1  We thank all reviewers for their constructive feedback and detailed comments. Our responses are provided below.

2  **To Reviewer 1**  Thanks for your supportive and helpful comments. We also believe that applying our methods on
3  various fusion settings is an interesting direction for future research (Conclusion).

4  **Q:** Other dependencies of input sources, e.g., correlated information.
5  **A:** Using multiple camera inputs from various viewpoints can simulate an extremely correlated case. In fact, our
6  experiments deal with correlated information. Object proposals or information about potential surrounding objects can
7  be extracted from both RGB and LIDAR data, and *shared information* is introduced to consider such correlation.

8  **Q:** How our methods scale with more than two input sources.
9  **A:** Considering its definition, our MAXSSN loss will still try to balance single source robustness for multiple sensory
10  inputs. On computational complexity, our TRAINSSN requires $n_s + 1$ forward passes and 1 back propagation. One can
11  further improve on speed by approximating the loss, e.g., mixing TRAINSSN and TRAINSSNALT.

12  **To Reviewer 2**  We would like to appreciate your valuable and encouraging comments.

13  **Q:** Degraded performance of our methods on clean data.
14  **A:** Retaining performance on clean data, is an idealistic design goal. In section 3, we solve problem (3) by optimizing
15  over flexible parameters $g_1$ and $g_2$. If the parts of input sources contributing to $z_3$ are known, then indeed we can
16  achieve this goal. In practice however, it is difficult to know which parts of an input source (or latent representation) are
17  related to shared information and which parameters are flexible, and also the design goal becomes a soft rather than a
18  hard constraint. Therefore there is minor degradation in performance, to pay for the added robustness. We will add this
19  discussion in our paper.

20  **Q:** More tasks/models to demonstrate the generalization ability of the work
21  **A:** AVOD was selected because it is the leading deep fusion model for 3D object detection, and this model has already
22  been favorably compared with prior approaches, so we did not want to take up space with additional comparisons. 3D
23  detection is both an important problem in self-driving cars and one where multiple sensors can contribute fruitfully by
24  providing both complementary and shared information. In contrast, models for 2D object detection heavily rely on
25  RGB data, which typically dominates other modalities. Trying our approaches on other tasks like audio-visual speech
26  recognition is worthwhile, but our paper is already at the page limit covering more basic aspects, e.g., comparison of
27  fusion methods, training algorithms, and corruption methods.
28  We set the hyper-parameter for $\ell_1$ constraint as 0.01, and is available in our source code (omitted for anonymity).

29  **To Reviewer 3**  Thanks for your constructive comments and valuable time in understanding our work.

30  **Q:** More analyses of Algorithm 1 and 2 and explain why TRAINSSNALT works?
31  **A:** We can certainly add more explanation. Essentially the approaches are motivated by findings from the simple linear
32  fusion model, as well as data augmentation. Alternating between clean data and corrupted data aims at increased
33  robustness without much degradation of performance on clean data (strategy is supported by experimental results).
34  We also tried *fine-tuning* only fusion related layers to preserve essential parts for normal data, but Algorithm 1 and
35  2 are selected for better performance. We discuss the reasons for superior performance of Algorithm 2 in Remark 3.
36  In particular, an element-wise mean operation restricts the features to be optimal for averaging, and hence updating
37  without the maximum loss helps balance the performance. Also note that when AVOD with *concatenation* fusion is
38  trained with TRAINSSNALT, it works a lot worse than TRAINSSN.

39  **Q:** The proposed latent ensemble layer is not a very novel approach
40  **A:** When we devised our LEL, we first set the three objectives (1$^{\text{st}}$ sentence of the 2$^{\text{nd}}$ paragraph of Section 4.2): (i) error
41  reduction of ensemble methods, (ii) admitting source-specific features to survive, and (iii) allowing different channel
42  depths. Then we wrote the equation in Figure 1 and noted that it can be implemented by using $1 \times 1$ convolution.
43  Without applying fancier approaches which could increase computational cost, our LEL showed appealing effectiveness
44  even with simple implementation.

45  **Q:** Importance of single source robustness and the generalization ability of the model.
46  **A:** Currently there are no formulations even for single source, and our research is the first step. This framework can
47  be extended to robustness against corruption in a subset of input sources, e.g., generalizing to robustness given $k$ of
48  $n_s$ corrupted sources. Self-driving cars using an RGB camera and ranging sensors like LIDAR and radar are exposed
49  to single source corruption. For example, Uber's fatal self-driving crash occurred in nighttime and a camera couldn't
50  detect the victim because of darkness, while LIDAR and radar was not affected. In adversarial settings too, single
51  source attacks seem more feasible. Situations with multiple failures (outside of combat situations) have not been much
52  documented in the automobile industry. Therefore we feel that single source robustness should be studied in depth prior
53  to more general cases.