# Fast Global Localization on Neural Radiance Field

**Mangyu Kong**[1], **Seongwon Lee**[2], **Jaewon Lee**[1], **Euntai Kim**[1*]

[1]Yonsei University  [2]Electronics and Telecommunications Research Institute (ETRI)

{mangyu0929,leejaewon,etkim}@yonsei.ac.kr, sungonce@etri.re.kr

**Abstract:** Neural Radiance Fields (NeRF) presented a novel way to represent scenes, allowing for high-quality 3D reconstruction from 2D images. Following its remarkable achievements, global localization within NeRF maps is an essential task for enabling a wide range of applications. Recently, Loc-NeRF demonstrated a localization approach that combines traditional Monte Carlo Localization with NeRF, showing promising results for using NeRF as an environment map. However, despite its advancements, Loc-NeRF encounters the challenge of a time-intensive ray rendering process, which can be a significant limitation in practical applications. To address this issue, we introduce Fast Loc-NeRF, which leverages a coarse-to-fine approach to enable more efficient and accurate NeRF map-based global localization. Specifically, Fast Loc-NeRF matches rendered pixels and observed images on a multi-resolution from low to high resolution. As a result, it speeds up the costly particle update process while maintaining precise localization results. Additionally, to reject the abnormal particles, we propose particle rejection weighting, which estimates the uncertainty of particles by exploiting NeRF's characteristics and considers them in the particle weighting process. Our Fast Loc-NeRF sets new state-of-the-art localization performances on several benchmarks, convincing its accuracy and efficiency. The code is available at this url

**Keywords:** Neural Radiance Fields, Visual Localization, Global Localization

## 1 Introduction

Global localization is one of the most important challenges in the field of mobile robotics. The detailed algorithm used for global localization is inherently tied to the type of map employed. Different mapping techniques require tailored localization algorithms to determine a robot's position effectively without knowledge about the initial pose. In the early days of mobile robotics, some relatively simple yet efficient maps were commonly used, for example, feature-based maps such as ORB maps [1, 2] or 2D occupancy grid maps [3], but there is a rapid shift towards utilizing a photo-realistic volumetric 3D environment map. Neural Radiance Fields (NeRF) [4] is certainly the most typical example of a realistic volumetric 3D map. The reason for this is that this type of 3D realistic volumetric map can be utilized not only in mobile robotics but also in immersive technologies such as AR/VR/MR. Therefore, global localization under the NeRF map emerges as one of the most urgent topics that should be considered in mobile robotics. Regarding the topic, several works have been reported. Gradient-based pose estimation methods, including iNeRF [5] and PI-NeRF [6], optimize the camera pose by inverting the NeRF. APR-based approaches [7, 8, 9] utilize NeRF to train absolute pose regression networks to predict the pose directly. Among them, recently Loc-NeRF [10] was developed as a groundbreaking approach by integrating NeRF with Monte Carlo Localization (MCL), a probabilistic method traditionally used in mobile robotics for estimating the pose of a robot.

The key idea of Loc-NeRF is to employ NeRF as a map model in the update step of the filter and robot dynamics for motion prediction. This integration allows for global localization and accurate

---

[*]Corresponding author

tracking on NeRF, capitalizing on the high-quality environmental cues captured by NeRF. However, while Loc-NeRF demonstrated impressive localization capabilities, it also revealed a critical bottleneck: the computational intensity of rendering a vast number of rays across numerous particles within the MCL framework. This challenge becomes particularly critical in scenarios requiring real-time localization or operating under computational constraints.

To address this problem, we introduce a new global localization method named Fast Loc-NeRF. The key idea of Fast Loc-NeRF is to divide the iterations of MCL into three distinct phases: the coarse phase, the intermediate phase, and the fine phase. As we progress through each phase, the matching between the rendered image and the observed image is performed at progressively higher resolutions: low resolution in the coarse step, medium resolution in the medium step, and high resolution in the fine step, shifting focus from exploration to exploitation. In response to progressively increasing the resolution, we progressively decrease the number of particles in MCL when we move from the coarse phase to the fine phase, keeping the amount of computation almost constant in all three distinct phases. The coarse phase aims to explore as many candidates as possible while spending a reasonable amount of computation, whereas the fine phase aims to exploit the most likely candidates as intensively as possible to achieve the best estimate. The key idea is summarized in Fig 1. Moreover, we propose an innovative particle rejection weighting scheme within Fast Loc-NeRF. This scheme estimates the uncertainty of each particle by exploiting the unique characteristics of NeRF's rendering process, integrating this uncertainty into the particle weighting mechanism. It is actually the incorporation of rejection sampling into the importance sampling in MCL. Such an approach not only improves the robustness of the localization but also accelerates our method, contributing to the overall performance gains of our method.

Through extensive evaluations, we demonstrate that Fast Loc-NeRF sets new benchmarks in localization performance, offering both efficiency and precision across various standard datasets. Our contributions are aimed at advancing the application of NeRF in visual localization, providing a robust and efficient framework that paves the way for future research and practical implementations.
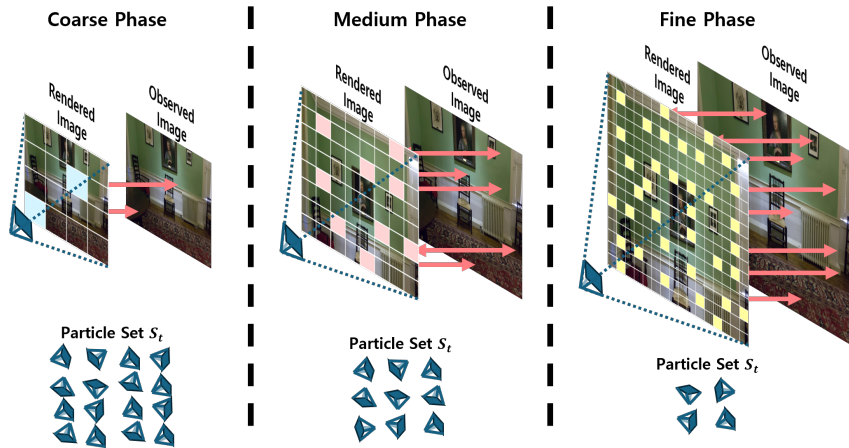


Figure 1: Overview of our proposed Fast Loc-NeRF. Our Fast Loc-NeRF introduces a coarse-to-fine multiscale matching framework to enhance the efficiency and accuracy of NeRF-based localization.

## 2 Related Works

**Neural Radiance Fields (NeRF)** Neural Radiance Fields (NeRF) [4] has emerged as a novel approach demonstrating that a neural network can be used to synthesize novel views of complex scenes. In detail, NeRF models the volumetric scene as a continuous function, using a deep network to map 5D coordinates (spatial location and viewing direction) to color and density. The success of NeRF has driven research aimed at more accurately representing environments, including addressing aliasing issues at various scales [11, 12, 13], achieving faster training and rendering

times [14, 15, 16, 17], and enabling the capability to learn larger scenes [18, 19, 20], as well as other improvements [21, 22, 23]. Additionally, it has extended into various applications such as SLAM [24, 25, 26], editing [27, 28, 29], dynamic scene reconstruction [30, 31, 32] and semantic understanding [33, 34].

