

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 OMNI-VIEW: UNLOCKING HOW GENERATION FA- CILITATES UNDERSTANDING IN UNIFIED 3D MODEL BASED ON MULTIVIEW IMAGES

006  
007     **Anonymous authors**  
008     Paper under double-blind review

## 010 011     ABSTRACT

013     This paper presents Omni-View, which extends the unified multimodal under-  
014     standing and generation to 3D scenes based on multiview images, exploring the  
015     principle that “generation facilitates understanding”. Consisting of understanding  
016     model, texture module, and geometry module, Omni-View jointly models scene  
017     understanding, novel view synthesis, and geometry estimation, enabling synergistic  
018     interaction between 3D scene understanding and generation tasks. By design, it  
019     leverages the spatiotemporal modeling capabilities of its texture module responsi-  
020     ble for appearance synthesis, alongside the explicit geometric constraints provided  
021     by its dedicated geometry module, thereby enriching the model’s holistic under-  
022     standing of 3D scenes. Trained with a two-stage strategy, Omni-View achieves a  
023     state-of-the-art score of 55.4 on the VSI-Bench benchmark, outperforming existing  
024     specialized 3D understanding models, while simultaneously delivering strong  
025     performance in both novel view synthesis and 3D scene generation.

## 026     1 INTRODUCTION

028     Recently, unified multimodal models (UMMs) (Team, 2024a; Xie et al., 2024; Deng et al., 2025)  
029     have emerged as a pivotal area of research. The primary goal is to empower multimodal large  
030     language models (MLLMs) to both understand and generate visual signals present in our world,  
031     laying the foundation for artificial general intelligence.

032     Numerous methods (Team, 2024a; Wu et al., 2024a;c; Xie et al., 2024; Chen et al., 2025a; Deng  
033     et al., 2025) have achieved coexistence between 2D image understanding and generation. Among  
034     them, some studies (Wu et al., 2024b; Tong et al., 2024; Yan et al., 2025) aim to improve model per-  
035     formance by exploring the synergies arising from the interaction between understanding and genera-  
036     tion. For example, Metamorph (Tong et al., 2024) has extensively validated the role of understand-  
037     ing in improving generation performance. However, the potential and effectiveness of generation to  
038     improve understanding capabilities within unified models remain underexplored.

039     Due to the intrinsic geometric and spatiotemporal nature of 3D scenes and their multiview images,  
040     generative tasks such as geometry estimation and novel view synthesis are particularly well-suited  
041     for facilitating understanding in the 3D domain. This is attributable to the fact that 3D under-  
042     standing tasks (Yang et al., 2024; Zhang et al., 2025a) inherently necessitate robust geometric and  
043     spatiotemporal modeling capabilities, which can be effectively acquired through these generative  
044     tasks. Moreover, biomedical evidence (Maus et al., 2013; Nortmann et al., 2015) suggests that hu-  
045     man understanding of 3D environments is governed by the capacity to generate and imagine future  
046     sensory and geometric data (Keller et al., 2012; Leinweber et al., 2017) in the observed scene. These  
047     findings provide us with a guidance: the “generation facilitates understanding” paradigm presents a  
048     promising approach to building a unified model for 3D scene understanding and generation.

049     Inspired by the above analysis, this paper aims to fully unlock and maximize the benefits of genera-  
050     tion to understanding, thereby constructing a unified model for 3D scene understanding and genera-  
051     tion, called **Omni-View**. We emphasize that geometry estimation and novel view synthesis can  
052     leverage their inherent geometric measurement and spatiotemporal modeling capabilities to improve  
053     3D scene understanding, localization, and spatial reasoning. Specifically, we achieve this through  
   two key aspects, *architectural design* and *training strategy*.

054 The *architectural design* aims at unifying 3D scene understanding and generation, leveraging geometry estimation and novel view synthesis to advance 3D scene understanding. Our Omni-View is  
 055 built on Bagel (Deng et al., 2025), a strong unified framework in which the interaction between  
 056 understanding and generation is facilitated by its shared multimodal self-attention. However, Bagel’s  
 057 generative capacity is limited to RGB images, thus only capturing texture information, whereas 3D  
 058 scene generation necessitates inclusion of both texture and geometric structure. In accordance with  
 059 the dual generative objectives, the generation model in Omni-View is split into two distinct mod-  
 060 ules: a texture module and a geometry module. The texture module receives reference images, a  
 061 list of targeted camera poses, and prompt tokens encoded by the understanding model to generate  
 062 novel views of the scene. Meanwhile, the geometry module employs the hidden features from the  
 063 understanding model and the latent output of the texture module to infer geometric information of  
 064 novel views, such as depth maps and camera poses. This dual-pathway architecture empowers the  
 065 model to develop both geometric and spatiotemporal modeling capabilities, which are essential for  
 066 3D scene understanding tasks.  
 067

068 The principal goal of *training strategy* is to comprehensively improve the performance of the model.  
 069 A two-stage training strategy is employed. The first stage aims to augment the benefits of generation  
 070 for understanding 3D scenes, as introduced by the proposed architecture. The subsequent stage  
 071 is intended to refine the generation performance. In stage 1, the understanding model, geometry  
 072 module, and texture module are trained simultaneously. Geometry estimation assists the model  
 073 in comprehending the relative positional relationships among objects, thus directly enhancing the  
 074 model’s capability to evaluate the relative distances and directions of objects in 3D scenes. The  
 075 autoregressive generation forces the understanding model to discern the spatiotemporal relations  
 076 between the generated novel views, thus improving its understanding capabilities. As iterations  
 077 progress in stage 1, the number of reference images gradually decreases. This progressive shift  
 078 from dense to sparse views supports a curriculum-like, easy-to-difficult training approach, ultimately  
 079 enhancing the performance of the understanding model. In stage 2, the understanding model is  
 080 forzen. The generation model is finetuned via RGB-Depth-Pose joint generation, thereby enhancing  
 081 its capabilities in 3D scene generation.  
 082

083 We evaluate Omni-View on scene understanding, spatial reasoning, and novel view synthesis tasks.  
 084 The model achieves an impressive score of 55.4 on the VSI-Bench, exceeding current MLLMs  
 085 designed for visual reasoning. It manifests particularly notable improvements in subtasks such as  
 086 Relative Distance and Appearance Order, which require spatiotemporal modeling and the estimation  
 087 of geometry acquired through generation tasks. It also exhibits superior performance compared to  
 088 existing 3D understanding MLLMs in the area of 3D question answering (Ma et al., 2023; Azuma  
 089 et al., 2022). Furthermore, it effectively narrows the performance gap between unified models and  
 090 specialized models focused on 3D visual localization (Chen et al., 2020). Furthermore, we attain  
 091 robust results in the domain of novel view synthesis and scene generation, with particular emphasis  
 092 on enhanced perceptual quality.  
 093

## 2 RELATED WORK

094 **Scene understanding.** Recent advances in understanding 3D scenes have been greatly provided  
 095 by incorporating 3D or video input and 3D reconstruction prior. LLaVA-3D (Zhu et al., 2024)  
 096 and GPT4Scene (Qi et al., 2025) function within the voxel space and BEV. Video3DLM (Zheng  
 097 et al., 2024) improves localization ability by encoding 3D coordinates as position embeddings,  
 098 while Ross3D (Wang et al., 2025a) improves 3D understanding through visual-centric reconstruc-  
 099 tion. However, the dependency on 3D input of these methods poses practical application challenges.  
 100 To alleviate this issue, VG-LLM (Zheng et al., 2025) and SpatialMLM (Wu et al., 2025a) use the  
 101 features of VGGT (Wang et al., 2025b) as input, embedding the 3D piror in the model.  
 102

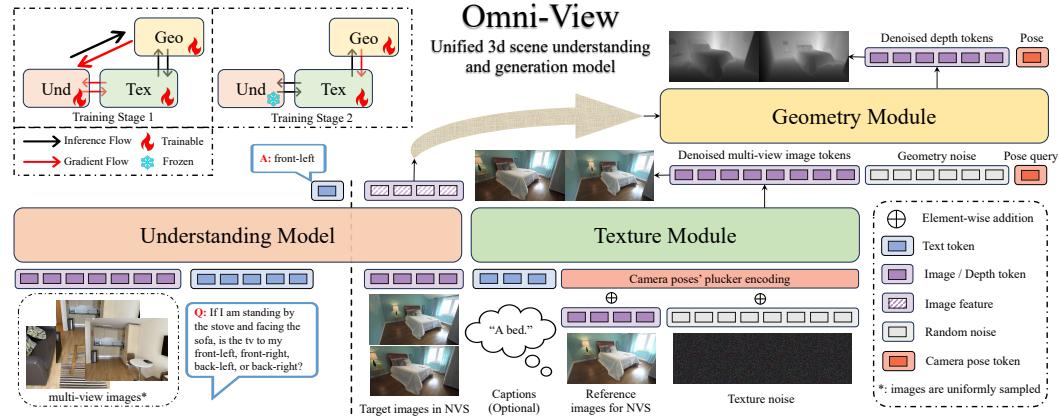
103 **Scene generation.** A vital element of scene generation methods is the presence of an efficient proxy  
 104 to represent the 3D scene, with panoramic images (Team et al., 2025), point clouds (Yu et al., 2024c),  
 105 and Gaussian Splatting (Yu et al., 2024a) among the viable options. In consideration of the advances  
 106 in image and video diffusion models (Rombach et al., 2022; Wan et al., 2025), contemporary 3D  
 107 scene generation strategies reconstruct the scene’s geometry from a single view, employing condi-  
 108 tioned video diffusion models to render the scene’s texture. ViewCrafter (Yu et al., 2024c) brings  
 109 explicit 3D information (point cloud) to generation models through iterative reconstruction. Won-  
 110

108 derJourney (Yu et al., 2024b) generates a comprehensive set of view sequences. Voyager (Huang  
109 et al., 2025) reconstructs the scene from a single view and uses it as conditions for inpainting.  
110

111 **Unified understanding and generation.** In the domain of 3D scenes, a unified model for under-  
112 standing and generation applicable in general scenarios remains absent. Hermes (Zhou et al., 2025)  
113 uses BEVs to design a unified understanding and generation model for autonomous driving. In  
114 contrast, significant progress has been made in 2D vision (Zhang et al., 2025b). These methods  
115 show primarily variations in architectures and training strategies. Chameleon (Team, 2024a) utilizes  
116 VQ-VAE for image tokenization, thereby improving generation competencies. However, its under-  
117 standing capabilities are inferior to those of Janus (Wu et al., 2024a), which employs SigLIP (Zhai  
118 et al., 2023) as visual understanding encoder. VILA-U (Wu et al., 2024c) integrates understanding  
119 and generation within the image encoder, bypassing multitask gradient conflicts during MLLM train-  
120 ing. BAGEL (Deng et al., 2025) implements task-based hard routing for MLLM, which also avoids  
121 gradient conflicts. Harmon (Wu et al., 2025b) takes advantage of MAE’s reconstruction ability and  
122 downstream understanding enhancement. BLIP3o (Chen et al., 2025a) introduces “understand first,  
123 then generation” training paradigm, thus achieving performance gains. Building on the progress in  
124 the 2D domain, this paper investigates a unified model in 3D scenes and explores how generation  
125 aids the understanding scheme within the framework.

