# GEOMETRIC ALGEBRA PLANES: CONVEX IMPLICIT NEURAL VOLUMES

## **Anonymous authors**

Paper under double-blind review

# Abstract

Volume parameterizations abound in recent literature, from the classic voxel grid 1 to the implicit neural representation and everything in between. While implicit 2 representations have shown impressive capacity and better memory efficiency 3 compared to voxel grids, to date they require training via nonconvex optimization. 4 This nonconvex training process can be slow to converge and sensitive to initial-5 ization and hyperparameter choices. We introduce a family of models, GA-Planes, 6 that is the first class of implicit neural volume representations that can be trained 7 by convex optimization. GA-Planes models include any combination of features 8 stored in tensor basis elements, followed by a neural feature decoder. They gen-9 eralize many existing representations and can be adapted for convex, semiconvex, 10 or nonconvex training as needed for different inverse problems. In the 2D setting, 11 we prove that GA-Planes is equivalent to a low-rank plus low-resolution matrix 12 factorization; we show that this approximation outperforms the classic low-rank 13 plus sparse decomposition for fitting a natural image. In 3D, we demonstrate 14 GA-Planes' competitive performance in terms of expressiveness, model size, and 15 optimizability across three volume fitting tasks: radiance field reconstruction, 3D 16 segmentation, and video segmentation. 17

# 18 1 INTRODUCTION

Volumes are everywhere—from the world we live in to the videos we watch to the organs and 19 tissues inside our bodies. In recent years tremendous progress has been made in modeling these 20 volumes using measurements and computation (Tewari et al., 2022), to make them accessible for 21 downstream tasks in applications including manufacturing (Intwala & Magikar, 2016; Šlapak et al., 22 2024), robotic navigation (Ming et al., 2024; Wijayathunga et al., 2023), entertainment and culture 23 (Liu et al., 2024; Croce et al., 2023), and medicine (Udupa & Herman, 1999; Masero et al., 2002; 24 Richter et al., 2024; Xu et al., 2024). All methods that seek to model a volume face a three-way 25 tradeoff between model size, which determines hardware memory requirements, expressiveness, 26 which determines how faithfully the model can represent the underlying volume, and optimizability, 27 which captures how quickly and reliably the model can learn the volume from measurements. Cer-28 tain applications place stricter requirements on model size (e.g. for deployment on mobile or edge 29 devices), expressiveness (e.g. resolution required for medical diagnosis or safe robotic navigation), 30 or optimizability (e.g. for interactive applications), but all stand to benefit from improvements to 31 this three-way pareto frontier. 32

Many existing strategies have been successfully applied at different points along this pareto frontier; some representative examples from computer vision are summarized in Appendix A.1. Our goal is to maintain or surpass the existing pareto frontier of model size and expressiveness while improving optimization stability through convex optimization.

Our approach introduces *convex* and *semiconvex* reformulations of the volume modeling optimization process that apply to a broad class of volume models we call *Geometric Algebra Planes*, or GA-Planes for short. We adopt the term *semiconvex* for Burer-Monteiro (BM) factorizations of a convex objective, as introduced in Sahiner et al. (2024), within the context of convex neural networks. BM factorized problems have the property that every local minimum is globally optimal (Sahiner et al., 2024).

GA-Planes is a mixture-of-primitives model that generalizes several existing volume models includ-43 44 ing voxels and tensor factorizations. Most importantly, most models in this family can be formulated for optimization by a convex program, as long as the objective function (to fit measurements of the 45 volume) is convex. At the same time, any GA-Planes model can also be formulated for nonconvex 46 optimization towards any objective, matching the range of applicability enjoyed by common mod-47 els. While only our convex and semiconvex models come with guarantees of convergence to global 48 optimality, all the models we introduce extend the pareto frontier of model size, expressiveness, and 49 optimizability on diverse tasks. 50

- 51 Concretely, we make the following contributions:
- We introduce GA-Planes, a mixture-of-primitives volume parameterization inspired by ge ometric algebra basis elements. GA-Planes combines any subset of line, plane, and volume
   features at different resolutions, with an MLP decoder. This GA-Planes family of parame terizations generalizes many existing volume and radiance field models.
- We derive convex and semiconvex reformulations of the GA-Planes training process for
   certain tasks and a large subset of the GA-Planes model family, to ensure our model opti mizes globally regardless of initialization.
- We analyze GA-Planes in the 2D setting and show equivalence to a low-rank plus low-resolution matrix approximation whose expressiveness can be directly controlled by design choices. We demonstrate that this matrix decomposition is expressive for natural images, outperforming the classic low-rank plus sparse approximation.
- We demonstrate convex, semiconvex, and nonconvex GA-Planes' high performance in terms of memory, expressiveness, and optimizability across three volume-fitting tasks: 3D radiance field reconstruction, 3D segmentation, and video segmentation.

# 66 2 RELATED WORK

Volume parameterization. Many volume parameterizations have been proposed and enjoy
widespread use across diverse applications. Here we give an overview of representative methods
used in computer vision, focusing on methods that parameterize an entire volume (rather than e.g. a
surface). These parameterizations achieve different tradeoffs between memory usage, representation
quality, and ease of optimization; richer descriptions are provided in Appendix A.1.

Coordinate MLPs like NeRF (Mildenhall et al., 2020) and Scene Representation Networks (Sitz-72 mann et al., 2019b) are representative of Implicit Neural Representations (INRs), which excel at 73 reducing model size (with decent expressiveness) but suffer from slow optimization. At the oppo-74 site end of the spectrum, explicit voxel grid representations like Plenoxels (Sara Fridovich-Keil and 75 Alex Yu et al., 2022) and Direct Voxel Grid Optimization (Sun et al., 2022) can optimize quickly but 76 require large model size to achieve good expressiveness (resolution). Many other methods (Chen 77 et al., 2022; Fridovich-Keil et al., 2023; Müller et al., 2022; Kerbl et al., 2023; Reiser et al., 2023; 78 Lombardi et al., 2021) find their niche somewhere in between, achieving tractable model size, good 79 expressiveness, and reasonably fast optimization time in exchange for some increased sensitivity (to 80 initialization, randomness, and prior knowledge) in the optimization process. GA-Planes matches or 81 exceeds the performance of strong baselines (Chen et al., 2022; Fridovich-Keil et al., 2023; Barron 82 et al., 2021) in terms of model size and expressiveness, while introducing the option to train by 83 convex or semiconvex optimization with guaranteed convergence to global optimality. 84

Radiance field modeling. Most of the works described above are designed for the task of model-85 ing a radiance field, in which the training measurements consist of color photographs from known 86 camera poses. The goal is then to faithfully model the optical density and view-dependent color of 87 light inside a volume so that unseen views can by rendered accurately. This task is also referred to as 88 novel view synthesis (Mildenhall et al., 2020; Sitzmann et al., 2019b). Although we do demonstrate 89 superior performance of GA-Planes in this setting, we note that the volumetric rendering formula 90 used in radiance field modeling (Max, 1995; Kajiya, 1986; Mildenhall et al., 2020) yields a noncon-91 vex photometric loss function, regardless of model parameterization. 92

3D segmentation. We test our convex and semiconvex GA-Planes parameterizations on fully con vex objectives, namely volume (xyz) segmentation with either indirect 2D tomographic supervision



**Figure 1:** Overview of the GA-Planes models we use in our experiments. Our nonconvex model (top) uses a standard MLP decoder and multiplication of features when the result yields a volume under geometric algebra; it also concatenates features across mult-resolution grids. Our semiconvex (middle) and convex (bottom) models use a single resolution for each feature grid, and avoid multiplication of features since that would induce nonconvexity. The pastel-colored grids inside the indicator function of the convex model are frozen at initialization and used as fixed ReLU gating patterns.  $\odot$  denotes concatenation and  $\circ$  denotes elementwise multiplication.

or direct supervision, as well as video (xyt) segmentation with direct 3D supervision. This 3D 95 (xyz) segmentation task has also been studied in recent work (Cen et al., 2023; Uy et al., 2023), 96 though these methods require additional inputs such as a pretrained radiance field model or monoc-97 ular depth estimator. Our setup is most similar to Mescheder et al. (2019), which uses an implicit 98 neural representation trained with cross-entropy loss and direct 3D supervision of the occupancy 99 100 function. Instead of having direct access to this 3D training data, we infer 3D supervision labels via Space Carving (Kutulakos & Seitz, 1999) from 2D image masks obtained by image segmentation 101 (via Kirillov et al. (2023)). 102

**Convex neural networks.** Recent work has exposed an equivalence between training a shallow (Pilanci & Ergen, 2020) or deep (Ergen & Pilanci, 2024) neural network and solving a convex program whose structure is defined by the architecture and parameter dimensions of the corresponding neural network. The key idea behind this convexification procedure is to enumerate (or randomly sample from) the possible activation paths through the neural network, and then treat these paths as a fixed dictionary whose coefficients may be optimized according to a convex program.

Geometric (Clifford) algebra. Geometric algebra (GA) is a powerful framework for modeling 109 geometric primitives and interactions between them (Dorst et al., 2009). The fundamental entity in 110 GA is the multivector, which is a sum of vectors, bivectors, trivectors, etc. In 3D GA, an example is 111 the multivector  $e_1e_2 + e_1e_2e_3$ , representing the sum of a bivector (a plane) and a trivector (a vol-112 ume). The geometric product in GA allows us to derive a volume element by multiplying a plane and 113 a line, e.g.  $(e_1e_2)e_3 = e_1e_2e_3$ . We use the shorthand  $e_{123} = e_1e_2e_3$ , and similarly for other multi-114 vector components throughout. Inspired by this framework, we define the GA-Planes model family 115 to include any volume parameterization that combines any subset (including the complete subset 116 and the empty subset) of the linear geometric primitives  $\{e_1, e_2, e_3\}$ , planar geometric primitives 117  $\{e_{12}, e_{13}, e_{23}\}$ , and/or volumetric primitive  $\{e_{123}\}$  with a (potentially convexified) MLP feature 118 decoder. We leverage geometric algebra to combine these primitives into a trivector (volume). To 119 our knowledge, this work is the first to use geometric algebra in neural volume models. 120

# 121 3 MODEL

## 122 3.1 THE GA-PLANES MODEL FAMILY

<sup>123</sup> A GA-Planes model represents a volume using a combination of geometric algebra features  $e_c$ <sup>124</sup> derived by interpolating the following parameter grids:

• Line (1-dimensional) feature grids  $\{g_1, g_2, g_3\}$ , where each grid has shape  $[r_1, d_1]$  with spatial resolution  $r_1$  and feature dimension  $d_1$ . • Plane (2-dimensional) feature grids  $\{g_{12}, g_{13}, g_{23}\}$ , where each grid has shape  $[r_2, r_2, d_2]$ with spatial resolution  $r_2$  and feature dimension  $d_2$ .

• A single volume feature grid  $\{\mathbf{g}_{123}\}$  with shape  $[r_3, r_3, r_3, d_3]$ .