**Visual Localization in NeRF** Visual localization, a crucial component in robotics and augmented reality, has been extensively studied, with methods ranging from particle-based approaches to optimization methods. The integration of NeRF into localization tasks represents a relatively new but rapidly growing area of research. The dense and continuous representation provided by NeRF offers a rich source of environmental cues that can be leveraged for more precise localization. Initial efforts in this space have focused on understanding how NeRF's environmental representations can be utilized to improve the accuracy of pose estimation. Several methods have emerged in this field, including gradient-based approaches [5, 6], APR-based methods [7, 8, 9], and Loc-NeRF [10], which utilizes Monte Carlo Localization. Gradient-based methods invert [5, 6] NeRF to estimate pose by rendering images from NeRF and comparing them to actual images at the pixel level. They refine the camera's pose using gradients derived from this comparison. APR-based approaches [7, 8, 9] train absolute pose regression networks, such as PoseNet [35], to directly predict the pose, leveraging the renderable properties of NeRF. However, this method requires extensive training time and resources to train the APR network. Loc-NeRF [10] is a notable example of combining NeRF with Monte Carlo Localization to achieve enhanced localization performance without any training process. However, the computational demands of this integration, particularly in the context of rendering multiple rays for numerous particles, present significant challenges. In this work, we proposed a novel approach named Fast-NeRF, which builds upon these foundations, addressing the computational challenges and enhancing the efficiency and accuracy of NeRF-based localization through coarse-to-fine localization framework and particle rejection weighting strategies.

## 3 Method

Our proposed method, Fast Loc-NeRF, is designed to address the computational challenges of NeRF-based localization, particularly focusing on the inefficiencies associated with the traditional Monte Carlo Localization approach when integrated with NeRF. Fast Loc-NeRF introduces a coarse-to-fine localization framework coupled with a novel particle rejection weighting mechanism, aimed at enhancing both the efficiency and accuracy of the localization process. The overall illustration of our methods is shown in Fig 1.

### 3.1 Monte Carlo Localization on NeRF

**Neural Radiance Fields** Neural radiance fields(NeRF) learn the visual and geometry of the environment and represent a three-dimensional environment that enables novel view synthesis. Given a set of $n$ images $\mathcal{I} = \{I_i\}_{i=1}^{n}$, and their corresponding camera poses and intrinsic parameters, NeRF encodes a function $F_{NeRF} : (x, d) \rightarrow (c, \sigma)$ from 3D coordinate $x = (x, y, z)$ and viewing direction $d = (\theta, \phi)$ to RGB color $c$ and density $\sigma$. The color and the densities of sampled points along the ray are predicted and integrated to render a novel view.

$$\hat{C}(r) = \int_{0}^{\infty} T(z)\sigma(z)c(z)dt, \tag{1}$$

where $T(z) = \exp(-\int_{0}^{z} \sigma(s)ds)$ is the transmittance checking occlusion by integrating along the ray $r(z) = o + zd$ between 0 to $z$. Since the density and radiance values are derived from the 5D coordinate of points on rays, NeRF techniques employ a sampling-based Riemann summation to approximate the integral. NeRF optimizes by the reconstruction error between the rendered colors and the corresponding ground-truth colors.

$$L_{color} = \sum_{r \in \mathcal{R}} \|\hat{C}(r) - C(r)\|^2, \tag{2}$$

where $\mathcal{R}$ is the set of rays $r$ with the known camera poses.

**Monte Carlo Localization leveraging NeRF as a map** With a well-trained Neural Radiance Field (NeRF) map denoted as $\mathcal{M}$, RGB images $I_t$ representing observations at time $t$, and motions $O_t$ at time $t$, a particle filter utilizing Monte Carlo localization is employed to determine the 6 Degrees of Freedom (DoF) camera pose $X_t$ for the corresponding time $t$. This process uses the particle filter to calculate the posterior probability $P(X_t|\mathcal{M}, I_{1:t}, O_{1:t})$, incorporating all images $I_{1:t}$ and motions $O_{1:t}$ recorded from the initial moment up to time $t$. Based on the derived posterior probability $P(X_t|\mathcal{M}, I_{1:t}, O_{1:t})$, the particle filter then samples a weighted set of $n$ particles to represent the distribution of possible camera poses as follows:

$$\mathcal{S}_t = \{\langle x_t^{[i]}, w_t^{[i]} \rangle \mid i = 1, \ldots, N\}, \tag{3}$$

where $x_t^{[i]}$ is the 6DoF pose of the $i$th particle and $w_t^{[i]}$ is the associated weight of $x_t^{[i]}$ normalized from 0 to 1. The particle filter updates the group of particles every time new images and measurements are observed. During the prediction step, the set of particles $\mathcal{S}_t$ at time $t$ is estimated from the set of particles $\mathcal{S}_{t-1}$ with the motion of the camera $O_t$ between time $t$ and time $t-1$ as the motion model $P(X_t|X_{t-1}, O_t)$. The pose $x_t^{[i]}$ of each particle is sampled through Gaussian noise as follows $x_t^{[i]} = x_{t-1}^{[i]} \cdot O_t \cdot x_\epsilon^{[i]}$, where $x_\epsilon \in SE(3)$ is the prediction noise modeled as an exponential map $\text{Exp}(\delta)$. Then, the importance weights of particles are updated based on the measurement likelihood $P(I_t|X_t, \mathcal{M})$, which indicates the probability of the pose $x_t^{[i]}$ given the observed image $I_t$ at time $t$. Loc-NeRF [10] employed a heuristic function to estimate the measurement likelihood as follows:

$$w_t^{[i]} = \left( \frac{B}{\sum_{j=1}^{B} \left( I_t(p_j) - \hat{C}(r(p_j, x_t^{[i]}), K) \right)^2} \right)^4, \tag{4}$$

where $B$ denotes the number of pixels rendered per particle. $\hat{C}(r(p_j, x_t^{[i]}), K)$ is the rendered color on the pixel $p_j$ with the pose $x_t^{[i]}$ and shared intrinsic camera parameter $K$. After importance weighting, $N$ particles are resampled with normalized weights $w_t^{[i]}$.

