

DetAny4D: Detect Anything 4D Temporally in a Streaming RGB Video

Supplementary Material

777

7. Dataset Composition

778 DA4D is a hybrid 4D detection dataset which includes 12
 779 sub-datasets. The supporting tasks include sequence 4D
 780 detection, monocular 3D detection, depth estimation, and
 781 reconstruction tasks. DA4D is built on six datasets in
 782 Omni3D [5] (ARKitScenes [3], Hypersim [27], KITTI [12],
 783 nuScenes [6], Objectron [1], and SUNRGBD [30]) and ex-
 784 panded with a sequence dimension. Six more datasets are
 785 introduced for sequence-wise tasks, including Replica [31],
 786 MP3D [7], HM3D [26], HSSD [16], Gibson [38], and Scan-
 787 net [9]. The format is standardized similar to the Omni3D
 788 structure with additional sequence information. Each se-
 789 quence compromises a list of frames. Each frame includes
 790 monocular RGB image, camera intrinsics, camera pose,
 791 depth map, and object information. The object information
 792 includes 3D b-box attribute ($[x, y, z, w, h, l, yaw]$), rotation
 793 pose, 2D b-box prompt, category, instance ID, and score.

794 **Dataset Composition.** The dataset compromises origi-
 795 nal single frame 3D detection data as 3D detection capac-
 796 ity validation and multi-frame sequences data for 4D tasks.
 797 As shown in Figure 7, DA4D consists of multi-frame se-
 798 quences and Omni3D 3D detection data considered as se-
 799 quences with length of 1.

800 **Dataset Split.** Sequences from Omni3D [5] follows the
 801 splitting strategy in prior works [5]. For newly compro-
 802 mised datasets to form multi-frame sequences, we split the
 803 scenes into the training and validation set, and separately
 804 record sequences from the selected scenarios, ensuring all
 805 scenes in the validation set have not been visited.

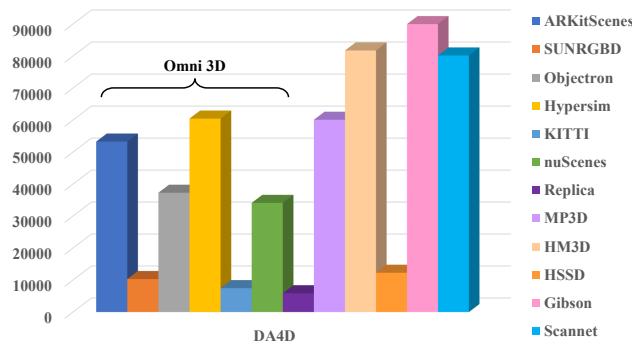


Figure 7. Visualization of the DA4D dataset composition.

806

8. Training Strategy Details

807 Here we illustrate the sequence crop and object padding
 808 strategy in Section 4.4 with more details. As shown in Fig-

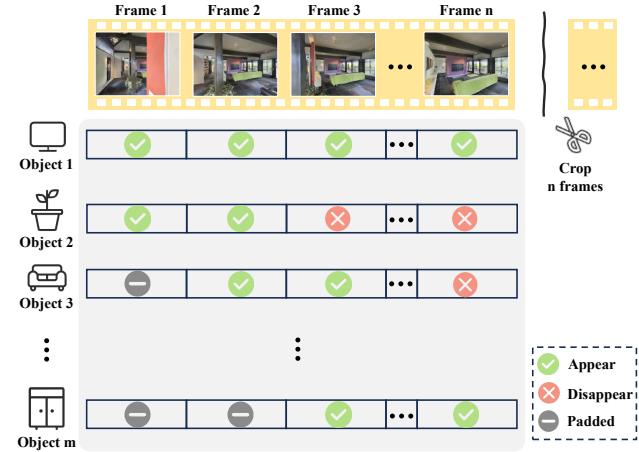


Figure 8. Visualization of the sequence crop and object padding strategy. The object query list maintains the objects and padding status. As frames in the sequence forecast, object status updates. The disappeared objects and padded objects do not contribute to the loss.

ure 8, we crop each clip under a fixed maximum length. For the cropped sequence, we count the total objects number in each sequence and pad the objects in each frame to this number. During training, this strategy ensures each frame has the same object query dimension, and predictions generated with the padded queries will be masked which do not contribute to the loss. During inference, the padded queries enable the model to manage newly appeared objects during forecasting, where new objects are registered to the padded queries and maintained if the prediction result has a high score and differs against objects in the query memory list.