### 126 3 METHOD

#### 127 3.1 ARCHTITECTURES



130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1348  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1508  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
191

162 **Texture Module.** The texture module in the generation model is tasked with novel view synthesis  
 163 using flow matching (Lipman et al., 2022). This module processes a textual description  $T_{des}$  and  
 164 some reference images  $I_{ref}$  of the current scene, together with a specified camera pose, to produce  
 165 consistent novel views of this scene. Within this module, the reference images  $I_{ref}$  are encoded  
 166 using the FLUX-VAE encoder (Labs, 2024)  $\varepsilon(\cdot)$ , and the vocabulary used to tokenize  $T_{des}$  aligns  
 167 with that used by the understanding model. This process incorporates the camera pose control as  
 168 delineated in MV-AR (Hu et al., 2025), embedding the Plucker-Ray encoding  $r_{i,j} = (o \times d, d)$  of  
 169 the camera pose as the absolute position encoding in the model, where  $o$  and  $d$  denote the origin  
 170 and direction of the ray,  $(i, j)$  represents the pixel coordinate. This camera pose embedding exhibits  
 171 adaptability to various image resolutions and demonstrates significant flexibility. The above process  
 172 can be described as:

$$F_{tex} = \text{TextureModule} ([LM\text{-}Head}(\tau(T_{des})); [\varepsilon(I_{ref}); N_{tex}] + r]), \quad (1)$$

173 where  $F_{tex}$  is the predicted image noise of the texture module,  $N_{tex}$  is the random input noise, and  
 174 LM Head is the text processing module of the understanding model.

175 **Geometry Module.** The geometry module in the generation model constructs the geometric aspects  
 176 of the generated images in the texture module. It synthesizes depth maps through flow matching  
 177 and employs a learnable query strategy to accurately estimate camera poses. It receives the latent of  
 178 novel view images  $F_{tex}$  from the texture module as input, which is concatenated with random depth  
 179 noise  $N_{dep}$  and a learnable camera pose query  $q_{cam}$  along the frame dimension. After processing,  
 180  $N_{dep}$  will be denoised as depth maps' latent, while  $q_{cam}$  will be decoded to reveal the intrinsic and  
 181 extrinsics of the camera. The features of novel view images  $F_{und}$ , output from the understanding  
 182 model's block, is also integrated into the geometry module through cross-attention. The process in  
 183 the geometry module can be described as follows.

$$[F_{dep}; \hat{g}] = \text{GeometryModule} ([F_{tex}; N_{dep}; q_{cam}], F_{und}), \quad (2)$$

187 where  $F_{dep}$  is the predicted depth noise of the geometry module.  $\hat{g}$  is the predicted intrinsics and  
 188 extrinsics of the camera, which is decoded by VGGT's camera decoder (Wang et al., 2025b).

189 The geometry module has independent parameters and maintains architectural connections with  
 190 both the texture module and the understanding model. Unlike the connection method in BAGEL,  
 191 the geometry module uniformly extracts features from the understanding model at the layer dimen-  
 192 sion as conditions. It utilizes cross-attention for fusion and ensures that the gradients of geometry  
 193 estimation can be backpropagated to the understanding model. Currently, the geometry module's  
 194 input only relies on the last-layer output latent of the texture module. This approach guarantees that  
 195 it acquires information closest to the image modality, thereby providing finely aligned features for  
 196 the estimation of depth map and camera pose.

### 200 3.2 TRAINING RECIPE

201 The Omni-View training recipe has two separate stages, as shown in the upper left of Figure 1.

202 **Stage 1: unify 3D understanding and generation.** In stage 1, we train the understanding model,  
 203 the texture module, and the geometry module simultaneously. It can leverage the fine-grained ge-  
 204 ometry estimation capability in the geometry module and combines it with the spatial-temporal  
 205 modeling ability within the texture module, thus improving the 3D understanding performance.

206 **Understanding model.** For the understanding model, we use the next token prediction to predict  
 207 the distribution of the answer text given the distribution of the multiview images and query text. The  
 208 loss function of the understanding model can be expressed as follows.

$$L_{und} = - \sum_{i=1}^T \log P_{\theta}(y_i | y_{<i}), \quad (3)$$

209 where  $y$  is the multimodal sequence that contains tokenized multiview images, query text, and an-  
 210 swer text.

216 **Texture module.** The loss function for the texture module is defined as the Mean Squared Error  
 217 (MSE) loss between each predicted texture noise  $F_{tex}$  generated by the texture module and the  
 218 provided texture noise  $N_{tex}$ .  
 219

$$220 \quad L_{tex} = \|F_{tex} - N_{tex}\|_2. \quad (4)$$

$$221$$

222 In contrast to the majority of novel view synthesis methods, the texture model employs a autoregres-  
 223 sive generation framework. Specifically, during the generation of the  $n$ -th frame, the model is ex-  
 224 posed to the visual data of the preceding  $n - 1$  frames, while excluding any subsequent frames. This  
 225 autoregressive methodology enables the model to fully grasp the concept of temporal sequences,  
 226 thereby enhancing its spatiotemporal modeling proficiency and effectively improving scene under-  
 227 standing. To optimize the 3D consistency of the sampled images, diffusion forcing (Chen et al.,  
 228 2024a) is employed when training the texture module.  
 229

230 To acquire this complex spatiotemporal generation capability incrementally, we gradually adjust  
 231 the reference images. As iterations progress, the reference images are systematically reduced in a  
 232 stepwise manner, transitioning from encompassing all input images, excluding the first image, to  
 233 including only the first image. This implies that for the model, the reference confidence transitions  
 234 from dense to sparse. We designate this progressive training approach as dense-to-sparse (D2S).  
 235 This strategy has been shown to be highly effective in facilitating improved understanding.  
 236

237 **Geometry module.** In stage 1, the geometry module is tasked with estimating the depth and pose  
 238 of the camera from both the provided images and those synthesized by the texture module. The loss  
 239 function for the geometry module comprises a sum of the depth estimation loss and the camera pose  
 240 estimation loss. In terms of depth estimation, we apply the MSE loss to the comparison between the  
 241 depth noise  $F_{dep}$  predicted by the geometry module and the given depth noise  $N_{dep}$ . The estimated  
 242 intrinsics and extrinsics of the camera  $\hat{g}$  are optimized directly through the Huber loss.  
 243

$$244 \quad L_{geo} = \|F_{dep} - N_{dep}\|_2 + \|\hat{g} - g_{gt}\|_\epsilon, \quad (5)$$

$$245$$

246 where  $g_{gt}$  is the ground truth of the camera pose and  $\|\cdot\|_\epsilon$  denotes the Huber loss.  
 247

248 Ultimately, a weighted summation of the aforementioned losses is conducted to derive the loss  
 249 function in stage 1  $L_{s1}$ .  
 250

$$251 \quad L_{s1} = \lambda_{und}L_{und} + \lambda_{tex}L_{tex} + \lambda_{geo}L_{geo}, \quad (6)$$

$$252$$

253 where  $\lambda_{und}$ ,  $\lambda_{tex}$ , and  $\lambda_{geo}$  represent the weighting coefficients and their default values are 1, 1, and  
 254 0.1 respectively.  
 255

256 **Stage 2: advance generation.** In stage 2, the texture module and the geometry modules are trained.  
 257 The RGB-Depth-Pose (RGBDP) joint learning methodology is used for training, capitalizing on the  
 258 geometry prior obtained from depth-pose estimation to enhance the ability to generate consistent  
 259 appearances for novel views.  
 260

261 In the texture module, the reference single-view image is used along with its depth map to recon-  
 262 struct the initial point cloud of the scene. The images rendered from different views, projected  
 263 through this point cloud, serve as conditions following Voyager (Huang et al., 2025). For the ge-  
 264 ometry module, it generates the depth map and camera pose from images synthesized by the texture  
 265 module. Concurrently, it no longer relies on the features of the understanding model as conditions  
 266 for cross-attention.  
 267

## 268 4 EXPERIENCE

$$269$$

### 270 4.1 EXPERIMENTAL SETUP

$$271$$

272 In our experiments, the understanding model and the texture module in the generation model are  
 273 initialized using the pre-trained BAGEL-7B (Deng et al., 2025). The geometry module in the gen-  
 274 eration model is configured to have the same dimensions as the texture module, with a depth of  
 275

270 four layers. For 3D scene understanding, a filtered dataset comprising 780k valid items, sourced  
 271 from SQA3D (Ma et al., 2023), ScanQA (Azuma et al., 2022), 3DOD (Zheng et al., 2025), Scan-  
 272 Refer (Chen et al., 2020), VLM-3R (Fan et al., 2025), SPAR (Zhang et al., 2025a) (234k subset)  
 273 and llava-hound4 (Zhang et al., 2024c) (64k subset), is meticulously curated for training. Regarding  
 274 novel view synthesis, 61k video clips are carefully selected from re10k (Zhou et al., 2018). The  
 275 corresponding depth maps are synthesized using the Voyager data pipeline (Huang et al., 2025) and  
 276 captions are synthesized using the QwenVLMax (Bai et al., 2025). During training, we do not use  
 277 images from the scene understanding task to train the generation model to fully demonstrate that the  
 278 understanding performance improvement brought by Omni-View does not come from memorizing  
 279 the data in understanding tasks.

280 Model training was performed using the AdamW optimizer, characterized by  $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ ,  
 281 the peak learning rate of  $1 \times 10^{-5}$ . The warm-up phase that constitutes 5% of the whole training  
 282 iterations. The training process is completed after one epoch of the understanding dataset. For the  
 283 understanding model, we do not rely on any 3D scene input to support both 3D scene understanding  
 284 and spatial reasoning tasks. In main comparisons, we use the same checkpoint for testing all tasks.

285

## 286 4.2 MAIN COMPARISONS

287

### 288 4.2.1 3D SCENE UNDERSTANDING

289

290 **Benchmarks.** We evaluate the 3D understanding performance of the model on question answering  
 291 (Ma et al., 2023; Azuma et al., 2022) and localization (Chen et al., 2020; Zheng et al., 2025).  
 292 During inference, we set the frame numbers as 32 following Video3DLM (Zheng et al., 2024).

293

294 **Comparison baselines.** We compare Omni-View with models specifically designed for 2D or 3D  
 295 visual understanding tasks (Team, 2024b; Wang et al., 2024a; Zhang et al., 2024c; Huang et al.,  
 296 2023; Zhang et al., 2024a; Chen et al., 2024b; Zhu et al., 2024; Zheng et al., 2024; Qi et al., 2025;  
 297 Wang et al., 2025a), as well as with some unified models applicable to video modalities (Deng et al.,  
 298 2025). Unified models are evaluated after fine-tuning with the same data. Within the results table,  
 299 we differentiate between models that incorporate explicit 3D input and those that do not. Although  
 300 the incorporation of explicit 3D input improves model performance, it concurrently restricts applicability  
 (Zheng et al., 2024).

