

# Improving Robustness for Joint Optimization of Camera Poses and Decomposed LowRank Tensorial Radiance Fields

Bo-Yu Cheng Wei-Chen Chiu Yu-Lun Liu

Code available at : <https://github.com/Nemo1999/Joint-TensoRF>



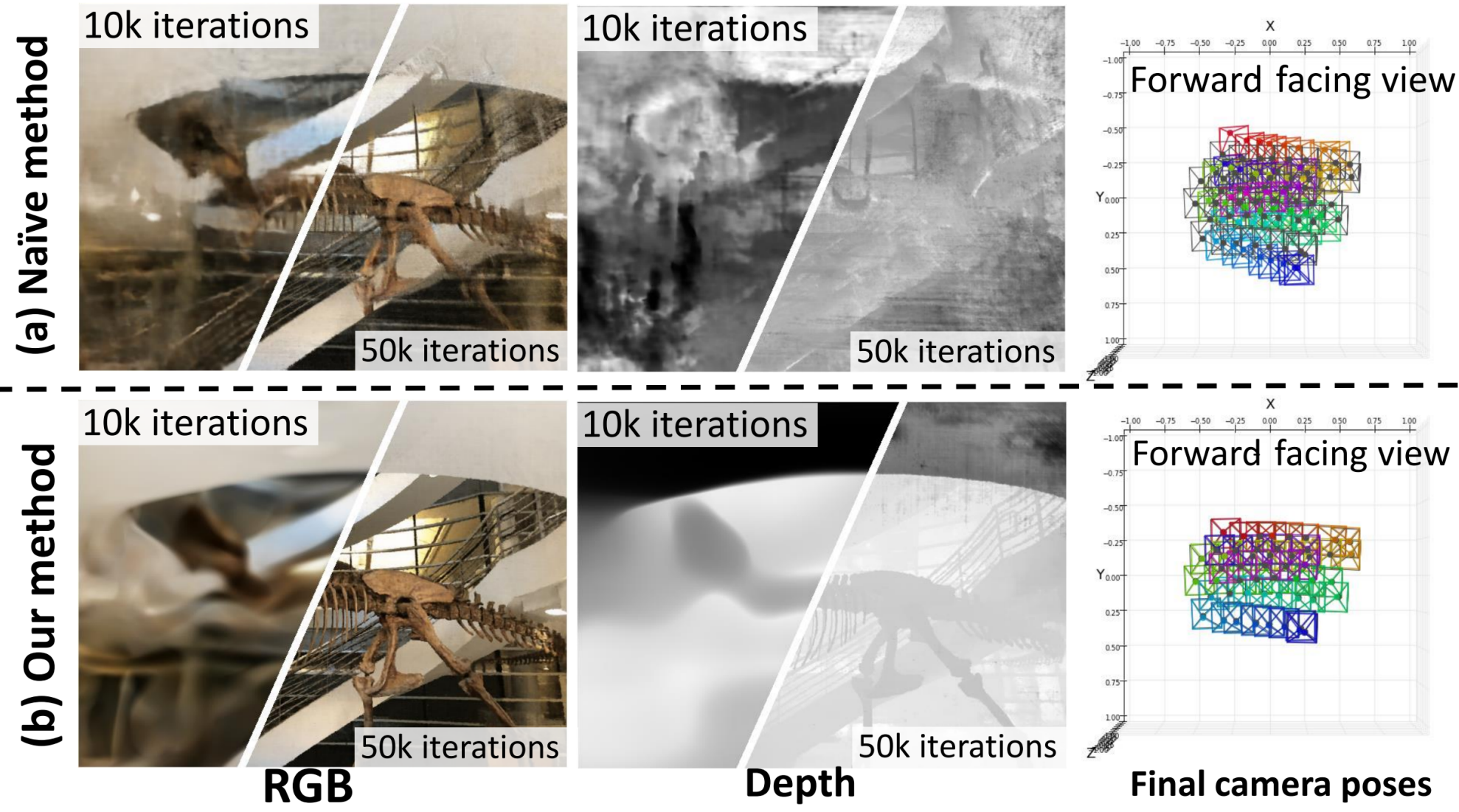
NATIONAL  
YANG MING CHIAO TUNG  
UNIVERSITY

