

# UMIFormer: Mining the Correlations between Similar Tokens for Multi-View 3D Reconstruction

Zhenwei Zhu<sup>†</sup> Liying Yang<sup>†</sup> Ning Li Chaozhao Jiang Yanyan Liang<sup>\*</sup>

Macau University of Science and Technology, Faculty of Innovation Engineering

## Abstract

In recent years, many video tasks have achieved breakthroughs by utilizing the vision transformer and establishing spatial-temporal decoupling for feature extraction. Although multi-view 3D reconstruction also faces multiple images as input, it cannot immediately inherit their success due to completely ambiguous associations between unordered views. There is not usable prior relationship, which is similar to the temporally-coherence property in a video. To solve this problem, we propose a novel transformer network for Unordered Multiple Images (UMIFormer). It exploits transformer blocks for decoupled intra-view encoding and designed blocks for token rectification that mine the correlation between similar tokens from different views to achieve decoupled inter-view encoding. Afterward, all tokens acquired from various branches are compressed into a fixed-size compact representation while preserving rich information for reconstruction by leveraging the similarities between tokens. We empirically demonstrate on ShapeNet and confirm that our decoupled learning method is adaptable for unordered multiple images. Meanwhile, the experiments also verify our model outperforms existing SOTA methods by a large margin.

## 1. Introduction

3D reconstruction, which lifts 2D view images to a 3D representation of an object, is a challenging problem. It plays an important role in numerous technologies, including intelligent driving, augmented reality and robotics. In the situation of single-view input, previous works have attempted to improve performance by strengthening the network capabilities [8, 12, 20, 30, 36] and leveraging the geometric information as priors knowledge [35, 39, 40]. However, for multi-view reconstruction, researchers concentrate on how to extract sufficient feature representation for the object shape from unordered multiple images

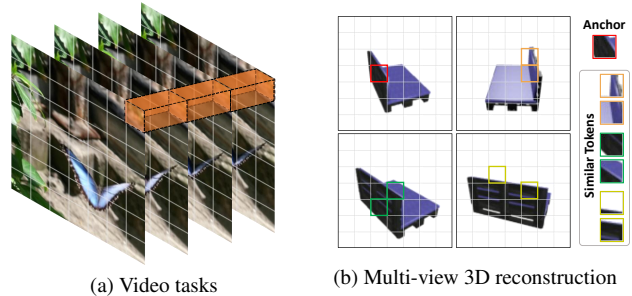


Figure 1. Comparison of the positional correspondence used for inter-image-decoupled feature extraction in (a) video tasks and (b) multi-view 3D reconstruction. Video tasks exploit the temporally-coherence property to establish the prior relationship as shown in (a). For multi-view reconstruction, each patch is treated as an anchor and associated with its similar tokens from other views to build the positional correspondence as shown in (b).

[4, 28, 33, 34, 38]. This paper is devoted to multi-view 3D reconstruction using the voxel representation.

In our investigation, the deep-learning-based algorithms for multi-view reconstruction typically involve two steps: feature extraction and shape reconstruction. The latter is generally accomplished using a 3D decoder module, while there are various solutions for the former including CNN-based and transformer-based methods.

CNN-based methods usually separate the feature extraction process into two stages. The first stage exploits a backbone network to encode on intra-view-dimension while the second stage processing on inter-view-dimension aggregate the features obtained from different views. The fusion method can be a pooling layer [9, 18, 21], a recurrent unit [4, 10, 15] or an attention operator [38]. In addition, Pix2Vox series [33, 34] put merger after decoder, which directly fusion the voxels predicted from different views, and also achieves good results. GARNet [44] sets up two fusion modules, which are located before and after the decoder respectively, to realize the adaptability of the merger to the global state.

Transformer-based methods [25, 28, 37] that can directly handle views as a sequence also attend global awareness.

<sup>†</sup>These authors contributed equally to this work.

<sup>\*</sup>Corresponding authors.

It takes natural advantage of its architecture to couple the procedures for intra-image and inter-image feature extraction. However, such approaches work poorly when facing few views as input since the size of the extracted feature is too small to hold enough information. In contrast, 3D-RETR [19] exploits transformer on the intra-view dimension and then aggregates the features from different views using an adaptive pooling layer. It is essentially the same method as before [21] but with a more advanced backbone network. As a result, it inherits the drawback of inadequate view association mining. Nevertheless, the success of this method reminds us that the power of vision transformer (ViT) [5] for the representation of views cannot be understated.

In video tasks that also face multiple images as input, recent works [1, 2, 13] have produced good performance using ViT as a spatially-decoupled feature extractor and additionally establish a temporally-decoupled feature extractor. They benefit from the fact that video frames are temporally coherent (as shown in Figure 1a). These approaches, however, cannot be directly transferred to our task, because multi-view reconstruction should deal with unordered multiple images without prior positional correspondence.

To address this problem, we propose a novel inter-view-decoupled block (IVDB) based on mining the correlation between similar patches from different views (as shown in Figure 1b). It can be inserted between the blocks of ViT to create a transformer encoder for unordered multiple images. This model maintains the advantages of ViT initialization pre-trained on large-scale datasets while alternating decoupling the intra- and inter-view encoding processes. Moreover, by clustering the tokens according to their similarities and exploiting a down-sampling transformer block, the tokens from all branches are compressed into a fixed-size compact representation, ensuring relatively steady performance for the varying number of views input.