301 Methods	3D Input	SQA3D <sub>test</sub>		ScanQA <sub>val</sub>				3DOD		ScanRefer	
		EM	EM-R	C	B-4	M	R	EM	F1	Acc@0.25	Acc@0.5
<i>Task-specific Models</i>											
LEO	✓	50.0	52.4	80.0	11.5	16.2	39.3	21.5	—	—	—
ChatScene	✓	54.6	57.5	87.7	14.3	18.0	41.6	21.6	—	55.5	50.2
Grounded 3D-LLM	✓	—	—	—	—	—	—	—	—	47.9	44.1
Video-3D-LLM	✓	58.6	—	102.1	16.4	20.0	49.3	30.1	—	58.1	51.7
GPT4Scene-HDM	✓	59.4	62.4	96.3	15.5	18.9	46.5	—	—	62.6	57.0
Ross3D	✓	63.0	65.7	107.0	17.9	20.9	50.7	30.8	—	61.1	54.4
LLaVA-3D	✓	60.1	—	103.1	16.4	20.8	49.6	30.6	—	50.1	42.7
InternVL2-8B	✗	33.0	45.3	62.5	3.3	14.5	34.3	—	—	—	—
Qwen2-VL-7B	✗	48.5	—	53.9	3.0	11.4	29.3	—	—	—	—
LLaVA-Video-7B	✗	48.5	—	88.7	3.1	17.7	44.6	—	—	—	—
SPAR-7B	✗	58.1	—	90.7	15.3	—	—	—	—	48.8 (31.9)	43.1 (12.4)
VG-LLM-4B	✗	—	—	—	—	—	—	—	47.0	53.5 (36.4)	47.5 (11.8)
SpatialMLLM-4B	✗	55.9	58.7	91.8	14.8	18.4	45.0	—	—	—	—
<i>Unified Models</i>											
BAGEL-7B-FT	✗	57.2	59.7	95.5	14.7	18.7	46.3	27.0	41.3	46.9 (28.0)	41.6 (7.7)
Omni-View-7B	✗	59.2	61.9	103.0	16.2	20.1	49.0	29.5	46.4	50.8 (32.5)	45.0 (9.9)

316 Table 1: **Evaluation of 3D scene understanding.** “—” indicates the number is not available for us.  
 317 **Bold** and underline denote the best and second-best models without 3D scene input, respectively. For  
 318 ScanRefer, the content in “()” indicates results without proposal refinement (Zhang et al., 2025a).

319

320 **Metrics.** For SQA3D, we use the EM metric to evaluate the accuracy, which stands for top-1 exact  
 321 match. EM-R means the refined EM following LEO (Huang et al., 2023). For ScanQA, we use  
 322 CIDEr (C), BLEU-4 (B-4), METEOR (M), ROUGE (R), and EM for more complete validation.  
 323 For 3DOD, we use the average F1 score (F1) as a metric to assess the correspondence between  
 the predicted and actual coordinates. For ScanRefer, we calculate the percentage of samples for

324 which the Intersection over Union (IoU) exceeds thresholds of 0.25 and 0.5, respectively, between  
 325 the predicted and true coordinates. The higher the above metrics, the better.

326 **Results.** As shown in Table 1, our analysis leads to four conclusions. (1) Our Omni-View exceeds  
 327 all current MLLM methods that do not depend on 3D scene input. Within the SQA3D test set,  
 328 Omni-View achieves an enhancement of 3.3 over SpatialMLLM in EM, while in the ScanQA val-  
 329 idation set, it surpasses SpatialMLLM by an increment of 11.2 in CIDEr. For the object detection  
 330 task, our unified model performs comparably with the 3D understanding model VG-LLM. (2) This  
 331 performance improvement is mainly attributed to the architectural design and training scheme that  
 332 we proposed. Compared with the fine-tuned BAGEL, our method still demonstrates substantial im-  
 333 provements. For example, on the SQA3D test set, Omni-View enhances EM by 2 points compared  
 334 to the fine-tuned BAGEL (Deng et al., 2025); on the ScanQA validation set, Omni-View improves  
 335 7.5 in CIDEr. (3) The efficacy of our approach in the question-answering task is equivalent to ad-  
 336 vanced MLLM methods that require 3D scene input. In particular, in the ScanQA validation set,  
 337 Omni-View achieved a performance comparable to Video3DLM (Zheng et al., 2024) and LLaVA-  
 338 3D (Zhu et al., 2024). (4) However, there remains a notable disparity between methodologies that do  
 339 not require 3D scene input and those that do, particularly in 3D grounding tasks. The experimental  
 340 results of VG-LLM (Zheng et al., 2025) also substantiate this observation.

#### 341 4.2.2 3D SPATIAL REASONING

343 **Benchmarks.** We evaluate the 3D spatial reasoning ability of the model in VSI-Bench (Yang et al.,  
 344 2024). During inference, we follow the VSI-Bench to set frame numbers ranging from 8 to 32 and  
 345 frame resolution to 640p.

346 **Comparison baselines.** We compare our Omni-View with models specifically designed for 2D or  
 347 3D visual understanding tasks (Xue et al., 2024; Zhang et al., 2024b; Bai et al., 2025; Ray et al.,  
 348 2024; Zhang et al., 2025a; Zheng et al., 2025; Wu et al., 2025a), as well as with some unified models  
 349 applicable to video modalities (OpenAI, 2024; Team et al., 2024; Deng et al., 2025; Xie et al., 2025).  
 350 Unified models are evaluated after fine-tuning with the same data we used.

352 Methods	Numerical Quesiton				Multiple-Choice Question				Avg.
	Obj. Cnt.	Abs. Dist.	Obj. Size	Room Size	Rel. Dist.	Rel. Dir.	Route Plan	Appr. Order	
<i>Task-specific Models</i>									
LongVILA-8B	29.1	9.1	16.7	0.0	29.6	30.7	32.5	25.5	21.6
LongVA-7B	38.0	16.6	38.9	22.2	33.1	43.3	25.4	15.7	29.2
LLaVA-OneVision-72B	43.5	23.9	57.6	37.5	42.5	39.9	32.5	44.6	40.2
LLaVA-Video-72B	48.9	22.8	57.4	35.3	42.4	36.7	35.0	48.6	40.9
Qwen2.5VL-7B	40.9	14.8	43.4	10.7	38.6	38.5	33.0	29.8	33.0
Qwen2.5VL-72B	25.1	29.3	54.5	38.8	38.2	37.0	34.0	28.9	37.0
SAT-LLaVA-Video-7B	–	–	–	47.3	41.1	37.1	<b>36.1</b>	40.4	–
SPAR-8B	–	–	–	–	–	–	–	–	41.1
VG-LLM-4B	<u>66.4</u>	<u>36.6</u>	55.2	<b>56.3</b>	40.8	43.4	30.4	39.5	46.1
Spatial-MLLM-4B	65.3	34.8	<u>63.1</u>	45.1	41.3	46.2	33.5	<u>46.3</u>	<u>48.4</u>
<i>Unified Models</i>									
GPT-4o (API)	46.2	5.3	43.8	38.2	37.0	41.3	31.5	28.5	34.0
Gemini-1.5 Pro (API)	56.2	30.9	64.1	43.6	<u>51.3</u>	46.3	<u>36.0</u>	34.6	45.4
BAGEL-7B-FT	62.8	36.3	56.4	49.7	46.1	<u>49.4</u>	26.8	43.1	46.3
Omni-View-7B	<b>70.3</b>	<b>46.4</b>	<b>68.6</b>	<u>54.7</u>	<b>65.9</b>	<b>54.4</b>	33.5	<b>49.0</b>	<b>55.4</b>

363 Table 2: **Evaluation of spatial reasoning on VSI-Bench.** “–” indicates the number is not available  
 364 for us. **Bold** and underline denote the best and second-best models, respectively.

370 **Results.** The results on spatial reasoning tasks more fully demonstrate Omni-View’s improvement  
 371 over previous methods in analyzing the relative or absolute position and orientation of spatial  
 372 objects. With an average score of 55.4, our Omni-View ranks first among existing spatial reasoning  
 373 MLLMs. Compared to existing Spatial-MLLM (Wu et al., 2025a) and VG-LLM (Zheng et al.,  
 374 2025), Omni-View improves Spatial-MLLM by 11.6, 9.6, and 24.6 in Absolute Distance (Abs.  
 375 Dist.), Room Size, and Relative Distance (Rel. Dist.), respectively. Omni-View improves VG-LLM  
 376 by 9.8, 13.4, 25.1, 11.0, and 9.5 in Abs. Dist., Object Size (Obj. Size), Rel. Dist., Relative Direction  
 377 (Rel. Dir.), and Appearance Order (Appr. Order), respectively. These tasks necessitate the model  
 378 to thoroughly predict the spatiotemporal state and measure the geometric properties of the scene it

378 observes (Yang et al., 2024). This demonstrates that our proposed method substantially enhances  
 379 the model’s pertinent capabilities.  
 380

#### 381 4.2.3 3D SCENE GENERATION 382

383 **Benchmarks.** The model’s generation capacity is validated under two tasks: novel view synthesis  
 384 (NVS) from a single view and scene generation. NVS from a single view necessitates that the model  
 385 generate the subsequent 25 frames from the first image. 3D scene generation evaluates scenes that  
 386 are reconstructed from the videos generated to 3DGS. The test scenes are randomly selected from  
 387 the Re10k test set following Chen et al. (2025b).  
 388

389 **Comparison baselines.** We compare with scene generation methods (Chung et al., 2023; Wang  
 390 et al., 2024b; Ma et al., 2025; Yu et al., 2024c; Chen et al., 2025b; Huang et al., 2025) and a unified  
 391 model (Deng et al., 2025). To evaluate the performance of scene generation, we use Dust3R to recon-  
 392 struct the generated videos, ensuring fairness in the comparison, according to Chen et al. (2025b).  
 393 The evaluation is only conducted on the first 25 frames generated by the model following Yu et al.  
 394 (2024c); Zhai et al. (2025).  
 395

396 **NVS from single view and scene generation.** Omni-View achieved the highest PSNR and SSIM,  
 397 and lowest LPIPS score, indicating that its image quality could surpass that of other methods. How-  
 398 ever, in terms of pixel-level fidelity, Omni-View shows only slight improvements over popular scene  
 399 generation models. This discrepancy may be attributed to Omni-View’s challenges in being pre-  
 400 cisely controlled via the camera pose. Visualization results are presented in the Appendix due to  
 401 page limit.  
 402

403 Methods	404 NVS from Single View			405 Scene Generation		
	406 PSNR ↑	407 SSIM ↑	408 LPIPS ↓	409 PSNR ↑	410 SSIM ↑	411 LPIPS ↓
<i>412 Task-specific Models</i>						
413 LucidDreamer	414 22.27	415 0.766	416 0.204	417 21.98	418 0.698	419 0.290
420 MotionCtrl	421 15.86	422 0.520	423 0.431	424 15.33	425 0.479	426 0.590
427 See3D	428 22.37	429 0.781	430 0.199	431 21.60	432 0.744	433 0.238
434 ViewCrafter	435 22.60	436 0.754	437 0.195	438 22.25	439 0.709	440 0.204
441 FlexWorld-5B	442 23.05	443 0.788	444 0.182	445 22.86	446 0.756	447 0.198
448 Voyager-13B	449 23.12	450 0.793	451 0.175	452 22.93	453 0.768	454 0.194
<i>455 Unified Models</i>						
456 BAGEL-7B-FT	457 21.76	458 0.703	459 0.288	460 21.04	461 0.599	462 0.403
463 Omni-View-7B	464 23.22	465 0.817	466 0.114	467 23.12	468 0.801	469 0.146

