# Expansive Supervision for Neural Radiance Field

Weixiang Zhang
SIGS, Tsinghua University
China
zhang-wx22@mails.tsinghua.edu.cn

Shuzhao Xie
SIGS, Tsinghua University
China
xsz24@mails.tsinghua.edu.cn

Shijia Ge
SIGS, Tsinghua University
China
gsj23@mails.tsinghua.edu.cn

Wei Yao
SIGS, Tsinghua University
China
w-yao22@mails.tsinghua.edu.cn

Chen Tang
The Chinese University of Hong Kong
China
chentang@link.cuhk.edu.hk

Zhi Wang
SIGS, Tsinghua University
China
wangzhi@sz.tsinghua.edu.cn

## ABSTRACT

Neural Radiance Fields have achieved success in creating powerful 3D media representations with their exceptional reconstruction capabilities. However, the computational demands of volume rendering pose significant challenges during model training. Existing acceleration techniques often involve redesigning the model architecture, leading to limitations in compatibility across different frameworks. Furthermore, these methods tend to overlook the substantial memory costs incurred. In response to these challenges, we introduce an expansive supervision mechanism that efficiently balances computational load, rendering quality and flexibility for neural radiance field training. This mechanism operates by selectively rendering a small but crucial subset of pixels and expanding their values to estimate the error across the entire area for each iteration. Compare to conventional supervision, our method effectively bypasses redundant rendering processes, resulting in notable reductions in both time and memory consumption. Experimental results demonstrate that integrating expansive supervision within existing state-of-the-art acceleration frameworks can achieve 69% memory savings and 42% time savings, with negligible compromise in visual quality.

## CCS CONCEPTS

• **Computing methodologies** → *Learning paradigms*; **Volumetric models**.

## 1 INTRODUCTION

Radiance field has emerged as a promising approach for representing 3D media content in the field of photorealistic novel view synthesis. Neural Radiance Fields (NeRFs) [30] employ a meticulously designed neural network $F(\Theta)$ to implicitly encode the scene. This neural network maps the position $\mathbf{x} = (x, y, z)$ and viewing direction $\mathbf{d} = (\theta, \varphi)$ to view-dependent color $\mathbf{c} = (r, g, b)$ and view-independent volumetric density $\tau$, i.e. $F(\Theta) : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$. With its powerful implicit neural scene representation, NeRFs leverage sampled query pairs (color $\mathbf{c}$ and density $\sigma$) along the ray for synthesizing and inferring the target pixel via volume rendering. As a consequence, NeRFs surpass traditional multi-view stereo methods in terms of visual quality. Despite the impressive performance of NeRF in novel view synthesis, the training speed of NeRFs remains a significant concern. In the original NeRF design, the rendering of each pixel (ray) requires sampling $N$ points to compute the color $\mathbf{c}$ and density $\sigma$. Considering a scene with dimensions $(h, w)$, this
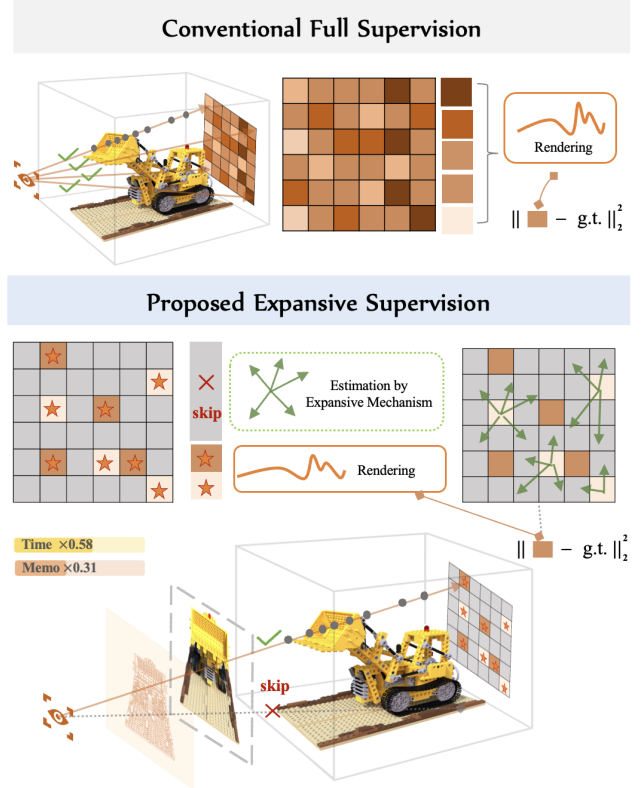


Figure 1: Overview of proposed method. Our approach adopts an expansive supervision technique to selectively render a subset of crucial pixels to estimate the error by expansive mechanism. Unlike conventional full supervision, which blindly renders all pixels, our method intelligently avoids redundant rendering processes, leading to significant reductions in training time and memory consumption.

process requires $h \cdot w \cdot N$ neural network forward passes, which can amount to over $10^6$ computations for rendering a view with a 1080p resolution. This computational burden substantially prolongs the training duration and negatively impacts the generation of this novel form of 3D media content.

Existing approaches for training acceleration predominantly rely on caching view-independent features using explicit representations, such as voxels [45], tensors [6] and hash table [33]. While

trading space for time achieves significant time savings, these acceleration methods suffer from compatibility limitations, as they are tailored to particular model architectures. With the emergence of more novel NeRF frameworks, the incompatibility issue of these existing methods has become increasingly evident. Furthermore, the memory cost associated with these acceleration techniques has often been overlooked, hindering the adaptation of the training process to more resource-limited devices.

To overcome these challenges, we introduce an expansive supervision mechanism for neural radiance field training. Our method is motivated by the observation that the distribution of training error exhibits long-tail characteristic and is highly consistent with the image content. During training, we selectively render a small but crucial subset of pixels $R' \subset R$ with image content prior $I$, and expand the error of these precisely rendered pixels to estimate the loss for the entire area in each iteration. By avoiding costly yet marginal renderings, our method can theoretically achieve $(1 - \frac{|R'|}{|R|})v\times$ time savings, where $v \in (0, 1)$ represents the proportion of rendering costs within the total training process. In our experiment, we can achieve 0.69× memory and 0.42× time savings by rendering only 30% of pixels to supervise the entire model.

In this paper, we observe that the long-tail distribution of training errors exhibits a strong correlation with the image content. As depicted in Figure 2 (column 3), the error map allows us to easily identify the image content. Moreover, regions with higher frequency display larger errors, while smoother areas exhibit smaller errors. Hence, leveraging image context to selectively omit a significant portion of the rendering process can effectively achieve substantial resource savings while maintaining rendering quality.

However, current NeRF training paradigm disrupts the connection between in-batch error and image content due to indiscriminate shuffling of training data. Simply removing the data shuffler can significantly compromise rendering quality due to the reduced entropy of the order-preserved training data. To address this, we propose a content-aware permutation that achieves maximum entropy within the constraints of expansive supervision. The effectiveness of permutation has been validated through theoretical analysis and empirical experiments.

With content-aware permutation, we satisfy the prerequisites for expansive supervision while preserving model performance. The selected set of pixels $R'$ comprises two areas in the batch $B$: the anchor area $A$ and the source area $S$. The anchor area are computed by the light-wight edge detector to displays prominent error patterns. And source area are sampled to expand its values to the reaming area. The final error estimate $\hat{L}$ is synthesized from both precise renderings $(A \cup S)$ and expanded estimation $(B \backslash (A \cup S))$. Subsequently, the model parameters are updated by $\Theta := \Theta - \eta \nabla \hat{L}$.

Extensive experiments have been conducted to validate the effectiveness of our method. In comparison to conventional full supervision, our expansive supervision approach achieves substantial time savings and memory usage reduction while maintaining rendering quality at a negligible loss. Importantly, our method exhibits unmatched compatibility with existing acceleration techniques, requiring no custom modifications for adaptation. Additionally, the

incorporation of content-aware permutation enriches the contextual information during loss computation, opening up possibilities for the development of more advanced loss functions.

Our contributions can be summarised as followed:

- We are the first to observe a strong correlation between error distribution and image content. To leverage this observation for accelerating NeRFs training, we introduce content-aware permutation to establish this connection while ensuring maximum model performance.
- We propose expansive supervision, a method that selectively renders a small yet crucial subset of pixels. By expanding the error values of these pixels, we estimate the overall loss in each iteration. This method effectively saves considerable time and memory during training by bypassing a significant number of redundant renderings.
- We conduct comprehensive experiments to validate the effectiveness of our method in both controlled test environments and real-world scenarios. Additionally, we analyze the trade-off between cost and quality.