9. 3D B-box annotation Pattern

This section illustrates the differences in the b-box annotation pattern for 3D detection task and 4D detection task. As shown in Figure 9, when observing an object from different views, 3D detection annotation pattern turns to consider each view as a singular observation and annotate the box referring only to the current view. On the contrary, for 4D detection, each object is registered with a b-box in the global coordinates and this constant b-box is projected to various views ensuring consistency globally.

Given the annotation formats of the 4D detection task and our goal to leverage powerful prior knowledge from 3D detection, specifically by utilizing pre-trained 3D detection models, we designed a specialized loss function (Section 4.4) to constrain the predictions to align with 4D task

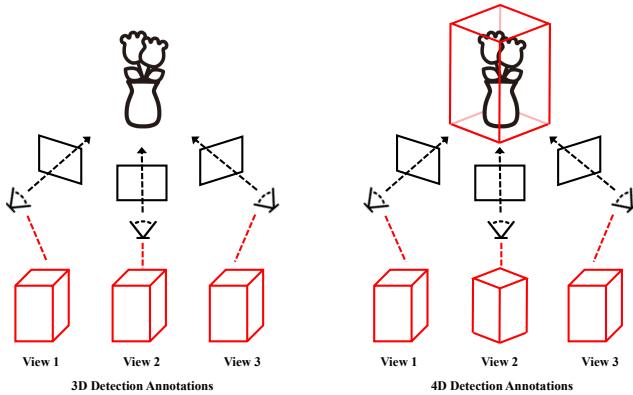


Figure 9. Visualization of the annotation pattern differences. 3D detection task turns to annotate each object referring only to the current view. 4D detection requires each view predicts objects in the global coordinates.

835 annotations. Our composite loss function, detailed in Section 4.4, incorporates constraints on the center, dimensions,
 836 and rotation angle of the 3D bounding boxes. Figure 10
 837 shows the supervision of the b-box. In contrast to the rigid
 838 constraints traditionally used in 3D detection, we employ a
 839 softened loss for the dimensions and rotation. This design
 840 is motivated by the potential misalignment between the out-
 841 puts of 3D pre-trained models and the 4D annotations, as
 842 shown in Figure 11. Directly applying hard constraints can
 843 lead to slow and difficult convergence, whereas the softened
 844 loss facilitates a more effective and stable alignment of the
 845 predictions to the 4D ground truth.
 846

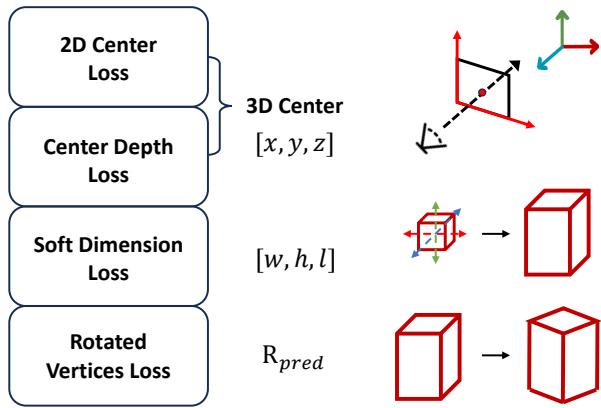


Figure 10. Visualization of the b-box supervision. The 3D bounding box is constrained sequentially by its center, dimensions, and rotation angle.

