

702 A EVALUATION METRICS

703
704 **Perception.** The evaluation for detection and tracking follows standard evaluation protocols Caesar
705 et al. (2020). For detection, we use mean Average Precision(mAP), mean Average Error of Transla-
706 tion(mATE), Scale(mASE), Orientation(mAOE), Velocity(mAVE), Attribute(mAAE) and nuScenes
707 Detection Score(NDS) to evaluate the model performance. For online mapping, we calculate the
708 Average Precision(AP) of three map classes: lane divider, pedestrian crossing and road boundary,
709 then average across all classes to get mean Average Precision(mAP).

710 **Planning.** We adopt commonly used L2 error and collision rate to evaluate the planning perfor-
711 mance. The evaluation of L2 error is aligned with VAD Jiang et al. (2023). For collision rate, there
712 are two drawbacks in previous Hu et al. (2023); Jiang et al. (2023) implementation, resulting in
713 inaccurate evaluation in planning performance. On one hand, previous benchmark convert obstacle
714 bounding boxes into occupancy map with a grid size of 0.5m, resulting in false collisions in certain
715 cases, e.g. ego vehicle approaches obstacles that smaller than a single occupancy map pixel Zhai
716 et al. (2023). (2) The heading of ego vehicle is not considered and assumed to remain unchanged Li
717 et al. (2024). To accurately evaluate the planning performance, we account for the changes in ego
718 heading by estimating the yaw angle through trajectory points, and assess the presence of a collision
719 by examining the overlap between the bounding boxes of ego vehicle and obstacles. We reproduce
720 the planning results on our benchmark with official checkpoints Hu et al. (2023); Jiang et al. (2023)
721 for a fair comparison.

722 B MORE ABLATION STUDY

723 **Necessity and Order of Object Selection.** Tab. 1 studies the necessity of agent and map selection
724 during the ego-centric hierarchical interaction. We can observe that agent selection contributes more
725 than the map selection, especially in the driving safety. And both of agent and map interactions are
726 conducted in a cascaded order is inferior than the parallel manner, where the updated ego query from
727 parallel outputs are concatenated for joint motion prediction.

730 Table 1: Effect of agent and map selection as well as interaction order in the hierarchical interaction
731 module.

Agent Selection	Map Selection	Cascade	Parallel	Planning L2 (m) ↓				Planning Coll. (%) ↓			
				1s	2s	3s	Avg.	1s	2s	3s	Avg.
✓	✗	-	-	0.16	0.34	0.64	0.38	0.03	0.05	0.22	0.10
✗	✓	-	-	0.17	0.35	0.63	0.38	0.02	0.06	0.28	0.12
✓	✓	✓	-	0.16	0.34	0.62	0.37	0.05	0.07	0.30	0.14
✓	✓	-	✓	0.16	0.33	0.59	0.35	0.00	0.04	0.18	0.07

739 **Effect of Interactive Score Fusion.** During the ego-centric query selection, both geometric and
740 classification scores are considered to ensure that the selected closest in-path queries are true positive
741 agents or maps, which are adopted for motion planner. Tab. 2 shows the effect of three types of scores
742 used for query ranking, namely attention, geometry and confidence scores. As described above,
743 interactive score S_{inter} obtained by multiplying these three scores can achieve the best selection
744 quality and planning performance. S_{inter} without confidence score fails to distinguish between
745 background and foreground queries, resulting in inferior performance.

747 Table 2: Effect of interactive score fusion process in the geometric attended selection step.

Attention Score	Geometric Score	Classification Score	Planning L2 (m) ↓				Planning Coll. (%) ↓			
			1s	2s	3s	Avg.	1s	2s	3s	Avg.
✓	✗	✗	0.18	0.36	0.66	0.39	0.09	0.11	0.28	0.16
✓	✓	✗	0.17	0.35	0.65	0.38	0.01	0.07	0.24	0.11
✓	✓	✓	0.16	0.33	0.59	0.35	0.00	0.04	0.18	0.07

754 **Effect of Iterative Refinement stages.** We continue to study the number of refinement stages in
755 Tab. 3. We can observe that our DiFSD can obtain superior planning performance with one addi-

tional refinement stage (36.3% collision rate reduction), which becomes saturated when introducing more stages. Hence, two-stage interacted motion planner is enough for achieving convincing results.

Table 3: Ablation for number of iterative refinement stages.