## 2 RELATED WORK

### 2.1 Neural Radiance Field

Neural radiance field [30] has revolutionized the field of 3D computer vision by leveraging multilayer perceptrons (MLPs) to implicitly represent the radiance field. Its outstanding performance in 3D reconstruction and novel view synthesis has inspired a plethora of research [1, 2, 21, 22, 27, 57]. In terms of rendering quality, several works have focused on addressing aliasing artifacts and improving multi-scale representation. Examples of such works include Mip-NeRF [1], Zip-NeRF [3], and Tri-MipRF [14]. For the adaption of NeRF to unbounded 360° scenes, several methods have implemented foreground-background separation and non-linear scene parameterization to optimize the realism of the scenes (e.g. NeRF++ [58], Mip-NeRF360 [2]). As a novel representation of 3D media content, NeRF has also sparked significant interest in various downstream applications, including segmentation [19, 25, 42], editing [13, 20, 31, 49, 54], and generation [16, 24, 48, 51, 52].

### 2.2 Efficient Training for NeRF

NeRFs have been challenged by the time-consuming training and rendering, primarily due to the intensive computation involved in volume rendering. Although recent advancements have enabled real-time rendering of NeRFs on mobile devices [5], the training process still demands a significant amount of time and effort. In order to accelerate training, various methods have been proposed. One of effective way is storing view-independent features in an explicit representation, which trades speed for space. Efficient explicit representations include octrees [23, 56], point cloud [53], voxel grids [10, 45], low rank tensors [6] and hash tables [33]. Another line of research focuses on employing decomposition schemes to reduce latency, such as DoNeRF [34] and KiloNeRF [39].

The key distinction between existing acceleration techniques [6, 8, 11, 12, 33, 38, 40] and our method lies in the elimination of costly renderings. While existing methods primarily focus on reducing computation for each rendering pixel, our method addresses the issue from a supervisory perspective by reducing the number of
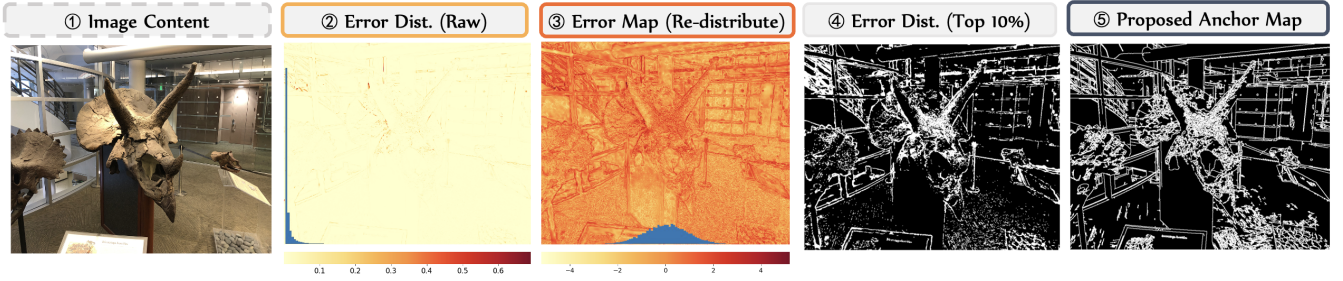
**Figure 2: Preliminary study for our observation. (#2) The blue histogram illustrates the distribution of errors after 1000 iterations, highlighting a pronounced long-tail characteristic. (#3) To enhance the discernibly of error data, we have transformed the data into a normal distribution, revealing the relationship between the redistributed error value and image content. (#4) The top 10% of errors identified during training are visualized, corresponding to regions with high-frequency details in the image content. (#5) The top10% error map generated by our expansive supervision exhibits a high correlation with the actual error distribution.**

rendering pixels required. To the best of our knowledge, our method is the first to achieve NeRF training acceleration through partial supervision.

## 3 METHODS

### 3.1 Problem Formulation

Neural Radiance Fields (NeRFs) learns a function $F(\Theta) : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$ using a multilayer perceptron (MLP), where $\mathbf{x} \in \mathbb{R}^3, \mathbf{d} \in \mathbb{R}^2$ represent the position and view direction of a point, while $\mathbf{c} \in \mathbb{R}^3$ and $\sigma \in \mathbb{R}$ represent the emitted color and density, respectively. Volume rendering allows computing the expected color $\mathbf{C}(\mathbf{r})$ in a novel view as:

$$\mathbf{C}(\mathbf{r}) = \int_0^t \mathcal{T}(t; \mathbf{r}) \cdot \tau(\mathbf{r}(t)) \cdot \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \tag{1}$$

where $\mathcal{T}(t; \mathbf{r}) = \exp\left(-\int_0^t \tau(\mathbf{r}(s)) ds\right)$ represents the accumulated transmittance along the ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$. In practice, numerical estimation of the rendering integral involves sampling $N$ points from partitioned bins along the ray, allowing the estimation of $\mathbf{C}(\mathbf{r})$ as:

$$\hat{\mathbf{C}}(\mathbf{r}) = \sum_{i=1}^N \mathcal{T}_i \left(1 - \exp\left(-\tau_i \delta_i\right)\right) \mathbf{c}_i, \tag{2}$$

where $\mathcal{T}_i = \exp\left(-\sum_{j=1}^{i-1} \tau_j \delta_j\right)$, and $\delta_i$ represents the distance between adjacent sampled points.

The training of radiance fields is computationally intensive due to the large number of neural network forward passes required, which amounts to $N\times$ batch_size for each iteration. Our proposed method address this issue by selectively rendering a subset of rays $R' \in R$ and utilizing image context $\mathcal{I}$ to estimate the error $L(R)$ with the expansive mechanism based on partial significant pixels $\hat{L}(R')$.

The implementation of expansive supervision requires a correlation between batch errors and corresponding image content. However, this link is disrupted by the arbitrary shuffling characteristic of standard NeRF training. To address this, we introduce a constraint for expansive supervision: pixels within the same batch

must derive from identical input views. Mathematically, this constraint is articulated as $C : B \cap I = B, \forall B \in \mathcal{B}, \exists I \in \mathcal{I}$, where $\mathcal{B}$ denotes the set of batch $B$, and $\mathcal{I}$ signifies the set of image $I$.

A straightforward solution is to align the data in a strict sequential order in image, as shown in left part of Figure 3. However, this approach results in a significant decrease in model performance due to the reduced entropy of the training data. This reduction in entropy negatively impacts the learning performance during each iteration [28]. Therefore, it becomes necessary to employ a permutation algorithm that maximizes the entropy of the training data while satisfying the constraint of expansive supervision.

This problem can be formulated as finding a permutation $P^* : \mathcal{D} \rightarrow \mathcal{B}$, which can be expressed as follows:

$$P^* = \arg\max_P H(P(\mathcal{D}))$$
$$\text{s.t. } C : B \cap I = B, \forall B \in \mathcal{B}, \exists I \in \mathcal{I}, \tag{3}$$

where $\mathcal{D} = g(\mathcal{B}) = g(\mathcal{I})$ and $g(\cdot)$ denotes a reshape function that maps a multi-dimensional set to a one-dimensional set while preserving the element order. $H(\cdot)$ represents entropy calculation. To solve this problem, we propose a content-aware ray shuffler, which is further elaborated in Section 3.2.

Once the constraint $C$ is satisfied, we can proceed with the implementing expansive supervision. The problem can be formulated as follows: given the shuffled batch set $\mathcal{B} = P^*(\mathcal{D})$ and the image set $\mathcal{I}$, our objective is to design a supervision mechanism that trains the model rendering only a subset of pixels $R' \subset R$. This mechanism allows for the conservation of both time and memory, resulting in a savings of $(1 - \frac{|R'|}{|R|})v\times$ computational resources, where $v \in (0, 1)$ represents the ratio of resources used by volume rendering in the total training process. Further details can be found in Section 3.3.