A GA-Planes model may include multiple copies of a given basis element with different resolution and feature dimensions, to effectively capture multi-resolution signal content. The x, y, and z spatial resolutions of each grid may differ in practice; for simplicity of notation we use isotropic resolutions.

We first extract features corresponding to  $q = (x, y, z) \in \mathbb{R}^3$  from each of our line, plane, and volume feature grids  $\mathbf{g}_c$  by linear, bilinear, and trilinear interpolation, respectively:

$$\mathbf{e}_c := \psi(\mathbf{g}_c, \pi_c(q)),\tag{1}$$

where  $\pi_c$  projects q onto the coordinates of the c'th feature grid  $\mathbf{g}_c$  and  $\psi$  denotes (bi/tri)linear 135 interpolation. The resulting feature  $\mathbf{e}_c$  is a vector of length  $d_1$  if  $c \in \{1, 2, 3\}, d_2$  if  $c \in \{12, 13, 23\}$ 136 or  $d_3$  if c = 123. We repeat this projection and interpolation procedure over each  $\mathbf{g}_c$ , and combine 137 the resulting feature vectors by any combination of elementwise multiplication ( $\circ$ ), addition (+), 138 139 and concatenation  $(\odot)$  along the feature dimension. Finally, the combined feature vector is decoded using an MLP decoder D. The decoder can take as input both the feature vector arising from the 140 feature grids as well as possible auxiliary inputs, such as (positionally encoded) viewing direction. 141 We consider any model that fits the above description to fall into the GA-Planes family. The specific 142 models we use for nonconvex, semiconvex, and convex optimization are illustrated in Figure 1. 143

Our experiments focus primarily on two specific GA-Planes models that exemplify some of the strongest convex and nonconvex representations in the GA-Planes family. For our experiments including convex optimization, namely 3D segmentation with 2D or 3D supervision, and video segmentation, we use the following GA-Planes model (illustrated in the second and third rows of Figure 1) which can be trained by either convex, semiconvex, or nonconvex optimization as described in the following subsections:

$$\mathbf{D}(\mathbf{e}_1 \odot \mathbf{e}_2 \odot \mathbf{e}_3 \odot \mathbf{e}_{12} \odot \mathbf{e}_{13} \odot \mathbf{e}_{23} \odot \mathbf{e}_{123}). \tag{2}$$

Here we use  $\odot$  to denote concatenation of features. For our radiance field experiments, since the

objective function is inherently nonconvex, we use the following nonconvex member of the GA-

Planes family (illustrated with multiresolution feature grids in the first row of Figure 1):

$$D((\mathbf{e}_1 \circ \mathbf{e}_2 \circ \mathbf{e}_3) \odot (\mathbf{e}_1 \circ \mathbf{e}_{23}) \odot (\mathbf{e}_2 \circ \mathbf{e}_{13}) \odot (\mathbf{e}_3 \circ \mathbf{e}_{12}) \odot \mathbf{e}_{123}), \tag{3}$$

which leverages geometric algebra to multiply ( $\circ$ ) lower-dimensional (vector and bivector) features together into 3D volume (trivector) features, but cannot be convexified because of this multiplication. We use multi-resolution copies of the line and plane feature grids  $\mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3, \mathbf{g}_{12}, \mathbf{g}_{13}, \mathbf{g}_{23}$ , but only a single resolution for the volume grid  $\mathbf{g}_{123}$  since it is already at lower resolution. For our nonconvex experiments the decoder D is a standard fully-connected ReLU neural network; decoder details for our semiconvex and convex models are presented in the following subsections.

### 159 3.2 SEMICONVEX GA-PLANES

For our segmentation experiments (with volumes and videos), we use the GA-Planes architecture in eq. (2), with concatenation instead of multiplication of features; we denote this concatenated feature vector as f(q), the input to the decoder. Our semiconvex formulation of this model uses a convex MLP (Pilanci & Ergen, 2020) as the decoder:

$$\tilde{y}(q) = \sum_{i=1}^{h} (W_i^{\top} f(q)) \mathbb{1}[\overline{W_i}^{\top} f(q) \ge 0].$$
(4)

Here W denotes the trainable hidden layer MLP weights, and  $\overline{W}$  denotes the same weights frozen 164 at initialization inside the indicator function. The indicator function, denoted as 1[\*], returns 1 if 165 the argument is true, and 0 otherwise. The indicator function here serves as a random gating pattern 166 that takes the place of the ReLU in a standard nonconvex MLP, where the gating pattern would be 167 optimized rather than fixed in a random pattern. Although this MLP decoder is fully convex, we 168 refer to this model as semiconvex (in particular biconvex; see (Sahiner et al., 2024)) because the 169 combined grid features f(q) are multiplied by the trainable MLP hidden layer weights W, though 170 the objective is separately convex in each of these parameters. 171

#### 172 3.3 CONVEX GA-PLANES

For our segmentation experiments (with volumes and videos), we also present a fully convex GA-Planes model that is similar to the semiconvex model described above, except that we fuse the learnable weights of the MLP decoder with the weights of the feature mapping, to remove the product of parameters (which is semiconvex but not convex). Our convex model is:

product of parameters (which is semiconvex but not convex). Our convex model is:  $\sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=1}$ 

$$\tilde{y}(q) = \sum_{c \in \{1,2,3,12,13,23,123\}} \mathbb{1}_{d(c)}^{+} (\mathbf{e}_{c} \circ \mathbb{1} [\overline{\mathbf{e}}_{c} \ge 0]),$$
(5)

where the features  $\mathbf{e}_c$  are interpolated from optimizable parameter grids with feature dimension  $d(c) \in \{d_1, d_2, d_3\}$ , whereas the gating variables  $\overline{\mathbf{e}}_c$  inside the indicator function are derived from the same grids frozen at their initialization values to preserve convexity. Here  $\circ$  denotes elementwise product of vectors. These indicator functions take the same role as the ReLU in a nonconvex MLP, using a sampling of random activation patterns based on the grid values at initialization.

# 182 4 THEORY

#### 183 4.1 Equivalence to Matrix Completion in 2D

In three dimensions, the complete set of geometric algebra feature grids are those that we include in the GA-Planes family:  $\{g_1, g_2, g_3, g_{12}, g_{13}, g_{23}, g_{123}\}$ . In two dimensions, this reduces to:  $\{g_1, g_2, g_{12}\}$ . In this 2D setting, we can analyze different members of our GA-Planes family and show equivalence to various formulations of the classic matrix completion problem.

**Notation.** As usual, we use  $\circ$  to denote elementwise multiplication and  $\odot$  to denote concatenation. 188 We use  $\mathbb{1}_{a \times b}$  to denote the all-ones matrix of size  $a \times b$  and  $\mathbb{1}[\cdot]$  to denote the indicator function, 189 which evaluates to 1 when its argument is positive and 0 otherwise. Our theorem statements consider 190 equivalence to a matrix completion problem with target matrix  $M \in \mathbb{R}^{m \times n}$  and low-rank compo-191 nents  $U \in \mathbb{R}^{m \times k}, V \in \mathbb{R}^{n \times k}$  to be optimized. We include theoretical results for 2D GA-Planes 192 models that combine features by addition (+), multiplication  $(\circ)$ , or concatenation  $(\odot)$  and decode 193 features using a linear decoder (as a warmup), a convex MLP, or a nonconvex MLP. The most il-194 luminating results are presented in the theorem statements that follow; the rest (and all proofs) are 195 deferred to Appendix A.2. 196

**Assumptions.** Our theorem statements assume that the line feature grids have the same spatial 197 resolution as the target matrix, and thus do not specify the type of interpolation. However, the results 198 hold even if the dimensions do not match and nearest neighbor interpolation is used; the empirical 199 performance is similar or even slightly improved in practice by using (bi)linear interpolation of 200 features (see Appendix A.3 and A.2 for a discussion of other interpolation methods). The theorems 201 assume that the optimization objective is to minimize the Frobenius norm of the error matrix; this 202 is equivalent to minimizing mean squared error measured directly in the representation space. This 203 objective function is the one we use for our convex experiments (video and volume segmentation 204 fitting), where we have access to direct supervision; our radiance field experiments instead use 205 indirect measurements (along rays) that are not equivalent to the setting of the theorems. 206

**Theorem 1.** The two-dimensional representation  $D(\mathbf{e}_1 + \mathbf{e}_2)$  with linear decoder  $D(f(q)) = \alpha^T f(q)$  is equivalent to a low-rank matrix completion model with the following structure:

$$\min_{U,V} \|M - (U\mathbb{1}_{k \times n} + \mathbb{1}_{m \times k} V^T)\|_F^2.$$
(6)

These two models are equivalent in the sense that  $U^* = \mathbf{g}_1^* \operatorname{diag}(\alpha^*)$  and  $V^* = \mathbf{g}_2^* \operatorname{diag}(\alpha^*)$  where  $U^*, V^*$  is the optimal solution to the low-rank matrix completion problem in eq. (6) and  $\mathbf{g}_1^*, \mathbf{g}_2^*, \alpha^*$ 

are the optimal grid features and linear decoder for the  $D(\mathbf{e}_1 + \mathbf{e}_2)$  model.

The two-dimensional representation  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  with the same linear decoder is equivalent to the standard low-rank matrix completion model:

$$\min_{U,V} \|M - UV^T\|_F^2.$$
(7)

These two models are equivalent in the same sense as above, except that  $V^* = \mathbf{g}_2^*$ .

**Remark.** Using a linear decoder reveals a dramatic difference in representation capacity between feature addition (or concatenation) and multiplication. Using addition, the maximum rank of the matrix approximation is 2 regardless of the feature dimension k. Using multiplication, the maximum rank of the approximation is k. With feature multiplication, the optimal values of the feature grids are identical to the rank-thresholded singular value decomposition (SVD) of M, where the feature grids  $\mathbf{g}_1$  and  $\mathbf{g}_2$  recover the left and right singular vectors and the decoder  $\alpha$  learns the singular values of M. This is the optimal rank k approximation of a matrix M.

**Theorem 2.** The two-dimensional representation  $D(\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_{12})$  with linear decoder  $D(f(q)) = \alpha^T f(q)$  is equivalent to a low-rank plus low-resolution matrix completion model with the following structure:

$$\min_{U,V,L} \|M - (U\mathbb{1}_{k \times n} + \mathbb{1}_{m \times k} V^T + \varphi(L))\|_F^2,$$
(8)

where  $L \in \mathbb{R}^{m_l \times n_l}$  is the low-resolution component to be learned, with upsampling (interpolation) function  $\varphi$ . These two models are equivalent in the sense that  $U^* = \mathbf{g}_1^* diag(\alpha^*)$ ,  $V^* = \mathbf{g}_2^* diag(\alpha^*)$ , and  $L^* = \mathbf{g}_{12}^* \alpha^*$ , where  $U^*, V^*, L^*$  is the optimal solution to the low-rank plus low-resolution matrix completion problem in eq. (8) and  $\mathbf{g}_1^*, \mathbf{g}_2^*, \mathbf{g}_{12}^*, \alpha^*$  are the optimal grid features and linear decoder for the  $D(\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_{12})$  model.

