

TK-Planes: Tiered K-Planes with High Dimensional Feature Vectors for Dynamic UAV-based Scenes

Christopher Maxey^{1,2}, Jaehoon Choi², Yonghan Lee², Hyungtae Lee³, Dinesh Manocha², and Heesung Kwon¹
¹DEVCOM Army Research Laboratory ²University of Maryland at College Park ³BlueHalo

In this paper, we present a new approach to bridge the domain gap between synthetic and real-world data for unmanned aerial vehicle (UAV)-based perception. Our formulation is designed for dynamic scenes, consisting of moving objects or human actions, where the goal is to recognize the pose or actions. We propose an extension of K-Planes Neural Radiance Field (NeRF), wherein our algorithm stores a set of tiered feature vectors. The tiered feature vectors are generated to effectively model conceptual information about a scene as well as an image decoder that transforms output feature maps into RGB images. Our technique leverages the information amongst both static and dynamic objects within a scene and is able to capture salient scene attributes of high altitude videos. We evaluate its performance on challenging datasets, including Okutama Action and UG2, and observe considerable improvement in accuracy over state of the art aerial perception algorithms.

I. INTRODUCTION

Synthetic data has been used for over a decade to train deep learning algorithms in a variety of tasks, from game play [24] to driving [37]. Data from sources such as game engines [27], [33], [2] and generative algorithms [13], [31] have made significant advancements in recent years, particularly in enhancing fidelity and bridging the domain gap between synthetic [35], [36] and real data [21]. Generative algorithms for instance can create realistic novel imagery and videos from a given prompt [31], [1], requiring an implicit model of real world objects and their codependencies. Neural Radiance Fields [23] and, more recently, Gaussian Splatting [16] have demonstrated remarkable capabilities in rendering both specific 3D static and dynamic scenes, providing novel views of a scene that were not available during the training process. Because of these advances, synthetic data have been playing a larger role in supplementing existing real world datasets for training data hungry perception algorithms [30], [8].

Specific domains, such as videos captured by UAVs [3], [17], [36], present unique challenges stemming from environmental conditions and a relative scarcity of datasets within these domains. In regard to UAVs, perception is not limited to static tasks such as object detection within a single frame but can also include more complex tasks like action recognition which can require multiple frames of a dynamic agent. The inclusion of dynamic objects and agents further compounds the difficulty of a task and accentuates the need for large amounts of data. Esoteric tasks, for example search and

rescue, add another layer of complication as more general datasets may not capture the requisite environmental detail and conditions or relevant agent behavior.

Given a dynamic UAV scene, several issues become pertinent when trying to model accurate dynamic behavior. First, because of the nature of high altitude data, target objects are often at a large distance from the camera and appear small relative to the background. Due to the more difficult task of modeling dynamic objects, this imbalance between static and dynamic elements of a scene can cause issues with separating static and dynamic elements within any given algorithm, in particular if there are inaccuracies present in camera pose information. Second, dynamic objects are often sparse in the scene at any given point of time. This further limits the amount of data that can afford training. Lastly, diverse poses at specific time stamps may be limited, often with just one camera pose per timestamp and limited variation in camera poses. This can lead to degeneracies in the output of a model the more poses and their corresponding timestamps diverge from the training data.

Main Results: Our method aims to rectify these issues with a novel Neural Radiance Field architecture [23]. We present a new approach to learn NeRF-based models for UAV-perception, that can create novel-view images from previously unseen camera positions capturing salient attributes of the scene. Our approach is used to augment current datasets by generating synthetic images for improved pose and action recognition.

We present a tiered version of the K-Planes algorithm [9], TK-Planes, that outputs and operates on feature vectors rather than RGB pixel values. These feature vectors can store conceptual information about a particular object or location within a scene and when output for multiple corresponding camera rays, form feature maps that are decoded to form a final image. Furthermore, taking advantage of the fact that many details of a scene, including dynamic objects, repeat throughout, the image decoder which operates on all feature vectors can learn how to model static and dynamic objects within the scene.

We have evaluated the performance on some challenging UAV datasets, including Okutama Action [3] and UG2 [39], which consist of high altitude videos. We show that our NeRF-based model based on TK-planes can generate supplementing UAV-based data that can improve the overall recognition accuracy for dynamic scenes, as compared to state of the art algorithms [9], [21]. To summarize, some novel contributions from our work include:

- Development of a novel grid-based NeRF architecture that utilizes multiple scales of feature space in order to render high fidelity regions of a scene.
- Extension of previous pipelines for training grid-based NeRF architectures to include an image decoder and address edge effects from rendering image patches.
- Showcasing high fidelity renderings over previous methods such as K-Planes [9] and Extended K-Planes [21].

