

# Coordinates are not lonely - Codebook Prior Helps Implicit Neural 3D Representations

## Supplementary Materials

Anonymous Author(s)

Affiliation

Address

email

### 1 Novel View Synthesis Videos

We visualize some of our experiments in a video, including sparse views and few views setting on the DTU [1], BlendedMVS [5] and H3DS [4] dataset (see CoCo-INR\_Sparse\_Views\_Setting.mp4 and CoCo-INR\_Few\_Views\_Setting.mp4). Specifically, two opposite views (i.e., front and behind) are selected as the start and end view, then an arc-shaped sequence of novel views with even adjacent view degree interval are then synthesized between the start and end views, constituting a total number of 60 views covering 180 degrees variation. Note that the DTU dataset is not selected from 360-degree surrounding views so it is not suitable for the above view generation method. For DTU dataset, we chose the leftmost and rightmost views to generate a arc-shaped view sequence between them and use masks to reduce the background ambiguity caused by the non-surround setting.

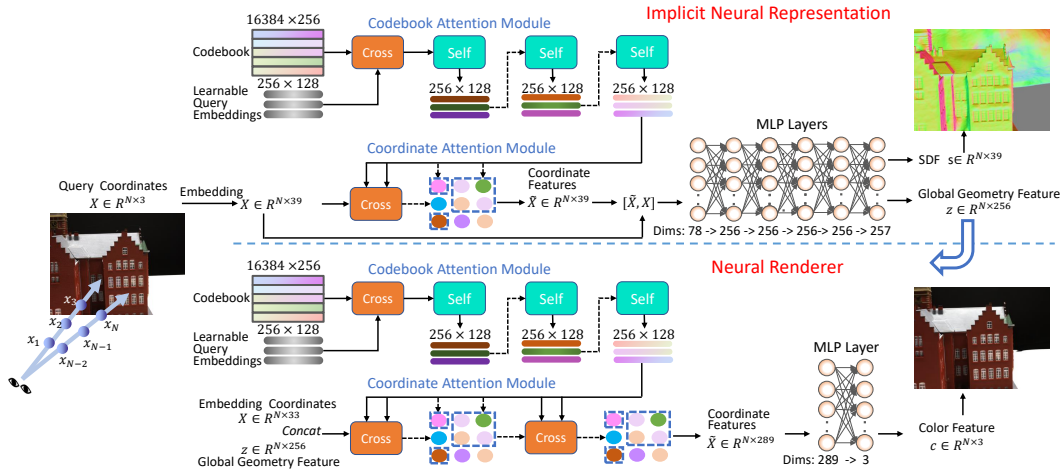


Figure 1: Pipeline of our method.

### 2 Details of Our CoCo-INR

We now present the pipeline and details of our proposed CoCo-INR. As shown in Fig. 1, it contains two modules: implicit neural representation and neural renderer, each of which consists of several codebook attention and coordinate attention modules and MLP layers. The structure of the codebook

attention and coordinate attention module has been described in our paper. Here is the workflow of our method.

### 3 Additional Experiments

In our paper, we have comprehensively evaluated the performance of our method and other state-of-the-art methods on DTU [1], BlendedMVS [5], and H3DS [4] datasets under different numbers of training views. In this section, we will present detailed data for qualitative comparisons, challenge fewer training views, and analyze limitations.

#### 3.1 Detailed Quantitative Results

Qualitative comparisons with NeRF [2], UNISURF [3], and VolSDF [6] have been shown in paper. Our CoCo-based method outperforms others in terms of PSNR, SSIM, and LPIPS whether sparse views or few views are available. This section further presents the evaluation results for each scan on DTU [1], BlendedMVS [5], and H3DS [4] datasets.

**For sparse views setting (16-32 in total)**, Table 1 shows the comparison of each scan on the DTU dataset, and Table 2 shows the comparison of each scan on the BlendedMVS dataset, and Table 3 shows the comparison of each scan on the H3DS dataset. It can be seen that our method outperforms other methods in most scans, which means that our method has good performance and robustness under sparse views.