## 3.2 Fast Loc-NeRF

**Framework of Fast Loc-NeRF** Our Fast Loc-NeRF consists of three phases: coarse, medium, and fine phases, shifting focus from exploration to exploitation. The key idea of Fast Loc-NeRF is to maintain the amount of computation represented by the product of the number of particles $N$ and the number of probings $B$ (renderings or queries) almost constant in all three phases, achieving both efficiency and precision. Specifically, as the filtering steps progress from coarse to fine, the focus shifts from utilizing a large number of particles $N$ to increasing the number of query pixels $B$ for precise localization. This adjustment narrows down the search space and focuses on the most promising solutions. During the coarse phase, fewer renderings challenge the capture of the image's overall shape,

---

**Algorithm 1:** Fast Loc-NeRF

**Input:** $I_t, u_t$
Initialize $X_0$; Initialize $X_0$;
$R = \{R_{\text{coarse}}, R_{\text{medium}}, R_{\text{fine}}\}$// Rendering scale
$B = \{B_{\text{coarse}}, B_{\text{medium}}, B_{\text{fine}}\}$// # of queries
$N = \{N_{\text{coarse}}, N_{\text{medium}}, N_{\text{fine}}\}$// # of particles
// Coarse Phase
**for** $t = 1$ *to* $T_{coarse}$ **do**
    $X_t =$
      MCL_NeRF$(X_{t-1}; R_{\text{coarse}}, B_{\text{coarse}}, N_{\text{coarse}})$
// Medium Phase
**for** $t = T_{coarse} + 1$ *to* $T_{medium}$ **do**
    $X_t =$
      MCL_NeRF$(X_{t-1}; R_{\text{medium}}, B_{\text{medium}}, N_{\text{medium}})$
// Fine Phase
**for** $t = T_{medium} + 1$ *to* $T_{fine}$ **do**
    $X_t = $ MCL_NeRF$(X_{t-1}; R_{\text{fine}}, B_{\text{fine}}, N_{\text{fine}})$
**return** $X_t$

---

so we adjust the rendering resolution $R$ in each phase to maintain the matching region's proportion over the image. Let us consider Algorithm 1. In the algorithm, $R$, $B$, and $N$ denote the rendering scale, the number of queries, and the number of particles, respectively. As we move from the coarse phase, we increase the rendering scale $R$ (for example, $R_{\text{coarse}}, R_{\text{medium}}, R_{\text{fine}} = 1/4, 1/2, 1$), and increase the number of queries $B$, ($B_{\text{coarse}}, B_{\text{medium}}, B_{\text{fine}} = 1/4, 1/2, 1$), but decrease the number of particles $N$, ($N_{\text{coarse}}, N_{\text{medium}}, N_{\text{fine}} = 1, 1/2, 1/4$).

**Detailed Structure** Fast Loc-NeRF divides three phases according to the initial distribution of particles and their convergence over time. The algorithm begins with the coarse phase, corresponding to the update at $t = 0$ from the initial particle distribution. After the initial exploration in the coarse phase, the algorithm proceeds to the medium phase, where the process is repeated until the particles converge. When the distribution of the 3D positions of the particles, $\{x_t^{[i]}\}$, falls within a certain threshold, $\delta_{\text{refine}}$, as $\text{Var}(\{x_t^{[i]}\}) < \delta_{\text{refine}}$, the algorithm enters to the fine phase. In this phase, a small number of particles is used to focus on precise localization by performing detailed rendering in high resolution, ensuring that the final estimated pose is as accurate as possible.

---

**Algorithm 2:** MCL-NeRF

**Input:** $X_{t-1}, u_t, I_t, R_t, B_t, N_t$

$\bar{X}_t = X_t = \emptyset$

**for** $i = 1$ *to* $N_{t-1}$ **do**

    sample $x_t^{[i]} \sim p(x_t \mid u_t, x_{t-1}^{[i]})$

    $I_t^{R_t} \leftarrow$ downscale$(I_t, R_t)$

    **for** $j = 1$ *to* $B_t$ **do**

        render $\hat{C}(r(p_j, x_t^{[i]}), R_t)$

    $w_t^{[i]} = p(I_t^{R_t} \mid x_t^{[i]}, \mathcal{M})$

    calculate $w_t^{[i]}$ through Eq (6)

    $\bar{X}_t = \bar{X}_t + \{\langle x_t^{[i]}, w_t^{[i]} \rangle\}$

**for** $i = 1$ *to* $N_t$ **do**

    draw $i$ with probability $\propto w_t^{[i]}$

    add $x_t^{[i]}$ to $X_t$

**return** $X_t$

---

**MCL-NeRF** In each phase, we run MCL-NeRF given in Algorithm 2 to update the particles. MCL-NeRF consists of sampling, importance weighting, and resampling, similarly performed as shown in Algorithm 2. To achieve importance weighting at the corresponding resolution $R_t$, the observed image $I_t$ is resized to $I_t^{R_t}$. The key difference between the vanilla MCL and MCL-NeRF is that we have to downsize the observed image to $I_t^{R_t}$ using the rendering scale $R_t$ and render the image $\hat{C}(r(p_j, x_t^{[i]}), R_t \cdot K)$ of query points $\{p_j\}_{j=1}^{B_t}$ to compute the importance weight $w^{[i]}$.

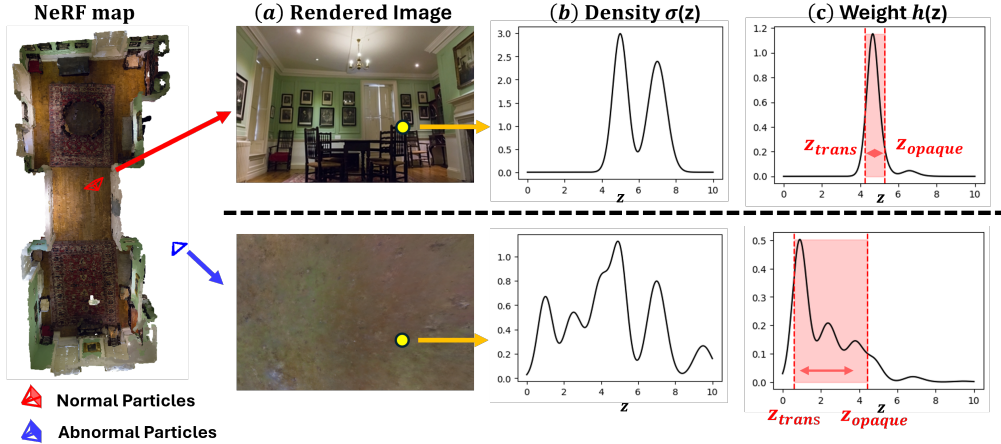### 3.3 Particle Rejection Weighting



Figure 2: **Particle Rejection Weighting** We plot the particle rejection weighting process. (a) Rendered image from the red and blue particle each, which is located inside and outside room (b) Density distribution $\sigma(z)$ over ray $r(z) = o + zd$ (c) Ray weight distribution $h(z) = T(z)\sigma(z)$ and red region indicate the range between $z_{\text{trans}}$ and $z_{\text{opaque}}$

When we apply a vanilla MCL to a NeRF map, we often encounter situations where early-stage particles are randomly or uniformly distributed throughout the environment, and some of them fall into regions where the NeRF is not trained at all. This occurs because we have no information about the structure of the environment. An example of this situation is shown in the leftmost figure in Fig 2. In the figure, the red particle (in the first row) falls into the "valid" region where the NeRF map is trained and successfully renders the corresponding image from the NeRF map. In contrast, the blue particle (in the second row) happens to fall outside the environment ("invalid

region") where the NeRF map is not trained at all. As a result, the blue particle renders a strange image, and such abnormal particles should be removed from subsequent processing. In this paper, we remove such abnormal particles by incorporating rejection sampling into the importance sampling framework. While importance sampling and rejection sampling are distinct techniques, they are related in that both methods use the ratio of the target distribution to the proposal distribution to draw samples. Specifically, we assign very low importance weights to particles that should be removed from further consideration, effectively rejecting them and focusing on more promising particles within the importance sampling context.

