , otherwise it is 0. As a result, for  $\lambda < 11.93$   
 970 the agent chooses action 1 with  $K = 20$ , and otherwise chooses action 2.

## 971 H

972 In this Appendix we derive an expression for the expected counterfactual harm in generalized additive  
 973 models. To calculate the expected counterfactual harm we derive a solution for a broad class of  
 974 SCMs, heteroskedastic additive noise models, which includes our GAM (11),

975 **Definition 12** (Heteroskedastic additive noise models). *For  $Y$ ,  $Pa(Y) = A \cup X$ , the mechanism*  
 976  *$y = f_Y(a, x)$  is a heteroskedastic additive noise model if  $Y$  is normally distributed with a mean and*  
 977 *variance that are functions of  $a, x$ ,*

$$y = \mu(a, x) + e^Y \sigma(a, x), \quad e^Y \sim \mathcal{N}(0, 1) \quad (26)$$

978 In Appendix 11 we show that the dose response model (11) can be parameterised as a heteroskedastic  
 979 additive noise model and calculate the expected counterfactual harm using the following theorem,

980 **Theorem 5** (Expected harm for heteroskedastic additive noise model). For  $Y = f_Y(a, x, e^Y)$  where  
 981  $f_Y$  is a heteroskedastic additive noise model (Definition 12) and default action  $A = \bar{a}$ , the expected  
 982 harm is

$$\mathbb{E}[h|a, x] = \frac{|\Delta\sigma|}{\sqrt{2\pi}} e^{-\frac{\Delta U^2}{2\Delta\sigma^2}} + \frac{\Delta U}{2} \left( \operatorname{erf} \left( \frac{\Delta U}{\sqrt{2}|\Delta\sigma|} \right) - 1 \right) \quad (27)$$

983 where  $\operatorname{erf}(\cdot)$  is the error function,  $\Delta U = \mathbb{E}[U|a, x] - \mathbb{E}[U|\bar{a}, x]$ ,  $\Delta\sigma = \sigma(a, x) - \sigma(\bar{a}, x)$ .

984 *Proof.* Note that if  $e^Y \sim \mathcal{N}(\mu, V)$  we can replace  $e^Y \rightarrow e'^Y = e^Y/\sqrt{V} - \mu$  and absorb these terms  
 985 into  $f(a, x)$  and  $\sigma(a, x)$ . Hence we need only consider zero-mean univariate noise. In the following  
 986 we use  $e^Y = \varepsilon \sim \mathcal{N}(0, 1)$  to denote the fact the the exogenous noise term is univariate normally  
 987 distributed. We also use the fact that there are no unobserved confounders between  $A$  and  $Y$  to give  
 988  $P(y|a, x) = P(y_a|x)$ . Calculating the expected counterfactual harm using gives

$$\mathbb{E}[h|a, x] = \int_y dy \int_{y^*} dy^* P(y, Y_{\bar{a}} = y^*, |a, x) \max(0, U(\bar{a}, x, y^*) - U(a, x, y)) \quad (28)$$

$$= \int_y dy \int_{y^*} dy^* P(Y_a = y, Y_{\bar{a}} = y^* | x) \max(0, U(\bar{a}, x, y^*) - U(a, x, y)) \quad (29)$$

$$= \int_{\varepsilon} P(\varepsilon) d\varepsilon \int_y dy \int_{y^*} dy^* P(Y_a = y, Y_{\bar{a}} = y^* | \varepsilon, a, x) \max(0, U(\bar{a}, x, y^*) - U(a, x, y)) \quad (30)$$

$$= \int_{\varepsilon} P(\varepsilon) d\varepsilon \int_y dy \int_{y^*} dy^* P(Y_a = y | \varepsilon, a, x) P(Y_{\bar{a}} = y^*, | \varepsilon, a, x) \max(0, U(\bar{a}, x, y^*) - U(a, x, y)) \quad (31)$$

989 Substituting in  $U(a, x, y) = y$  and  $P(y|\varepsilon, a, x) = \delta(y - f(a, x) - \varepsilon\sigma(a, x))$  gives,

$$\mathbb{E}[h|a, x] = \int d\varepsilon P(\varepsilon) \max\{0, f(\bar{a}, x) - f(a, x) + \varepsilon(\sigma(\bar{a}, x) - \sigma(a, x))\} \quad (32)$$

$$= \int d\varepsilon P(\varepsilon) \max(0, -(\mathbb{E}[U|a, x] - \mathbb{E}[U|\bar{a}, x]) - \varepsilon(\sigma(a, x) - \sigma(\bar{a}, x))) \quad (33)$$

990 where we have used the fact that  $\mathbb{E}[U|a, x] = \int d\varepsilon P(\varepsilon) (f(a, x) + \varepsilon\sigma(a, x)) = f(a, x)$ . For ease  
 991 of notation we use  $\Delta U = E[U|a, x] - E[U|\bar{a}, x]$ ,  $\Delta\sigma = \sigma(a, x) - \sigma(\bar{a}, x)$ . Next, we remove the  
 992  $\max()$  by incorporating it into the bounds for the integral. If  $\Delta U > 0$  and  $\Delta\sigma > 0$ , this is equivalent  
 993 to  $\varepsilon < -\Delta U/\Delta\sigma$  and hence,

$$\mathbb{E}[h|a, x] = \int_{\varepsilon < -\Delta U/\Delta\sigma} d\varepsilon P(\varepsilon) (-\Delta U - \varepsilon\Delta\sigma) \quad (34)$$

$$= -\Delta U \int_{-\infty}^{-\Delta U/\Delta\sigma} P(\varepsilon) d\varepsilon - \Delta\sigma \int_{-\infty}^{-\Delta U/\Delta\sigma} \varepsilon P(\varepsilon) d\varepsilon \quad (35)$$

$$(36)$$

994 Using the standard Gaussian integrals

$$\int_a^b P(\varepsilon) d\varepsilon = \frac{1}{2} \left[ \operatorname{erf} \left( \frac{b}{\sqrt{2}} \right) - \operatorname{erf} \left( \frac{a}{\sqrt{2}} \right) \right] \quad (37)$$