Table 1: The results of different methods on the DTU dataset for each scan with sparse views (16-32) and without object masks.  $\uparrow$  means the higher, the better.

Scan	PSNR $\uparrow$				SSIM $\uparrow$				LPIPS $\downarrow$			
	NeRF	UNISURF	VolSDF	Ours	NeRF	UNISURF	VolSDF	Ours	NeRF	UNISURF	VolSDF	Ours
24	10.720	20.060	23.100	<b>23.896</b>	0.478	0.704	0.676	<b>0.791</b>	0.377	0.277	0.246	<b>0.227</b>
37	19.789	20.065	<b>22.576</b>	22.022	0.628	0.650	0.607	<b>0.703</b>	<b>0.227</b>	0.277	0.254	0.230
40	9.424	12.185	21.705	<b>22.803</b>	0.364	0.512	<b>0.593</b>	0.569	0.338	0.313	0.234	<b>0.231</b>
55	20.323	22.264	<b>24.364</b>	23.804	0.734	0.783	<b>0.835</b>	0.821	0.329	0.379	0.333	<b>0.318</b>
63	24.492	24.539	25.378	<b>26.568</b>	<b>0.826</b>	0.811	0.699	0.705	<b>0.184</b>	0.229	0.233	0.200
65	18.161	23.859	<b>27.904</b>	26.121	0.780	0.837	<b>0.870</b>	0.850	<b>0.262</b>	0.310	0.271	0.282
69	23.478	24.427	25.434	<b>25.935</b>	0.867	0.882	0.910	<b>0.915</b>	0.302	0.339	<b>0.274</b>	0.283
83	25.733	25.808	26.165	<b>27.975</b>	0.922	0.905	0.925	<b>0.944</b>	<b>0.270</b>	0.370	0.360	0.391
97	22.119	21.833	23.765	<b>24.029</b>	0.854	0.845	0.887	<b>0.896</b>	<b>0.297</b>	0.392	0.308	0.339
105	24.746	24.240	27.257	<b>27.577</b>	0.895	0.888	<b>0.926</b>	0.924	<b>0.268</b>	0.372	0.363	0.300
106	27.629	26.919	30.394	<b>30.813</b>	0.904	0.907	0.932	<b>0.940</b>	0.381	0.389	0.345	<b>0.338</b>
110	25.353	26.046	28.084	<b>30.190</b>	0.916	0.910	0.924	<b>0.939</b>	<b>0.369</b>	0.426	0.392	<b>0.369</b>
114	24.297	24.394	28.039	<b>28.707</b>	0.863	0.875	<b>0.907</b>	0.902	0.348	0.387	0.340	<b>0.338</b>
118	29.082	28.484	<b>32.218</b>	31.358	0.919	0.921	<b>0.948</b>	0.943	<b>0.312</b>	0.377	0.341	0.329
122	27.091	28.125	<b>32.765</b>	29.270	0.901	0.921	<b>0.952</b>	0.942	0.334	0.387	0.350	<b>0.294</b>
<b>Mean</b>	22.162	23.549	26.609	<b>26.738</b>	0.790	0.823	0.839	<b>0.852</b>	0.306	0.348	0.309	<b>0.298</b>

**For few views setting (5-8 in total)**, Table 4 shows the comparison of each scan on the DTU dataset, and Table 5 shows the comparison of each scan on the BlendedMVS dataset, and Table 6 shows the comparison of each scan on the H3DS dataset. VolSDF [6] achieves better performance than other state-of-the-art methods in the sparse views setting, so we only compare with VolSDF in this more difficult setting. It can be seen that in each scan, using fewer training views will inevitably lead to performance degradation compared with sparse views, but our method still maintains better robustness and performance for the majority of scans.

Table 2: The results of different methods on the BlendedMVS dataset for each scan with sparse views (16-32) and without object masks.  $\uparrow$  means the higher, the better.