847 10. Open-Set Validation

848 To provide a detailed assessment of our model’s open-set
 849 detection capabilities, we separately evaluated the 4D de-

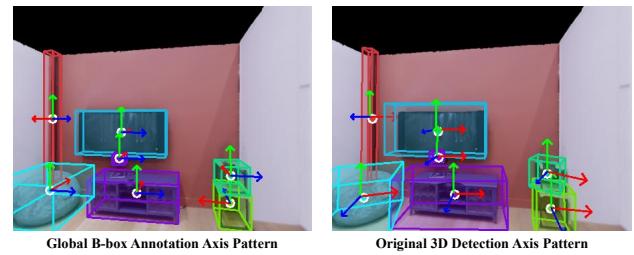


Figure 11. Visualization of the axis pattern comparison between global annotations and 3D prediction results. Global annotated b-box dimensions and rotation differs from the predicted results of pre-trained 3D detection model.

Datasets	Categories	Objects Num	Open-Set		None-Open-Set	
	Open-Set / All	Open-Set / All	AP _{3D} ↑	F1 _↑	AP _{3D} ↑	F1 _↑
Replica	18 / 62	0.3k / 16k	27.9	46.8	28.0	47.1
MP3D	20 / 330	3k / 98k	24.7	41.9	25.0	43.2
HM3D	75 / 698	11k / 344k	27.2	44.4	26.7	43.7

Table 4. Evaluation of the open-set performance. We separately evaluate the AP and F1 score on three sub-datasets and also report the ratio of the categories and numbers of open-set objects.

tention metrics on objects from categories that were within the training set (seen categories) and those that were outside of it (unseen categories) in the validation set. Table 4 shows the AP and F1 score under threshold IoU@0.5 (the same metrics illustrated in Section 5.1) of open-set categories and none-open-set categories in the validation set. It can be seen that the model performance on open-set highly keeps align with the none-open-set performance.

11. More Results

We provide more visualization comparison to show the prior performance of our proposed DetAny4D on the spatiotemporal consistency validation.

850
851
852
853
854
855
856
857

858
859
860
861



Figure 12. More qualitative comparison results.

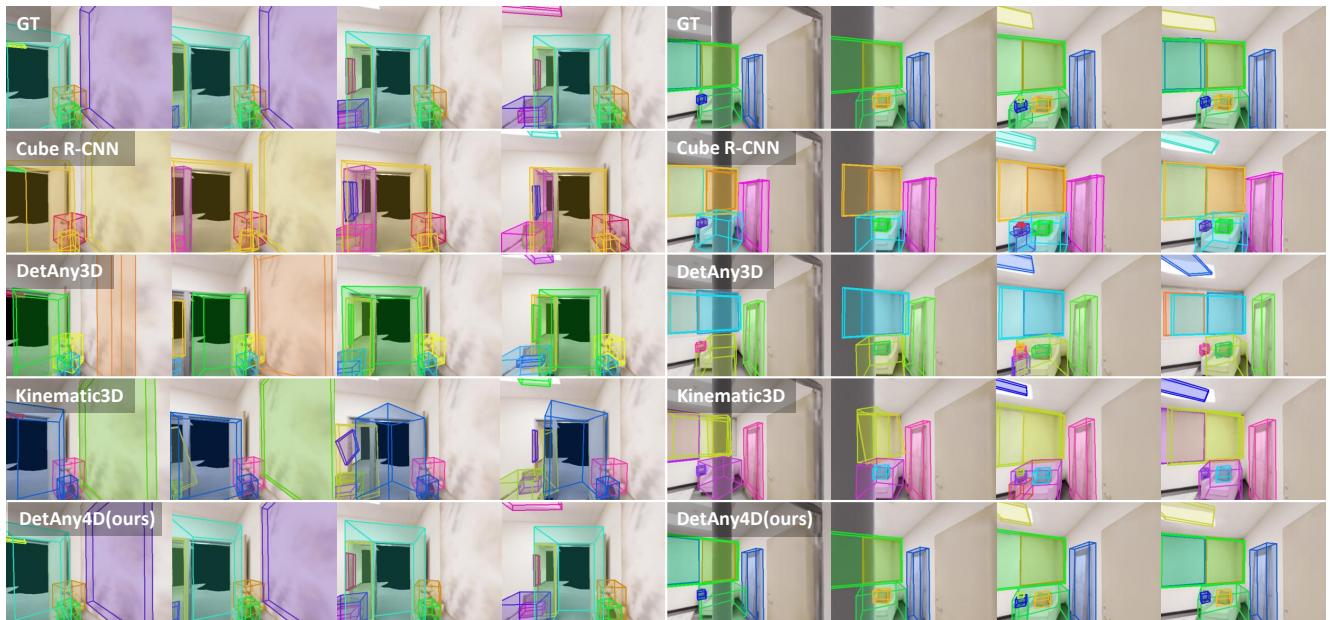


Figure 13. More qualitative comparison results.