In detail, our contributions are as follows:

- To our best knowledge, we are the first to propose a transformer network that alternates decoupled intra- and inter-view feature extraction for multi-view 3D reconstruction, a problem facing unordered multiple images as input.
- Leveraging the correlations between similar tokens, we propose a novel inter-view-decoupled block (IVDB) that rectifies the tokens according to the related information from other views and a similar-token merger (STM) that compresses the features from all branches.
- Experiments on ShapeNet [31] verify that our model achieves a performance better than previous SOTA methods by a large margin.

## 2. Related Works

### 2.1. Multi-View 3D Reconstruction

In the early years, the traditional algorithms, e.g. SFM [17] and SLAM [7], build mappings from 2D pixels to 3D positions based on feature matching. They are hard to deal with complex situation of view images. At present, neural network algorithms become the mainstream for solving multi-view reconstruction.

Among them, CNN-based methods usually extract features from each view in parallel and then these features are aggregated as the representation of the shape. Most of the research works focus on the fusion approach. [9, 18, 21] employ pooling based fusion method that concatenates the feature maps from different views and then compresses them to a specified size by a maximum pooling or an average pooling layer. Despite being straightforward, it performs poorly because it lacks learnable parameters. 3D-R2N2 series [4, 15] and LSM [10] use recurrent neural network (RNN)-based fusion method that treats features from images as a sequence. However, a recurrent unit cannot satisfy invariant to permutations and is not suitable for facing a large number of views input due to limited long-term memory. Attsets [38], Pix2Vox series [33, 34] and GAR-Net [44] exploits attention-based fusion method that accumulates the features from different views weighted according to the score maps predicted by an extra branch.

To learn the relatively complex latent correlation between different views, transformer-based methods are proposed. EVoIT [28], LegoFormer [37] and 3D-C2FT [25] treats the input as a sequence on the inter-view-dimension. However, their reconstruction quality is terrible when facing few view images due to the insufficient size of the feature. 3D-RETR [19] deals with tokens on intra-view-dimension. It utilizes ViT [5] to extract features from each view and fuses them using the pooling-based fusion. Although the strong representation learning ability of the transformer for images is utilized, the potential information between views is not fully discovered.

### 2.2. Transformer for Multiple Images

The transformer paradigm is proposed by [27] for natural language processing. ViT [5] widely extends it to computer vision and mainly works on single-image. Some research about video tasks introduces the transformer network to solve the problems facing multiple images as input. Leveraging prior relationships based on spatial and temporal, TimeSformer [2] proposes decoupled spatial-temporal attention where the two kinds of attention operation coexist in a transformer block and DSTT [13] decouples the spatial and temporal encoding into separate transformer blocks which are used alternately. ViViT [1] factorises the multi-head dot-product attention operation to execute the two

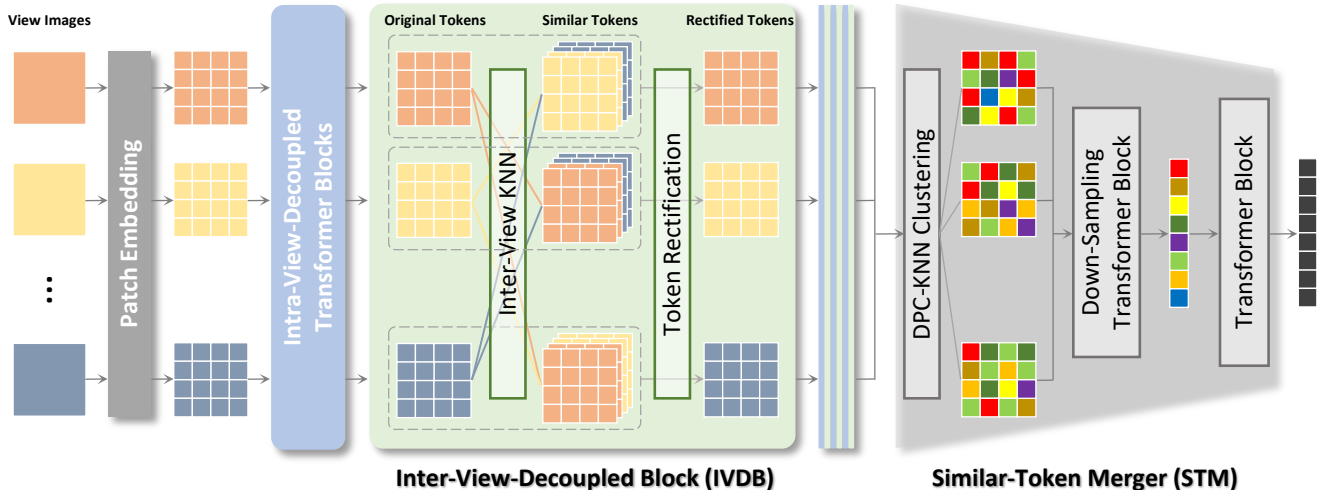


Figure 2. The architecture used for feature extraction in UMIFormer. The network encodes unordered multiple images utilizing the intra-view-decoupled transformer block and the inter-view-decoupled block alternately. Then, the feature is compressed to a compact representation by the similar-token merger.

kinds of decoupled attention in parallel. In addition, DeViT [3], FGT [43] and E<sup>2</sup>FGVI [11] further mine the relationship in spatial or temporal to acquire a better representation.

These decoupling methods achieve good representation learning capability because they exploit spatial-temporal relations. However, they cannot be transferred to solve multi-view reconstruction because there is no inter-image-coherent in the problem of processing unordered multiple images.