Scan	PSNR $\uparrow$				SSIM $\uparrow$				LPIPS $\downarrow$			
	NeRF	UNISURF	VolSDF	Ours	NeRF	UNISURF	VolSDF	Ours	NeRF	UNISURF	VolSDF	Ours
1	19.179	15.586	20.004	<b>20.038</b>	0.675	0.626	<b>0.725</b>	<b>0.725</b>	0.256	0.295	0.215	<b>0.211</b>
2	14.220	14.848	19.297	<b>20.951</b>	0.630	0.651	0.757	<b>0.810</b>	0.320	0.351	0.234	<b>0.203</b>
3	16.684	14.717	16.808	<b>17.263</b>	0.647	0.600	0.671	<b>0.682</b>	0.293	0.336	0.268	<b>0.265</b>
4	<b>23.684</b>	15.200	23.214	23.007	0.838	0.755	0.847	<b>0.854</b>	0.170	0.258	0.154	<b>0.148</b>
5	17.160	14.361	16.744	<b>17.528</b>	0.745	0.700	0.758	<b>0.784</b>	0.244	0.304	0.215	<b>0.207</b>
6	15.575	13.739	19.775	<b>21.741</b>	0.670	0.630	0.788	<b>0.819</b>	0.282	0.304	0.175	<b>0.148</b>
7	8.722	12.086	17.972	<b>17.997</b>	0.528	0.610	<b>0.721</b>	<b>0.721</b>	0.349	0.333	0.237	<b>0.223</b>
8	17.851	17.187	18.924	<b>19.827</b>	0.743	0.737	0.821	<b>0.837</b>	0.183	0.171	0.140	<b>0.128</b>
9	15.634	15.019	17.747	<b>17.997</b>	0.535	0.540	0.639	<b>0.640</b>	0.356	0.346	0.284	<b>0.279</b>
<b>Mean</b>	16.523	14.749	18.942	<b>19.594</b>	0.667	0.649	0.747	<b>0.764</b>	0.272	0.299	0.213	<b>0.201</b>

Table 3: The results of different methods on the H3DS dataset for each scan with sparse views (16-32) and without object masks.  $\uparrow$  means the higher, the better.