The two-dimensional representation  $D(\mathbf{e}_1 \circ \mathbf{e}_2 + \mathbf{e}_{12})$  with the same linear decoder is equivalent to a low-rank plus low-resolution matrix completion model:

$$\min_{U,V,L} \|M - (UV^T + \varphi(L))\|_F^2.$$
(9)

These two models are equivalent in the same sense as above, except that  $V^* = \mathbf{g}_2^*$ .

**Remark.** Theorem 2 describes the behavior of a 2D, linear-decoder version of our GA-Planes model, both the version with addition/concatenation of features (eq. (2)) and the version with multiplication of features (eq. (3)). Extending the same idea to 3D, we can interpret GA-Planes as a low-rank plus low-resolution approximation of a 3D tensor (volume). We can understand this model as first fitting a low-resolution volume and then finding a low-rank approximation to the high-frequency residual volume. When we use multiplication of features, the low-rank residual approximation is optimal and analogous to the rank-thresholded SVD.

**Theorem 3.** The two-dimensional representation  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  with a two-layer convex MLP decoder  $D(f(q)) = \sum_{i=1}^{h} (W_i^{\top} f(q)) \mathbb{1}[\overline{W_i}^{\top} f(q) \ge 0]$  is equivalent to a masked low-rank matrix completion model:

$$\min_{U,V,W} \left\| M - \sum_{i,j} W_{i,j} U_j V_j^\top \circ B_i \right\|_F^2, \tag{10}$$

where  $W \in \mathbb{R}^{h \times k}$  contains the trainable weights of the convex MLP decoder, with indices  $j = 1, \ldots, k$  for the input dimension and  $i = 1, \ldots, h$  for the hidden layer dimension.  $B_i \in \mathbb{R}^{m \times n}$ denotes the binary masking matrix formed by random, fixed gates of the convex MLP decoder;  $B_i = \mathbb{1}[\sum_j \overline{W}_{i,j} U_j V_j^\top \ge 0]$ , where  $\overline{W}$  denotes the weight matrix W with values fixed at random initialization.

This matrix completion model and our GA-Planes model  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  with convex MLP decoder are equivalent in the sense that  $U^* = \mathbf{g}_1^*$ ,  $V^* = \mathbf{g}_2^*$ , and  $W^* = W^*$ , where  $U^*$ ,  $V^*$ ,  $W^*$  is the optimal solution to the masked low-rank matrix completion problem eq. (10) and  $\mathbf{g}_1^*, \mathbf{g}_2^*, W^*$  are the optimal grid features and convex MLP decoder weights for the  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  model. The optimal mask matrices  $B_i^*$  are defined by the fixed random weight initialization  $\overline{W}$  and the optimal singular vector matrices  $U^*, V^*$ .

**Remark.** We can interpret the matrix completion model of eq. (10) as a sum of h different low-254 rank approximations, where the matrices within each of the h groups are constrained to share the 255 same singular vectors  $U_i, V_i$ . The binary masks  $B_i$  effectively allow each of these h low-rank 256 approximations to attend to (or complete) a different part of the matrix M before being linearly 257 combined through the weights (singular values)  $W_{i,j}$ . The upper limit of the rank of this matrix 258 approximation is thus  $\min(n, m)$ , because the mask matrices can arbitrarily increase the rank beyond 259 the constraint faced by models with a linear decoder. Note that if the feature grids  $g_1$  and  $g_2$  have 260 spatial resolution  $r_1$  less than  $\min(n, m)$ , the maximum rank will be  $r_1$ . 261

**Theorem 4.** The two-dimensional representation  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  with a standard two-layer MLP decoder  $D(f(q)) = \alpha^T (Wf(q))_+$  is equivalent to a low-rank matrix completion model with the following structure:

$$\min_{U,V,W,\alpha} \left\| M - \sum_{i=1}^{h} \alpha_i \Big( \sum_{j=1}^{k} W_{i,j} U_j V_j^{\top} \Big)_+ \right\|_F^2,$$
(11)

where  $W \in \mathbb{R}^{h \times k}$  is the weight matrix for the MLP decoder's hidden layer (with width h) and  $\alpha \in \mathbb{R}^{h}$  is the weight vector of the MLP decoder's output layer.

This matrix completion model and our GA-Planes model  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  with nonconvex MLP decoder are equivalent in the sense that  $U^* = \mathbf{g}_1^*$ ,  $V^* = \mathbf{g}_2^*$ ,  $W^* = W^*$ , and  $\alpha^* = \alpha^*$ , where  $U^*, V^*, W^*, \alpha^*$  is the optimal solution to the masked low-rank matrix completion problem eq. (11) and  $\mathbf{g}_1^*, \mathbf{g}_2^*, W^*, \alpha^*$  are the optimal grid features and MLP decoder for the  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  model.

**Remark.** The upper limit of the rank of this matrix approximation is  $\min(n, m, r_1)$ , the same as with a convex MLP decoder.

We summarize the maximum attainable ranks of different 2D models in Table 1 (see Appendix A.2.3 and A.2.4 for matrix representations of  $D(e_1 + e_2)$  and  $D(e_1 \odot e_2)$  with convex and nonconvex MLP decoders). Experimental validation of these theoretical results on the task of 2D image compression is provided with a comparison of interpolation schemes in Figure 5 in the appendix.

Model	Linear decoder	convex MLP decoder	MLP decoder
$D(\mathbf{e}_1 + \mathbf{e}_2)$	2	$r_1$	$r_1$
$D(\mathbf{e}_1\odot\mathbf{e}_2)$	2	$r_1$	$r_1$
$D(\mathbf{e}_1 \circ \mathbf{e}_2)$	k	$r_1$	$r_1$

**Table 1:** Maximum attainable ranks of different 2D GA-Planes models, using only line features. Here k is the feature dimension and  $r_1$  is the spatial dimension of the features, which need never exceed  $\min(m, n)$ . Replacing a linear decoder with a convex or nonconvex MLP can dramatically increase the rank of the representation.

## 277 4.2 INTERPRETATION: LOW RANK + LOW RESOLUTION

Combining multiple parameterization strategies with complementary representation capacities is a 278 time-honored strategy in signal processing. A classic example is the combination of sparse and 279 low-rank models used to represent matrices in the compressive sensing literature (Chandrasekaran 280 et al., 2009). As shown in Theorem 2, we can view the GA-Planes family as following a similar 281 strategy with a combination of low-rank and low-resolution approximations. This low-rank plus 282 low-resolution parameterization is generally easier to train, because sparse models must either store 283 large numbers of empty values (high memory) or store the locations of nonzero entries and suffer 284 from poorly-conditioned spatial gradients (difficult optimization). Indeed there are volume param-285 eterizations that utilize sparsity, such as point clouds, surface meshes, surfels, and Gaussian splats, 286 287 but these tend to be more challenging to optimize (e.g. requiring high memory (Sara Fridovich-Keil and Alex Yu et al., 2022) or heuristic updates and good initialization (Kerbl et al., 2023)). 288

We illustrate this low-rank plus low-resolution interpretation in Figure 2 with a simple experiment, in which we approximate a grayscale image (the *astronaut* image from SciPy) using either a sum of low-rank and low-resolution components (similar to the structure of GA-Planes) or the classic sum of low-rank and sparse components. In this experiment we compute the optimal low-rank components of each model type using the SVD, which corresponds to multiplication of features.

## 294 5 EXPERIMENTS

## 295 5.1 RADIANCE FIELD MODELING

Our experiments for the radiance field reconstruction task are built on the NeRFStudio framework (Tancik et al., 2023) and use all 8 scenes from NeRF-Blender (Mildenhall et al., 2020). We train each volume representation based on the photometric loss that is standard in the NeRF literature (Max,



for Image Compression

PSNR 29.60

**PSNR 26.26** 

Figure 2: For a natural image, approximation as a sum of low rank and low resolution components (green points and subfigure b) achieves higher fidelity compared to the classic matrix decomposition as a sum of low rank and sparse components (blue points and subfigure c), with the same parameter budget (18.75% of the original image size, for subfigures b and c). The GA-Planes model family generalizes the idea of a low rank plus low resolution approximation to three dimensions.

1995; Kajiya, 1986; Mildenhall et al., 2020). Our results are summarized in Figure 3, which reports 299 PSNR, SSIM (Wang et al., 2004), and LPIPS (Zhang et al., 2018) for each model as a function of its 300

size. We provide per-scene pareto-optimal curves and renderings in Appendix A.6 and A.7. 301



Figure 3: Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on all 8 scenes from the Blender dataset, and the average results are shown.

Because the photometric loss function is inherently nonconvex, we use our most expressive GA-302 Planes parameterization as defined in eq. (3). We compare this with several popular models as 303 implemented in NeRFStudio: Mip-NeRF (Barron et al., 2021), TensoRF (Chen et al., 2022), and K-304 Planes (Fridovich-Keil et al., 2023). For K-Planes, we include versions with and without proposal 305 sampling, a strategy for efficiently allocating ray samples during training and rendering. Proposal 306

sampling is the default for K-Planes, but we include a version without proposal sampling because 307 308 none of the other models in this experiment use proposal sampling. We also include several ablations of our GA-Planes model: one with only the volume features (similar to a multiresolution version 309 of DVGO (Sun et al., 2022) or Plenoxels (Sara Fridovich-Keil and Alex Yu et al., 2022)), one 310 with only the line features (similar to a multiresolution TensoRF-CP), and one with only the line-311 plane products (similar to a multiresolution TensoRF-VM (Chen et al., 2022)). At large model sizes 312  $(\sim 10 \text{ million parameters})$  most models including GA-Planes perform well. As model size shrinks, 313 only GA-Planes and its line-only ablation reach comparable metrics as the larger models. Detailed 314 descriptions of parameter allocations for each model are provided in Appendix A.8. 315

#### 316 5.2 3D SEGMENTATION

Our experiments for 3D segmentation use the opacity masks from the NeRF-Blender *lego* scene (Mildenhall et al., 2020). The task is to "lift" these 2D segmentation masks to 3D. We compare the linear and nonlinear feature combination versions of GA-Planes and Tri-Planes Chan et al. (2022); Fridovich-Keil et al. (2023) (models with plane features only).

**2D Supervision.** In this experiment we use 2D tomographic supervision, in which we minimize the mean squared error between the ground truth 2D object segmentation masks and the average ray density at each viewpoint. In Table 2 we report the intersection over union (IOU) metric on thresholded projections from our trained model on test views. These results validate that the GA-Planes architecture in eq. (2) works well regardless of convex, semiconvex, or nonconvex formulation– whereas the Tri-Planes model requires a nonconvex decoder for good performance.

	Convex	Semiconvex	Nonconvex
GA-Planes (with $\circ$ )	-	-	0.877
GA-Planes (with $\odot$ )	0.875	0.883	0.880
Tri-Planes (planes with $\circ$ , like K-Planes)	-	-	0.877
Tri-Planes (planes with +, like (Chan et al., 2022))	0.681	0.868	0.863

**Table 2:** Intersection over union (IOU) for recovering novel view object segmentation masks from segmentation mask training with 2D tomographic supervision.