470 Table 3: **Evaluation of novel view synthesis from single view and scene generation on Re10k.**  
 471

#### 472 4.3 ABLATION STUDIES 473

474 We perform ablation studies on the proposed architecture and training strategy to ascertain their ef-  
 475 ficacy. The data used and the hyperparameters applied during the ablation studies remain consistent.  
 476

477 **Effect of two modules in the generation model.** Both the texture module and the geometry module  
 478 can improve understanding performance, but they focus on different aspects. The results presented  
 479 in Table 4 demonstrate that the integration of the texture module facilitates notable advancements in  
 480 tasks dependent on spatiotemporal modeling, exemplified by an increase of 4.1 points in Appr. Order.  
 481 Conversely, the introduction of the geometry module markedly enhances performance in tasks  
 482 contingent upon relative position information, notably in Rel. Dist. However, because of the lacking  
 483 absolute metric in the synthesized depth maps, improvements in tasks pertaining to absolute metric  
 484 comprehension, such as Abs. Dist., are constrained. Incorporating the task of accurately predicting  
 485 camera pose can mitigate the reduction resulting from the imprecise depth prediction. Furthermore,  
 486 our results corroborate that the segregation of the texture and geometry modules results in superior  
 487 understanding performance compared to employing a unified architecture that concurrently learns  
 488 both texture and geometry.  
 489

490 **Effect of the autoregressive generation.** Autoregressive generation significantly improves tasks  
 491 requiring spatiotemporal modeling, like Appr. Order. Table 5 compares the understanding perfor-  
 492 mance with and without autoregressive generation. Although bidirectional generation offers some  
 493

Texture	Geometry		SQA3D		ScanQA		ScanRefer		VSI-Bench (subset)				
	depth	camera	EM	C	B-4	Acc@0.25	Obj. Cnt.	Abs. Dist.	Obj. Size	Rel. Dist.	Appr. Order		
✗	✗	✗	57.2	95.5	14.7	46.9 (28.0)	62.8	36.3	56.4	46.1	43.1		
✓	✗	✗	58.3	97.4	14.7	48.8 (31.5)	67.7	44.6	59.0	53.2	47.2		
✓	✓	✗	58.9	100.5	16.0	48.2 (31.2)	69.0	44.9	69.7	63.0	47.9		
✓	✓	✓	58.7	99.2	15.0	49.0 (31.6)	69.5	45.9	67.8	63.8	48.2		
✓	✓	✓	59.2	103.0	16.2	50.8 (32.5)	70.3	46.4	68.6	65.9	49.0		

Table 4: **Ablation on modules in the generation model.** The gray row denotes that, in this experiment, both the texture module and the geometry module use the same architecture and parameters.

enhancement, autoregressive generation provides greater improvements. Specifically, it boosts performance by 5.8 and 4.4 points in Abs. Dist. and Appr. Order tasks, respectively, compared to when not used. This demonstrates that autoregressive generation strengthens spatiotemporal modeling and enhances related scene understanding tasks.

Autoregressive	SQA3D		ScanQA		ScanRefer		VSI-Bench (subset)				
	EM	B-4	EM	Acc@0.25	Obj. Cnt.	Abs. Dist.	Obj. Size	Rel. Dist.	Appr. Order		
None	57.2	95.5	14.7	46.9 (28.0)	62.8	36.3	56.4	46.1	43.1		
✗	57.2	98.8	15.0	47.0 (28.0)	68.3	40.6	69.0	60.6	44.6		
✓	59.2	103.0	16.2	50.8 (32.5)	70.3	46.4	68.6	65.9	49.0		

Table 5: Ablation on the autoregressive generation in stage 1. “None” means we only train understanding model in this experiment.

**Effect of the D2S (dense-to-sparse).** The efficacy of D2S training is substantiated during stage 1. The results show that D2S significantly improves the model understanding performance, exceeding the visual reconstruction method advanced in Ross3D (Wang et al., 2025a) for scene understanding tasks necessitating spatiotemporal ordering. In this ablation, the geometry module is kept trainable.

Condition in S1	SQA3D		ScanQA		ScanRefer		VSI-Bench (subset)				
	EM	B-4	EM	Acc@0.25	Obj. Cnt.	Abs. Dist.	Obj. Size	Rel. Dist.	Appr. Order		
None	57.2	95.5	14.7	46.9 (28.0)	62.8	36.3	56.4	46.1	43.1		
random mask	58.9	98.9	15.5	49.4 (31.7)	65.7	38.9	55.5	51.6	46.5		
dense	57.3	94.7	15.0	47.2 (28.3)	67.7	42.3	60.1	65.6	46.6		
sparse	57.9	97.3	15.7	50.1 (32.0)	69.0	44.9	65.5	63.1	48.3		
dense → sparse (Ours)	59.2	103.0	16.2	50.8 (32.5)	70.3	46.4	68.6	65.9	49.0		

Table 6: **Ablation on the D2S mechanism in stage 1 (S1).** “None” means we don’t train generation model in this experiment. The “random mask” denotes the visual reconstruction method in Ross3D (Wang et al., 2025a).

**Effect of training stage 2.** The results in Table 7 show that stage 2 can significantly improve the scene generation performance of the model.

Stage 2	NVS from single view			Scene Generation		
	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓
✗	21.90	0.705	0.265	21.44	0.683	0.307
✓	23.22	0.817	0.114	22.93	0.768	0.194

Table 7: **Ablation of the stage 2.**

## 5 CONCLUSION

We introduce Omni-View, a unified 3D scene understanding and generation model. By decomposing the generation model into distinct texture and geometry components, we establish the feasibility and efficacy of employing generation processes to enhance 3D scene understanding and spatial reasoning. Our method illustrates that a unified understanding and generation model is capable of achieving performance on par with leading specialized understanding models. We posit that Omni-View can serve as a foundational model across 3D and multiview domains, thereby advancing the development of downstream applications such as spatial and intelligence.

486 **Ethics Statement.** This paper aims to develop unified models to understand and generate 3D scenes.  
 487 In light of the ongoing advancements in scene generation technology, we emphasize the importance  
 488 of preventing its misapplication, such as the fabrication of deceptive scenes or the creation of scenes  
 489 with nefarious intents.

490 **Reproducibility Statement.** We state that Omni-View is highly reproducible. Implementation  
 491 details on our main experiences are provided in Section 4.1. It is anticipated that these descriptions  
 492 can sufficiently demonstrate the reproducibility of Omni-View. We plan to open-source the code and  
 493 weight files after the paper passes peer review.  
 494

495 **REFERENCES**

496 Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question an-  
 497 swering for spatial scene understanding. In *Proceedings of the IEEE/CVF Conference on Com-*  
 498 *puter Vision and Pattern Recognition (CVPR)*, pp. 19129–19139, 2022.

499 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,  
 500 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang  
 501 Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie,  
 502 Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl tech-  
 503 nical report. *ArXiv*, abs/2502.13923, 2025. URL <https://api.semanticscholar.org/CorpusID:276449796>.  
 504

505 Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitz-  
 506 mann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in*  
 507 *Neural Information Processing Systems*, 37:24081–24125, 2024a.

508 Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. Scanrefer: 3d object localization in  
 509 rgb-d scans using natural language. In *European Conference on Computer Vision (ECCV)*, pp.  
 202–221. Springer, 2020.

510 Juhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi  
 511 Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal  
 512 models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025a.

513 Luxi Chen, Zihan Zhou, Min Zhao, Yikai Wang, Ge Zhang, Wenhao Huang, Hao Sun, Ji-Rong Wen,  
 514 and Chongxuan Li. Flexworld: Progressively expanding 3d scenes for flexible-view synthesis.  
 515 *arXiv preprint arXiv:2503.13265*, 2025b.

516 Yilun Chen, Shuai Yang, Haifeng Huang, Tai Wang, Runsen Xu, Ruiyuan Lyu, Dahua Lin, and  
 517 Jiangmiao Pang. Grounded 3d-llm with referent tokens. *arXiv preprint arXiv:2405.10370*, 2024b.

518 Jaeyoung Chung, Suyoung Lee, Hyeongjin Nam, Jaerin Lee, and Kyoung Mu Lee. Luciddreamer:  
 519 Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384*, 2023.

520 Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Wei-  
 521 hao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified  
 522 multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025.

523 Zhiwen Fan, Jian Zhang, Renjie Li, Junge Zhang, Runjin Chen, Hezhen Hu, Kevin Wang, Huaizhi  
 524 Qu, Dilin Wang, Zhicheng Yan, et al. Vlm-3r: Vision-language models augmented with  
 525 instruction-aligned 3d reconstruction. *arXiv preprint arXiv:2505.20279*, 2025.

526 JiaKui Hu, Yuxiao Yang, Jialun Liu, Jinbo Wu, Chen Zhao, and Yanye Lu. Auto-regressively gen-  
 527 erating multi-view consistent images. *arXiv preprint arXiv:2506.18527*, 2025.

528 Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puahao Li, Yan Wang, Qing Li,  
 529 Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world.  
 530 *arXiv preprint arXiv:2311.12871*, 2023.

531 Tianyu Huang, Wangguandong Zheng, Tengfei Wang, Yuhao Liu, Zhenwei Wang, Junta Wu, Jie  
 532 Jiang, Hui Li, Rynson WH Lau, Wangmeng Zuo, et al. Voyager: Long-range and world-consistent  
 533 video diffusion for explorable 3d scene generation. *arXiv preprint arXiv:2506.04225*, 2025.

540 Georg B Keller, Tobias Bonhoeffer, and Mark Hübener. Sensorimotor mismatch signals in primary  
 541 visual cortex of the behaving mouse. *Neuron*, 74(5):809–815, 2012.

542

543 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.

544

545 Marcus Leinweber, Daniel R Ward, Jan M Sobczak, Alexander Attinger, and Georg B Keller. A  
 546 sensorimotor circuit in mouse cortex for visual flow predictions. *Neuron*, 95(6):1420–1432, 2017.

547

548 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
 549 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

550

551 Baorui Ma, Huachen Gao, Haoge Deng, Zhengxiong Luo, Tiejun Huang, Lulu Tang, and Xinlong  
 552 Wang. You see it, you got it: Learning 3d creation on pose-free videos at scale. In *Proceedings  
 553 of the Computer Vision and Pattern Recognition Conference*, pp. 2016–2029, 2025.

554

555 Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao Liang, Song-Chun Zhu, and Siyuan Huang.  
 556 Sq3d: Situated question answering in 3d scenes. In *International Conference on Learning Rep-  
 557 resentations (ICLR)*, 2023.

