

NeRO: Neural Road Surface Reconstruction

Ruibo Wang¹, Song Zhang¹, Ping Huang², Donghai Zhang², and Haoyu Chen²

Z-one Technology Co., LTD. Shanghai, China

{wangruibo01,zhangsong05,huangping01,zhangdonghai,chenhaoyu03}@saicmotor.com

Abstract. Accurately reconstructing road surfaces is pivotal for various applications especially in autonomous driving. This paper introduces a position encoding Multi-Layer Perceptrons (MLPs) framework to reconstruct road surfaces, with input as world coordinates x and y , and output as height, color, and semantic information. The effectiveness of this method is demonstrated through its compatibility with a variety of road height sources like vehicle camera poses, LiDAR point clouds, and SFM point clouds, robust to the semantic noise of images like sparse labels and noise semantic prediction, and fast training speed, which indicates a promising application for rendering road surfaces with semantics, particularly in applications demanding visualization of road surface, 4D labeling, and semantic groupings.

Keywords: Road Reconstruction · Multi-resolution Hash Positional Encoding · Positional Encoding · Semantic Label

1 Introduction

The evolving landscape of 3D reconstruction has led to significant advancements, especially in reconstructing complex urban environments like road surfaces, which is useful in 4D labeling and semantic groupings. While effective, traditional approaches often grapple with challenges such as computational intensity, low-quality rendering, and semantic noise to be improved.

NeRF presents a method for synthesizing novel views of complex scenes by modeling the volumetric scene function using a fully connected deep neural network. While NeRF’s methodology provides high-quality 3D reconstructions, its voxel representations are redundant for surface modeling and computationally intensive. Tesla claimed on its 2021 AI Day that it uses implicit MLPs to reconstruct road surface color, semantics, and elevation without using any explicit meshes or point clouds but offering no further information. RoME [14] uses explicit mesh objects from Pytorch3D [19] lib to represent the color and semantics and uses the implicit MLPs network with position encoding only to calculate the road height. After learning road height with supervised road ground truth like SFM point clouds or lidar point clouds, it utilizes ray-tracing based method in Pytorch3D lib to optimize the color and semantics, accelerating the render speed than NeRF.

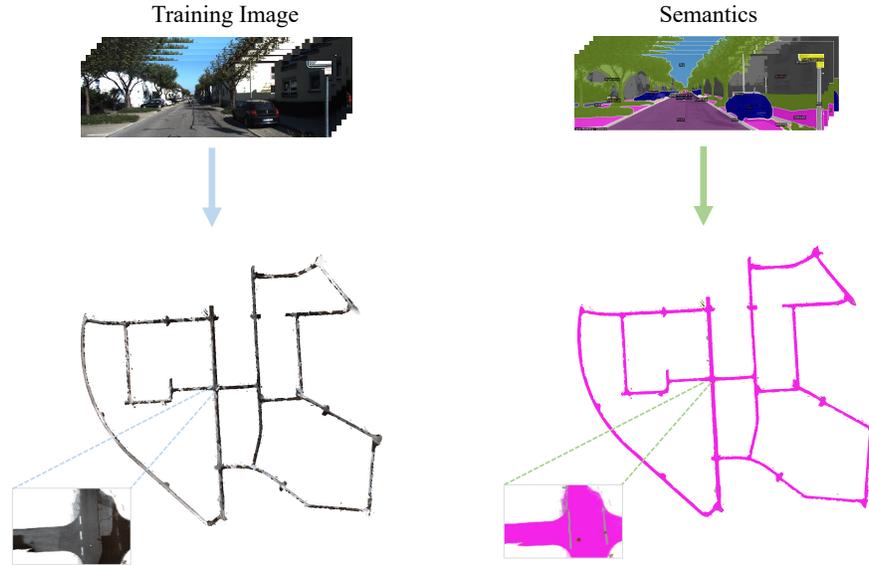


Fig. 1: The result of our method in reconstructing an entire segment of the road from the KITTI Odometry Sequence 00.