### 3.2 Content-aware Permutation

The permutation that satisfies the constraint $C$ can be defined as $\hat{P}$. Our object is to find a permutation $P^* = \arg\max H(\hat{P}(\mathcal{D}))$. To ensure $C$: $\exists I \in \mathcal{I} \rightarrow B \cap I = B$ for $\forall B \in \mathcal{B}$, we partition $\hat{P}(\mathcal{D})$ into $\hat{P}_{\text{intra}}^{\mathcal{I}}(B)$ and $\hat{P}_{\text{inter}}^{\mathcal{B}}(\mathcal{D})$. Here, $\hat{P}_{\text{intra}}^{\mathcal{I}}(B)$ represents the intra-batch permutation from the same input view, and $P_{\text{inter}}^{\mathcal{B}}(\mathcal{D})$ represents the inter-batch permutation. Consequently, the entropy of $P^*(\mathcal{D})$
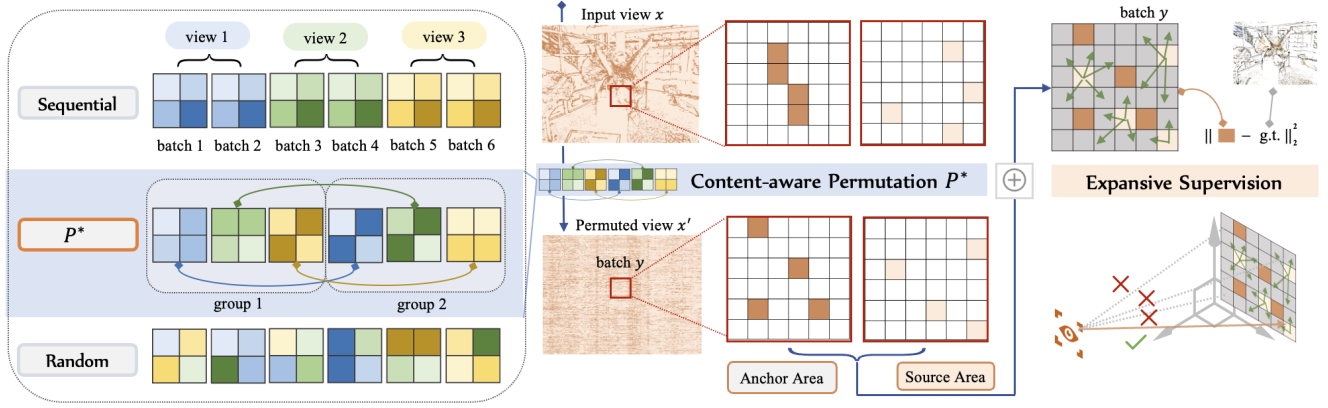
**Figure 3: Pipeline of expansive supervision. The training process begins with the application of content-aware permutation to ensure that the data within the same batch originate from the same view. Subsequently, we exclusively render the crucial pixels, which consist of the pre-computed anchor area and sampled source areas, to estimate the loss. This estimation is accomplished through the expansive strategy described in Section 3.3. Our expansive supervision method results in significant time and memory savings while maintaining negligible compromise in visual quality.**

can be expressed as follows:

$$
\begin{aligned}
H(P^*(\mathcal{D})) &= H(P^{\mathcal{I}}_{intra}(B), P^{\mathcal{B}}_{inter}(\mathcal{D})) \\
&= H(P^{\mathcal{I}}_{intra}(B)) + H(P^{\mathcal{B}}_{inter}(\mathcal{D})|P^{\mathcal{I}}_{intra}(B)).
\end{aligned}
\tag{4}
$$

Given that $P^{\mathcal{B}}_{inter}(\mathcal{D})$ is independent of the specific permutation for $\forall B \in \mathcal{B}$, i.e. $P^{\mathcal{I}}_{intra}(\cdot) \perp P^{\mathcal{B}}_{inter}(\cdot)$, we can maximize each component of $H(P^*(\mathcal{D}))$ separately based on Equation 4:

$$
\begin{aligned}
\max\{H(P^*(\mathcal{D}))\} = &\max\{H(P^{\mathcal{I}}_{intra}(B))\} \\
&+ \max\{H(P^{\mathcal{B}}_{inter}(\mathcal{D}))\}.
\end{aligned}
\tag{5}
$$

The entropy of intra-batch permutation $H(P^{\mathcal{I}}_{intra}(B))$ for $\forall B \in \mathcal{B}$ can be represent as:

$$
H(P^{\mathcal{I}}_{intra}(B)) = - \sum_{i=0}^{\sqrt{|B|}-1} \sum_{j=0}^{\sqrt{|B|}-1} p(i,j) \log p(i,j),
\tag{6}
$$

and

$$
p(i,j) = \frac{1}{MN} \sum_{m=0}^{M} \sum_{n=0}^{N} \mathcal{P}(b_{i,j} = \hat{b}_{m,n}),
\tag{7}
$$

where $p(i,j) \in [0,1]$ denotes that the predictability of entry $b_{i,j} \in B$. $\hat{b} \in \hat{B}$ denotes the element in the batch in natural sequence of image, as Sequential permutation shown in the left part of Figure 3. $\mathcal{P}$ represents probability measure.

According to the principle of maximum entropy [15], maximum entropy is attained when the distribution is uniform, i.e. $p(i,j) = \frac{1}{|\mathcal{B}|}$. This can be achieved by employing a uniformly random permutation, which is self-evident. Hence we have $P^{\mathcal{I}}_{intra}(B) := P_{random}(B)$.

For inter-batch permutation $H(P^{\mathcal{B}}_{inter}(\mathcal{D}))$, we consider the inherent correlation of the pixels in a image and group the batches by input views $I$. Based on the assumption that the correlation coefficient [18] between input views are negligible ($\rho(I_{inter}) \ll$

$\rho(I_{intra})$), we can represent the entropy of inter-batch permutation as follows:

$$
\begin{aligned}
H(P^{\mathcal{B}}_{inter}(\mathcal{D})) &= - \sum_{k=0}^{|\mathcal{I}|-1} \sum_{l=0}^{\lceil \frac{|\mathcal{D}|}{|B||\mathcal{I}|} \rceil - 1} p(k,l) \log p(k,l) \\
&= - |\mathcal{I}| \sum_{l=0}^{\lceil \frac{|\mathcal{D}|}{|B||\mathcal{I}|} \rceil - 1} p(l) \log p(l).
\end{aligned}
\tag{8}
$$

Similarly, it is evident that when $p(l) = 1/\lceil \frac{|\mathcal{D}|}{|B||\mathcal{I}|} \rceil$, $H(P^{\mathcal{B}}_{inter}(\mathcal{D}))$ achieve its maximum. The maximum can be attained by employing a uniformly random permutation within a content-aware group. Each group consists of $|\mathcal{I}|$ batches randomly selecting from different input views. It is important to note that each view can only be selected once within the same group, as illustrated in Figure 3 (left part). This process is repeated until all of training data is involved.

In summary, our content-aware permutation scheme achieves the maximum entropy of training batches while ensuring that all data in a batch are from the same input views. Experimental results demonstrate that our permutation scheme achieves comparable training performance to pure random permutation, as discussed in Section 4.3.

## 3.3 Expansive Supervision

The objective of expansive supervision is to utilize only a small subset of supervision to guide the training of the radiance field, while minimizing any degradation in rendering quality and achieving significant time and memory savings.

Our design is based on the observation that the error distribution exhibits strong correlations with image content. Specifically, areas of high frequency in the image are expected to be more challenging to train and exhibit larger errors compared to other areas. Furthermore, similar patterns in the image should demonstrate similar error distribution.

The observation has been demonstrated through a preliminary study, as illustrated in Figure 2. We observed that the error generated during standard training exhibits a clear long-tail phenomenon, where approximately 99.4% of the data with the least errors only contribute to around 10% of the overall importance. Furthermore, our observations indicate that higher error values are predominantly concentrated in the high-frequency regions of the image, corresponding to edges and areas rich in texture. This pattern is clearly illustrated in column 3 of Figure 2. Based on above validation, the error distribution can be estimated by part of rendering and the global loss could be expansively calculated.

The pipeline of expansive supervision is demonstrated in Figure 3. The training process begins with the application of content-aware permutation to ensure that the data within the same batch originate from the same view. Subsequently, we exclusively render the crucial pixels to estimate the loss. The selected set of pixels $R'$ are obtained from two distinct areas within the given input view $I$:

**Anchor area** $A \subset I$. We define $A$ as the area within $I$ where patterns exhibit larger errors. This is computed using the anchor extractor function $\mathcal{F}_A(\cdot)$. In other words, we have $A = \mathcal{F}_A(I, \beta_A)$, where $\beta_A$ controls the size of the anchor area. The cardinality of $A$ is determined by $\beta_A$, such that $|A| = \beta_A|I|$.

**Source area** $S \subset I\backslash A$. We define $S$ as the leftover area after excluding the anchor set. The source set is composed of sampled points, and the error is estimated based on these source points, which expand to cover all remaining areas. Similarly, we have $|S| = \beta_S|I|$.