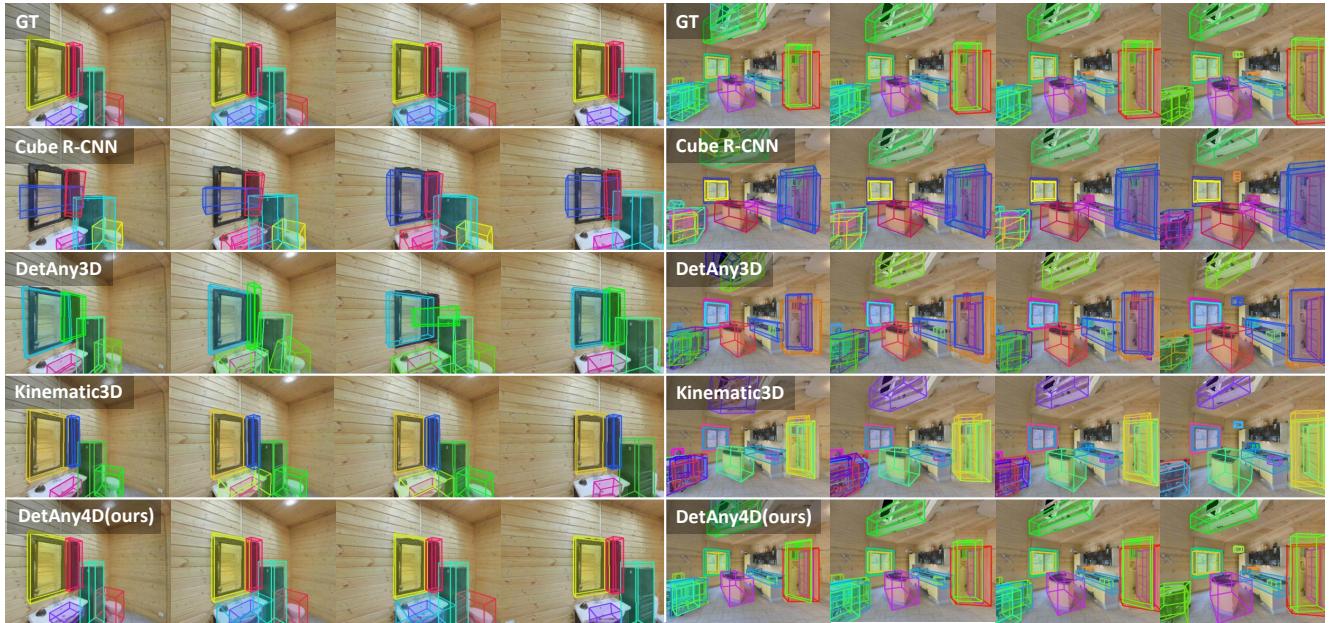


Figure 14. More qualitative comparison results.



Figure 15. More qualitative comparison results.

