
NerfBaselines: Consistent and Reproducible Evaluation of Novel View Synthesis Methods

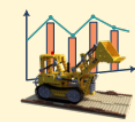
Jonas Kulhanek
CTU in Prague and ETH Zurich
Czech Republic and Switzerland
jonas.kulhanek@cvut.cz

Torsten Sattler
CTU in Prague
Czech Republic
torsten.sattler@cvut.cz

Implemented datasets

- Standardized dataset format and coordinate systems
- Implements common loaders, e.g. COLMAP, NerfStudio, Bundler, etc.
- Standardized evaluation protocols matching the original papers

NerfBaselines



- Unified API to access methods
- Better reproducibility
- Simplified installation
- Online benchmark



Interactive viewer

Integrated methods

- Unified interface for NeRFs and 3DGS
- Original codebases are used such that reported metrics match the papers
- Each method installed into isolated encapsulated environment

Isolated encapsulated environments



Online benchmark

Abstract

1 Novel view synthesis is an important problem with many applications, including
2 AR/VR, gaming, and simulations for robotics. With the recent rapid development of
3 Neural Radiance Fields (NeRFs) and 3D Gaussian Splatting (3DGS) methods, it is
4 becoming difficult to keep track of the current state of the art (SoTA) due to methods
5 using different evaluation protocols, codebases being difficult to install and use, and
6 methods not generalizing well to novel 3D scenes. Our experiments support this
7 claim by showing that tiny differences in evaluation protocols of various methods
8 can lead to inconsistent reported metrics. To address these issues, we propose
9 a framework called NerfBaselines, which simplifies the installation of various
10 methods, provides consistent benchmarking tools, and ensures reproducibility. We
11 validate our implementation experimentally by reproducing numbers reported in
12 the original papers. To further improve the accessibility, we release a web platform
13 where commonly used methods are compared on standard benchmarks.
14 **Web:** <https://jkulhanek.com/nerfbaselines>

15 1 Introduction

16 Due to the explosive growth in the number of NeRF-based and 3DGS-based works, with about 10
17 new papers citing NeRF [37] and 3 papers citing 3DGS [19] each day, there isn't a single evaluation
18 framework and the individual works are based on different codebases for evaluation, resulting in
19 numbers that sometimes are not directly comparable. Different methods sometimes use different
20 evaluation protocols for the same datasets and often use specific data processing or hyperparameters
21 which do not translate well to novel datasets (not used in the original papers). Furthermore, methods
22 are often difficult to install with conflicting dependencies, and the codebases differ significantly,

23 which makes custom visualization or benchmarking laborious. To address these issues, we propose a
24 framework designed to standardize and simplify the evaluation and comparison of these methods.

25 Not all papers follow the same evaluation protocol (partly because it is not fully specified) and
26 sometimes innocuous changes can have a significant impact on the overall results. Examples
27 include methods using different image resolutions for evaluation, different parameters for image
28 downscaling, or different parameters for computing metrics. In our experiments, we show that
29 even small discrepancies can cause substantial differences in the metrics, validating the need for a
30 framework such as NerfBaselines. Our framework standardizes evaluation protocols, following best
31 practices for each benchmark dataset.

32 Even with access to original source codes, reproducing paper results can be challenging due to
33 codebases evolving over time or some codebases becoming obsolete where it is no longer possible to
34 install the dependencies. To tackle this, for each method, we provide a consistent and reproducible way
35 of installation (by encapsulating the runtime environment) which enables us to release checkpoint
36 that can be used to render from custom camera positions without the need to rerun the training.
37 Furthermore, we release a website with benchmarks comparing various methods on multiple datasets.

38 Many existing codebases cannot be easily applied to novel datasets due to having simplifying as-
39 sumptions in the code, e.g., not implementing various camera models, requiring uniform camera
40 intrinsic, or image sizes being equal. Various methods also use different interfaces, dataset formats,
41 and coordinate systems, making integration and custom rendering laborious. Therefore, in NerfBase-
42 lines, we have integrated some important and influential codebases and implemented the missing
43 functionality. Thankfully, most codebases are derived from the few original repositories (often with
44 minor modifications), and therefore, they can also be integrated easily. Our experiments show that
45 our integrated methods match the original ones with sufficient precision.

46 Furthermore, most methods compare using a set of metrics (e.g., PSNR) computed on a set of
47 images taken as a subset of the training trajectory, which does not necessarily correlate well with
48 the performance outside of the training trajectory [29]. Rendering videos from custom trajectories is
49 often much better at demonstrating multiview consistency. Therefore, we release a web-based camera
50 trajectory editor, based on Viser [53], to design camera trajectories and standardize their format. The
51 common interface allows all integrated methods to render novel images from these trajectories.

52 In summary, to simplify benchmarking and improve reproducibility, we propose:

- 53 • A unified interface for NeRF and 3DGS methods, standardizing dataset formats and evaluation
54 protocols.
- 55 • The NerfBaselines framework, which installs each method in an isolated environment to manage
56 dependencies and ensure reproducibility. We experimentally verify our integrated methods
57 reproduce results from original papers.
- 58 • A web-based camera trajectory editor and tools to render different camera trajectories, offering a
59 more comprehensive assessment of performance beyond traditional metrics like PSNR, SSIM [6],
60 or LPIPS [70].
- 61 • A web platform for comparing the performance of various methods on different datasets,
62 providing checkpoints, and enabling interactive viewing of results.

63 2 Literature review

64 In recent years, progress in novel view synthesis has enabled real-time photo-realistic rendering of
65 images from novel camera viewpoints. There has been a surge of interest, first with the advent of
66 neural radiance field methods (NeRFs) [37, 30] which enabled photo-realistic results and a wide range
67 of application. The second wave of attention came with the introduction of 3D Gaussian Splatting
68 [19], matching NeRFs in terms of the rendering quality while enabling real-time rendering.

69 **Ray-based methods.** NeRFs represent a scene using a continuous volumetric radiance field, originally
70 parameterized by a fully connected neural network [37, 2]. While this method enables high-quality

71 rendering for forward-facing or bounded and object-centric scenes, it initially suffered from aliasing
 72 issues, which were addressed in [2] and extended to unbounded scenes in [3, 69]. However, the MLPs
 73 used in NeRFs posed a bottleneck in terms of both training and rendering speeds. To address this,
 74 caching of the radiance field was proposed to enable faster rendering [16, 45], though this came at the
 75 cost of large storage requirements. Consequently, several methods have investigated replacing MLPs
 76 with alternatives such as *sparse grids* [14, 52], *point clouds* [61], *tetrahedral meshes* [24], *tensorial*
 77 *factorizations* [9, 15, 46], or *hash-grids* [39, 4, 56], often combined with tiny MLPs on top of the
 78 internal representation. While each approach has its advantages and disadvantages, hash-grids [39]
 79 have become the most popular representation due to their speed and scalability, with Zip-NeRF [4]
 80 achieving SoTA performance by addressing the aliasing issues of grid-based representations.

