

DistillNeRF: Data-Efficient Initialization of Neural Radiance Fields using Knowledge Distillation

Anonymous ACCV 2022 submission

Paper ID ***

Abstract. Neural Radiance Fields (NeRF) learn a high-quality continuous 3D implicit representation of a scene given multiple views. While the approach has gained popularity in Novel View Synthesis (NVS), its vanilla implementation is not suited for real-time applications. This paper presents DistillNeRF, a data-efficient method for initializing and breeding smaller models using Knowledge Distillation (KD). Smaller models naturally benefit from lower inference times but decrease in perceptual quality. DistillNeRF optimizes training time and builds a priority grid from a teacher network as a data-efficient proxy for sampling better training examples reducing the quality loss while using half the memory.

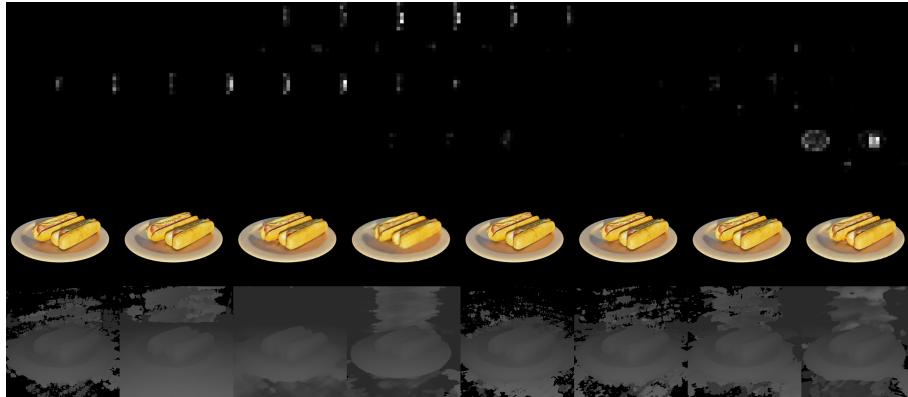


Fig. 1: DistillNeRF initializes smaller capacity NeRF models using KD from a teacher model. The figure is a depiction of the hot dog scene from the Blender Synthetic Dataset. The top row shows the occupancy voxel grid sampled from the pretrained NeRF. The middle row are renders from NeRF variations, and the bottom are their corresponding depth maps. The models used are NeRF, TinyNeRF, MicroNeRF, NanoNeRF, DistillNeRF, DistillTinyNeRF, DistillMicroNeRF, and DistillNanoNeRF respectively from left to right. The Distillation initialization reduces the perceptual quality loss observed in smaller models.

045

1 Introduction

 045

046

047 **Motivation.** Novel View Synthesis (NVS) is a challenging problem for Com- 047
 048 puter Science as it requires an understanding of the 3D scene structure, image 048
 049 reasoning, and a method to render its internal representation into target view- 049
 050 points. Current advances in Deep Learning allowed neural approaches to gain 050
 051 popularity and enabled applications that were impossible to achieve without a 051
 052 significant loss in quality. Neural NVS is becoming crucial to graphics pipeline 052
 053 based on real-world imagery such as image and video viewpoint editing [1, 2], 053
 054 and photo-realistic facial puppetry in Virtual Reality (VR) chat applications [3, 054
 055 4]. To this end, Neural Radiance Fields (NeRF) [5], learns a compact continuous 055
 056 3D implicit representation of the scene and can generate high-quality novel views 056
 057 while dealing with thin, transparent, semi-transparent, and reflective surfaces. 057

058

059 **Problem.** NeRF optimizes a large capacity Multi Layer Perceptron (MLP) to 059
 060 regress a volume density and view-dependent RGB radiance from a 5D vector 060
 061 (x, y, z, θ, ϕ). In its vanilla implementation, NeRF uses two neural networks to 061
 062 perform a hierarchical sampling of the scene, a coarse for approximating the 062
 063 geometry and a fine to refine its prediction. 063

064

064 While this formulation of neural NVS allows unprecedented render quality, 064
 065 it suffers from long training and inference time, thus failing at providing real- 065
 066 time use-cases. These factors can be explained by the number of network queries 066
 067 needed to render a frame and by NeRF’s two large capacity internal neural 067
 068 networks. Reducing both factors often results in lower quality renders and cannot 068
 069 be done without additional countermeasures. While the current literature has 069
 070 focused on reducing the number of queries [6, 7] and the use of render acceleration 070
 071 structures [8–13], this paper explores how to minimize NeRF’s network capacity 071
 072 while limiting the quality degradation. 072

073

073 **Solution.** We introduce DistillNeRF, a data-efficient method to initialize and 073
 074 breed lower capacity NeRF models using Knowledge Distillation (KD). Distill- 074
 075 NeRF uses only one network for both coarse and fine steps, naturally reducing 075
 076 the number of network queries and diving memory usage by half. Our method 076
 077 relies on extracting a priority occupancy voxel grid from a pretrained model, 077
 078 the teacher, to generate better data samples for training a smaller model, the 078
 079 student, and reduce the quality degradation due to the lower weight capacity. 079

080

081 **Findings.** We evaluate DistillNeRF on the Synthetic Blender NVS benchmark 081
 082 proposed in the original work [5] by distilling three flavors of NeRF students 082
 083 we named TinyNeRF, MicroNeRF, and NanoNeRF in decreasing order of depth 083
 084 and wideness. Our approaches reduces the quality degradation due to the lower 084
 085 capacity network on Mean Squared Error (MSE), Peak Signal to Noise Ratio 085
 086 (PSNR), and Structural Similarity (SSIM), using half the memory. 086

087

088

089

087

088

089

090

Contributions.

- A data-efficient method, DistillNeRF, to initialize and breed lower capacity NeRF models using KD. We observe reduced memory usage and lower quality degradation compared to previous initialization methods. 092
 - A study comparing three low capacity flavours of NeRF architectures we named TinyNeRF, MicroNeRF, and NanoNeRF against the vanilla training procedure, a meta-learning initialization scheme called Reptile on a Blender Synthetic Dataset, and an ablation study to justify our design choices. 095
 - An open-source C++ renderer using LibTorch, Cuda kernels, and OpenGL textures to limit memory transfers between the RAM and the VRAM. 099

103

Implication. Prior work focused on the acceleration of the rendering pipeline via the use of proxies such as acceleration structures, reformulation of the ray sampling strategies, and the radiance representation.

105

DistillNeRF is orthogonal to these contributions and focuses on reducing the size of the NeRF inner MLP. Our method can be applied in conjunction with sample reduction methods and acceleration structures. The reduction of the neural network size naturally accelerates the inference and rendering.

109

Such improvements would result in the use of neural NVS for real-time applications on consumer-grade hardware, the emergence of new artistic 3D workflows, and faster creative iterations.

112

Reproduction. For replicability and transparency, we published our implementation and experimentation at <https://github.com/after/publication>.