The remaining problem is how to compute the importance weight so that the abnormal particles are removed in the context of importance sampling. When a normal particle located in the valid region renders an image by querying a NeRF with rays (for example, a yellow circle in Fig 2), the density $\sigma(z)$ is quite low for empty spaces but abruptly increases for objects such as walls, as shown in the first row of Fig 2 (b). In this case, if we define $z_{\text{trans}}$ and $z_{\text{opaque}}$ as the points on the ray until the ray remains transparent and the point from which the ray becomes opaque, respectively, and compute them by:

$$
\begin{aligned}
z_{\text{trans}} &= \min \left\{ z \mid \int_0^z T(s)\sigma(s)\,ds \geq \alpha \right\}, \\
z_{\text{opaque}} &= \max \left\{ z \mid \int_0^z T(s)\sigma(s)\,ds \leq 1 - \alpha \right\}.
\end{aligned}
\tag{5}
$$

We can see that the difference between the two points $z_{\text{opaque}}$ and $z_{\text{trans}}$ is relatively small, as shown in the first row of Fig 2 (c), where $\alpha$ is a predefined threshold close to zero. On the contrary, when the abnormal particle located on the invalid region renders an image by querying a NeRF with rays, the density $\sigma(z)$ tends to have meaningless values, lacking a clear distinction between the empty space and the object, as shown in the second row of Fig 2 (b). Consequently, the corresponding ray weight distribution $h(z) = T(z)\sigma(z)$ does not have a clear distinction between the transparent point $z_{\text{trans}}$ and the opaque point $z_{\text{opaque}}$, as shown in the second row of Fig 2 (c). In this case, the difference between the two points $z_{\text{opaque}}$ and $z_{\text{trans}}$ is relatively large. To remove such abnormal particles in the context of importance sampling, we weight each particle by:

$$
w_t^{[i]} = \left( \frac{B_t}{\sum_{j=1}^{B_t} \left( I_t^{R_t}(p_j) - \hat{C}(r(p_j, x_t^{[i]}), R_t \cdot K) \right)^2 \cdot F(r(p_j, x_t^{[i]}))} \right)^4,
\tag{6}
$$

where the additional term $F(r(p_j)) = \max(z_{\text{opaque}} - z_{\text{trans}}, \tau)$ penalizes the abnormal particles using the gap between the transparent point $z_{\text{trans}}$ and the opaque point $z_{\text{opaque}}$. $\tau$ is a minimum bound of $F(r(p_j))$, and it simply ensures the algorithm stability.

## 4  Experiments

Section 4.1 focuses on benchmarking Fast Loc-NeRF against Loc-NeRF [10], iNeRF [5], and PI-NeRF (Parallel Inversion of NeRF) [6] in global localization settings to demonstrate its superior accuracy and speed. Section 4.2 explores the benefits of incorporating multi-scale representation, utilizing Zip-NeRF [13] as a map. In addition to comparisons on the real-world forward-facing dataset [36], as done in previous works, this section also showcases the localization performance in a larger real-world indoor dataset [37]. Analysis of module ablation and initial particle number are included in the supplementary materials.

### 4.1  Global Localization

**Experiment Setting:** For our evaluation, we conducted global localization experiments across all 8 forward-facing scenes of the LLFF dataset [36]. We generated maps using a basic NeRF model [4]

with NeRF-pytorch [38] same as comparisons, and trained for 20k iterations to represent each scene accurately. To ensure a fair comparison with Loc-NeRF under identical conditions, we randomly selected 5 test images from each scene and performed 5 localization trials per image. This resulted in a total of 25 experiments per scene, aiming to minimize randomness by increasing the number of trials and averaging the results. For our method's setting, observed images were downscaled to $R = \{1/4, 1/2, 1\}$, the number of particles is $N = \{9600, 600, 100\}$ and the number of rendered rays are $B = \{8, 16, 32\}$ at each phase.

| Method | Fern | Flower | Leaves | Horns | Fortress | Room | Orchids | Trex | Time |
|---|---|---|---|---|---|---|---|---|---|
| iNeRF [5] | 0.75/167 | 0.92/75.0 | 1.07/145 | 1.12/138 | 1.26/136 | 1.47/111 | 1.22/152 | 1.06/132 | **0.208** |
| PI-NeRF [6] | 1.02/103 | 0.98/96.8 | 0.97/101 | 1.07/93.7 | 1.03/92.4 | 1.05/96.9 | 1.00/98.6 | 1.08/70.0 | - |
| Loc-NeRF [10] | 0.03/0.64 | 0.56/18.1 | 0.89/4.58 | 0.39/3.84 | 0.03/**0.96** | 0.06/0.94 | 1.43/35.0 | 0.64/5.68 | 0.724 |
| Ours | **0.02/0.59** | **0.01/0.19** | **0.21/0.72** | **0.02/0.23** | **0.03**/1.01 | **0.05/0.94** | **0.38/3.47** | **0.14/1.81** | 0.327 |

Table 1: Comparison of Fast Loc-NeRF and other state-of-the-art methods on the LLFF dataset, with position error (cm), rotation error (deg), and the time for each update step (sec).
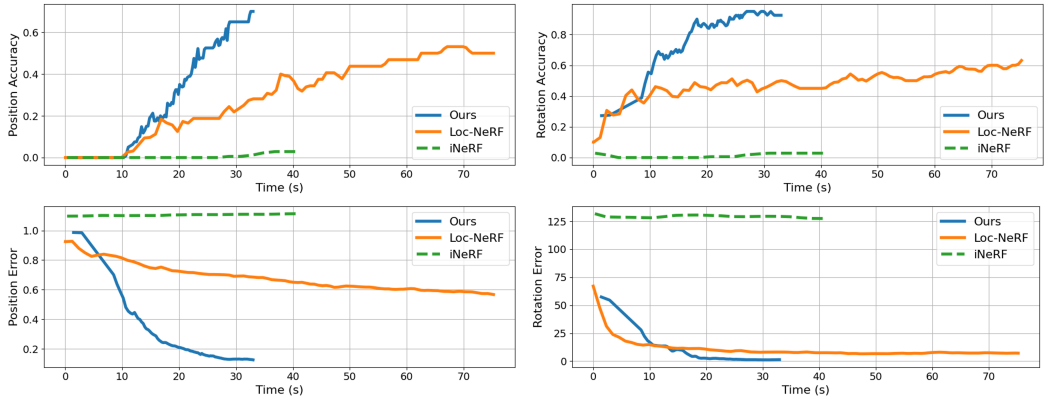


Figure 3: Quantitative Comparison with Loc-NeRF and iNeRF. The first column presents the comparison result of position accuracy and error. The second column presents the comparison result of rotation accuracy and error.