## Introduction

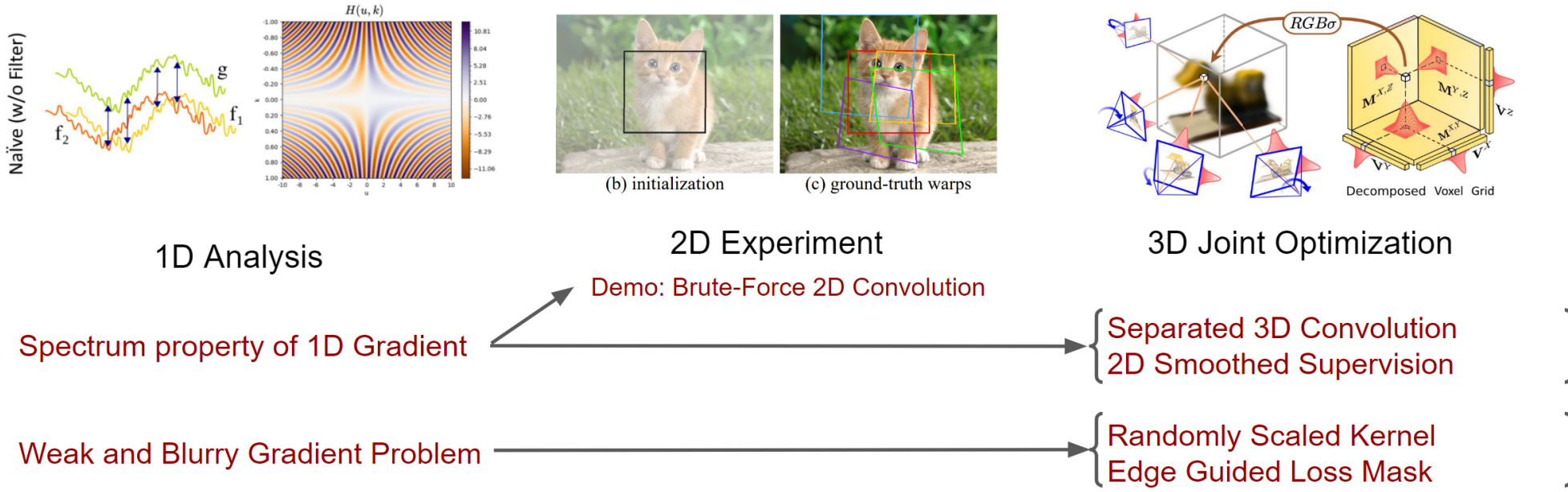
**Goal:** Following BARF, we enable joint optimization of camera pose on Tensorial Radiance Field, accelerating the joint optimization and getting better quality.

**Challenge:** Unlike MLP architecture used in BARF, voxel-based architectures lack spectral bias and is unstable in joint optimization.

**Contribution:** We solve the overfitting problem of naive method, and enable joint optimization of camera pose on TensoRF.

