

1 We thank the reviewers for their comments. All reviewers agree that our approach on exploring computational
2 complexity as a way of comparing model interpretability would be of interest to the NeurIPS community, which is
3 very encouraging for us. **Reviewer 1** and **Reviewer 2** also welcomed our formal approach to it (**R1**: “*Classifying*
4 *the complexity of explanations is a natural step. Therefore this paper is a welcome contribution*”, **R2**: “*a thoughtful*
5 *approach to an important problem*”). Both reviewers mention as a strong point the use of parameterized complexity
6 when analyzing explanations for MLPs. **R1** even encourages us to present in more detail the classes that we use in our
7 results as they “*could have other uses in deep learning*” (because of space limitations it was difficult to do so in the
8 current version of the paper, but we would definitely add it in an extended version). We deal with specific comments by
9 the reviewers below.

10 **Reviewer 4** appreciates that we take the risk of formally studying the elusive notion of interpretability (**R4**: “*Important*
11 *goal of pinpointing an elusive concept*”), also praising our theoretical results (**R4**: “*Novel and sound complexity results*”).
12 Nevertheless, the reviewer questions the significance of the approach, which can be summarized in the following
13 comment: “*There does not seem to be a well-founded basis for [linking model interpretability with computational*
14 *complexity], since interpretability is, to the best of my knowledge, an inherently subjective concept.*” We agree with
15 the reviewer: interpretability is largely a subjective concept and the community is far from having a definitive answer
16 for what interpretability exactly means. But we actually consider this subjectivity as our main motivation: we wanted
17 to explore ways into which some formal ground can be given to this subjective concept, even if in a preliminary
18 form. This concern can thus be reduced to our choice of computational complexity as such a possible ground.
19 **R4** considers this choice as “*not well-founded*” while the other two reviewers consider the study of the computational
20 complexity of explanations as a plausible choice/step. We think that our results in the paper are a preliminary proof
21 of this plausibility/adequacy, as the complexity of explanation queries actually correlates with the informal views on
22 interpretability that can be found in the literature.

23 **R4** also questions our position in the following comment: “[...] *perhaps it would be more interesting to set out from a*
24 *more neutral position of looking for a correlation between [complexity and interpretability].*” We disagree with this
25 comment. Neutrality was a central concern when we were wording our submission and we were very careful in not
26 claiming that computational complexity is the right or only way to understand interpretability. This is why we phrased
27 our results explicitly as “*correlations*” between complexity and interpretability (exactly as the reviewer suggests). For
28 instance, in the Abstract (line 4) we write that “*We make a step towards such a notion by studying whether folklore*
29 *interpretability claims have a correlate in terms of computational complexity*”, and also in the Introduction with a
30 similar comment (line 46). We were also careful in naming our notion as *complexity-based interpretability* when we
31 formalize it (Section 2, line 110) to make it explicit that we are not formalizing a general notion of interpretability, but
32 exploring a (new) narrow yet formal one. Although we disagree with the reviewer’s view, this comment pushes us to
33 make our position even more explicit in the paper. For the next version of the paper, we will include an additional
34 discussion about this in the **Limitations** section and move this section to be right after the Introduction (it is the last
35 section in the current submission).

36 We now focus on the comments by **R1** and **R2**. Both reviewers point out that some of the problems we treat, or slight
37 variations of them, have already been discussed in the literature over different settings and under different names.
38 We agree that our paper would improve by adding a discussion of the different names and contexts under which the
39 literature has discussed similar problems, and explicitly mentioning their differences with our setting. We were under
40 space constraints but in case of acceptance we will use the additional space to improve our Related Work section.
41 **R1** and **R2** also suggest that our paper would benefit from being more explicit about the novelty and difficulty of some
42 of the proofs by distinguishing those that involve new insights (e.g., parameterized complexity) from those that follow
43 by already used techniques for the models that we consider. We will also implement this change in the next version.
44 Thank you very much for the comments and the pointers. We appreciate **R1** pointing out certificate complexity and the
45 study of boolean sensitivity as related subjects that could make for interesting lines of future research.

46 **R2** also comments about our choice of minimum sufficient reason (MSR) as one of the explainability queries as “*it’s*
47 *not clear why it would need to be minimum*”. We consider the minimality as a desirable property since an explanation
48 (as a sufficient reason) can always be padded with superfluous information while remaining valid. Our rationale was
49 that, among the different minimal sufficient reasons, it is arguably better to provide explanations that are as succinct as
50 possible (observe that an input x is always a sufficient reason for its own classification, but this is not really interesting).
51 This comment also pushes us to provide a more technical discussion of the relevance of the chosen queries in the final
52 version of our paper. Unfortunately due to space constraints we cannot address other interesting comments by **R2**,
53 like the relationship between monotone DNFs and MSR, that suggests a further study of asymmetric models (for
54 which MSR for positive and negative examples would have different complexities), or the need of using SHORTEST
55 IMPLICANT CORE instead of just SHORTEST IMPLICANT. But we thank the reviewer as both comments give us the
56 opportunity to improve our paper.