### 3. Methods

According to an arbitrary number of view images  $\mathcal{I} = \{I_1, I_2, \dots, I_n\}$  with sizes of  $224^2 \times 3$  of an object, our model is to generate the corresponding voxel representation  $V$  with a size of  $32^3$ . To begin, views are used to extract a feature representation  $f$  (described in section 3.1). Then, the binary voxel  $V$  will be constructed according to the feature (described in section 3.2). The entire process is formulated as:

$$V = \text{UMIFormer}(\mathcal{I}) = R(E(I_1, I_2, \dots, I_n)), \quad (1)$$

where  $E$  and  $R$  donate the processes of feature extraction and shape reconstruction respectively.

#### 3.1. Feature Extraction

The feature extractor is designed based on ViT [5]. It contains four types of blocks: patch embedding, intra-view-decoupled transformer block, inter-view-decoupled block (IVDB) and similar-token merger (STM). The first two of them are derived from the ViT structure. Patch embedding

is consist of splitting images into patches, linearly mapping them and adding position embeddings. The foundation structure of ViT is actually assembled by all the intra-view-decoupled transformer blocks.

IVDBs (elaborated in Section 3.1.1), which construct the relationships between various views, are periodically inserted into the ViT backbone. Thus, it is feasible to alternate intra- and inter-view-dimension encoding. Comparing to the approaches that separate the two encoding modes, which extract features from each view and then fuse them, our method mines richer correlations between different views. Comparing to the approaches that couple these two dimensions, our method is equivalent to providing prior knowledge to reduce the complexity of representation learning.

STM (elaborated in Section 3.1.2), at the end of the extractor, downsamples the feature obtained from all branches into a compact representation. Note that STM is utilized to compress the feature, as opposed to being employed for aggregating like the merger blocks in previous works. Because the connection between various views has been created by IVDBs and the transformer-based decoder can handle sequences with variable lengths, it is not necessary to set a specific fusion function. STM is designed for enabling the extractor to provide fixed-size features for reconstruction when receiving varying numbers of views, ensuring relatively stable performance.

The process of feature extraction is shown in Figure 2. The patched view images are alternately decoupled encoding by the transformer blocks and IVDBs and then compressed into a compact representation using STM.

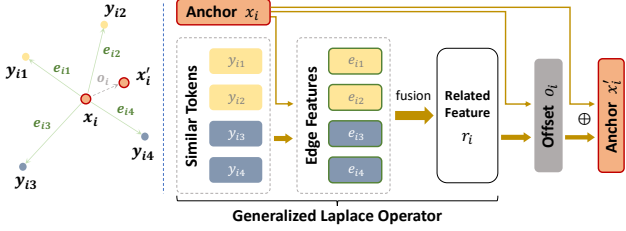


Figure 3. Visual illustration of the token rectification used in IVDB.

### 3.1.1 Inter-View-Decoupled Block

Since the unordered images have no prior positional correspondence, we consider that it can be modeled by the relationship between similar tokens from different views. Therefore, we adopt an inter-view KNN layer, which takes each token as an anchor  $x_i$  and matches it with the nearest  $k$  tokens  $y_i$  from each other view in Euclidean space. A token represents a particular region of the object and similar tokens are the representation of related regions. Therefore, we hold that the relationship created by similar tokens can be used to replace the lacking position correspondence.

To exploit the above relationships, we propose token rectification for mining the inter-view correlation, illustrated in Figure 3. In essence, the tokens can be treated as point clouds in a high-dimensional manifold space. Inspired by [29], we utilize similar tokens to support the anchor providing more accurate feature representations.

Firstly, a Generalized Laplace Operator extracts the related feature  $r_i$  of the anchor  $x_i$ . This operator embeds the edge features based on the related positional relationship between anchor and its similar tokens in feature space and aggregate them through an attention-based fusion [38]. The definition is as follows:

$$e_{ij} = \text{MLP}_{\text{edge}}(y_{ij} - x_i), \quad (2)$$

$$r_i = \text{Fusion}(e_{i1}, e_{i2}, \dots, e_{i(k(n-1))}), \quad (3)$$

where  $\text{MLP}_{\text{edge}}$  is a MLP (multilayer perceptron) with non-linear activation function and  $\text{Fusion}$  indicates the attention-based fusion approach. Then, the related feature help to rectify the anchor by predicting the feature offset  $o_i$ . Compared with mapping directly to update the anchor, adding offset better preserves the advantages of ViT initialization pre-trained on large-scale datasets due to keeping the token in the original feature space. We propose two strategies for predicting offsets:

$$o_i = \text{MLP}_{\text{os}}(x_i, r_i), \quad (4)$$

$$o_i = \text{MLP}_{\text{ow}}(x_i, r_i) \times x_i, \quad (5)$$

where  $\text{MLP}_{\text{os}}$  is an MLP with a ReLU function to predict the offset directly, while  $\text{MLP}_{\text{ow}}$  is an MLP with a tanh

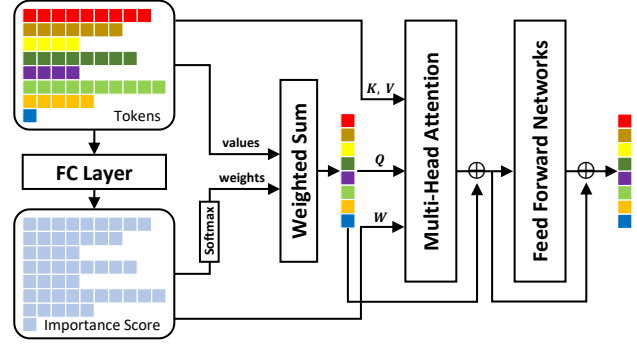


Figure 4. Visual illustration of the down-sampling transformer block used in STM.

function to predict the offset weight. We empirically verify that the latter is more favorable for multi-view reconstruction. Finally, the rectified anchor token is produced by a straightforward addition operation:

$$x'_i = x_i + o_i. \quad (6)$$

### 3.1.2 Similar-Token Merger

The previous network fetches a large number of tokens from views. Several of these tokens contain relatively close information that may cause information redundancy, especially when facing extremely high input amounts. Moreover, the view images are composed of both foreground and background while only the foreground information supports reconstruction. As a result, we propose the feature compression method that maximizes diversity while minimizing irrelevant information for the preserved feature.