We define $B^*$ as the batch data after content-aware permutation $P^*$, and $A^*$ as the corresponding anchor area of $B^*$. The source are $S$ can then be randomly sampled from $B^*\backslash A^*$. The global estimated error can be represented as follow:

$$\hat{L} = \frac{1}{|A^*|} \sum_{r_A \in A^*} ||\hat{C}(r_A) - C(r_A)||_2^2 + \frac{1}{|S|}(\frac{1}{\beta_A + \beta_S} - 1) \sum_{r_S \in S} ||\hat{C}(r_S) - C(r_S)||_2^2, \quad (9)$$

where $C(\cdot)$, $\hat{C}(\cdot)$ denotes ground truth and predicted RGB colors from given rays. At the end of the iteration, thee radiance field parameter $\Theta$ would then be updated by $\Theta := \Theta - \eta\nabla\hat{L}$.

Compared to state-of-the-art full supervision, expansive supervision only renders a subset of rays $A^* \cup S^*$ to guide model learning process. This selective rendering theoretically saves $(1 - \beta_A - \beta_S)v\times$ computation resources.

**Details of anchor area extractor.** we provide a detailed description of the anchor extractor function $\mathcal{F}_A(\cdot)$ design here. Given that the objective of expansive supervision is to accelerate training, we chose lightweight canny [4] as the basic extractor. To ensure that each anchor area has the same intensity of $\beta_A|I|$, we designed a simple progressive adjuster for the threshold of edge detector. For iteration $i$, the threshold $T_i$ is updated with $T_i = 1 + \mu(\mathcal{E}(T_{i-1}) - \beta_A|I|)$, where $\mathcal{E}(T_{i-1})$ is the sum of the output edge map generated by the edge detector with threshold $T_{i-1}$. The step rate is denoted as $\mu$. The iterative process continues until the condition $0.8 \leqslant \frac{\mathcal{E}(T_i)}{\beta_A|I|} \leqslant 1.2$ is satisfied.

---

**Algorithm 1** Expansive Supervision Training

---

**Input:** Input View Set $\mathcal{I}$, $\beta_A$, $\beta_S$, $\mathcal{A} = \varnothing$
**for** $I$ in $\mathcal{I}$ **do**
    $A := \mathcal{F}_A(I, \beta_A)$
    $\mathcal{A} := \mathcal{A} \cup \{A\}$
**end for**
**Content-aware Permutation** [Section 3.2]
$\mathcal{B}^* := P^*(g(\mathcal{I})), \mathcal{A}^* := P^*(g(\mathcal{A}))$
**repeat**
    **for** $B^*$ **in** $\mathcal{B}^*$ **do**
        $S := $ Random Sampling$(B^*\backslash A^*, \beta_S)$
        $\hat{L}_A := \frac{1}{|A^*|} \sum_{r_A \in A^*} ||\hat{C}(r_A) - C(r_A)||_2^2$
        $\hat{L}_S := \frac{1}{|S|}(\frac{1}{\beta_A + \beta_S} - 1) \sum_{r_S \in S^*} ||\hat{C}(r_S) - C(r_S)||_2^2$
        $\hat{L} := \hat{L}_A + \hat{L}_S$ [Eq.9]
        $\Theta := \Theta - \eta\nabla\hat{L}$
    **end for**
**until** Training End

---

## 4 EXPERIMENTS

### 4.1 Implementation Details

To demonstrate the compatibility of our method with the state-of-the-art NeRF acceleration framework, we utilized TensoRF as the underlying backbone model for our experiments. The parameter settings were aligned with the default configuration in [6]. We use $\beta = \beta_A + \beta_S = \frac{|R'|}{|R|}$ to represent the ratio of pixel to supervise and $\beta_A = \beta_S$ was set empirically. We conducted a total of 30,000 iterations using a batch size of 4,096. The step rate for the anchor area extractor was set to 15.

We adopted the following metrics for evaluation: PSNR [9], SSIM [50], and LPIPS [59]. Specifically, we utilized L(A) and L(V) to represent the VGG [43] and AlexNet [17] versions of LPIPS, respectively. Our extensive experiments encompassed Synthetic-NeRF [30] dataset, which consisted of synthetic data, as well as real-world forward-facing scenes from the LLFF [29]. In the absence of specific explanations, we selected the Synthetic-NeRF dataset for analysis purposes.

Our training were conducted on four GPUs equipped with NVIDIA RTX 3090(24.58GB VRAM) on Pytorch Framework[36]. We conducted our experiments in an ideal test environment as well as real-world scenarios to measure computational resources. The CPU used was AMD EPYC 7542, and the RAM capacity was 512 GB. More details can be found in Section 4.4.
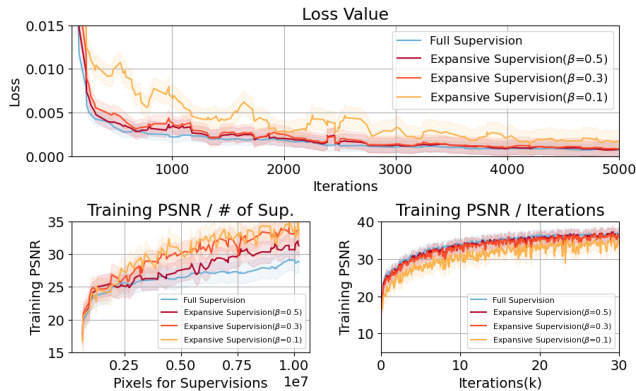
### 4.2 Comparison of Different Supervision Mechanisms

As illustrated in Table 1, we conducted a comparison between our method and various supervision mechanisms. In the table header, **Sup.Rays** represents the number of rendered rays used for supervision in each iteration. Traditional training techniques employ full supervision with a batch size $|B|$, whereas expansive supervision only requires $\beta|B|$ rendered rays. The indicator Sup.Rays serves as a measure of algorithm efficiency, which can be controlled by adjusting the batch size or the parameter $\beta$.

**Table 1: Quantitative comparison of different supervisions. Expansive Sup. † denotes the default version of Expansive Supervision, which strikes a balance between rendering quality and computational savings.**

| Methods | Sup.Rays | $\beta$ | Batch Size ↓ | Memory ↓ | Time ↓ | Synthetic-NeRF | | | | LLFF | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | PSNR (dB)↑ | SSIM ↑ | L(A) ↓ | L(V) ↓ | PSNR (dB)↑ | SSIM ↑ | L(A) ↓ | L(V) ↓ |
| Full Sup. | 4096 | - | 4096 | 1× | 1× | 32.73 | 0.961 | 0.030 | 0.051 | 26.68 | 0.835 | 0.113 | 0.201 |
| Full Sup. | 2048 | - | 2048 | 0.46× | 0.69× | 32.08 | 0.956 | 0.036 | 0.059 | 26.64 | 0.833 | 0.126 | 0.212 |
| Expansive Sup. | 2048 | 0.5 | 4096 | 0.46× | 0.69× | **32.36** | 0.959 | 0.033 | 0.056 | **26.68** | 0.827 | 0.124 | 0.206 |
| Full Sup. | 1229 | - | 1229 | 0.31× | 0.58× | 31.36 | 0.951 | 0.043 | 0.066 | 26.29 | 0.826 | 0.129 | 0.218 |
| Random Sup.(30%) | 1229 | - | 4096 | 0.31× | 0.58× | 31.24 | 0.942 | 0.043 | 0.068 | 25.90 | 0.812 | 0.147 | 0.243 |
| **Expansive Sup. †** | 1229 | 0.3 | 4096 | 0.31× | 0.58× | **32.20** | 0.956 | 0.035 | 0.058 | **26.36** | 0.825 | 0.143 | 0.232 |
| Full Sup. | 410 | - | 410 | 0.10× | 0.67× | 29.45 | 0.933 | 0.066 | 0.093 | 25.67 | 0.802 | 0.190 | 0.266 |
| Random Sup. (10%) | 410 | - | 4096 | 0.10× | 0.67× | 29.40 | 0.93 | 0.065 | 0.093 | 25.58 | 0.799 | 0.192 | 0.270 |
| Expansive Sup. | 410 | 0.1 | 4096 | 0.10× | 0.67× | **30.51** | 0.940 | 0.053 | 0.081 | 25.47 | 0.786 | 0.208 | 0.292 |

To demonstrate the compatibility of our method with the state-of-the-art NeRF acceleration framework, we utilized TensoRF (VM-192) as the underlying backbone model for our experiments.



**Figure 4: Convergence performance of expansive supervision. Our method achieves precise error estimation comparable to full supervision(upper) and exhibits faster convergence as the number of supervised pixels increases(lower left).**

To demonstrate the superiority of our approach, we implemented various mechanisms to achieve an equivalent number of rendered rays per batch with comparable memory and time consumption. These mechanisms included reducing the batch size with full supervision and randomly selecting a subset of rays for rendering (represented as "Random Sup."). The effectiveness of our method are validated with both synthetic and real dataset. Our method with $\beta = 0.3$ achieved the highest rendering quality with comparable efficiency, thus we selected it as our default setting, which will be further detailed in Section 4.5. Compared to standard full supervision, our method only utilizes 0.31× of the memory and 0.58× of the training time to achieve nearly the same visual quality.