Scan	PSNR $\uparrow$				SSIM $\uparrow$				LPIPS $\downarrow$			
	NeRF	UNISURF	VolSDF	Ours	NeRF	UNISURF	VolSDF	Ours	NeRF	UNISURF	VolSDF	Ours
*a287	21.826	14.314	25.283	<b>26.070</b>	0.896	0.801	<b>0.941</b>	<b>0.941</b>	0.121	0.198	0.082	<b>0.076</b>
*42a8	19.783	20.170	22.437	<b>23.785</b>	0.825	0.832	0.863	<b>0.887</b>	0.179	0.202	0.138	<b>0.121</b>
*1d54	23.779	18.028	25.420	<b>25.933</b>	0.878	0.810	0.893	<b>0.904</b>	0.130	0.184	0.114	<b>0.107</b>
*0854	20.302	20.236	22.819	<b>24.274</b>	0.840	0.835	0.865	<b>0.875</b>	0.151	0.151	0.127	<b>0.124</b>
*0c89	22.241	21.371	23.658	<b>25.116</b>	0.911	0.898	0.932	<b>0.942</b>	0.117	0.122	0.087	<b>0.084</b>
*0226	22.207	12.104	25.582	<b>25.742</b>	0.861	0.762	<b>0.907</b>	<b>0.907</b>	0.130	0.246	<b>0.088</b>	<b>0.088</b>
*4baa	21.458	15.378	24.949	<b>26.755</b>	0.884	0.820	0.909	<b>0.931</b>	0.139	0.187	0.106	<b>0.094</b>
*512f	21.009	20.250	24.217	<b>25.141</b>	0.830	0.826	0.878	<b>0.895</b>	0.183	0.223	0.141	<b>0.127</b>
*3924	21.108	18.979	23.774	<b>26.090</b>	0.865	0.832	0.897	<b>0.918</b>	0.133	0.201	0.112	<b>0.091</b>
*e5bc	23.274	12.605	24.051	<b>24.395</b>	0.908	0.768	0.918	<b>0.926</b>	0.099	0.210	0.084	<b>0.079</b>
*ee0b	21.128	12.636	22.629	<b>24.136</b>	0.886	0.763	0.911	<b>0.921</b>	0.131	0.217	0.100	<b>0.091</b>
*e187	24.039	20.294	25.256	<b>26.738</b>	0.893	0.846	0.913	<b>0.923</b>	0.114	0.146	0.092	<b>0.084</b>
*d436	19.579	18.488	21.114	<b>21.356</b>	0.817	0.799	0.847	<b>0.851</b>	0.168	0.237	0.160	<b>0.136</b>
*b9c0	21.514	20.944	23.069	<b>27.221</b>	0.854	0.855	0.877	<b>0.916</b>	0.153	0.182	0.129	<b>0.101</b>
*9c4e	21.905	17.540	25.714	<b>27.557</b>	0.911	0.881	0.943	<b>0.950</b>	0.110	0.144	0.080	<b>0.073</b>
*7be3	20.288	14.227	20.190	<b>22.983</b>	0.835	0.787	0.843	<b>0.866</b>	0.138	0.198	0.138	<b>0.106</b>
*fd85	22.270	9.058	24.448	<b>25.250</b>	0.873	0.658	0.900	<b>0.910</b>	0.140	0.296	0.114	<b>0.102</b>
*c1e6	21.110	15.816	25.507	<b>26.204</b>	0.889	0.832	0.932	<b>0.938</b>	0.132	0.193	0.090	<b>0.082</b>
*4b87	21.769	22.298	24.960	<b>25.244</b>	0.897	0.901	0.931	<b>0.934</b>	0.126	0.116	0.094	<b>0.092</b>
*2a8f	11.296	15.006	24.133	<b>25.439</b>	0.697	0.821	0.928	<b>0.931</b>	0.285	0.180	0.094	<b>0.093</b>
*244e	22.966	14.668	23.207	<b>25.295</b>	0.833	0.758	0.853	<b>0.875</b>	0.155	0.226	0.123	<b>0.106</b>
*deb0	22.809	20.577	25.772	<b>28.006</b>	0.898	0.862	0.931	<b>0.942</b>	0.139	0.141	0.114	<b>0.097</b>
*2091	19.607	18.199	22.038	<b>22.685</b>	0.839	0.838	0.864	<b>0.871</b>	0.159	0.176	0.134	<b>0.122</b>
<b>Mean</b>	21.185	17.095	23.922	<b>25.279</b>	0.861	0.816	0.898	<b>0.911</b>	0.144	0.190	0.110	<b>0.098</b>

Table 4: The results of different methods on the DTU dataset for each scan with few views (5-8) and without object masks.  $\uparrow$  means the higher, the better.

Scan	PSNR $\uparrow$		SSIM $\uparrow$		LPIPS $\downarrow$	
	VolSDF	Ours	VolSDF	Ours	VolSDF	Ours
24	17.476	<b>17.867</b>	0.544	<b>0.617</b>	<b>0.274</b>	0.285
37	10.783	<b>15.372</b>	0.554	<b>0.615</b>	0.310	<b>0.271</b>
40	<b>14.489</b>	13.627	<b>0.574</b>	0.496	<b>0.286</b>	0.294
55	14.840	<b>18.014</b>	0.641	<b>0.720</b>	0.393	<b>0.382</b>
63	7.661	<b>18.391</b>	0.482	<b>0.572</b>	0.361	<b>0.226</b>
65	19.173	<b>19.265</b>	0.769	<b>0.771</b>	0.319	<b>0.307</b>
69	<b>16.769</b>	16.268	<b>0.720</b>	0.705	0.387	<b>0.379</b>
83	<b>18.746</b>	17.119	<b>0.866</b>	0.864	<b>0.375</b>	0.391
97	18.321	<b>18.821</b>	<b>0.827</b>	<b>0.827</b>	0.382	<b>0.378</b>
105	<b>18.914</b>	18.165	<b>0.873</b>	0.861	0.360	<b>0.344</b>
106	<b>23.132</b>	23.103	<b>0.845</b>	0.836	<b>0.377</b>	0.407
110	21.763	<b>25.954</b>	0.843	<b>0.900</b>	0.437	<b>0.436</b>
114	<b>24.510</b>	22.408	<b>0.861</b>	0.834	0.372	<b>0.352</b>
118	<b>28.324</b>	26.602	<b>0.913</b>	0.905	0.397	<b>0.368</b>
122	<b>25.488</b>	23.385	<b>0.901</b>	0.836	0.400	<b>0.372</b>
<b>Mean</b>	18.693	<b>19.624</b>	0.748	<b>0.757</b>	0.362	<b>0.346</b>