81 NeRFs have been extended in various ways, such as handling dynamic scenes [41, 44], modeling
 82 different image appearances [33, 15], and accommodating data with imprecise poses [5, 42]. They
 83 have also been applied to model semantics [21, 22] and for style transfer [54, 17]. To better model
 84 geometry, some methods replace the radiance field’s density with a signed distance function (SDF)
 85 [55, 57, 63]. Finally, another line of work focuses on extending these methods to generative [11, 43]
 86 or sparse-view settings [66, 8, 26], where the method is trained on a class of environments and then
 87 performs novel view synthesis given a few context images.

88 **Rasterization methods**, which project 3D geometry onto 2D image planes [19, 13], can be fast
 89 thanks to hardware optimization (rasterization pipelines have been dominant in the past). Recently,
 90 Gaussian Splatting (3DGS) [19] has gained popularity by matching the rendering quality of NeRFs
 91 while rendering at real-time speeds. 3DGS represents scenes using 3D Gaussians, which are projected
 92 onto the image plane [72], rasterized, and alpha-composited to render a view.

93 Like NeRFs, 3DGS has been modified to handle aliasing is-
 94 sues [67, 31], extended to larger scenes [20], and optimized to
 95 reduce the size of the representation [40] and to improve its ge-
 96 ometric accuracy [18]. A minor change to the adaptive density
 97 control (ACD) was proposed in [65, 68] – where the sum of
 98 gradients was replaced with a sum of its norm, and a faster ini-
 99 tialization method was proposed in [12]. Additionally, 3DGS
 100 has been applied to various domains, including SLAM [34] and
 101 physics simulation [59]. It has been extended to handle dynamic
 102 scenes [32] and semantics [71] using an N-D rasterizer. Genera-
 103 tive 3DGS has also been proposed [10].

104 Unlike ray-based approaches which sample pixels randomly from
 105 training views, 3DGS optimizes on full images, requiring differ-
 106 ent training batch formation. In NerfBaselines, we do not enforce
 107 any specific way of constructing training batches, allowing each
 108 method to use the format it needs.

109 2.1 Existing codebases

110 Most current methods are based on a few core repositories: Nerf-
 111 Studio [53], Multi-NeRF [2, 3], Instant-NGP [39] for NeRFs,
 112 and Gaussian Splatting (INRIA) [19] for 3DGS, typically with
 113 moderate modifications. Therefore, in NerfBaselines, we focused
 114 on integrating these core repositories, as it simplifies the subse-
 115 quent integration of derivative works. In Figure 1, we illustrate the relationships between popular
 116 repositories and highlight those currently integrated with NerfBaselines.

117 **NerfStudio** is a popular framework that introduced the modular separation of NeRFs into components
 118 such as ray samplers, radiance field heads, etc. It supports various dataloaders, camera types, and
 119 export formats. Unfortunately, the rapid evolution of NerfStudio introduces frequent breaking changes.
 120 Therefore, in NerfBaselines, we freeze each method to ensure reproducibility.

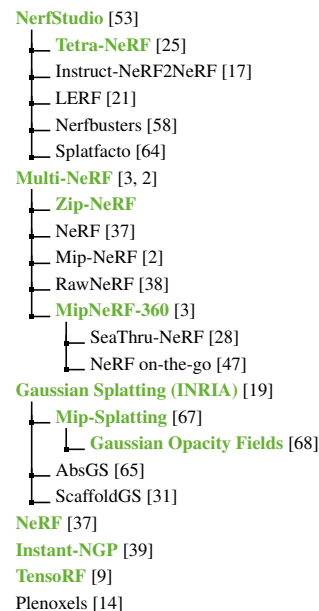


Figure 1: **Existing codebases.** Integrated methods are **bold green**.

121 **Multi-NeRF** is fully implemented in JAX and does not have custom CUDA kernels (unlike Instant-
122 NeRF or NerfStudio), making it easy to install, but slower. Early versions based on Mip-NeRF360
123 supported only a single camera per dataset and a single image size; therefore, in NerfBaselines, we
124 have extended these methods to handle more complex datasets.

125 **Instant-NGP** is a highly optimized implementation that has inspired numerous follow-up works
126 [4, 53, 21, 22]. However, it is a less popular choice as a codebase due to its C++ training code
127 being more difficult to extend. Since the repository does not natively support rendering with custom
128 (distorted) cameras, in NerfBaselines, we have implemented this functionality.

129 **Gaussian Splatting (INRIA)** is the most popular choice for 3DGS methods, primarily because (until
130 recently) it was unmatched in performance. The repository only supports pinhole cameras with their
131 centre of projection being in center of the image. We have extended it to work with arbitrary camera
132 models, performing undistortion/distortion for more complicated camera models.

133 3 Framework design

134 When designing the NerfBaselines framework, we aimed to address common issues when bench-
135 marking novel view synthesis methods. These include: **different interfaces**, where each codebase
136 has a different structure, making it difficult to interface with it. To this end, we propose a unified
137 interface that all methods can share. **Installation challenges** arise as various methods often require
138 specific and potentially difficult-to-install dependencies. To simplify the process, we install each
139 method in a separate environment where we fully control the dependencies. For **fair evaluation**, we
140 implement a standardized evaluation protocol for all methods, closely matching the original protocols
141 proposed with the datasets. **Reproducibility** is ensured by carefully storing and tracking checkpoints
142 used to generate predictions and compute metrics. Additionally, we have implemented a viewer with
143 a trajectory editor and a public website to compare existing methods on standard benchmark datasets.