$$\int_a^b \varepsilon P(\varepsilon) d\varepsilon = P(a) - P(b) \quad (38)$$

995 where  $P(\varepsilon) = e^{-\varepsilon^2/2}/\sqrt{2\pi}$  and  $\text{erf}(z)$  is the error function, we recover

$$\mathbb{E}[h|a, x] = \frac{-\Delta U}{2} \left[ \text{erf}\left(\frac{-\Delta U}{\sqrt{2}\Delta\sigma}\right) - \text{erf}(-\infty) \right] - \Delta\sigma [P(-\infty) - P(-\Delta U/\Delta\sigma)] \quad (39)$$

$$= \frac{\Delta U}{2} \left[ \text{erf}\left(\frac{\Delta U}{\sqrt{2}\Delta\sigma}\right) - 1 \right] + \frac{\Delta\sigma}{\sqrt{2\pi}} e^{-\frac{\Delta U^2}{2\Delta\sigma^2}} \quad (40)$$

996 where we have used  $\text{erf}(-z) = -\text{erf}(z)$  and  $P(-z) = P(z)$ . Similarly, if  $\Delta U > 0$ ,  $\Delta\sigma < 0$  then  
 997 the  $\max()$  in (33) can be replaced with a definite intergral over  $\varepsilon > \Delta U/\Delta\sigma$  giving,

$$\mathbb{E}[h|a, x] = \int_{\varepsilon > \Delta U/\Delta\sigma}^{\infty} d\varepsilon P(\varepsilon) (-\Delta U - \varepsilon\Delta\sigma) \quad (41)$$

$$= -\Delta U \int_{-\Delta U/\Delta\sigma}^{\infty} P(\varepsilon) d\varepsilon - \Delta\sigma \int_{-\Delta U/\Delta\sigma}^{\infty} \varepsilon P(\varepsilon) d\varepsilon \quad (42)$$

$$= -\frac{\Delta U}{2} \left[ \text{erf}(\infty) - \text{erf}\left(\frac{-\Delta U}{\sqrt{2}\Delta\sigma}\right) \right] - \Delta\sigma \left[ P\left(\frac{-\Delta U}{\sqrt{2}\Delta\sigma}\right) - P(\infty) \right] \quad (43)$$

$$= \frac{\Delta U}{2} \left[ \text{erf}\left(\frac{\Delta U}{\sqrt{2}|\Delta\sigma|}\right) - 1 \right] + \frac{|\Delta\sigma|}{\sqrt{2\pi}} e^{-\frac{\Delta U^2}{2\Delta\sigma^2}} \quad (44)$$

998 Next, if  $\Delta U < 0$  and  $\Delta\sigma > 0$  we recover the same integral as (35), and if  $\Delta U < 0$  and  $\Delta\sigma < 0$  we  
 999 recover the same integral as (41). Hence the general solution for all  $\Delta\sigma$  is (44).

1000

□

## 1001 I

1002 In this Appendix we present the GAM dose response model including parameter values, and show  
 1003 that it corresponds to a heteroskedastic additive noise model and calculate the expected harm for a  
 1004 given dose.

1005 We follow the set-up described in (18), where outcome  $Y$  denotes the level of improvement in the  
 1006 symptoms of schizoaffective patients following treatment and compared to pre-treatment levels,  
 1007 measured in terms of the Positive and Negative Syndrome Scale (PANSS) (44). The response of  $Y$   
 1008 w.r.t dose  $A$  (Aripiprazole mg/day) is determined using a generalized additive model fit with a cubic  
 1009 splines regression and random effects,

$$y = \theta_1 a + \theta_2 f(a) + \varepsilon_0 \quad (45)$$

1010 where the parameters  $\theta_i$  are random variables  $\theta_i \sim \mathcal{N}(\hat{\theta}_i, V_i)$ ,  $\varepsilon_0 \sim \mathcal{N}(0, V_0)$  is the sample noise,  
 1011 and the spline function  $f(a)$  is given by,

$$f(a) = \frac{(a - k_1)_+^3 - \frac{k_3 - k_1}{k_3 - k_2} (a - k_2)_+^3 + \frac{k_2 - k_1}{k_3 - k_2} (a - k_3)_+^3}{(k_3 - k_1)^2} \quad (46)$$

1012 where  $k_1, k_2, k_3$  are the knots at  $a = 0, 10$  and  $30$  respectively, with  $(u)_+ = \max\{0, u\}$ . In the  
 1013 following we assume for simplicity that  $\theta_1$  and  $\theta_2$  are independent. This hierarchical model can be  
 1014 expressed as an SCM with the mechanism for  $Y$  given by,

$$y = \left( \hat{\theta}_1 a + \hat{\theta}_2 f(a) \right) + \varepsilon_1 a + \varepsilon_2 f(a) + \varepsilon_0 \quad (47)$$

1015 where  $\varepsilon_i \sim \mathcal{N}(0, V_i)$ . We will now reparameterise this as an equivalent SCM that is an additive  
 1016 heteroskedastic noise model. Using the identifies  $Z = kY$ ,  $Y \sim \mathcal{N}(0, 1) \implies Z \sim \mathcal{N}(0, k^2)$ ,