We use a unified position encoding MLPs-based network architecture Neural Road Surface Reconstruction (**NeRO**) to reconstruct height, color, and semantics in a function like $(\text{height}, \text{color}, \text{semantics}) = \text{MLP}(x, y)$. The ground truth of height could be derived from various sources like vehicle camera pose, LiDAR point clouds, and SFM point clouds. The ground truth of color and semantics are queried by sampling millions of world coordinates road surface 3D points around each camera and projecting them back to related image coordinates using camera extrinsic and intrinsic, without utilizing ray-tracing based method. The main contributions of this paper include:

- We introduce a position encoding MLPs-based road surface reconstruction method to represent road surfaces with color, semantics, and height. Our method uses the same software libraries as NeRF without any further complex software, and it may be the closest implementation to Tesla’s implicit road surface reconstruction.
- We have verified the performance of our architecture in road height reconstruction under different input sources, such as vehicle camera pose, SFM point clouds, and LiDAR point clouds. We have also tested its capability to reconstruct slopes (using square waves as a substitute) and its ability to complete sparse point cloud inputs or incomplete height information.
- Additionally, we validated our architecture’s robustness against semantic noise. Experiments show that semantic noise in single-frame images can be optimized to some extent by aggregating multi-frame, multi-view semantic information, which helps improve the accuracy of 4D road surface labeling.

Our code will be released at: <https://github.com/ToeleoT/NeRO>

2 Related Works

2.1 Multi-view 2D Image to 3D Reconstruction

A notable contribution in this realm is the rendering technique outlined in [21, 22]. Those methods demonstrate a robust capability to reconstruct three-dimensional geometries from two-dimensional images, leveraging known camera poses. In the context of urban mapping, [18] presents a technique that focuses on extracting road line segmentation data from images, followed by an inverse projection of this information to generate three-dimensional maps. A similar method [27] introduces a method for creating textured meshes from images. [11] explores depth map creation from various viewpoints, which fuse them into a cohesive three-dimensional structure. [23, 32] capitalizes on the advantages of planar structures, thereby refining the efficiency and accuracy of three-dimensional reconstructions in certain contexts.

2.2 NeRF-based Reconstruction

Neural Radiance Fields (NeRF) [15] represents an implicit approach, utilizing MLPs and positional encoding method divergent by inputting viewing angles to reconstruct three-dimensional environments. Subsequent research has expanded upon the foundational principles of NeRF, adapting it for extensive, unbounded scenes. Specifically, NeRF++ [34] employs a methodology of normalizing the unbounded scene into a unit sphere to render the environments. Further advancements are seen in BlockNeRF [24], which deconstructs large-scale scenes into smaller blocks with each NeRF framework, incorporating factors of exposure and visibility to reconstruct these extensive environments. Conversely, Mip-NeRF 360 [1] introduces an innovative wrapping technique for parameterizing Gaussians. [10] renders the background with a single MLP, and other dynamic objects are learned by a set of MLPs to present a panoptic street view. [3, 12, 20, 31] uses the depth map, LiDAR, and point cloud data to help NeRF to build more accurate 3D geometry.

2.3 NeRF with Semantic

There has been much work integrating semantic information into 3D scene reconstruction in the evolving landscape of the implicit method. A pivotal development in this domain is Semantic-NeRF [35], which incorporates an additional layer within the MLPs of NeRF to render the semantic details within three-dimensional environments. Expanding upon this concept, PanopticNeRF [6] merges 3D semantic data with three-dimensional bounding boxes to enhance the geometric fidelity of semantic representations. Another innovative approach is presented in NeRF-SOS [4], which leverages self-supervised learning techniques

within the NeRF framework to calculate semantic information. Furthermore, NeSF (Neural Semantic Fields) [26] demonstrates a unique approach by extracting density information from NeRF outputs, which is processed through a 3D U-Net architecture.

2.4 NeRF-based Road surface Reconstruction

RoME [14] significantly contributes to the reconstruction of road surfaces by using explicit mesh objects from Pytorch3D [19] lib to represent the color and semantics and using the implicit MLPs network with position encoding to calculate the road height. After learning road height with supervised road ground truth like SFM and lidar, it utilizes ray-tracing based method in Pytorch3D lib to go on optimizing the color and semantics. MV-Map [30] adopts a voxel-based approach within the NeRF framework like [5, 33], paving the way for the construction of high-definition maps. Moreover, PlaNeRF [28] presents a plane regulation methodology grounded in the decomposition (SVD) technique to reconstruct scenes efficiently. StreetSurf [8] distinguishes itself by segmenting scenes within images and applying varied scales of multi-resolution hash positional encoding.