The visual quality comparison is depicted in Figure 5. Although our method achieved high efficiency with minimal impact on quantitative metrics, the visual quality remains indistinguishable. Thanks to our emphasis on high-frequency areas, our method demonstrates

**Table 2: Comparison of different permutations with full supervision.**

| | PSNR↑ | SSIM↑ | L(A)↓ | L(V)↓ | $C$ |
|---|---|---|---|---|---|
| Random | 32.73 | 0.961 | 0.013 | 0.026 | × |
| Sequential | 26.42 | 0.914 | 0.085 | 0.101 | √ |
| Intra(16) | 28.20 | 0.932 | 0.062 | 0.086 | √ |
| Intra(4) | 28.41 | 0.935 | 0.056 | 0.080 | √ |
| Intra(1) | 28.35 | 0.936 | 0.052 | 0.074 | √ |
| Intra(1)+inter | 28.38 | 0.936 | 0.052 | 0.074 | √ |
| $P^*$ (proposed) | 32.55 | 0.956 | 0.031 | 0.052 | √ |

superior rendering quality in terms of details compared to full supervision with small batch size. Furthermore, the convergence performance of our method is visualized in Figure 4, confirming the effectiveness of expansive supervision. Our method with $\beta = 0.3$ and $\beta = 0.5$ achieves precise error estimation comparable to full supervision (as shown in the upper and lower right sub-figures). Additionally, the expansive supervision exhibits faster convergence as the number of supervised pixels increases.

## 4.3 Comparison of Different Permutations

To validate the effectiveness of content-aware permutation, we conducted a comparative analysis of different permutation methods and their impact on model performance, as shown in Table 2. To eliminate the effects of permutation from those of expansive supervision, we utilized standard full supervision across all evaluations. These experiments were conducted on the Synthetic-NeRF dataset.

We established a set of permutation methods, each satisfying the prerequisite for implementing expansive supervision, to benchmark against our approach. The **Sequential** method does not involve any permutation and adheres strictly to the natural order of images. **Intra(n)** denotes a random shuffle within each batch based on a sequential permutation, where $n$ signifies the size of the patch ($n \times n$)
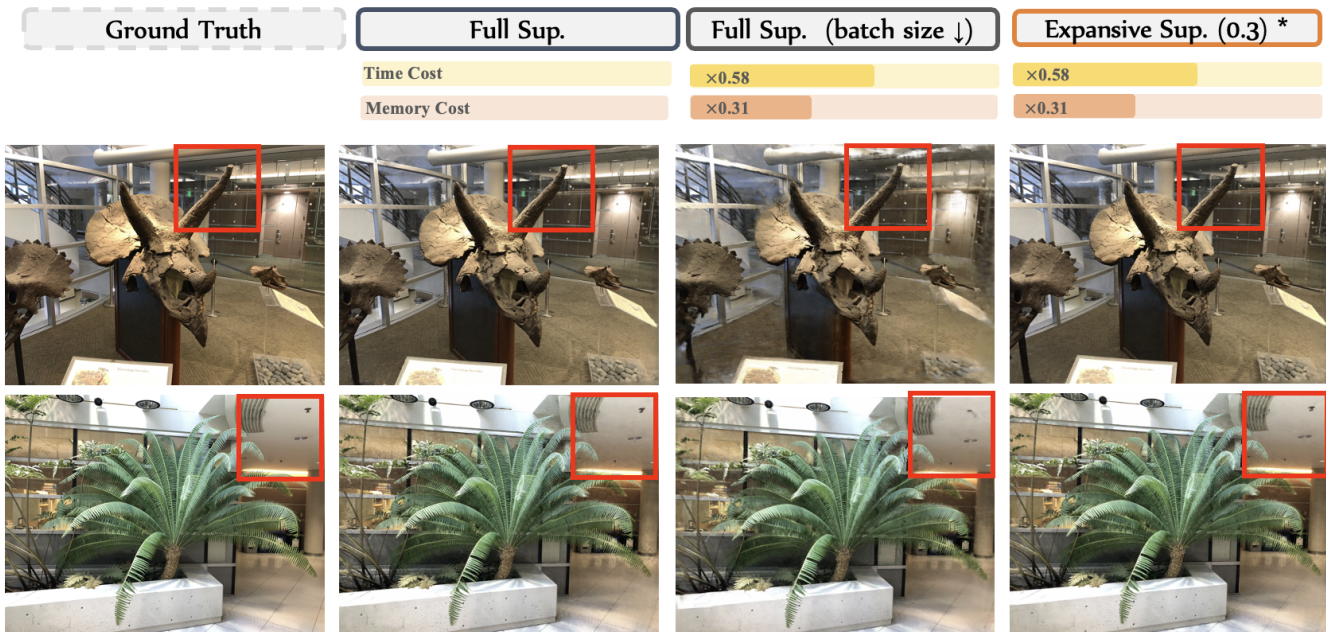
**Figure 5: Visual quality comparison with standard full supervision. In contrast to full supervision, expansive supervision exhibits no noticeable artifacts effectively reducing training time and memory usage. Under the same constrained computational resources, expansive supervision demonstrates higher quality reconstruction compared to full supervision.**
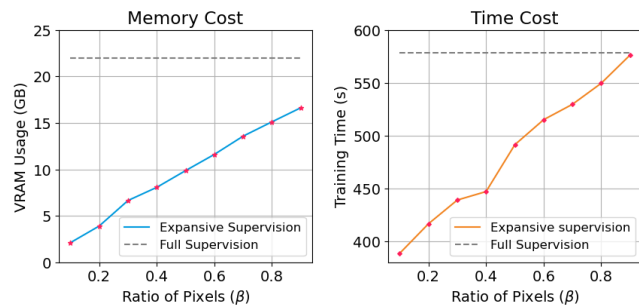


**Figure 6: Memory and training time cost of expansive supervision.**

whose order is preserved. As the value of $n$ increases, the entropy of Intra($n$) is expected to decrease. **Intra(1)+Inter** introduces an additional layer of random permutation across batches.

Random permutation was used as a reference for performance comparison. The results in Table 2 suggest a positive correlation between the entropy of permutation and the rendering quality of the model. The content-aware permutation $P^*$ closely approaches the upper bound, with a marginal loss of only 0.3 in PSNR and approximately 0.01 in LPIPS.

## 4.4 Analysis of Resources Savings

In order to validate the theoretical resource savings outlined in Section 3.1, we assessed the practical memory and time savings

resulting from the implementation of expansive supervision. Directly measuring the training time for each case may not accurately reflect resource savings due to environmental variability, such as other processes and I/O operations. To ensure a fair comparison, we designed the experiment settings to ensure consistent running conditions. We conducted the measurements in both a test environment and real-world applications. In the test environment, we cleared all other processes on the server and conducted experiments one by one. We measured the computation cost with 10 different $\beta$ settings ranging from 0.1 to 1.0. The quantitative results are presented in Table 3, and the visualization of memory and time costs can be found in Figure 6. For real-world applications, we conducted all experiments simultaneously without any specific settings to simulate limited computational resources. The results are shown in Table 4.
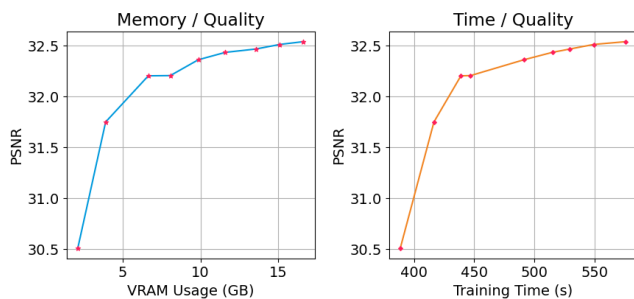
The result in Table 3 indicates that volume rendering accounts for the majority of training time, averaging between 55.86% and 83.64% across our experiments. Given the high cost of rendering pixels for supervision, reducing the number of rendering pixels can lead to significant resource savings. Specifically, by rendering only 30% of pixels for model supervision, we observed savings of 69% in memory and 25% in training time in the test environment. As shown in Table 4, the performance in terms of time savings is better in real-world scenarios with limited computational resources, where it can save nearly half of the training time. It can be concluded that expansive supervision achieves greater time savings when computational resources are limited.