II. RELATED WORKS

Synthetic Data Generation. Collecting real-world dynamic dataset can be expensive, especially for challenging environments such as unmanned aerial vehicle (UAV)-based scenes. One way of generating synthetic dataset is utilization of game engines or simulators [28], [33], [2]. While these rendering pipelines can generate images in a fully controlled environment setting, they require large amount of efforts and skills for setup. Meanwhile, recent advances in neural rendering [44], [5] show surprising performance in novel view synthesis, which shows their feasibility to augment sparse image datasets with arbitrary rendering in a novel view [21], [12] and appearance [45], [41], or even with injected objects in the scene [41].

Neural Rendering for Dynamic Scenes. To represent 3D dynamic scenes, various NeRF research explores 4D neural radiance fields representation with an additional temporal dimension [6], [43], [10], [40]. Instead of embedding continuous time dimension into the radiance fields, some of variants uses time-varying latent code [18], [26]. While these 4D NeRF show capability to represent complicated dynamic scenes, they generally suffer from inherent representational redundancy when they learn 4D contents from 2D imagery only. The other major variants of NeRF incorporates canonical scene and a deformation field [25], [29]. However, these methods demonstrate limited capability to learn topological changes or large motion in the scene due to their heavy dependence on the canonical scene.

NeRF with Explicit Grid Representation Integration of 4D NeRF models with explicit voxel grid structures [14], [20], [26], [7] has been widely studied to boost slow training and rendering speed of implicit NeRF models. However, voxel-based methods have limitation in representing large-scale scenes due to their inefficient use of memory storage. One way to achieve compact explicit representation is factoring 4D dynamic volumes into multiple planar features planes or volumes [9], [4], [34], [11], [22], [42]. These tensor factorization techniques show good trade-off between compactness and high-quality representation. However, all of these methods primarily target dynamic scenes observed from ground-level perspectives, typically featuring a single person exhibiting transient motions close to the camera. UAV-Sim [21] shares a similar motivation as ours, utilizing extended K-planes for dynamic scene representation and addressing multiple human actors captured by UAVs. We enhance their method [21] by introducing higher-level features to augment the low-level detail of multiple human actors.

III. METHODOLOGY

A. Novel View Synthesis for Dynamic Scenes

Given a difficult domain such as UAV video with limited available training data and a small dynamic to static pixel ratio, our approach utilizes NeRF in feature space in order to better capture and render objects of interest within a scene, e.g. dynamic objects. We employ a grid-based Neural Radiance Field, the grids of which are similar to those from K-Planes [9]. A key difference being that the grids in our algorithm store feature vectors that do not directly code for RGB values but rather for more abstract high level concepts within the scene, such as patches of ground, trees, and people. Thus, our NeRF algorithm outputs a feature vector and density value for a given ray rather than an RGB and density value. The overall output is a feature map rather than an RGB image. Output feature maps are decoded by a convolutional neural network in order to generate a final rendered image.

B. Feature Plane NeRF Overview

To understand our algorithm in more detail, a brief overview of Neural Radiance Fields follows. Neural radiance fields (NeRF) are a type of volumetric rendering that utilize neural networks to output a color and density given an input camera pose. From a designated camera pose, rays are cast into a scene and incremental points along each ray are sampled. Typically, two Multi-layer Perceptrons (MLP) are employed, one for spatial density and the other for a pixel RGB value. Each sampled point’s location and direction are fed as input into the MLPs. The first MLP takes as input a position within 3D space and outputs a density value as well as a feature vector. The feature vector and the point’s direction are fed into the second MLP which outputs a RGB color value. The density and color at each sampled point along the ray are then multiplied, rectified and summed along the ray to get a final pixel color value for the given ray. This version of NeRF relies on the MLPs learning an implicit 3D structure of a scene.

Alternative to the implicit method, some NeRF algorithms implement an explicit representation of a scene stored in a data structure such as a matrix [15], [38] or a grid [9], [4]. These data structures contain feature vectors that are input into the spatial density MLP and learn explicit details wherein the data structure itself represents the spatial or temporal layout of a scene. By storing feature vectors in an explicit data structure, the MLPs used to output density and color values can be smaller and lightweight which coincides with reduced training and inference times.

C. Tiered Planes in Feature Space

Our main approach operates in a conceptual feature space rather than an RGB space with respect to the output of the NeRF algorithm. The purpose of focusing on feature space rather than RGB space is to learn feature vectors that represent cohesive objects of a scene, such as people, as opposed to learning solely a chromaticity and luminance value. This helps to alleviate the issue of poor fidelity of

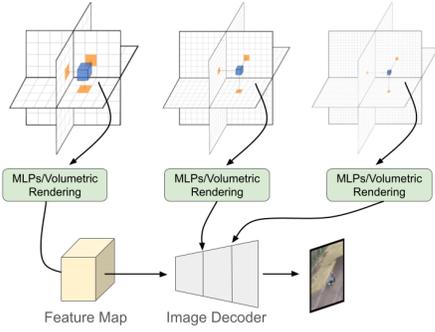


Fig. 1: Overview of our tiered planes algorithm: Three sets of grids are shown, representing three different scales of feature space, with larger scales capturing more abstract scene details. Note that the algorithm uses nine grids in total per scale, three static spatial, three dynamic spatial and three dynamic spatio-temporal. Feature vectors are processed via concurrent MLPs to output a density value and final feature vector. These are accumulated with volumetric rendering into a feature map. Feature maps from each set of grids are input into subsequent stages of the image decoder.

high interest objects across different poses and time-slices. We use the same grid representation as seen by K-Planes [9] with changes influenced from the extended version as seen in [21]. Figure 1 shows our system in full with varying scale feature space grids as well as the image decoder.