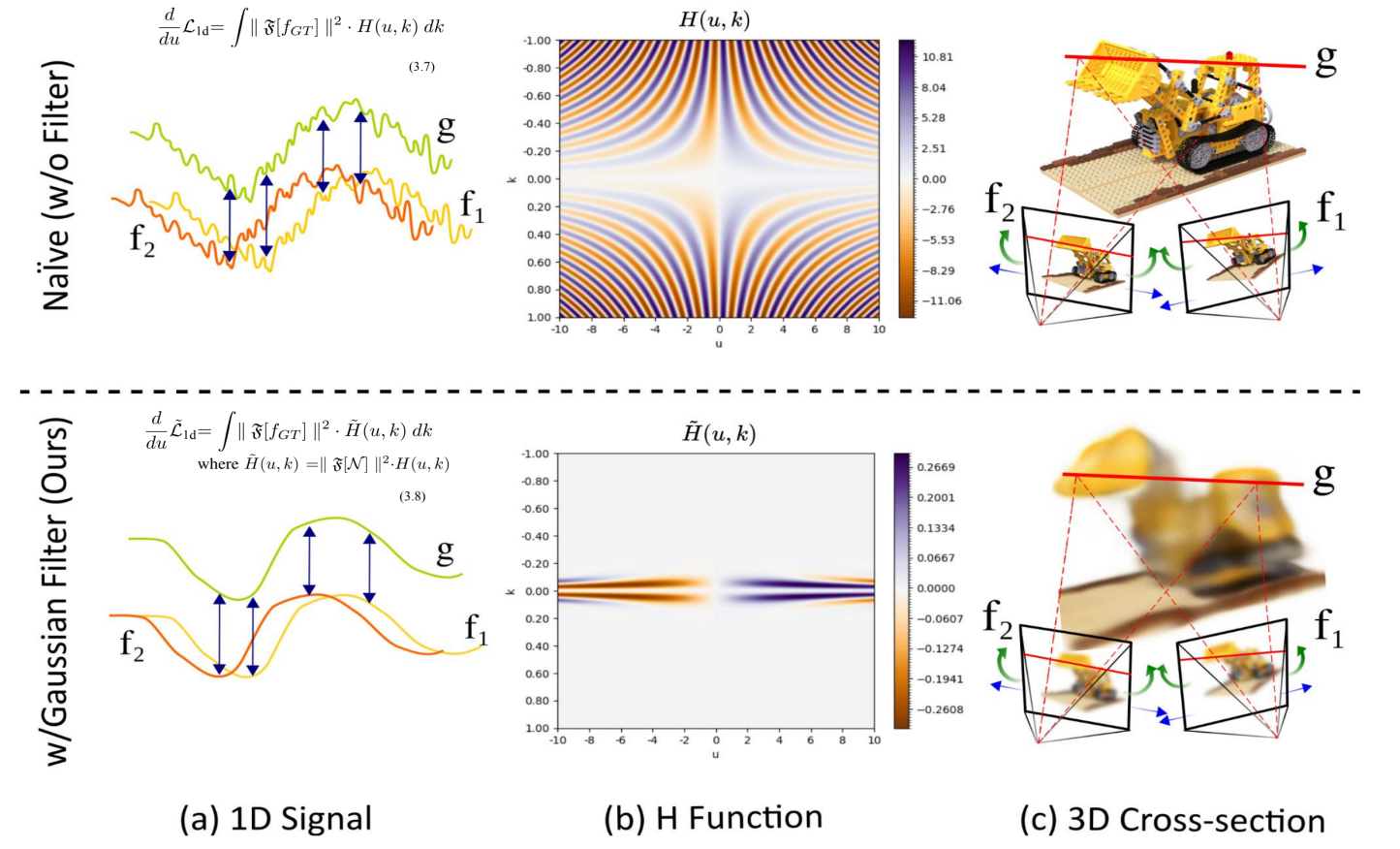


**Proposed Methods:** We start with 1D pilot study that discusses the effect of filtering strength on the joint optimization, from which we propose various methods that are proven effective in 2D and 3D experiment.



## 1D Analysis

We analyze how the spectral property of 1D signal can affect joint alignment gradient. We also explain how filtering helps finding global minima, and discuss the dilemma of filtering strength for weak gradient problems