144 **Unified method API.** The different codebases [53, 3, 19, 9, 19] have very different code structures
145 making it difficult to interface with them in a unified way. While studying the codebases, we
146 identified a common structure which can be shared between all existing methods, regardless if they
147 are raycasting-based (as in the case of NeRF codebases [53, 3, 19, 9, 60]), or rasterization-based (as
148 in 3DGS methods [19, 20, 67, 68, 64]). For each codebase, we define a class called `method`, which
149 encapsulates all that is needed to train the method and render new images from an existing checkpoint.
150 Each method implements (among others) methods: `train_iteration` - which performs one training
151 step, and `render` - which renders images from input cameras. Additionally, if the method supports
152 appearance conditioning [33, 53], it implements the functionality to obtain the embeddings from
153 unseen images. The detailed interface is given in the *supp. mat.*

154 **Environment isolation.** Installing methods is often difficult due to a set of hard-to-install and
155 sometimes obsolete dependencies, e.g., CUDA, CuDNN, GCC, or OpenGL, where the version
156 must match. The dependencies are sometimes in conflict with each other or in conflict with the
157 user’s environment. Therefore, we provide a level of isolation, where each method is installed in
158 its own isolated environment, allowing for different/conflicting dependencies without polluting the
159 user’s environment. The user then communicates with the isolated environments using interprocess
160 communication. NerfBaselines then handles the installation of all necessary dependencies the first
161 time a user uses a method. We provide three levels of isolation: Anaconda [1], Docker [35], and
162 Apptainer [27]: **1)** While Anaconda [1] is commonly used in research and is easy to install even on
163 HPC clusters, it does not offer a good level of isolation, as some dependencies are inherited from the
164 system, e.g., glibc, OpenGL, and, therefore, may fail on some systems. **2)** Docker [35] provides the
165 best level of isolation, fully encapsulating the user environment (except for the CUDA drivers needed
166 to control GPUs). Unfortunately, it is more difficult to install and is not commonly present on HPC
167 clusters. Therefore, we also support the alternative **3)** Apptainer [27] with a better level of isolation
168 than Anaconda that is more commonly present on HPC clusters.

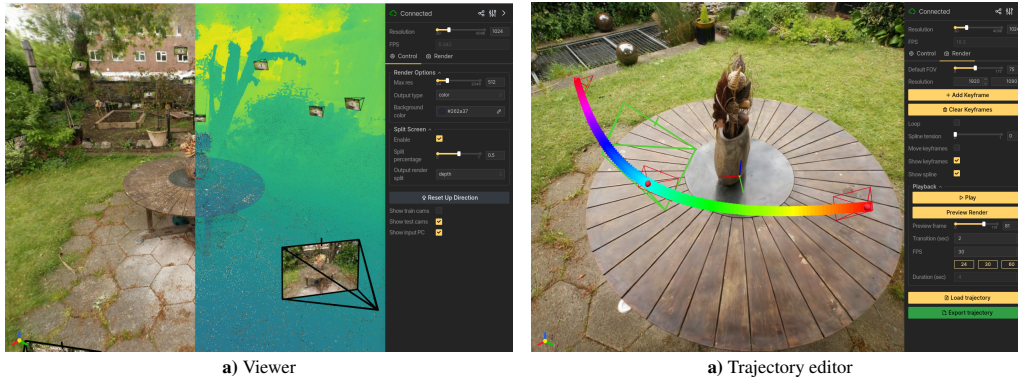


Figure 2: The **NerfBaselines Viewer** enables interactive rendering, shows train/test cameras, and input point cloud. It also has a trajectory editor (shown in the left figure).

169 **Unified dataloader.** Different codebases store and load data in various ways, complicating the
 170 process of integrating and evaluating methods across different datasets. For example, they use
 171 different downscaling algorithms and resolutions, different background colors, etc. The camera poses
 172 processing pipeline is also important, where methods like Instant-NGP [39] and NerfStudio [53] are
 173 sensitive to the scale of the scene, which depends on the camera placement relative to the scene’s
 174 geometry. We propose a unified interface for dataset loading and processing to address these issues
 175 and to facilitate the transfer to new datasets. We implement support for commonly used dataset
 176 formats such as COLMAP [49, 48], NeRF’s transforms.json [37], NerfStudio’s transforms.json
 177 [53], Bundler [51, 50] (also supported by RealityCapture [7]), NeRF in the Wild Photo Tourism
 178 [33, 50] splits format, Tanks and Temples [23], and the LLFF dataset format [36]. The loaders ensure
 179 consistent data processing and provide a uniform interface and camera format for all datasets, thereby
 180 simplifying integration and evaluation across diverse datasets. It also makes the development of new
 181 methods easier as they do not need to implement their own dataloaders anymore. Furthermore, to
 182 simplify the ease of use, we support an automatic download of commonly used benchmark datasets
 183 such as the Blender dataset [37] and the Mip-NeRF 360 dataset [3].

184 **Viewer & camera trajectory rendering.** For novel view synthesis methods, evaluation is often
 185 conducted on a set of test images that are typically a subsequence of the training views. Unfortunately,
 186 this approach does not adequately demonstrate the method’s robustness and multiview consistency in
 187 capturing the 3D scene [29]. Rendering a video from a custom trajectory from truly novel viewpoints
 188 farther away from the training images provides a more insightful evaluation. While the unified
 189 interface presented in Section 3 ‘Unified method API’ significantly simplifies the process of rendering
 190 such trajectories, we further enhance this capability by providing an interactive viewer. This viewer
 191 allows users to inspect how the method performs outside the training camera distribution and includes
 192 a camera trajectory editor for designing custom trajectories. The viewer, based on NerfStudio’s
 193 viewer [53] and utilizing the Viser platform [53] for the web interface, is depicted in Figure 2, where
 194 the viewer interface is shown on the left and the trajectory editor on the right. The trajectory editor
 195 also enables exporting the trajectory in a format suitable for subsequent rendering.

196 **Fair evaluation & Reproducibility** Evaluating and comparing different methods is challenging
 197 because original codebases often use different evaluation protocols. For example, some methods use
 198 different parameters for SSIM [6] and different LPIPS [70] architectures (like AlexNet vs. VGG).
 199 Some methods even use different image resolutions [60, 62], or compute metrics on the raw image
 200 (float range) as opposed to rounding to the uint8 range – which cannot be reproduced from the
 201 commonly saved images. Therefore, we standardize the evaluation protocol by fixing the parameters
 202 of the metrics, resolution of images, and other aspects of the evaluation.