1017 and  $Z = X + Y$ ,  $X \sim \mathcal{N}(0, V_X)$ ,  $Y \sim \mathcal{N}(0, V_Y) \implies Z \sim \mathcal{N}(0, V_X + V_Y)$  (where  $V_X$  is the  
1018 variance of  $X$  and likewise for  $V_Y, Y$ ), we can replace  $\varepsilon_1 a + \varepsilon_2 f(a) \rightarrow \varepsilon g(a)$  where  $\varepsilon \sim \mathcal{N}(0, 1)$   
1019 and  $g(a) = \sqrt{a^2 V_1 + f(a)^2 V_2}$ . We can therefore reparameterise the mechanism for  $Y$  as

$$y = \mathbb{E}[U|a] + g(a)\varepsilon + \varepsilon_0 \quad (48)$$

1020 where we have used  $U(a, x, y) = U(a, y) = y$  and the fact that  $\varepsilon, \varepsilon_0$  are mean zero to give  
1021  $\mathbb{E}[U|a] = \theta_1 a + \theta_2 f(a)$ . Finally, we note that the sample noise term  $\varepsilon_0$  cancels in the expression for  
1022 the harm,

$$\mathbb{E}[h|a] = \int_y dy \int_{y^*} dy^* P(y, Y_{\bar{a}} = y^*, |a) \max(0, U(\bar{a}, y^*) - U(a, y)) \quad (49)$$

$$= \int_y dy \int_{y^*} dy^* P(Y_a = y, Y_{\bar{a}} = y^*) \max(0, U(\bar{a}, y^*) - U(a, y)) \quad (50)$$

$$= \int_{\varepsilon} P(\varepsilon) d\varepsilon \int_{\varepsilon_0} P(\varepsilon_0) d\varepsilon_0 \int_y dy \int_{y^*} dy^* P(Y_a = y, Y_{\bar{a}} = y^* | \varepsilon, \varepsilon_0, a) \max(0, U(\bar{a}, y_{\bar{a}}^*) - U(a, y)) \quad (51)$$

1023

$$= \int_{\varepsilon} P(\varepsilon) d\varepsilon \int_{\varepsilon_0} P(\varepsilon_0) d\varepsilon_0 \int_y dy \int_{y^*} dy^* P(y, | \varepsilon, \varepsilon_0, a) P(Y_{\bar{a}} = y^*, | \varepsilon, \varepsilon_0) \max(0, U(\bar{a}, y^*) - U(a, y)) \quad (52)$$

1024 Substituting in  $P(Y_a = y | \varepsilon, \varepsilon_0) = \delta(y - f(a) + g(a)\varepsilon + \varepsilon_0)$  gives,

$$\mathbb{E}[h|a] = \int_{\varepsilon} P(\varepsilon) d\varepsilon \int_{\varepsilon_0} P(\varepsilon_0) d\varepsilon_0 \max(0, f(\bar{a}) + g(\bar{a})\varepsilon + \varepsilon_0 - f(a) - g(a)\varepsilon - \varepsilon_0) \quad (53)$$

$$= \int_{\varepsilon} P(\varepsilon) d\varepsilon \int_{\varepsilon_0} P(\varepsilon_0) d\varepsilon_0 \max(0, f(\bar{a}) - f(a) + (g(\bar{a}) - g(a))\varepsilon) \quad (54)$$

$$= \int_{\varepsilon} P(\varepsilon) d\varepsilon \max(0, f(\bar{a}) - f(a) + (g(\bar{a}) - g(a))\varepsilon) \quad (55)$$

1025 Therefore we can ignore the sample noise term when calculating the expected harm, instead calcu-  
1026 lating the expected harm for the model  $Y = f(a) + g(a)\varepsilon$ . This is a heteroskedastic additive noise  
1027 model, and therefore by Theorem 5 the expected harm is,

$$\mathbb{E}[h|a] = \frac{\Delta U}{2} \left[ \operatorname{erf} \left( \frac{\Delta U}{\sqrt{2}\Delta\sigma} \right) - 1 \right] + \frac{\Delta\sigma}{\sqrt{2\pi}} e^{-\Delta U^2/2\Delta\sigma^2} \quad (56)$$

1028 where  $\Delta U = \mathbb{E}[U|a] - \mathbb{E}[U|\bar{a}]$ ,  $\Delta\sigma = g(a) - g(\bar{a})$  and  $g(a) = \sqrt{a^2 V_1 + f(a)^2 V_2}$

1029 The resulting curves presented in Figure 2 are calculated using (56) and the parameter values  
1030 taken from [18] (Table 1), which are fitted in a meta-analysis of the dose-responses reported in  
1031 [19, 43, 57, 70, 86].

## 1032 J

1033 In this Appendix we present proofs of Theorems 2, 3 and 4. First, we prove Theorem 2.

1034

1035 **Theorem 2:** For any utility functions  $U$ , environment  $\mathcal{M}$  and default action  $A = \bar{a}$  the expected  
1036 HPU is never a harmful objective for  $\lambda > 0$ .

Table 1: Parameters for the hierarchical generalized additive dose-response model reported in [18]

Parameter	Value
$\hat{\theta}_1$	0.937
$\hat{\theta}_2$	-1.156
$V_1$	0.03
$V_2$	0.10

1037 *Proof.* Let  $a_{\max} = \arg \max_a \{\mathbb{E}[U|a, x] - \lambda \mathbb{E}[h|a, x; \mathcal{M}]\}$ . If  $\exists a' \neq a_{\max}$  such that  $\mathbb{E}[U|a', x] \geq$   
1038  $\mathbb{E}[U|a_{\max}, x]$  and  $\mathbb{E}[h|a', x; \mathcal{M}] < \mathbb{E}[h|a_{\max}, x; \mathcal{M}]$ , then  $\mathbb{E}[U|a_{\max}, x] + \lambda \mathbb{E}[h|a_{\max}, x; \mathcal{M}] <$   
1039  $\mathbb{E}[U|a', x] + \lambda \mathbb{E}[h|a', x; \mathcal{M}] \forall \lambda > 0$  and so  $a_{\max} \neq \arg \max_a \{\mathbb{E}[U|a, x] - \lambda \mathbb{E}[h|a, x; \mathcal{M}]\}$ .  $\square$

1040 Next, we prove theorems [3] and [4] by example, constructing distributional shifts that reveal if an  
1041 objective function is harmful. To do this we make use of a specific family of structural causal  
1042 models—counterfactually independent models.

1043 **Definition 13** (counterfactual independence (CFI)). *Y is counterfactually independent in with respect*  
1044 *to A in M if,*

$$P(y_{a^*}^*, y_a | x) = \begin{cases} P(y_a | x) \delta(y_a - y_{a^*}^*) & a = a^* \\ P(y_{a^*}^* | x) P(y_a | x) & \text{otherwise} \end{cases} \quad (57)$$

1045 Counterfactually independent models (CFI models) are those for which the outcome  $Y_a$  is independent  
1046 to any counterfactual outcome  $Y_{a'}$ . Next we show that there is always a CFI model that can induce  
1047 any factual outcome statistics.