3 Method

We show our framework and pipeline in Fig. 2. **NeRO** processes the input to calculate the height along the vertical z-axis, RGB, and semantic.

3.1 NeRO Network Structure

NeRO takes the x and y coordinates in the world coordinate system, $\mathbf{X} = (x, y)$, as input. Before entering the network layers, our input \mathbf{X} is normalized to between [-1,1] to facilitate the calculation of encoding methods. We feed the normalized input \mathbf{X}' into position encoding functions. Then, the outputs from the position encoding methods are processed by three different MLPs, resulting in outputs for road surface height \mathbf{z} , colour output $\mathbf{c} = (r, g, b)$, and semantic output \mathbf{s} . NeRO can also choose different position encoding method for different MLPs. By the way, only MLPs without any position encoding fails to reconstruct road height, color and semantic through all of our experiments.

3.2 Encoding methods

Positional Encoding From [15, 16, 25], we need to leverage Positional Encoding, which is used to learn high-frequency details with the formula:

$$PE(X) = (\sin(2^0\pi X), \cos(2^0\pi X), \dots, \sin(2^{L-1}\pi X), \cos(2^{L-1}\pi X)) \quad (1)$$

Where L is the hyperparameter that decides the length of the function.

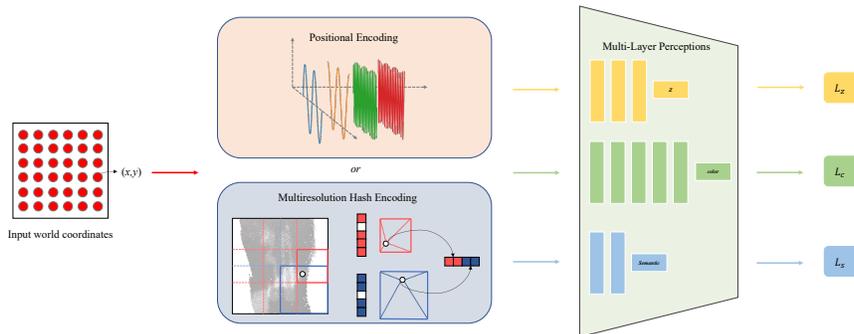


Fig. 2: NeRO Overview. The world coordinates $\mathbf{X} = (x, y)$ are encoded by the Positional Encoding or Multiresolution Hash Positional Encoding. Then, the encoded information is passed into three different MLPs responsible for calculating the height z , color, and semantic.

Multi-Resolution Hash Positional Encoding In [16], a learnable encoding method, Multi-Resolution 3D Hash Positional Encoding is introduced. Our approach modify a Multi-Resolution 2D hash Positional Encoding based on it.

3.3 Reconstruction methods

Z-axis Reconstruction We employ ground truth height from three sources: vehicle camera pose, LiDAR point clouds, and SfM point clouds. In vehicle camera pose, it is assumed that the ground plane near the corresponding pose is flat, and each pose will sample the points within a length * width certain area to form pseudo point clouds. The different representations in the encoding method will affect the result of the height value.

Color Reconstruction In color reconstruction, we sample millions of 2D world coordinates as network input $\mathbf{X} = (x, y)$ for each pose. Then, we use those coordinates to obtain height z and color \mathbf{c} by the complete learned height network and the color network separately. After that, we combine the road surface height z with $\mathbf{X} = (x, y)$ to obtain the 3D world coordinates $\mathbf{W} = (x, y, z)$, then project them into the pixel coordinate system (\mathbf{u}, \mathbf{v}) with camera’s extrinsic and intrinsic to get the corresponding ground truth pixel color \mathbf{c}' to optimize network output color \mathbf{c} .

Semantic Reconstruction The network semantic output \mathbf{s} from the encoding methods is used to render the semantic information for the network input $\mathbf{X} = (x, y)$, which use the same method in color reconstruction to obtain the ground truth semantic \mathbf{s}' .