**327 3D Supervision.** Our second set of 3D segmentation experiments leverages direct 3D supervision via Space Carving (Kutulakos & Seitz, 1999), using the principle that if any ray passing through a given 3D coordinate is transparent, the density at that 3D coordinate must be zero. Our results with 3D supervision, in Table 3, parallel those with 2D supervision: GA-Planes performs well regardless of convex, semiconvex, or nonconvex formulation, whereas the Tri-Planes model performs decently with a nonconvex decoder but much worse with convex or semiconvex formulation.

	Convex	Semiconvex	Nonconvex
GA-Planes (with $\circ$ )	-	-	0.926
GA-Planes (with $\odot$ )	0.932	0.957	0.964
Tri-Planes (planes with ○, like K-Planes)	-	-	0.881
Tri-Planes (planes with +, like (Chan et al., 2022))	0.642	0.636	0.941

**Table 3:** Intersection over union (IOU) for recovering novel view object segmentation masks from segmentation mask training with 3D Space Carving supervision.

#### 333 5.3 VIDEO SEGMENTATION

Our video segmentation task is similar to volume segmentation with 3D supervision: here the volume dimensions are x, y, t rather than x, y, z, and the supervision is performed directly in 3D using segmentation masks for a subset of the video frames (every third frame is held out for testing). Our dataset preparation pipeline uses the skateboarding video and preprocessing steps described at Labelbox.com, which involves first extracting a bounding box with YOLOv8 (Jocher et al., 2023) and then segmenting the skateboarder with SAM (Kirillov et al., 2023). This is essentially a temporal superresolution task on segmentation masks. Our results are summarized in Figure 4 and Table 4. We find that GA-Planes performs well across convex, semiconvex, and nonconvex formulations, though its performance is slightly reduced under fully convex training, perhaps because the convex model size is slightly reduced due to fusing the decoder parameters into the feature grids. In contrast, the simpler Tri-Plane models perform poorly on this task regardless of training strategy: they fail to learn the temporal sequence of the video,

producing masks that focus only on the skateboarder's less-mobile core.



**Figure 4:** Intersection over union (IOU) for predicting segmentation masks for unseen frames within a video of a segmented skateboarder.

	Convex	Semiconvex	Nonconvex
GA-Planes (with $\circ$ )	-	-	0.974
GA-Planes (with $\odot$ )	0.913	0.975	0.981
Tri-Planes (planes with $\circ$ , like K-Planes)	-	-	0.727
Tri-Planes (planes with +, like Chan et al. (2022))	0.557	0.647	0.732

**Table 4:** Intersection over union (IOU) for temporal superresolution of segmentation masks in a video, computed on held-out test frames. Models that involve multiplication of features can only be trained by nonconvex optimization.

# 347 6 DISCUSSION

In this work we introduce GA-Planes, a family of volume parameterizations that generalizes many 348 existing representations (see Appendix A.1). We specifically focus on two members of the GA-349 Planes family (with concatenation versus multiplication of features), and offer both theoretical in-350 terpretation and empirical evaluation of these models. Our nonconvex GA-Planes model shows 351 pareto-optimal performance in terms of model size and quality on fitting a 3D radiance field, while 352 our convex and semiconvex GA-Planes formulations are effective on several 3D linear inverse prob-353 lems. In 2D, we show connections between GA-Planes and a low-rank plus low-resolution matrix 354 completion model, and derive how this model's rank capacity is affected by various design decisions. 355

Limitations. Here we focus on 3D (or smaller) representations, rather than higher dimensions (e.g. dynamic volumes), and we demonstrate GA-Planes for reconstruction rather than generation tasks. Both of these extensions are promising avenues for extending GA-Planes. In our experiments, we use the same first-order optimization algorithm for all models. However, our convex GA-Planes formulation is compatible with any convex solver (e.g. cvxpy), and we expect its performance may improve by leveraging these efficient convex optimization algorithms. We demonstrate preliminary benefits of convexity and semiconvexity in terms of training stability in Appendix A.5.

# 363 REFERENCES

Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and

- Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields.
   In *ICCV*, pp. 5835–5844. IEEE, 2021. doi: 10.1109/ICCV48922.2021.00580. URL https:
- 367 //doi.org/10.1109/ICCV48922.2021.00580.
- Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Chen Yang, Wei Shen, Lingxi Xie, Dongsheng Jiang, Xiaopeng Zhang, and Qi Tian. Segment anything in 3d with nerfs. In *NeurIPS*, 2023.
- Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello,
- Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, 2022.
- Venkat Chandrasekaran, Sujay Sanghavi, Pablo A Parrilo, and Alan S Willsky. Sparse and low-rank
   matrix decompositions. *IFAC Proceedings Volumes*, 42(10):1493–1498, 2009.
- Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022.
- V Croce, G Caroti, L De Luca, A Piemonte, and P Véron. Neural radiance fields (nerf): Review and
   potential applications to digital cultural heritage. *The International Archives of the Photogram- metry, Remote Sensing and Spatial Information Sciences*, 48:453–460, 2023.
- Leo Dorst, Daniel Fontijne, and Stephen Mann. *Geometric algebra for computer science (revised edition): An object-oriented approach to geometry.* Morgan Kaufmann, 2009.
- Tolga Ergen and Mert Pilanci. Path regularization: A convexity and sparsity inducing regularization for parallel relu networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo
   Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *CVPR*, 2023.
- Aditya M Intwala and Atul Magikar. A review on process of 3d model reconstruction. In 2016
   *International Conference on Electrical, Electronics, and Optimization Techniques (ICEEOT)*, pp. 2851–2855. IEEE, 2016.
- Glenn Jocher, Jing Qiu, and Ayush Chaurasia. Ultralytics YOLO, January 2023. URL https: //github.com/ultralytics/ultralytics.
- James T Kajiya. The rendering equation. In *Proceedings of the 13th annual conference on Computer graphics and interactive techniques*, pp. 143–150, 1986.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics, 42(4), July 2023.
   URL https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
   Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.
   Segment anything, 2023. URL https://arxiv.org/abs/2304.02643.
- K.N. Kutulakos and S.M. Seitz. A theory of shape by space carving. In *Proceedings of the Seventh IEEE International Conference on Computer Vision*, volume 1, pp. 307–314 vol.1, 1999. doi:
   10.1109/ICCV.1999.791235.
- Labelbox.com. Using meta's segment anything (sam) model on video with labelbox's modelassisted labeling. https://labelbox.com/guides/using-metas-segment-
- 404 anything-sam-model-on-video-with-labelbox-model-assisted-
- 405 labeling/. Accessed: 2024-10-01.
- Jian Liu, Xiaoshui Huang, Tianyu Huang, Lu Chen, Yuenan Hou, Shixiang Tang, Ziwei Liu, Wanli
   Ouyang, Wangmeng Zuo, Junjun Jiang, et al. A comprehensive survey on 3d content generation.
   *arXiv preprint arXiv:2402.01166*, 2024.

Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Trans. Graph.*, 40 (4), jul 2021. ISSN 0730-0301. doi: 10.1145/3450626.3459863. URL https://doi.org/ 10.1145/3450626.3459863.

Valentín Masero, JUAN M LEÓN-ROJAS, and José Moreno. Volume reconstruction for health
care: a survey of computational methods. *Annals of the New York Academy of Sciences*, 980(1):
198–211, 2002.

<sup>416</sup> Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and* <sup>417</sup> *Computer Graphics*, 1(2):99–108, 1995.

Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4460–4470, 2019.

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. In *ECCV*, 2020.

Yuhang Ming, Xingrui Yang, Weihan Wang, Zheng Chen, Jinglun Feng, Yifan Xing, and Guofeng
 Zhang. Benchmarking neural radiance fields for autonomous robots: An overview. *arXiv preprint arXiv:2405.05526*, 2024.

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

- doi: 10.1145/3528223.3530127. URL https://doi.org/10.1145/3528223.3530127.
- Alan V Oppenheim. *Discrete-time signal processing*. Pearson Education India, 1999.

Mert Pilanci and Tolga Ergen. Neural networks are convex regularizers: Exact polynomial-time
convex optimization formulations for two-layer networks. In Hal Daumé III and Aarti Singh
(eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 7695–7705. PMLR, 13–18 Jul 2020. URL

434 https://proceedings.mlr.press/v119/pilanci20a.html.

Christian Reiser, Richard Szeliski, Dor Verbin, Pratul P. Srinivasan, Ben Mildenhall, Andreas
Geiger, Jonathan T. Barron, and Peter Hedman. Merf: Memory-efficient radiance fields for realtime view synthesis in unbounded scenes. *SIGGRAPH*, 2023.

Alexander Richter, Till Steinmann, Jean-Claude Rosenthal, and Stefan J Rupitsch. Advances in real-time 3d reconstruction for medical endoscopy. *Journal of Imaging*, 10(5):120, 2024.

Arda Sahiner, Tolga Ergen, Batu Ozturkler, John M Pauly, Morteza Mardani, and Mert Pilanci.
 Scaling convex neural networks with burer-monteiro factorization. In *ICLR*, 2024.

Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo
 Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.

Vishwanath Saragadam, Daniel LeJeune, Jasper Tan, Guha Balakrishnan, Ashok Veeraraghavan,
 and Richard G Baraniuk. Wire: Wavelet implicit neural representations. In *Conf. Computer Vision and Pattern Recognition*, 2023.

Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Confer- ence on Computer Vision and Pattern Recognition (CVPR)*, 2016.

Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise
 view selection for unstructured multi-view stereo. In *European Conference on Computer Vision* (*ECCV*), 2016.

Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael
 Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2019a.

Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Con tinuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Pro- cessing Systems*, 2019b.

- Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Proc. NeurIPS*, 2020.
- Eugen Šlapak, Enric Pardo, Matúš Dopiriak, Taras Maksymyuk, and Juraj Gazda. Neural radiance
   fields in the industrial and robotics domain: applications, research opportunities and use cases.
   *Robotics and Computer-Integrated Manufacturing*, 90:102810, 2024.
- Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022.

Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan,
Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020.

Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang,
 Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, David McAllister, and
 Angjoo Kanazawa. Nerfstudio: A modular framework for neural radiance field development.
 In ACM SIGGRAPH 2023 Conference Proceedings, SIGGRAPH '23, 2023.

- Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Wang Yifan,
  Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. Advances in neural rendering. In *Computer Graphics Forum*, volume 41, pp. 703–735. Wiley Online
  Library, 2022.
- <sup>476</sup> Jayaram K Udupa and Gabor T Herman. *3D imaging in medicine*. CRC press, 1999.

Mikaela Angelina Uy, Ricardo Martin-Brualla, Leonidas Guibas, and Ke Li. Scade: Nerfs from
 space carving with ambiguity-aware depth estimates. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.

Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:
from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.