1048 **Lemma 1.** *For any desired outcome distribution  $P(y|a, x)$  there is a choice of exogenous noise*  
1049 *distribution  $P(e^Y)$  and causal mechanism  $f_Y(a, x, e^Y)$  such that Y is counterfactually independent*  
1050 *with respect to A*

1051 *Proof.* Consider the causal mechanism  $y = f_Y(a, x, e^Y)$  for some fixed  $X = x$ , and exogenous  
1052 noise distribution  $P(E^Y = e^Y)$ . Let the noise term be described by the random field  $E^Y =$   
1053  $\{E^Y(a, x) : a \in A, x \in X\}$ , with  $P(E^Y = e^y) = \times_{a \in A, x \in X} P(E^y(a, x) = e^y(a, x))$  and with  
1054  $\text{dom}(E^Y(a, x)) = \text{dom}(Y) \forall A = a, X = x$ . I.e. we choose the noise distribution to be joint state  
1055 over mutually independent noise variables, one for every action  $A = a$  and context  $X = x$ , and  
1056 where each of these variables has the same domain as  $Y$ . Next, we choose the causal mechanism,

$$f_Y(a, x, e^Y) = e^Y(a, x) \quad (58)$$

1057 i.e. the value of  $Y$  for action  $A = a$  and context  $X = x$  is the state of the independent noise variable  
1058  $E^Y(a, x)$ . By construction this is a valid SCM, and we note that the factual distributions (calculated  
1059 with [4]) are given simply by,

$$P(y|a, x) = P(E^Y(a, x) = y) \quad (59)$$

1060 Likewise applying our choice of mechanism and noise distribution to [4] gives (for  $a \neq a'$ ) the  
1061 counterfactual distribution,

$$P(Y_a = y, Y_{a'} = y' | x) = P(E^Y(a, x) = y) P(E^Y(a', x) = y') \quad (60)$$

$$= P(Y_a = y | x) P(Y_{a'} = y' | x) \quad (61)$$

1062 and likewise gives  $P(y_a | x) \delta(y_a - y'_{a'})$  for  $a = a'$ . Finally, we note that we can choose any  
1063  $P(y_a | x) = P(E^Y(a, x) = y)$ , hence there is a CFI model that induces any factual outcome  
1064 distribution we desire.  $\square$

1065 Next, we show that in counterfactually independent models there are outcome distributional shifts  
1066 that only change the expected harm of individual actions, without changing any other factual or  
1067 counterfactual statistics.



1068 **Lemma 2.** For  $\mathcal{M}$  and (context-dependent) default action  $A = \bar{a}(x)$ , if  $U$  is outcome dependent for  
 1069 the default action  $\bar{a}(x)$  and some other action  $a \neq \bar{a}(x)$ , then there are three outcome distributionally  
 1070 shifted environments  $\mathcal{M}_0$ ,  $\mathcal{M}_+$  and  $\mathcal{M}_-$  such that;

- 1071 1.  $\mathbb{E}[h|a, x; \mathcal{M}_-] < \mathbb{E}[h|a, x; \mathcal{M}_0] < \mathbb{E}[h|a, x; \mathcal{M}_+]$   
 1072 2.  $\mathbb{E}[h|b, x; \mathcal{M}_-] = \mathbb{E}[h|b, x; \mathcal{M}_0] = \mathbb{E}[h|b, x; \mathcal{M}_+] \forall b \neq a$   
 1073 3.  $P(y|a', x; \mathcal{M}_0) = P(y|a', x; \mathcal{M}_+) = P(y|a', x; \mathcal{M}_-) \forall a' \in A$ , including  $a, \bar{a}(x)$

1074 *Proof.* In the following we suppress the notation  $\bar{a}(x) = \bar{a}$ . To construct the environment  $\mathcal{M}_0$  we  
 1075 restrict to a binary outcome distribution for each action such that  $P(y_a|x)$  is completely concentrated  
 1076 on the highest and lowest utility outcomes,

$$Y_a = 1 \implies Y_a = \arg \max_y U(a, x, y) \quad (62)$$

$$Y_a = 0 \implies Y_a = \arg \min_y U(a, x, y) \quad (63)$$

$$1 = P(Y_a = 1|x; \mathcal{M}_0) + P(Y_a = 0|x; \mathcal{M}_0) \quad (64)$$

1077 Note that we abuse notation as the variables  $Y_a = 1$  and  $Y_b = 1$  will not be in the same state  
 1078 in general, and the states 1, 0 denote the max/min utility states under any given action, rather  
 1079 than a fixed state of  $Y$ . By Lemma 1 we can choose  $Y_a$  to be counterfactually independent with  
 1080 respect to  $A$ . Recalling our parameterization of CFI models in Lemma 1 with noise distribution  
 1081  $P(E^Y = e^Y) = \times_{a \in A, x \in X} P(E^Y(a, x) = e^Y(a, x))$ ,  $\text{dom}(E^Y(a, x)) = \text{dom}(Y)$ , and causal  
 1082 mechanism  $f_Y(a, x, e^Y) = e^Y(a, x)$ , therefore  $E^Y(a, x) \in \{0, 1\} \forall a, x$ . The expected harm for  
 1083 action  $\text{do}(A = a)$  is,

$$\mathbb{E}[h|a, x; \mathcal{M}_0] = \sum_{y_a=0}^1 \sum_{y_{\bar{a}}=0}^1 P(y_{\bar{a}}|x) P(y_a|x) \max\{0, U(\bar{a}, x, y_{\bar{a}}) - U(a, x, y_a)\} \quad (65)$$

1084 where we have used the fact that  $P(y_{\bar{a}}^*, y|a, x) = P(y_{\bar{a}}^*, y_a|x)$  and used counterfactual independence.  
 1085  $U(a, x, 0) < U(\bar{a}, x, 1)$  and so if we choose non-deterministic outcome distributions for  $P(y_a|x)$   
 1086 and  $P(y_{\bar{a}}|x)$  then (65) is strictly greater than 0.