558

559 Gerrit W Maus, Jason Fischer, and David Whitney. Motion-dependent representation of space in  
 560 area mt+. *Neuron*, 78(3):554–562, 2013.

561

562 Nora Nortmann, Sascha Rekauzke, Selim Onat, Peter König, and Dirk Jancke. Primary visual cortex  
 563 represents the difference between past and present. *Cerebral Cortex*, 25(6):1427–1440, 2015.

564

565 OpenAI. Hello GPT-4o, 2024. URL <https://openai.com/index/hello-gpt-4o>.

566

567 Zhangyang Qi, Zhixiong Zhang, Ye Fang, Jiaqi Wang, and Hengshuang Zhao. Gpt4scene: Un-  
 568 derstand 3d scenes from videos with vision-language models. *arXiv preprint arXiv:2501.01428*,  
 569 2025.

570

571 Arijit Ray, Jiafei Duan, Ellis Brown, Reuben Tan, Dina Bashkirova, Rose Hendrix, Kiana Ehsani,  
 572 Aniruddha Kembhavi, Bryan A Plummer, Ranjay Krishna, et al. Sat: Dynamic spatial aptitude  
 573 training for multimodal language models. *arXiv preprint arXiv:2412.07755*, 2024.

574

575 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
 576 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
 577 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.

578

579 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint  
 580 arXiv:2405.09818*, 2024a.

581

582 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,  
 583 Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal under-  
 584 standing across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.

585

586 HunyuanWorld Team, Zhenwei Wang, Yuhao Liu, Junta Wu, Zixiao Gu, Haoyuan Wang, Xuhui  
 587 Zuo, Tianyu Huang, Wenhuan Li, Sheng Zhang, et al. Hunyuanworld 1.0: Generating immersive,  
 588 exploratory, and interactive 3d worlds from words or pixels. *arXiv preprint arXiv:2507.21809*,  
 589 2025.

590

591 OpenGVLab Team. InternVL2: Better than the Best—Expanding Performance Boundaries of  
 592 Open-Source Multimodal Models with the Progressive Scaling Strategy, 2024b. URL <https:////internvl.github.io/blog/2024-07-02-InternVL-2.0/>.

593

594 Shengbang Tong, David Fan, Jiachen Zhu, Yunyang Xiong, Xinlei Chen, Koustuv Sinha, Michael  
 595 Rabbat, Yann LeCun, Saining Xie, and Zhuang Liu. Metamorph: Multimodal understanding and  
 596 generation via instruction tuning. *arXiv preprint arXiv:2412.14164*, 2024.

597

598 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,  
 599 Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative  
 600 models. *arXiv preprint arXiv:2503.20314*, 2025.

594 Haochen Wang, Yucheng Zhao, Tiancai Wang, Haoqiang Fan, Xiangyu Zhang, and Zhaoxiang  
 595 Zhang. Ross3d: Reconstructive visual instruction tuning with 3d-awareness. *arXiv preprint*  
 596 *arXiv:2504.01901*, 2025a.

597 Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David  
 598 Novotny. Vggt: Visual geometry grounded transformer. In *Proceedings of the IEEE/CVF Con-*  
 599 *ference on Computer Vision and Pattern Recognition*, 2025b.

600 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,  
 601 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the  
 602 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024a.

603 Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo,  
 604 and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. In  
 605 *ACM SIGGRAPH 2024 Conference Papers*, pp. 1–11, 2024b.

606 Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu,  
 607 Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified  
 608 multimodal understanding and generation. *arXiv preprint arXiv:2410.13848*, 2024a.

609 Diankun Wu, Fangfu Liu, Yi-Hsin Hung, and Yueqi Duan. Spatial-mllm: Boosting mllm capabilities  
 610 in visual-based spatial intelligence. *arXiv preprint arXiv:2505.23747*, 2025a.

611 Junfeng Wu, Yi Jiang, Chuofan Ma, Yuliang Liu, Hengshuang Zhao, Zehuan Yuan, Song Bai, and  
 612 Xiang Bai. Liquid: Language models are scalable and unified multi-modal generators. *arXiv*  
 613 *preprint arXiv:2412.04332*, 2024b.

614 Size Wu, Wenwei Zhang, Lumin Xu, Sheng Jin, Zhonghua Wu, Qingyi Tao, Wentao Liu, Wei Li,  
 615 and Chen Change Loy. Harmonizing visual representations for unified multimodal understanding  
 616 and generation. *arXiv preprint arXiv:2503.21979*, 2025b.

617 Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang, Dacheng Li, Yunhao Fang, Ligeng  
 618 Zhu, Enze Xie, Hongxu Yin, Li Yi, et al. Vila-u: a unified foundation model integrating visual  
 619 understanding and generation. *arXiv preprint arXiv:2409.04429*, 2024c.

620 Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin,  
 621 Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer  
 622 to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024.

623 Jinheng Xie, Zhenheng Yang, and Mike Zheng Shou. Show-o2: Improved native unified multimodal  
 624 models. *arXiv preprint arXiv:2506.15564*, 2025.

625 Fuzhao Xue, Yukang Chen, Dacheng Li, Qinghao Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian  
 626 Tang, Shang Yang, Zhijian Liu, Ethan He, Hongxu Yin, Pavlo Molchanov, Jan Kautz, Linxi Fan,  
 627 Yuke Zhu, Yao Lu, and Song Han. Longvila: Scaling long-context visual language models for  
 628 long videos. *ArXiv*, abs/2408.10188, 2024. URL <https://api.semanticscholar.org/CorpusID:271903601>.

629 Zhiyuan Yan, Kaiqing Lin, Zongjian Li, Junyan Ye, Hui Han, Zhendong Wang, Hao Liu, Bin Lin,  
 630 Hao Li, Xue Xu, et al. Can understanding and generation truly benefit together—or just coexist?  
 631 *arXiv preprint arXiv:2509.09666*, 2025.

632 Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in  
 633 space: How multimodal large language models see, remember, and recall spaces. *arXiv preprint*  
 634 *arXiv:2412.14171*, 2024.

635 Hong-Xing Yu, Haoyi Duan, Charles Herrmann, William T Freeman, and Jiajun Wu. Wonderworld:  
 636 Interactive 3d scene generation from a single image. *arXiv preprint arXiv:2406.09394*, 2024a.

637 Hong-Xing Yu, Haoyi Duan, Junhwa Hur, Kyle Sargent, Michael Rubinstein, William T Freeman,  
 638 Forrester Cole, Deqing Sun, Noah Snavely, Jiajun Wu, et al. Wonderjourney: Going from any-  
 639 where to everywhere. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
 640 *Pattern Recognition*, pp. 6658–6667, 2024b.

648 Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao, Tien-  
 649 Tsin Wong, Ying Shan, and Yonghong Tian. Viewcrafter: Taming video diffusion models for  
 650 high-fidelity novel view synthesis. *arXiv preprint arXiv:2409.02048*, 2024c.  
 651

652 Shangjin Zhai, Zhichao Ye, Jialin Liu, Weijian Xie, Jiaqi Hu, Zhen Peng, Hua Xue, Danpeng Chen,  
 653 Xiaomeng Wang, Lei Yang, et al. Stargen: A spatiotemporal autoregression framework with video  
 654 diffusion model for scalable and controllable scene generation. In *Proceedings of the Computer  
 Vision and Pattern Recognition Conference*, pp. 26822–26833, 2025.  
 655

656 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language  
 657 image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*,  
 658 pp. 11975–11986, 2023.  
 659

660 Jiahui Zhang, Yurui Chen, Yanpeng Zhou, Yueming Xu, Ze Huang, Jilin Mei, Junhui Chen, Yu-  
 661 Jie Yuan, Xinyue Cai, Guowei Huang, et al. From flatland to space: Teaching vision-language  
 662 models to perceive and reason in 3d. *arXiv preprint arXiv:2503.22976*, 2025a.  
 663

664 Jiawei Zhang, Chejian Xu, and Bo Li. Chatscene: Knowledge-enabled safety-critical scenario gener-  
 665 ation for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
 and Pattern Recognition (CVPR)*, pp. 15459–15469, 2024a.  
 666

667 Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue  
 668 Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to  
 669 vision. *ArXiv*, abs/2406.16852, 2024b. URL <https://api.semanticscholar.org/CorpusID:270703489>.  
 670

671 Xinjie Zhang, Jintao Guo, Shanshan Zhao, Minghao Fu, Lunhao Duan, Jiakui Hu, Yong Xien Chng,  
 672 Guo-Hua Wang, Qing-Guo Chen, Zhao Xu, et al. Unified multimodal understanding and genera-  
 673 tion models: Advances, challenges, and opportunities. *arXiv preprint arXiv:2505.02567*, 2025b.  
 674

675 Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video  
 676 instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024c.  
 677

678 Duo Zheng, Shijia Huang, and Liwei Wang. Video-3d llm: Learning position-aware video repre-  
 679 sentation for 3d scene understanding. *arXiv preprint arXiv:2412.00493*, 2024.  
 680

681 Duo Zheng, Shijia Huang, Yanyang Li, and Liwei Wang. Learning from videos for 3d world:  
 682 Enhancing mllms with 3d vision geometry priors. *arXiv preprint arXiv:2505.24625*, 2025.  
 683

684 Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification:  
 685 Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018.  
 686

687 Xin Zhou, Dingkang Liang, Sifan Tu, Xiwu Chen, Yikang Ding, Dingyuan Zhang, Feiyang Tan,  
 688 Hengshuang Zhao, and Xiang Bai. Hermes: A unified self-driving world model for simultaneous  
 689 3d scene understanding and generation. *arXiv preprint arXiv:2501.14729*, 2025.  
 690

691 Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang, and Xihui Liu. Llava-3d: A simple  
 692 yet effective pathway to empowering lmms with 3d-awareness. *arXiv preprint arXiv:2409.18125*,  
 693 2024.  
 694

695

696

697

698

699

700

701

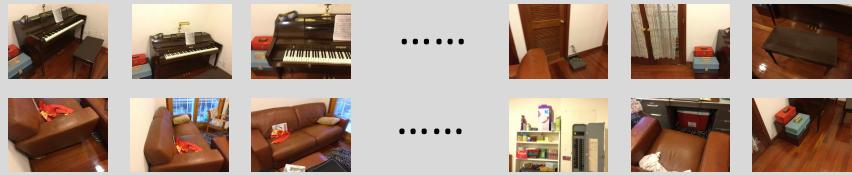
702 **Limitations and future work.** Although Omni-View has demonstrated a unified advance in the  
 703 domains of question answering, spatial reasoning, and novel view synthesis, its grounding capabilities  
 704 remain to be substantiated. Additionally, its generation model currently lacks the capability for  
 705 long-range world generation. Future efforts will be concentrated on using reinforcement learning  
 706 to augment Omni-View’s performance in 3D visual grounding and long-range generation. At the  
 707 same time, since the training data used by its geometry module is a synthetic depth map, its actual  
 708 geometry prediction ability may not be accurate enough.