203 Sharing checkpoints often greatly improve research reproducibility, enabling validation of predicted
 204 test images and rendering 3D scenes from novel viewpoints. To enable this, NerfBaselines stores

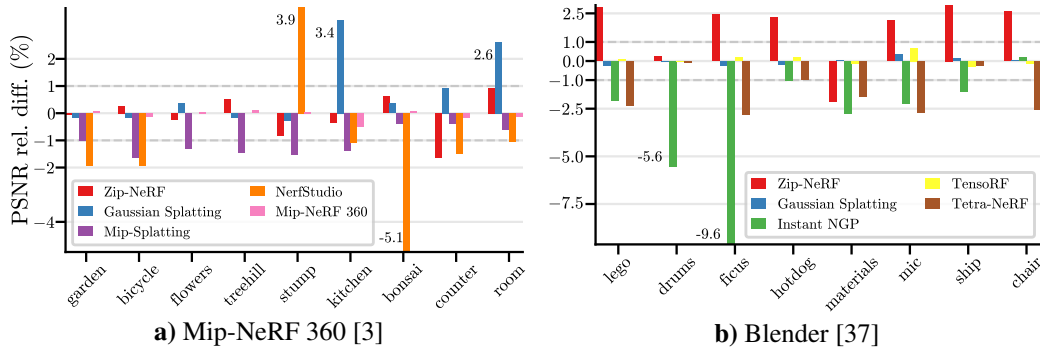


Figure 3: **Mip-NeRF 360 [3] and Blender [37] results** comparing PSNRs obtained via NerfBaselines with those reported in the original papers. We show the relative difference in %. In most cases, the difference is $< 1\%$. Instant-NGP [39] and Mip Splatting [67] are consistently underperforming because different evaluation protocols were used in the papers.

205 checkpoints to ensure models can be restored later. During method evaluation, we store the checkpoint
 206 SHA and verify it matches the released checkpoint. For all methods, we fix random seeds to match
 207 the official implementations. Additionally, for each integrated method, we provide a set of tests to
 208 confirm that loading the checkpoint genuinely reproduces the results.

209 **Web platform & Extensibility.** As part of the NerfBaselines toolbox, we release a website where
 210 methods are compared and their numbers (measured with NerfBaselines) reported. It also enables
 211 the checkpoints and the predictions to be downloaded for all scenes. The web can be found at
 212 <https://jkulhanek.com/nerfbaselines>.

213 Finally, NerfBaselines is designed in such that it makes it easy to extend it with new methods and
 214 datasets. Therefore we provide extension points to register new datasets, new evaluation protocols,
 215 and add new loggers (currently, we implement wandb and tensorboard loggers). The new method can
 216 be added by implementing the protocol described in Section 3 ‘Unified method API’ and providing a
 217 specification file, containing the installation script and method’s metadata.

218 4 Evaluation

219 In our experiments, we **1)** verify that the integrated methods match the original metrics, **2)** show
 220 the tradeoff between quality and training cost, **3)** motivate the need for standardized evaluation
 221 protocols by demonstrating that small changes in evaluation protocol can bias results, **4)** demonstrate
 222 transferability to novel datasets, and **5)** qualitatively compare methods outside training trajectories.
 223 All our experiments used NVIDIA A100 GPUs. A single GPU was used for all but Mip-NeRF 360 [3],
 224 which used four GPUs. For comparisons, we use PSNR, SSIM [6], and LPIPS [70] (AlexNet). In the
 225 main paper, we mostly report PSNR, while other results can be found in the *supp. mat*.

226 4.1 Reproducing published results

227 First, to validate our framework, we reevaluate important methods on the standard benchmark datasets:
 228 Mip-NeRF 360 [3] and Blender [37]. We use the same evaluation protocol for all methods. The
 229 results are compared to the original numbers as published in the papers in Figure 3, with detailed
 230 numbers given in the *supp. mat*. Note, that in the figures and tables, we only compare with methods
 231 that released their numbers on the datasets in the corresponding publications.

232 **Mip-NeRF 360 results.** As shown in Figure 3, NerfBaselines reproduces the original results with a
 233 deviation of less than 1% for most scenes. For Mip-Splatting [67] and 3DGS [19], the difference in
 234 the numbers was caused by the different evaluation protocols used as discussed in Section 4.3. In
 235 the case of NerfStudio [53] and Tetra-NeRF [25], the codebase evolved since the time of the release
 236 which is likely the cause of the difference.

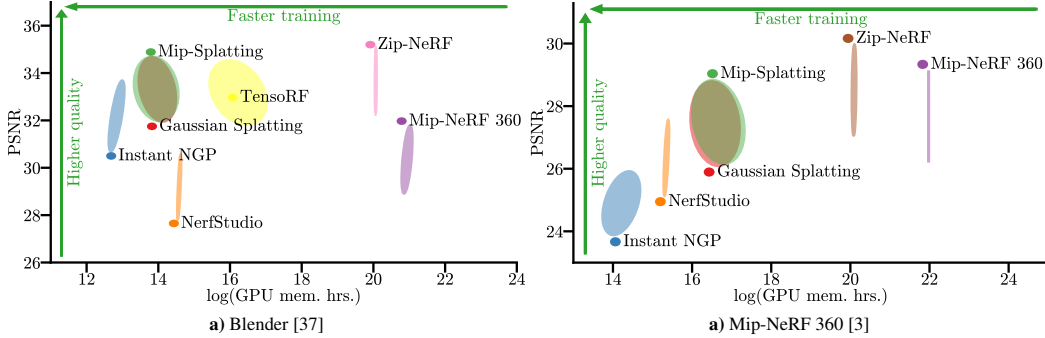


Figure 4: **Performance vs. training cost tradeoff.** We compare different methods’ rendering quality (PSNR \uparrow) and rendering cost (GPU mem. hrs. \downarrow). The variance across different scenes is visualized by the ellipse scale. While Zip-NeRF [4] performs the best, it is more costly to train than 3DGS[19].

	garden kitchen	bicycle bonsai	flowers counter	treehill room	stump
NerfStudio [53]	26.21/25.96/+0.96% 29.87/29.96/-0.28%	24.06/23.61/+1.90% 30.51/30.52/-0.02%	21.43/21.12/+1.45% 26.91/26.80/+0.43%	23.23/22.85/+1.66% 30.79/30.56/+0.76%	26.09/25.75/+1.33%
Zip-NeRF [4]	28.58/28.18/+1.40% 32.57/32.39/+0.55%	26.35/25.87/+1.85% 34.89/34.67/+0.64%	22.69/22.34/+1.56% 29.14/28.90/+0.82%	24.43/24.01/+1.73% 33.22/32.95/+0.80%	27.70/27.32/+1.37%
Gaussian Splatting [19]	27.75/27.37/+1.40% 31.65/31.36/+0.92%	25.62/25.20/+1.65% 32.29/32.10/+0.60%	21.84/21.60/+1.14% 28.97/28.97/+0.03%	22.75/22.46/+1.30% 31.38/31.43/-0.17%	26.82/26.48/+1.29%
Mip-Splatting [4]	27.85/27.48/+1.34% 31.55/31.12/+1.38%	25.68/25.30/+1.51% 32.39/32.18/+0.65%	21.93/21.64/+1.35% 29.07/29.04/+0.10%	22.99/22.64/+1.52% 31.68/31.55/+0.42%	26.88/26.52/+1.33%

Table 1: **Mip-NeRF 360 [3] image downscale protocol.** We compare the manual downscaling with the default protocol using released downsampled images (*middle number*). While NeRFs use released downsampled images, 3DGS-based methods [19] downscale images internally without compression. *Indoor scenes* and *outdoor scenes* have downscale factors of 2 and 4, respectively.