In particular, we add the additional planes for spatial features that are dynamic as well as the cosine similarity loss function that encourages separation of static and dynamic features in the spatial planes. The resulting static and dynamic features are described as follows:

$$f_s(\mathbf{q}) = \prod_{c \in C_s} f_s(\mathbf{q})_c, \quad f_d(\mathbf{q}) = \prod_{c \in C_d} f_d(\mathbf{q})_c. \quad (1)$$

$$f(\mathbf{q}) = f_s(\mathbf{q}) \oplus f_d(\mathbf{q}) \quad (2)$$

wherein f_s is a function of the static spatial planes and f_d is a function of the dynamic spatial and temporal planes and \prod represents the Hadamard product. In our version, each feature vector $f_*(\mathbf{q})_c$ is a D dimensional vector wherein D is a hyperparameter, which differs from [21] in that it does not include a learned mask value for the temporal feature vectors. In focusing on rendering quality over data generation, we can omit the mask dimension to simplify the model. From 2, a cohesive feature vector $f(\mathbf{q})$ is calculated by concatenating f_s and f_d . $f(\mathbf{q})$ is then input into the MLP networks to output a density value and another feature vector. Volumetric rendering along the ray returns a final feature vector value per pixel of an output feature map at a given feature space scale, d_s . Furthermore, each set of grids reduces the dimensionality of their feature vectors by two compared to the previous grids. In our experiments we choose a starting D of 256 such that vectors in the first stage of grids can learn complex and abstract details of a scene. The feature map from these grids is fed into the start of the image decoder. Subsequent grid sets then have feature vectors of size 128 and 64 and are

fed into later stages of the decoder, as explained below. The purpose of this is to focus on more fine grained detail in the later grids which does not require large vector lengths. This also allows to grid feature map outputs to integrate with an image decoder in which the spatial dimensions are doubled each stage, but the depth dimension is halved.

K-Planes [9] and its extended version [21] use multi resolution grids offering different levels of scene details. The resulting feature vector from each level can either be summed or concatenated before being processed by the MLPs and then accumulated via volumetric rendering. In these cases, the volumetric rendering is the final step. However, in our algorithm it is more intuitive to separate the different levels of detail as input into appropriate stages of the decoder rather than as input entirely into the front of the decoder. Convolutional image decoders operate on high-level abstract features to low-level detail features as their stages progress and typically benefit from feature map input at varying stages [32]. This allows the grids at varying scales of feature space to learn appropriate details with respect to rendering an output image.

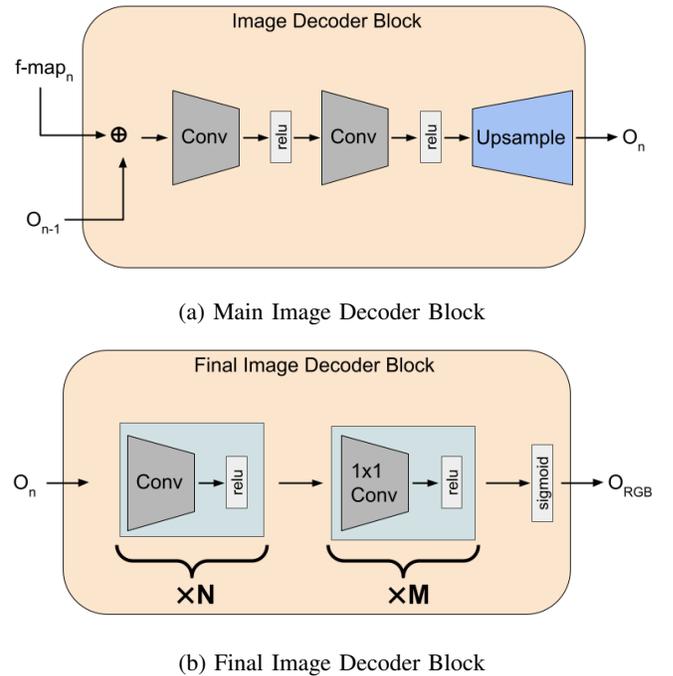


Fig. 2: a.) A diagram of a block within the image decoder. The n^{th} block accepts two inputs, a feature map from the n^{th} set of grids after volumetric rendering and the output O_{n-1} from the previous block. At $n = 0$, there is no input from a previous block, only the feature map. b.) The final image decoder block with xN 3×3 convolutional layers followed by xM 1×1 convolutional layers and a sigmoid activation layer to rectify the raw output into RGB values.