## 2D Experiment

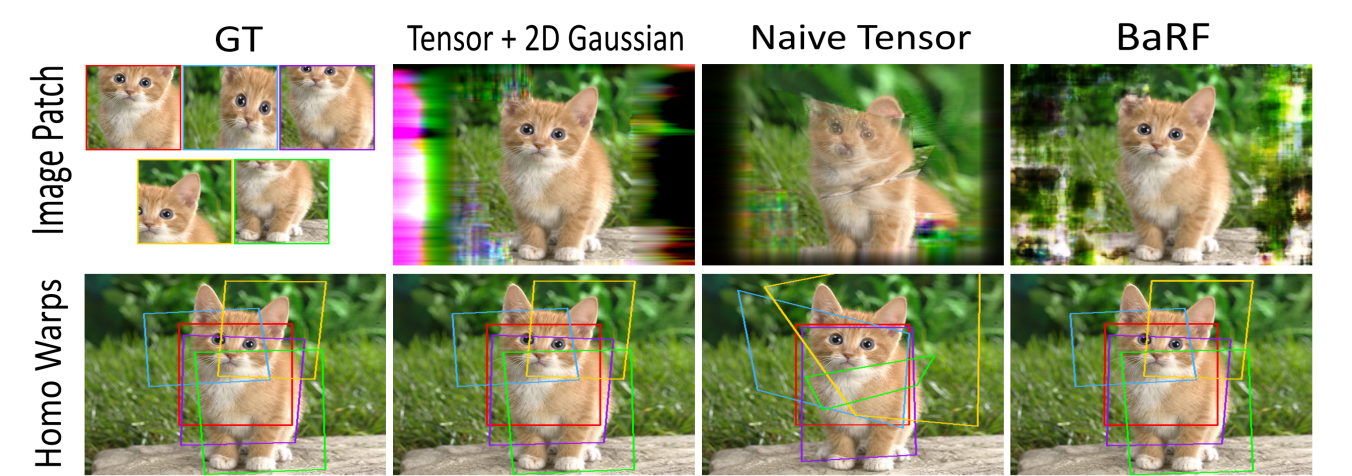
$$L_{2D}(P_{2D}, \mathcal{P}_{2D}) = \sum_{i=1}^L \sum_{u \in \mathcal{U}_{2D}} \|F_{2D}(W_{2D}(P_i, u)) - J_{1D}\|^2, \quad (9)$$

$$F_{2D}(x) = (\mathcal{N}_{2D} *_{2D} \mathcal{T}_{2D})(x) = (\mathcal{N}_{2D} *_{2D} \sum_{v \in \mathcal{V}} v^x \otimes v^y)(x), \quad (10)$$

In 2D example, brute force 2D convolution outperforms BARF in both quality and efficiency

Methods	$s(3)$ error ↓	patch PSNR ↑
BARF	0.0105	35.19
Naïve 2D TensoRF	0.5912	20.80
2D TensoRF + 2D Gaussian	<b>0.0023</b>	<b>40.70</b>

Table 1: Quantitative results of planar image alignment.



**Fast Convergence:** *Seperable Component-Wise Convolution* allows efficient 3D spectrum control on tensorial field, which in terms prevent the need for MLP PE control in BARF and sequential multi-resolution grids learning in HASH (Heo et al. 2023).