Liyana Wijayathunga, Alexander Rassau, and Douglas Chai. Challenges and solutions for au tonomous ground robot scene understanding and navigation in unstructured outdoor environ ments: A review. *Applied Sciences*, 13(17):9877, 2023.

Mengya Xu, Ziqi Guo, An Wang, Long Bai, and Hongliang Ren. A review of 3d reconstruction
 techniques for deformable tissues in robotic surgery. *arXiv preprint arXiv:2408.04426*, 2024.

Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
 effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.

# 490 A APPENDIX

#### 491 A.1 CONTEXT FOR GA-PLANES

Table 5 summarizes some representative volume models popular in computer vision, and how they relate to GA-Planes along the threeway pareto frontier of model size, expressiveness, and optimizability.

Implicit Neural Representations (INRs) or Co-497 ordinate Neural Networks (Mildenhall et al., 498 2020; Sitzmann et al., 2019b; Tancik et al., 499 2020; Sitzmann et al., 2020; Saragadam et al., 500 2023) parameterize the volume implicitly 501 through the weights of a neural network, typi-502 cally a multilayer perceptron (MLP) with some 503 modification to overcome spectral bias and rep-504 resent high frequency content. These models 505 tend to provide decent expressiveness with very 506 507 small model size; their main drawback is slow optimization. 508

Voxel grids (Kutulakos & Seitz, 1999; Sara 509 Fridovich-Keil and Alex Yu et al., 2022; Sun 510 et al., 2022; Sitzmann et al., 2019a) are perhaps 511 the most traditional parameterization of a vol-512 ume, where each parameter denotes the func-513 tion value (density, color, a latent feature, etc.) 514 at a specific grid cell location within the vol-515 ume. These voxel values can then be combined 516 into a continuous function over the 3D space 517 by some form of interpolation, following the 518 standard Nyquist sampling and reconstruction 519 paradigm of digital signal processing (Oppen-520 heim, 1999). Voxels offer direct control over 521

522 expressivity (via resolution) and are easily op-

	Model Size	Expressiveness	Optimizability
Coordinate MLP (NeRF, SRN)	<ul> <li>✓</li> </ul>	2	X
Voxels (Space Carving, Plenoxels, DVGO)	X	1	1
Tensor Factorization (TensoRF, K-Planes)	~	~	~
Hash Embedding (Instant-NGP)	~	1	~
Point Cloud / Splat (3D Gaussian Splatting)	~	1	~
Mixture of Primitives (MVP, MERF)	~	1	~
GA-Planes (Nonconvex)	~	1	~
GA-Planes (Convex)	~	1	1
GA-Planes (Semi-Convex)	~	1	1

**Table 5: Context.** All volume models face a tradeoff between memory efficiency, expressiveness, and optimizability. The qualitative categorizations here are based on the tradeoffs achieved by representative example methods listed in each category. *Model Size* denotes memory usage during training; other methods exist to compress trained models, e.g. for rendering on mobile hardware. *Optimizability* denotes both speed and stability of optimization/training. For example, Coordinate MLPs tend to train slowly, while Splats train quickly but are sensitive to initialization.

timized; their main drawback is memory usage because the number of parameters grows cubically with the spatial resolution.

Tensor factorizations (Chen et al., 2022; Chan et al., 2022; Fridovich-Keil et al., 2023) parameterize a 3D volume as a combination of lower-dimension objects, namely vectors and matrices (lines and planes). Tensor factorizations tend to balance the three attributes somewhat evenly, offering decent expressiveness and optimizability while using more memory than an INR but less than a high resolution voxel grid.

Hash embeddings (Müller et al., 2022; Tancik et al., 2023) are similar to voxels, but replace the
 explicit voxel grid in 3D with a multiresolution 3D hash function followed by a small MLP decoder
 to disambiguate hash collisions. They can optimize very quickly and with better memory efficiency
 compared to voxels; quality is mixed with good high-resolution details but also some high-frequency
 noise likely arising from unresolved hash collisions or sensitivity to random initialization.

Point clouds / splats (Kerbl et al., 2023; Schönberger & Frahm, 2016; Schönberger et al., 2016) represent a volume as a collection of 3D points or blobs, where the points need not be arranged on a regular grid. They are highly expressive and less memory-intensive than voxels (but still more so than some other methods). They can optimize very quickly but often require heuristic or discrete optimization strategies that result in sensitivity to initialization.

540 *Mixture of primitives* (Reiser et al., 2023; Lombardi et al., 2021) models combine multiple of the 541 above representation strategies to balance their strengths and weaknesses. For example, combining 542 low resolution voxels with a high resolution tensor factorization is an effective strategy to improve on the expressiveness of tensor factorizations without resorting to the cubic memory requirement of a high resolution voxel grid; this strategy underlies both MERF (Reiser et al., 2023) and GA-Planes.

We emphasize that all of these existing methods (except perhaps voxels) require nonconvex optimization, often for a feature decoder MLP, and thus risk getting stuck in suboptimal local minima depending on the randomness of initialization and the trajectory of stochastic gradients. In practice, as described above, some of the prior methods exhibit greater optimization stability than others, though none (except voxels in limited settings) come with guarantees of convergence to global optimality. In contrast, both the convex and semiconvex GA-Planes formulations come with guarantees that all local optima are also global (Sahiner et al., 2024).

**Relation to Prior Models.** Without any convexity restrictions, the GA-Planes family includes many previously proposed models as special cases:

- NeRF (Mildenhall et al., 2020): D
- Plenoxels (Sara Fridovich-Keil and Alex Yu et al., 2022), DVGO (Sun et al., 2022):  $D(e_{123})$
- TensoRF (Chen et al., 2022):  $D((\mathbf{e}_1 \circ \mathbf{e}_{23}) \odot (\mathbf{e}_2 \circ \mathbf{e}_{13}) \odot (\mathbf{e}_3 \circ \mathbf{e}_{12}))$
- Tri-Planes (Chan et al., 2022):  $D(e_{12} + e_{13} + e_{23})$
- K-Planes (Fridovich-Keil et al., 2023):  $D(\mathbf{e}_{12} \circ \mathbf{e}_{13} \circ \mathbf{e}_{23})$
- MERF (Reiser et al., 2023):  $D(e_{12} + e_{13} + e_{23} + e_{123})$

Of these, all except for TensoRF and K-Planes are compatible with convex optimization towards any convex objective. Note that different models may use different decoder architectures for D, including both linear and MLP decoders and additional decoder inputs such as encoded viewing direction and/or positionally-encoded coordinates.

Convex Neural Networks. Given a data matrix  $X \in \mathbb{R}^{n \times d}$  and labels  $y \in \mathbb{R}^n$ , a 2-layer nonconvex ReLU MLP approximates y as

$$y \approx \sum_{j=1}^{m} (XU_j)_+ \alpha_j, \tag{12}$$

where m is the number of hidden neurons and U and  $\alpha$  are the first and second linear layer weights, respectively. Pilanci & Ergen (2020) proposed to instead approximate y as

$$y \approx \sum_{i=1}^{P} D_i X(v_i - w_i), \tag{13}$$

subject to  $(2D_i - I_n)Xv_i \ge 0$  and  $(2D_i - I_n)Xw_i \ge 0$  for all *i*. The parameters v and w in 569 eq. (13) replace the first and second layer weights U and  $\alpha$  from the nonconvex formulation in 570 eq. (12) (optimal values of U and  $\alpha$  can be recovered from optimal values of v and w). Here  $D_i$ 571 represent different possible activation patterns of the hidden neurons as  $\{D_i\}_{i=1}^P := \{\text{Diag}(\mathbb{1}[Xu \geq i])\}$ 572  $0]): u \in \mathbb{R}^d$ , which is the finite set of hyperplane arrangement patterns obtained for all possible 573  $u \in \mathbb{R}^d$ . We can sample different u's to find all distinct activation patterns  $\{D_i\}_{i=1}^P$ , where P is 574 the number of regions in the partitioned input space. Enumerating all such patterns would yield an 575 exact equivalence with the global minimizer of the nonconvex ReLU MLP in eq. (12), but may be 576 complicated or intractable due to memory limitations. Subsampling  $\hat{P}$  patterns results in a convex 577 program with tractable size, whose solution is one of the stationary points of the original non-convex 578 problem (Pilanci & Ergen, 2020). We apply this idea to create convex and semiconvex GA-Planes 579 models by convexifying the feature decoder MLP according to this procedure. 580

## 581 A.2 PROOF OF THEOREMS

A general note on proofs. In order to represent a matrix  $M \in \mathbb{R}^{m \times n}$  with an implicit model, we compute D(f(q)) for q = (k, l),  $\forall k \in \{1, ..., m\}$ ,  $\forall l \in \{1, ..., n\}$ . Considering line feature grids with resolutions matching m, n; the features will become  $\mathbf{e}_1 = (\mathbf{g}_1)_k$ ,  $\mathbf{e}_2 = (\mathbf{g}_2)_l$  for q = (k, l)

otherwise they will be  $\mathbf{e}_1 = \varphi(\mathbf{g}_1)_k$ ,  $\mathbf{e}_2 = \varphi(\mathbf{g}_2)_l$  where  $\varphi(\mathbf{g}_1)$ ,  $\varphi(\mathbf{g}_2)$  now have resolutions m, 585 n after interpolation through  $\varphi$ . Here  $\varphi$  can be any interpolation scheme with linear weighting of 586 inputs, e.g. nearest neighbor, (bi)linear, (bi)cubic, spline, Gaussian, sinc, etc. For the simplicity 587 of notation, we omit  $\varphi$  in line feature grids, and only apply it to the plane feature grid  $\mathbf{g}_{12}$ , which 588 has lower resolution by design. The proofs consider  $\mathbf{g}_1, \mathbf{g}_2 \in \mathbb{R}^{r_1 \times d_1}$  and  $\mathbf{g}_{12} \in \mathbb{R}^{r_2 \times r_2 \times d_1}$  (equal 589 feature dimensions, different resolutions), resulting in the matrix representation  $\hat{M} \in \mathbb{R}^{r_1 \times r_1}$  (the 590 case where  $m = n = r_1$ ). Note that if the interpolation is done by a method other than nearest 591 neighbor, this may allow a (convex or nonconvex) MLP decoder to increase the rank beyond  $r_1$ . We 592 derive expressions for M implied by different GA-Planes variations in the parts that follow. The 593 coordinate-wise optimization objective (in the case of a directly supervised mean-square-error loss) 594 corresponds to minimizing the Frobenius norm of the ground truth matrix M and its approximation 595 Â. 596

# 597 A.2.1 PROOF OF THEOREM 1

The forward mapping of the model  $D(e_1 + e_2)$  is:

$$\tilde{y}(q) = \mathbf{D}(\mathbf{e}_1 + \mathbf{e}_2) = \alpha^{\top}(\mathbf{e}_1 + \mathbf{e}_2) = \alpha^{\top}((\mathbf{g}_1)_k + (\mathbf{g}_2)_l) = \sum_{i=1}^{d_1} \alpha_i((\mathbf{g}_1)_{k,i} + (\mathbf{g}_2)_{l,i}), \quad (14)$$

- 599 where  $(\mathbf{g}_1)_k, (\mathbf{g}_1)_l \in \mathbb{R}^{d_1 \times 1}$ .
- 600 In matrix form,

$$\hat{M} = \sum_{i=1}^{d_1} \alpha_i ((\mathbf{g}_1)_i \mathbb{1}^\top + \mathbb{1}(\mathbf{g}_2)_i^\top) = \sum_{i=1}^{d_1} \alpha_i (\mathbf{g}_1)_i \mathbb{1}^\top + \sum_{i=1}^{d_1} \alpha_i \mathbb{1}(\mathbf{g}_2)_i^\top.$$
(15)

Defining  $U := \mathbf{g}_1 \operatorname{diag}(\alpha), U \in \mathbb{R}^{r_1 \times d_1}$  and  $V := \mathbf{g}_2 \operatorname{diag}(\alpha), V \in \mathbb{R}^{r_1 \times d_1}$ , this can be expressed as

$$\hat{M} = U \mathbb{1}_{r_1 \times d_1}^{\top} + \mathbb{1}_{r_1 \times d_1} V^{\top}.$$
(16)

Note that the resulting matrix  $\hat{M} \in \mathbb{R}^{r_1 \times r_1}$  has rank at most 2—very limited expressivity regardless of the resolution  $r_1$ . This is because the all-ones matrix is rank 1, and a product of matrices cannot have higher rank than either of its factors.