1087 We can construct the desired  $\mathcal{M}_{\pm}$  by keeping the causal mechanism but changing the factorized  
 1088 exogenous noise distribution in  $\mathcal{M}$  to be,

$$P'(E^Y = e^Y; \mathcal{M}_+) = P(E^Y = e^Y; \mathcal{M}_0) + (-1)^{e^Y(a, x) - e^Y(\bar{a}, x)} \phi_+ \quad (66)$$

$$P'(E^Y = e^Y; \mathcal{M}_-) = P(E^Y = e^Y; \mathcal{M}_0) + (-1)^{e^Y(a, x) - e^Y(\bar{a}, x)} \phi_- \quad (67)$$

1089 where  $\phi_{\pm} \in \mathbb{R}$  are constants that satisfy the bounds  $\max\{-P(Y_{\bar{a}} = 1|x)P(Y_a = 1|x), -P(Y_{\bar{a}} =$   
 1090  $0|x)P(Y_a = 0|x)\} \leq \phi_{\pm} \leq \min\{P(Y_{\bar{a}} = 1|x)P(Y_a = 0|x), P(Y_{\bar{a}} = 0|x)P(Y_a = 1|x)\}$ . It is  
 1091 simple to check that for any  $\phi$  that satisfies these bounds we recover  $\sum_{e^Y} P'(E^Y = e^Y) = 1$ ,  
 1092  $P'(E^Y = e^Y) \geq 0 \forall e^Y$ , and therefore  $P'$  is a valid noise distribution. Keeping the same causal  
 1093 mechanism  $f_Y$  is  $\mathcal{M}_{\pm}$  as in  $\mathcal{M}_0$  gives  $P(y_a|x; \mathcal{M}_0) = P(y_a|x; \mathcal{M}_+) = P(y_a|x; \mathcal{M}_-)$  as,

$$P'(y_i|x) = \sum_{e^Y(0, x)=0}^1 \dots \sum_{e^Y(i-1, x)=0}^1 \sum_{e^Y(i+1, x)=0}^1 \dots \sum_{e^Y(|A|, x)=0}^1 \left[ \prod_{j=1}^{|A|} P(e^Y(j, x)) + (-1)^{e^Y(i, x) - e^Y(\bar{a}, x)} \phi_{\pm} \right] \quad (68)$$

$$= P(e^Y(i, x)) + (-1)^{e^Y(i, x) - 0} \phi_{\pm} + (-1)^{e^Y(i, x) - 1} \phi_{\pm} \quad (69)$$

$$= P(e^Y(i, x)) = P(y_i|x) \quad (70)$$

1094 and likewise for  $i = \bar{a}$ . This implies that for any desired outcome statistics  $P(y_a|x)$  there is a model  
 1095 where  $Y_a \perp Y_{a'} \forall (a, a')$  where  $a \neq a'$  except for the pair  $a, \bar{a}$ , so long as  $P(y_{\bar{a}}|x)$  and  $P(y_a|x)$  are  
 1096 non-deterministic (if they are deterministic,  $\phi_{\pm} = 0$  and  $\mathcal{M}_0 = \mathcal{M}_{\pm}$ ). Because  $Y_{a'} \perp Y_{\bar{a}} \forall a' \neq a$ ,  
 1097 then  $H(a', x; \mathcal{M}_0) = H(a', x; \mathcal{M}_{\pm}) \forall a' \neq a$ . Also note that  $H(\bar{a}, x; \mathcal{M}) = 0$  for any  $U$  or  $\mathcal{M}$  if

1098  $P(a|x) = \delta(a - \bar{a})$ , as  $P(Y_{\bar{a}} = i, Y_a = k) = 0$  if  $i \neq k$  and if  $i = k$  (factual and counterfactual  
1099 outcomes are identical) then the expected harm is zero. The only difference between  $\mathcal{M}_0$  and  $\mathcal{M}_{\pm}$  is  
1100  $P(y_{\bar{a}}, y_a|x; \mathcal{M}_+) \neq P(y_{\bar{a}}, y_a|x; \mathcal{M}_-) \neq P(y_{\bar{a}}, y_a|x; \mathcal{M}_0)$ , which differ for  $\phi_+ \neq 0$ ,  $\phi_- \neq 0$  and  
1101  $\phi_+ \neq \phi_-$ . Substituting (66) and (67) into our expression for the expected harm as using the notation  
1102  $\Delta_{y,y'} = \max\{0, U(\bar{a}, x, y) - U(a, x, y')\}$  gives,

$$\mathbb{E}[h|a, x; \mathcal{M}_{\pm}] = \mathbb{E}[h|a, x; \mathcal{M}_0] + \phi_{\pm} [\Delta_{00} + \Delta_{11} - \Delta_{10} - \Delta_{01}] \quad (71)$$

$$\mathbb{E}[h|a', x; \mathcal{M}_{\pm}] = \mathbb{E}[h|a', x; \mathcal{M}_0], \quad a' \neq a \quad (72)$$

1103 Now, as  $\max_y U(a, x, y) > \min_y U(\bar{a}, x, y)$  then  $\Delta_{01} = 0$ . For the coefficient of  $\phi_{\pm}$  in (71) to be  
1104 zero, we would therefore require that  $\Delta_{00} + \Delta_{11} = \Delta_{10}$ . We know  $\Delta_{10} > 0$  because otherwise  
1105  $\min_y U(a, x, y) > \max_y U(\bar{a}, x, y)$ , therefore the minimal value of  $\Delta_{10}$  is  $\max_y U(\bar{a}, x, 1) -$   
1106  $\min_y U(a, x, y)$ . If  $\Delta_{00} \neq 0$  and  $\Delta_{11} \neq 0$  then  $\Delta_{00} + \Delta_{11} \geq \Delta_{10}$  implies  $\min_y U(\bar{a}, x, y) \geq$   
1107  $\max_y U(a, x, y)$  which violates our assumptions, therefore  $\Delta_{00} + \Delta_{11} < \Delta_{10}$ . If  $\Delta_{00} = 0$  clearly  
1108 we cannot have  $\Delta_{11} = \Delta_{10}$  as  $\min_y U(a, x, y) < \max_y U(\bar{a}, x, y)$  by our assumptions, and likewise  
1109 if  $\Delta_{11} = 0$  we cannot have  $\Delta_{00} = \Delta_{10}$  as this would imply  $\min_y U(\bar{a}, x, y) = \max_y U(\bar{a}, x, y)$   
1110 which violates our assumptions. Therefore we can conclude that the coefficient in (71) is greater than  
1111 zero.