**Table 3: Analysis of time/memory cost and rendering quality.**

| | Memory Cost ↓ (GB) | | Training Time (s) ↓ | | | | Rendering Quality | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | Total | Rendering | Backward | Others | PSNR ↑ | SSIM ↑ | L(A) ↓ | L(V) ↓ |
| Full Sup. | 21.51 | ×1.00 | 578.54 | 323.16 | 237.35 | 18.03 | 32.73 | 0.961 | 0.030 | 0.051 |
| Expansive Sup. $\beta = 0.9$ | 16.62 | ×0.77 | 576.17 | 319.70 | 235.29 | 21.18 | 32.53 | 0.959 | 0.031 | 0.053 |
| Expansive Sup. $\beta = 0.7$ | 13.57 | ×0.63 | 529.47 | 289.64 | 221.42 | 21.21 | 32.46 | 0.959 | 0.032 | 0.054 |
| Expansive Sup. $\beta = 0.5$ | 9.86 | ×0.46 | 491.32 | 250.40 | 218.03 | 22.89 | 32.36 | 0.958 | 0.033 | 0.056 |
| Expansive Sup. $\beta = 0.3$ | 6.64 | ×0.31 | 438.91 | 215.46 | 201.57 | 21.88 | 32.20 | 0.956 | 0.035 | 0.058 |
| Expansive Sup. $\beta = 0.1$ | 2.11 | ×0.10 | 388.46 | 177.01 | 190.13 | 21.32 | 30.51 | 0.940 | 0.053 | 0.081 |

**Table 4: Training time in test and real environment.**

| $\beta$ | Test Environment | | Real Environment | |
|---|---|---|---|---|
| | Total (s) | Rendering (s) | Total (s) | Rendering (s) |
| Full. | 578.54 (×1.00) | 323.16 | 2873.89 (×1.00) | 2403.25 |
| 0.5 | 491.32 (×0.85) | 250.40 | 1987.17 (×0.69) | 1622.27 |
| 0.3 | 438.91 (×0.76) | 215.46 | 1678.29 (×0.58) | 1334.48 |
| 0.1 | 388.46 (×0.67) | 177.01 | 1382.53 (×0.48) | 1171.61 |



**Figure 7: Memory and training time cost of expansive supervision. As the supervised pixel ratio increases, the margin of resource savings decreases.**

## 4.5 Analysis of Time-Quality Trade-off

Based on our measurements in Section 4.4, we investigated the impact of resource savings on model performance. The configurations of $\beta$ control the time and memory costs, and the resource-quality curves are depicted in Figure 7.

We observed that the optimal reduction of supervised pixels occurs at approximately $\beta = 0.3$, where $\frac{dC(\text{Time})}{d\beta}$ and $\frac{dC(\text{Memo.})}{d\beta}$ reach their maximum values. As the supervised pixel ratio increases beyond this point, the margin of resource savings decreases. Furthermore, when setting with a lower supervised pixel ratio($\beta = 0.1$), we observed noticeable artifacts that are not feasible in our mechanism. Therefore, we have determined that the default setting of $\beta = 0.3$ strikes a balance between rendering quality and computational savings.

## 4.6 Ablation Studies

To validate the effectiveness of each components within our proposed expansive supervision mechanism, we conducted a series of

**Table 5: Ablation experiments for expansive supervision.**

| | PSNR↑ | SSIM↑ | L(A)↓ | L(V)↓ |
|---|---|---|---|---|
| Baseline ($\beta = 0.3$) | 32.20 | 0.956 | 0.035 | 0.058 |
| w/o $P^*$ | 26.42 | 0.912 | 0.086 | 0.112 |
| w/o Anchor Sup. | 30.10 | 0.940 | 0.0058 | 0.083 |
| w/o Source Sup. | 31.53 | 0.949 | 0.043 | 0.068 |
| w/o Adjuster | 31.51 | 0.948 | 0.045 | 0.087 |
| w/ Pre-Recovery | 32.14 | 0.956 | 0.036 | 0.057 |
| w/ Post-Recovery | 32.18 | 0.956 | 0.036 | 0.058 |

ablation experiments. The result are presented in Table 5. Expansive supervision consists three primary components: content-aware permutation $P^*$, supervision from anchor set $A$ and source set $S$. We conducted ablation studies for each of these components. We also tested the effectiveness of adjuster coefficient in Equation 9.

We used the default expansive supervision as the baseline. Our findings revealed that omitting the content-aware permutation $P^*$ resulted in a significant performance degradation, with a decrease of 5.78 dB in PSNR. Removing the supervision from the anchor set $A$, the source set $S$, and the adjuster led to declines of 2.10 dB, 0.67 dB, and 0.79 dB in PSNR, respectively. Furthermore, it may seem intuitive to implement full supervision at the start or end stage of training to recover performance. However, our experiments showed that there was no improvement with these recovery approaches, and thus we eliminated them from our pipeline.

## 5 CONCLUSION

In this paper, we introduce an expansive mechanism designed to enhance the efficiency of neural radiance field training. Our approach is motivated by the observation that the long-tail distribution of training errors exhibits a strong correlation with the image content. To establish this correlation while preserving the maximal entropy of the data, we employ a content-aware permutation technique. By selectively rendering subsets of rays and leveraging the image context for expansive error estimation, our method achieves significant savings in both memory and training time. Compared to existing methods for NeRF training acceleration, our approach offers substantial savings in memory usage and unparalleled compatibility with minimal implementation effort.

# REFERENCES

[1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. 2021. Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields. In ICCV. IEEE, 5835–5844.

[2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2022. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. In CVPR. IEEE, 5460–5469.

[3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2023. Zip-NeRF: Anti-Aliased Grid-Based Neural Radiance Fields. In ICCV. IEEE, 19640–19648.

[4] John F. Canny. 1986. A Computational Approach to Edge Detection. IEEE Trans. Pattern Anal. Mach. Intell. 8, 6 (1986), 679–698.

[5] Junli Cao, Huan Wang, Pavlo Chemerys, Vladislav Shakhrai, Ju Hu, Yun Fu, Denys Makoviichuk, Sergey Tulyakov, and Jian Ren. 2023. Real-Time Neural Light Field on Mobile Devices. In CVPR. IEEE, 8328–8337.

[6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. 2022. TensoRF: Tensorial Radiance Fields. In ECCV (32) (Lecture Notes in Computer Science, Vol. 13692). Springer, 333–350.

[7] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, and Abhinav Shrivastava. 2021. Nerv: Neural representations for videos. Advances in Neural Information Processing Systems 34 (2021), 21557–21568.

[8] Zhiqin Chen, Thomas A. Funkhouser, Peter Hedman, and Andrea Tagliasacchi. 2023. MobileNeRF: Exploiting the Polygon Rasterization Pipeline for Efficient Neural Field Rendering on Mobile Architectures. In CVPR. IEEE, 16569–16578.

[9] Ahmet M. Eskicioglu and Paul S. Fisher. 1995. Image quality measures and their performance. IEEE Trans. Commun. 43, 12 (1995), 2959–2965.

[10] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. 2022. Plenoxels: Radiance Fields without Neural Networks. In CVPR. IEEE, 5491–5500.

[11] Quankai Gao, Qiangeng Xu, Hao Su, Ulrich Neumann, and Zexiang Xu. 2023. Strivec: Sparse Tri-Vector Radiance Fields. In ICCV. IEEE, 17523–17533.

[12] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien P. C. Valentin. 2021. FastNeRF: High-Fidelity Neural Rendering at 200FPS. In ICCV. IEEE, 14326–14335.

[13] Ayaan Haque, Matthew Tancik, Alexei A. Efros, Aleksander Holynski, and Angjoo Kanazawa. 2023. Instruct-NeRF2NeRF: Editing 3D Scenes with Instructions. In ICCV. IEEE, 19683–19693.

[14] Wenbo Hu, Yuling Wang, Lin Ma, Bangbang Yang, Lin Gao, Xiao Liu, and Yuewen Ma. 2023. Tri-MipRF: Tri-Mip Representation for Efficient Anti-Aliasing Neural Radiance Fields. In ICCV. IEEE, 19717–19726.

[15] Edwin T Jaynes. 1957. Information theory and statistical mechanics. Physical review 106, 4 (1957), 620.

[16] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. 2023. LERF: Language Embedded Radiance Fields. In ICCV. IEEE, 19672–19682.

[17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In NIPS. 1106–1114.

[18] Joseph Lee Rodgers and W Alan Nicewander. 1988. Thirteen ways to look at the correlation coefficient. The American Statistician 42, 1 (1988), 59–66.

[19] Hao Li, Dingwen Zhang, Yalun Dai, Nian Liu, Lechao Cheng, Jingfeng Li, Jingdong Wang, and Junwei Han. 2023. GP-NeRF: Generalized Perception NeRF for Context-Aware 3D Scene Understanding. CoRR abs/2311.11863 (2023).