## 710 A QUALITATIVE RESULTS

712 We show some qualitative results for 3D scene understanding and generation tasks in the appendix.

### 714 A.1 3D SCENE UNDERSTANDING

#### 718 Multiview images

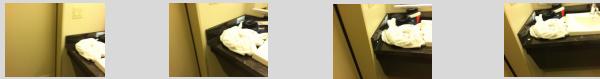


725 Question: I am standing by the ottoman on my right facing a couple of toolboxes.

726 What instrume in fro of me is ebony and ivory?

727 Answer: **piano**

#### 730 Multiview images



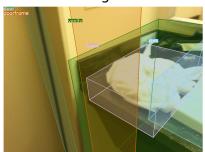
734 Question: Detect the 3D bounding boxes in the camera coordinate system of the first frame.

735 Answer: `{"label": "door", "bbox_3d": [-0.29, 0.21, 0.63, 0.07, 0.6, 1.69, 2.0, 1.14, -3.09]},`  
 736 `{"label": "vanity", "bbox_3d": [0.6, 0.47, 0.53, 0.59, 1.36, 0.7, 0.43, 1.14, -3.09]},`  
 737 `{"label": "towel", "bbox_3d": [0.31, 0.07, 0.79, 0.25, 0.35, 0.09, 2.0, 1.14, -3.09]},`  
 738 `{"label": "sink", "bbox_3d": [0.57, 0.21, 0.53, 0.42, 0.63, 0.15, 0.43, 1.14, -3.09]},`  
 739 `{"label": "doorframe", "bbox_3d": [-0.28, 0.19, 0.59, 0.2, 0.85, 1.63, 2.0, 1.14, -3.09]},`  
 740 `{"label": "switch", "bbox_3d": [0.16, -0.23, 0.66, 0.02, 0.08, 0.11, 2.0, 1.14, -3.09]},`  
 741 `{"label": "mirror", "bbox_3d": [0.86, -0.14, 0.44, 0.08, 1.06, 0.13, 0.43, 1.14, -3.09]}`

741 Image 0



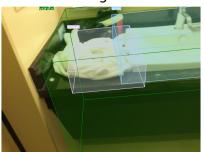
742 Image 1



743 Image 2

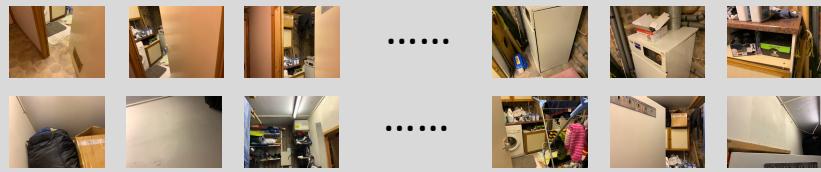


744 Image 3

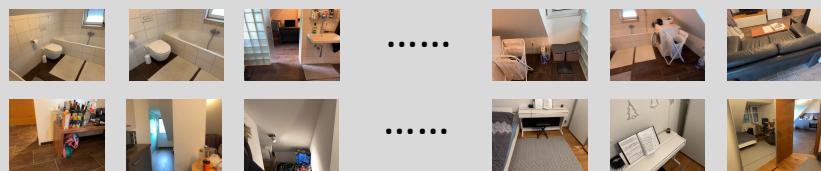


756 A.2 SPATIAL REASONING  
757  
758  
759

## 760 Multiview images

761 Question: What is the size of this room (in square meters)?  
762  
763  
764  
765  
766767 If multiple rooms are shown, estimate the size of the combined space.  
768  
769770 Answer: **10.3**771 Ground Truth: **10.5**772 (a) Results on room size estimation.  
773  
774  
775  
776  
777  
778

## 779 Multiview images

780 Question: What will be the first-time appearance order of the following categories in the video:  
781  
782  
783  
784785 A. ceiling light, heater, tv, printer  
786 B. ceiling light, tv, heater, printer  
787 C. heater, tv, ceiling light, printer  
788 D. ceiling light, heater, printer, tv789 Answer: **B**790 (b) Results on appearance order.  
791  
792  
793  
794  
795  
796  
797  
798

## 799 Multiview images

800 Question: Measuring from the closest point of each object,  
801 which of these objects (stove, tv, table, sofa) is the closest to the stool?  
802 Options:803 A. stove  
804 B. tv  
805 C. table  
806 D. sofa807 Answer with the option's letter from the given choices directly.  
808  
809810 Answer: **D.**

811 GT: D.

812 (c) Object relative distance.  
813  
814  
815  
816  
817  
818  
819

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

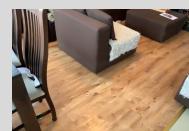
860

861

862

863

### Multiview images



Question: How many chair(s) are in this room?

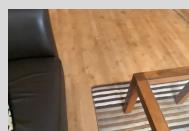
Please answer the question using a single word or phrase.

Answer: **4.**

GT: 4.

(d) Object counting.

### Multiview images



Question: What is the length of the longest dimension (length, width, or height) of the sofa, measured in centimeters?

Please answer the question using a single word or phrase.

Answer: **184.**

GT: 173.

(e) Object size estimation.

### Multiview images



Question: Measuring from the closest point of each object,

what is the distance between the tv and the bed (in meters)?

Please answer the question using a single word or phrase.

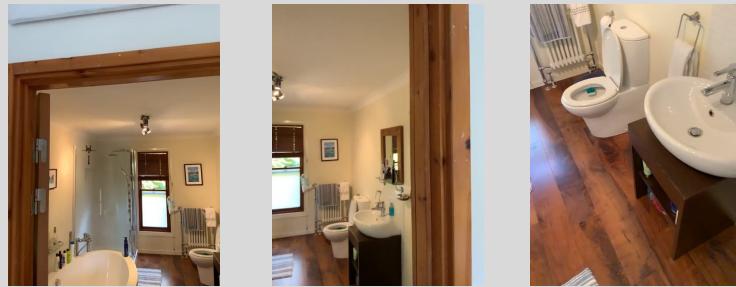
Answer: **1.0.**

GT: 1.1.

(f) Object absolute distance.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876

Multiview images



877 **Question:** If I am standing by the bathtub and facing the toilet,  
878 is the table to my front-left, front-right, back-left, or back-right?  
879 The directions refer to the quadrants of a Cartesian plane  
880 (if I am standing at the origin and facing along the positive y-axis).  
881 Options:  
882 A. back-left  
883 B. front-right  
884 C. front-left  
885 D. back-right  
886 Answer with the option's letter from the given choices directly.

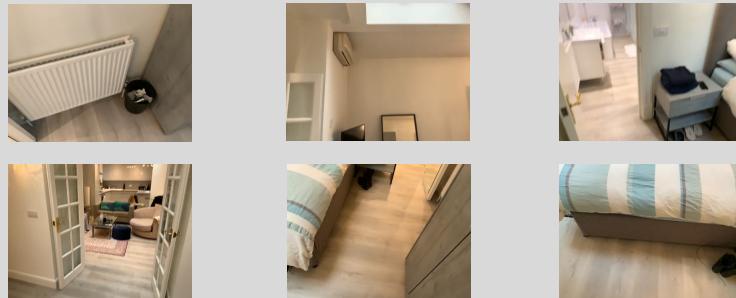
Answer: **B.**

GT: B.

887 (g) Object relative direction.

888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899

Multiview images



900 **Question:** You are a robot beginning at the tv facing the bed.  
901 You want to navigate to the trash bin.  
902 You will perform the following actions  
903 (Note: for each [please fill in], choose either 'turn back,' 'turn left,' or 'turn right.'):  
904 1. [please fill in]  
905 2. Go forward until the cabinet  
906 3. [please fill in]  
907 4. Go forward until the trash bin is on your right.  
908 You have reached the final destination.  
909 Options:  
910 A. Turn Left, Turn Left  
911 B. Turn Right, Turn Left  
912 C. Turn Back, Turn Left  
913 D. Turn Right, Turn Right  
914 Answer with the option's letter from the given choices directly.

Answer: **B.**

GT: B.

(h) Route planning.

915  
916  
917

918 A.3 SCALE CONSISTENCY ACROSS SCENES  
919

920 We analyze the scale consistency across different scenes. To this end, we selected some test samples  
921 from the SPAR-7M dataset and used Omni-View to perform absolute distance prediction based on  
922 multi-view images. This SPAR-7M consists of scenes from three datasets: ScanNet and ScanNet++,  
923 where the average scene depth ranges from 3 to 5 meters with a maximum depth of 12 meters;  
924 and Structured3D, which features larger-scale environments with average depths between 4 and 6  
925 meters and maximum depths reaching up to 20 meters. Qualitative visualizations demonstrate that  
926 Omni-View can predict absolute depth in scenes with different metric scale.

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

## Multiview images



Question: The Object63 (red point at Frame-0, [0.14, 0.67]) is positioned at 3.4 meters depth.

Predict the depth of the Object22 (blue point at Frame-0, [0.23, 0.65]).

Calculate or judge based on the 3D center points of these objects.

Answer: **Object22's central depth is estimated to be about 2.0 meters.**

GT: With a central depth of 2.2 meters, Object22 is referenced here.

(a) Absolute depth estimation in structured3d, scene\_03201\_80.

## Multiview images



Question: The Object26 (at Frame-1, [0.10, 0.39]) lies at 2.7 meters depth.

Predict the depth of the Object30 (at Frame-1, [0.10, 0.70]).

Calculate or judge based on the 3D center points of these objects.

Answer: **Object30 is centered at about 4.0 meters depth.**

GT: At its center, object Object30 has a depth of 4.3 meters.

(b) Absolute depth estimation in structured3d, scene\_02053\_25.

## Multiview images



Question: If the plant (Frame-1, [0.30, 0.41]) is positioned at 2.7 meters,

estimate the depth of the laptop (at Frame-1, [0.16, 0.63]).

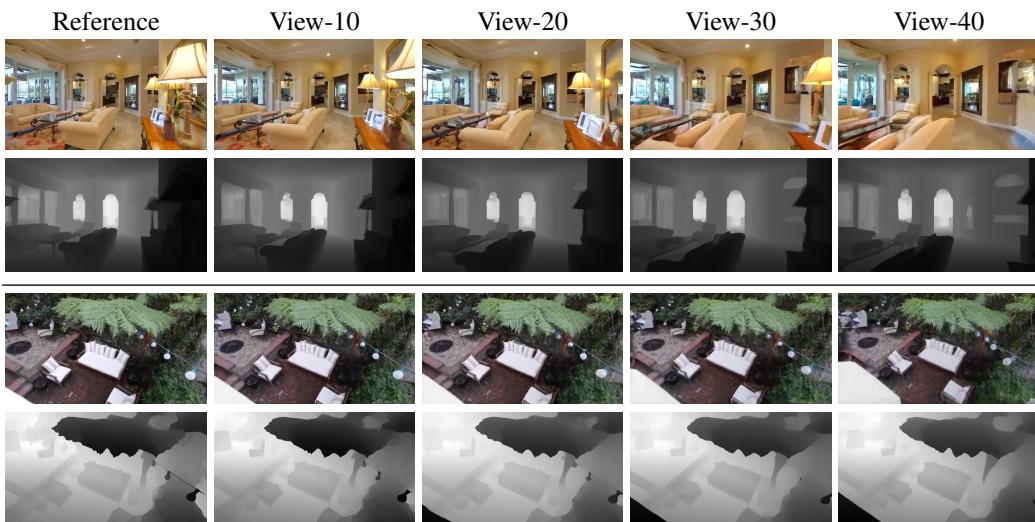
Calculate or judge based on the 3D center points of these objects.