Number of stages	Planning L2 (m) ↓				Planning Coll. (%) ↓			
	1s	2s	3s	Avg.	1s	2s	3s	Avg.
1	0.16	0.33	0.61	0.37	0.01	0.08	0.23	0.11
2	0.16	0.33	0.59	0.35	0.00	0.04	0.18	0.07
3	0.16	0.33	0.60	0.36	0.01	0.40	0.22	0.09
4	0.16	0.33	0.61	0.36	0.00	0.04	0.20	0.08

Effect of Uncertainty Denoising. We also validate the effectiveness of uncertainty denoising strategy including position-level motion diffusion and trajectory-level planning denoising. As shown in Tab. 4, motion diffusion can improve the prediction stability with uncertain agent positions, while the planning denoising can also strengthen the trajectory regression precision of ego-vehicle.

Table 4: Ablation for uncertainty denoising procedure.

Position Diffusion	Trajectory Denoising	Planning L2 (m) ↓				Planning Coll. (%) ↓			
		1s	2s	3s	Avg.	1s	2s	3s	Avg.
✗	✗	0.16	0.34	0.64	0.38	0.07	0.07	0.17	0.10
✓	✗	0.16	0.34	0.63	0.37	0.02	0.04	0.15	0.07
✓	✓	0.16	0.33	0.59	0.35	0.00	0.04	0.18	0.07

C ANALYSIS & DISCUSSION

The GroundTruth future state distribution of ego-vehicle on nuScenes validation set is illustrated in Fig. 1, which is calculated with fixed time interval (1s) between consecutive predicted waypoints. And we also compare the output ego-state distribution of different popular end-to-end methods based on planned trajectories respectively, as show in Fig. 2. We can observe that without ego-centric design, the optimized end-to-end model is unable to handle various emergencies appearing in the driving scenarios, where the absolute values of Δv and Δa are larger than normal situations. Under this circumstance, the output planned trajectories cannot conform to the expert routes as expected. However, our DiFSD performs consistently better in planning the future ego states with variable speed and acceleration, owing to the ego-centric hierarchical interaction and selection mechanism, thus the iterative motion planner can focus on the interactive agents rather than irrelevant objects.

D VISUALIZATION

As show in Fig. 3, 4 and 5, we provide additional visualization results to illustrate the generalizability of DiFSD on various driving scenarios under different commands.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

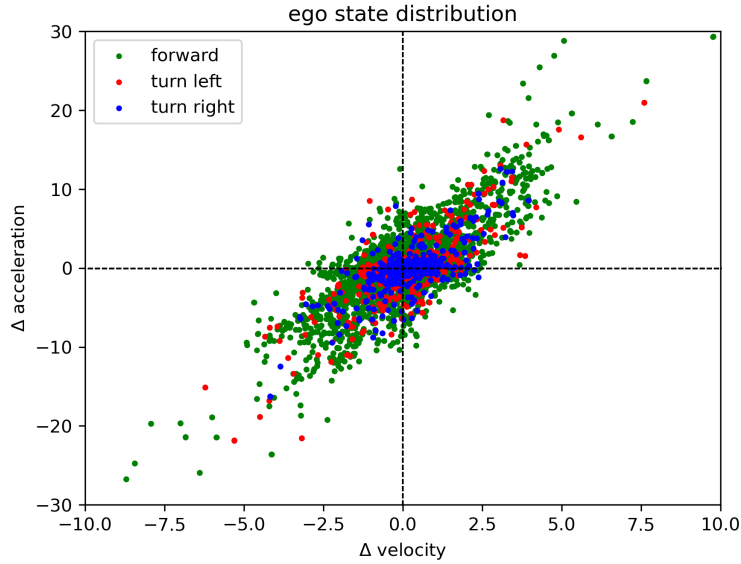


Figure 1: Distribution of GroundTruth future ego states (Δv vs. Δa) with different driving commands on the nuScenes val set.

