

000 001 002 003 004 005 006 007 UNI4D-LLM: A UNIFIED SPATIOTEMPORAL-AWARE 008 VLM FOR 4D UNDERSTANDING AND GENERATION 009 010 011 —Supplementary Material— 012 013 014 015 016 017

018 Anonymous authors 019 Paper under double-blind review 020 021 022 023 024 025 026

In this supplementary material, we provide additional technical details and experimental validations on our training details, method, and comparison. For our training details, we introduce the detailed training pipeline and datasets in Sec. 1. For our method, we conduct extensive experiments on ablation study and discussion in Sec. 2, including applicability to different camera settings in Sec. 2.1, task prompt in Sec. 2.2, textual coordinate embedding in Sec. 2.3, and the potential solution to the limitations in Sec. 2.4. Regarding our comparison, we provide more quantitative and qualitative results of 2D/3D/4D scene understanding and generation in Sec. 3.1-3.3.

027 1 DETAILS OF TRAINING PIPELINE AND DATASET 028 029 030 031 032 033 034 035

In Table 1, we provide details of our training datasets and training strategy. We integrate diverse 2D, 3D, and 4D vision-language datasets for both understanding and generation, which are utilized at different training stages. The entire training pipeline is divided into three stages as follows:

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
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
0100
0101
0102
0103
0104
0105
0106
0107
0108
0109
0110
0111
0112
0113
0114
0115
0116
0117
0118
0119
0120
0121
0122
0123
0124
0125
0126
0127
0128
0129
0130
0131
0132
0133
0134
0135
0136
0137
0138
0139
0140
0141
0142
0143
0144
0145
0146
0147
0148
0149
0150
0151
0152
0153
0154
0155
0156
0157
0158
0159
0160
0161
0162
0163
0164
0165
0166
0167
0168
0169
0170
0171
0172
0173
0174
0175
0176
0177
0178
0179
0180
0181
0182
0183
0184
0185
0186
0187
0188
0189
0190
0191
0192
0193
0194
0195
0196
0197
0198
0199
0200
0201
0202
0203
0204
0205
0206
0207
0208
0209
0210
0211
0212
0213
0214
0215
0216
0217
0218
0219
0220
0221
0222
0223
0224
0225
0226
0227
0228
0229
0230
0231
0232
0233
0234
0235
0236
0237
0238
0239
0240
0241
0242
0243
0244
0245
0246
0247
0248
0249
0250
0251
0252
0253
0254
0255
0256
0257
0258
0259
0260
0261
0262
0263
0264
0265
0266
0267
0268
0269
0270
0271
0272
0273
0274
0275
0276
0277
0278
0279
0280
0281
0282
0283
0284
0285
0286
0287
0288
0289
0290
0291
0292
0293
0294
0295
0296
0297
0298
0299
0300
0301
0302
0303
0304
0305
0306
0307
0308
0309
0310
0311
0312
0313
0314
0315
0316
0317
0318
0319
0320
0321
0322
0323
0324
0325
0326
0327
0328
0329
0330
0331
0332
0333
0334
0335
0336
0337
0338
0339
0340
0341
0342
0343
0344
0345
0346
0347
0348
0349
0350
0351
0352
0353
0354
0355
0356
0357
0358
0359
0360
0361
0362
0363
0364
0365
0366
0367
0368
0369
0370
0371
0372
0373
0374
0375
0376
0377
0378
0379
0380
0381
0382
0383
0384
0385
0386
0387
0388
0389
0390
0391
0392
0393
0394
0395
0396
0397
0398
0399
0400
0401
0402
0403
0404
0405
0406
0407
0408
0409
0410
0411
0412
0413
0414
0415
0416
0417
0418
0419
0420
0421
0422
0423
0424
0425
0426
0427
0428
0429
0430
0431
0432
0433
0434
0435
0436
0437
0438
0439
0440
0441
0442
0443
0444
0445
0446
0447
0448
0449
0450
0451
0452
0453
0454
0455
0456
0457
0458
0459
0460
0461
0462
0463
0464
0465
0466
0467
0468
0469
0470
0471
0472
0473
0474
0475
0476
0477
0478
0479
0480
0481
0482
0483
0484
0485
0486
0487
0488
0489
0490
0491
0492
0493
0494
0495
0496
0497
0498
0499
0500
0501
0502
0503
0504
0505
0506
0507
0508
0509
0510
0511
0512
0513
0514
0515
0516
0517
0518
0519
0520
0521
0522
0523
0524
0525
0526
0527
0528
0529
0530
0531
0532
0533
0534
0535
0536
0537
0538
0539
0540
0541
0542
0543
0544
0545
0546
0547
0548
0549
0550
0551
0552
0553
0554
0555
0556
0557
0558
0559
0560
0561
0562
0563
0564
0565
0566
0567
0568
0569
0560
0561
0562
0563
0564
0565
0566
0567
0568
0569
0570
0571
0572
0573
0574
0575
0576
0577
0578
0579
0580
0581
0582
0583
0584
0585
0586
0587
0588
0589
0590
0591
0592
0593
0594
0595
0596
0597
0598
0599
0600
0601
0602
0603
0604
0605
0606
0607
0608
0609
0610
0611
0612
0613
0614
0615
0616
0617
0618
0619
0620
0621
0622
0623
0624
0625
0626
0627
0628
0629
0630
0631
0632
0633
0634
0635
0636
0637
0638
0639
0640
0641
0642
0643
0644
0645
0646
0647
0648
0649
0650
0651
0652
0653
0654
0655
0656
0657
0658
0659
0660
0661
0662
0663
0664
0665
0666
0667
0668
0669
0660
0661
0662
0663