237 **Blender results.** From the results, the discrepancy is again small for most methods. However, for
 238 the Instant-NGP method [39], we can notice larger differences in PSNR, especially for ‘drums’,
 239 and ‘ficus’. Note, that Instant-NGP [39] uses a black background for training and evaluation. We
 240 run additional experiments confirming this to be the case of the difference in the metrics. Since
 241 Tetra-NeRF [25] uses NerfStudio [53] and the codebase evolved since the time of the release, we can
 242 notice a slight drop in the performance.

243 4.2 Quality vs. computational cost

244 When comparing methods, image quality is important, but we also need to consider the computational
 245 cost of training. Some methods perform better but require more computational power, while others
 246 achieve good results with less computation. To compare and visualize this, we plot the performance
 247 of different methods against their computational resources in Figure 4. We measure computational
 248 cost as the training time multiplied by GPU memory use. The reasoning is that most methods can be
 249 sped up with more GPU memory or larger batch sizes, while low GPU memory usage allows training
 250 on cheaper GPUs. We also show the variance in performance and computational cost. For both
 251 the Mip-NeRF 360 [3] and Blender [37] datasets, performance variance is similar across methods.
 252 However, due to its adaptive nature, computational cost variance is much higher for methods based
 253 on Gaussian Splatting [19] compared to NeRF-based methods [37, 3, 4].

254 4.3 Mip-NeRF 360 evaluation protocol details

255 As stressed in the introduction, ensuring the same evaluation protocol is crucial in benchmarking. In
 256 this section, we demonstrate this by comparing two evaluation protocols that were used for the Mip-
 257 NeRF 360 dataset [3]. In Section 4.1, we have seen that NerfBaselines achieves consistently lower
 258 PSNR for Mip Splatting [67]. The reason is that while NeRFs train/evaluate on released downsampled

	barr	caterpillar	truck	highhouse	playground	train	auditorium	ballroom	courtroom	museum	palace	temple
PSNR↑	Training Data			Intermediate			Advanced					
Instant NGP [39]	25.90	21.72	22.85	21.65	23.33	20.01	20.67	21.62	19.44	15.19	19.09	17.84
NerfStudio [53]	26.40	21.71	23.37	20.85	24.69	20.43	20.77	22.68	20.24	17.84	17.68	17.06
Zip-NeRF [4]	29.26	23.94	25.09	23.07	27.13	22.19	24.52	25.45	22.17	19.34	19.11	20.58
Gaussian Splatting [19]	27.51	23.38	24.25	22.11	25.37	21.67	24.13	24.07	23.12	20.92	19.63	20.85
Mip-Splatting [67]	27.75	23.42	24.36	22.25	25.87	21.82	24.41	24.15	23.00	20.88	19.63	20.55
Gaussian Opacity Fields [68]	25.72	21.78	22.33	21.80	23.89	19.69	23.20	22.84	21.15	19.92	16.46	20.29
SSIM↑												
Instant NGP [39]	0.772	0.633	0.770	0.765	0.696	0.657	0.761	0.652	0.640	0.471	0.668	0.689
NerfStudio [53]	0.794	0.666	0.797	0.768	0.755	0.693	0.771	0.705	0.673	0.648	0.640	0.678
Zip-NeRF [4]	0.884	0.802	0.864	0.849	0.880	0.814	0.877	0.835	0.790	0.746	0.718	0.805
Gaussian Splatting [19]	0.852	0.791	0.853	0.843	0.848	0.791	0.871	0.824	0.790	0.764	0.736	0.806
Mip-Splatting [67]	0.855	0.790	0.857	0.844	0.861	0.795	0.872	0.826	0.791	0.768	0.731	0.805
Gaussian Opacity Fields [68]	0.866	0.791	0.860	0.833	0.869	0.796	0.871	0.818	0.781	0.761	0.683	0.794
LPIPS↑												
Instant NGP [39]	0.271	0.360	0.216	0.281	0.343	0.334	0.429	0.352	0.448	0.606	0.440	0.424
NerfStudio [53]	0.215	0.302	0.167	0.245	0.249	0.261	0.330	0.261	0.336	0.311	0.452	0.392
Zip-NeRF [4]	0.083	0.152	0.081	0.131	0.095	0.119	0.153	0.113	0.153	0.159	0.317	0.183
Gaussian Splatting [19]	0.160	0.190	0.108	0.156	0.170	0.171	0.193	0.101	0.165	0.160	0.350	0.222
Mip-Splatting [67]	0.161	0.197	0.109	0.159	0.155	0.172	0.196	0.098	0.165	0.158	0.354	0.226
Gaussian Opacity Fields [68]	0.140	0.187	0.099	0.181	0.142	0.164	0.194	0.107	0.168	0.152	0.443	0.234

Table 2: **Tanks & Temples [23] results.** We show the PSNR, SSIM [6], and LPIPS [70] (AlexNet) of various implemented methods. The **first**, **second**, and **third** values are highlighted.

259 images, the 3DGS [19] and Mip Splatting codebases [67] downscales from the large original images
260 during training and evaluation (without storing them as JPEGs, thus avoiding compression artifacts).
261 From the results presented in Table 1, we can see that while the difference is not large and the
262 relative ordering is preserved, not using compression consistently gives results with larger PSNR.
263 The difference is larger for outdoor scenes, where the downscale factor of 4 was used as opposed to
264 indoor scenes with a downscale factor of 2.

265 4.4 Tanks & Temples evaluation

266 To demonstrate how NerfBaseline simplifies the transfer of existing methods to new datasets, we
267 evaluate various integrated methods on the Tanks and Temples [23] dataset. For the dataset, we
268 run COLMAP reconstruction [48, 49] with a simple radial camera model shared for all images.
269 Afterward, we undistorted and downscaled images by a factor of 2. For NerfStudio [53], we run
270 the Mip-NeRF 360 configuration (which from our experiments performs better on the dataset). The
271 results are given in Table 2. As we can see the reconstructions are dominated by Zip-NeRF [4] for
272 easier scenes, while for ‘Advanced’, there is no single best-performing method. We believe this is
273 caused by NeRFs with its fixed capacity does not scale as well to larger scenes as 3DGS, where the
274 capacity is adaptively increased.