Inspired by [41], we establish the similarity relationship of tokens again. The tokens are divided into  $g$  groups by the DPC-KNN clustering [6] and then fed into the down-sampling transformer block (illustrated in Figure 4). We fuse the features from each group to obtain an aggregated token set using the attention-based fusion method same as  $\text{Fusion}$  in Equation 3. The set is entered into multi-head attention (MHA) as  $Q$  to extract information from the ungrouped tokens which provide  $K$  and  $V$ . In contrast to the general MHA, we introduce the extra weights  $W$ , which reuse the importance score predicted by the additional branch in the attention-based fusion process, to ensure that tokens with different importance have different effects on the result. This MHA is following [41] and defined as:

$$\text{Attention}(Q, K, V, W) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}} + W\right)V, \quad (7)$$

where  $d_k$  is the dimensions of  $Q$ ,  $K$  and  $V$ . At the end of STM, the down-sampled feature is further processed by

Methods		1 view	2 views	3 views	4 views	5 views	8 views	12 views	16 views	20 views
CNN-Based	<b>3D-R2N2</b> [4]	0.560 / 0.351	0.603 / 0.368	0.617 / 0.372	0.625 / 0.378	0.634 / 0.382	0.635 / 0.383	0.636 / 0.382	0.636 / 0.382	0.636 / 0.383
	<b>AttSets</b> [38]	0.642 / 0.395	0.662 / 0.418	0.670 / 0.426	0.675 / 0.430	0.677 / 0.432	0.685 / 0.444	0.688 / 0.445	0.692 / 0.447	0.693 / 0.448
	<b>Pix2Vox++</b> [34]	0.670 / <b>0.436</b>	0.695 / 0.452	0.704 / 0.455	0.708 / 0.457	0.711 / 0.458	0.715 / 0.459	0.717 / 0.460	0.718 / 0.461	0.719 / 0.462
	<b>GARNet</b> [44]	0.673 / 0.418	0.705 / 0.455	0.716 / 0.468	0.722 / 0.475	0.726 / 0.479	0.731 / 0.486	0.734 / 0.489	0.736 / 0.491	0.737 / 0.492
	<b>GARNet+</b>	0.655 / 0.399	0.696 / 0.446	0.712 / 0.465	0.719 / 0.475	0.725 / 0.481	0.733 / 0.491	0.737 / 0.498	0.740 / 0.501	0.742 / 0.504
Transformer-Based	<b>EVolt</b> [28]	- / -	- / -	- / -	0.609 / 0.358	- / -	0.698 / 0.448	0.720 / 0.475	0.729 / 0.486	0.735 / 0.492
	<b>Legoformer</b> [37]	0.519 / 0.282	0.644 / 0.392	0.679 / 0.428	0.694 / 0.444	0.703 / 0.453	0.713 / 0.464	0.717 / 0.470	0.719 / 0.472	0.721 / 0.472
	<b>3D-C2FT</b> [25]	0.629 / 0.371	0.678 / 0.424	0.695 / 0.443	0.702 / 0.452	0.702 / 0.458	0.716 / 0.468	0.720 / 0.475	0.723 / 0.477	0.724 / 0.479
	<b>3D-RETR</b> (3 views)	0.674 / -	0.707 / -	0.716 / -	0.720 / -	0.723 / -	0.727 / -	0.729 / -	0.730 / -	0.731 / -
	<b>3D-RETR</b> <sup>†</sup> [19]	0.680 / -	0.701 / -	0.716 / -	0.725 / -	0.736 / -	0.739 / -	0.747 / -	0.755 / -	0.757 / -
	<b>UMIFormer</b>	<b>0.6802</b> / 0.4281	<b>0.7384</b> / <b>0.4919</b>	<b>0.7518</b> / <b>0.5067</b>	0.7573 / <b>0.5127</b>	0.7612 / 0.5168	0.7661 / 0.5213	0.7682 / 0.5232	0.7696 / 0.5245	0.7702 / 0.5251
	<b>UMIFormer+</b>	0.5672 / 0.3177	0.7115 / 0.4568	0.7447 / 0.4947	<b>0.7588</b> / 0.5104	<b>0.7681</b> / <b>0.5216</b>	<b>0.7790</b> / <b>0.5348</b>	<b>0.7843</b> / <b>0.5415</b>	<b>0.7873</b> / <b>0.5451</b>	<b>0.7886</b> / <b>0.5466</b>

Table 1. Evaluation and comparison of the performance on ShapeNet using IoU  $\uparrow$  / F-Score@1%  $\uparrow$ . The best results are highlighted in bold. <sup>†</sup> The results in this row are derived from models that train individually for the various number of input views.

a transformer block into the ultimate feature representation extracted from the multiple image views.

It is a great promotion for the compactness of the compressed representation that similar features are stored in one token or a few tokens. Besides, even tokens corresponding to background information account for a large proportion of the input, this block can significantly reduce their occupancy in the final feature.