1112 Therefore if we choose any  $0 < \phi_+ < \min\{P(Y_{\bar{a}} = 1|x)P(Y_a = 0|x), P(Y_{\bar{a}} = 0|x)P(Y_a =$   
1113  $1|x)\}$  we get  $\mathbb{E}[h|a, x; \mathcal{M}_+] > \mathbb{E}[h|a, x; \mathcal{M}_0]$ , and any  $\max\{P(Y_{\bar{a}} = 1|x)P(Y_a = 1|x), P(Y_{\bar{a}} =$   
1114  $0|x)P(Y_a = 0|x)\} < \phi_- < 0$ , we get  $\mathbb{E}[h|a, x; \mathcal{M}_-] < \mathbb{E}[h|a, x; \mathcal{M}_0]$ .  $\square$

1115 **Lemma 3.** For (context dependent) default action  $A = \bar{a}(x)$ ,  $\mathbb{E}[h|\bar{a}(x), x; \mathcal{M}] = 0 \forall \mathcal{M}$

1116 *Proof.* In the following we suppress the notation  $\bar{a}(x) = \bar{a}$ .

$$\mathbb{E}[h|\bar{a}, x; \mathcal{M}] = \int_{y^*, y} P(Y_{\bar{a}} = y^*, Y = y|\bar{a}, x; \mathcal{M}) \max\{0, U(\bar{a}, x, y^*) - U(\bar{a}, x, y)\} \quad (73)$$

$$= \int_{y^*, y} P(Y_{\bar{a}} = y^*, Y_{\bar{a}} = y|x; \mathcal{M}) \max\{0, U(\bar{a}, x, y^*) - U(\bar{a}, x, y)\} \quad (74)$$

$$= \int_{y^*, y} P(Y_{\bar{a}} = y|x; \mathcal{M}) \delta(y^* - y) \max\{0, U(\bar{a}, x, y^*) - U(\bar{a}, x, y)\} \quad (75)$$

$$= \int_y P(Y_{\bar{a}} = y|x; \mathcal{M}) \max\{0, U(\bar{a}, x, y) - U(\bar{a}, x, y)\} \quad (76)$$

$$= 0 \quad (77)$$

1117  $P(y|a, x) = P(y_a|x)$ .

1118  $\square$

1119 **Theorem 3:** For any (context dependent) default action  $A = \bar{a}(x)$ , if there is a context  $X = x$  where  
1120 the user's utility function is outcome dependent for  $\bar{a}(x)$  and some other action  $a \neq \bar{a}(x)$ , then there  
1121 is an outcome distributional shift such that  $U$  is harmful in the shifted environment.

1122 *Proof.* For the expected utility to not be harmful by Definition 6, it must be that  $\mathbb{E}[h|a, x] > \mathbb{E}[h|b, x]$   
1123  $\implies \mathbb{E}[U|a, x] < \mathbb{E}[U|b, x]$ . Given our assumption of outcome dependence, we know there  
1124 is a context  $X = x$  such that the utility functions for  $\bar{a}(x)$  and  $a \neq \bar{a}(x)$  overlap, that is  
1125  $\min_y U(a, x, y) < \max_y U(\bar{a}(x), x, y)$  and  $\max_y U(a, x, y) > \min_y U(\bar{a}(x), x, y)$ . In the fol-  
1126 lowing we drop the notation  $\bar{a}(x) = \bar{a}$ . We can restrict our agent to choose between these two  
1127 actions and construct an outcome distributional shift such that; i) The outcomes  $Y_a$  and  $Y_{\bar{a}}$  are  
1128 binary with one outcome maximizing the utility for that action and the other minimizing the utility,  
1129 i.e.  $Y_a \in \{\max_y U(a, x, y), \min_y U(a, x, y)\}$  and  $Y_{\bar{a}} \in \{\max_y U(\bar{a}, x, y), \min_y U(\bar{a}, x, y)\}$ , ii)  
1130  $\mathbb{E}[U|a, x] = \mathbb{E}[U|\bar{a}, x]$ , iii)  $P(y_a|x)$  and  $P(y_{\bar{a}}|x)$  are non-deterministic. This follows from the fact  
1131 that the set of possible expected utility values for an action  $a$  is the set of mixtures over  $U(a, x, y)$

1132 with respect to  $y$ , and as  $Y_a = 0, 1$  are the extremal points of this convex set, the expected utility for  
 1133 action  $a$  in context  $x$  can be written as  $P(Y_a = 0|x)U(a, x, 0) + P(Y_a = 1|x)U(a, x, 1)$ . Then, as  
 1134 the utility functions for  $a$  and  $\bar{a}$  overlap there is point in the intersection of these convex sets that is  
 1135 non-extremal (and hence, a non-deterministic mixture).

1136 By Lemma 3 the default action causes zero expected harm. By Lemma 2 we can construct a shifted  
 1137 environment  $\mathcal{M}_0$  where the non-default action  $a \neq \bar{a}$  has non-zero harm for any non-deterministic  
 1138  $P(y_a|x)$ . We can therefore construct  $\mathcal{M}_0$  such that i)  $\mathbb{E}[Y_a|x] = \mathbb{E}[Y_{\bar{a}}|x]$ , and ii)  $\mathbb{E}[h|a, x] >$   
 1139  $\mathbb{E}[h|\bar{a}, x]$ , violating our requirement that  $\mathbb{E}[h|a, x] > \mathbb{E}[h|b, x] \implies \mathbb{E}[U|a, x] < \mathbb{E}[U|b, x]$ .