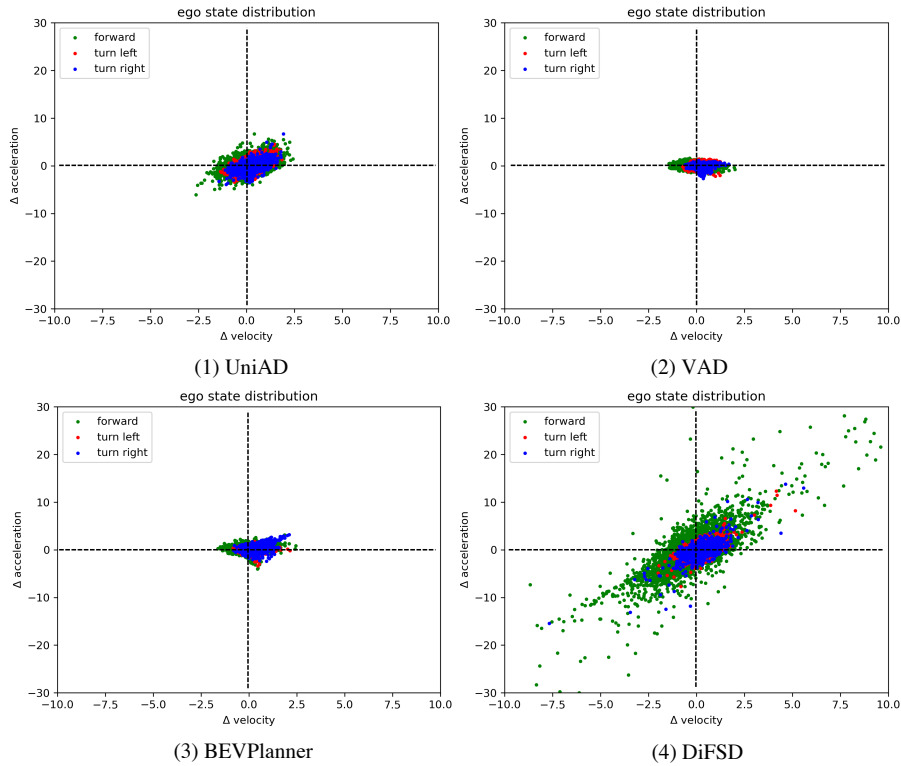


Figure 2: Comparison of predicted future ego states of different end-to-end methods on the validation set of nuScenes.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917



Figure 3: Qualitative results of DiFSD under “Go Straight” driving command in interactive scenes. In the first row, the pedestrian and the construction vehicle are selected as the closest in-path agents for motion prediction and interactive planning, thus DiFSD adjusts the planned trajectory from afar to avoid a collision. In the second row, DiFSD notices the pedestrian in the distance and plans the future trajectory taking the pedestrian intention into consideration. In the third row, DiFSD completes interactive decision-making in the “Cut-in” scenario, and outputs the planned trajectory constrained by the lane divider.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