### 3.2. Shape Reconstruction

We employ a decoder composed of a transformer stage and a CNN stage for shape reconstruction, which shares the same structure as [19]. The transformer stage contains 8 transformer decoder blocks while excluding any upsampling layers. A feature map with a size of  $64 \times 768$  is generated. After reshaping to  $4^3 \times 768$ , it entered into the CNN stage and upsample to  $32^3$  voxel gradually.

### 3.3. Loss Function

The task to reconstruct the shape of an object can be seen as a voxel-level segmentation for occupied or empty. Consequently, the loss function is defined as Dice loss [16] between predicting volume and the ground truth (GT). The previous work [19] indicates that it is suitable for 3D reconstruction, the problem with an extremely unbalanced amount of samples between categories. Mathematically, this loss function is defined as:

$$\mathcal{L} = 1 - \frac{\sum_{i=1}^{32^3} p_i g t_i}{\sum_{i=1}^{32^3} p_i + g t_i} - \frac{\sum_{i=1}^{32^3} (1 - p_i) (1 - g t_i)}{\sum_{i=1}^{32^3} 2 - p_i - g t_i} \quad (8)$$

where  $p$  and  $gt$  indicate the confidence of the grids on the reconstructed volume and GT.

## 4. Experiments

### 4.1. Datasets and Implementation Details

Following [4], our experiments are primarily carried out on a subset of ShapeNet [31] to evaluate the ability of multi-view 3D reconstruction. The subset includes 13 categories and 43,783 objects with a 3D representation and rendered images from 24 random poses. Moreover, single-view reconstruction experiments on the chair category from Pix3D [23] dataset including 2,894 data with untruncated and unoccluded view image are supplemented to verify that our model is capable of handling real-world data. The reconstruction results are measured using both 3D Intersection over Union (IoU) and F-Score@1% [24, 34].

We adopt the pre-training model of DeiT-B [26], a variant of ViT, to initialize the intra-view-decoupled transformer blocks in our model. To facilitate visualization and analysis, the cls token and distillation token are removed. The model contains 12 transformer blocks and we insert the IVDB with  $k = 5$  after every third block. For STM,  $k_{dpc}$  in DPC-KNN clustering and  $g$  are defined as 15 and 196. For multi-view 3d reconstruction, we eventually provide two models with the same structure called UMIFormer and UMIFormer+, whose input view numbers are respectively fixed to 3 and 8 during training. The models are trained by an AdamW optimizer [14] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  with a batch size of 32 for 150 epochs. The learning rate is defined as 1e-4 and decreases by 0.1 after 50 and 120 epochs sequentially. UMIFormer is trained on 2 Tesla V100 for 2 days and UMIFormer+ is trained on 8 Tesla V100 for 2.5 days. The fixed threshold for binarizing the probabilities is set to 0.5 for UMIFormer and 0.4 for UMIFormer+.



Figure 5. Multi-view reconstruction results on the test set of ShapeNet when facing 5 views, 10 views, 15 views and 20 views as input.

IVDB	Fusion	3 views	5 views	8 views	12 views	16 views	20 views
✗	PBM	0.7325	0.7406	0.7447	0.7472	0.7487	0.7493
✗	ABM	0.7394	0.7479	0.7522	0.7545	0.7560	0.7566
✗	STM	0.7477	0.7557	0.7587	0.7598	0.7606	0.7606
✓	PBM	0.7372	0.7453	0.7488	0.7514	0.7530	0.7536
✓	ABM	0.7412	0.7503	0.7548	0.7574	0.7588	0.7593
✓	STM	<b>0.7518</b>	<b>0.7612</b>	<b>0.7661</b>	<b>0.7682</b>	<b>0.7696</b>	<b>0.7702</b>

Table 2. The ablation experiments on ShapeNet evaluated by IoU about IVDB and STM. Among them, STM is compared with two mainstream fusion methods: pooling-based merger (PBM) and attention-based merger (ABM).

## 4.2. Multi-view 3D Reconstruction Results

The performance qualification results of methods are shown in Table 1. Undoubtedly, UMIFormer has a significant advantage over the previous methods in almost all metrics. It outperforms current SOTA methods by a large margin. Even training 3D-RETR models separately for different input view numbers (the row marked in gray) trails our model by a big gap. Furthermore, UMIFormer+ boasts a more powerful capability for multi-view reconstruction. Whereas, it has a somewhat limited capacity for single-view reconstruction.

Figure 5 shows two examples of the reconstruction results when facing various view amounts as inputs. Com-

pared with the other methods, our two models produce more accurate results for rifle reconstruction. In addition, the texture on predicted volumes is improved significantly with more views. It not only demonstrates the effectiveness of our algorithm but also verifies that our model can continue to mine information from the increasing input. For lamp reconstruction, our models, especially UMIFormer+, realize a relatively complete representation for the intermediate bracket, which is difficult for the other methods. Certainly, it also demonstrates the strong learning ability of our feature extractor for details.

## 4.3. Ablation Experiments

Ablation analysis on IVDB and STM is based on experimental results as shown in Table 2.