275 4.5 Off-trajectory qualitative comparison

276 While test-set metrics enable an effective way of comparing different methods, they are insufficient
277 to fully evaluate the perceived quality [29]. Rendering images from poses with varying distances
278 from the training camera’s trajectory provides a lot of insight into the robustness of the learned
279 representations. Therefore, NerfBaselines provides a viewer and a renderer to enable visualising
280 methods and rendering images/videos outside the train trajectory. In Figure 5, we compare various
281 methods by rendering trained scenes both close to the training camera trajectory and far from it.
282 Notice how in the second row Instant-NGP [39], and 3DGS methods [19, 67] cannot fill the sparsely
283 observed sky, while NerfStudio [53] and Zip-NeRF [4] can achieve it thanks to space contraction.
284 Also, notice how 3DGS methods [19, 67] are more blurred in less observed regions.

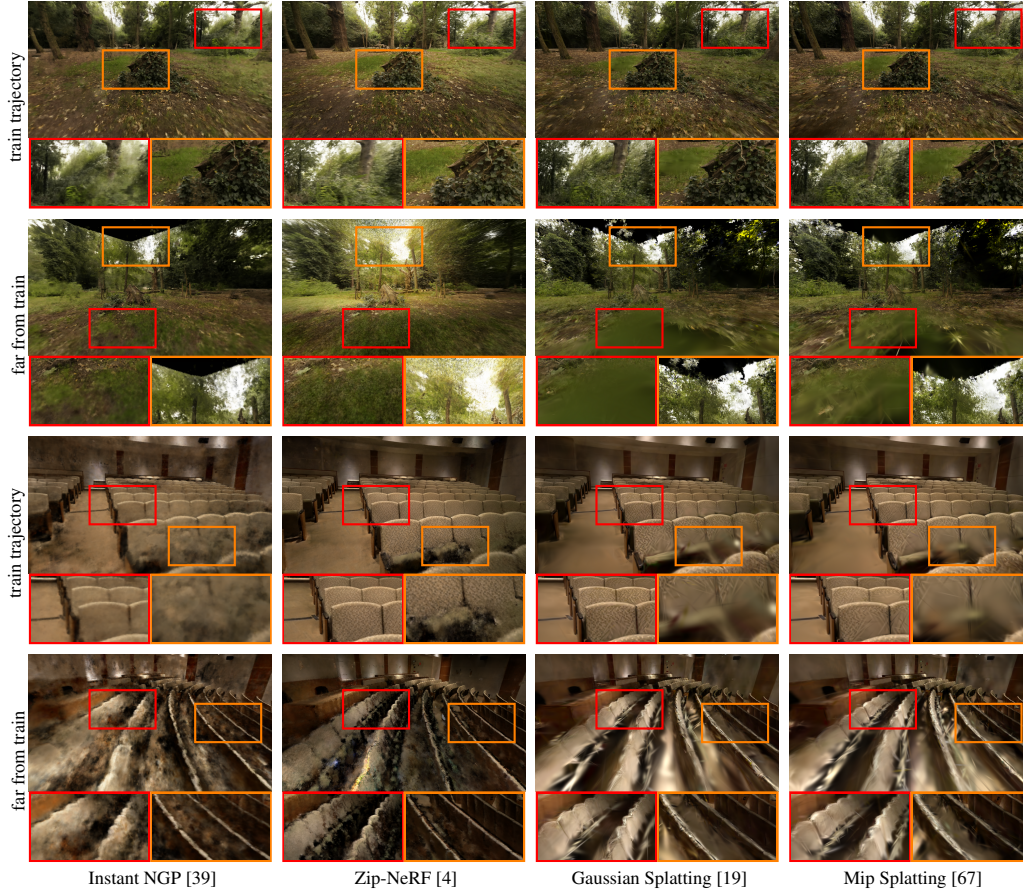


Figure 5: **Qualitative results.** We compare methods on views close and far from the training trajectory. **Top:** MipNeRF360/stump scene, **bottom:** T&T/Auditorium.

285 5 Conclusion

286 In conclusion, our NerfBaselines framework addresses the major challenges in evaluating novel view
 287 synthesis methods, e.g. NeRFs and 3DGS. By standardizing evaluation protocols and designing a
 288 unified interface, we enable fair comparisons and scalability to novel datasets. The camera trajectory
 289 editor enables multiview consistency evaluation, while the NerfBaselines framework ensures smooth
 290 installation and reproducibility by using isolated environments. Additionally, our web platform
 291 displays benchmark results, comparing various methods across different datasets. The NerfBaselines
 292 framework thus improves the fairness and effectiveness of novel view synthesis method evaluations.

293 **License.** This project is released under the **MIT** license. The integrated methods may be licensed
 294 under various licenses, and it is the user’s responsibility to conform to the conditions before using
 295 a method. Note that while NerfBaselines enables access to various methods, the user still uses the
 296 original official codebases (with some wrapper code to interface with it).

297 **Limitations.** While NerfBaselines offers consistent and reproducible benchmarking, it requires
 298 methods to expose the same interface either directly or by writing a wrapper script. While we
 299 integrated some well-known methods and will gradually add more, our hope is that the scientific
 300 community could collaborate and adopt the interface for future methods.

301 **Broader impact.** While training and evaluating all methods require significant computational
 302 resources and associated carbon emissions, our release of reproducible checkpoints minimizes the
 303 need for repeated training by other researchers, thereby reducing overall environmental impact.

304 Acknowledgments and Disclosure of Funding

305 We want to thank Brent Yi and the NerfStudio Team for helpful discussions regarding the NerfStudio
306 codebase and for releasing the Viser platform. This work was supported by the Czech Science
307 Foundation (GAČR) EXPRO (grant no. 23-07973X), the Grant Agency of the Czech Technical
308 University in Prague (grant no. SGS24/095/OHK3/2T/13), and by the Ministry of Education, Youth
309 and Sports of the Czech Republic through the e-INFRA CZ (ID:90254)