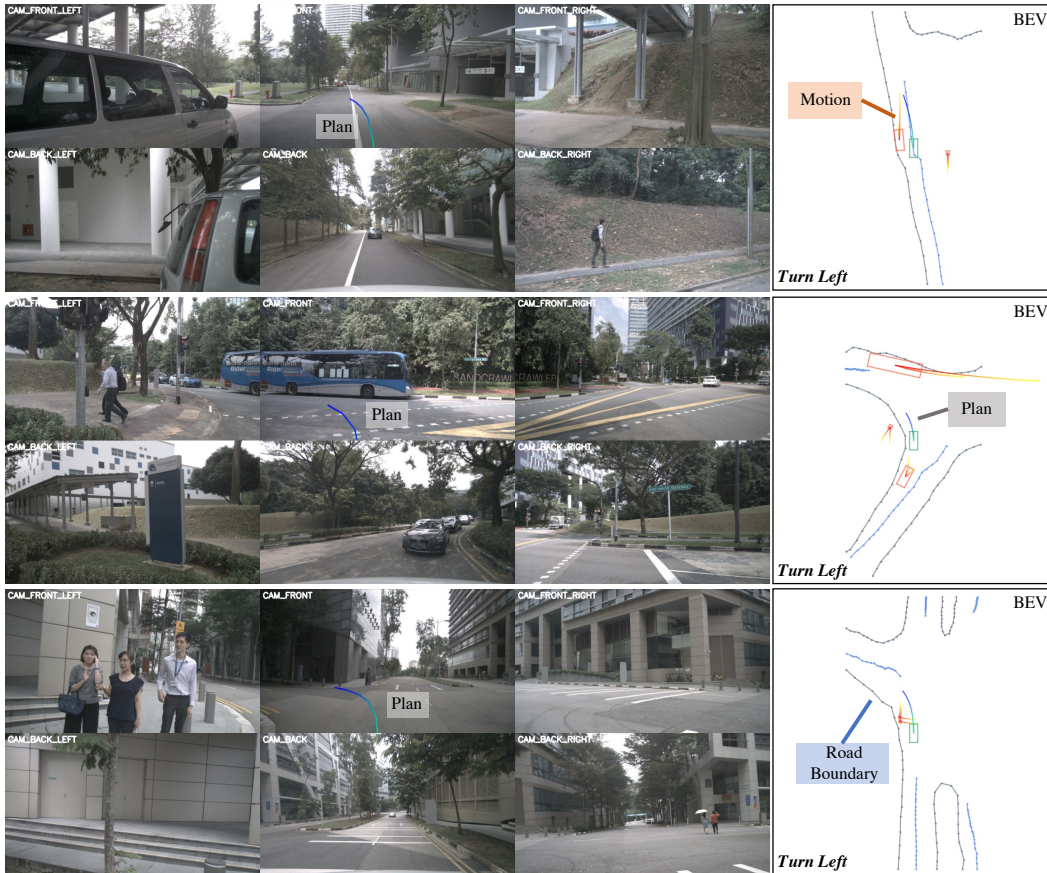


Figure 4: Qualitative results of DiFSD under “*Turn Left*” driving command in diverse scenarios. In the first scenario, DiFSD makes an “*Overtaking*” decision from the ride side of the front vehicle, considering the motions of both target vehicle and neighboring pedestrian to ensure driving safety. In the latter two intersection scenarios, DiFSD detects the pedestrians waiting at the crossing and the opposite bus passing the intersection, then decelerates to make a turning decision.

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

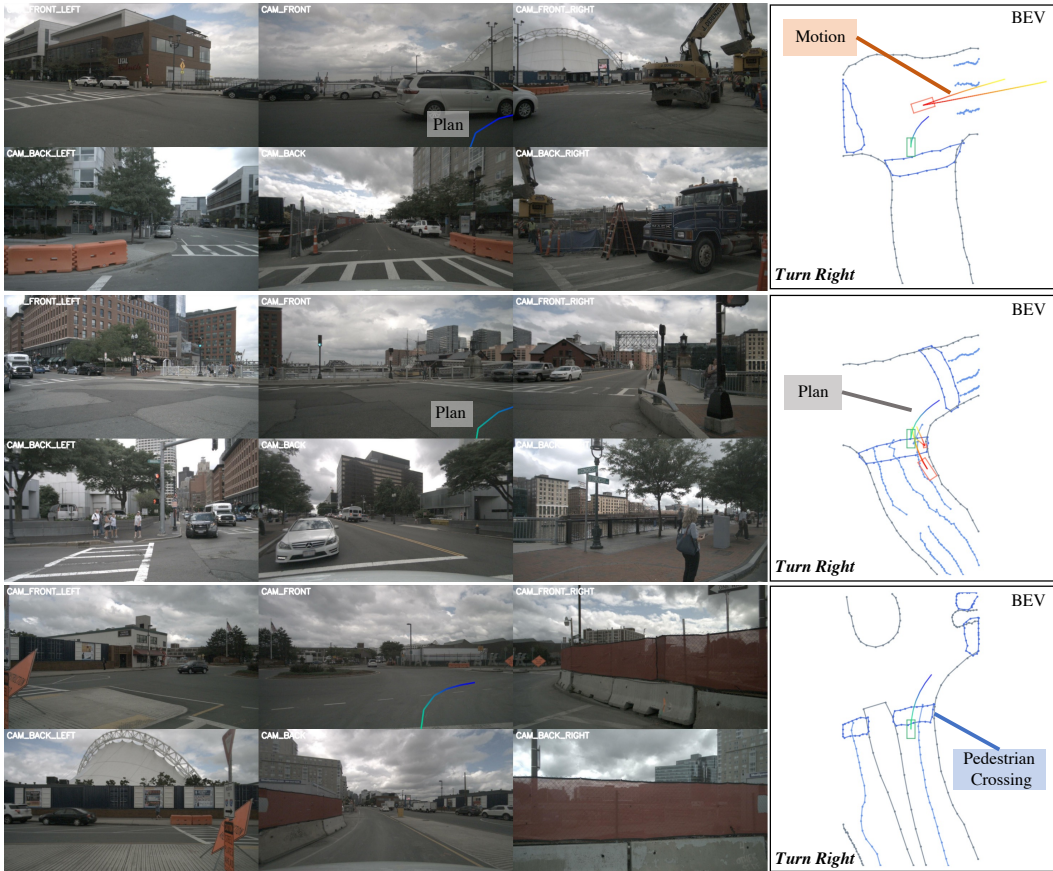


Figure 5: Qualitative results of DiFSD under “Turn Right” driving command at both interactive and non-interactive intersections. Joint motion prediction of agents and ego-vehicle is essential for DiFSD especially in the turning scenarios at intersections. The first two rows illustrate the interactive scenarios either inside and outside the intersection. And the last row presents a non-interactive intersection where DiFSD plans the future trajectory merely based on the detected pedestrian crossing.