D. Image Decoder

The set of feature maps d_s at varying scales, d , are input to an image decoder that consists of convolutional blocks

followed by bilinear upsampling. Figure 2 details a block of the image decoder. Each block takes as input a corresponding feature map and, if applicable, the output of the previous block. These two inputs are concatenated together before further processing. The output of each block has spatial dimensions upsampled by two and feature dimensions downsampled by two. The number of blocks is equivalent to the number of scaling levels of the grids, in our case three. A final block consisting of four convolutional layers further processes the resultant output from all feature maps without further input. This is followed by a series of 1x1 convolutions that reduce the feature dimensions of each pixel to a raw RGB value and then output a rectified color per pixel via a sigmoid activation layer.

E. Training Details

Due to the use of an image decoder, rays must be processed in patches rather than sampled individually throughout the image. This leads to a potential complication from edge effects within the image decoder because of the down-sampling nature of convolutional layers and the padding that may compensate for it. Because we are rendering an image based on patches and not in its entirety due to memory constraints, a simple padding scheme should not be used when training or during inference. If using padding in the convolutional layers, discontinuities will appear during training and inference that affect the quality of the final images. In order to adjust for this, we over-sample our patches and do not utilize padding within the convolutional layers of the image decoder. This over-sampling effectively serves as padding for the desired final shape after decoding except that actual grid values are used rather than approximations such as zero-padding or reflectance. In this case the over-sampling must be tailored to the architecture of the image decoder itself such that the final dimensions of the output image are appropriate. Given our architecture, we over-sample by eight "pixels" in each spatial dimension for each feature map.

We use a NVIDIA RTX 3090 GPU for training with a learning rate of $1e - 5$ for both the feature grids as well as the image decoder. The training is run for 500,000 iterations with a patch from one image per training iteration. Patches are sized so as to be an integer divisible by the width and height of the actual image while also being evenly divisible by two for each scale. In our case we chose a patch size of 360x640 for every experiment except Okutama video 1.1.1, in which case we chose a smaller patch size of 64x72. In most cases the larger patch size showed considerable performance improvements over using smaller patch sizes III.

IV. RESULTS

A. Datasets

1) *Okutama-Action Dataset*: The Okutama-Action dataset [3] is collected using unmanned aerial vehicles (UAVs). It comprises scenes captured at various altitudes, angles, and time periods, containing different humans in 12 different action classes. This dataset presents challenges due

TABLE I: PSNR Comparison on Okutama-Action and UG2 Validation Subsets

Okutama	K-Planes	Extended K-Planes	TK-Planes
1.1.1	18.23	32.61	31.17
1.2.2	14.45	32.63	32.93
1.2.6	14.37	27.89	28.97
UG2	K-Planes	Extended K-Planes	TK-Planes
UAV-133	17.11	27.56	27.88

TABLE II: Dynamic PSNR Comparison on Okutama-Action and UG2 Validation Subsets

Okutama	K-Planes	Extended K-Planes	TK-Planes
1.1.1	20.55	24.07	27.36
1.2.2	15.63	24.52	26.47
1.2.6	16.31	20.81	21.76
UG2	K-Planes	Extended K-Planes	TK-Planes
UAV-133	15.54	17.51	16.91

to dynamic transitions between actions, humans of multiple concurrent actions, and actors labeled with multiple actions. We train on the three different scenes: 1.1.1, 1.2.2, 1.2.6. Each of these scenes offers distinct views of a baseball diamond at different times of day with varying numbers of people present, see Figures 7a, 7b, 7c.

2) *UG2 Dataset*: We also sample a video from the UG2 dataset [39] which consists of a large variety of videos from varying styles of aerial vehicles, including piloted aircraft and UAVs such as fixed wing gliders and quadcopters. Although many of the videos within UG2 are not applicable for this given application, video 133 within the test set for UAV videos offers a challenging view of a cricket match in progress. This particular video contrasts well with the scenes from Okutama-Action, offering smaller and more dynamic targets moving with less precise and scripted actions 7d.