**Effect of IVDB.** We observe that employing IVDB can consistently improve the reconstruction performance for various amounts of view inputs. Table 3 presents experimental results related to several rectification strategies (discussed in Section 3.1.1) used in IVDB. Token rectification by mapping directly performs terribly. Because it disrupts the prior knowledge of the backbone network learned during pre-training by mapping tokens to a new feature space. As a result, the benefits of pre-training are significantly diminished. For the other two strategies, predicting offset works better for the case of few input views and predict-

Rectification Strategy	3 views	5 views	8 views	12 views	16 views	20 views
FC Mapping	0.6935	0.7022	0.7049	0.7038	0.7028	0.7022
Offset	<b>0.7528</b>	<b>0.7614</b>	0.7648	0.7664	0.7674	0.7678
Offset Weight	0.7518	0.7612	<b>0.7661</b>	<b>0.7682</b>	<b>0.7696</b>	<b>0.7702</b>

Table 3. Comparison of performance on ShapeNet evaluated by IoU when using different rectification strategies in IVDB. FC mapping refers to using a fully connected layer to map the concatenation of an anchor and its related features to the rectified token directly. Offset and offset weight respectively correspond to the strategies defined in Equation 4 and 5.

ing offset weight is suitable for processing a large number of input views. In this paper, we adopt predicting offset weight for token rectification.

**Effect of STM.** All three types of merger — pooling-based merger [21], attention-based merger [38] and our proposed STM — compress features from all branches to a fixed size of  $196 \times 768$  in our network. Notably, the model using STM can achieve better reconstruction performance. It verifies that STM preserves richer information than other compression methods.

Furthermore, experimental results indicate that using IVDB and STM together performs much better than using them alone.

#### 4.4. Evaluation on Real-World Dataset

Pix3D dataset [23] is usually used as the testing set for evaluating the performance of single-view reconstruction. Most of the view images in it are real-world images with complex backgrounds. Therefore, we attempt to validate the effectiveness of UMIFormer on this dataset to verify that it works on various domains of image. Following previous works, the training set uses the data from the chair category in ShapeNet and the view images are re-synthesized by Render for CNN [22] with random backgrounds from the SUN database [32]. Among them, each object includes 60 view images.

IVDB is not used in this experiment because it is irrelevant for single-view input. In addition,  $g$  in STM is defined as 32. Table 4 shows the performance qualification results of UMIFormer and other SOTA methods which can be used for both single-view reconstruction and multi-view reconstruction. In comparison, our model slightly outperforms them. Figure 6 shows that our method performs better on the restoration of texture, particularly for the thin strip shapes.

#### 4.5. Visualization of Similar Tokens

In our algorithm, mining the correlations between similar tokens is involved several times. To further investigate the behavior of our algorithm, we attempt to visualize the pertinent procedures. Taking the UMIFormer dealing

IoU $\uparrow$			
Pix2Vox++	3D-RETR	GARNet	UMIFormer
0.279	0.297	0.291	<b>0.300</b>
F-Score@1% $\uparrow$			
Pix2Vox++	3D-RETR	GARNet	UMIFormer
0.113	0.125	0.116	<b>0.129</b>

Table 4. Evaluation and comparison of the performance for single-view reconstruction on Pix3D.

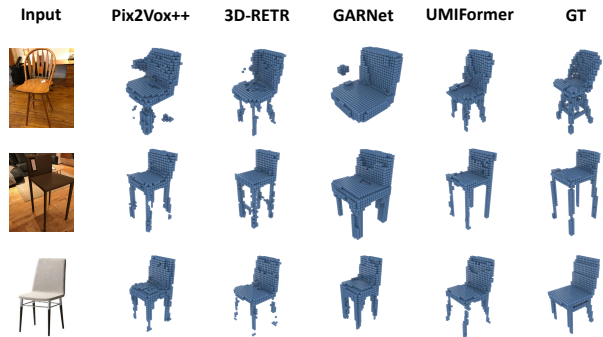


Figure 6. Single-view reconstruction results on real-world data (Pix3D test set).

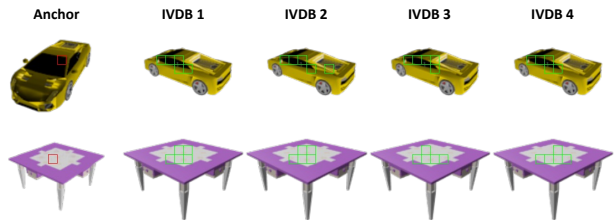


Figure 7. Visualization of several anchors (marked with red frames) and their similar tokens (marked with green frames) paired by inter-view KNN in IVDBs during multi-view reconstruction processing.

with 2 view inputs from ShapeNet as an example, Figure 7 shows the results that an anchor (marked with a red frame) finds its similar tokens (marked with green frames) from another view image through the inter-view KNN layers in the 4 IVDBs. We observe that regions relating to the anchor content are surrounded by the matched tokens. Therefore, there is indeed a semantic relationship between them, which can be used as the position correspondence for decoupling.

Figure 8 shows some examples of token grouping by the clustering layer in STM. As anticipated, the patches corresponding to the foreground regions are divided finely, whereas the background patches are only assigned to a few groups. Nevertheless, cluster maps do not clearly distinguish subject edges since the features for clustering have been highly abstracted by the previous layers. It does not affect STM to maximize diversity while minimizing irrelevant information at the feature level.