310 References

- 311 [1] Anaconda. Anaconda software distribution, 2020. URL <https://docs.anaconda.com/>.
- 312 [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P
313 Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, pages
314 5855–5864, 2021.
- 315 [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360:
316 Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022.
- 317 [4] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Zip-nerf:
318 Anti-aliased grid-based neural radiance fields. In *ICCV*, pages 19697–19705, 2023.
- 319 [5] Wenjing Bian, Zirui Wang, Kejie Li, Jia-Wang Bian, and Victor Adrian Prisacariu. Nope-nerf: Optimising
320 neural radiance field with no pose prior. In *CVPR*, pages 4160–4169, 2023.
- 321 [6] Dominique Brunet, Edward R. Vrscay, and Zhou Wang. On the mathematical properties of the structural
322 similarity index. *IEEE TIP*, 2012.
- 323 [7] Capturing Reality. Reality capture 1.4, 2024. URL <https://www.capturingreality.com/>. Software.
- 324 [8] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf:
325 Fast generalizable radiance field reconstruction from multi-view stereo. In *ICCV*, pages 14124–14133,
326 2021.
- 327 [9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensorRF: Tensorial radiance fields. In
328 *ECCV*, 2022.
- 329 [10] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen
330 Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. *arXiv*,
331 2024.
- 332 [11] Congyue Deng, Chiyu Jiang, Charles R Qi, Xinchen Yan, Yin Zhou, Leonidas Guibas, Dragomir Anguelov,
333 et al. Nerd: Single-view nerf synthesis with language-guided diffusion as general image priors. In *CVPR*,
334 pages 20637–20647, 2023.
- 335 [12] Zhiwen Fan, Wenyan Cong, Kairun Wen, Kevin Wang, Jian Zhang, Xinghao Ding, Danfei Xu, Boris
336 Ivanovic, Marco Pavone, Georgios Pavlakos, Zhangyang Wang, and Yue Wang. InstantSplat: Unbounded
337 sparse-view pose-free gaussian splatting in 40 seconds. *arXiv*, 2024.
- 338 [13] Linus Franke, Darius Rückert, Laura Fink, and Marc Stamminger. Trips: Trilinear point splatting for
339 real-time radiance field rendering. *CGF*, 43(2), 2024. doi: <https://doi.org/10.1111/cgf.15012>.
- 340 [14] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa.
341 Plenoxels: Radiance fields without neural networks. In *CVPR*, pages 5501–5510, 2022.
- 342 [15] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa.
343 K-planes: Explicit radiance fields in space, time, and appearance. In *CVPR*, 2023.
- 344 [16] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf:
345 High-fidelity neural rendering at 200fps. In *ICCV*, pages 14346–14355, 2021.
- 346 [17] Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-
347 nerf2nerf: Editing 3d scenes with instructions. In *ICCV*, 2023.
- 348 [18] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2D gaussian splatting for
349 geometrically accurate radiance fields. *arXiv*, 2024.

- 350 [19] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for
351 real-time radiance field rendering. *ACM TOG*, 2023.
- 352 [20] Bernhard Kerbl, Andreas Meuleman, Georgios Kopanas, Michael Wimmer, Alexandre Lanvin, and George
353 Drettakis. A hierarchical 3d gaussian representation for real-time rendering of very large datasets. *ACM*
354 *TOG*, 43(4), July 2024.
- 355 [21] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lurf: Language
356 embedded radiance fields. In *ICCV*, pages 19729–19739, 2023.
- 357 [22] Chung Min Kim, Mingxuan Wu, Justin Kerr, Ken Goldberg, Matthew Tancik, and Angjoo Kanazawa.
358 Garfield: Group anything with radiance fields. *arXiv*, 2024.
- 359 [23] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking
360 large-scale scene reconstruction. *ACM TOG*, 36(4), 2017.
- 361 [24] Jonas Kulhanek and Torsten Sattler. Tetra-nerf: Representing neural radiance fields using tetrahedra. In
362 *ICCV*, 2023.
- 363 [25] Jonas Kulhanek and Torsten Sattler. Tetra-nerf: Representing neural radiance fields using tetrahedra. In
364 *ICCV*, pages 18458–18469, 2023.
- 365 [26] Jonáš Kulháněk, Erik Derner, Torsten Sattler, and Robert Babuška. Viewformer: Nerf-free neural rendering
366 from few images using transformers. In *ECCV*, pages 198–216. Springer, 2022.
- 367 [27] Gregory M. Kurtzer, cclerget, Michael Bauer, Ian Kaneshiro, David Trudgian, and David Godlove.
368 hpcng/singularity: Singularity 3.7.3, 2021.
- 369 [28] Deborah Levy, Amit Peleg, Naama Pearl, Dan Rosenbaum, Derya Akkaynak, Simon Korman, and Tali
370 Treibitz. Seathru-nerf: Neural radiance fields in scattering media. In *CVPR*, pages 56–65, 2023.
- 371 [29] Hanxue Liang, Tianhao Wu, Param Hanji, Francesco Banterle, Hongyun Gao, Rafal Mantiuk, and Cengiz
372 Oztireli. Perceptual quality assessment of nerf and neural view synthesis methods for front-facing views.
373 In *Eurographics*, 2023.
- 374 [30] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh.
375 Neural volumes: learning dynamic renderable volumes from images. *ACM TOG*, 38(4):1–14, 2019.
- 376 [31] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-GS:
377 Structured 3D gaussians for view-adaptive rendering. *arXiv*, 2023.
- 378 [32] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking
379 by persistent dynamic view synthesis. In *3DV*, 2024.
- 380 [33] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and
381 Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*,
382 2021.
- 383 [34] Hidenobu Matsuki, Riku Murai, Paul H. J. Kelly, and Andrew J. Davison. Gaussian Splatting SLAM.
384 *CVPR*, 2024.
- 385 [35] Dirk Merkel. Docker: lightweight linux containers for consistent development and deployment. *Linux*
386 *journal*, 2014(239):2, 2014.
- 387 [36] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi,
388 Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling
389 guidelines. *ACM TOG*, 38(4), 2019.
- 390 [37] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng.
391 NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- 392 [38] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P Srinivasan, and Jonathan T Barron. Nerf
393 in the dark: High dynamic range view synthesis from noisy raw images. In *CVPR*, pages 16190–16199,
394 2022.
- 395 [39] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives
396 with a multiresolution hash encoding. *ACM TOG*, 2022.
- 397 [40] Panagiotis Papantonakis, Georgios Kopanas, Bernhard Kerbl, Alexandre Lanvin, and George Drettakis.
398 Reducing the memory footprint of 3d gaussian splatting. *ACM CGIT*, 2024.