116

2 Related Work

110

Novel View Synthesis. Computer Graphics (CG) methods are defined as the set of formulations used to represent a scene geometry and rendering techniques used to generate their image representation. As such, Novel View Synthesis (NVS) methods are categorized by these criteria.

123

Mesh [14–17] and Point-Cloud [18–21] methods relies on preprocessing steps such as Structure from Motion (SfM) or RGB-D sensors to reconstruct the scene geometry from a set of multiview images. Multi-Plane Images (MPI) approaches [22–25] represents the scene as a stack of image layers allowing high quality renders in low range viewpoint displacements. Voxel-based pipelines [26–31] make use of a voxel occupancy grid as a low resolution scene representation. The use of voxels for NVS is often limited by the lack of structural definition due to a lack of resolution.

131

These four methods are often referred to as classical frameworks and present limitations that implicit neural representation can alleviate. Implicit differentiable rendering methods [32, 5, 33–36] performs optimization in the image space by formulating the rendering pipeline as a differentiable rasterization process.

135 Such approaches require accurate geometry masks and cannot handle thin, trans- 135
 136 parent, nor semi-transparent objects. In contrast, the use of differentiable volu- 136
 137 metric rendering in conjunction with Signed Distance Fields (SDF) [37, 38, 34] 137
 138 or Neural Radiance Fields (NeRF) [5] to represent the scene geometry does not 138
 139 require the use of any mask and can handle such complex objects. The use of 139
 140 Alpha-Blending enables those methods to recover thin, transparent, and semi- 140
 141 transparent materials. 141

142 Neural implicit representations and volumetric rendering are at the core of 142
 143 the current state of the art in geometric Computer Vision (CV) tasks such as 143
 144 3D shape reconstruction [39–41], scene relighting [42–45], object, human, and 144
 145 camera pose estimation [46–48], and more. 145

146
 147 **Training Stability.** Recent contributions and application of NeRF have ex- 147
 148 posed the limitations of its vanilla implementation. 148

149 The NeRF internal neural network is sensible to initialization and may suffer 149
 150 from instabilities during training on certain scenes resulting in exploding or 150
 151 vanishing values. In their work [49], Tancik et al. proposed the use of a meta- 151
 152 learning approach called Reptile [50] to learn better initializations. 152

153 In Mip-NeRF [51], Barron et al. not only present a ray casting formulation 153
 154 to reduce aliasing when the camera position differs from the dataset but also 154
 155 propose a set of stabilization tricks. They replace the softplus activation with 155
 156 a shifted one, the sigmoid activation a widened one to avoid common failure 156
 157 modes and yield smoother optimizations. 157

158
 159 **Fast Rendering.** Vanilla NeRF is not suited for realtime application (60 158
 160 Frames per Second (FPS)). In follow-up works, acceleration structures and op- 159
 161 timization tricks are used to reduce inference and accelerate rendering. 160

162 DeRF [8] splits the scene volume into sections for which a specific NeRF sub- 162
 163 network is trained on. Similarly, KiloNeRF [10] pushes the concept further by 163
 164 training thousands of tiny Multi Layer Perceptron (MLP) on bounded regions 164
 of the scene. 164

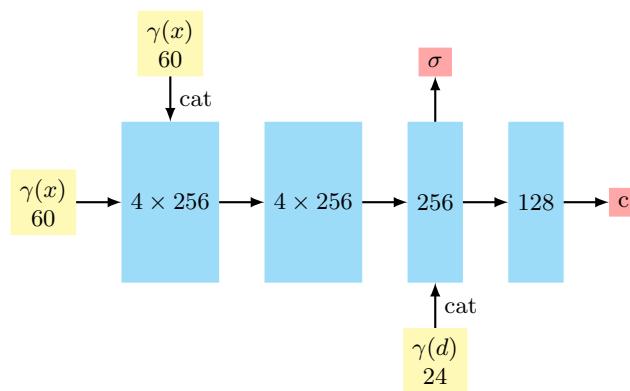
165 Neural sparse voxel fields [9] use a sparse voxel grid fitted to the scene geo- 165
 166 metry to cache and compute ray voxel intersections, thus avoiding sampling empty 166
 167 volumes. Hadman et al. [11] propose a similar approach by baking the NeRF 167
 168 model into an acceleration structure in combination with the use of differend 168
 169 rendering to accelerate the rendering performances. Yu et al. [12] build a plenoptic 169
 170 octree fitted onto the scene content and use spherical harmonics decompositon 170
 171 to encode the scene appearance. 171

172 DONeRF [6] trains a depth oracle in an end-to-end framework with the NeRF 172
 173 model to predict the surface location when raycasting. The number of samples 173
 174 per ray is then chosen accordingly to this prediction to reduce the amount of 174
 175 evaluation associated with empty space. 175

176
 177 **Extended Applications.** Since publication, NeRF has gained attraction in the 177
 178 field of 3D scene reconstruction and NVS, and has been adapted to a multitude 178
 179 of use-cases. 179

180 NeRF++ [52] extends its applications to 360 captures of large-scale un- 180
181 bounded scenes, NeRF-W [53] to in the wild monument photographies, and 181
182 NeRF in the Dark to High Dynamic Range (HDR) using noisy raw images. Oth- 182
183 ers have applied NeRF to text guided generation [39], free view realistic facial 183
184 animation [54] and controllable human appearance and pose articulations [47]. 184

185 Each of these contributions inherently suffers from NeRF vanilla formulation. 185
186 Faster training, inference, and stability would extend their usage to real-world 186
187 and real-time applications. 187



202 Fig. 2: Schematic of the NeRF internal MLP architecture. The first section of the 202
203 network is in charge of learning valuable feature representations of the scene con- 203
204 tent given an positional encoded position vector $\gamma(x)$ with residual connections. 204
205 These features are then used to regress a density value σ and a RGB radiance c 205
206 given a positional encoded view direction $\gamma(d)$. 206

210 3 Background 210

212 **Geometry Representation.** Neural Radiance Fields (NeRF) encodes a scene 212
213 content given multiple views into a neural network parameterized by θ , f_θ . The 213
214 network takes a 5D input vector (x, y, z, θ, ϕ) composed of a 3D position vector 214
215 and a view direction to encode view-dependent visual properties such as reflec- 215
216 tion and specularity. Both the position and view direction are encoded using a 216
217 deterministic positional encoding γ to recover high-frequency details as demon- 217
218 strated in a follow-up work by the authors [55]. NeRF's internal neural network 218
219 is a simple 8 layer deep and 256 neurons wide Multi Layer Perceptron (MLP). 219
220 A complete depiction of the architecture is shown in Figure 2. 220

222 **Volumetric Rendering.** NeRF uses a classical volume rendering pipeline with 222
223 alpha blending to render a frame. Rays are cast from the camera origin towards 223
224 the center of each pixel. To this end, 64 positions are uniformly sampled along 224