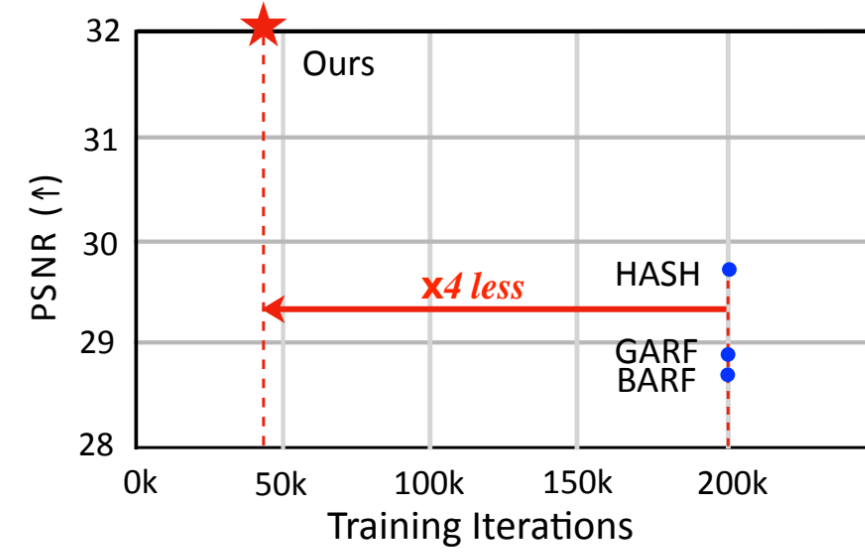
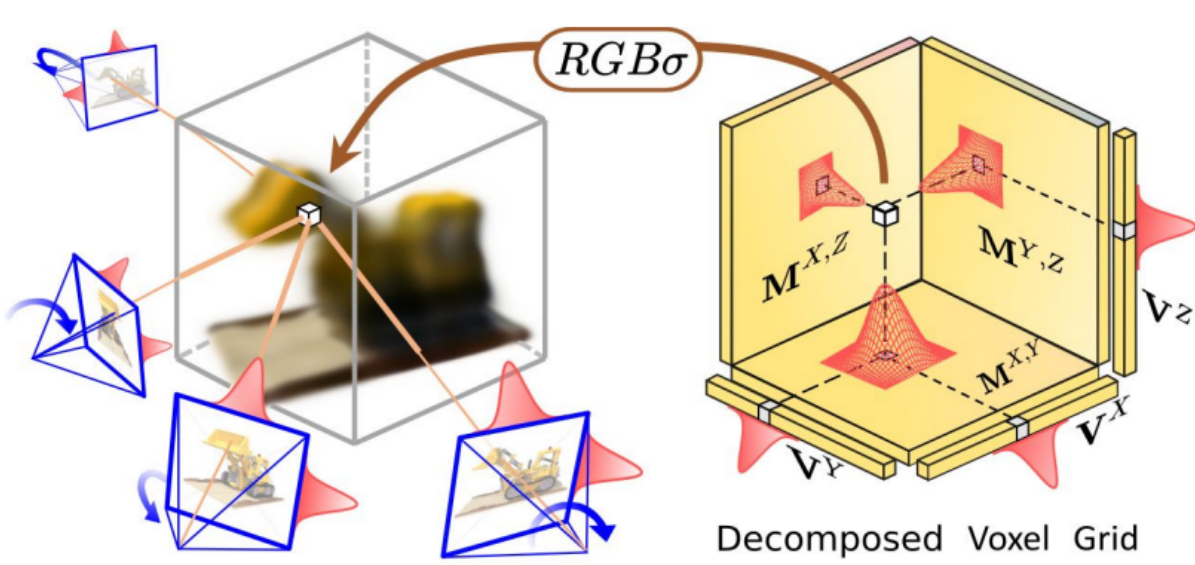
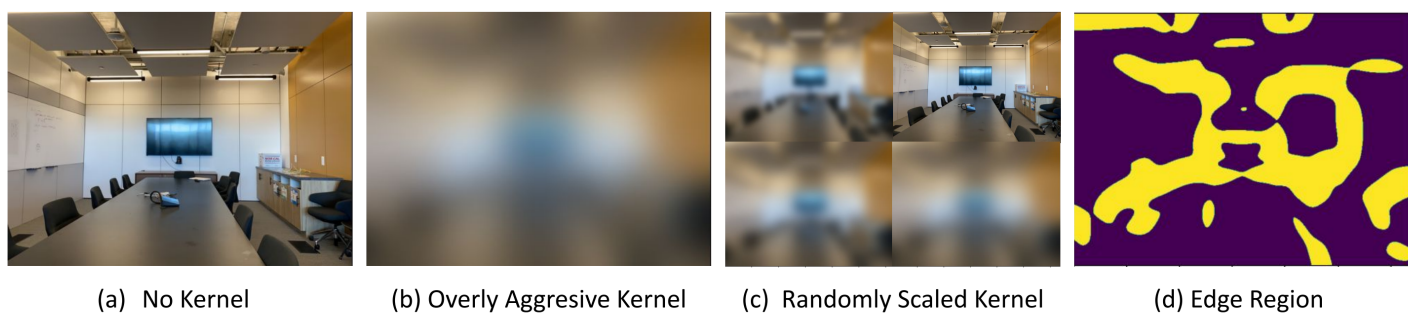


Figure 7: PSNR and training iterations comparison.