B. Analysis

Because we are focused on dynamic portions of a scene, we calculate not only overall PSNR values but also what we call Dynamic-PSNR values which represent the average PSNR between each bounding box of people within each frame. With respect to tasks such as object detection, pose recognition and action recognition, the fidelity of the target object is key and thus the most pertinent for validation. As can be seen from Tables I and II, the overall performance of TK-Planes beats stock and extended K-Planes in most experiments. Furthermore, Table III shows how important

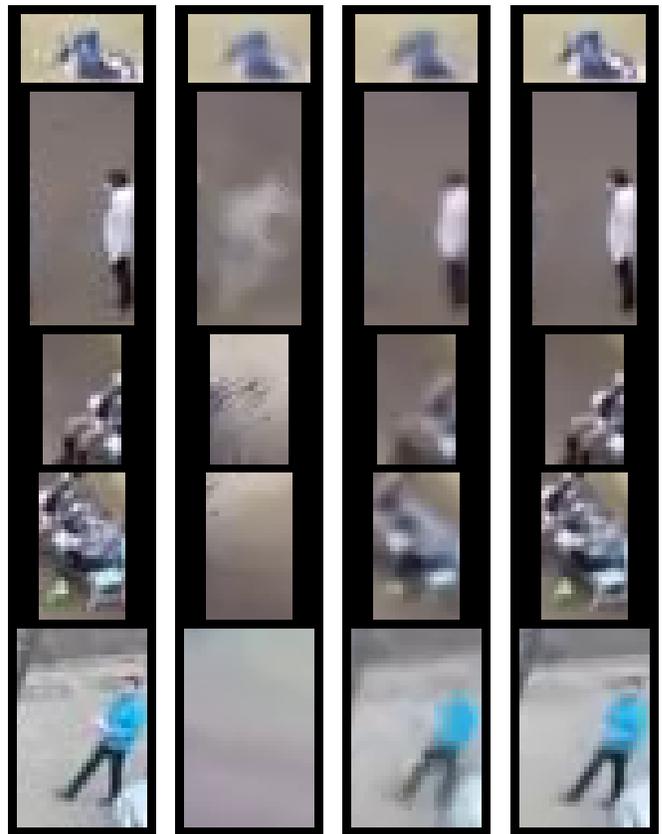
TABLE III: Dynamic PSNR for Small and Large Patch-Size on Okutama-Action and UG2 Validation Subsets for TK-Planes

Okutama	64x72	360x640
1.1.1	27.36	25.20
1.2.2	22.59	26.47
1.2.6	18.40	21.76
UG2	64x72	360x640
UAV-133	15.38	16.91

hyperparameter tuning is with respect to TK-Planes. In all but the first experiment, larger patch-sizes provides upwards of 3.88 points of improvement for Dynamic-PSNR. Given that TK-Planes has more hyperparameters to adjust than stock or extended K-Planes, there is a larger space for exploration in achieving improved results which indicates that TK-Planes may further improve PSNR values compared to stock and extended K-Planes.

The ability to utilize an image decoder and to store more abstract feature vectors within the grids allows for the network to learn *how* to render a person versus rendering individual pixel values. In the latter case such as seen in Extended K-Planes, a higher level of noise is captured in the renderings while also losing some of the detail and sharpness of smaller portions of the people. Figures 3d, 4d, and 5d in particular shows enhanced detail from TK-Planes with respect to limbs and clothing detail. One of the main benefits from this is in rendering synthetic data for the previously mentioned tasks that require specific limb placement and outlines, such as action recognition.

Samples from the Okutama experiments and Tables I and II shows that TK-Planes can ultimately render higher fidelity imagery when compared to Extended K-Planes, and especially when compared to standard K-Planes which fails to render details in many cases. However, the final experiment from UG2 of the UAV cricket scene shows more mixed results. The overall PSNR value for TK-Planes is higher when compared to stock and extended K-Planes, but the Dynamic-PSNR value lower when compared to extended K-Planes. This scene is particularly difficult with small dynamic portions and large unbounded backgrounds. Figure 6 compares renderings of people from all three methods with that of the ground truth image from UG2. Although the background is largely captured, the dynamic portions of the scene are very small resulting in a lack of detail with respect to the people in the scene. The top two rows show extended K-Planes rendering more accurate people while the bottom two rows show TK-Planes renderings to be more accurate. Due to the difficulty of this scene, none of the methods perform particularly well, but TK-Planes has shown that hyperparameter tuning can play a large role in its overall



(a) Okutama (b) K-Planes (c) Extended (d) TK-Planes

Fig. 3: A comparison of dynamic regions from a subsection of video 1.1.1 in Okutama-Action. a.) Okutama ground truth b.) stock K-Planes c.) Extended K-Planes d.) TK-Planes. We highlight the improved rendering quality generated by TK-Planes over the other methods on these challenging dynamic scenes.

performance which leads to the intuition that TK-Planes is not optimized for this particular scene and may yet achieve better performance.