0664
0665
0666
0667
0668
0669
0670
0671
0672
0673
0674
0675
0676
0677
0678
0679
0680
0681
0682
0683
0684
0685
0686
0687
0688
0689
0690
0691
0692
0693
0694
0695
0696
0697
0698
0699
0700
0701
0702
0703
0704
0705
0706
0707
0708
0709
07010
07011
07012
07013
07014
07015
07016
07017
07018
07019
07020
07021
07022
07023
07024
07025
07026
07027
07028
07029
07030
07031
07032
07033
07034
07035
07036
07037
07038
07039
07040
07041
07042
07043
07044
07045
07046
07047
07048
07049
07050
07051
07052
07053
07054
07055
07056
07057
07058
07059
07060
07061
07062
07063
07064
07065
07066
07067
07068
07069
07070
07071
07072
07073
07074
07075
07076
07077
07078
07079
07080
07081
07082
07083
07084
07085
07086
07087
07088
07089
07090
07091
07092
07093
07094
07095
07096
07097
07098
07099
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140
070141
070142
070143
070144
070145
070146
070147
070148
070149
070150
070151
070152
070153
070154
070155
070156
070157
070158
070159
070160
070161
070162
070163
070164
070165
070166
070167
070168
070169
070170
070171
070172
070173
070174
070175
070176
070177
070178
070179
070180
070181
070182
070183
070184
070185
070186
070187
070188
070189
070190
070191
070192
070193
070194
070195
070196
070197
070198
070199
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140
070141
070142
070143
070144
070145
070146
070147
070148
070149
070150
070151
070152
070153
070154
070155
070156
070157
070158
070159
070160
070161
070162
070163
070164
070165
070166
070167
070168
070169
070170
070171
070172
070173
070174
070175
070176
070177
070178
070179
070180
070181
070182
070183
070184
070185
070186
070187
070188
070189
070190
070191
070192
070193
070194
070195
070196
070197
070198
070199
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140
070141
070142
070143
070144
070145
070146
070147
070148
070149
070150
070151
070152
070153
070154
070155
070156
070157
070158
070159
070160
070161
070162
070163
070164
070165
070166
070167
070168
070169
070170
070171
070172
070173
070174
070175
070176
070177
070178
070179
070180
070181
070182
070183
070184
070185
070186
070187
070188
070189
070190
070191
070192
070193
070194
070195
070196
070197
070198
070199
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140
070141
070142
070143
070144
070145
070146
070147
070148
070149
070150
070151
070152
070153
070154
070155
070156
070157
070158
070159
070160
070161
070162
070163
070164
070165
070166
070167
070168
070169
070170
070171
070172
070173
070174
070175
070176
070177
070178
070179
070180
070181
070182
070183
070184
070185
070186
070187
070188
070189
070190
070191
070192
070193
070194
070195
070196
070197
070198
070199
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140
070141
070142
070143
070144
070145
070146
070147
070148
070149
070150
070151
070152
070153
070154
070155
070156
070157
070158
070159
070160
070161
070162
070163
070164
070165
070166
070167
070168
070169
070170
070171
070172
070173
070174
070175
070176
070177
070178
070179
070180
070181
070182
070183
070184
070185
070186
070187
070188
070189
070190
070191
070192
070193
070194
070195
070196
070197
070198
070199
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140
070141
070142
070143
070144
070145
070146
070147
070148
070149
070150
070151
070152
070153
070154
070155
070156
070157
070158
070159
070160
070161
070162
070163
070164
070165
070166
070167
070168
070169
070170
070171
070172
070173
070174
070175
070176
070177
070178
070179
070180
070181
070182
070183
070184
070185
070186
070187
070188
070189
070190
070191
070192
070193
070194
070195
070196
070197
070198
070199
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140
070141
070142
070143
070144
070145
070146
070147
070148
070149
070150
070151
070152
070153
070154
070155
070156
070157
070158
070159
070160
070161
070162
070163
070164
070165
070166
070167
070168
070169
070170
070171
070172
070173
070174
070175
070176
070177
070178
070179
070180
070181
070182
070183
070184
070185
070186
070187
070188
070189
070190
070191
070192
070193
070194
070195
070196
070197
070198
070199
070100
070101
070102
070103
070104
070105
070106
070107
070108
070109
070110
070111
070112
070113
070114
070115
070116
070117
070118
070119
070120
070121
070122
070123
070124
070125
070126
070127
070128
070129
070130
070131
070132
070133
070134
070135
070136
070137
070138
070139
070140