Encoder		1 view	2 views	3 views	4 views	5 views	8 views	12 views	16 views	20 views
Independent Branches		<b>0.6923</b>	0.7292	0.7394	0.7445	0.7479	0.7522	0.7545	0.7560	0.7566
Video Transformer	Joint Attention [42]	0.6771	0.7112	0.7201	0.7241	0.7264	0.7275	0.7269	0.7261	0.7252
	Factorised Transformer Block [13]	0.6809	0.7179	0.7288	0.7337	0.7365	0.7403	0.7424	0.7437	0.7441
	Factorised Attention [2]	0.6918	0.7287	0.7400	0.7451	0.7485	0.7528	0.7552	0.7566	0.7572
	Factorised Dot-Product [1]	0.6684	0.7139	0.7275	0.7331	0.7373	0.7423	0.7445	0.7457	0.7463
Ours		0.6802	<b>0.7384</b>	<b>0.7518</b>	<b>0.7573</b>	<b>0.7612</b>	<b>0.7661</b>	<b>0.7682</b>	<b>0.7696</b>	<b>0.7702</b>

Table 5. Comparison of the reconstruction performance exploiting different decoupling strategies in feature extractor. To control the variables, all of them are based on the structure of ViT. For more setup details refer to the supplementary material.

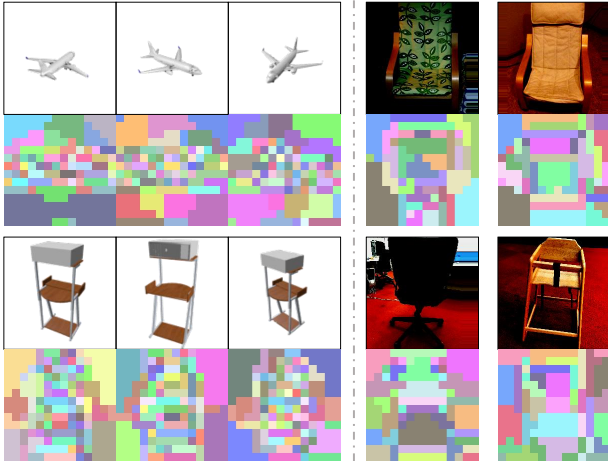


Figure 8. Visualization of clustering results in STM, including the case of multi-view reconstruction on ShapeNet (left) and single-view reconstruction on Pix3D (right).

## 5. Discussion

In this work, we expect to leverage ViT and decoupled encoding strategy to establish a robust representation learning model for multiple image inputs. The effectiveness of the model has been verified by several works [1, 2, 13] about video transformer algorithms. Although both video tasks and multi-view reconstruction are oriented to multi-image input problems, their inter-image relationships are quite different. The inter-view-decoupled encoder established in this way is theoretically not suitable for processing unordered multiple images.

In Table 5, we compare the performance of multi-view 3D reconstruction using different feature extraction methods. Encoding for independent branches using ViT as shown in the first row treats as the baseline. To control the variables, the four video transformer methods are all implemented based on ViT, which may differ from the architecture in the original paper. Encoder with joint attention means that all tokens from various views are processed uniformly without decoupled encoding. We observe that the performance is significantly worse than the baseline and the reconstruction results are worse when dealing with more

than 8 views as input. Since the positional encoding for each image in a multi-view reconstruction is consistent, the prior association between tokens in the attention layer is completely lost. It is very challenging to mine the correlation solely relying on the adaptive capability of the network. This problem becomes more serious as the number of views increases. The other three video transformer approaches respectively use factorised transformer block, factorised attention layers and factorised dot products in the self-attention layer to establish decoupled encoding. These methods increase the connection between branches relative to the baseline while they do not conform to the nature of unordered multiple images. Actually, these kinds of inter-view-decouple feature extraction methods do not simplify the multi-view reconstruction problem. Therefore, these methods cannot achieve the desired performance. Summarizing these results, we confirm that decoupling based on mining the correlation between similar tokens is most suitable for handling unordered multiple images.

## 6. Conclusion

**Limitations.** 1) Our model requires large memory occupation since it is based on the parallel transformer network and devoted to a 3D reconstruction task. Therefore, it is hard to generalize to high-resolution voxel reconstruction under existing mainstream hardware devices. 2) The computational consumption of both the inter-view KNN layer and the DPC-KNN clustering layer grows exponentially with the increasing number of input views. Therefore, the predicting efficiency of our algorithm is not dominant when facing an extremely high number of view inputs.

In this paper, we propose a transformer-based method for multi-view 3D reconstruction which achieves brilliant performance. The feature extractors alternate decoupled intra- and inter-view encoding for unordered multiple images by mining the correlation between similar tokens. In future work, we expect to overcome its limitation on higher-resolution reconstruction by compressing the model and alleviate the inference efficiency problem by accelerating KNN and DPC-KNN clustering algorithm. Furthermore, this encoding mode may also be extended to other issues involving unordered multiple images.



