# 3D Gated Recurrent Fusion for Semantic Scene Completion

Yu Liu$^{1\dagger}$, Jie Li$^{2\dagger}$, Qingsen Yan$^3$, Xia Yuan$^2$, Chunxia Zhao$^2$, Ian Reid$^1$ and Cesar Cadena$^4$ *

## Abstract

*This paper tackles the problem of data fusion in the semantic scene completion (SSC) task, which can simultaneously deal with semantic labeling and scene completion. RGB images contain texture details of the object(s) which are vital for semantic scene understanding. Meanwhile, depth images capture geometric clues of high relevance for shape completion. Using both RGB and depth images can further boost the accuracy of SSC over employing one modality in isolation. We propose a 3D gated recurrent fusion network (GRFNet), which learns to adaptively select and fuse the relevant information from depth and RGB by making use of the gate and memory modules. Based on the single-stage fusion, we further propose a multi-stage fusion strategy, which could model the correlations among different stages within the network. Extensive experiments on two benchmark datasets demonstrate the superior performance and the effectiveness of the proposed GRFNet for data fusion in SSC. Code will be made available.*

## 1. Introduction

Understanding the surroundings is a fundamental capability for many real-world applications such as augmented reality [2], robot grasping [37], or autonomous navigation [9]. Different abstractions are possible, and even complementary. Semantic labeling of the scene allows for a high level reasoning, while 3D geometry completion enables basic spatial capabilities. Semantic scene completion (SSC) aims at solving both simultaneously.

An RGB-D sensor allows acquiring depth information from the scene along side the RGB image. On the one hand, RGB image contains rich details about the color and texture, which are the primary cues for the semantic scene understanding. On the other hand, depth carries more clues about the object geometry and distance information, which are much reliable in reflecting the position, shape, and occlusion relationship between objects within the scene. Many vision applications have already benefit from using both modalities in their tasks, such as object detection [16, 3], video segmentation [12, 36, 10], action recognition [21, 20, 42], or visual SLAM [23, 40, 29]. Recent studies [13, 25] in SSC also demonstrate that employing both, RGB image and depth, can outperform using only one modality [35].

However, fusing the information from RGB and depth is still an unsolved problem, and becomes an obstacle which hinders the performance of SSC. Albeit some recent works conduct data fusion between RGB and depth, they usually employ some, "manually" set, basic operation to fuse the data. Those includes *sum fusion* [25, 18], *max fusion* [22], *concatenate fusion* [7, 15], *transform fusion* [38] and *bilinear fusion* [28]. Nevertheless, RGB and depth data are not equivalent quantities, while still providing complementary yet redundant information. Therefore, we propose to extract the information in a selective manner from both modalities, and fuse them accordingly with respect to the specific task.

We present the Gated Recurrent Fusion (GRF) block, which can provide adaptive selection and aggregation of RGB and depth information. On the one hand, the *gate* component in the GRF fusion block selects, in an adaptive manner, various positions of different importance in aligned RGB-D frames regarding to the contribution from both modalities. The *gate* effectively selects valid information while filters out the irrelevant one. On the other hand, the *memory* component in the GRF fusion block effectively preserves the complementary information, which can compensate the missing or ambiguous details of the data obtained from different modalities.

Furthermore, the GRF fusion block offers the flexibility to be cascaded to a multi-stage configuration that combines high-level and low-level features. GRF fusion block is extended from the Gated Recurrent Unit (GRU) [5]. Based on the GRF fusion block, we build the GRFNet for the semantic scene completion, and provide single- and multi-stage

---

*Yu Liu and Jie Li contributed equally to this work.

