

1 We would like to thank the reviewers for their comments and feedback. We are aware that in a largely conceptual paper
2 like ours there are subtleties, and highly appreciate the time and effort that the reviewers are putting in to digest these.

3 **Reviewer #1:** Causal Shapley values (SVs) are defined in Section 2. These do *not* coincide with what [9] and others
4 call the interventional SVs (marginal SVs in our terminology). Janzing et al. [9] write down the same equation, but
5 then choose to ignore any dependencies between the features in the real world (e.g., that in summer it tends to be
6 warmer than in winter). We do choose to incorporate these dependencies and hence cannot simplify to $P(\mathbf{X}_{\bar{S}}|do(\mathbf{X}_S =$
7 $\mathbf{x}_S)) = P(\mathbf{X}_{\bar{S}})$, but keep $P(\mathbf{X}_{\bar{S}}|do(\mathbf{X}_S = \mathbf{x}_S))$ in our definition of the causal SVs. We will follow the reviewer’s
8 suggestion to make this more explicit in Section 2. This distinction then hopefully also resolves the reviewer’s issue
9 about the indirect effect: it indeed vanishes for marginal SVs, but need not vanish for causal and conditional SVs. See
10 also the examples in Section 4 (Figure 1). The decomposition for conditional SVs follows by replacing “conditioning
11 by intervention” with “conditioning by observation”, i.e., by replacing $do(\mathbf{X} = \mathbf{x})$ with \mathbf{x} on the righthand side of the
12 bar. The decomposition is introduced in Section 3 to assist our illustration of how the different SVs attribute a model’s
13 prediction to the features involved in this prediction in Section 4 for different causal models. Here we also discuss in
14 which cases (most notably the fork and the confounder) conditional SVs fail to provide an intuitive causal attribution.

15 Causal chain graphs are introduced as a means to compute causal SVs (whether symmetric or asymmetric) when users
16 are willing/able to specify a (partial) causal ordering, but not a full-fledged causal model. The asymmetric SVs of [6]
17 indeed rely on the same information. On top of [6] we offer a formalization in terms of causal chain graphs and show
18 that, with “conditioning by intervention” instead of “by observation” as in [6], there is no need for asymmetry in the
19 SVs. Unlike conditional (asymmetric) SVs, causal SVs provide the right intuition in the case of common confounding.

20 **Reviewer #2:** W.r.t. the novelty in comparison to [6]: asymmetric (conditional) SVs as defined in [6] in some cases
21 coincide with symmetric or asymmetric (causal) SVs, but are different in general. See also the previous paragraph.

22 Section 4 aims to illustrate the behavior of the various SVs in simple cases that can be analyzed analytically and then to
23 argue which is the most intuitive, indeed also linking to psychological literature when appropriate. Here one prominent
24 theory, dating back to [15], states that humans sample over different possible scenarios to judge causation. Translating
25 this to a situation in which there are two possible causes, X_1 and X_2 , where it is unknown which one is intervened
26 upon first, may suggest that the natural interpretation is to consider both options and average over them.

27 We fully agree that quantifying causal influence is a difficult topic and any method has its weaknesses, but causal
28 SVs appear to fare better than the reviewer suggests. Discontinuity w.r.t. arrows with zero strength is an issue for the
29 asymmetric SVs, but not for the symmetric SVs that consider all orderings, not just those consistent with the causal
30 DAG. After averaging over all these orderings, the indirect effect already does incorporate all possible paths (so we do
31 not see how or why it needs to be generalized), but of course in the game-specific way inherent to the Shapley value
32 approach. We will add comments and disclaimers to clarify this and adapt our description of Janzing et al. and related
33 work as suggested by the reviewer. Our statement ‘not every causal query need be identifiable (see e.g., [24])’ did not
34 presume DAGs with all variables observed, but more general causal structures possibly including latent variables.

35 **Reviewer #3:** W.r.t. the scope, see our answer to Reviewer #1 (third paragraph) and the beginning of Section 5: causal
36 SVs are generally applicable when a user is willing/able to specify a causal model among the features that are used as
37 input to the model and when all causal queries are indeed identifiable. Specifying when this is the case is a topic on its
38 own: we will add more references (see also the supplement). Causal chain graphs are “just” proposed as a practical
39 approach to handle partial causal knowledge. In causal chain graphs, all causal queries are guaranteed to be identifiable
40 and can be answered based on the available observational data. These graphs allow for handling cycles, confounders,
41 etc (see Figure 2). In fact, all examples in Figure 1 are easily translated to causal chain graphs. An illustrative example
42 for the fork could be predicting hotel occupation (Y), based on season (X_2) and temperature (X_1).

43 We miss the point the reviewer tries to make w.r.t. counterfactual analysis. As far as we can tell, the counterfactual
44 question posed by the reviewer (assuming all features are known) can be answered simply by reading off the output of
45 the model. Our analysis can be interpreted as counterfactual (third rung) reasoning to analyze what the model prediction
46 would have been had we not known some of the input features (see second paragraph of Section 4). Counterfactual
47 *explanations* as in e.g. [33] may be improved with similar techniques, but are beyond the scope of the current paper.

48 Causal relationships are indeed asymmetric, but that does not prevent the causal SVs from being symmetric according
49 to the standard symmetry axiom for SVs (see the definition in Section 2 and the elaborate discussion in [9], Section 3 in
50 response to Sundarajan and Najmi, 2019). We chose not to repeat this argumentation, but will add a reference.

51 Figure 4 is meant to illustrate the difference between the various SVs (asymmetric SVs focus on the root cause, marginal
52 SVs on the direct effect, symmetric causal SVs consider both), not necessarily to claim that one is always better than
53 the other. We will extend the supplement with additional empirical analyses, e.g., on (deep) neural networks.

54 (7) indeed should have been (6). We will fix the other minor issues, also those rightfully indicated by **Reviewer #4**.