## References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6836–6846, 2021.
- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, volume 2, page 4, 2021.
- [3] Jiayin Cai, Changlin Li, Xin Tao, Chun Yuan, and Yu-Wing Tai. Devit: Deformed vision transformers in video inpainting. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 779–789, 2022.
- [4] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *European conference on computer vision*, pages 628–644. Springer, 2016.
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [6] Mingjing Du, Shifei Ding, and Hongjie Jia. Study on density peaks clustering based on k-nearest neighbors and principal component analysis. *Knowledge-Based Systems*, 99:135–145, 2016.
- [7] Jorge Fuentes-Pacheco, José Ruiz-Ascencio, and Juan Manuel Rendón-Mancha. Visual simultaneous localization and mapping: a survey. *Artificial intelligence review*, 43(1):55–81, 2015.
- [8] Rohit Girdhar, David F Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *European Conference on Computer Vision*, pages 484–499. Springer, 2016.
- [9] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2821–2830, 2018.
- [10] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. *Advances in neural information processing systems*, 30, 2017.
- [11] Zhen Li, Cheng-Ze Lu, Jianhua Qin, Chun-Le Guo, and Ming-Ming Cheng. Towards an end-to-end framework for flow-guided video inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17562–17571, 2022.
- [12] Caixia Liu, Dehui Kong, Shaofan Wang, Jinghua Li, and Baocai Yin. A spatial relationship preserving adversarial network for 3d reconstruction from a single depth view. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 18(4):1–22, 2022.
- [13] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Decoupled spatial-temporal transformer for video inpainting. *arXiv preprint arXiv:2104.06637*, 2021.
- [14] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [15] Tingsong Ma, Ping Kuang, and Wenhong Tian. An improved recurrent neural networks for 3d object reconstruction. *Applied Intelligence*, 50(3):905–923, 2020.
- [16] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. IEEE, 2016.
- [17] Onur Özyeşil, Vladislav Voroninski, Ronen Basri, and Amit Singer. A survey of structure from motion\*. *Acta Numerica*, 26:305–364, 2017.
- [18] Despoina Paschalidou, Osman Ulusoy, Carolin Schmitt, Luc Van Gool, and Andreas Geiger. Raynet: Learning volumetric 3d reconstruction with ray potentials. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3897–3906, 2018.
- [19] Zai Shi, Zhao Meng, Yiran Xing, Yunpu Ma, and Roger Wattenhofer. 3d-retr: End-to-end single and multi-view 3d reconstruction with transformers. *arXiv preprint arXiv:2110.08861*, 2021.
- [20] Edward J Smith and David Meger. Improved adversarial systems for 3d object generation and reconstruction. In *Conference on Robot Learning*, pages 87–96. PMLR, 2017.
- [21] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015.
- [22] Hao Su, Charles R Qi, Yangyan Li, and Leonidas J Guibas. Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views. In *Proceedings of the IEEE international conference on computer vision*, pages 2686–2694, 2015.
- [23] Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B Tenenbaum, and William T Freeman. Pix3d: Dataset and methods for single-image 3d shape modeling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2974–2983, 2018.
- [24] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3405–3414, 2019.
- [25] Leslie Ching Ow Tiong, Dick Sigmund, and Andrew Beng Jin Teoh. 3d-c2ft: Coarse-to-fine transformer for multi-view 3d reconstruction. *arXiv preprint arXiv:2205.14575*, 2022.

- [26] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [28] Dan Wang, Xinrui Cui, Xun Chen, Zhengxia Zou, Tianyang Shi, Septimiu Salcudean, Z Jane Wang, and Rabab Ward. Multi-view 3d reconstruction with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5722–5731, 2021.
- [29] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019.
- [30] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Advances in neural information processing systems*, 29, 2016.
- [31] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.
- [32] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010.
- [33] Haozhe Xie, Hongxun Yao, Xiaoshuai Sun, Shangchen Zhou, and Shengping Zhang. Pix2vox: Context-aware 3d reconstruction from single and multi-view images. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2690–2698, 2019.
- [34] Haozhe Xie, Hongxun Yao, Shengping Zhang, Shangchen Zhou, and Wenxiu Sun. Pix2vox++: Multi-scale context-aware 3d object reconstruction from single and multiple images. *International Journal of Computer Vision*, 128(12):2919–2935, 2020.
- [35] Zhen Xing, Yijiang Chen, Zhixin Ling, Xiangdong Zhou, and Yu Xiang. Few-shot single-view 3d reconstruction with memory prior contrastive network. In *European Conference on Computer Vision*, pages 55–70. Springer, 2022.
- [36] Zhen Xing, Hengduo Li, Zuxuan Wu, and Yu-Gang Jiang. Semi-supervised single-view 3d reconstruction via prototype shape priors. In *European Conference on Computer Vision*, pages 535–551. Springer, 2022.
- [37] Farid Yagubbayli, Alessio Tonioni, and Federico Tombari. Legoformer: Transformers for block-by-block multi-view 3d reconstruction. *arXiv preprint arXiv:2106.12102*, 2021.
- [38] Bo Yang, Sen Wang, Andrew Markham, and Niki Trigoni. Robust attentional aggregation of deep feature sets for multi-view 3d reconstruction. *International Journal of Computer Vision*, 128(1):53–73, 2020.
- [39] Shuo Yang, Min Xu, Haozhe Xie, Stuart Perry, and Jiahao Xia. Single-view 3d object reconstruction from shape priors in memory. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3152–3161, 2021.
- [40] Yang Yang, Junwei Han, Dingwen Zhang, and Qi Tian. Exploring rich intermediate representations for reconstructing 3d shapes from 2d images. *Pattern Recognition*, 122:108295, 2022.
- [41] Wang Zeng, Sheng Jin, Wentao Liu, Chen Qian, Ping Luo, Wanli Ouyang, and Xiaogang Wang. Not all tokens are equal: Human-centric visual analysis via token clustering transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11101–11111, 2022.
- [42] Yanhong Zeng, Jianlong Fu, and Hongyang Chao. Learning joint spatial-temporal transformations for video inpainting. In *European Conference on Computer Vision*, pages 528–543. Springer, 2020.
- [43] Kaidong Zhang, Jingjing Fu, and Dong Liu. Flow-guided transformer for video inpainting. *arXiv preprint arXiv:2208.06768*, 2022.
- [44] Zhenwei Zhu, Liying Yang, Xuxin Lin, Chaohao Jiang, Ning Li, Lin Yang, and Yanyan Liang. Garnet: Global-aware multi-view 3d reconstruction network and the cost-performance tradeoff. *arXiv preprint arXiv:2211.02299*, 2022.