$^1$Y. Liu and I. Reid are with School of Computer Science, The University of Adelaide, 5005, North Terrace, SA `yu.liu04@adelaide.edu.au`

$^2$J. Li, X. Yuan and C. Zhao are with School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing, 210094, China `jieli_cn@163.com`

$^3$ Q. Yan is with School of Computer Science and Engineering, Northwestern Polytechnical University, Xi'an, 710072, China `yqs@mail.nwpu.edu.cn`

$^4$C. Cadena is with Autonomous Systems Lab, ETH Zurich, Leonhardstrasse 21, 8092, Zurich `cesarc@ethz.ch`
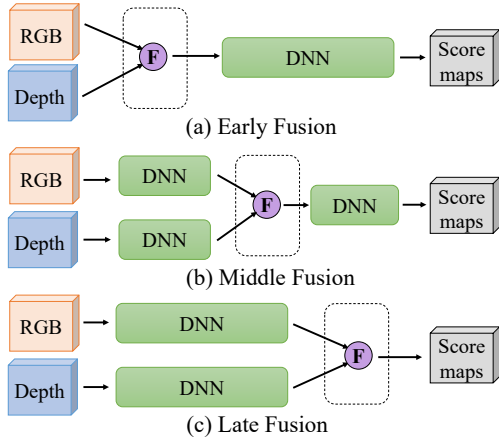
Figure 1. Fusion at different stages. Early fusion contains more low-level features, while the input of the late fusion contains more abstract high-level features.

fusion versions of GRFNet. In the single-stage fusion version, depth and RGB features are fed into the same GRF fusion block individually. In the multi-stage fusion version, depth and RGB features of different stages form an interleaved sequence and are input into the same GRF fusion block consecutively.

The multi-stage version takes advantage of both low-level and high-level features, and achieves better performance than the single-stage version.

In summary, the contributions of this work are mainly two-fold:

- An end-to-end 3D-GRF based network, GRFNet, is presented for fusing RGB and depth information in the SSC task, through employing *gate* and *memory* components, the selection and fusion between two modalities can be conducted effectively. To the best of our knowledge, this is the first time that gated recurrent network is employed for data fusion in the SSC task.

- Within the framework of GRFNet, single-stage and multi-stage fusion strategies are proposed. While outperforming existing fusing strategies in the SSC task already with the single stage, the multi-stage fusion proves to give the best results.

Extensive experiments demonstrate that the proposed GRFNet achieves superior performance on NYU [33] and NYUCAD [11] datasets.

## 2. Related Work

In this section, we briefly look through the deep learning based methods for SSC, with emphasis on the discussion of the existing multi-modality fusion strategies.

### 2.1. Semantic Scene Completion

The goal of SSC is to produce a complete 3D voxel representation for a scene from a single-view input. Specifi-
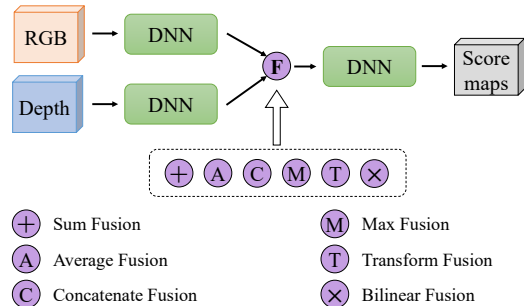


Figure 2. Several typical single-stage fusion methods.

cally, Song *et al*. ([35]) propose an end-to-end 3D convolutional network (SSCNet) which is based on the single-view depth as input that can simultaneously predict the results of scene completion and semantic labeling. SSCNet has high computational costs due to the adoption of 3D convolutions. Zhang *et al*. ([43]) introduce spatial group convolution (SGC) into SSC for accelerating the computation of 3D dense prediction task. Meanwhile, Han *et al*. ([17]) employ the long short-term memory (LSTM) to recover missing parts of 3D shapes. Dai *et al*. ([8]) use a coarse-to-fine strategy to handle large scenes with varying spatial extent. Although the depth-based approach has made significant progress, the absence of texture details prevents improving SSC.

In order to incorporating the color information, TS3D [13] introduces the RGB image into SSC and uses a 2D network to acquire semantic segmentation results. Semantic outputs of the RGB stream are concatenated with inputs of the depth stream to obtain the completed 3D scene. DDR-SSC [25] uses two parallel feature extraction branches with the same structure to obtain information from RGB and depth simultaneously. A multi-stage structure with element-wise addition is employed to perform feature fusion. Thanks to the semantic information provided by RGB, the semantic labeling accuracy of both TS3D and DDR-SSC has significantly been improved compared to SSCNet. However, none of these methods takes into account the selective fusion of multi-modal information, which limit those algorithms to achieve better performance.