**Results on LLFF dataset** In this section, we demonstrate the outcomes of Fast Loc-NeRF, Loc-NeRF [10], and gradient-based methods [5, 6], showcasing the comparison of position accuracy and error, alongside rotation error and accuracy over time, as depicted in Fig 3. In NeRF Localization, pose accuracy is defined as having a position error within 5 cm and a rotation error under 5 degrees. From the results illustrated in Fig 3, iNeRF [5] and PI-NeRF [6] struggle to meet the target pose from initial poses, failing to demonstrate convergence. Unlike gradient-based methods [5, 6], Loc-NeRF [10], utilizing Monte Carlo Localization (MCL), estimates the true position and rotation in many cases despite challenging starting points, ultimately achieving a localization accuracy of 53%. However, achieving this level of position accuracy with Loc-NeRF takes 74 seconds. Our approach significantly reduces this time, reaching similar accuracy in just 24 seconds. Moreover, when compared at equivalent update iterations, Fast Loc-NeRF exhibits a 55% improvement in average time per update step and a 19% higher performance in final accuracy, highlighting the efficiency and effectiveness of our proposed method in NeRF-based global localization tasks.

## 4.2 Localization on Multi-Scale Representation

In this section, we integrate Zip-NeRF [13] as a NeRF map within our Fast Loc-NeRF framework, employing multi-scale representation which allows for accurate rendering at various resolutions without aliasing. We expect the synergy between multi-scale matching and multi-scale representation to enable more accurate particle weighting.

**Experiment Setup**: We utilized multi-scale representation via Zip-NeRF [13] to conduct comparative analyses between Loc-NeRF and Fast Loc-NeRF. The experimental dataset consists of all 8

| | Fern | Leaves | Flower | Horns | Fortress | Room | Orchids | Trex | Time per update |
|---|---|---|---|---|---|---|---|---|---|
| Loc-NeRF | 1.54/57.00 | 0.63/14.5 | 0.37/33.00 | 0.64/12.00 | 0.28/1.59 | 0.34/5.60 | **1.10**/133.00 | 0.54/4.45 | 0.205 |
| Ours | **0.67/16.40** | **0.32/4.30** | **0.10/21.00** | **0.02/0.22** | **0.02/0.49** | **0.02/0.38** | 1.80/**84.00** | **0.06/0.73** | **0.118** |

(a) Position (cm) /Rotation (deg) Error

| | Fern | Leaves | Flower | Horns | Fortress | Room | Orchids | Trex | Time per update |
|---|---|---|---|---|---|---|---|---|---|
| Loc-NeRF | 0.00/0.00 | 0.10/0.70 | 0.55/0.65 | 0.50/0.75 | 0.55/0.95 | 0.35/0.60 | 0.00/0.00 | 0.70/0.75 | 0.205 |
| Ours | **0.68/0.68** | **0.28/0.80** | **0.88/0.88** | **1.00/1.00** | **1.00/1.00** | **0.96/1.00** | **0.20/0.20** | **0.80/0.96** | **0.118** |

(b) Position/Rotation Accuracy

Table 2: Comparison on LLFF dataset [36] using Zip-NeRF [13] as a map.

scenes from the LLFF dataset [36], consistent with Section 4.1, and is extended to include two additional scenes from the Deep Blending dataset [37]: $playroom$ and $Drjohnson$. The Deep Blending dataset [37], unlike the LLFF dataset [36] which comprises primarily forward-facing scenes, showcases entire rooms, making it particularly suitable for demonstrating global localization for the large real world.

**Evaluation on LLFF dataset** We present the results of employing Zip-NeRF [13] as a map on our method and Loc-NeRF in Table 2. Our approach, Fast Loc-NeRF demonstrated superior performance against baseline across most of the forward-facing scenes. While Loc-NeRF [10] had difficulty reaching a 0.8 accuracy threshold in any scene, our Fast Loc-NeRF method consistently attained this benchmark across all scenes except for $fern$ and $orchids$ scenes, still demonstrating the highest performance.

**Evaluation on Deep Blending dataset [37]** In the evaluation conducted on the Deep Blending dataset, we observed that each scene, representing a single room, presents a more challenging environment for global localization compared to the forward-facing scenes in the LLFF dataset. Our method significantly outperforms Loc-NeRF shown on Table 3. This enhancement suggests that our method efficiently understands and explores extensive spatial environments, such as larger rooms, indicating its adaptability and efficiency in addressing more complex localization challenges.

| | Drjohnson | Playroom | Average |
|---|---|---|---|
| Loc-NeRF | 0.83/54.40 | 0.97/35.65 | 0.90/45.03 |
| Ours | **0.42/34.04** | **0.37/33.40** | **0.40/33.72** |

(a) Position (cm)/Rotation (deg) Error

| | Drjohnson | Playroom | Average |
|---|---|---|---|
| Loc-NeRF | 0.10/0.10 | 0.10/0.15 | 0.10/0.13 |
| Ours | **0.35/0.25** | **0.25/0.40** | **0.30/0.33** |

(b) Position/Rotation Accuracy

Table 3: Comparison on Deep Blending Dataset [37] using Zip-NeRF [13] as a map.

# 5   Conclusion

In summary, we introduce Fast Loc-NeRF, which exploits particle filtering in a coarse-to-fine manner to enable more efficient and accurate localization on a NeRF map. Leveraging coarse-to-fine matching, Fast Loc-NeRF speeds up the expansive particle weight update process, while progressively matching particle states and measurements from low to high resolution. Furthermore, to deal with particles distributed in irrelevant locations, we propose particle rejection weighting, which estimates ray uncertainty by analyzing the distribution over each ray without any additional training step. By not simply using the RGB values obtained from rendering to weight the particles required for the update step, but also applying the uncertainty of each ray, we reduced the exploration of inappropriate space, enabling faster localization. As a result, our method speeds up the localization process by 55% and improves localization accuracy by 19% with the evaluation of various NeRF datasets and multi-scale NeRF models.

**Limitation and Future Work** Our Fast Loc-NeRF improves both the performance and efficiency of global localization on NeRF maps. However, it still faces limitations due to the inherent inference speed constraints of the NeRF model itself even utilizing a fast inference model like Zip-NeRF [13]. In future work, we plan to integrate and evaluate the concept of Fast Loc-NeRF with 3D renderable maps that offer faster rendering speeds, such as 3D Gaussian splitting [39].

# References

[1] R. Mur-Artal and J. D. Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE transactions on robotics*, 33(5):1255–1262, 2017.