6 ACCV-22 submission ID ***

225 the ray direction. This step is referred to as the coarse step, which is part of the 225
 226 hierarchical sampling scheme proposed by Mildenhall et al. [5]. Every sample 226
 227 x_i is evaluated given the ray direction d by a coarse MLP f_c to output the 227
 228 corresponding density σ_i and the RGB radiance c_i . 228

229

$$230 \quad (\sigma_i, c_i) = f_\theta(\gamma(x_i), \gamma(d)) \quad 1 \leq i \leq 64 \quad (1) \quad 230 \\ 231$$

232 The ray evaluations are then composed using a differentiable formulation of 232
 233 alpha blending to output the coarse pixel color \hat{c}_c . 233

234

$$235 \quad \hat{c} = \sum_{i=1}^{64} T_i \alpha_i c_i \quad T_i = \prod_{j=1}^{64} 1 - \alpha_j \quad \alpha_i = 1 - \exp(-\sigma_i \delta_i) \quad (2) \quad 235 \\ 236 \\ 237 \\ 238$$

239 T describes the transmittance and δ the distance between adjacent sample 239
 240 positions $\delta_i = \|x_{i+1} - x_i\|$. The coarse evaluation is finally used to compute a 240
 241 probability density function along the ray to generate 64 new fine samples that 241
 242 are more likely to be close to the scene geometry. A second network, the fine 242
 243 network f_f is used to evaluate both coarse and fine samples to obtain a refined 243
 244 evaluation of the final pixel color \hat{c}_f . 244

245 In total $(64 + 128) \times H \times W$ network evaluations are performed during the 245
 246 generation of one frame, thus $122M$ evaluation when considering a 800×800 246
 247 image. Only the fine network is used during inference. 247

248

249 **Optimization.** NeRF optimizes its coarse and fine networks in an end-to-end 249
 250 fashion by computing an L_2 loss in image space against ground truth. 250

251

$$252 \quad \mathcal{L} = \frac{1}{P} \sum_{p=1}^P \|c_p - \hat{c}_{cp}\|_2 + \|c_p - \hat{c}_{fp}\|_2 \quad (3) \quad 252 \\ 253 \\ 254 \\ 255 \\ 256$$

257 4 Proposed Method 257

258

259 Our proposed method, DistillNeRF, aims at breeding smaller capacity Neural 259
 260 Radiance Fields (NeRF) Multi Layer Perceptron (MLP)s by initializing them 261
 261 from pretrained high capacity models. We propose to use Knowledge Distilla- 262
 262 tion (KD) to transfer the knowledge captured by the weights of a large capacity 263
 263 network called the teacher to a low capacity network called the student. During 264
 264 this KD phase, the student is trained on samples produced by the teacher net- 265
 265 work to match its output distribution. The quality of the student is thus bound 266
 266 by the teacher's examples quality. 267

268 This section describes how the teacher model is trained, how to generate a 268
 269 student candidate, and how to sample good teacher examples. 269

4.1 Teacher Model

The teacher model follows the same architecture as the vanilla NeRF described in Figure 2 with modified activation functions as proposed by Barron et al. in Mip-NeRF [51]. The softplus activation is replaced by a shifted version, Equation 4, and the sigmoid activation by a widened version, Equation 5.

$$\text{softplus}(x) = \log(1 + \exp(x - 1))$$

$$\text{sigmoid}(x) = \frac{1 + 2\epsilon}{1 + \exp(-x))} - \epsilon \quad \epsilon = 1e-3$$

Contrary to the original work [5], we train the same network f_θ for both coarse and fine steps. Using only one network during training reduces memory usage by half and network evaluation from $64 + 128 = 192$ to 128 per pixel. Training is thus less resource-intensive and faster.

Table 1: Models architectures. The width describes the number of neurons in each layer of the NeRF internal neural network, the depth the number of hidden layers, and residual the presence or absence of a residual connection in the model. The models are displayed in descending order of size (MB). Smaller models naturally demonstrate higher render frame rates (FPS). However, low capacity results in perceptual degradation (PSNR). Results are averaged over the 8 scenes of the Blender Synthetic Dataset.

Model	Width	Depth	Residual	PSNR \uparrow	FPS \downarrow	MB \downarrow
NeRF	256	8	✓	30.01	0.28	2.38
TinyNeRF	128	4	✓	26.37	0.94	0.37
MicroNeRF	64	4	✗	24.09	1.71	0.10
NanoNeRF	32	2	✗	23.45	2.16	0.03

4.2 Student Model

The architecture of NeRF internal neural network is a sample MLP described in Figure 2. Reducing its inference time requires either a reduction of capacity, the number of parameters θ , or the use of faster operation routes in the network graph.

DistillNeRF generates student networks by reducing three hyperparameters: the MLP width, depth, and an eventual residual connection. Following this recipe, we propose three new flavors of NeRF: TinyNeRF, MicroNeRF, and NanoNeRF in decreasing order of size. Their respective hyperparameters are displayed in Table 1. The parameters are chosen to reflect a wide range of student

315 variations. Smaller models naturally benefit from faster inference and rendering 315
 316 time. NeRF renders 400×400 images at 0.28 Frames per Second (FPS) com- 316
 317 pared to 2.16 FPS for NanoNeRF without the use of any additional acceleration 317
 318 technique. 318
 319 319

320 4.3 Knowledge Distillation 320

321 The student candidates are initialized using KD. Their weights are distilled 322
 322 from the pretrained NeRF teacher model. The teacher examples are randomly 323
 323 sampled from a Bounding Box (BBox) centered on the scene geometry. The 324
 324 BBox is described by a set of six parameters $((x_a, x_b), (y_a, y_b), (z_a, z_b))$, two for 325
 325 each axis initialized to a value of $(-20, 20)$. Random positions and direction 326
 326 vectors are then sampled inside this structure. 327

327 The student models are trained to match the teacher’s predictions. We opti- 328
 328 mize the student’s parameters to minimize a mixture of L_2 losses for density σ 329
 329 and radiance c . The alpha blending values α are computed instead of the density 330
 330 output σ using $\delta = 3.125e-2$, a value close to the ones encountered during train- 331
 331 ing. This computational change, introduced by Reiser et al. in KiloNeRF [10], 332
 332 reduces the emphasis on small differences between big densities. 333

$$\mathcal{L}_{KD} = \frac{1}{2}(\|\hat{c}_t - \hat{c}_s\|_2 + \|\hat{\alpha}_t - \hat{\alpha}_s\|_2) \quad (6)$$



338 Fig. 3: Visualization of the occupancy voxel grid slices normalized between 0 and 338
 339 1 for a vanilla NeRF model trained on the Lego scene from the Blender Synthetic 339
 340 Dataset. 340
 341 341
 342 342
 343 343