### 2.2. Fusion Schemes

The RGB-D information fusion is important to many vision applications. In general, the fusion scheme can be divided into three categories, *e.g.* early fusion [7], middle fusion [31] and late fusion [41, 34], as shown in Figure 1. According to the stages of fusion, these schemes can also be divided into single-stage fusion and multi-stage fusion [18, 30].

**Single-Stage Fusion** There are several general patterns for single-stage fusion as shown in Figure 2. Specifically, Sum fusion [25, 18] computes the sum of the two feature maps
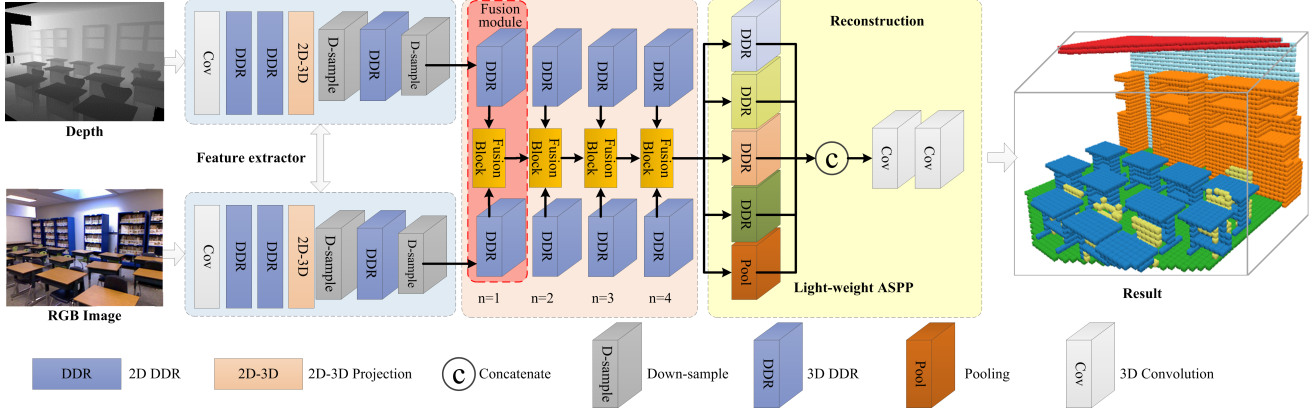
Figure 3. The network architecture of GRFNet is extended from Dimensional Decomposition Residual (DDR) network [25]. GRFNet has two feature extractors to capture the features from depth and RGB images respectively. The feature extractor contains a projection layer to map 2D feature to 3D space. The GRF fusion block (denoted by yellow boxes in the middle) replaces the original fusion unit to take full advantage of the multi-modal information. With two DDR plus their corresponding GRF fusion block to form a fusion module, also named single-stage fusion module (denoted by the red box in one column). GRFNet is composed of a multi-stage ( 4-stage here) fusion module. Then we use light-weight ASPP to obtain multiple receptive fields information. Different colors of the DDR block denote various receptive fields. Then the network uses two 3D convolutions to predict occupancies and object labels simultaneously.

at the same spatial locations. Average fusion is essentially a weighted sum fusion with equal weights. Max fusion [22] takes the feature with the maximal value from multiple feature maps. Concatenate fusion stacks the features with channels [7, 15]. Wang *et al*. ([38]) propose an encoder-decoder architecture which exchanges the information of multi-modal data in the latent space. Bilinear fusion [28] computes an outer matrix product of the two features at each pixel location.

There are a few methods that consider the complementary and selectivity of data fusion. Specifically, Li *et al*. ([26]) develop a novel LSTM model to fuse scene contexts adaptively. Cheng *et al*. ([4]) use the concatenated feature maps of RGB and depth to learn an array $G$ to weight the contribution of one input modality and $1 - G$ to weight the other input modality. Wang *et al*. ([39]) use the same strategy as [4] to fuse the feature maps from RGB and Depth in saliency detection. However, Li *et al*. ([26]) only consider the complementarity of information but ignore the selectivity of the data, the other two methods only consider the selectivity of information and cannot guarantee the complementarity of information. Moreover, these methods are single-stage fusion and lack of scalability.