[20] Shaoxu Li and Ye Pan. 2023. Interactive geometry editing of neural radiance fields. arXiv preprint arXiv:2303.11537 (2023).

[21] Zhihao Li, Kexue Fu, Haoran Wang, and Manning Wang. 2023. PI-NeRF: A Partial-Invertible Neural Radiance Fields for Pose Estimation. In ACM Multimedia. ACM, 7826–7836.

[22] Zhihao Li, Kexue Fu, Haoran Wang, and Manning Wang. 2023. PI-NeRF: A Partial-Invertible Neural Radiance Fields for Pose Estimation. In ACM Multimedia. ACM, 7826–7836.

[23] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. 2020. Neural Sparse Voxel Fields. In NeurIPS.

[24] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. 2023. One-2-3-45: Any Single Image to 3D Mesh in 45 Seconds without Per-Shape Optimization. In NeurIPS.

[25] Yichen Liu, Benran Hu, Chi-Keung Tang, and Yu-Wing Tai. 2023. SANeRF-HQ: Segment Anything for NeRF in High Quality. CoRR abs/2312.01531 (2023).

[26] Zhen Liu, Hao Zhu, Qi Zhang, Jingde Fu, Weibing Deng, Zhan Ma, Yanwen Guo, and Xun Cao. 2023. FINER: Flexible spectral-bias tuning in Implicit NEural Representation by Variable-periodic Activation Functions. arXiv preprint arXiv:2312.02434 (2023).

[27] Hongming Luo, Fei Zhou, Zehong Zhou, Kin-Man Lam, and Guoping Qiu. 2023. Restoration of Multiple Image Distortions using a Semi-dynamic Deep Neural Network. In ACM Multimedia. ACM, 7871–7880.

[28] Qi Meng, Wei Chen, Yue Wang, Zhi-Ming Ma, and Tie-Yan Liu. 2019. Convergence analysis of distributed stochastic gradient descent with shuffling. Neurocomputing 337 (2019), 46–57.

[29] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. 2019. Local light field fusion: practical view synthesis with prescriptive sampling guidelines. ACM Trans. Graph. 38, 4 (2019), 29:1–29:14.

[30] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In ECCV (1) (Lecture Notes in Computer Science, Vol. 12346). Springer, 405–421.

[31] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Konstantinos G. Derpanis, Jonathan Kelly, Marcus A. Brubaker, Igor Gilitschenski, and Alex Levinshtein. 2023. SPIn-NeRF: Multiview Segmentation and Perceptual Inpainting with Neural Radiance Fields. In CVPR. IEEE, 20669–20679.

[32] Thomas Müller. 2021. Tiny CUDA Neural Network Framework. https://github.com/nvlabs/tiny-cuda-nn.

[33] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant neural graphics primitives with a multiresolution hash encoding. ACM Trans. Graph. 41, 4 (2022), 102:1–102:15.

[34] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. 2021. DONeRF: Towards Real-Time Rendering of Compact Neural Radiance Fields using Depth Oracle Networks. Comput. Graph. Forum 40, 4 (2021), 45–59.

[35] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019. Deepsdf: Learning continuous signed distance functions for shape representation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 165–174.

[36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In NeurIPS. 8024–8035.

[37] Sameera Ramasinghe and Simon Lucey. 2022. Beyond periodicity: Towards a unifying framework for activations in coordinate-mlps. In European Conference on Computer Vision. Springer, 142–158.

[38] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. 2021. DeRF: Decomposed Radiance Fields. In CVPR. Computer Vision Foundation / IEEE, 14153–14161.

[39] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. 2021. KiloNeRF: Speeding up Neural Radiance Fields with Thousands of Tiny MLPs. In ICCV. IEEE, 14315–14325.

[40] Christian Reiser, Richard Szeliski, Dor Verbin, Pratul P. Srinivasan, Ben Mildenhall, Andreas Geiger, Jonathan T. Barron, and Peter Hedman. 2023. MERF: Memory-Efficient Radiance Fields for Real-time View Synthesis in Unbounded Scenes. ACM Trans. Graph. 42, 4 (2023), 89:1–89:12.

[41] Vishwanath Saragadam, Daniel LeJeune, Jasper Tan, Guha Balakrishnan, Ashok Veeraraghavan, and Richard G Baraniuk. 2023. Wire: Wavelet implicit neural representations. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 18507–18516.

[42] Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulò, Norman Müller, Matthias Nießner, Angela Dai, and Peter Kontschieder. 2023. Panoptic Lifting for 3D Scene Understanding with Neural Fields. In CVPR. IEEE, 9043–9052.

[43] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In ICLR.

[44] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. 2020. Implicit neural representations with periodic activation functions. Advances in neural information processing systems 33 (2020), 7462–7473.

[45] Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2022. Direct Voxel Grid Optimization: Super-fast Convergence for Radiance Fields Reconstruction. In CVPR. IEEE, 5449–5459.

[46] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. 2020. Fourier features let networks learn high frequency functions in low dimensional domains. Advances in neural information processing systems 33 (2020), 7537–7547.

[47] Jiaxiang Tang. 2022. Torch-ngp: a PyTorch implementation of instant-ngp. https://github.com/ashawkey/torch-ngp.

[48] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. 2023. Make-It-3D: High-Fidelity 3D Creation from A Single Image with Diffusion Prior. In ICCV. IEEE, 22762–22772.

[49] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. 2022. NeRF-Art: Text-Driven Neural Radiance Fields Stylization. CoRR abs/2212.08070 (2022).

[50] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. IEEE Trans. Image Process. 13, 4 (2004), 600–612.

[51] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. 2023. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation. In NeurIPS.

[52] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying Shan, Xiaohu Qie, and Shenghua Gao. 2023. Dream3D: Zero-Shot Text-to-3D Synthesis Using 3D Shape Prior and Text-to-Image Diffusion Models. In CVPR. IEEE, 20908–20918.

[53] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. 2022. Point-NeRF: Point-based Neural Radiance Fields. In CVPR. IEEE, 5428–5438.

[54] Zijiang Yang, Zhongwei Qiu, Chang Xu, and Dongmei Fu. 2023. MM-NeRF: Multimodal-Guided 3D Multi-Style Transfer of Neural Radiance Field. arXiv preprint arXiv:2309.13607 (2023).

[55] Lin Yen-Chen. 2020. NeRF-pytorch. https://github.com/yenchenlin/nerf-pytorch/.

[56] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. 2021. PlenOctrees for Real-time Rendering of Neural Radiance Fields. In ICCV. IEEE, 5732–5741.

[57] Junyi Zeng, Chong Bao, Rui Chen, Zilong Dong, Guofeng Zhang, Hujun Bao, and Zhaopeng Cui. 2023. Mirror-NeRF: Learning Neural Radiance Fields for Mirrors with Whitted-Style Ray Tracing. In ACM Multimedia. ACM, 4606–4615.

[58] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. 2020. NeRF++: Analyzing and Improving Neural Radiance Fields. CoRR abs/2010.07492 (2020).

[59] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In CVPR. Computer Vision Foundation / IEEE Computer Society, 586–595.

## A SUPPLEMENTARY EXPERIMENTS

### A.1 Compatibility with Current NeRF Variants

Expansive supervision is a plug-and-play method that seamlessly integrates with all learnable radiance field frameworks without requiring custom modifications. We successfully implemented our method on two widely-used NeRF acceleration frameworks, namely INGP [33] and TensoRF [6]. In order to assess the performance of our method, we conducted a comparative analysis with various NeRF variants [6, 30, 30, 45]. The results of this comparison can be found in Table 6. To ensure a fair comparison, we measured the memory cost by considering its maximum usage during training. Additionally, we recorded only the time taken for rendering, backward computation, and loss computation as the training time, eliminating the effects of data processing and validation. It is important to note that our time measurements were conducted in an ideal testing environment, as outlined in Section 4.4. For the implementation of these frameworks, we utilized the respective PyTorch versions [47, 55]. Our main focus throughout the comparison was solely on the algorithm design, thus we excluded general engineering accelerations such as model quantization and TCNN backend [32].