Similarly for the multiplicative representation  $D(e_1 \circ e_2)$ , the mapping is

$$\tilde{y}(q) = \mathsf{D}(\mathbf{e}_1 \circ \mathbf{e}_2) = \alpha^{\top}(\mathbf{e}_1 \circ \mathbf{e}_2) = \alpha^{\top}((\mathbf{g}_1)_k \circ (\mathbf{g}_2)_l) = \sum_{i=1}^{d_1} \alpha_i(\mathbf{g}_1)_{k,i}(\mathbf{g}_2)_{l,i}, \quad (17)$$

where  $(\mathbf{g}_1)_k, (\mathbf{g}_1)_l \in \mathbb{R}^{d_1 \times 1}$ . In matrix form,

$$\hat{M} = \sum_{i=1}^{d_1} \alpha_i(\mathbf{g}_1)_i(\mathbf{g}_2)_i^{\top} = \mathbf{g}_1 \operatorname{diag}(\alpha) \mathbf{g}_2^{\top}.$$
(18)

Defining  $U := \mathbf{g}_1 \operatorname{diag}(\alpha), U \in \mathbb{R}^{r_1 \times d_1}$  and  $V := \mathbf{g}_2, V \in \mathbb{R}^{r_1 \times d_1}$ , this can be expressed as

$$\hat{M} = UV^{\top},\tag{19}$$

- which is the optimal rank- $d_1$  decomposition.
- 610 A.2.2 PROOF OF THEOREM 2

ĵ

The forward mapping of the model  $D(e_1 + e_2 + e_{12})$  becomes:

$$\tilde{p}(q) = \mathbf{D}(\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_{12}) = \alpha^{\top}(\mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_{12}) = \alpha^{\top}((\mathbf{g}_1)_k + (\mathbf{g}_2)_l + \varphi(\mathbf{g}_{12})_{k,l})$$
(20)

$$=\sum_{i=1}^{a_1} \alpha_i((\mathbf{g}_1)_{k,i} + (\mathbf{g}_2)_{l,i}) + \sum_{i=1}^{a_1} \alpha_i \varphi(\mathbf{g}_{12})_{k,l,i},$$
(21)

612 where  $(\mathbf{g}_1)_k, (\mathbf{g}_1)_l, \varphi(\mathbf{g}_{12})_{k,l} \in \mathbb{R}^{d_1 \times 1}$ . In matrix form,

$$\hat{M} = \sum_{i=1}^{d_1} \alpha_i ((\mathbf{g}_1)_i \mathbb{1}^\top + \mathbb{1}(\mathbf{g}_2)_i^\top) + \sum_{i=1}^{d_1} \alpha_i \varphi(\mathbf{g}_{12})_i.$$
(22)

Noting that the first term is the same as in eq. (15) and defining  $L := \mathbf{g}_{12}\alpha, L \in \mathbb{R}^{r_2 \times r_2}$ , we reach the expression

$$\hat{M} = U \mathbb{1}_{d_1 \times r_1} + \mathbb{1}_{r_1 \times d_1} V^\top + \varphi(L),$$
(23)

since  $\sum_{i=1}^{d_1} \alpha_i \varphi(\mathbf{g}_{12})_i = \varphi(\sum_{i=1}^{d_1} \alpha_i(\mathbf{g}_{12})_i) = \varphi(\mathbf{g}_{12}\alpha)$ , following the linearity of the upsampling function  $\varphi$ . Note that in the definition of L there is a tensor-vector product that effectively takes a dot product along the last (feature) dimension.

618 Similarly for the multiplicative representation  $D(e_1 \circ e_2 + e_{12})$ , the mapping is

$$\tilde{y}(q) = \mathbf{D}(\mathbf{e}_1 \circ \mathbf{e}_2 + \mathbf{e}_{12}) = \alpha^\top (\mathbf{e}_1 \circ \mathbf{e}_2 + \mathbf{e}_{12})$$
(24)

$$= \alpha^{\top} ((\mathbf{g}_{1})_{k} \circ (\mathbf{g}_{2})_{l} + \varphi(\mathbf{g}_{12})_{k,l}) = \sum_{i=1}^{a_{1}} \alpha_{i}(\mathbf{g}_{1})_{k,i}(\mathbf{g}_{2})_{l,i} + \sum_{i=1}^{a_{1}} \alpha_{i}\varphi(\mathbf{g}_{12})_{k,l,i}.$$
 (25)

619 In matrix notation, we have

$$\hat{M} = \sum_{i=1}^{d_1} \alpha_i ((\mathbf{g}_1)_i (\mathbf{g}_2)_i^{\top}) + \sum_{i=1}^{d_1} \alpha_i \varphi(\mathbf{g}_{12})_i.$$
(26)

Following eq. (18) and using the same definition of L, the final expression becomes

$$\hat{M} = UV^{\top} + \varphi(L). \tag{27}$$

621 A.2.3 PROOF OF THEOREM 3

 $\tilde{y}$ 

For a 2-layer convex MLP with hidden size h, denote the trainable first layer weights as  $W \in \mathbb{R}^{h \times d_1}$ and the gating weights as  $\overline{W} \in \mathbb{R}^{h \times d_1}$  (which are fixed at random initialization). We will handle three different cases for merging the interpolated features: multiplication ( $\circ$ ), addition (+), and concatenation ( $\odot$ ).

The forward mapping of the multiplicative model using a convex MLP,  $D(e_1 \circ e_2)$  at q = (k, l) is

$$\tilde{y}(q) = \mathbb{1}_{h}^{\top} \left( \left( W((\mathbf{g}_{1})_{k} \circ (\mathbf{g}_{2})_{l}) \right) \circ \mathbb{1} \left[ \overline{W}((\mathbf{g}_{1})_{k} \circ (\mathbf{g}_{2})_{l}) \ge 0 \right] \right)$$

$$(28)$$

$$=\sum_{i=1}^{h} \left( \sum_{j=1}^{d_1} W_{i,j}(\mathbf{g}_1)_{k,j}(\mathbf{g}_2)_{l,j} \right) \mathbb{1} \left[ \sum_{j=1}^{d_1} \overline{W}_{i,j}(\mathbf{g}_1)_{k,j}(\mathbf{g}_2)_{l,j} \ge 0 \right],$$
(29)

where  $\circ$  denotes elementwise multiplication (Hadamard product). The resulting matrix decomposition can then be written as

$$\hat{M} = \sum_{i=1}^{h} \left( \sum_{j=1}^{d_1} W_{i,j}(\mathbf{g}_1)_j(\mathbf{g}_2)_j^\top \right) \circ \mathbb{1} \left[ \sum_{j=1}^{d_1} \overline{W}_{i,j}(\mathbf{g}_1)_j(\mathbf{g}_2)_j^\top \ge 0 \right].$$
(30)

Now, we define the masking matrix  $B_i = \mathbb{1}\left[\sum_{j=1}^{d_1} \overline{W}_{i,j}(\mathbf{g}_1)_j(\mathbf{g}_2)_j^\top \ge 0\right]$  and the eigenvectors  $U_j = (\mathbf{g}_1)_j, V_j = (\mathbf{g}_2)_j$  to reach the expression from the theorem statement:

$$\hat{M} = \sum_{i,j} W_{i,j} U_j V_j^\top \circ B_i.$$
(31)

When the model uses additive features as in  $D(e_1 + e_2)$ , and D is a convex MLP, the prediction is

$$(q) = \mathbb{1}_{h}^{\top} \left( \left( W((\mathbf{g}_{1})_{k} + (\mathbf{g}_{2})_{l} \right) \right) \circ \mathbb{1} \left[ \overline{W}((\mathbf{g}_{1})_{k} + (\mathbf{g}_{2})_{l} \right) \geq 0 \right] \right)$$
(32)

$$=\sum_{i=1}^{h} \left( \sum_{j=1}^{d_1} W_{i,j}((\mathbf{g}_1)_{k,j} + (\mathbf{g}_2)_{l,j}) \right) \mathbb{1} \left[ \sum_{j=1}^{d_1} \overline{W}_{i,j}((\mathbf{g}_1)_{k,j} + (\mathbf{g}_2)_{l,j}) \ge 0 \right].$$
(33)

<sup>632</sup> The resulting matrix decomposition can then be written as

$$\hat{M} = \sum_{i=1}^{h} \left( \sum_{j=1}^{d_1} W_{i,j}((\mathbf{g}_1)_j \mathbb{1}_{r_1}^\top + \mathbb{1}_{r_1}(\mathbf{g}_2)_j^\top) \right) \mathbb{1} \left[ \sum_{j=1}^{d_1} \overline{W}_{i,j}((\mathbf{g}_1)_j \mathbb{1}_{r_1}^\top + \mathbb{1}_{r_1}(\mathbf{g}_2)_j^\top) \ge 0 \right].$$
(34)

633 Defining  $B_i = \mathbb{1}\left[\sum_{j=1}^{d_1} \overline{W}_{i,j}((\mathbf{g}_1)_j \mathbb{1}_{r_1}^\top + \mathbb{1}_{r_1}(\mathbf{g}_2)_j^\top) \ge 0\right], U_j = (\mathbf{g}_1)_j, V_j = (\mathbf{g}_2)_j$ , we reach the 634 final expression:

$$\hat{M} = \sum_{i,j} W_{i,j} \left( U_j \mathbb{1}_{r_1}^\top + \mathbb{1}_{r_1} V_j^\top \right) \circ B_i.$$
(35)

<sup>635</sup> Finally, we show that concatenation of features results in a very similar expression to eq. (35).