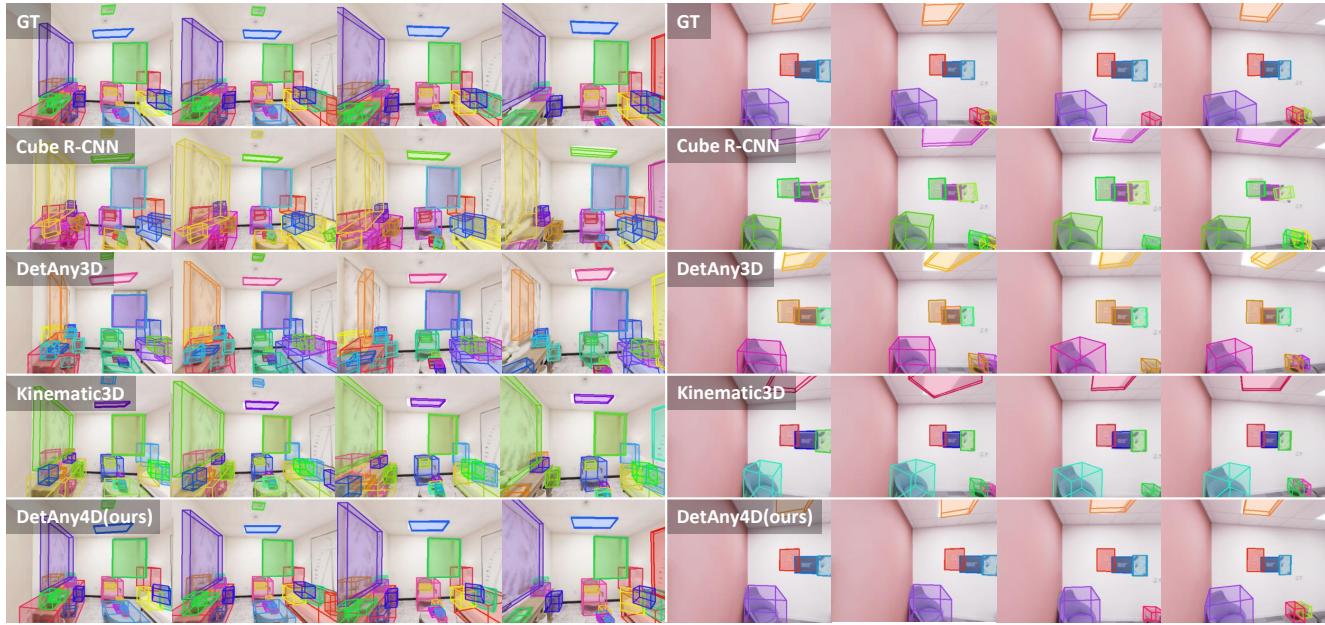


Figure 16. More qualitative comparison results.