(Zhang et al., 2023) and ScanRef (Chen et al., 2020). We also introduce a small portion of 4D scene understanding datasets related to dense caption from Chat4D (Zhou & Lee, 2025). Additionally, we integrate 3D generation datasets, such as CO3D (Reizenstein et al., 2021), Objaverse (Deitke et al., 2023), RealEstate10k (Zhou et al., 2018) and MVIImgNet (Yu et al., 2023b) to improve the generation quality, where the conditional text can also serve as the caption for scene understanding tasks. In total, these datasets contain 980.6K samples. We employ hybrid datasets to align fine-grained spatiotemporal information between visual and linguistic representations for both understanding and generation tasks. For the training strategy, we update the trainable parameters of spatiotemporal embedding, adaptive cross-attention fusion, higher layers of the LLM, and multi-task heads while keeping the remaining modules frozen.

Stage 3: 4D Task Instruction Fine-Tuning. This stage is to improve the generalization of our model for understanding and generation tasks in more complex 4D dynamic scenes. For the datasets, we introduce the typical 4D understanding dataset Chat4D (Zhou & Lee, 2025), and integrate 4D generation datasets including 4DNex-10M (Chen et al., 2025) and DyCheck (Gao et al., 2022). In total, these datasets contain 160.7K samples. We employ these 4D multimodal datasets to perform fine-tuning to adapt our model to the intricacies of complex 4D scenarios. For the training strategy, all trainable parameters are optimized through LoRA adapters. The vision encoder-decoder and geometry encoder remain frozen.

2 ABLATION STUDY AND DISCUSSION

2.1 APPLICABILITY TO DIFFERENT INPUT SETTINGS

Although we introduce the entire model using multi-view videos as an input example, we also discuss the applicability of our model to other input settings such as single-view video. We use single-view videos collected from the 4D vision datasets of Chat4D (Zhou & Lee, 2025) and the corresponding instruction-following texts as ground truths for evaluation of the understanding task. As shown in Table 2, the performance metrics of our model with single-view videos as input are comparable to our model with multi-view videos as input. This shows that our model is applicable to different input settings and is practical for the real world.

Table 2: Ablation study on the effect of various vision inputs on scene understanding performance.

Vision input	C↑	SAcc@0.5↑	TAcc↑
w/ Single-view video	93.6	57.9	54.5
w/ Multi-view video	93.8	58.2	54.6

2.2 IMPACT OF TASK PROMPT

In Table 3, we analyze the impact of the task prompt on model performance. The results show that introducing the task prompt significantly improves the model performance on both understanding and generation for 4D scenes. The main reason is that the task prompt not only distinguishes the features of understanding and generation within the visual representation, but also guides the subsequent attention mask to dynamically regulate the information flow of different tasks in our model. These mechanisms enhance the capability of our model for multi-task prediction.

Table 3: Ablation study on the effect of task prompt on model performance.