<sup>636</sup> When the model uses concatenated features as in  $D(\mathbf{e}_1 \odot \mathbf{e}_2)$ , and D is a convex MLP (with trainable <sup>637</sup> weights  $W \in \mathbb{R}^{h \times 2d_1}$  and fixed gates  $\overline{W} \in \mathbb{R}^{h \times 2d_1}$ ), the prediction at a point q is

$$\tilde{y}(q) = \mathbb{1}_{h}^{\top} \left( \left( W((\mathbf{g}_{1})_{k} \odot (\mathbf{g}_{2})_{l} \right) \right) \circ \mathbb{1} \left[ \overline{W}((\mathbf{g}_{1})_{k} \odot (\mathbf{g}_{2})_{l} \right) \geq 0 \right] \right).$$
(36)

<sup>638</sup> Denoting the weights and gates each as a concatenation of 2 matrices,  $W = (W_1 \odot W_2)$ ,  $\overline{W} = (\overline{W}_1 \odot \overline{W}_2)$ , where  $W_1, W_2, \overline{W}_1, \overline{W}_2 \in \mathbb{R}^{h \times d_1}$ , we have the following expression:

$$\tilde{y}(q) = \sum_{i=1}^{h} \left( \sum_{j=1}^{d_1} W_{1i,j}(\mathbf{g}_1)_{k,j} + W_{2i,j}(\mathbf{g}_2)_{l,j} \right) \mathbb{1} \left[ \sum_{j=1}^{d_1} \overline{W}_{1i,j}(\mathbf{g}_1)_{k,j} + \overline{W}_{2i,j}(\mathbf{g}_2)_{l,j} \ge 0 \right].$$
(37)

<sup>640</sup> Following similar steps as for the additive case, we express the matrix decomposition as

$$\hat{M} = \sum_{i,j} \left( W_{1i,j} U_j \mathbb{1}_{r_1}^\top + W_{2i,j} \mathbb{1}_{r_1} V_j^\top \right) \circ B_i,$$
(38)

641 where  $B_i = \mathbb{1}\left[\sum_{j=1}^{d_1} \overline{W}_{1i,j}(\mathbf{g}_1)_j \mathbb{1}_{r_1}^\top + \overline{W}_{2i,j} \mathbb{1}_{r_1}(\mathbf{g}_2)_j^\top \ge 0\right], U_j = (\mathbf{g}_1)_j, V_j = (\mathbf{g}_2)_j.$ 

In all these representations, a low-rank matrix is multiplied elementwise with a binary mask  $B_i$ , which makes the maximum attainable rank  $r_1$ . Thus, with a convex MLP decoder, rank of  $\hat{M}$  is limited by the resolution of the feature grids.

# 645 A.2.4 PROOF OF THEOREM 4

For a standard 2-layer ReLU MLP with hidden size h, denote the trainable first and second layer weights as  $W \in \mathbb{R}^{h \times d_1}$ ,  $\alpha \in \mathbb{R}^{h \times 1}$ . We will handle three different cases for merging the interpolated features: multiplication ( $\circ$ ), addition (+), and concatenation ( $\odot$ ).

The forward mapping of the multiplicative model using a standard nonconvex MLP,  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  is:

$$\tilde{y}(q) = \alpha^{\top} \left[ W((\mathbf{g}_1)_k \circ (\mathbf{g}_2)_l) \right]_+$$
(39)

$$=\sum_{i=1}^{h} \alpha_{i} \left[ \sum_{j=1}^{d_{1}} W_{i,j}(\mathbf{g}_{1})_{k,j}(\mathbf{g}_{2})_{l,j} \right]_{+}.$$
 (40)

<sup>650</sup> The resulting matrix decomposition can then be written as

$$\hat{M} = \sum_{i=1}^{h} \alpha_i \left( \sum_{j=1}^{d_1} W_{i,j}(\mathbf{g}_1)_j(\mathbf{g}_2)_j^\top \right)_+ = \sum_{i=1}^{h} \alpha_i \left( \sum_{j=1}^{d_1} W_{i,j}U_jV_j^\top \right)_+,$$
(41)

651 with  $U_j = (\mathbf{g}_1)_j, V_j = (\mathbf{g}_2)_j$ .

When the model uses additive features as in  $D(e_1 + e_2)$ , the prediction is

$$\tilde{y}(q) = \alpha^{\top} \left[ W((\mathbf{g}_1)_k + (\mathbf{g}_2)_l) \right]_+$$
(42)

$$=\sum_{i=1}^{h} \alpha_{i} \left[ \sum_{j=1}^{d_{1}} W_{i,j} ((\mathbf{g}_{1})_{k,j} + (\mathbf{g}_{2})_{l,j}) \right]_{+}.$$
(43)

<sup>653</sup> The resulting matrix decomposition can then be written as

$$\hat{M} = \sum_{i=1}^{h} \alpha_i \left( \sum_{j=1}^{d_1} W_{i,j} ((\mathbf{g}_1)_j \mathbb{1}_{r_1}^\top + \mathbb{1}_{r_1} (\mathbf{g}_2)_j^\top) \right)_+ = \sum_{i=1}^{h} \alpha_i \left( \sum_{j=1}^{d_1} W_{i,j} (U_j \mathbb{1}_{r_1}^\top + \mathbb{1}_{r_1} V_j^\top) \right)_+$$
(44)

- again with  $U_j = (\mathbf{g}_1)_j, V_j = (\mathbf{g}_2)_j$ .
- <sup>655</sup> Finally, we show that concatenation of features results in a very similar expression.
- <sup>656</sup> When the model uses concatenated features as in  $D(\mathbf{e}_1 \odot \mathbf{e}_2)$  and D is a standard nonconvex MLP <sup>657</sup> (with trainable weights  $W \in \mathbb{R}^{h \times 2d_1}$  and  $\alpha \in \mathbb{R}^{h \times 1}$ ), the prediction at a point q is

$$\tilde{y}(q) = \alpha^{\top} \left[ \left( W((\mathbf{g}_1)_k \odot (\mathbf{g}_2)_l) \right)_+ \right].$$
(45)

Denoting the hidden layer weights as a concatenation of 2 matrices,  $W = (W_1 \odot W_2)$ , where  $W_1, W_2 \in \mathbb{R}^{h \times d_1}$ , we have the following expression:

$$\tilde{y}(q) = \sum_{i=1}^{h} \alpha_i \left( \sum_{j=1}^{d_1} W_{1i,j}(\mathbf{g}_1)_{k,j} + W_{2i,j}(\mathbf{g}_2)_{l,j} \right)_+.$$
(46)

<sup>660</sup> Following similar steps, we express the matrix decomposition as

$$\hat{M} = \sum_{i=1}^{h} \alpha_i \left( \sum_{j=1}^{d_1} W_{1i,j}(\mathbf{g}_1)_j \mathbb{1}_{r_1}^\top + W_{2i,j} \mathbb{1}_{r_1}(\mathbf{g}_2)_j^\top \right)_+$$
(47)

$$=\sum_{i=1}^{h} \alpha_{i} \left( \sum_{j=1}^{d_{1}} W_{1i,j} U_{j} \mathbb{1}_{r_{1}}^{\top} + W_{2i,j} \mathbb{1}_{r_{1}} V_{j}^{\top} \right)_{+},$$
(48)

661 where  $U_j = (\mathbf{g}_1)_j, V_j = (\mathbf{g}_2)_j$ .

By a similar argument to Appendix A.2.3, the maximum attainable rank of all three representations derived here is limited by  $r_1$ .

## 664 A.3 INTERPOLATION COMPARISON

We present 2D image fitting experiments with the *astronaut* image from SciPy, validating ma-665 trix completion analysis summarized in Table 1. We compare 2D GA-Planes models of the form 666  $D(\mathbf{e}_1 \circ \mathbf{e}_2)$  (solid colorful lines) and  $D(\mathbf{e}_1 + \mathbf{e}_2)$  (dotted colorful lines) with the optimal low-rank 667 approximation provided by singular value decomposition (solid black line) in Figure 5. The same 668 experiment is repeated for linear interpolation into the vector (line) features (left) versus nearest 669 neighbor interpolation (right, same as theorems). In this experiment we find qualitatively similar 670 results regardless of the type of interpolation, with slightly better performance using linear interpo-671 lation; in our 3D experiments we use (bi/tri)linear interpolation. 672

As expected, we find that a linear decoder model with multiplication dramatically outperforms its additive counterpart, which does not improve with increasing model size. We also find that 2D GA-Planes models with MLP decoders can match or exceed the compression performance of the optimal low-rank representation found by singular value decomposition (SVD), especially when using a nonconvex MLP. This is a testament to the capacity of an MLP decoder to increase representation rank using fewer parameters than a traditional low-rank decomposition, as well as to the resolution compressibility of natural images.



**Figure 5:** 2D image fitting experiments matching the setting of our theoretical results, with a GA-Planes version using only vector (line) features and a decoder as specified in the legend. Solid lines denote features combined by multiplication; dotted lines use addition. Left: linear interpolation of features; Right: nearest neighbor interpolation of features. For the SVD baseline we use low-rank factors whose resolution matches the target image, so no interpolation is needed (this, and the use of a nonlinear decoder, is why in some cases 2D GA-Planes can outperform SVD).

## 680 A.4 LOWER BOUNDS

Based on the matrix completion theorems and their summary in Table 1, we present lower bounds on the Frobenius norm errors of each 2D GA-Planes model. We denote the optimal fitting error of the linear and MLP decoder models by  $E_{linear}(D(f(q)))$  and  $E_{MLP}(D(f(q)))$  for different feature combinations f(q). For models with a linear decoder,

$$E_{linear}(\mathbf{D}(\mathbf{e}_1 + \mathbf{e}_2)) \ge \sigma_2(M) \tag{49}$$

$$E_{linear}(\mathbf{D}(\mathbf{e}_1 \circ \mathbf{e}_2)) \ge \sigma_k(M) \tag{50}$$

$$E_{linear}(\mathbf{D}(\mathbf{e}_1 \circ \mathbf{e}_2 + \mathbf{e}_{12})) \ge \sigma_k(M - \varphi(L^*)), \tag{51}$$

where  $L^*$  is a downsampled version of the target M, at the same resolution as the feature grid  $g_{12}$ .

686 For models with convex or nonconvex MLP decoders,

$$E_{MLP}(\mathbf{D}(\mathbf{e}_1 + \mathbf{e}_2)) \ge \sigma_{r_1}(M) \tag{52}$$

$$E_{MLP}(\mathbf{D}(\mathbf{e}_1 \circ \mathbf{e}_2)) \ge \sigma_{r_1}(M) \tag{53}$$

$$E_{MLP}(\mathbf{D}(\mathbf{e}_1 \circ \mathbf{e}_2 + \mathbf{e}_{12})) \ge \sigma_{r_1}(M - \varphi(L^*)).$$
(54)

We can see from these bounds that the approximation error of a model can be reduced dramatically by the introduction of a convex or nonconvex MLP decoder, depending on the singular value decay of the target image M.