Table 5: The results of different methods on the BlendedMVS dataset for each scan with few views (5-8) and without object masks.  $\uparrow$  means the higher, the better.

Scan	PSNR $\uparrow$		SSIM $\uparrow$		LPIPS $\downarrow$	
	VolSDF	Ours	VolSDF	Ours	VolSDF	Ours
1	13.313	<b>16.597</b>	0.562	<b>0.630</b>	0.369	<b>0.265</b>
2	15.037	<b>15.277</b>	0.635	<b>0.637</b>	0.311	<b>0.303</b>
3	<b>12.312</b>	10.965	<b>0.556</b>	0.490	<b>0.358</b>	0.398
4	14.132	<b>14.614</b>	0.722	<b>0.735</b>	0.278	<b>0.254</b>
5	13.781	<b>14.358</b>	0.675	<b>0.699</b>	0.299	<b>0.268</b>
6	11.662	<b>12.114</b>	<b>0.579</b>	0.575	0.375	<b>0.366</b>
7	14.708	<b>15.955</b>	0.647	<b>0.659</b>	0.281	<b>0.269</b>
8	<b>17.438</b>	16.384	<b>0.760</b>	0.759	<b>0.171</b>	0.178
9	12.876	<b>13.491</b>	0.478	<b>0.487</b>	0.415	<b>0.386</b>
<b>Mean</b>	13.918	<b>14.417</b>	0.624	<b>0.630</b>	0.317	<b>0.299</b>

Table 6: The results of different methods on the H3DS dataset for each scan with few views (5-8) and without object masks.  $\uparrow$  means the higher, the better.