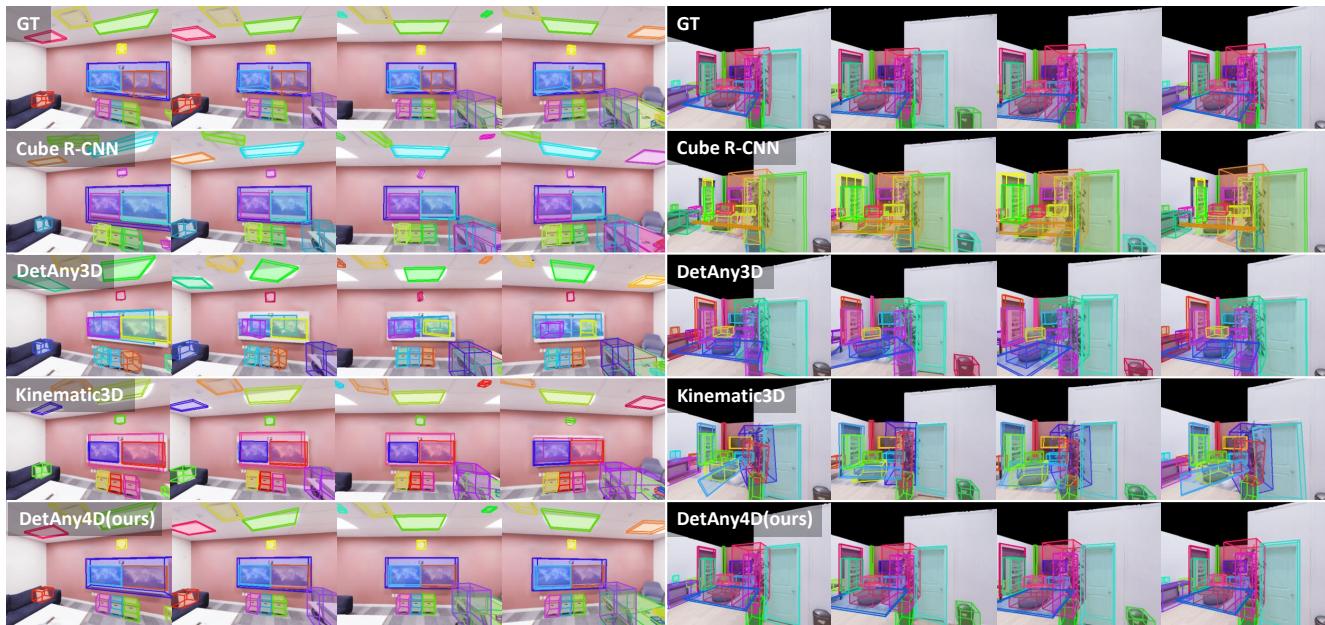


Figure 17. More qualitative comparison results.

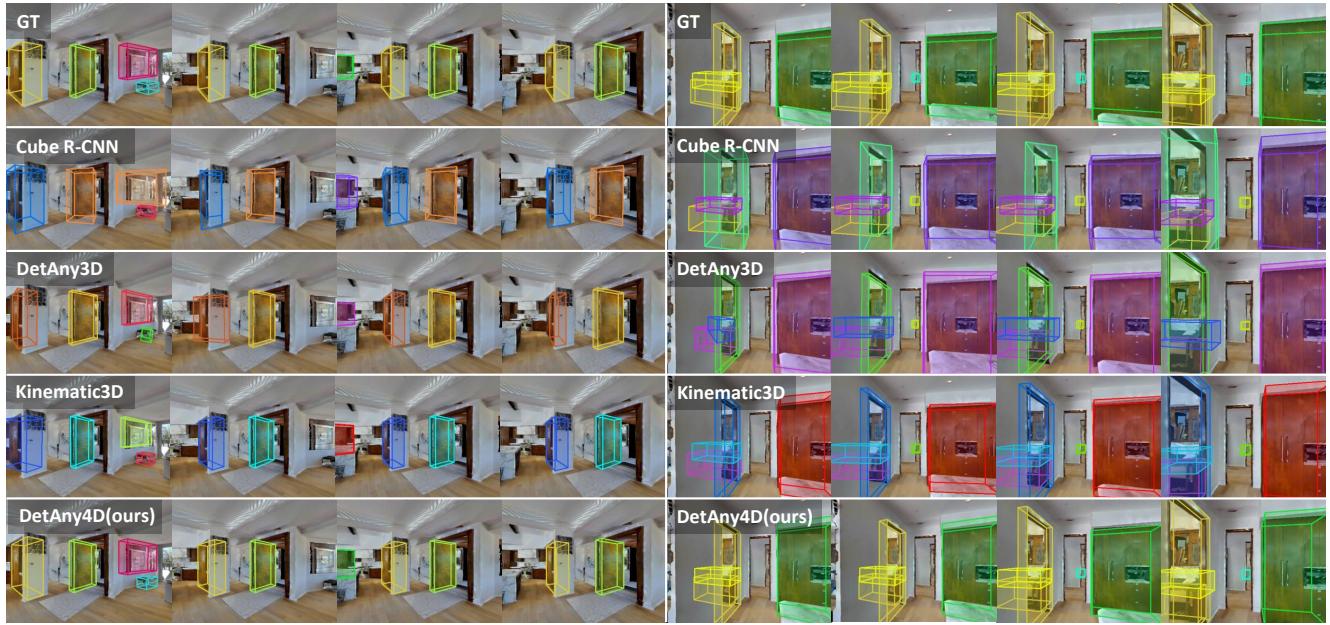


Figure 18. More qualitative comparison results.

