**Example 1.** *Two assassins ($Assassin_1$ and $Assassin_2$), in place as snipers, shoot and kill Victim, with each of the bullets fatally piercing Victim's heart at exactly the same moment. Although neither of them could have prevented the outcome, each of them is clearly responsible for Victim's death.*

Let $V$ stand for Victim's death ($V = 1$) or survival ($V = 0$), and let $A_1, A_2$ stand for the actions of the two assassins, where $A_i = 1$ if and only if $Assassin_i$ shoots. We can then capture this example with the single equation $V = A_1 \lor A_2$, and a context $\vec{u}$ such that $A_1 = 1$ and $A_2 = 1$.

Does the BvH definition (Definition 5) succeed in establishing that each of the assassins is responsible for Victim's death? To find out, we first need to evaluate whether $A_1 = 1$ (resp. $A_2 = 1$) directly NESS-causes $V = 1$ (Def. 4). We can choose $\vec{W} = \varnothing$ to get the desired result, as follows.

AC1 is fulfilled because $A_1 = 1$ and $V = 1$ actually happened. AC2 is established by verifying that the following two claims hold: $(M, \vec{u}) \vDash [A_1 \leftarrow 1, A_2 \leftarrow 1]V = 1$ and $(M, \vec{u}) \vDash [A_1 \leftarrow 1, A_2 \leftarrow 0]V = 1$. Since $\vec{W} = \varnothing$, verifying AC3 is easy: we need to find a single intervention on the variables other than $V$ such that they result in $V = 0$. The intervention $[A_1 \leftarrow 0, A_2 \leftarrow 0]$ does the job.

To evaluate the **Epistemic Condition** requires making some assumptions about the assassins's probability attributions. It sounds reasonable to assume that, without evidence to the contrary, each assassin attributed a higher probability to them shooting causing the outcome than them not shooting causing the outcome. Therefore the **Epistemic Condition** is also fulfilled for each assassin, and thus the BvH definition arrives at the right verdict for this example.

We continue with the approach pursued by Halpern Kleiman-Weiner (HK) [17], which uses the modified Halpern & Pearl definition of causation [16]:

**Definition 6 (HP).** *$\vec{X} = \vec{x}$ HP-causes $Y = y$ w.r.t. $(M, \vec{u})$ if there exists a $\vec{W} = \vec{w}$ so that the following conditions hold:*

AC1. *$(M, \vec{u}) \vDash \vec{X} = \vec{x} \land \vec{W} = \vec{w} \land Y = y$.*

AC2. *There is a setting $\vec{x}'$ such that $(M, \vec{u}) \vDash [\vec{X} \leftarrow \vec{x}', \vec{W} \leftarrow \vec{w}]Y \neq y$.*

AC3. *$\vec{X}$ is minimal; there is no strict subset $\vec{X}'$ of $\vec{X}$ such that $\vec{X}' = \vec{x}''$ satisfies AC2, where $\vec{x}''$ is the restriction of $\vec{x}$ to the variables in $\vec{X}'$.*

Note that, contrary to the direct NESS definition, the HP definition allows for conjunctive causes $\vec{X} = \vec{x}$, instead of merely atomic causes $X = x$. The minimality condition (AC3) is there to prevent irrelevant events to be added to such conjuncts. We can retrieve a definition of causation for atomic events by simply considering any conjunct $X = x$ that appears in an HP-cause $\vec{X} = \vec{x}$ to be a cause as well, which is indeed what Halpern suggests himself repeatedly [16].

The heart of the HP definition is AC2: it states that the outcome $Y = y$ counterfactually depends on the cause $\vec{X} = \vec{x}$ given that we intervene to hold fixed a suitably chosen set of variables $\vec{W}$ at their actual values $\vec{w}$. To see how this definition works, let us apply it to Example 1.

First we try substituting $\vec{X} = \vec{x}$ with $A_1 = 1$. Alas, this will not allow us to get $A_1 = 1$ as a cause of $V = 1$. We start with choosing $\vec{W} = \varnothing$, and we get that $(M, \vec{u}) \vDash [A_1 \leftarrow 0]V = 1$. The reason is that $\vec{u}$ encodes the actual context, in which $A_2 = 1$, and thus also $V = 1$. Yet what is required for AC2 would be $(M, \vec{u}) \vDash [A_1 \leftarrow 0]V = 0$. The only other choice for $\vec{W} = \vec{w}$ would be $A_2 = 1$, and that does not work either: $(M, \vec{u}) \vDash [A_1 \leftarrow 0, A_2 \leftarrow 1]V = 1$.

Second we try $A_1 = 1 \land A_2 = 1$. If this works, then AC3 is satisfied due to the fact that neither of the conjuncts themselves satisfied AC2. $\vec{W}$ has to be $\varnothing$, since there are no other variables. Thus what remains is to find counterfactual values for $A_1$ and $A_2$. As they are binary, the only option is to consider $A_1 = 0 \land A_2 = 0$. Clearly, for this choice AC2 is satisfied, as $(M, \vec{u}) \vDash [A_1 \leftarrow 0, A_2 \leftarrow 0]V = 0$. Therefore $A_1 = 1$ is an HP-cause of $V = 1$.

We can now formulate a definition of responsibility that is closely inspired by HK.

**Definition 7 (HK Responsibility).** *An agent who performs $A = a$ is responsible for outcome $O = o$ w.r.t. a responsibility setting $(M, \vec{u}, \mathcal{E})$ if:*

- **(Causal Condition)** *$A = a$ HP-causes $O = o$ w.r.t. $(M, \vec{u})$.*