Scan	PSNR $\uparrow$		SSIM $\uparrow$		LPIPS $\downarrow$	
	VolSDF	Ours	VolSDF	Ours	VolSDF	Ours
*a287	18.014	<b>18.804</b>	0.864	<b>0.878</b>	0.152	<b>0.136</b>
*42a8	18.207	<b>18.538</b>	<b>0.811</b>	<b>0.811</b>	0.202	<b>0.195</b>
*1d54	<b>22.271</b>	16.674	<b>0.872</b>	0.761	<b>0.139</b>	0.202
*0854	18.837	<b>19.950</b>	0.809	<b>0.827</b>	0.187	<b>0.171</b>
*0c89	<b>20.037</b>	14.813	<b>0.892</b>	0.779	<b>0.129</b>	0.226
*0226	20.999	<b>21.706</b>	0.863	<b>0.869</b>	<b>0.132</b>	0.134
*4baa	18.800	<b>20.792</b>	0.847	<b>0.875</b>	0.174	<b>0.148</b>
*512f	19.587	<b>19.780</b>	0.829	<b>0.832</b>	<b>0.190</b>	0.200
*3924	<b>19.950</b>	19.362	<b>0.857</b>	0.851	<b>0.155</b>	0.157
*e5bc	<b>22.159</b>	20.907	<b>0.913</b>	0.900	<b>0.113</b>	0.120
*ee0b	15.451	<b>16.959</b>	0.817	<b>0.841</b>	0.209	<b>0.187</b>
*e187	20.125	<b>21.774</b>	0.858	<b>0.891</b>	0.154	<b>0.128</b>
*d436	16.528	<b>17.806</b>	0.780	<b>0.794</b>	0.211	<b>0.193</b>
*b9c0	14.505	<b>19.288</b>	0.783	<b>0.843</b>	0.230	<b>0.188</b>
*9c4e	19.285	<b>20.137</b>	<b>0.894</b>	<b>0.894</b>	<b>0.121</b>	0.123
*7be3	19.439	<b>20.038</b>	0.825	<b>0.864</b>	<b>0.151</b>	<b>0.151</b>
*fd85	<b>19.874</b>	19.391	<b>0.858</b>	0.850	<b>0.168</b>	<b>0.168</b>
*c1e6	20.311	<b>20.470</b>	0.891	<b>0.894</b>	0.135	<b>0.122</b>
*4b87	<b>20.192</b>	19.527	<b>0.890</b>	0.880	0.136	<b>0.135</b>
*2a8f	<b>17.871</b>	16.654	<b>0.859</b>	0.831	<b>0.173</b>	0.196
*244e	18.218	<b>18.508</b>	0.797	<b>0.825</b>	<b>0.179</b>	0.182
*deb0	<b>20.181</b>	19.616	<b>0.865</b>	<b>0.865</b>	<b>0.152</b>	0.171
*2091	14.767	<b>17.450</b>	0.789	<b>0.816</b>	0.220	<b>0.170</b>
<b>Mean</b>	18.939	<b>19.085</b>	0.846	<b>0.847</b>	0.166	<b>0.165</b>

### 3.2 Setting with Fewer Views

In this section we attempt a challenge: only use three training views to train our CoCo-based network. Specifically, we select three representative views for each scan as training views (usually left, right, and top views), and the remaining views are used as test views to evaluate the performance of our method under extremely limited views. In this challenge, three training views mean that it is difficult for the network to learn background information, so we use masks to make the network more focus on the foreground. We conduct experiments on the DTU dataset and compare with mask-based VolSDF [6].

As shown in Table 7, our method has significantly improved over VolSDF with a higher mean PSNR, SSIM, and LPIPS respectively. Especially for the LPIPS, VolSDF is nearly 25% higher than our method, which means that the new view images generated by our method perform better at the semantic level. We visualize part of the experimental results in Fig. 2, and it can be seen that our CoCo-INR can still render satisfactory RGB images and surface features under the extremely limited 3 training views.

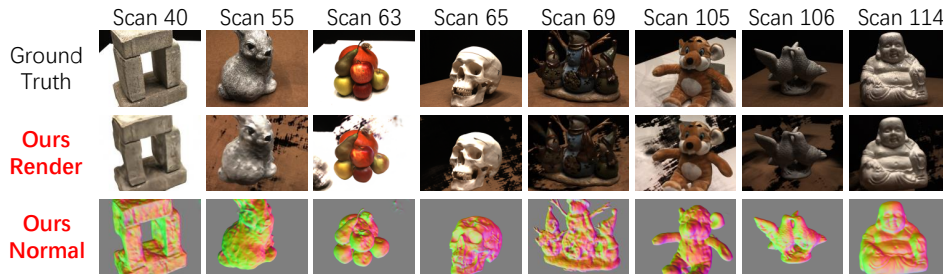


Figure 2: Qualitative visualization results (zoom-in for the best of views) on the DTU dataset with 3 extremely limited views.

Table 7: The results of different methods on the DTU dataset for each scene with 3 training views and objects masks.  $\uparrow$  means the higher, the better.