C. Failure Modes and Limitations

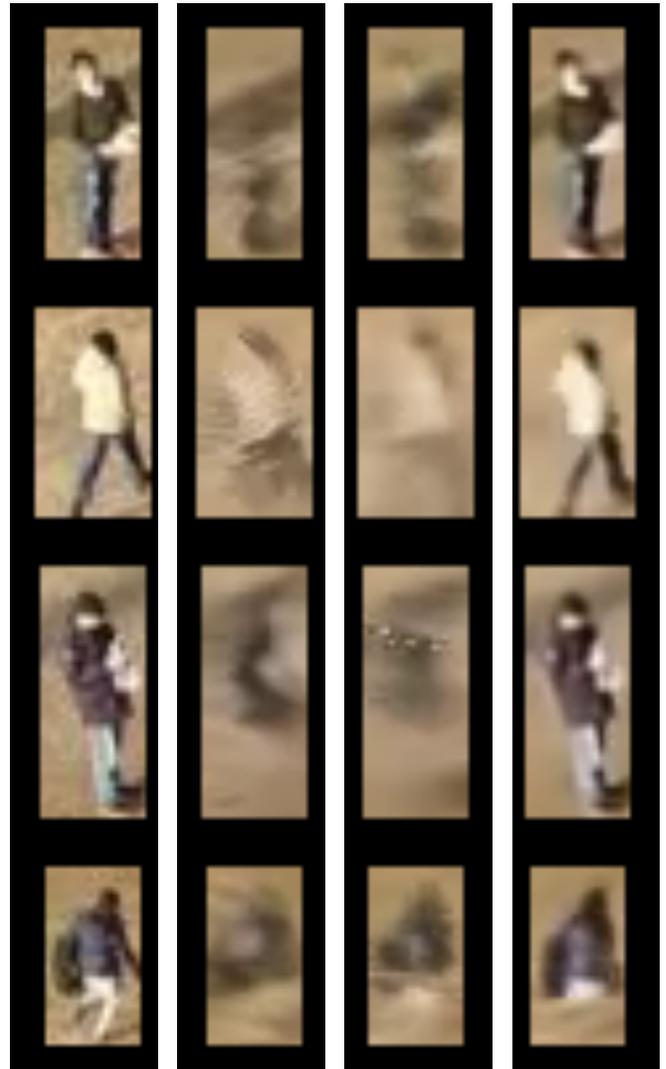
In order to explain the failure cases such as seen in Figure 6d, there are a variety of factors that contribute to this dichotomy between Extended and TK-Planes. First, the training method for TK-Planes has a larger number of hyperparameters, the most likely of which to affect specific frames having portions of lower fidelity being the sampling rate of dynamic portions of the scene and the overall patch size per image during training. Second, each scene for Okutama and UG2 does not provide ground truth camera poses which are a necessity for many NeRF algorithms. To compensate, we run advanced structure from motion algorithms [19] to estimate the camera poses for each image within a scene. It is possible that these estimated poses are large enough confounding factors between scenes with respect to the overall PSNR of dynamic portions of a scene.



(a) Okutama (b) K-Planes (c) Extended (d) TK-Planes

Fig. 4: A comparison of dynamic regions from a subsection of video 1.2.2 in Okutama-Action. a.) Okutama ground truth b.) stock K-Planes c.) Extended K-Planes d.) TK-Planes. We highlight the improved rendering quality generated by TK-Planes over the other methods on these challenging dynamic scenes.

A further limitation that is also seen in Extended K-Planes is that of limited novelty beyond training poses and time-slices. In Extended K-Planes, the further from the training poses and time-slices a rendering pose and time departed, the more distorted the rendered result would be. A similar limitation is also present within TK-Planes. This is ultimately due to the difficult restrictions of the UAV datasets that have tight correlations between poses and time-slices, as well as limited diversity in overall poses at varying time-slices.



(a) Okutama (b) K-Planes (c) Extended (d) TK-Planes

Fig. 5: A comparison of dynamic regions from a subsection of video 1.2.6 in Okutama-Action. a.) Okutama ground truth b.) stock K-Planes c.) Extended K-Planes d.) TK-Planes. We highlight the improved rendering quality generated by TK-Planes over the other methods on these challenging dynamic scenes.

V. CONCLUSION AND FUTURE WORK

We have presented a novel method for synthetic data generation. Our tiered K-Planes algorithm has shown that it can improve the fidelity of dynamic objects within difficult to render scenes such as those from UAVs. In certain cases this improvement is upwards of 3 points in Dynamic-PSNR. Although there are failure cases in some instances, this can easily be explained by a lack of proper hyperparameter grid search with respect to patch size, feature vector dimension and grid resolution, amongst other parameters. Overall, the increase in fidelity of dynamic objects will help with tasks such as object detection, pose recognition and action recognition wherein the details of the subjects are paramount.