1140 □

1141 **Theorem 4:** For any (context dependent) default action  $A = \bar{a}(x)$ , if there is a context  $X = x$  where  
 1142 the user’s utility function is outcome dependent for  $\bar{a}(x)$  and two other actions  $a_1, a_2 \neq \bar{a}(x)$ , then  
 1143 for any factual objective function  $J$  there is an outcome distributional shift such that maximizing the  
 1144  $J$  is harmful in the shifted environment.

1145 *Proof.* By assumption there is a context  $X = x$  for which the utility functions for  $a_1, a_2$  and  $\bar{a}(x)$   
 1146 overlap. In the following we drop the notation  $\bar{a}(x) = \bar{a}$ . There is a choice of non-deterministic  
 1147 outcome distributions  $P(y_{\bar{a}}|x)$ ,  $P(y_{a_1}|x)$  and  $P(y_{a_2}|x)$  such that all three actions have the same  
 1148 expected utility. By Lemma 2 for any non-deterministic outcome distribution we can choose  $\mathcal{M}_0$   
 1149 such that  $\mathbb{E}[h|a_1, x; \mathcal{M}_0] > 0$ , and  $\mathbb{E}[h|a_2, x; \mathcal{M}_0] > 0$ , and by Lemma 3  $\mathbb{E}[h|\bar{a}, x; \mathcal{M}] = 0 \forall \mathcal{M}$ .  
 1150 Therefore  $\exists \mathcal{M}_0$  that is an outcome distributional shift of the original environment  $\mathcal{M}$  such that  
 1151  $\bar{a}, a_1, a_2$  have the same expected utility,  $\bar{a}$  has zero expected harm and  $a_1, a_2$  have non-zero expected  
 1152 harm.

1153 If  $\mathbb{E}[h|a_1, x; \mathcal{M}_0] = \mathbb{E}[h|a_2, x; \mathcal{M}_0]$  then by Lemma 2 there are outcome-shifted environments  $\mathcal{M}_{\pm}$   
 1154 such that  $\bar{a}, a_1$  and  $a_2$  have the same factual statistics as in  $\mathcal{M}_0$  and  $\mathbb{E}[h|a_2, x; \mathcal{M}_{\pm}] = \mathbb{E}[h|a_2, x; \mathcal{M}_0]$ ,  
 1155 but the harm caused by  $a_1$  is increased(decreased) by some non-zero amount. Therefore in  $\mathcal{M}_+$   
 1156  $a_1$  and  $a_2$  have the same expected utility but  $a$  has a strictly higher expected harm, and in order  
 1157 to be non-harmful it must be that  $\mathbb{E}[J|a_1, x; \mathcal{M}_+] < \mathbb{E}[J|a_2, x; \mathcal{M}_+]$ . Likewise in  $\mathcal{M}_-$   $a_1$  and  $a_2$   
 1158 have the same expected utility but the expected harm for  $a_1$  is strictly lower than for  $a_2$ , therefore  
 1159 in order to be non-harmful it must be that  $\mathbb{E}[J|a_1, x; \mathcal{M}_-] > \mathbb{E}[J|a_2, x; \mathcal{M}_-]$ . Finally we note that  
 1160  $\mathbb{E}[J|a, x; \mathcal{M}_+] = \mathbb{E}[J|a, x; \mathcal{M}_-] = \mathbb{E}[J|a, x; \mathcal{M}_0] \forall a \in A$  as the factual statistics are identical  
 1161 in  $\mathcal{M}_0, \mathcal{M}_{\pm}$ , i.e.  $P(y_a|x; \mathcal{M}_+) = P(y_a|x; \mathcal{M}_-) = P(y_a|x; \mathcal{M}_0)$ . Therefore any  $J$  must be  
 1162 harmful in either  $\mathcal{M}_+$  and  $\mathcal{M}_-$ , and therefore there is an outcome distributional shift  $\mathcal{M} \rightarrow \mathcal{M}_+$  or  
 1163  $\mathcal{M} \rightarrow \mathcal{M}_-$  such that  $J$  is harmful in the shifted environment.

1164 If  $\mathbb{E}[h|a_1, x; \mathcal{M}_0] \neq \mathbb{E}[h|a_2, x; \mathcal{M}_0]$ , assume without loss of generality that  $\mathbb{E}[h|a_1, x; \mathcal{M}_0] >$   
 1165  $\mathbb{E}[h|a_2, x; \mathcal{M}_0]$ . As  $\bar{a}, a_1$  and  $a_2$  have the equal expected utilities then so does any mixture of  
 1166 these actions, in  $\mathcal{M}_0$  and  $\mathcal{M}_{\pm}$ . Restrict the agent to choose between action  $a_2$  and a mixture  
 1167 of actions  $\bar{a}$  and  $a_1$ —i.e. a stochastic or ‘soft’ intervention [17, 66], which involves replacing  
 1168 the causal mechanism for  $A$  with a mixture  $\tau := q[A = a_1] + (1 - q)[A = \bar{a}]$  where  $q$   
 1169 is an independent binary noise term. By linearity the expected utility for this mixed action is  
 1170  $\mathbb{E}[U_{\tau}|x] = q\mathbb{E}[U_{a_1}|x] + (1 - q)\mathbb{E}[U_{\bar{a}}|x] = \mathbb{E}[U_{a_1}|x]$  as all three actions have the same expected utility,  
 1171 and has an expected harm  $\mathbb{E}[h|\tau, x; \mathcal{M}_0] = q\mathbb{E}[h|a_1, x; \mathcal{M}_0] + (1 - q)\mathbb{E}[h|\bar{a}, x; \mathcal{M}_0] = q\mathbb{E}[h|a_1, x; \mathcal{M}_0]$   
 1172 as  $\mathbb{E}[h|\bar{a}, x; \mathcal{M}] = 0 \forall \mathcal{M}$ . Therefore as  $\mathbb{E}[h|a_1, x; \mathcal{M}_0] > 0$  and  $\mathbb{E}[h|a_2, x; \mathcal{M}_0] > 0$  we can choose  
 1173  $p > 0$  such that  $\mathbb{E}[h|\tau, x; \mathcal{M}_0] = \mathbb{E}[h|a_2, x; \mathcal{M}_0]$ . Therefore in  $\mathcal{M}_0$ ,  $a_2$  and  $\tau$  have the same ex-  
 1174 pected harm and utility, and in  $\mathcal{M}_+$  they have the same expected utility but  $\tau$  is more harmful than  
 1175  $a_2$  as  $\mathbb{E}[h|a_1, x; \mathcal{M}_+] > \mathbb{E}[h|a_1, x; \mathcal{M}_0]$  and  $p > 0$ , and in  $\mathcal{M}_-$  they have the same expected utility  
 1176 but  $a_2$  is more harmful than  $\tau$ . As the factual statistics  $P(y_a|x)$  are identical for  $\mathcal{M}_0$  and  $\mathcal{M}_{\pm}$ , so is  
 1177 the value of any factual objective function across all three environments. Hence, any factual objective  
 1178 function must be harmful in either  $\mathcal{M}_+$  or  $\mathcal{M}_-$ . □