350 4.4 Priority Occupancy Voxel Grid 350

351 NeRF learns to handle empty space in the early stage of its training process. 352
 352 Thus, we optimize the data-efficiency of our distillation process by introducing 353
 353 a priority occupancy voxel grid to emphasize non-empty space. The occupancy 354
 354 scores act as a rough description of the scene content distribution. An example 355
 355 visualization is displayed in Figure 3. 356

356 To fill this grid, the BBox is first refined to make it fit the geometry bounds 357
 357 more tightly. The structure is turned into a $256 \times 256 \times 256$ voxel grid with 358
 358 256 bins on each axis. Each voxel is then filled with a 32 samples Monte Carlo 359

360 approximation of the mean density σ evaluated through the teacher network. 360
 361 The grid is normalized to values between 0 and 1. At this stage of the process, 361
 362 the 3D grid cells contain scores describing the occupancy of the scene. The BBox 362
 363 is finally tightened to the closest cell with a score greater or equal to a selected 363
 364 threshold $\tau = 1e - 2$. 364

365 The exact same process is used to generate a $16 \times 16 \times 16$ low-resolution 365
 366 priority occupancy voxel grid with 16 bins. The probability score of each cell 366
 367 is clamped to a minimum value of 0.2 to avoid forgetting empty space during 367
 368 the distillation process. During KD, the position vectors are generated by draw- 368
 369 ing random cells taking into account their priority, their score, and randomly 369
 370 offset by the voxel size. Finally, the ray directions are generated by sampling a 370
 371 normalized unit sphere. 371

372 372

373 5 Implementation 373

374 374
 375 DistillNeRF is implemented using the PyTorch Python library. While this ap- 375
 376 proach allows fast prototyping and training, a naive implementation is not suited 376
 377 for production and real-time applications. This sections describes a set of strate- 377
 378 gies we employ to reduce inference time and accelerate volume rendering. 378
 379

380 **Inference.** For inference, we reduce the Multi Layer Perceptron (MLP) weights 380
 381 to half-precision, FP16. This precision reduction has the benefit of halving mem- 381
 382 ory usage at training and inference time. The model is then exported as a Torch- 382
 383 Script module to be used in a C++ environment thanks to the LibTorch API. 383
 384 The batch size is also maxed out during inference, 8,192 on an Nvidia GeForce 384
 385 RTX 3090. 385

386 386

387 **Rendering.** We also implement a set of specialized Cuda kernels to handle 387
 388 the volumetric rendering process, including the ray casting, the alpha blending 388
 389 computation, and accumulation. The final pixel colors are then written into 389
 390 Cuda memory to the VRAM and bound to an OpenGL 2D texture. The texture 390
 391 is finally projected on a quad fitting the screen. Every computation is done on 391
 392 the GPU to limit RAM to VRAM memory transfers and maximize utilization. 392

393 393

394 6 Evaluation Setup 394

395 395
 396 **Data.** We evaluate and train our method on the Blender Synthetic Dataset 396
 397 proposed by Neural Radiance Fields (NeRF)'s authors [5]. The dataset gathers 397
 398 a set of 8 scenes: chair, drums, ficus, hotdog, lego, materials, mic, and ship, each 398
 399 containing 100 views of the scene for training, 100 for validation, and 200 for 399
 400 testing. The scenes are rendered using the Blender Cycles ray tracing engine, and 400
 401 the views are sampled from a centered hemisphere surface with the view direction 401
 402 targeting the center. The rendered images are of size 800×800 and contain 402
 403 various challenging materials, including thin, transparent, and semi-transparent 403
 404 objects. The generated images are reduced to 400×400 using nearest pixel 404

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

10 ACCV-22 submission ID ***

405 values to make the training of every model variation and initialization methods 405
 406 affordable. 406

407 **Metrics.** The NeRF literature reports evaluations with a set of objective met- 407
 408 rics such as the Mean Squared Error (MSE), the Peak Signal to Noise Ratio 408
 409 (PSNR), Structural Similarity (SSIM) [56] and the Learned Perceptual Image 409
 410 Patch Similarity (LPIPS) [57]. The MSE evaluates the mean squared distance of 410
 411 the generated and the target image at a pixel level. The PSNR is derived from 411
 412 the MSE metric and measures the ratio between the maximum possible power of 412
 413 the image and the power of corrupting noise affecting the fidelity of its approxi- 413
 414 mation. The SSIM measures the similarity between a pair of images at different 414
 415 window sizes. The LPIPS makes use of a pretrained foundation model, AlexNet, 415
 416 and measures the perceptual feature differences between the two images. 416

417 In their contribution, Barron et al. [51], demonstrate a strong correlation 417
 418 between the LPIPS and SSIM. The LPIPS is not measured in this paper to 418
 419 limit redundancy. 419

420 We also report the models' sizes in MB and their rendering performance 420
 421 in Frames per Second (FPS) including the ray generation, casting, the model 421
 422 inference, the ray radiance and alpha blending aggregation. 422

423 **Baseline.** DistillNeRF is an initialization method used to breed smaller NeRF 423
 424 models for faster inference. In this setup, we evaluate our method against two 424
 425 baselines, the vanilla NeRF [5], and the meta-learning Reptile initialization [49]. 425
 426 NeRF, TinyNeRF, MicroNeRF, and NanoNeRF are benchmarked on the three 426
 427 methods. 427

428 **Training.** The models are trained using 64 coarse and 64 fine samples with 428
 429 half-precision enabled. We optimize NeRF using the same regime as the original 429
 430 implementation for 20 epochs using the Adam optimizer with a $5e - 4$ learning 430
 431 rate which is decayed to $5e - 6$ over time using a log scheduler. Before training, 431
 432 the model's weights are initialized for 1 epoch with 16 samples using the SGD 432
 433 optimizer with a $5e - 2$ learning rate for Reptile, and 100,000 steps using the 433
 434 Adam optimizer with a $5e - 4$ learning rate for Knowledge Distillation (KD). The 434
 435 training hyperparameters are chosen to match the initialization time required 435
 436 for both methods. 436

437 **Measurements.** Training and inference measurements are realized on a modern 437
 438 station equipped with an AMD Ryzen 9 5900X 12-Core Processor, 32 Go of 438
 439 DDR4 RAM, and an NVidia GeForce RTX 3090 with 24 Go of VRAM. 439

440 7 Results 440

441 We quantitatively in Table 2 and qualitatively in Figure 5, 6, and 4 show that 441
 442 our initialization method DistillNeRF outperforms the meta-learning Reptile 442
 443 initialization in Mean Squared Error (MSE), Peak Signal to Noise Ratio (PSNR), 443
 444 and Structural Similarity (SSIM). We also provide an ablation study in Table 3 444
 445 to justify our current design strategy. 445

450
 451 Table 2: Benchmark on Blender Synthetic Dataset re-scaled to 400×400 . Re-
 452 results are averaged on the 8 scenes of the dataset. Our approach, highlighted in
 453 light gray, improves the performances on MSE, PSNR, and SSIM compared to
 454 ReptileNeRF for similar training regimes and presents lower quality degradation
 455 on smaller models.