Task prompt	Chat4D			DyCheck		
	C↑	SAcc@0.5↑	TAcc↑	PSNR↑	FVD↓	CLIP-C↑
w/o Prompt	85.1	47.5	45.3	17.54	216.1	0.91
w/ Prompt	93.8	58.2	54.6	21.38	152.3	0.97

108
109 2.3 EFFECT OF TEXTUAL COORDINATE EMBEDDING110
111 In Table 4, we study the effect of textual coordinate embedding, *i.e.* pos/time encoding via special
112 token embedding (Li et al., 2025) on understanding performance. First, textual coordinates as
113 instructions improve the fine-grained spatiotemporal understanding of our model. Second, textual
114 coordinate embedding further improves the upper limit of 4D spatiotemporal understanding. This is
115 because textual coordinate embedding helps minimize the risk that large language models misinterpret
116 coordinate values.

117 Table 4: Discussion on the effect of textual coordinate embedding for scene understanding.

118
119

Text instruction		C↑	SAcc@0.5↓	TAcc↑
w/o Coordinate		84.2	39.6	19.3
w/ Coordinate	w/o Embedding	90.5	55.1	52.7
	w/ Embedding	93.8	58.2	54.6

124 2.4 LIMITATION AND POTENTIAL SOLUTION

125
126 Despite strong performance in most short-term scene understanding and generation, our Uni4D-LLM
127 struggles with long-term dynamics. Capturing such variations requires memory-based reasoning to
128 model cross-spatiotemporal interactions and causal relations. However, the current framework relies
129 on short-range attention, which primarily captures correlations within limited temporal windows. As a
130 result, the model lacks explicit mechanisms to retain or propagate motion information across extended
131 sequences and thus makes it difficult to reason about long-horizon dependencies. In the future, we
132 plan to integrate a world model (Ha & Schmidhuber, 2018) to enable long-term spatiotemporal
133 reasoning and extend scene understanding and generation to longer temporal horizons. This will
134 improve the practicality of our model.

136 3 QUANTITATIVE AND QUALITATIVE RESULTS

137 3.1 QUANTITATIVE COMPARISON ON 2D UNDERSTANDING BENCHMARK

138
139 In Table 5, we compare the basic visual understanding capability of our model with other 2D, 3D and
140 4D VLMs on several typical 2D understanding benchmarks, including VQAv2 (Goyal et al., 2017),
141 MMBench (Liu et al., 2024b), MME (Fu et al., 2024), MM-Vet (Yu et al., 2023a). The results show
142 that our method can maintain the same level of image understanding as other competing VLMs. This
143 also indirectly proves that our model can obtain the initial multimodal understanding capability from
144 the 2D datasets.

146 Table 5: Quantitative results of VLMs for scene understanding on 2D zero-shot benchmarks.

147
148

Methods		VQAv2	MMBench	MME	MM-Vet
2D	MobileVLM (Chu et al., 2023)	47.5	59.6	1289	–
	Qwen-VL (Bai et al., 2023)	63.8	38.2	–	–
	Qwen-VL-Chat (Bai et al., 2023)	61.5	60.6	1488	–
	LLaMA-VID (Li et al., 2024)	–	65.1	1521	–
	LLaVA-1.5 (Liu et al., 2024a)	58.2	65.2	1511	31.1
3D	LLaVA-3D (Zhu et al., 2024)	57.8	65.0	1502	30.9
4D	Uni4D-LLM (Ours)	58.8	64.5	1506	30.7

158 3.2 QUALITATIVE COMPARISON ON 3D UNDERSTANDING BENCHMARK

159
160 In Fig. 1, we provide more visual comparisons on 3D scene understanding and generation, where we
161 provide some typical scenes as examples. In 3D understanding, our Uni4D-LLM is more fine-grained
than the competing 3D VLMs in representing spatial characteristics of the scene. In 3D generation,