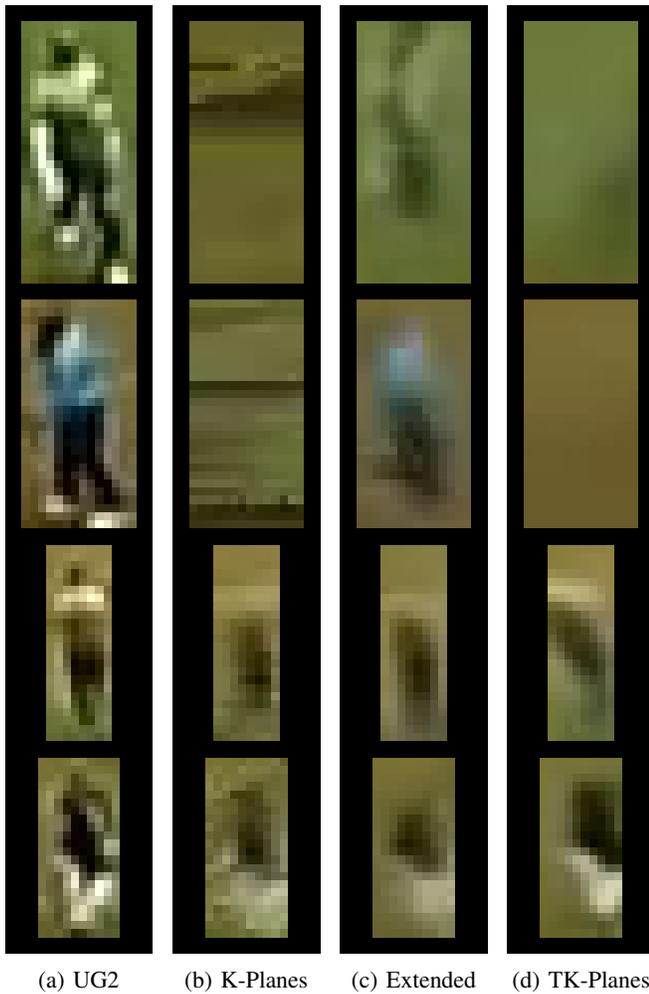


Fig. 6: A comparison of dynamic regions from a subsection of video UAV-133 in UG2. a.) UG2 ground truth b.) stock K-Planes c.) Extended K-Planes d.) TK-Planes. The UG2 dataset is particularly challenging, extended K-Planes performs better in some cases while TK-Planes in others.

Future work with respect to TK-Planes involves a more robust grid search of hyperparameters given particular scenes and environmental conditions. Furthermore, the use of abstract feature vectors can be exploited to take advantage of redundant information within a scene and improve to a greater degree the fidelity of images and novelty of poses and time-slices of dynamic scenes. It may be useful to design some hybrid methods that can combine Extended K-Planes and TK-Planes to further improve the recognition accuracy.

REFERENCES

[1] Sora: Creating video from text, 2024. <https://openai.com/research/video-generation-models-as-world-simulators>.

[2] Brendan Alvey, Derek T Anderson, Andrew Buck, Matthew Deardorff, Grant Scott, and James M Keller. Simulated photorealistic deep learning framework and workflows to accelerate computer vision and unmanned aerial vehicle research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3889–3898, 2021.

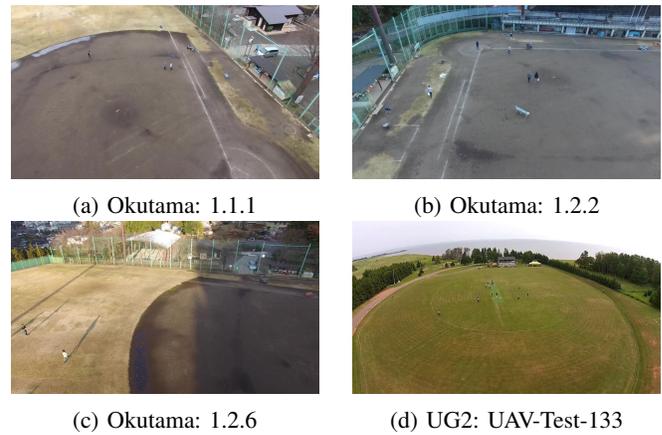


Fig. 7: A example frame from each of the four scenes used to validate our experiments. Each scene from Okutama-Action offers a unique viewpoint, time-of-day and number of people. The scene from UG2 offers a challenging viewpoint from a quadcopter of a cricket match in session.

[3] Mohammadamin Barekatin, Miquel Martí, Hsueh-Fu Shih, Samuel Murray, Kotaro Nakayama, Yutaka Matsuo, and Helmut Prendinger. Okutama-action: An aerial view video dataset for concurrent human action detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 28–35, 2017.

[4] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2023.

[5] Guikun Chen and Wenguan Wang. A survey on 3d gaussian splatting. *arXiv preprint arXiv:2401.03890*, 2024.

[6] Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14324–14334, 2021.

[7] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia 2022 Conference Papers*, 2022.

[8] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 203–213, 2020.

[9] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023.

[10] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.

[11] Quankai Gao, Qiangeng Xu, Hao Su, Ulrich Neumann, and Zexiang Xu. Strivec: Sparse tri-vector radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17569–17579, 2023.