456 Model	456 Train (h:m:s) ↓	456 MSE ↓	456 PSNR ↑	456 SSIM ↑
457 NeRF	02:44:58	1.24e-03	30.01	0.94
458 ReptileNeRF	03:59:56	1.17e-03	30.90	0.95
459 DistillNeRF	03:02:49	1.15e-03	31.02	0.95
460 TinyNeRF	01:08:40	2.15e-03	26.37	0.91
461 ReptileTinyNeRF	01:35:57	1.98e-03	28.29	0.91
462 DistillTinyNeRF	01:25:46	1.96e-03	28.39	0.92
463 MicroNeRF	00:39:51	2.98e-03	24.09	0.87
464 ReptileMicroNeRF	00:52:24	2.80e-03	26.51	0.88
465 DistillMicroNeRF	00:56:55	2.63e-03	26.82	0.89
466 NanoNeRF	00:32:27	3.97e-03	23.45	0.85
467 ReptileNanoNeRF	00:40:56	3.87e-03	25.01	0.85
468 DistillNanoNeRF	00:49:37	3.53e-03	25.36	0.86

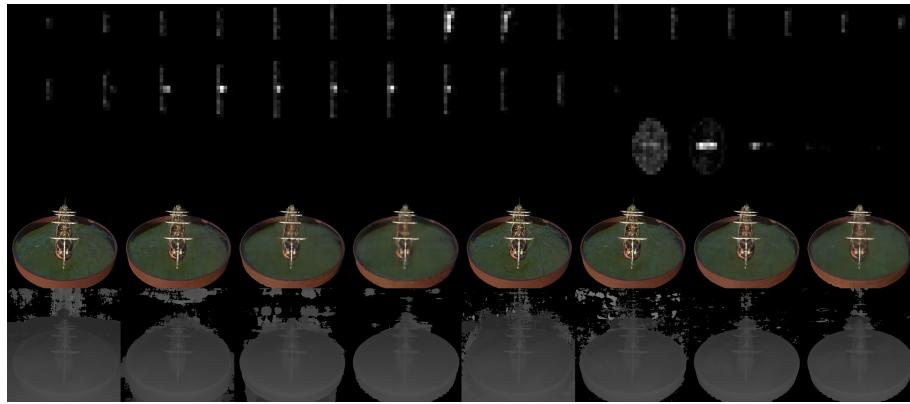
470
 471
 472 **Benchmark.** For similar training regimes and approximately the same training
 473 duration, our initialization method DistillNeRF outperforms the Reptile meta
 474 initialization procedure on average on MSE, PSNR and SSIM on the Blender
 475 Synthetic Dataset.

476
 477
 478
 479 Table 3: Ablation study on the Blender Synthetic Dataset hot dog scene rescaled
 480 to 400×400 . The table studies the implication of the occupancy voxel grid used
 481 of importance sampling and its refinement. The ablation is realised on the most
 482 challenging NeRF variation, NanoNeRF.

483 Model	483 MSE ↓	483 PSNR ↑	483 SSIM ↑
485 NanoNeRF + KD	1.41e-3	28.66	0.92
486 NanoNeRF + KD + Voxel Grid	1.21e-3	29.25	0.92
487 NanoNeRF + KD + Voxel Grid + Refined	1.05e-3	29.89	0.93

488
 489
 490
 491
 492 **Ablation.** The ablation study evaluates the impact of our design choices. We
 493 report the MSE, PSNR and SSIM on the Blender Synthetic hot dog scene for
 494 NanoNeRF, our most challenging model variation. The results shows that the

495 use of a well fitted occupancy voxel grid for importance sampling improves the 495
 496 standard Knowledge Distillation (KD) strategy.
 497



500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 Fig. 4: Renders of the Blender Synthetic Dataset ship scene. The top row shows
 512 the occupancy voxel grid, the middle are renders from NeRF variations, and
 513 the bottom are their corresponding depth maps. The models used are NeRF,
 514 TinyNeRF, MicroNeRF, NanoNeRF, DistillNeRF, DistillTinyNeRF, DistillMi-
 515 croNeRF, and DistillNanoNeRF respectively from left to right.
 516

520 8 Conclusion

521
 522 Our work DistillNeRF addresses the need for faster rendering performances in
 523 Neural Radiance Fields (NeRF) methods. While the current state of the art
 524 has focused on reducing the number of sample evaluations performed to train
 525 such models, we explored the reduction of the inner Multi Layer Perceptron
 526 (MLP) size. Smaller models naturally exhibit faster inference but often result
 527 in a quality drop. We demonstrate that the use of Knowledge Distillation (KD)
 528 in conjunction with an occupancy voxel grid as a proxy for the importance-
 529 sampling of teacher examples reduces the perceptual loss induced by a small
 530 capacity network. Our method can be combined with other approaches such as
 531 acceleration structure caching, sampling oracles, and others to allow real-time
 532 rendering at a smaller footprint.
 533

534 **Future Work** The use of KD for initialization requires training a teacher net-
 535 work which constitutes a consequent overhead. In future work, we want to explore
 536 the use of other proxies such as the information contained in the training images
 537 to initialize NeRF approaches at a smaller cost.
 538