**Table 6: Comparison of current NeRF variants**

|  | Resources Cost ↓ | | Rendering Quality (PSNR ↑) | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
|  | Memo.(GB) | Time(s) | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship | Mean |
| Vanilla NeRF [30] | 4.54 | 19943.33 (~5.5 h) | 33.81 | 25.76 | 29.03 | 36.92 | 31.51 | 29.35 | 32.89 | 28.46 | 30.97 |
| DVGO [45] | 11.27 | 510.27 | 34.10 | 25.46 | 32.67 | 36.71 | 34.60 | 29.51 | 33.13 | 29.13 | 31.91 |
| TensoRF [6] | 21.51 | 578.54 | 34.96 | 25.88 | 33.44 | 37.13 | 36.03 | 29.79 | 34.16 | 30.47 | 32.73 |
| TensoRF + E.S. | 6.64 | 438.91 | 34.50 | 25.58 | 32.78 | 36.53 | 35.45 | 28.98 | 33.80 | 29.96 | 32.20 |
| INGP [33] | 8.78 | 340.15 | 33.12 | 25.17 | 30.95 | 35.37 | 33.36 | 28.02 | 33.75 | 28.89 | 31.08 |
| INGP + E.S. | 4.79 | 183.17 | 32.72 | 25.04 | 30.78 | 35.28 | 32.90 | 27.65 | 33.32 | 28.21 | 30.74 |

In our experiments, we set the number of iterations to 200,000 for vanilla NeRF and 30,000 for all other NeRF variants. Furthermore, we chose a value of $\beta = 0.3$ for the expansive supervision in both TensoRF and INGP backbones. Additionally, we made adjustments to the default settings [47] by setting the parameter "bound" to 1.1 for the ship and hotdog scenes, as we observed a performance drop when using the default settings. As shown in Table 6, the implementation of expansive supervision results in significant improvements in the memory and time efficiency of training NeRF models, while maintaining a minimal loss in reconstruction quality. Notably, our method can be seamlessly integrated into any NeRF framework without the need for custom modifications. When compared to TensoRF, our method demonstrates superior time efficiency, achieving a 46% reduction in training time when applied to the INGP framework, with only a minor decrease in performance (0.34dB in PSNR).

### A.2 Compatibility with Other Modality of Implicit Neural Representations

Implicit Neural Representations (INR) have gained significant popularity for their efficient memory usage and potential for various downstream tasks. INR leverages neural networks to parameterize signals through implicit continuous functions. It has shown remarkable progress in representing multimedia content, including images [44], videos [7], and 3D shapes [35]. NeRF, a special variant of INR, utilizes neural networks to parameterize the radiance field and implicitly encode 3D scenes. While our method is primarily designed to accelerate the training of NeRF, it also exhibits high compatibility with other forms of INR. Experimental results demonstrate that expansive supervision can be effectively extended to various forms of INR, offering significant time savings in the production of novel implicit media content representations.

We conducted experiments to evaluate the effectiveness of our method by selecting four state-of-the-art INR frameworks, namely SIREN [44], Gauss [37], WIRE [41], and the recent FINER [26], as backbones. These INR frameworks were used to fit 2D images and learn a function: $f : \mathbb{R}^2 \to \mathbb{R}^3$. The input to the function is the pixel location $(x, y)$, and the output is the corresponding pixel color $(r, g, b)$. For our experiments, we utilized natural datasets [46] with a resolution of 512 ×512. The parameters were kept consistent with FINER, and we set $\beta = 0.5$ for this particular experiment. The results are presented in Table 7.

The results presented in Table 7 demonstrate that our expansive supervision method can be seamlessly integrated into any INR framework, improving training efficiency without the need for custom settings. In the context of implicit image representation, our method exhibits even greater time savings compared to its implementation on NeRFs. By utilizing 50% of the pixels for supervision, we were able to achieve approximately a 47% reduction in training time. Furthermore, the performance degradation observed in terms of PSNR is minimal, with only a compromise of 0.19 dB. It is worth noting that in certain specific tests, our method even outperformed the baseline in PSNR (WIRE backbone in "Island" and "Mushroom"). In summary, the high flexibility and compatibility of expansive supervision allow its application to extend beyond NeRF and its variants. Other modalities of INR also hold great potential for its implementation. Our method effectively addresses the challenges posed by the high computational burden in INR training and contributes to the advancement of this novel multimedia representations.

### Table 7: Comparison of current INRs backbone

| | Resources Cost ↓ | | Rendering Quality (PSNR ↑) | | | | | | | | |
| | Memo.(GB) | Time(s) | Archway | Market | Island | Mushroom | Colosseo | Topview | Wolf | Seaside | Mean |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SIREN [44] | 3.59 | 178.16 | 33.72 | 40.03 | 37.53 | 40.21 | 38.43 | 34.45 | 38.82 | 38.49 | 37.71 |
| SIREN + E.S. | 3.07 | 88.95 | 32.63 | 39.20 | 36.36 | 39.50 | 37.11 | 33.26 | 37.38 | 37.77 | 36.65 |
| Gauss [37] | 3.83 | 244.61 | 31.44 | 35.01 | 34.77 | 35.50 | 35.89 | 31.78 | 35.59 | 35.15 | 34.39 |
| Gauss + E.S. | 3.20 | 123.42 | 30.96 | 34.99 | 35.07 | 35.70 | 35.31 | 31.65 | 34.90 | 34.98 | 34.20 |
| WIRE [41] | 3.59 | 563.07 | 28.72 | 31.10 | 29.92 | 30.75 | 32.85 | 28.94 | 31.65 | 30.07 | 30.50 |
| WIRE + E.S. | 3.21 | 279.73 | 27.92 | 30.88 | 29.71 | 30.03 | 31.58 | 28.60 | 31.16 | 28.81 | 29.84 |
| FINER [26] | 4.60 | 243.46 | 35.96 | 42.51 | 39.68 | 42.15 | 40.75 | 36.87 | 41.80 | 39.79 | 39.94 |
| FINER + E.S. | 4.09 | 154.32 | 35.25 | 41.67 | 38.80 | 41.51 | 39.72 | 35.93 | 40.83 | 39.25 | 39.12 |

## A.3 Details of Time/Memory Cost and Rendering Quality

As presented in Table 8, we provide detailed data corresponding to Table 3. This data forms the foundation for the analysis of the time-quality trade-off in Section 4.5. Figure 7 is plotted based on the information provided in this table.

These records adhere to the settings of the test environment. To ensure fairness, we cleared all other processes on the server and conducted the experiments sequentially. The computation cost was measured using 10 different $\beta$ settings ranging from 0.1 to 1.0. It is worth noting that the additional execution time of our method (mainly attributed to the pre-processed anchor area extractor in Section 3.3) is $1.91 \pm 0.2$s, which can be considered negligible compared to the total training time.

### Table 8: Analysis of time/memory cost and rendering quality.

| | Memory Cost ↓ (GB) | | Training Time (s) ↓ | | | | Rendering Quality | | | |
| | | | Total | Rendering | Backward | Others | PSNR ↑ | SSIM ↑ | L(A) ↓ | L(V) ↓ |
|---|---|---|---|---|---|---|---|---|---|---|
| Full Sup. | 21.51 | ×1.00 | 578.54 | 323.16 | 237.35 | 18.03 | 32.73 | 0.961 | 0.030 | 0.051 |
| Expansive Sup. $\beta = 0.9$ | 16.62 | ×0.77 | 576.17 | 319.70 | 235.29 | 21.18 | 32.53 | 0.959 | 0.031 | 0.053 |
| Expansive Sup. $\beta = 0.8$ | 15.11 | ×0.70 | 549.53 | 303.28 | 225.62 | 20.60 | 32.51 | 0.959 | 0.031 | 0.053 |
| Expansive Sup. $\beta = 0.7$ | 13.57 | ×0.63 | 529.47 | 289.64 | 221.42 | 21.21 | 32.46 | 0.959 | 0.032 | 0.054 |
| Expansive Sup. $\beta = 0.6$ | 11.61 | ×0.54 | 515.01 | 273.48 | 219.57 | 22.03 | 32.43 | 0.958 | 0.032 | 0.054 |
| Expansive Sup. $\beta = 0.5$ | 9.86 | ×0.46 | 491.32 | 250.40 | 218.03 | 22.89 | 32.36 | 0.958 | 0.033 | 0.056 |
| Expansive Sup. $\beta = 0.4$ | 8.06 | ×0.37 | 446.94 | 225.13s | 200.58 | 21.23 | 31.90 | 0.954 | 0.038 | 0.061 |
| Expansive Sup. $\beta = 0.3$ | 6.64 | ×0.31 | 438.91 | 215.46 | 201.57 | 21.88 | 32.20 | 0.956 | 0.035 | 0.058 |
| Expansive Sup. $\beta = 0.2$ | 3.90 | ×0.18 | 416.46 | 194.72 | 200.19 | 21.56 | 31.75 | 0.952 | 0.041 | 0.065 |
| Expansive Sup. $\beta = 0.1$ | 2.11 | ×0.10 | 388.46 | 177.01 | 190.13 | 21.32 | 30.51 | 0.940 | 0.053 | 0.081 |