Scan	PSNR $\uparrow$		SSIM $\uparrow$		LPIPS $\downarrow$	
	VolSDF	Ours	VolSDF	Ours	VolSDF	Ours
24	<b>12.118</b>	11.266	<b>0.518</b>	0.443	0.352	<b>0.332</b>
37	9.840	<b>10.413</b>	<b>0.464</b>	0.387	0.337	<b>0.334</b>
40	<b>10.743</b>	9.963	<b>0.443</b>	0.372	<b>0.309</b>	0.317
55	14.709	<b>14.743</b>	<b>0.630</b>	0.619	<b>0.401</b>	0.414
63	<b>11.516</b>	11.179	<b>0.428</b>	0.374	0.294	<b>0.276</b>
65	15.411	<b>15.762</b>	0.674	<b>0.689</b>	0.337	<b>0.324</b>
69	<b>17.798</b>	17.472	<b>0.771</b>	0.761	<b>0.364</b>	0.375
83	10.316	<b>10.681</b>	0.675	<b>0.715</b>	0.430	<b>0.405</b>
97	11.936	<b>12.302</b>	0.707	<b>0.714</b>	0.419	<b>0.412</b>
105	<b>11.296</b>	11.261	0.716	<b>0.727</b>	0.410	<b>0.405</b>
106	17.311	<b>19.219</b>	0.513	<b>0.746</b>	0.713	<b>0.427</b>
110	<b>16.703</b>	16.367	<b>0.497</b>	<b>0.497</b>	0.769	<b>0.730</b>
114	16.763	<b>19.308</b>	0.524	<b>0.754</b>	0.672	<b>0.400</b>
118	15.928	<b>20.354</b>	0.281	<b>0.770</b>	0.780	<b>0.396</b>
122	16.549	<b>21.122</b>	0.445	<b>0.799</b>	0.789	<b>0.407</b>
<b>Mean</b>	13.929	<b>14.761</b>	0.552	<b>0.624</b>	0.491	<b>0.396</b>

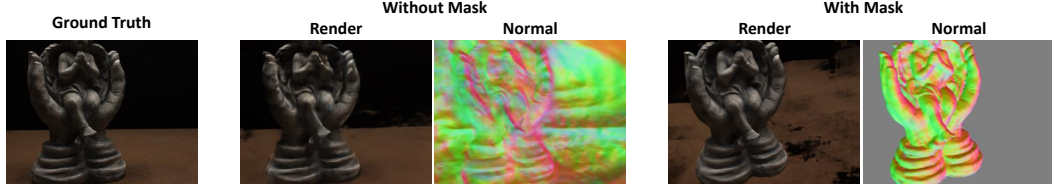


Figure 3: A failure case of our method with extremely limited views.

### 53 3.3 Limitations

54 For settings where the number of training views is extremely limited, it is possible that the limited  
55 views can not cover the entire background. So it is difficult to distinguish objects and background  
56 when they have similar color, as shown in Fig. 3. To improve the robustness of our method, for the  
57 setting of three training views, we have to use extra masks to reduce the background interference. In  
58 fact, we can fuse unsupervised methods and implicit neural representations to learn masks to reduce  
59 the manual annotation. Therefore, we will study self-supervised multi-task assisted CoCo-based  
60 methods to further improve the effectiveness under restricted conditions. Moreover, if we replace the  
61 per-scene learnable query embeddings in codebook attention modules with the image/pose-dependent  
62 feature tokens, our framework could be extended to a generalizable network across scenes rather than  
63 the current per-scene optimization.

## 64 References

- 65 [1] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view  
66 stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
67 pages 406–413, 2014.
- 68 [2] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng.  
69 Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer*  
70 *vision*, pages 405–421. Springer, 2020.
- 71 [3] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and  
72 radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on*  
73 *Computer Vision*, pages 5589–5599, 2021.
- 74 [4] Eduard Ramon, Gil Triginer, Janna Escur, Albert Pumarola, Jaime Garcia, Xavier Giro-i Nieto, and Francesc  
75 Moreno-Noguer. H3d-net: Few-shot high-fidelity 3d head reconstruction. In *Proceedings of the IEEE/CVF*  
76 *International Conference on Computer Vision*, pages 5620–5629, 2021.
- 77 [5] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan.  
78 Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *Proceedings of the*  
79 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1790–1799, 2020.
- 80 [6] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces.  
81 *Advances in Neural Information Processing Systems*, 34, 2021.