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181

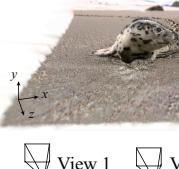
	
User What is the object at [2.16, 1.85, 0.68]?	User Please describe this scene.
PQ3D A bed located near the wall, horizontally aligned with pillows and a dark blanket on top.	PQ3D The scene contains a bed, sofa, chairs, and a desk with monitors, arranged along the walls of a compact room.
LLaVA-3D This is a bed placed at the back of the room, with a dark blanket and pillow, next to the wall and under a window.	LLaVA-3D This indoor scene includes a bed, central desk with monitors, two chairs, and a sofa opposite the bed.
Video-3D LLM The object is bed covered in dark bedding, placed against the wall and adjacent to a small nightstand and window.	Video-3D LLM This is a small bedroom scene with a bed on the right, a central desk with monitors, and a green sofa and chairs on the left.
Uni4D-LLM The object is a bed covered with dark blue sheets, stationary at the center of the room, aligned with the wall, and surrounded by a nightstand and a desk.	Uni4D-LLM This is a furnished room with a bed beneath the window, a desk with two monitors placed in front, chairs neatly tucked in under the desk, and a green sofa positioned facing the bed across the room.

(a) Global 3D Scene Understanding

(b) Local 3D Scene Understanding

182
183
184
185

Figure 1: Visual comparison of VLMs on 3D scene understanding.

	
User What happened in this scene?	
LLaVA-3D A seal is lying on a sandy beach, turning its head to the right.	
LLaVA-4D A spotted seal lying on the beach slowly lifted its head and turned from side to side, appearing to look around curiously across several moments. Faint curved trails behind its body suggest it had recently moved forward along the wet sand.	
Uni4D-LLM A spotted seal lies on the sandy beach with its body stretched forward. It slowly lifts its head and turns from side to side, appearing to look around curiously. Behind the seal, faint curved trails in the damp sand suggest it had recently moved forward.	

201
202
203
204

Figure 2: Visual comparison of VLMs on 4D scene understanding.

205
206
207
208
209
210
211
212
213
214
215

the results generated by our model are on par with those of 3D generation models. These results demonstrate the effectiveness of our model on both 3D scene understanding and generation.

3.3 QUALITATIVE COMPARISON ON 4D UNDERSTANDING AND GENERATION BENCHMARK

In Fig. 2 and 3, we also provide more visual comparisons on 4D scene understanding and generation. In 4D understanding, 3D VLMs struggle to capture temporal dynamics, while our model demonstrates strong spatiotemporal reasoning on par with recent 4D VLMs. In 4D generation, our Uni4D-LLM produces sharp and coherent results that rival those of advanced 4D diffusion models. These results demonstrate the superiority of our Uni4D-LLM in 4D understanding and generation, underscoring its potential as a unified multi-task framework for the physical world.

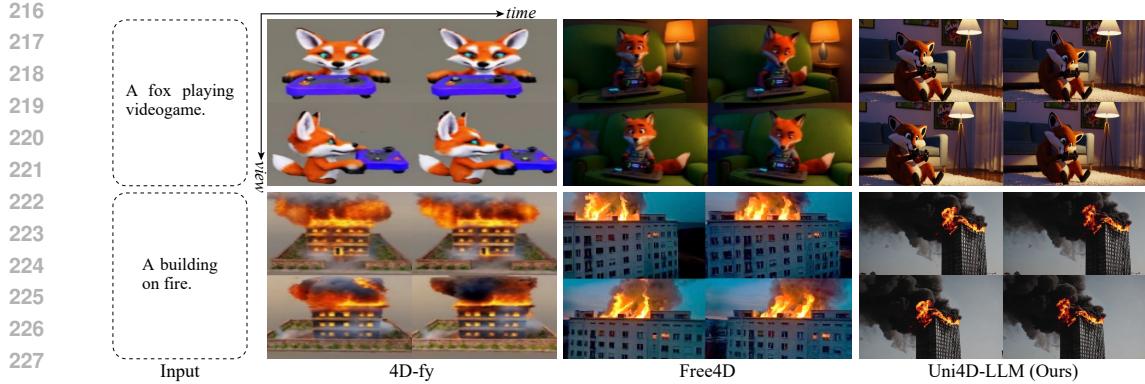


Figure 3: Visual comparison on 4D scene generation.

REFERENCES

Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 19129–19139, 2022.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 1(2):3, 2023.

Max Bain, Arsha Nagrani, GüL Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1728–1738, 2021.

Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. In *European conference on computer vision*, pp. 202–221. Springer, 2020.

Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13320–13331, 2024.

Zhaoxi Chen, Tianqi Liu, Long Zhuo, Jiawei Ren, Zeng Tao, He Zhu, Fangzhou Hong, Liang Pan, and Ziwei Liu. 4dnex: Feed-forward 4d generative modeling made easy. *arXiv preprint arXiv:2508.13154*, 2025.

Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X Chang. Scan2cap: Context-aware dense captioning in rgb-d scans. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 3193–3203, 2021.

Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 1(2):3, 2023.

Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 13142–13153, 2023.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 248–255. Ieee, 2009.

Chaoyou Fu, Yuhang Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024.

270 Hang Gao, Rui long Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Monocular
 271 dynamic view synthesis: A reality check. *Advances in Neural Information Processing Systems*, 35:
 272 33768–33780, 2022.

273

274 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa
 275 matter: Elevating the role of image understanding in visual question answering. In *Proceedings of
 276 the IEEE conference on computer vision and pattern recognition*, pp. 6904–6913, 2017.

277

278 David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Adv. Neural
 279 Inform. Process. Syst.*, 31, 2018.

280

281 De-An Huang, Shijia Liao, Subhashree Radhakrishnan, Hongxu Yin, Pavlo Molchanov, Zhiding Yu,
 282 and Jan Kautz. Lita: Language instructed temporal-localization assistant. In *European Conference
 283 on Computer Vision*, pp. 202–218. Springer, 2024.

284

285 Hongyu Li, Jinyu Chen, Ziyu Wei, Shaofei Huang, Tianrui Hui, Jialin Gao, Xiaoming Wei, and Si Liu.
 286 Llava-st: A multimodal large language model for fine-grained spatial-temporal understanding.
arXiv preprint arXiv:2501.08282, 2025.

287

288 Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language
 289 models. In *European Conference on Computer Vision*, pp. 323–340. Springer, 2024.

290

291 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
 292 tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 pp. 26296–26306, 2024a.

293

294 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi
 295 Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player?
 296 In *European conference on computer vision*, pp. 216–233. Springer, 2024b.

297

298 Ruipu Luo, Zi Wang Zhao, Min Yang, Junwei Dong, Da Li, Pengcheng Lu, Tao Wang, Linmei Hu,
 299 Minghui Qiu, and Zhongyu Wei. Valley: Video assistant with large language model enhanced
 300 ability. *arXiv preprint arXiv:2306.07207*, 2023.

301

302 Ruiyuan Lyu, Jingli Lin, Tai Wang, Xiaohan Mao, Yilun Chen, Runsen Xu, Haifeng Huang, Chenming
 303 Zhu, Dahua Lin, and Jiangmiao Pang. Mmscan: A multi-modal 3d scene dataset with hierarchical
 304 grounded language annotations. *Advances in Neural Information Processing Systems*, 37:50898–
 305 50924, 2024.

306

307 Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao Liang, Song-Chun Zhu, and Siyuan Huang.
 308 Sq3d: Situated question answering in 3d scenes. *arXiv preprint arXiv:2210.07474*, 2022.

309

310 Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham
 311 Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel
 312 grounding large multimodal model. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 13009–13018,
 313 2024.

314

315 Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David
 316 Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category
 317 reconstruction. In *Int. Conf. Comput. Vis.*, pp. 10901–10911, 2021.

318

319 Zhongwei Ren, Zhicheng Huang, Yunchao Wei, Yao Zhao, Dongmei Fu, Jiashi Feng, and Xiaojie
 320 Jin. Pixellm: Pixel reasoning with large multimodal model. In *Proceedings of the IEEE/CVF
 321 Conference on Computer Vision and Pattern Recognition*, pp. 26374–26383, 2024.

322

323 Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan
 324 Chen, Yaohui Wang, et al. Internvid: A large-scale video-text dataset for multimodal understanding
 325 and generation. *arXiv preprint arXiv:2307.06942*, 2023.

326

327 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,
 328 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv
 329 preprint arXiv:2308.02490*, 2023a.

324 Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan,
325 Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimgnet: A large-scale dataset of
326 multi-view images. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 9150–9161, 2023b.
327

328 Yiming Zhang, ZeMing Gong, and Angel X Chang. Multi3drefer: Grounding text description to
329 multiple 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer
330 Vision*, pp. 15225–15236, 2023.

331 Hanyu Zhou and Gim Hee Lee. Llava-4d: Embedding spatiotemporal prompt into lmms for 4d scene
332 understanding. *arXiv preprint arXiv:2505.12253*, 2025.

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

335 Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang, and Xihui Liu. Llava-3d: A simple
336 yet effective pathway to empowering lmms with 3d-awareness. *arXiv preprint arXiv:2409.18125*,
337 2024.

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