The answer should be a single number, nothing more. .

Answer: **2.4**

GT: 2.3

(c) Absolute depth estimation in scannetpp, 9d7e20fbec\_6022.

972 A.4 NVS FROM SINGLE VIEW  
973974 **Indoor scenes.** We present the two visual results of NVS in indoor scenes. Omni-View can reasonably  
975 imagine unseen areas while maintaining the texture and structure of observed objects, like the  
976 followers on the table in the first row.  
977986 **Outdoor scenes.** We present the two visual results of NVS in outdoor scenes. When the camera  
987 movement is small, Omni-View can consistently generate new views.  
988997 **Depth estimation.** We present the depth prediction results of Omni-View in indoor and outdoor  
998 scenes. In indoor scenes, the output of Omni-View's geometry module is more accurate.  
9991016 **Failure case.** We claim that the existing Omni-View is inadequate for effectively processing outdoor  
1017 scenes with substantial camera movement. As illustrated in the figure below, the result images  
1018 from Omni-View exhibit significant artifacts, highlighted within the red dashed boxes. Future work  
1019 will focus on resolving the following challenges to enhance the handling of outdoor scenes with  
1020 extensive camera motion: (1) the development of more precise camera control mechanisms, and (2)  
1021 the improvement of inter-frame texture consistency stability.  
1022

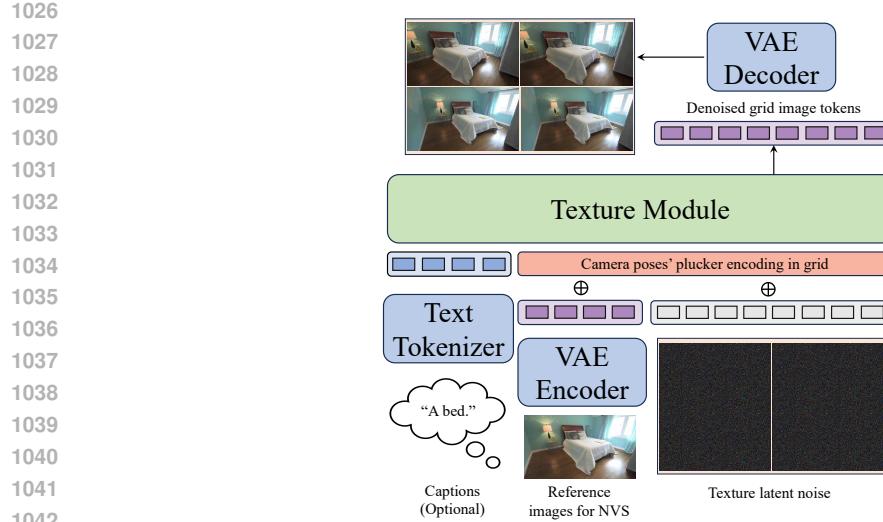
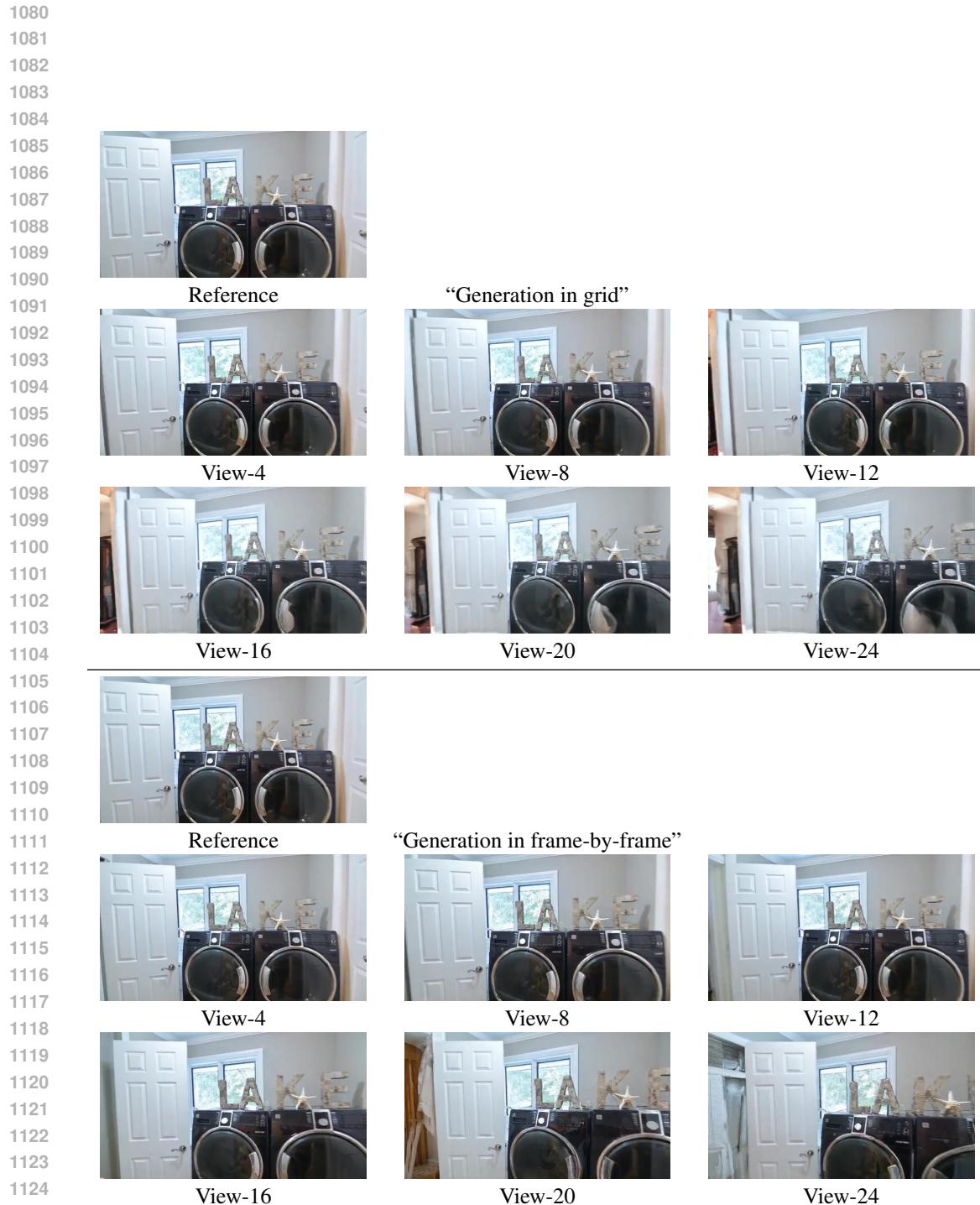


Figure 2: Generation in grid.

## A.5 EFFECT OF GENERATION IN GRID

**Generation in grid.** We explored small improvements to enhance inter-frame consistency and support long-sequence scene generation: “generation in grid”. Specifically, we arrange 4 consecutive views into a single frame by organizing them in a grid layout, and perform autoregressive generation over the resulting sequence of grid-organized frames.

**Results.** As shown in the figures on the next page, we demonstrate that after using “generation in grid”, Omni-View can generate more consistent novel view images.



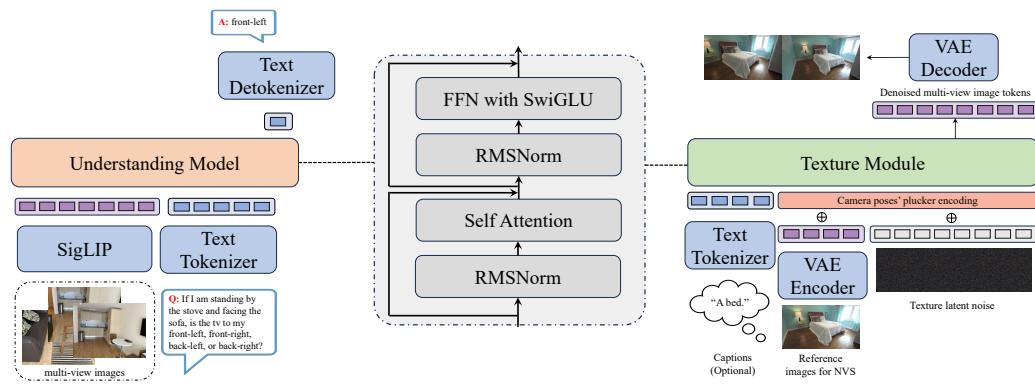
1134 **B TECHNICAL DETAILS**  
11351136 **B.1 DATASET**  
1137

1138 **3D scene understanding and spatial reasoning.** We curate a filtered training dataset containing 1139 780k valid samples by consolidating data from multiple sources, including SQA3D, ScanQA, 1140 3DOD, ScanRefer, VLM-3R, a 234k subset of SPAR from VG-LLM, and a 64k subset of 1141 llava-hound4 from VG-LLM. To ensure data quality, we perform two main filtering steps: (i) deduplication 1142 across the combined dataset, and (ii) removal of samples with invalid bounding box annotations. 1143 Specifically, we convert all bounding boxes to the  $[x_1, y_1, x_2, y_2]$  format and exclude any instance 1144 where  $x_1 < 0, y_1 < 0, x_2 > \text{width}$ , or  $y_2 > \text{height}$ .  
1145

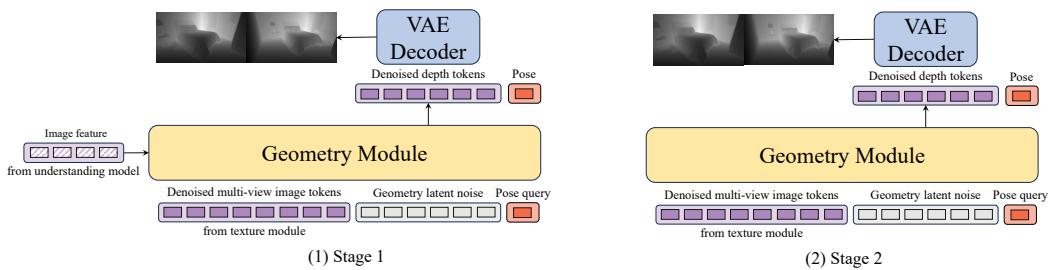
1146 **Novel view synthesis.** We select 61k video clips from RealE10K train set, excluding those with 1147 fewer than 32 frames or significant motion blur that could degrade training stability. Depth maps are 1148 synthesized using the Voyager data pipeline, while captions are generated via QwenVLMMax.  
1149

1150 **B.2 ARCHITECTURE.**  
1151

1152 **Understanding model.** The detailed architecture is shown in Figure 4. The text tokenizer inherits 1153 from the vocabulary used by Qwen2. The image tokenizer for understanding is SigLIP. The backbone 1154 has layers as 28, attention head as 28, and hidden size as 3584, totaling 7B parameters. 1155 Key architectural choices include SwiGLU as the activation function, RMSNorm for normalization, 1156 FlexAttention (via PyTorch) for efficient self-attention computation, and position encoding adapted 1157 from Qwen2.  
1158