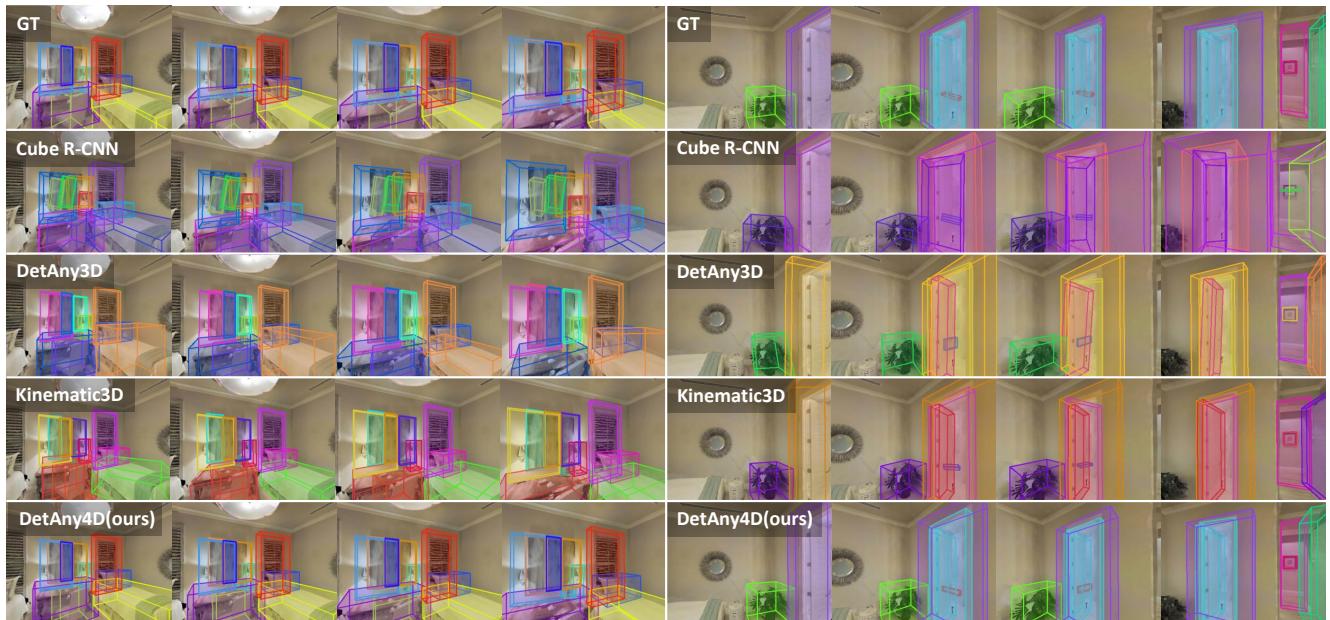


Figure 19. More qualitative comparison results.