## 1179 K

1180 In this appendix we discuss related works; counterfactual fairness [48] and path-specific objectives  
 1181 [26], as well as discussing some deep learning implementations that are capable of supporting  
 1182 counterfactual inferences of the type used to estimate counterfactual harm. For the sake of generality  
 1183 our results are derived in the SCM framework, and so taken at face value they assume knowledge of

1184 the SCM for the data generating process. Often in complex domains we will not have access to the  
 1185 true SCM that describes the data generating process but some approximation. However, there have  
 1186 been several recent proposals for performing counterfactual inference using deep learning methods,  
 1187 with promising results in diverse complex domains including learning deep structural causal models  
 1188 for medical imaging [65], visual question answering [61], and vision-and-language navigation in  
 1189 robotics [64] and text generation [54]. These studies evidence that deep learning algorithms can learn  
 1190 to make good counterfactual inferences that can be used to support decision making. This is achieved  
 1191 often without perfect knowledge of the underlying SCM (one notable exception being when the  
 1192 environment is simulated). This is somewhat analogous to the fact that human decision making often  
 1193 utilizes counterfactual reasoning for various cognitive tasks [23] (for example, it is important for legal  
 1194 and ethical reasoning [50]). This is in spite of the fact that humans clearly do not having access to  
 1195 perfect structural causal models of their environments, but have to learn good enough approximations  
 1196 through heuristics and inductive biases. While it is known that counterfactuals cannot be identified  
 1197 from data alone [79] but are only defined up to a structural causal models of the environment, clearly  
 1198 humans [32] and increasingly AI systems are capable of learning good structural causal models of  
 1199 real-world environments and using these to make counterfactual inferences capable of guiding actions  
 1200 and reasoning about harm.

## 1201 K.1 Related work

1202 Counterfactual fairness deals with prediction tasks  $\hat{Y} : X \rightarrow Y$  where the desire is to have a predictor  
 1203  $\hat{Y}$  that is not unfairly influenced by a protected attribute  $A$  such as gender or race. Note  $A$  is a  
 1204 feature that typically cannot be intervened on, whereas in our setup  $A$  denotes an agent’s action.  
 1205 Counterfactual fairness quantifies this unfair influence causally, using the counterfactual constraint,

$$P(\hat{Y}_a = y | X = x, A = a) = P(\hat{Y}_{a'} = y | X = x, A = a) \quad \forall a' \in A, y \in \hat{Y} \quad (78)$$

1206 which states that the probability of predicting any given outcome should not be caused on average  
 1207 by the protected attribute  $A$ , where type causation is established using the counterfactual  $P(\hat{Y}_{a'} =$   
 1208  $y | X = x, A = a)$  which is the probability of  $\hat{Y}$  given  $A = a$  if  $A$  had been equal to  $a'$ . Note  
 1209 that the counterfactual in (78) does not deal with the joint statistics of the factual outcome  $\hat{Y}_a$  and  
 1210 the counterfactual outcome  $\hat{Y}_{a'}$ , as so is an example of type causality compared to harm which an  
 1211 example of actual causality [34]. Harm is conceptually distinct from fairness—for example, it is  
 1212 possible to apply a needlessly harmful action fairly—but the two measures can be used in tandem.  
 1213 For example, one could quantify if a action or decision was unfair, and whether or not the user was  
 1214 harmed due to this unfair action.

1215 Another perhaps more related use of counterfactual inference for ethical AI is path-specific objectives  
 1216 [26]. This work similarly refines expected utility theory in the CID framework to take into account  
 1217 the fact that we often want to maximize utility via specific causal pathways due to ethical constraints.  
 1218 For example we can consider a simple model where the agent’s action  $A$  influences user feedback  
 1219  $Y$  (and utility  $U(y)$ ) but also effects the users preferences  $H$  where  $A \rightarrow Y$ ,  $A \rightarrow H$  and  $H \rightarrow Y$ .  
 1220 To maximize utility without intentionally manipulating the user we must maximize along the causal  
 1221 pathway (1) :  $A \rightarrow Y$  without including contributions to the expected utility from the mediator  
 1222 pathway (2) :  $A \rightarrow H \rightarrow Y$ . This involves replacing the expected utility with its path-specific  
 1223 equivalent, much as our path-specific harm (Definition 9) generalizes our path-independent definition  
 1224 of harm (Definition 3). As such the path-specific expected utility is still agnostic to harm just as  
 1225 the expected utility is, although it could be combined with the path specific harm in [26] to give a  
 1226 path-specific variant of the HPU (Definition 4). This would allow for harm averse decision making  
 1227 where the necessary degree of harm-aversion  $\lambda_{(i)}$  differs depending on the causal path ( $i$ )—for  
 1228 example, if we desire agents that have a high aversion for being directly harmful, but a lower degree  
 1229 of harm-aversion for indirect harm mediated by the actions of other agents (as described in Appendix  
 1230 B).