1160 Figure 4: Architecture of understanding model and texture module.  
1161

1162 **Texture module.** This module uses FLUX-VAE as the image tokenizer, sharing the same backbone 1163 architecture as the understanding model.  
1164

1165 Figure 5: Architecture of geometry module in two stages.  
1166

1167 **Geometry module.** The detailed architecture is shown in Figure 5. Depth maps are also encoded 1168 by FLUX-VAE. Its backbone is basically the same as the block architecture, but reduced to 4 layers, 1169 containing about 1B parameters. Unlike the others, the architecture of the geometry module changes 1170 across different training stages. In stage 1, it receives features from the understanding model for 1171

1188 cross-attention. In stage 2, it no longer receives outputs from the understanding model. The reason  
 1189 for this design can be found in the discussion of stage 2 in Section 3.2.  
 1190

1191 **B.3 TRAINING AND COMPUTATIONAL COST.**  
 1192

1193 All training phases use the AdamW optimizer with  $\beta_1 = 0.9, \beta_2 = 0.95$ , a peak learning rate of  
 1194  $1 \times 10^{-5}$ , and a linear warm-up schedule covering the first 5% of iterations.  
 1195

1196 **Stage 1.** We train on the 3D scene understanding dataset for 10,000 iterations with a packed se-  
 1197 quence length of 50k, using 32 H100 GPUs, which takes approximately 160 hours.  
 1198

1199 **Stage 2.** For the generation task, we perform 20,000 iterations on the novel view synthesis dataset  
 1200 with a packed sequence length of 32k, using 32 H100 GPUs, requiring approximately 40 hours in  
 1201 total.  
 1202

1203 **Understanding inference.** The model takes approximately 2.5 seconds on average to process a  
 1204 32-frame multi-view scene understanding query using a single H100 GPU.  
 1205

1206 **Generation inference.** Generating a single image at resolution  $640 \times 352$  takes about 2.2 seconds  
 1207 on average with one H100 GPU.  
 1208

1209 **B.4 CONVERGENCE BEHAVIOR.**  
 1210

1211 **Stage 1.** Overall, most losses exhibit smooth and stable convergence. However, we observe spikes in  
 1212 the camera pose loss during training, particularly in later epochs. We hypothesize that this behavior  
 1213 may stem from instable optimizion by learnable queries in the camera pose estimation task. How-  
 1214 ever, these fluctuations do not prevent the model from converging to effective solutions, as evidenced  
 1215 by strong performance in 3D scene understanding, spatial reasoning, and novel view generation.  
 1216

1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

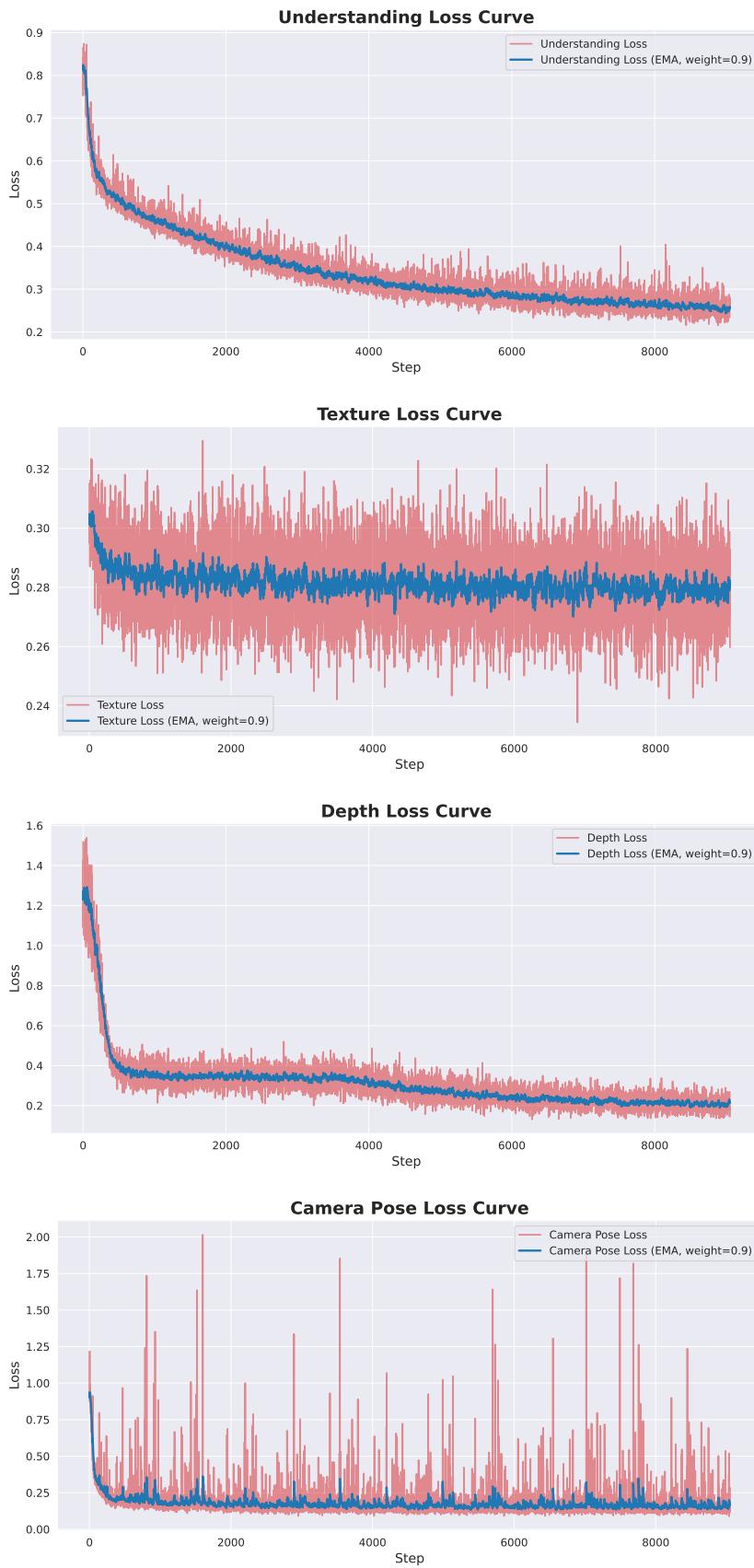
1291

1292

1293

1294

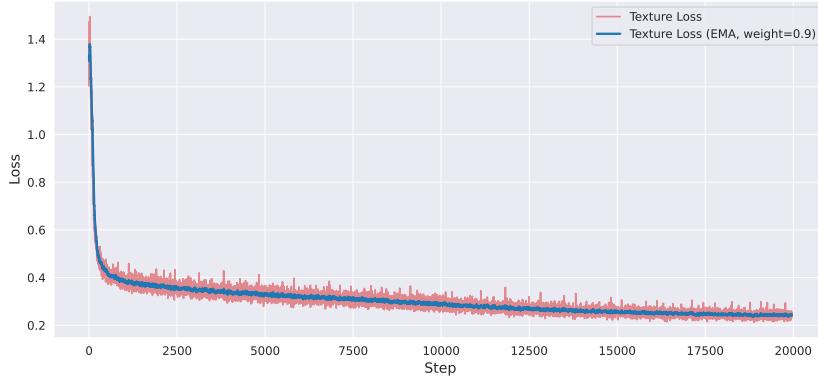
1295



1296 **Stage 2.** Since we used a grid generation format in stage 2, the per-image prediction method in stage  
 1297 1 cannot be directly transferred to the “generation in grid” paradigm in stage 2. However, due to  
 1298 the generative prior learned in stage 1, Omni-View can quickly converge to a lower loss in stage 2.  
 1299 However, it was observed that the geometry loss showed a slight increase at around 8000 iterations,  
 1300 followed by a decrease, suggesting that it might be possible to further reduce  $\lambda_{geo}$  to improve the  
 1301 training stability.

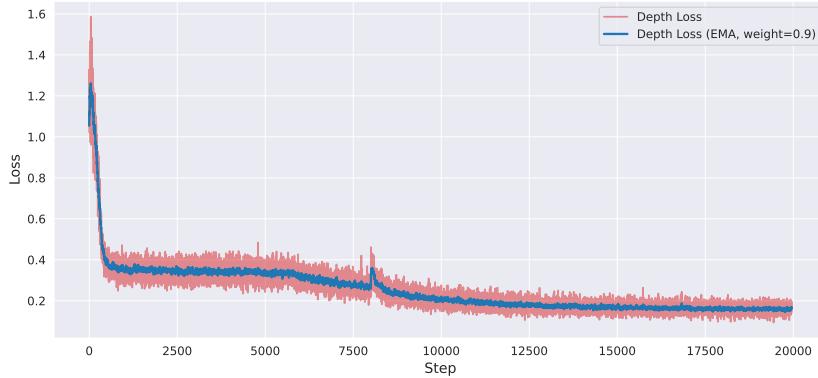
1302

1303

**Texture Loss Curve**

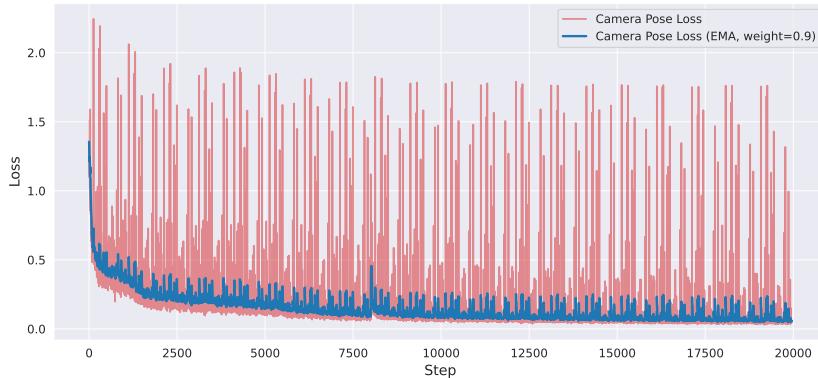
1316

1317

**Depth Loss Curve**

1329

1330

**Camera Pose Loss Curve**

1343

1344

## B.5 ACTIVATION VISUALIZATION

1346

We visualize the activation maps of Bagel and Omni-View when they perform spatial reasoning tasks, as shown in the figure below. In this example, we want the model to locate how many cabinets are in the current scene, using the prompt: “How many cabinet(s) are in this room?”. It can be seen that Bagel mainly focuses on the first two images that are irrelevant to

1350 the question, whereas Omni-View is able to attend to each image as much as possible. The wider  
 1351 attention may be the reason why Omni-View performs better on 3D scene understanding tasks.  
 1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403



Figure 6: Activation map. Bagel-FT vs. Omni-View.

1404 **C THE USE OF LARGE LANGUAGE MODELS (LLMs)**  
1405

1406 LLMs are used to correct potential grammatical inaccuracies in the manuscript. LLMs do not par-  
1407 ticipate in research ideation.  
1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457