[2] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós. Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021.

[3] G. Grisetti, C. Stachniss, and W. Burgard. Improved techniques for grid mapping with rao-blackwellized particle filters. *IEEE transactions on Robotics*, 23(1):34–46, 2007.

[4] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.

[5] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin. inerf: Inverting neural radiance fields for pose estimation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1323–1330. IEEE, 2021.

[6] Y. Lin, T. Müller, J. Tremblay, B. Wen, S. Tyree, A. Evans, P. A. Vela, and S. Birchfield. Parallel inversion of neural radiance fields for robust pose estimation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9377–9384. IEEE, 2023.

[7] S. Chen, Z. Wang, and V. Prisacariu. Direct-posenet: Absolute pose regression with photometric consistency. In *2021 International Conference on 3D Vision (3DV)*, pages 1175–1185. IEEE, 2021.

[8] A. Moreau, N. Piasco, D. Tsishkou, B. Stanciulescu, and A. de La Fortelle. Lens: Localization enhanced by nerf synthesis. In *Conference on Robot Learning*, pages 1347–1356. PMLR, 2022.

[9] S. Chen, X. Li, Z. Wang, and V. A. Prisacariu. Dfnet: Enhance absolute pose regression with direct feature matching. In *European Conference on Computer Vision*, pages 1–17. Springer, 2022.

[10] D. Maggio, M. Abate, J. Shi, C. Mario, and L. Carlone. Loc-nerf: Monte carlo localization using neural radiance fields. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4018–4025. IEEE, 2023.

[11] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.

[12] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.

[13] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *arXiv preprint arXiv:2304.06706*, 2023.

[14] T. Müller, A. Evans, C. Schied, and A. Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.

[15] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.

[16] A. Chen, Z. Xu, A. Geiger, J. Yu, and H. Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022.

[17] C. Sun, M. Sun, and H.-T. Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. *CVPR*, 2022.

[18] H. Turki, D. Ramanan, and M. Satyanarayanan. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12922–12931, 2022.

[19] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022.

[20] Z. Mi and D. Xu. Switch-nerf: Learning scene decomposition with mixture of experts for large-scale neural radiance fields. In *International Conference on Learning Representations (ICLR)*, 2023. URL https://openreview.net/forum?id=PQ2zoIZqvm.

[21] K. Deng, A. Liu, J.-Y. Zhu, and D. Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022.

[22] W. Bian, Z. Wang, K. Li, J.-W. Bian, and V. A. Prisacariu. Nope-nerf: Optimising neural radiance field with no pose prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4160–4169, 2023.

[23] E. Sandström, Y. Li, L. Van Gool, and M. R. Oswald. Point-slam: Dense neural point cloud-based slam. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18433–18444, 2023.

[24] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12786–12796, 2022.

[25] A. Rosinol, J. J. Leonard, and L. Carlone. Nerf-slam: Real-time dense monocular slam with neural radiance fields. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3437–3444. IEEE, 2023.

[26] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6229–6238, 2021.

[27] Y.-J. Yuan, Y.-T. Sun, Y.-K. Lai, Y. Ma, R. Jia, and L. Gao. Nerf-editing: geometry editing of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18353–18364, 2022.

[28] A. Mirzaei, T. Aumentado-Armstrong, K. G. Derpanis, J. Kelly, M. A. Brubaker, I. Gilitschenski, and A. Levinshtein. Spin-nerf: Multiview segmentation and perceptual inpainting with neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20669–20679, 2023.

[29] M. Kong, S. Lee, and E. Kim. Roomnerf: Representing empty room as neural radiance fields for view synthesis. In *British Machine Vision Conference*, 2023. URL https://api.semanticscholar.org/CorpusID:267000673.

[30] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021.

[31] S. Fridovich-Keil, G. Meanti, F. R. Warburg, B. Recht, and A. Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023.

[32] Z. Li, Q. Wang, F. Cole, R. Tucker, and N. Snavely. Dynibar: Neural dynamic image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4273–4284, 2023.

[33] S. Zhi, T. Laidlow, S. Leutenegger, and A. J. Davison. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15838–15847, 2021.

[34] J. Cen, Z. Zhou, J. Fang, W. Shen, L. Xie, D. Jiang, X. Zhang, Q. Tian, et al. Segment anything in 3d with nerfs. *Advances in Neural Information Processing Systems*, 36:25971–25990, 2023.

[35] A. Kendall, M. Grimes, and R. Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015.

[36] B. Mildenhall, P. P. Srinivasan, R. Ortiz-Cayon, N. K. Kalantari, R. Ramamoorthi, R. Ng, and A. Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019.

[37] P. Hedman, J. Philip, T. Price, J.-M. Frahm, G. Drettakis, and G. Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)*, 37(6):1–15, 2018.

[38] L. Yen-Chen. Nerf-pytorch. https://github.com/yenchenlin/nerf-pytorch/, 2020.

[39] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.

[40] W. Hu, Y. Wang, L. Ma, B. Yang, L. Gao, X. Liu, and Y. Ma. Tri-miprf: Tri-mip representation for efficient anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19774–19783, 2023.

# Appendix

## A   Detailed Experiment Setting

For the global localization experiments in Section 4.1, we followed the initial distribution settings of Loc-NeRF [10], generating offset pose from the ground truth pose within a range of [-1m, 1m] for each axis. Particles were then uniformly sampled within a 2m × 2m × 2m cube centered on these offset poses. Additionally, we sample yaw angles from a uniform distribution within the range of [-180 °, 180°] and rotate the sampled particle. In comparison to our method, Loc-NeRF was configured with a batch size (ray per particle) of 32 and utilized 600 particles. We also applied *particle annealing*, reducing the number of particles to 100. In the case of iNeRF [5] and PI-NeRF [6], the experimental setting follows the original methodology of each paper. The offset pose is used as the initial pose for the experiments on gradient-based optimization comparisons. For iNeRF [5], we use $batch\_size = 2024$ per each gradient step, and interest region sampling as a point sampling method. Adam optimizer with an exponentially decaying learning rate ($\beta_1 = 0.9$ and $\beta_2 = 0.999$) is used to optimize the camera pose. For PI-NeRF [6], which utilized instant-NGP [14] as a map, we trained 8 scenes of LLFF dataset [36] on instant-NGP to evaluate the global localization performance. In the optimization process, 64 camera samples were utilized as parameters, and the L2 loss function was applied as the loss function. Since the update time varies depending on the rendering speed of the map, the time per update step for PI-NeRF [6] was not included in Table 1, as instant-NGP was used as the map instead of the basic NeRF [4].

## B   Ablation Studies

In this section, we perform a module ablation on Fast Loc-NeRF, analyzing the impact of each contribution on global localization on NeRF: coarse-to-fine framework and particle rejection weighting. Additionally, we present analyses of particle rejection weighting and the effect of the initial number of particles.