3.4 Loss Function

We calculate the Mean Squared Error loss for the height MLPs output \mathbf{z} with the ground truth \mathbf{z}' , which is the \mathcal{L}_z . Then, we use the true colour \mathbf{c}' obtained from the pixel coordinate system and compare it with the network output colour \mathbf{c} . After that, it will perform Mean Squared Error loss calculation \mathcal{L}_c . Finally, we use the semantic ground truth \mathbf{s}' and the semantic output \mathbf{s} to perform a Cross-entropy loss calculation \mathcal{L}_s .

The total loss for our method is that:

$$\mathcal{L} = \mathcal{L}_z + \mathcal{L}_c + \mathcal{L}_s \quad (2)$$

$$= \sum_{n=1}^N \left[\|z - z'\|_2^2 + \|c - c'\|_2^2 + \sum_{l=1}^L s \log(s') \right] \quad (3)$$

Where N is the number of points inside the road surface that appear in the road part of the training images, and L is the number of semantic labels used in the training process.

4 Experiment

4.1 Experiment Settings

Datasets Our experiment uses Kitti odometry [7] sequence 00 datasets to test our method, and we use images from the left camera as our training datasets. As the Kitti odometry datasets do not have any semantic information, we use the state-of-the-art semantic prediction method Mask2Former [2] with Swin-L [13] backbone to acquire that information. The semantic categories we used are road, traffic lane and Manhole, which represent the main information appears in road. The SfM points are obtained by COLMAP, which uses the provided but not perfect ground truth camera pose in kitti odometry datasets for reconstruction.

Evaluation Metrics Our experiment compares the pixel color from the ground images, and use PSNR as metric to verify the rendering quality. For semantic label accuracy, the metric is mIoU.

Environment Setup We implements NeRO method by Pytorch [17] framework, which uses the Nvidia A100 80G for the experiments, while a GPU with 8G memory is possible. The training images are the original fixed size, 1241 x 376, with batch size 1. The training optimizer we selected is Adam [9], with a learning rate of 5e-4.

4.2 Road Surfaces Height Reconstruction with different PE Method

To illustrate the road height reconstruction performance, an experimental setup was set utilizing a square wave-like road structure, as depicted in the Fig. 3. The

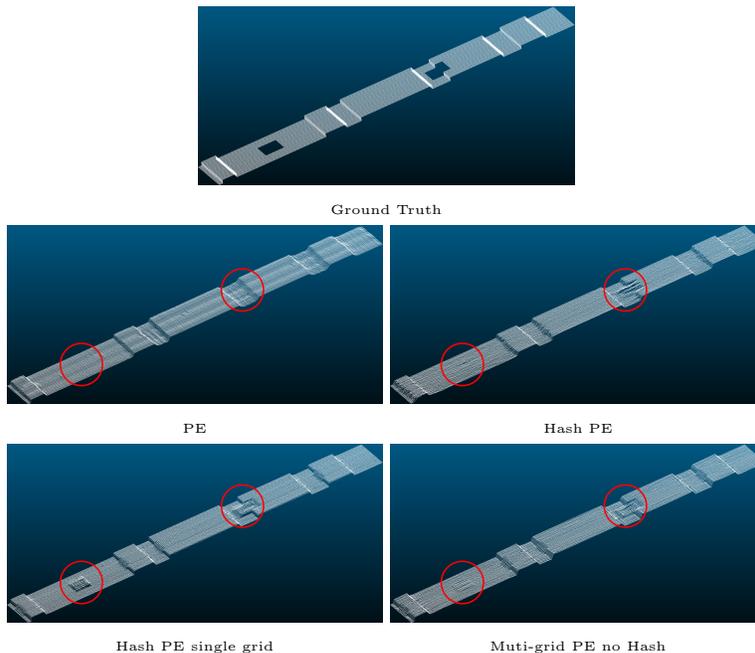


Fig. 3: Comparison in reconstructing the incomplete road. The red circle indicates the hole in the horizontal and corner of the road surface.

positional encoding method successfully fills these holes, accurately replicating the correct shape. This effectiveness is attributed to its inherent periodic property, which provides prior information. However, a limitation of this approach is the lack of smoothness on the road surface. Multi-resolution hash positional encoding excels in rendering the road surface with remarkable smoothness. It fills the horizontal holes smoothly and attempts to fill the corner holes with values from adjacent areas. This is due to the absence of prior information in the multi-resolution hash positional encoding approach, leading to a less smooth appearance. We also examine the ability of multi-resolution hash positional encoding when it has only a single resolution or devoid hash function. The figure shows that only a single hash function without multiresolution cannot fill horizontal and corner holes. Conversely, employing a multiresolution approach without a hash function demonstrates a filling capability similar to the multi-resolution hash positional encoding method.