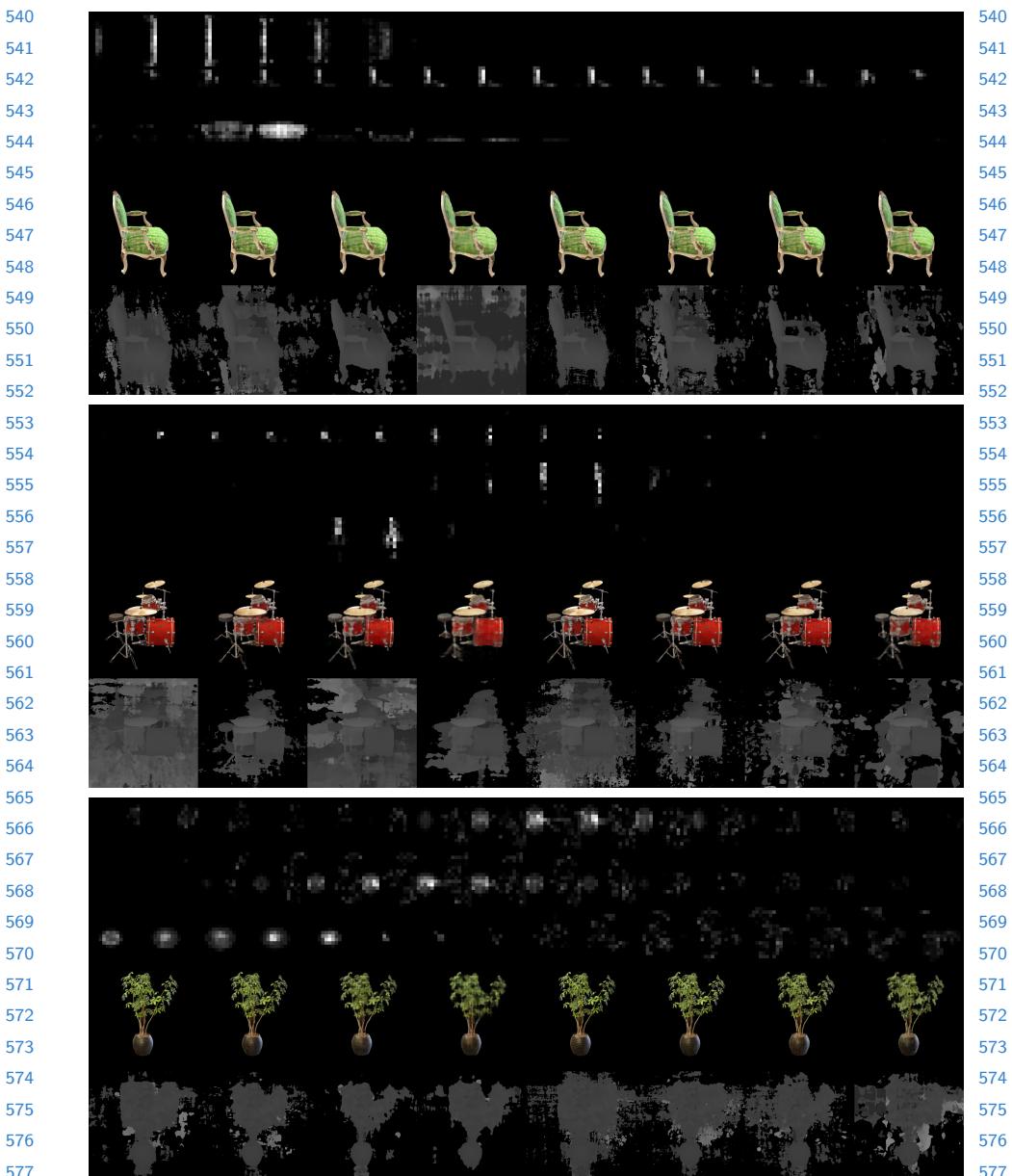


Fig. 5: Renders of the Blender Synthetic Dataset chair, drums and ficus scenes. The top row shows the occupancy voxel grid, the middle are renders from NeRF variations, and the bottom are their corresponding depth maps. The models used are NeRF, TinyNeRF, MicroNeRF, NanoNeRF, DistillNeRF, DistillTinyNeRF, DistillMicroNeRF, and DistillNanoNeRF respectively from left to right.

584

579

580

581

582

583

584

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

14 ACCV-22 submission ID ***

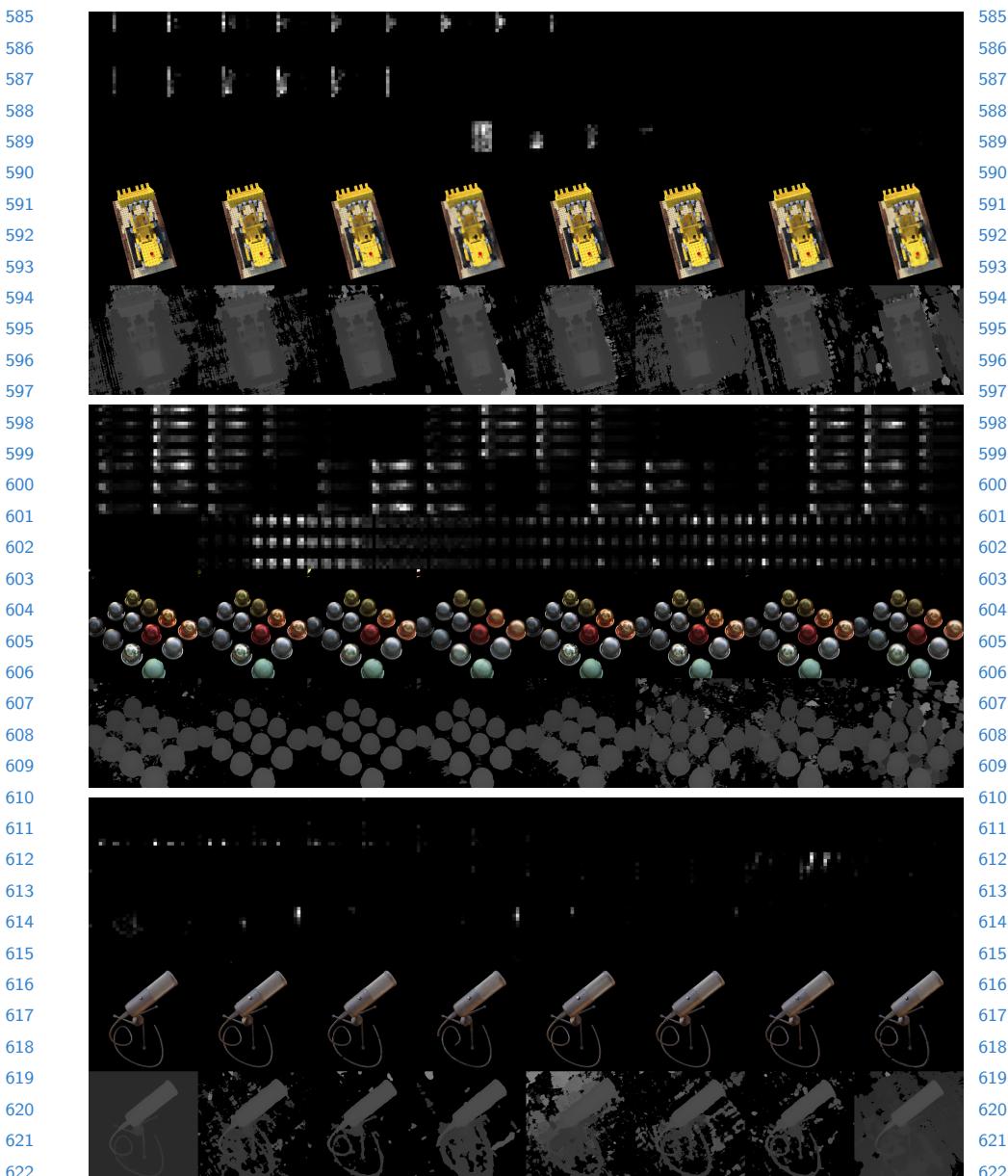


Fig. 6: Renders of the Blender Synthetic Dataset lego, materials and mic scenes. The top row shows the occupancy voxel grid, the middle are renders from NeRF variations, and the bottom are their corresponding depth maps. The models used are NeRF, TinyNeRF, MicroNeRF, NanoNeRF, DistillNeRF, DistillTinyNeRF, DistillMicroNeRF, and DistillNanoNeRF respectively from left to right.