#### 690 A.5 BENEFITS OF CONVEXITY

In most of our experiments, all models (convex, semiconvex, and nonconvex) are large enough that 691 they are able to optimize well. However, we highlight a benefit of our convex and semiconvex mod-692 els that they enjoy more stable optimization even with very small model sizes. In Figure 6 we com-693 pare test intersection-over-union (IoU) curves for very small models for our video fitting task (hid-694 den dimension 4 in the decoder MLP, and feature dimensions  $[d_1, d_2, d_3] = [4, 4, 2]$  and resolutions 695  $[r_1, r_2, r_3] = [32, 32, 16]$  for line, plane, and volume features, respectively). We repeat optimization 696 with 10 different random seeds used to initialize the optimizable parameters (gating weights for the 697 convex and semiconvex models are fixed). While the convex and semiconvex models enjoy stable 698 training curves across random seeds, we find that the nonconvex model experiences much more 699 volatile training behavior (completely failing to optimize with some of the random seeds). 700



**Figure 6:** Test performance throughout training on our video fitting task, for a very small GA-Planes model across 10 random seeds. We find that the favorable optimization landscape of our convex and semiconvex models enables reliable training across seeds, whereas the nonconvex model fails to fit any test frame with some of the seeds.

On average, the semiconvex model performs best, followed by the convex model, with the nonconvex model performing worst on average. We can also see from the error bars that the convex model is extremely stable, with standard deviations that are visually imperceptible; this is also the case for the semiconvex model in the latter half of training. In contrast, the nonconvex model is highly unstable, with very large standard deviations throughout training. Although with a lucky seed the nonconvex model can outperform the (semi)convex ones, the opposite is true on average.

707 A.6 RESULTS FOR ALL NERFSTUDIO-BLENDER SCENES

708 A.7 EXAMPLE RENDERINGS

In this section, we provide qualitative rendering comparisons on various scenes from the Blender dataset. We highlight the superior performance of GA-Planes with limited number of parameters by comparing the smallest K-Planes, TensoRF and GA-Planes models.

712 A.8 MODEL CONFIGURATIONS USED FOR EXPERIMENTS

713 A.8.1 RADIANCE FIELD MODELING

In our NeRFStudio experiments, we find that because GA-Planes contains features with different dimensionalities (line, plane, and volume), it experiences little loss in quality over a wide range of model sizes. For small model sizes, we allocate most of the model memory to the line features, since their spatial resolution grows linearly with parameter count. As the parameter budget grows, we allocate more parameters to the plane features, whereas the performance of the line-only model stagnates with increasing size. Similarly, as model size grows even further we allocate more parameters to the volume features, whose memory footprint grows cubically with spatial resolution.

The original K-planes model uses 2 proposal networks with different resolutions (as noted in Ta-721 ble 6) and a fixed channel dimension of 8 for both. The resolutions and channel dimensions 722 for either K-planes model (with vs. without proposal sampling) refer to  $r_2$  and  $d_2$ , respectively. 723 TensoRF model resolutions and channel dimensions can be interpreted in a similar way, since 724 their feature combination dictates that  $d_1 = d_2$  and they initialize the line and plane grids with 725 the same resolution. The only nuance is that TensoRF constructs separate features for color and 726 density decoding. Hence, the channel dimensions for density and color features are listed. In-727 stead of a multiresolution scheme, TensoRF starts from the base resolution of  $r_1 = r_2 = 128$ 728



**Figure 7:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *lego* scene.

and upsamples the grids to reach the final resolutions on Table 6. Resolutions listed under GA-729 Planes should be interpreted as  $[r_1, r_2, r_3]$ ; channel dimensions as  $[d_1, d_2, d_3]$ . For all mod-730 els that use the multiresolution scheme, the base resolutions (i.e.  $[r_1, r_2, r_3]$ ) are multiplied 731 with the upsampling factors. For instance, a base resolution  $[r_1, r_2, r_3]$  with channel dimensions 732  $[d_1, d_2, d_3]$  and multiresolution copies  $[m_1, m_2, m_3]$  will generate the grids of GA-Planes as fol-733 lows: Linear feature grids  $\mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3$  will have the shapes  $\{[m_1r_1, d_1], [m_2r_1, d_1], [m_3r_1, d_1]\}$ , plane 734 grids  $\mathbf{g}_{12}, \mathbf{g}_{23}, \mathbf{g}_{13}$  will have the shapes  $\{[m_1r_2, m_1r_2, d_2], [m_2r_2, m_2r_2, d_2], [m_3r_2, m_3r_2, d_2]\}, [m_3r_2, m_3r_2, d_2]\}$ 735 and the volume grid  $\mathbf{g}_{123}$  will have the shape  $[r_3, r_3, r_3, d_3]$ . Although we don't use multires-736 olution copies for the volume grid in GA-Planes, we do use multiresolution for the volume-737 only GA-Planes ablation. The resolution for that model refers to  $r_3$ , and the channel dimen-738 sions are also allowed to vary for each resolution (unlike other variants with multiresolution, 739 where the feature dimension is fixed across resolutions). If we denote these varying feature di-740 mensions as  $[d_{3a}, d_{3b}, d_{3c}]$ , the multiresolution copies of the volume grids will have the shapes 741  $\{[m_1r_3, m_1r_3, m_1r_3, d_{3a}], [m_2r_3, m_2r_3, m_2r_3, d_{3b}], [m_3r_3, m_3r_3, m_3r_3, d_{3c}]\}.$ 742

## 743 A.8.2 3D SEGMENTATION

GA-Planes model uses feature dimensions  $[d_1, d_2, d_3] = [36, 24, 8]$  (with  $\odot$ ) or  $[d_1, d_2, d_3] = [25, 25, 8]$  (with  $\circ$ ) and resolutions  $[r_1, r_2, r_3] = [128, 32, 24]$ . Multiresolution grids are not used for this task since density prediction can be achieved by a simpler architecture. The model size is 0.22 M. Tri-Planes model has the feature dimension  $d_2 = 4$ , and resolution  $r_2 = 128$  resulting in a total number of parameters of 0.2 M. Note that we fix these sizes across (non/semi)convex formulations, which causes slight variations in the size of the decoder, however, the grids constitute the most number of parameters, making this effect negligible.

#### 751 A.8.3 VIDEO SEGMENTATION

GA-Planes model uses feature dimensions  $[d_1, d_2, d_3] = [32, 16, 8]$  and resolutions  $[r_1, r_2, r_3] = [128, 128, 64]$ . When the features are combined by multiplication in the nonconvex model,  $d_1 = 154$   $d_2 = 16$ . Multiresolution grids are not used for this task. The model size is 2.9 M. Tri-Planes



**Figure 8:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *chair* scene.

model has the feature dimension  $d_2 = 59$ , and resolution  $r_2 = 128$  resulting in a total number of parameters of 2.9 M.



**Figure 9:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *mic* scene.



**Figure 10:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *drums* scene.



**Figure 11:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *hotdog* scene.



**Figure 12:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *materials* scene.



**Figure 13:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *ficus* scene.



**Figure 14:** Results on radiance field reconstruction. Nonconvex GA-Planes (with feature multiplication) offers the most efficient representation: when the model is large it performs comparably to the state of the art models, but when model size is reduced it retains higher performance than other models. Here all models are trained for the same number of epochs on the *ship* scene.



**Figure 15:** Rendering comparison for the *chair* scene: TensoRF on the left (0.32 M parameters), K-Planes in the middle (0.39 M parameters), GA-Planes on the right (0.25 M parameters).



**Figure 16:** Rendering comparison for the *mic* scene: TensoRF on the left (0.32 M parameters), K-Planes in the middle (0.39 M parameters), GA-Planes on the right (0.25 M parameters).



**Figure 17:** Rendering comparison for the *lego* scene: TensoRF on the left (0.32 M parameters), K-Planes in the middle (0.39 M parameters), GA-Planes on the right (0.25 M parameters).



**Figure 18:** Rendering comparison for the *materials* scene: TensoRF on the left (0.32 M parameters), K-Planes in the middle (0.39 M parameters), GA-Planes on the right (0.25 M parameters).

Model	Resolutions	Channel Dimen-	Multiresolution	Proposal Network	Number of model
		sions		Resolutions	parameters (M)
	32	4	[1, 2, 4]	[32, 64]	0.390
	32	8	[1, 2, 4]	[32, 64]	0.649
	64	4	[1, 2, 4]	[64, 128]	1.533
K-plane	64	8	[1, 2, 4]	[128, 256]	4.041
	64	16	[1, 2, 4]	[128, 256]	6.107
	128	8	[1, 2, 4]	[128, 256]	10.234
	128	16	[1, 2, 4]	[128, 256]	18.493
	32	4	[1, 2, 4]	-	0.298
	40	4	[1, 2, 4]	-	0.444
	64	4	[1, 2, 4]	-	1.073
K-plane without proposal sampling	64	8	[1, 2, 4]	-	2.108
	100	6	[1, 2, 4]	-	3.822
	128	10	[1, 2, 4]	-	10.367
	129	16	[1, 2, 4]	-	16.824
	128	[2, 4]	-	-	0.320
	256	[2, 4]	-	-	1.207
	256	[6, 8]	-	-	2.786
TangaBE	256	[12, 12]	-	-	4.760
Tensokr	300	[12, 16]	-	-	7.609
	300	[16, 24]	-	-	10.860
	300	[32, 32]	-	-	17.362
	300	[16, 48]	-	-	17.364
	[200, 4, 4]	[32, 32, 4]	[1, 2, 4]	-	0.254
	[200, 8, 4]	[32, 32, 4]	[1, 2, 4]	-	0.351
	[200, 16, 8]	[32, 32, 4]	[1, 2, 4]	-	0.740
GA-plane	[200, 32, 8]	[16, 16, 4]	[1, 2, 4]	-	1.164
_	[100, 100, 16]	[6, 6, 8]	[1, 2, 4]	-	3.874
	[200, 128, 32]	[10, 10, 8]	[1, 2, 4]	-	10.681
	[200, 128, 32]	[16, 16, 8]	[1, 2, 4]	-	16.908
	32	4	[1, 2, 4]	-	0.301
	64	4	[1, 2, 4]	-	1.078
GA-plane ablation-VM	128	4	[1, 2, 4]	-	4.180
	128	6	[1, 2, 4]	-	6.251
	128	12	[1, 2, 4]	-	12.465
CA-plane ablation-CP	200	32	[1, 2, 4]	-	0.196
GA-plane ablation-CF	200	64	[1, 2, 4]	-	0.355
	18	[3, 5, 6]	[1, 2, 3]	-	1.236
CA-plane ablation-volume	24	[4, 4, 6]	[1, 2, 3]	-	2.778
GA-plane ablauon-volume	32	[4, 4, 6]	[1, 2, 3]	-	6.529
	32	[4, 6, 8]	[1, 2, 3]	-	8.824

**Table 6:** Model configurations used for the radiance field modeling task on the Blender dataset.