**Improving Robustness:** We propose *Randomly Scaled Kernel* and *Edge Guided Loss Mask* to improve the robustness of joint optimization, the former prevents local minima by randomized combination of 2D and 3D filtering, and the latter amplify gradient signal in edge regions which are critical for alignment.



**Algorithm :** Our proposed 3D method combines  
 1. Separated 3D Convolution on tensorial field.  
 2. Smoothed 2D supervision.  
 3. Randomly Scaled Kernels for both 2D and 3D  
 4. Edge Guided Loss for amplifying more useful gradients

Algorithm 1: Conceptual Training Process for Our Proposed 3D Joint Optimization Training

```

 $\mathcal{T}_\sigma, \mathcal{T}_c \leftarrow$  Initialize Voxel Grid
 $\mathcal{P} \leftarrow$  Initialize Camera Poses
for  $s = 1$  to train_iters do
 $i \leftarrow 2D\_kernel\_sched(s)$ 
 $\sigma, c \leftarrow 3D\_kernel\_sched(s)$ 
 $u_\sigma \leftarrow$  randomly sample density kernel scale
 $u_I \leftarrow$  randomly sample 2D kernel scale
 $\mathcal{N}_\sigma \leftarrow$  generate gaussian kernel with variance  $(\sigma \cdot u_\sigma)^2$ 
 $\mathcal{N}_c \leftarrow$  generate gaussian kernel with variance  $c^2$ 
 $\mathcal{N}_I \leftarrow$  generate gaussian kernel with variance  $(i \cdot u_I)^2$ 
 $\mathcal{L}_{joint} \leftarrow \mathcal{L}_{joint}(\mathcal{T}_\sigma, \mathcal{T}_c, \mathcal{P}, \mathcal{N}_\sigma, \mathcal{N}_c, \mathcal{N}_I)$  (with randomly selected pixels in all training views)
 $\mathcal{L}_{L1} = \mathcal{L}_{L1}(v_{\sigma,r}, M_{c,r}, v_{c,r}, M_{\sigma,r})$ 
 $\mathcal{L}_{TV} = \mathcal{L}_{TV}(v_{\sigma,r}, M_{c,r}, v_{c,r}, M_{\sigma,r})$ 
 $\mathcal{L}_{3d} = w_1 \cdot \mathcal{L}_{joint} + w_2 \cdot \mathcal{L}_{L1} + w_3 \cdot \mathcal{L}_{TV}$ 
back propagation
update  $\mathcal{T}_\sigma, \mathcal{T}_c, \mathcal{P}$ 
end for
    
```

**Ablation :** We show the importance of each proposed components and also demonstrate the necessity by showing the methods proposed by BARF or GARF are not applicable to tensorial radiance field.

	3D Gauss.	2D Gauss.	Random Kernel	Edge Guided	Rot. ↓	Trans. ↓	PSNR ↑
(a)	✓	✓	✓	✓	<b>0.72</b>	<b>0.33</b>	<b>25.36</b>
(b)	✓	✓	✓	✓	1.00	0.37	25.25
(c)	✓	✓	✓	✓	1.91	0.93	25.12
(d)	✓	✓	✓	✓	33.00	12.7	20.10
(e)	✓	✓	✓	✓	26.25	8.9	19.73
(d)	✓	✓	✓	✓	23.29	9.4	23.97

Table 4: Ablation study of the components of the proposed method on the real-world LLFF dataset.