[12] Yunhao Ge, Harkirat Behl, Jiashu Xu, Suriya Gunasekar, Neel Joshi, Yale Song, Xin Wang, Laurent Itti, and Vibhav Vineet. Neural-sim: Learning to generate training data with nerf. In *European Conference on Computer Vision*, pages 477–493. Springer, 2022.

[13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

[14] Xiang Guo, Guanying Chen, Yuchao Dai, Xiaoqing Ye, Jiadai Sun, Xiao Tan, and Errui Ding. Neural deformable voxel grid for fast optimization of dynamic view synthesis. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*, 2022.

[15] Xiang Guo, Guanying Chen, Yuchao Dai, Xiaoqing Ye, Jiadai Sun, Xiao Tan, and Errui Ding. Neural deformable voxel grid for fast

- optimization of dynamic view synthesis. In *Proceedings of the Asian Conference on Computer Vision*, pages 3757–3775, 2022.
- [16] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.
- [17] Tianjiao Li, Jun Liu, Wei Zhang, Yun Ni, Wenqian Wang, and Zhiheng Li. Uav-human: A large benchmark for human behavior understanding with unmanned aerial vehicles. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16266–16275, 2021.
- [18] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhaoyang Lv. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5521–5531, 2022.
- [19] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-perfect structure-from-motion with featuremetric refinement. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5987–5997, 2021.
- [20] Jia-Wei Liu, Yan-Pei Cao, Weijia Mao, Wenqiao Zhang, David Junhao Zhang, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Devrf: Fast deformable voxel radiance fields for dynamic scenes. *arXiv preprint arXiv:2205.15723*, 2022.
- [21] Christopher Maxey, Jaehoon Choi, Hyungtae Lee, Dinesh Manocha, and Heesung Kwon. Uav-sim: Nerf-based synthetic data generation for uav-based perception. *arXiv preprint arXiv:2310.16255*, 2023.
- [22] Marko Mihajlovic, Sergey Prokudin, Marc Pollefeys, and Siyu Tang. Resfields: Residual neural fields for spatiotemporal signals. *International Conference on Learning Representations*, 2024.
- [23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [24] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [25] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021.
- [26] Sida Peng, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Representing volumetric videos as dynamic mlp maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [27] Justin Plowman. *3D Game Design with Unreal Engine 4 and Blender*. Packt Publishing Ltd, 2016.
- [28] Justin Plowman. *3D Game Design with Unreal Engine 4 and Blender*. Packt Publishing Ltd, 2016.
- [29] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [30] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [32] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015.
- [33] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and Service Robotics: Results of the 11th International Conference*, pages 621–635. Springer, 2018.
- [34] Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boming Liu, Hongwen Zhang, and Yebin Liu. Tensor4d: Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16632–16642, 2023.
- [35] Yi-Ting Shen, Hyungtae Lee, Heesung Kwon, and Shuvra S Bhattacharyya. Progressive transformation learning for leveraging virtual images in training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 835–844, 2023.
- [36] Yi-Ting Shen, Yaesop Lee, Heesung Kwon, Damon M Conover, Shuvra S Bhattacharyya, Nikolas Vale, Joshua D Gray, G Jeremy Leong, Kenneth Evensen, and Frank Skirlo. Archangel: A hybrid uav-based human detection benchmark with position and pose metadata. *IEEE Access*, 2023.
- [37] Zhihang Song, Zimin He, Xingyu Li, Qiming Ma, Ruibo Ming, Zhiqi Mao, Huaxin Pei, Lihui Peng, Jianming Hu, Danya Yao, et al. Synthetic datasets for autonomous driving: A survey. *IEEE Transactions on Intelligent Vehicles*, 2023.
- [38] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022.
- [39] Rosaura G Vidal, Sreya Banerjee, Klemen Grm, Vitomir Struc, and Walter J Scheirer. Ug²: A video benchmark for assessing the impact of image restoration and enhancement on automatic visual recognition. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1597–1606. IEEE, 2018.
- [40] Chaoyang Wang, Ben Eckart, Simon Lucey, and Orazio Gallo. Neural trajectory fields for dynamic novel view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [41] Liao Wang, Qiang Hu, Qihan He, Ziyu Wang, Jingyi Yu, Tinne Tuytelaars, Lan Xu, and Minye Wu. Neural residual radiance fields for streamably free-viewpoint videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 76–87, 2023.
- [42] Minye Wu, Zehao Wang, Georgios Kourou, and Tinne Tuytelaars. Tettrif: Temporal tri-plane radiance fields for efficient free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [43] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9421–9431, 2021.
- [44] Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond. In *Computer Graphics Forum*, volume 41, pages 641–676. Wiley Online Library, 2022.
- [45] Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. Unisim: A neural closed-loop sensor simulator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1389–1399, June 2023.