### B.1   Module Ablation

Fast Loc-NeRF proposes two methods: a coarse-to-fine localization framework utilizing MCL, which focuses on exploration to exploitation performing at progressively higher resolution, and a particle rejection weighting method that assigns low weights to abnormal particles outside the searching area. Table B.1 shows the results of applying each method separately. By applying a coarse-to-fine localization framework, the overall time per update step is halved while maintaining localization accuracy as shown in Table B.1. Additionally, the implementation of particle rejection weighting not only improves position accuracy by 13% but also accelerates convergence by assigning lower weights to abnormal particles, leading to a reduction in overall computation time.

| Coarse-to-Fine | Rejection | Position | | Rotation | | Time per step (sec) |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | Accuracy(%) | Error(cm) | Accuracy(%) | Error(∘) | |
| | | 52.5 | 0.50 | 65.0 | 5.84 | 0.7450 |
| ✓ | | 55.0 | 0.45 | 75.0 | 3.18 | 0.3836 |
| | ✓ | 65.0 | 0.33 | 85.0 | 8.12 | 0.5056 |
| ✓ | ✓ | **75.0** | **0.11** | **97.5** | **1.09** | **0.3267** |

Table B.1: Results of the module ablation on Fast Loc-NeRF. The table presents the position error (cm) and accuracy (%), rotation error (deg) and accuracy (%), and time per update step (sec) for different combinations of coarse-to-fine multi-scale matching and particle rejection weighting methods in LLFF dataset [36].

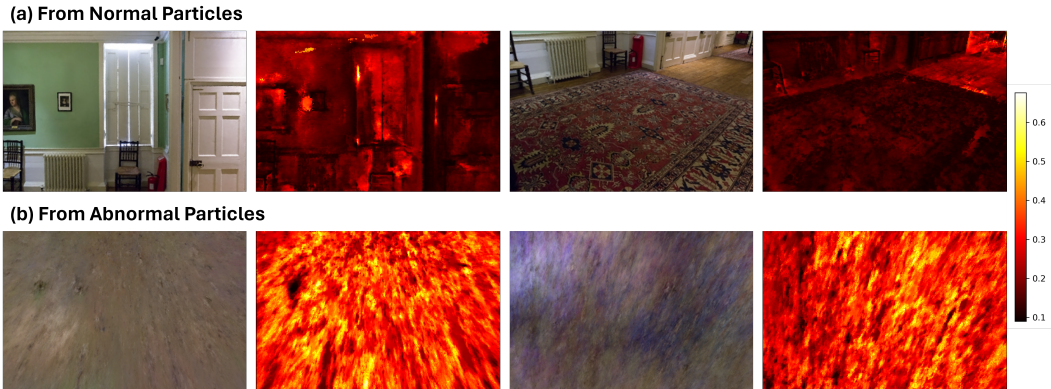**(a) From Normal Particles**

**(b) From Abnormal Particles**

Figure B.1: **Uncertainty Visualization** Rendered image and visualization of the uncertainty ($z_{\text{opaque}} - z_{\text{trans}}$) of each pixel from normal particles and abnormal particles on $Drjohnson$ scene of Deep Blending dataset [37]. The first row shows the properly rendered image from normal particles located on a valid region and the corresponding uncertainty for each pixel. The second row shows the strange rendered image from abnormal particles located on an invalid region and the corresponding uncertainty for each pixel, showing high values overall.

## B.2 Analysis on Particle Rejection Weighting

**Ray Uncertainty Visualization** The particle rejection weighting method is designed to filter out abnormal particles by incorporating particle rejection sampling into importance sampling. As we described in Equation 5, $z_{\text{trans}}$ and $z_{\text{opaque}}$ represent the distances along a ray at which the ray remains transparent and at which it becomes opaque, respectively. By utilizing the distance between $z_{\text{trans}}$ and $z_{\text{opaque}}$, we estimate the uncertainty of each ray. Fig B.1 visualizes the rendered images and the uncertainty of each pixel from particles. The first row shows a rendered image from normal particles on a valid region and visualizes the corresponding uncertainty for each pixel. The second row displays a rendered image and the associated uncertainty from abnormal particles on an invalid region. With normal particles, proper images well-representing scenes are rendered, while abnormal particles produce strange images. Also, the uncertainty values ($z_{\text{opaque}} - z_{\text{trans}}$) from these abnormal particles are generally high, compared to those from normal particles. This indicates that, through particle rejection weighting, abnormal particles receive lower weights, increasing the likelihood of their rejection during the resampling process. Furthermore, we show the uncertainty visualization on Fig C.1 across various scenes in the LLFF dataset [36].

**Ablation Studies on Threshold of Particle Rejection Weighting**

Our particle rejection weighting method stabilizes by setting a minimum bound $\tau$ on the gap between the last transparent point $z_{\text{trans}}$ and the opaque point $z_{\text{opaque}}$ as $F(r(p_j)) = \max(z_{\text{opaque}} - z_{\text{trans}}, \tau)$ on Equation 6. If the gap between $z_{\text{trans}}$ and $z_{\text{opaque}}$ is very small, such as 0.01, it results in assigning too high weights to particles with rays having a small gap between these two points, rather than serving the purpose of rejecting particles on the invalid region. To prevent overweighting and ensure stability, a minimum bound $\tau$ is imposed. Fig B.2 presents the results for various minimum bounds
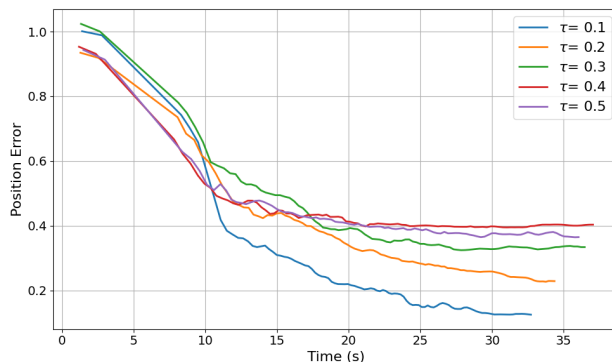


Figure B.2: **Ablation results on minimum bound $\tau$.** The graph shows position errors (cm) over time (sec) for different $\tau$ values evaluated on LLFF dataset [36].

13

$\tau$. As observed from Fig B.2, smaller values of $\tau$ make the particle rejection weighing more effective in reducing position errors. Additionally, when $\tau$ exceeds 0.4, the results become similar to those obtained without particle rejection weighting. The reason is that when the minimum value becomes large, the uncertainty no longer significantly influences the importance weighting process. Fig B.2 shows that our proposed particle rejection weighting effectively eliminates abnormal particles, resulting in improved global localization performance and enhanced convergence to the valid region.

## B.3  Analysis on # of initial particle

In this section, we address the number of particles used during the coarse phase of the Fast Loc-NeRF framework. Since no prior information is provided on the NeRF map $\mathcal{M}$, we initially sample the particles randomly or uniformly across the searching region. The randomly sampled particles in the initial distribution play a crucial role in the exploration process. Table B.2 illustrates the results of the global localization for varying numbers of initial particles, from 600 to 38400. Increasing the number of initial particles enhances exploration, leading to reduced position error and accuracy.