- 399 [41] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and
400 Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, pages 5865–5874, 2021.
- 401 [42] Keunhong Park, Philipp Henzler, Ben Mildenhall, Jonathan T Barron, and Ricardo Martin-Brualla. Camp:
402 Camera preconditioning for neural radiance fields. *ACM TOG*, 42(6):1–11, 2023.
- 403 [43] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion.
404 In *ICLR*, 2022.
- 405 [44] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural
406 Radiance Fields for Dynamic Scenes. In *CVPR*, 2020.
- 407 [45] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance
408 fields with thousands of tiny mlps. In *ICCV*, pages 14335–14345, 2021.
- 409 [46] Christian Reiser, Rick Szeliski, Dor Verbin, Pratul Srinivasan, Ben Mildenhall, Andreas Geiger, Jon Barron,
410 and Peter Hedman. Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded
411 scenes. *ACM TOG*, 42(4):1–12, 2023.
- 412 [47] Weining Ren, Zihan Zhu, Boyang Sun, Jiaqi Chen, Marc Pollefeys, and Songyou Peng. Nerf on-the-go:
413 Exploiting uncertainty for distractor-free nerfs in the wild. In *CVPR*, 2024.
- 414 [48] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016.
- 415 [49] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view
416 selection for unstructured multi-view stereo. In *ECCV*, 2016.
- 417 [50] Noah Snavely, Steven M Seitz, and Richard Szeliski. Photo tourism: exploring photo collections in 3d.
418 *ACM TOG*, 2006.
- 419 [51] Noah Snavely, Steven M Seitz, and Richard Szeliski. Modeling the world from internet photo collections.
420 *IJCV*, 80:189–210, 2008.
- 421 [52] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for
422 radiance fields reconstruction. In *CVPR*, 2022.
- 423 [53] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Terrance Wang, Alexander Kristoffersen,
424 Jake Austin, Kamyar Salahi, Abhik Ahuja, et al. Nerfstudio: A modular framework for neural radiance
425 field development. In *ACM TOG*, 2023.
- 426 [54] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Nerf-art:
427 Text-driven neural radiance fields stylization. *IEEE TVCG*, 2023.
- 428 [55] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning
429 neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 34:27171–27183,
430 2021.
- 431 [56] Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and
432 Wenping Wang. F2-nerf: Fast neural radiance field training with free camera trajectories. In *CVPR*, pages
433 4150–4159, 2023.
- 434 [57] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2:
435 Fast learning of neural implicit surfaces for multi-view reconstruction. In *ICCV*, pages 3295–3306, 2023.
- 436 [58] Frederik Warburg, Ethan Weber, Matthew Tancik, Aleksander Holynski, and Angjoo Kanazawa. Nerf-
437 busters: Removing ghostly artifacts from casually captured nerfs. In *ICCV*, pages 18120–18130, 2023.
- 438 [59] Tianyi Xie, Zeshun Zong, Yuxin Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. Physgaussian:
439 Physics-integrated 3d gaussians for generative dynamics. *arXiv*, 2023.
- 440 [60] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann.
441 Point-nerf: Point-based neural radiance fields. In *CVPR*, pages 5438–5448, 2022.
- 442 [61] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann.
443 Point-nerf: Point-based neural radiance fields. In *CVPR*, pages 5438–5448, 2022.
- 444 [62] Yifan Yang, Shuhai Zhang, Zixiong Huang, Yubing Zhang, and Mingkui Tan. Cross-ray neural radiance
445 fields for novel-view synthesis from unconstrained image collections. In *ICCV*, 2023.

- 446 [63] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces.
447 *NeurIPS*, 34:4805–4815, 2021.
- 448 [64] Vickie Ye, Matias Turkulainen, and the Nerfstudio team. gsplat, 2024. URL [https://github.com/
449 nerfstudio-project/gsplat](https://github.com/nerfstudio-project/gsplat).
- 450 [65] Zongxin Ye, Wenyu Li, Sidun Liu, Peng Qiao, and Yong Dou. Absgs: Recovering fine details for 3d
451 gaussian splatting. *arXiv*, 2024.
- 452 [66] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or
453 few images. In *CVPR*, pages 4578–4587, 2021.
- 454 [67] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d
455 gaussian splatting. In *CVPR*, 2024.
- 456 [68] Zehao Yu, Torsten Sattler, and Andreas Geiger. Gaussian opacity fields: Efficient high-quality compact
457 surface reconstruction in unbounded scenes. *arXiv*, 2024.
- 458 [69] Kai Zhang, Gernot Riegler, Noah Snively, and Vladlen Koltun. Nerf++: Analyzing and improving neural
459 radiance fields. *arXiv*, 2020.
- 460 [70] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
461 effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- 462 [71] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suya
463 You, Zhangyang Wang, and Achuta Kadambi. Feature 3dgs: Supercharging 3d gaussian splatting to enable
464 distilled feature fields. *CVPR*, 2024.
- 465 [72] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Surface splatting. In *ACM TOG*,
466 2001.

467 **Checklist**

- 468 1. For all authors...
- 469 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
470 contributions and scope? [Yes] Claims reflect paper’s contribution and scope
- 471 (b) Did you describe the limitations of your work? [Yes] Limitations described in Section 5
- 472 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Broader
473 impact described in Section 5
- 474 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
475 them? [Yes] The paper conforms to the ethics review guidelines
- 476 2. If you are including theoretical results...
- 477 (a) Did you state the full set of assumptions of all theoretical results? [N/A] No theoretical
478 results
- 479 (b) Did you include complete proofs of all theoretical results? [N/A] No theoretical results
- 480 3. If you ran experiments (e.g. for benchmarks)...
- 481 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
482 mental results (either in the supplemental material or as a URL)? [Yes] Code included
483 including the instructions needed to reproduce the results.
- 484 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
485 were chosen)? [Yes] Details specified in the experimental section.
- 486 (c) Did you report error bars (e.g., with respect to the random seed after running ex-
487 periments multiple times)? [No] Not possible due to computation cost of a single
488 run
- 489 (d) Did you include the total amount of compute and the type of resources used (e.g., type
490 of GPUs, internal cluster, or cloud provider)? [Yes] Amount of compute and type of
491 resourced explicitly mentioned in the main paper.
- 492 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 493 (a) If your work uses existing assets, did you cite the creators? [Yes] All used assets are
494 properly cited.
- 495 (b) Did you mention the license of the assets? [Yes] License mentioned in Section 5.
496 Further details are included in *supp. mat.*
- 497 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
498 We include the source code for the framework and URL to the webpage.
- 499 (d) Did you discuss whether and how consent was obtained from people whose data you’re
500 using/curating? [N/A] We do not release a dataset. Public datasets are used with
501 authors’ consents
- 502 (e) Did you discuss whether the data you are using/curating contains personally identifiable
503 information or offensive content? [N/A] No dataset is released and the public data do
504 not contain offensive or personally identifiable information.
- 505 5. If you used crowdsourcing or conducted research with human subjects...
- 506 (a) Did you include the full text of instructions given to participants and screenshots, if
507 applicable? [N/A] No humans involved
- 508 (b) Did you describe any potential participant risks, with links to Institutional Review
509 Board (IRB) approvals, if applicable? [N/A]
- 510 (c) Did you include the estimated hourly wage paid to participants and the total amount
511 spent on participant compensation? [N/A] No humans involved