4.3 Color and Semantic Reconstruction with Different PE methods and Ground Truth Height

In this part, as shown in Fig. 4, we show our experiment results in color and semantic results in different positional encoding methods for our method **NeRO**

Table 1: Quantitive result for our method in reconstructing the road surface in color and semantics with the road height from vehicle camera pose, LiDAR, and SfMs using multi-resolution hash positional encoding or positional encoding method.

Road height	Positional Encoding Type	PSNR	mIoU
Vehicle Camera Pose	PE	17.81	0.704
Vehicle Camera Pose	Hash PE	25.73	0.988
LiDAR	PE	18.87	0.701
LiDAR	Hash PE	29.20	0.994
SfM-Dense	PE	18.76	0.784
SfM-Dense	Hash PE	25.81	0.975
SfM-Sparse	PE	18.38	0.780
SfM-Sparse	Hash PE	24.37	0.967

with different input datasets, and in Tab. 1, we offer our quantitive metric for each result.

The results show that road surface height reconstruction with lidar point clouds performs best over other road height sources. However, the road surface obtained by SFM reconstruction is similar to that of vehicle pose, probably caused by the fact that the Kitti dataset is collected on flat roads. We also show that the performs of hash PE is better than PE.

4.4 Road Semantic Reconstruction with Noise Labels

This section shows how our method is applied in some applications, like sparse semantic labels and semantic label denoising. In NeRF-based models, like [29, 35] offer its ability on that kind of application. We only chose the LiDAR road height for this process because it has relatively accurate road height values, and in Tab. 2 and Tab. 3, we offer our quantitive metric for each result.

Sparse Label We hypothesize a scenario where a significant portion of ground truth semantic images is unavailable, leaving only a minimal subset for use. From [35], it indicated that using less than 10% of semantic images will significantly decrease the rendering quality. Our investigation will set the threshold to only 10% of images to evaluate whether our method could reconstruct the road surface. Fig. 5 presents the outcomes by using positional encoding or multi-resolution hash positional encoding. The results demonstrate that, despite the limited availability of training semantics, both methods can reconstruct the road surfaces. However, the distinct advantage of multi-resolution hash positional encoding is its ability to render detailed features, such as manholes, which positional encoding fails to replicate with the same level of detail. Tab. 2 shows the quantitive result in the mIoU. Here, it is evident that multi-resolution hash positional encoding outperforms positional encoding, showcasing superior mIoU values.