Figure B.3: **Ablation results on # of initial particles.** The graph shows position errors (cm) over time (sec) for different # of initial particles evaluated on LLFF dataset [36].

However, beyond a certain number of particles, the localization inference time during the coarse phase increases exponentially. Therefore, in the coarse phase, the number of particles was set to 9600, considering the tradeoff between inference time and performance. This parameter can be adjusted based on the size of the searching space.
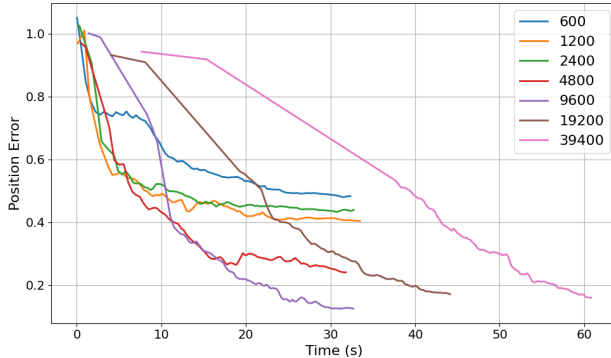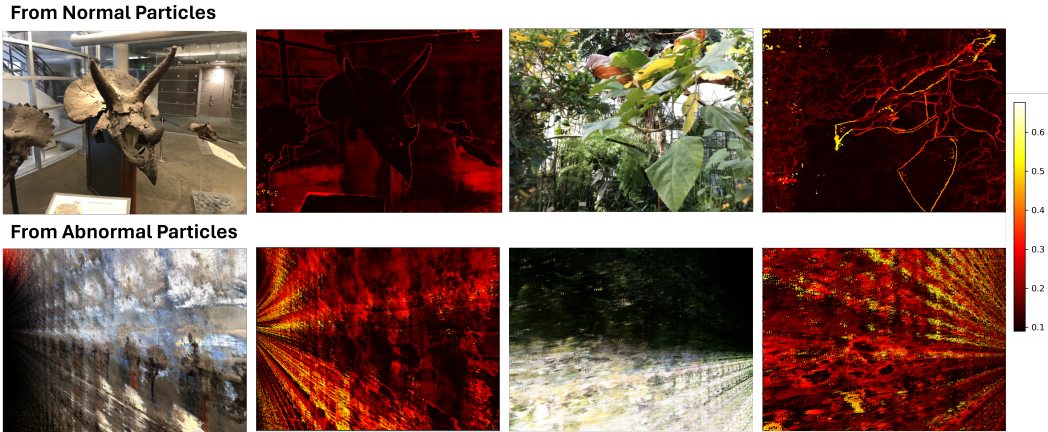
| # of Initial Particle | Position | | Rotation | | Time per step (sec) |
|---|---|---|---|---|---|
| | Accuracy(%) | Error(cm) | Accuracy(%) | Error(∘) | |
| 600 | 42.5 | 0.51 | 67.5 | 9.30 | 0.3818 |
| 1200 | 45.0 | 0.51 | 72.5 | 10.4 | 0.3645 |
| 2400 | 60.0 | 0.30 | 80.0 | 3.56 | **0.2996** |
| 4800 | 62.5 | 0.20 | 90.0 | 10.9 | 0.3329 |
| 9600 | 75.0 | 0.11 | **97.5** | **1.09** | 0.3237 |
| 19200 | 77.5 | 0.16 | 92.5 | 8.12 | 0.4371 |
| 38400 | **77.5** | **0.08** | 95.0 | 1.11 | 0.6020 |

Table B.2: Ablation study on the effect of the initial number of particles in Fast Loc-NeRF. The table presents the position error (cm) and rotation error (deg) for different initial particle counts.
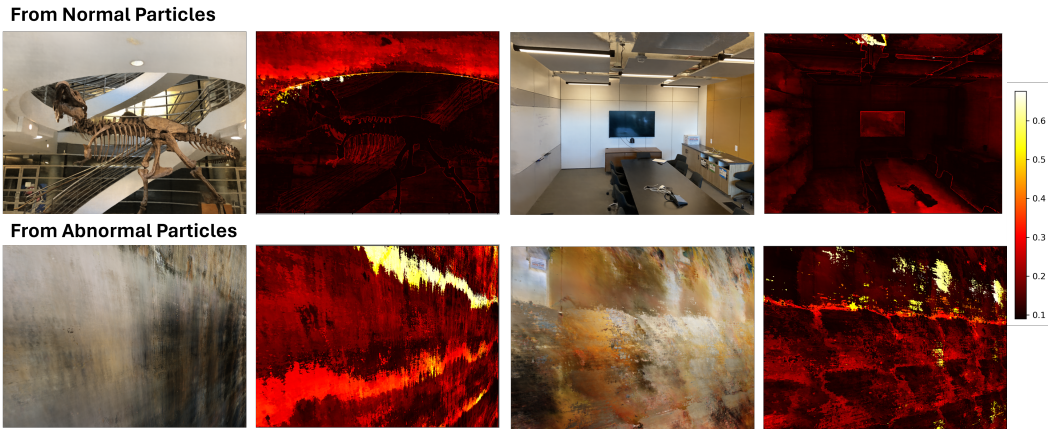
## C  Implementation Detail on Zip-NeRF

Fast loc-NeRF progressively increases the resolution of rendered rays in a coarse-to-fine strategy. Our weighting method matches rendered images with observed images across various resolutions. Therefore, a 3D model that can accurately represent and render without aliasing issues at different resolutions is suitable for our proposed coarse-to-fine multiscale matching. Among various multiscale representations [11, 12, 13, 40], we utilized Zip-NeRF [13], a grid-based approach with fast rendering speed, to serve as a map for global localization, as evaluated in Section 4.2.

To train the multi-scale representation using Zip-NeRF [13] as a map, we converted each training image into a set of four bicubically down-sampled images. By utilizing 4 sets of bicubically down-sampled images as training data, we generated a generalized model capable of effectively handling various scales. We trained the Zip-NeRF model on 8 scenes from the LLFF dataset [36] and 2 scenes from the Deep Blending dataset [37], each for 25k iterations. Zip-NeRF was implemented with two proposal MLPs and 1 final NeRF MLP. We use 64 samples for two rounds of proposal sampling each and 32 samples for the final NeRF sampling.



(a) Rendered image and uncertainty visualization from $Horns$ and $Orchichs$ scenes of LLFF dataset.



(b) Rendered image and uncertainty visualization from $Trex$ and $Room$ scenes of LLFF dataset

Figure C.1: **Uncertainty Visualization** Rendered image and visualization of the uncertainty $(z_{\text{opaque}} - z_{\text{trans}})$ of each pixel through heat map from normal particles and abnormal particles on LLFF dataset [36].