624
625
626
627
628
629

References

- 630

631 References

632 1. Li, T., Slavcheva, M., Zollhofer, M., Green, S., Lassner, C., Kim, C., Schmidt, T., Lovegrove, S., Goesele, M., Lv, Z.: Neural 3d video synthesis. arXiv e-prints (2021) arXiv–2103

633 2. Zhang, J., Liu, X., Ye, X., Zhao, F., Zhang, Y., Wu, M., Zhang, Y., Xu, L., Yu, J.: Editable free-viewpoint video using a layered neural representation. ACM Transactions on Graphics (TOG) **40** (2021) 1–18

634 3. Gafni, G., Thies, J., Zollhofer, M., Nießner, M.: Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2021) 8649–8658

635 4. Xu, H., Alldieck, T., Sminchisescu, C.: H-nerf: Neural radiance fields for rendering and temporal reconstruction of humans in motion. Advances in Neural Information Processing Systems **34** (2021)

636 5. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M., eds.: Computer Vision – ECCV 2020, Cham, Springer International Publishing (2020) 405–421

637 6. Neff, T., Stadlbauer, P., Parger, M., Kurz, A., Mueller, J.H., Chaitanya, C.R.A., Kaplanyan, A., Steinberger, M.: Donerf: Towards real-time rendering of compact neural radiance fields using depth oracle networks. arXiv preprint arXiv:2103.03231 (2021)

638 7. Arandjelović, R., Zisserman, A.: Nerf in detail: Learning to sample for view synthesis. arXiv preprint arXiv:2106.05264 (2021)

639 8. Rebain, D., Jiang, W., Yazdani, S., Li, K., Yi, K.M., Tagliasacchi, A.: Derf: Decomposed radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2021) 14153–14161

640 9. Liu, L., Gu, J., Lin, K.Z., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. arXiv preprint arXiv:2007.11571 (2020)

641 10. Reiser, C., Peng, S., Liao, Y., Geiger, A.: Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. arXiv preprint arXiv:2103.13744 (2021)

642 11. Hedman, P., Srinivasan, P.P., Mildenhall, B., Barron, J.T., Debevec, P.: Baking neural radiance fields for real-time view synthesis. arXiv preprint arXiv:2103.14645 (2021)

643 12. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenoctrees for real-time rendering of neural radiance fields. arXiv preprint arXiv:2103.14024 (2021)

644 13. Garbin, S.J., Kowalski, M., Johnson, M., Shotton, J., Valentin, J.: Fastnerf: High-fidelity neural rendering at 200fps. arXiv preprint arXiv:2103.10380 (2021)

645 14. Buehler, C., Bosse, M., McMillan, L., Gortler, S., Cohen, M.: Unstructured lumigraph rendering. In: Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques. SIGGRAPH ’01, New York, NY, USA, Association for Computing Machinery (2001) 425–432

646 15. Debevec, P.E., Taylor, C.J., Malik, J.: Modeling and rendering architecture from photographs: A hybrid geometry- and image-based approach. In: Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques. SIGGRAPH ’96, New York, NY, USA, Association for Computing Machinery (1996) 11–20

647 16. Waechter, M., Moehrle, N., Goesele, M.: Let there be color! large-scale texturing of 3d reconstructions. In: ECCV. (2014)

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

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

16 ACCV-22 submission ID ***

- 675 17. Wood, D.N., Azuma, D.I., Aldinger, K., Curless, B., Duchamp, T., Salesin, D.H., 675
676 Stuetzle, W.: Surface light fields for 3d photography. In: Proceedings of the 676
677 27th Annual Conference on Computer Graphics and Interactive Techniques. SIG- 677
678 GRAPH '00, USA, ACM Press/Addison-Wesley Publishing Co. (2000) 287–296 678
- 679 18. Aliev, K.A., Sevastopolsky, A., Kolos, M., Ulyanov, D., Lempitsky, V.: Neural 679
680 point-based graphics. In Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M., eds.: 680
681 Computer Vision – ECCV 2020, Cham, Springer International Publishing (2020) 681
682 696–712
- 683 19. Riegler, G., Koltun, V.: Free view synthesis. In: European Conference on Computer 682
684 Vision, Springer (2020) 623–640
- 685 20. Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Proceedings 684
686 of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 685
687 (2016)
- 688 21. Song, Z., Chen, W., Campbell, D., Li, H.: Deep novel view synthesis from colored 3d 687
689 point clouds. In Vedaldi, A., Bischof, H., Brox, T., Frahm, J., eds.: Computer Vision 688
690 - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23–28, 2020, 689
691 Proceedings, Part XXIV. Volume 12369 of Lecture Notes in Computer Science., 690
692 Springer (2020) 1–17
- 693 22. Flynn, J., Broxton, M., Debevec, P., DuVall, M., Fyffe, G., Overbeck, R., Snavely, 691
694 N., Tucker, R.: Deepview: View synthesis with learned gradient descent. In: Pro- 692
695 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog- 693
696 nition. (2019) 2367–2376
- 697 23. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, 695
698 R., Ng, R., Kar, A.: Local light field fusion: Practical view synthesis with prescrip- 696
699 tive sampling guidelines. ACM Trans. Graph. **38** (2019)
- 700 24. Srinivasan, P.P., Mildenhall, B., Tancik, M., Barron, J.T., Tucker, R., Snavely, 697
701 N.: Lighthouse: Predicting lighting volumes for spatially-coherent illumination. 698
702 In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 699
703 Recognition. (2020) 8080–8089
- 704 25. Zhou, T., Tucker, R., Flynn, J., Fyffe, G., Snavely, N.: Stereo magnification: 701
705 Learning view synthesis using multiplane images. arXiv preprint arXiv:1805.09817 702
706 (2018)
- 707 26. Kutulakos, K., Seitz, S.: A theory of shape by space carving. In: Proceedings of 704
708 the Seventh IEEE International Conference on Computer Vision. Volume 1. (1999) 705
709 307–314 vol.1
- 710 27. Lombardi, S., Simon, T., Saragih, J., Schwartz, G., Lehrmann, A., Sheikh, Y.: 707
711 Neural volumes: Learning dynamic renderable volumes from images. ACM Trans. 708
712 Graph. **38** (2019)
- 713 28. Penner, E., Zhang, L.: Soft 3d reconstruction for view synthesis. ACM Trans. 709
714 Graph. **36** (2017)
- 715 29. Seitz, S., Dyer, C.: Photorealistic scene reconstruction by voxel coloring. In: Pro- 711
716 ceedings of IEEE Computer Society Conference on Computer Vision and Pattern 712
717 Recognition. (1997) 1067–1073
- 718 30. Sitzmann, V., Thies, J., Heide, F., Nießner, M., Wetstein, G., Zollhofer, M.: 714
719 Deepvoxels: Learning persistent 3d feature embeddings. In: Proceedings of the 715
720 IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2019) 2437– 716
721 2446
- 722 31. Szeliski, R., Golland, P.: Stereo matching with transparency and matting. In: 717
723 Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271). 718
724 (1998) 517–524