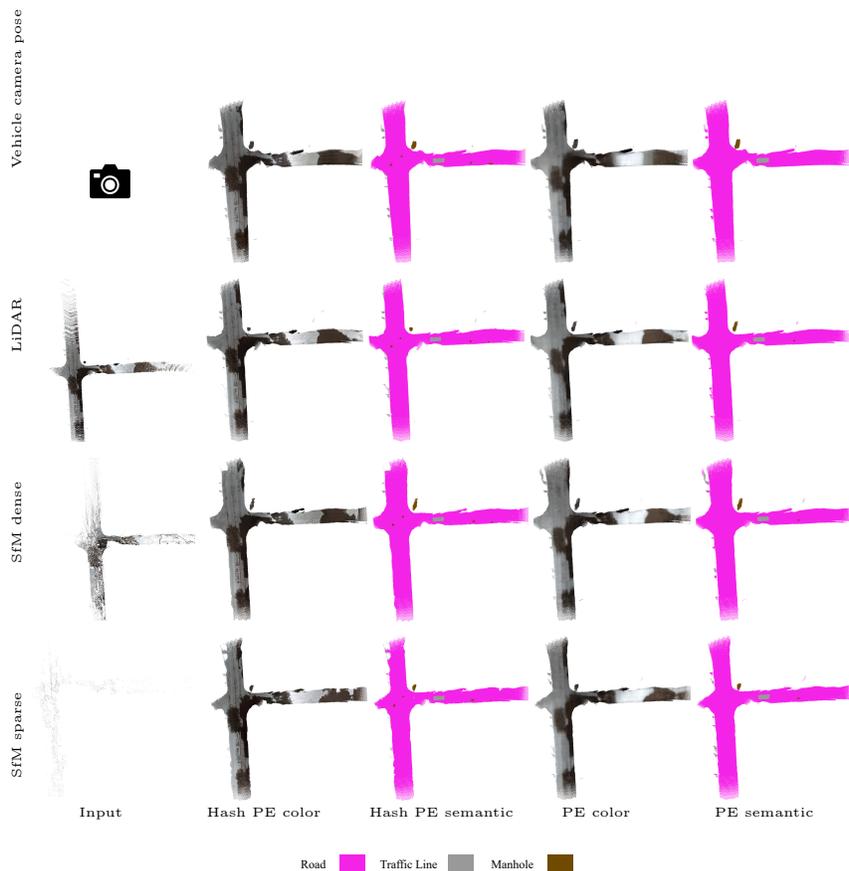
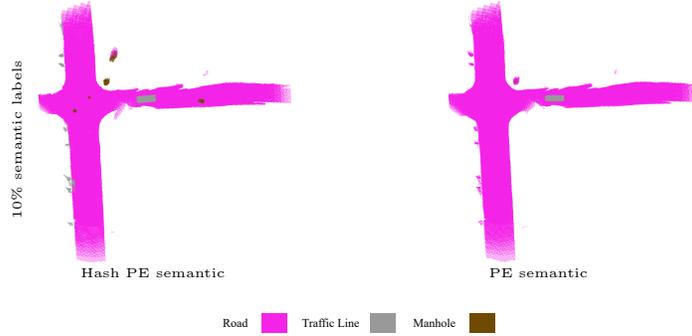


Fig. 4: In qualitative comparison with the PE and Hash PE in the datasets of the vehicle camera pose, LiDAR points, SfMs dense points, and SfMs sparse points. The first column shows each input dataset, and the rest illustrates the reconstruction result. Hash PE has a better render quality in each result than in color and semantics.

Noise Label In the actual scenarios, the road datasets have noise content, so we simulate some pixel noises to the ground truth by using the method in [35], and the noised semantic labels are used to examine our method’s denoise ability. Fig. 6 illustrates the outcomes of applying denoising reconstruction to noisy labels using positional encoding or multi-resolution hash positional encoding. In scenarios with 50% noise in the labels, it is observed that both multi-resolution hash positional encoding and positional encoding can reconstruct the road surface. Still, the positional encoding method fails to reconstruct some specific items, like manholes in the middle of the road. When we add the ratio of noise labels to 90%, multi-resolution hash positional encoding cannot denoise the label due to its accurate learning ability. Its precision in learning leads to the unintentional incorporation of noise labels as road surface components. In contrast,

Table 2: Quantitive result for our method for sparse labels, when using only 10% training semantic dataset.

Sparse Ratio	Positional Encoding Type	mIoU
0.1	PE	0.670
0.1	Hash PE	0.823

**Fig. 5:** Comparison with the positional encoding and multi-resolution hash positional encoding in the sparse semantic labels. Both methods render the whole structure of the road, but the Hash PE on the left gives more detail in the result.

the positional encoding method continues to denoise the labels. However, it can still not render detailed features like manholes in the middle of the road. The quantitative analysis in Tab. 3 shows that the mIoU value of positional encoding is higher than multi-resolution hash positional encoding, indicating a better semantic denoise ability. Still, positional encoding needs to be more comprehensive in rendering the details. Further investigations were conducted to determine the noise ratio threshold of the denoise ability of multi-resolution hash positional encoding. Visual results indicate that multi-resolution hash positional encoding begins to integrate noise labels into the road surface when the noise ratio exceeds 0.6.

Table 3: The quantitative result for our method is denoising the noisy semantic labels, which use positional encoding or multi-resolution hash positional encoding method to denoise the label with the noise ratio when 50% or 90% labels are noised in LiDAR datasets.

Noise Ratio	Positional Encoding Type	mIoU
0.5	PE	0.693
0.5	Hash PE	0.625
0.9	PE	0.420
0.9	Hash PE	0.292



Fig. 6: Comparison with the positional encoding and multi-resolution hash positional encoding in the noisy semantic labels. The first column shows the input noisy datasets, the second column records the multi-resolution hash positional encoding result, and the final column gives the positional encoding result.