	Rot. ↓	Trans. ↓	PSNR ↑	SSIM ↑	LPIPS ↓
TensoRF + BARF	45.47	0.17	20.71	0.630	0.314
TensoRF + GARF	73.92	0.29	10.47	0.287	0.679
Ours	0.43	0.003	26.92	0.872	0.104

Table 5: Ablation on Directly Applying BARF and GARF on TensoRF (Potential Baseline)

## 3D Experiments

**Quantitative Results :** Quantitative results shows superior synthesis quality compared to previous methods.

Scene	Camera Pose Registration				View Synthesis Quality			
	Rotation (°) ↓		Translation ↓		PSNR ↑		SSIM ↑	
	GARF	BARF	HASH	Ours	GARF	BARF	HASH	Ours
Chair	0.113	0.096	<b>0.085</b>	0.874	0.549	0.428	<b>0.365</b>	3.501
Drum	0.052	0.043	0.041	<b>0.037</b>	0.232	0.225	0.214	<b>0.118</b>
Ficus	0.081	0.085	0.079	<b>0.050</b>	0.461	0.474	0.479	<b>0.173</b>
Hotdog	0.235	0.248	0.229	<b>0.105</b>	1.123	1.308	1.123	<b>0.499</b>
Lego	0.101	0.082	0.071	<b>0.049</b>	0.299	0.291	0.272	<b>0.100</b>
Materials	<b>0.842</b>	0.844	0.852	0.854	<b>2.688</b>	2.692	2.743	2.690
Mic	0.070	0.071	<b>0.068</b>	1.177	0.293	0.301	<b>0.287</b>	5.000
Ship	0.073	0.075	0.079	<b>0.058</b>	0.310	0.326	0.287	<b>0.167</b>
Mean	0.195	0.193	<b>0.189</b>	0.400	0.744	0.756	<b>0.722</b>	1.533

Table 2: Quantitative results on the NeRF-Synthetic dataset. Our method achieves the best average novel-view synthesis quality and the best pose error in 5 out of 8 scenes. Notice that our method converges within 40k iterations, while all previous methods train for 200k iterations.

Scene	Camera Pose Registration				View Synthesis Quality			
	Rotation (°) ↓		Translation ↓		PSNR ↑		SSIM ↑	
	GARF	BARF	HASH	Ours	GARF	BARF	HASH	Ours
Fern	0.470	0.191	<b>0.110</b>	0.472	0.250	<b>0.102</b>	0.102	0.199
Flower	0.460	<b>0.251</b>	0.301	1.375	0.220	0.224	<b>0.211</b>	0.389
Fortress	<b>0.030</b>	0.479	0.211	0.449	0.270	0.364	<b>0.241</b>	0.419
Horns	<b>0.030</b>	0.304	0.049	0.386	0.210	0.222	<b>0.209</b>	0.251
Leaves	<b>0.130</b>	1.272	0.840	1.990	0.230	0.249	<b>0.228</b>	0.397
Orchids	0.330	0.627	0.399	<b>0.279</b>	0.410	0.404	0.386	<b>0.340</b>
Room	0.270	0.320	0.271	<b>0.188</b>	0.200	0.270	0.213	0.191
T-Rex	<b>0.420</b>	1.138	0.894	0.523	<b>0.360</b>	0.720	0.474	0.416
Mean	<b>0.280</b>	0.573	0.384	0.709	0.269	0.331	<b>0.258</b>	0.325

Table 3: Quantitative results on the LFF dataset. Our method achieves the best average novel-view synthesis quality and best LPIPS in 7 out of 8 scenes. Our method converges within 50k iterations, while all previous methods train for 200k iterations.

## Conclusion :

**1. Theoretically,** we provide insights into the impact of scene properties on the convergence of joint optimization beyond the coarse-to-fine heuristic discussed in prior research, proposing a filtering based strategy for improving the joint optimization.

**2. Algorithmically,** we introduce ( and prove the equivalence of ) an effective method for applying the filtering based strategy on decomposed low-rank tensor, notice that the proposed *Seperable Component-Wise Convolution* is unique and more efficient than traditional separable methods in the sense that we (aside from the separated kernel) additionally utilize the separability of input signal. We additionally proposed techniques (i.e. *Randomly Scaled Kernel* and *Edge Guided Loss Mask*) for improving the robustness of joint optimization.

**3. Comprehensive Evaluations** demonstrates our proposed framework's state-of-the-art performance and rapid convergence.