- 720 32. Liu, S., Chen, W., Li, T., Li, H.: Soft rasterizer: Differentiable rendering for unsu- 720
721 pervised single-view mesh reconstruction. arXiv preprint arXiv:1901.05567 (2019) 721
- 722 33. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumet- 722
723 ric rendering: Learning implicit 3d representations without 3d supervision. In: 723
724 Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). 724
(2020)
- 725 34. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning 725
726 continuous signed distance functions for shape representation. In: Proceedings of 726
727 the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2019) 727
165–174
- 728 35. Park, K., Mousavian, A., Xiang, Y., Fox, D.: Latentfusion: End-to-end differ- 728
729 entiable reconstruction and rendering for unseen object pose estimation. In: Pro- 729
730 ceedings of the IEEE/CVF conference on computer vision and pattern recognition. 730
731 (2020) 10710–10719
- 732 36. Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Basri, R., Lipman, Y.: 732
733 Multiview neural surface reconstruction by disentangling geometry and appear- 733
734 ance. arXiv preprint arXiv:2003.09852 (2020)
- 735 37. Chibane, J., Mir, A., Pons-Moll, G.: Neural unsigned distance fields for implicit 735
736 function learning. arXiv preprint arXiv:2010.13938 (2020)
- 737 38. Liu, S., Zhang, Y., Peng, S., Shi, B., Pollefeys, M., Cui, Z.: Dist: Rendering deep 737
738 implicit signed distance function with differentiable sphere tracing. In: Proceedings 738
739 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. (2020) 739
2019–2028
- 740 39. Jain, A., Mildenhall, B., Barron, J.T., Abbeel, P., Poole, B.: Zero-shot text-guided 740
741 object generation with dream fields. arXiv preprint arXiv:2112.01455 (2021)
- 742 40. Jang, W., Agapito, L.: Codenerf: Disentangled neural radiance fields for object cat- 742
743 egories. In: Proceedings of the IEEE/CVF International Conference on Computer 743
744 Vision. (2021) 12949–12958
- 745 41. Saito, S., Huang, Z., Natsume, R., Morishima, S., Kanazawa, A., Li, H.: Pifu: Pixel- 745
746 aligned implicit function for high-resolution clothed human digitization. In: Pro- 746
747 ceedings of the IEEE/CVF International Conference on Computer Vision. (2019) 747
2304–2314
- 748 42. Boss, M., Braun, R., Jampani, V., Barron, J.T., Liu, C., Lensch, H.: Nerd: 748
749 Neural reflectance decomposition from image collections. In: Proceedings of the 749
750 IEEE/CVF International Conference on Computer Vision. (2021) 12684–12694
- 751 43. Rudnev, V., Elgharib, M., Smith, W., Liu, L., Golyanik, V., Theobalt, C.: Neural 751
752 radiance fields for outdoor scene relighting. arXiv preprint arXiv:2112.05140 (2021)
- 753 44. Srinivasan, P.P., Deng, B., Zhang, X., Tancik, M., Mildenhall, B., Barron, J.T.: 753
754 Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In: 754
755 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog- 755
756 nition. (2021) 7495–7504
- 757 45. Zhang, X., Srinivasan, P.P., Deng, B., Debevec, P., Freeman, W.T., Barron, J.T.: 757
758 Nerfactor: Neural factorization of shape and reflectance under an unknown illumina- 758
759 nation. arXiv preprint arXiv:2106.01970 (2021)
- 760 46. Guo, J., Yang, Z., Lin, X., Zhang, Q.: Template nerf: Towards modeling dense 760
761 shape correspondences from category-specific object images. arXiv preprint 761
arXiv:2111.04237 (2021)
- 762 47. Su, S.Y., Yu, F., Zollhöfer, M., Rhodin, H.: A-nerf: Articulated neural radiance 762
763 fields for learning human shape, appearance, and pose. Advances in Neural Infor- 763
764 mation Processing Systems **34** (2021)

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

18 ACCV-22 submission ID ***

- 765 48. Yen-Chen, L., Florence, P., Barron, J.T., Rodriguez, A., Isola, P., Lin, T.Y.: [765](#)
766 inerf: Inverting neural radiance fields for pose estimation. arXiv preprint [766](#)
767 arXiv:2012.05877 (2020) [767](#)
- 768 49. Tancik, M., Mildenhall, B., Wang, T., Schmidt, D., Srinivasan, P.P., Barron, J.T., [768](#)
769 Ng, R.: Learned initializations for optimizing coordinate-based neural represen- [769](#)
770 tations. In: Proceedings of the IEEE/CVF Conference on Computer Vision and [770](#)
771 Pattern Recognition. (2021) 2846–2855 [771](#)
- 772 50. Nichol, A., Achiam, J., Schulman, J.: On first-order meta-learning algorithms. [771](#)
773 arXiv preprint arXiv:1803.02999 (2018) [772](#)
- 774 51. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srin- [773](#)
775 ivasan, P.P.: Mip-nerf: A multiscale representation for anti-aliasing neural radiance [774](#)
776 fields. arXiv preprint arXiv:2103.13415 (2021) [775](#)
- 777 52. Zhang, K., Riegler, G., Snavely, N., Koltun, V.: Nerf++: Analyzing and improving [776](#)
778 neural radiance fields. arXiv preprint arXiv:2010.07492 (2020) [777](#)
- 779 53. Martin-Brualla, R., Radwan, N., Sajjadi, M.S., Barron, J.T., Dosovitskiy, A., Duck- [778](#)
780 worth, D.: Nerf in the wild: Neural radiance fields for unconstrained photo col- [779](#)
781 lections. In: Proceedings of the IEEE/CVF Conference on Computer Vision and [780](#)
782 Pattern Recognition. (2021) 7210–7219 [781](#)
- 783 54. Athar, S., Shu, Z., Samaras, D.: Flame-in-nerf: Neural control of radiance fields [781](#)
784 for free view face animation. arXiv preprint arXiv:2108.04913 (2021) [782](#)
- 785 55. Tancik, M., Srinivasan, P.P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Sing- [783](#)
786 hal, U., Ramamoorthi, R., Barron, J.T., Ng, R.: Fourier features let networks learn [784](#)
787 high frequency functions in low dimensional domains. NeurIPS (2020) [785](#)
- 788 56. Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E.: Image quality assessment: from [786](#)
789 error visibility to structural similarity. IEEE Transactions on Image Processing **13** [787](#)
790 (2004) 600–612 [788](#)
- 791 57. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable [789](#)
792 effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE [790](#)
793 conference on computer vision and pattern recognition. (2018) 586–595 [791](#)
- 794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809