4.5 Training speed

During the networks' training phase, the loss reduction is shown in the Fig. 7. When examining the vertical z-axis loss, both positional encoding and multi-resolution hash positional encoding methods exhibit comparable training speeds, ultimately converging on their global minima. Multi-resolution hash positional encoding demonstrates an accelerated color and semantic loss training speed, efficiently attaining the global minima. Conversely, the positional encoding method experiences a slower decrease in loss, facing challenges in achieving convergence to the global minima.

5 Limitation

The limitation of our method is the difficulty in achieving a balance between the ability to fit detailed information and model complexity on road surfaces height reconstruction. Most road surface is flat and could be described with low parameter multivariate equations, and there is no need to use MLPs with big parameters. But for the sake of some little detailed information like holes in the road we have to use position encoding and MLPs with huge parameters. In

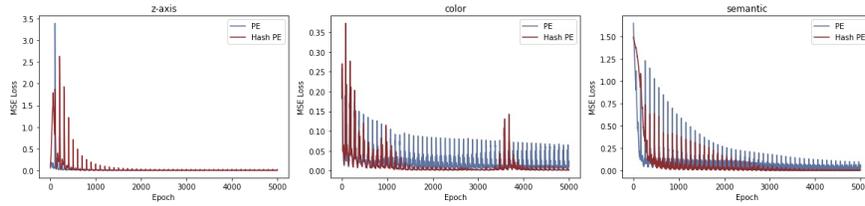


Fig. 7: Comparison in the training loss decay between positional encoding and multi-resolution hash positional encoding in the z-axis, color, and semantic. Those spikes appear because we train in sequence, not in random.

future research, we propose to refine our approach by a hybrid representation of road surface.

6 Conclusion

In conclusion, we introduce an position encoding MLPs-based neural road reconstruction method that accepts the various sources of round truth of height to render the color and semantic information with road height output. Our experiment shows the success of our rendering ability in either color or semantic information. In the semantic applications, it also shows that it can handle sparse labels and noise labels. We also compare the Positional Encoding and Muti-resolution Positional Encoding methods to deliver each performance. This indicates that the Muti-resolution Positional Encoding method performs better in rendering the road surfaces in quality and speed.

References

1. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5470–5479 (2022) [3](#)
2. Cheng, B., Misra, I., Schwing, A.G., Kirillov, A., Girdhar, R.: Masked-attention mask transformer for universal image segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 1290–1299 (2022) [6](#)
3. Deng, K., Liu, A., Zhu, J.Y., Ramanan, D.: Depth-supervised nerf: Fewer views and faster training for free. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12882–12891 (2022) [3](#)
4. Fan, Z., Wang, P., Jiang, Y., Gong, X., Xu, D., Wang, Z.: Nerf-sos: Any-view self-supervised object segmentation on complex scenes. arXiv preprint arXiv:2209.08776 (2022) [3](#)
5. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5501–5510 (2022) [4](#)
6. Fu, X., Zhang, S., Chen, T., Lu, Y., Zhu, L., Zhou, X., Geiger, A., Liao, Y.: Panoptic nerf: 3d-to-2d label transfer for panoptic urban scene segmentation. In: 2022 International Conference on 3D Vision (3DV). pp. 1–11. IEEE (2022) [3](#)
7. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: 2012 IEEE conference on computer vision and pattern recognition. pp. 3354–3361. IEEE (2012) [6](#)
8. Guo, J., Deng, N., Li, X., Bai, Y., Shi, B., Wang, C., Ding, C., Wang, D., Li, Y.: Streetsurf: Extending multi-view implicit surface reconstruction to street views. arXiv preprint arXiv:2306.04988 (2023) [4](#)
9. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014) [6](#)
10. Kundu, A., Genova, K., Yin, X., Fathi, A., Pantofaru, C., Guibas, L.J., Tagliasacchi, A., Dellaert, F., Funkhouser, T.: Panoptic neural fields: A semantic object-aware neural scene representation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12871–12881 (2022) [3](#)
11. Leroy, V., Franco, J.S., Boyer, E.: Shape reconstruction using volume sweeping and learned photoconsistency. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 781–796 (2018) [3](#)
12. Li, Z., Li, L., Ma, Z., Zhang, P., Chen, J., Zhu, J.: Read: Large-scale neural scene rendering for autonomous driving. arxiv preprint [2022-12-11] (2022) [3](#)
13. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 10012–10022 (2021) [6](#)
14. Mei, R., Sui, W., Zhang, J., Zhang, Q., Peng, T., Yang, C.: Rome: Towards large scale road surface reconstruction via mesh representation. arXiv preprint arXiv:2306.11368 (2023) [1](#), [4](#)
15. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM **65**(1), 99–106 (2021) [3](#), [4](#)

16. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* **41**(4), 1–15 (2022) [4](#), [5](#)
17. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32** (2019) [6](#)
18. Qin, T., Zheng, Y., Chen, T., Chen, Y., Su, Q.: A light-weight semantic map for visual localization towards autonomous driving. In: 2021 IEEE International Conference on Robotics and Automation (ICRA). pp. 11248–11254. IEEE (2021) [3](#)
19. Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501* (2020) [1](#), [4](#)
20. Rematas, K., Liu, A., Srinivasan, P.P., Barron, J.T., Tagliasacchi, A., Funkhouser, T., Ferrari, V.: Urban radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12932–12942 (2022) [3](#)
21. Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4104–4113 (2016) [3](#)
22. Schönberger, J.L., Zheng, E., Frahm, J.M., Pollefeys, M.: Pixelwise view selection for unstructured multi-view stereo. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III 14. pp. 501–518. Springer (2016) [3](#)
23. Sun, S., Zheng, Y., Shi, X., Xu, Z., Liu, Y.: Phi-mvs: Plane hypothesis inference multi-view stereo for large-scale scene reconstruction. *arXiv preprint arXiv:2104.06165* (2021) [3](#)
24. Tancik, M., Casser, V., Yan, X., Pradhan, S., Mildenhall, B., Srinivasan, P.P., Barron, J.T., Kretschmar, H.: Block-nerf: Scalable large scene neural view synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8248–8258 (2022) [3](#)
25. Tancik, M., Srinivasan, P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Singhal, U., Ramamoorthi, R., Barron, J., Ng, R.: Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems* **33**, 7537–7547 (2020) [4](#)
26. Vora, S., Radwan, N., Greff, K., Meyer, H., Genova, K., Sajjadi, M.S., Pot, E., Tagliasacchi, A., Duckworth, D.: Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes. *arXiv preprint arXiv:2111.13260* (2021) [4](#)
27. Waechter, M., Moehrle, N., Goesele, M.: Let there be color! large-scale texturing of 3d reconstructions. In: Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13. pp. 836–850. Springer (2014) [3](#)
28. Wang, F., Louys, A., Piasco, N., Bennehar, M., Roldão, L., Tsishkou, D.: Planerf: Svd unsupervised 3d plane regularization for nerf large-scale scene reconstruction. *arXiv preprint arXiv:2305.16914* (2023) [4](#)
29. Wang, R., Zhang, S., Huang, P., Zhang, D., Yan, W.: Semantic is enough: Only semantic information for nerf reconstruction. In: 2023 IEEE International Conference on Unmanned Systems (ICUS). pp. 906–912. IEEE (2023) [8](#)
30. Xie, Z., Pang, Z., Wang, Y.X.: Mv-map: Offboard hd-map generation with multi-view consistency. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8658–8668 (2023) [4](#)

31. Xu, Q., Xu, Z., Philip, J., Bi, S., Shu, Z., Sunkavalli, K., Neumann, U.: Point-nerf: Point-based neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5438–5448 (2022) [3](#)
32. Xu, Q., Tao, W.: Planar prior assisted patchmatch multi-view stereo. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 12516–12523 (2020) [3](#)
33. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenotrees for real-time rendering of neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5752–5761 (2021) [4](#)
34. Zhang, K., Riegler, G., Snavely, N., Koltun, V.: Nerf++: Analyzing and improving neural radiance fields. arXiv preprint arXiv:2010.07492 (2020) [3](#)
35. Zhi, S., Laidlow, T., Leutenegger, S., Davison, A.J.: In-place scene labelling and understanding with implicit scene representation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15838–15847 (2021) [3](#), [8](#), [9](#)