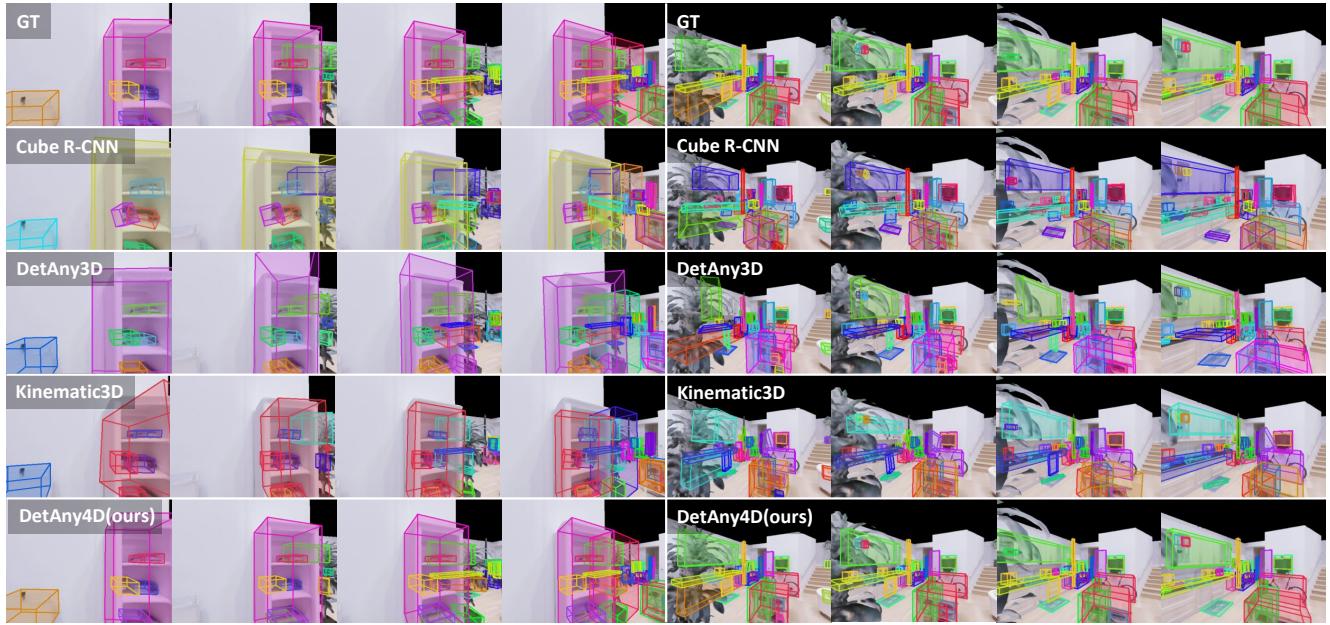


Figure 20. More qualitative comparison results.

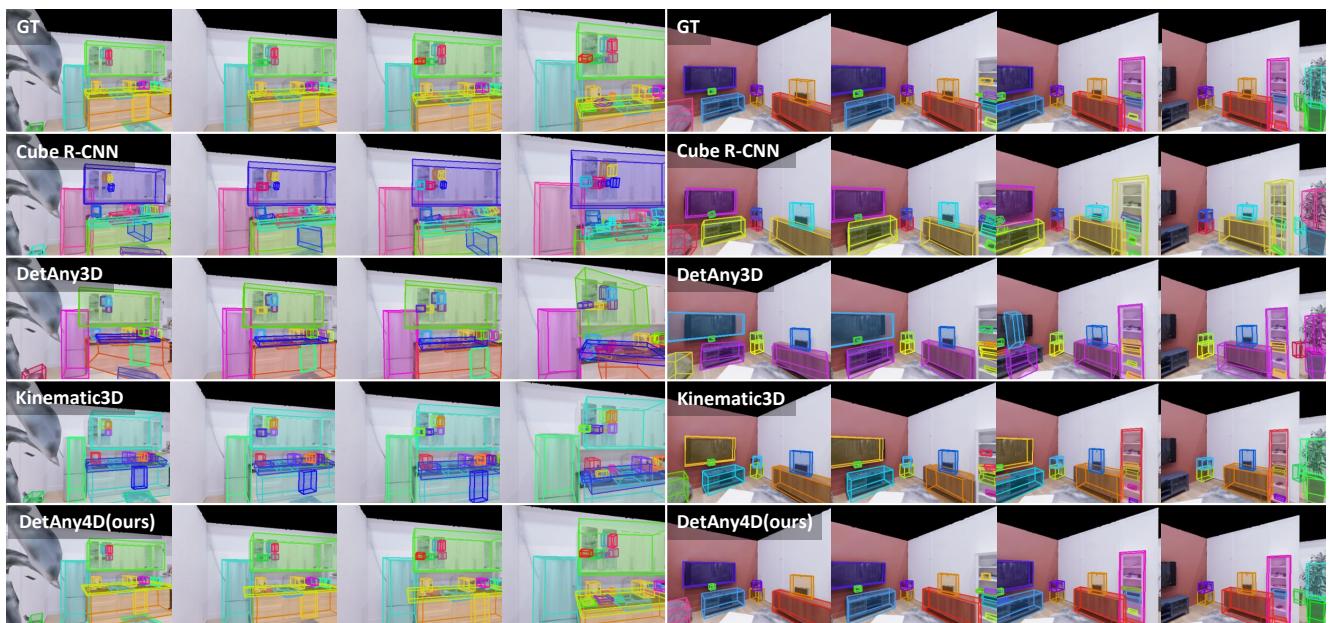


Figure 21. More qualitative comparison results.