1 We appreciate all the reviewer comments. First, we summarize our contributions. We use a novel $2D^2CCA$ loss to
2 assure similar semantics are captured from the RGB and depth domains; we design a new model structure for sparse
3 depth completion, which exploits the relationship between RGB and depth data; our method achieves the state of the art
4 (SOTA) on several indoor/outdoor datasets.
5 **Response to common concerns**: First, we will release codes including the data processing tools. Second, we fixed the
6 symbol/annotation inconsistency and added symbols on Figure 2. We will carefully check grammar in the whole paper.
7 **Response for reviewer #1**:
8 1. "The difference between using the complementary RGB image and the dense RGB image is quite small." - The
9 difference between known and predicted depth is small (see ablation study, line 229-233 of paper). However, Table 1
10 validates that using only sparse and complementary information for depth completion is sufficient and can yield better
11 performance than using dense RGB information.
12 2. "For Table 3 a version of CFCNet with 0 points is missing (only RGB information)" - The core elements of our
13 CFCNet are $2D^2CCA$ and a range transformer. With only RGB information, the network structure would be just an
14 image encoder-decoder, and reduced to a normal FCN. This contradicts our purpose of multi-modal learning. Table 1
15 with only image information is for the ablation study. It shows that our core elements can improve performance.
16 3. Need to explain the result reported in the abstract. - The numbers in the abstract are compared with CNN-SLAM.
17 The "13.03" is the number for that sequence. The "58.89" is a typo - the correct number is "64.34" in Table 5. We agree
18 that using the result of a single sequence for comparison could be biased and unfair. Hence, we instead calculate the
19 average performance on SLAM RGBD datasets. We attain a +194% improvement on the datasets.
20 4. "Which dataset was used for the results in Table 1? Were the networks trained for the specific input configuration in
21 Table 1?" - We used the NYUv2 dataset, as described in the supplementary material. NYUv2 dataset was pre-processed
22 by a 100-point stereo sparsifier and then used as the training data. We will add the information to the body of the paper.
23 5. "Do you backpropagate to the parameters of the sparse depth branch from the 2D CCA loss?" - Yes, we use Eq. (4)
24 in the paper to backpropagate gradients through both the image branch and sparse depth branch. The whole network is
25 end-to-end trainable.
26 6. For the K->K experiment the number of sparse depth points is missing. - We state 100 depth points in the title of the
27 table. We will add it to the table body as well.
28 **Response for reviewer #2**:
29 We would like to clarify that our proposed loss is CCA-based but not SSA-based.
30 1. "False color visualizations, eg Fig 6, are only very qualitative and do allow to really compare performances of the
31 methods." - Most of the related works also use false color to visualize results. We follow this practice and use color
32 codes in the same way. Visualizing depth maps with false color is indeed qualitative. Hence we also provide numerical
33 results in paper to compare with others. We will add the numeric data of each example shown in Fig. 6 in the paper.
34 2. "Many tables (eg 1,2,3,4,5) have no units." - Meters for MAE and RMSE. Percentage for all $\delta$. We will add this info
35 in our final version.
36 3. "Model seems to work but is just another model." & "The model is a combination of existing elements but reasonably
37 engineered." - We respectfully disagree. Our method is not a simple combination of existing elements. Rather, it is a
38 multi-modal learning framework effectively uses $2D^2CCA$ to exploit the relationship between different types of data.
39 We want to demonstrate to other researchers a meaningful direction, in which the correlation between multi-modal
40 features for depth completion is fully utilized. Since our novel loss is independent from existing elements such as the
41 VGG-16 architecture, future researchers may use our loss along with other models to boost performance.
42 4. "The paper misses an important previous work: [P1] which should be reviewed and included into the comparison."
43 - We thought it might be inappropriate to compare our results with theirs because the study objectives are different.
44 Our study aims to complete depth when observable measurements are highly limited, while [P1] aims to complete
45 missing data with significantly more observable depth measurements. The authors note that their performance would be
46 negatively impacted if the observed depth samples are fewer than 2000. Our experiments use samples typically fewer
47 than 500.
48 **Response for reviewer #3**:
49 1. Meaning of many "directional"s. - "Directional" is synonymous to "dimensional" here, used in some statistical
50 communities. Indeed this may be confusing. We changed "directional" to "dimensional" to avoid confusion.
51 2. "The authors already demonstrate results when sampling 100 points from LiDAR data, but I am still interested in
52 what will happen when using all LiDAR measurements." - Thank you for your interest. We show numerical results
53 for training and evaluation on KITTI train/val dataset with an additional visual result here (left to right: image, depth
completion using full LiDAR points (MAE=0.2822), depth completion using 500 LiDAR points (MAE=0.7871)).

| Sample# | MAE | RMSE | $\delta_1$ | $\delta_2$ | $\delta_3$ |
|---|---|---|---|---|---|
| 500 (from paper,Table 2) | 1.197 | 2.964 | 94.0 | 98.0 | 99.3 |
| full points | 0.596 | 1.568 | 97.5 | 99.3 | 99.8 |