**Multi-stage Fusion** According to the way for aggregating multi-modal information, this paper divides multi-stage fusion algorithms into merge fusion, cross fusion, and external fusion. Hazirbas *et al*. ([18]) adopt the merge fusion structure to fuse the two branches of features extracted from RGB and depth images. The feature maps from depth are fused into the RGB branch by stages with an element-wise summation. Wang *et al*. ([38]) use cross fusion to merge the

common features of RGB and depth, and keep the modality specific features separated from each other. Both Park *et al*. ([30]) and Li *et al*. ([25]) use an external fusion mechanism. Specifically, Li *et al*. ([25]) capture features of RGB and depth image at different levels, these features at each level are fused separately and then assembled all at once before the reconstruction part. Park *et al*. ([30]) propose RDFNet to fuse multi-modal features separately by multiple fusion blocks, and refine the fused features one by one through a set of refine blocks. In RDFNet, each fusion introduces an additional fusion block with a new set of extra parameters. The artificially designed fusion blocks are complex and require multiple parameters that are not easy to migrate to other applications. These multi-stage fusion methods use high-level and low-level features achieving high accuracy. However, each fusion block within the multi-stage mostly adopts concatenation or summation, ignoring the adaptive selection of the multi-modal data.

On the contrary, our proposed GRF fusion block extends the standard gated recurrent unit (GRU), where the gate and the memory structures can adaptively select and preserve valid information. Besides, GRFNet adopts the form of a recurrent network. When performing multi-stage fusion, GRF modules exploit parameter sharing. And experiments show that both of the proposed single- and multi-stage GRFNets achieve better accuracy than previous methods.

## 3. Methodology

### 3.1. Overview

Our proposal, GRFNet, extends the network architecture of DDR-SSC [25], and focuses on improving the fusion

strategy. We subtly adopt the gate structure and memory mechanism in GRU unit to form a multi-modal feature GRF fusion block with the power of autonomous selectivity and adaptive memory preservation. Moreover, taking advantage of its recurrent nature, we further propose a multi-stage fusion strategy to utilize both low-level and high-level features with introducing insignificant parameters.

In the feature extractor part, the network uses dimensional decomposition residual (DDR) blocks to extract the local textures and the geometry information. A projection layer is employed to connect the 2D and 3D parts. The multi-stage fusion module consists of four single-stage fusion modules that can effectively combines the RGB features and depth features. The fused features are fed into the subsequent light-weight atrous spatial pyramid pooling (LW-ASPP [25]). After that, another two point-wise convolutional layers are used to predict the semantic labels for each voxel in the 3D volume.

The network maps each voxel to one of the labels $C = c_0, c_1, \cdots c_N$, where $N$ is the number of semantic classes, and $c_0$ represents the empty voxel.

## 3.2. Gated Recurrent Unit

Gated recurrent unit (GRU) [5] is a popular model in recurrent neural networks (RNN) and has an outstanding performance in many natural language processing (NLP) tasks [6, 24]. A GRU has two gate structures, a memory structure, and can be reused recurrently. However, few researchers have explored the power of GRU in the field of 3D vision, especially for feature fusion. We find that GRU highly aligns with our requirements for an effective multi-modal fusion strategy in SSC.

The gate structure in GRU enables the selective fusion of multi-modal features. The memory structure ensures that valid information can be retained for future fusion purpose. The characteristics of its recurrent network enable GRU to be reused in the multi-stage fusion while sharing the same set of parameters. Compared to GRU, the structure of Long short-term memory (LSTM) [19] is more complicated and has an extra forget gate with more parameters. ConvGRU [1] is a convolutional version of GRU. We extend ConvGRU to 3D convolutional in our GRFNet, and modify it to fit the feature fusion in SSC task.

## 3.3. Gated Recurrent Fusion Block of RGB-D Features

As shown in Figure 4, at fist step (left), gated recurrent fusion (GRF) block takes one of the RGB-D features as input. The outputs of this step will be used as the hidden state. Then, in the second step (right), GRF fusion block takes the features of another modality as input. These two steps reuse the GRF fusion block and share the same set of parameters. Next, we will use the first step with input $f^d$ as an example
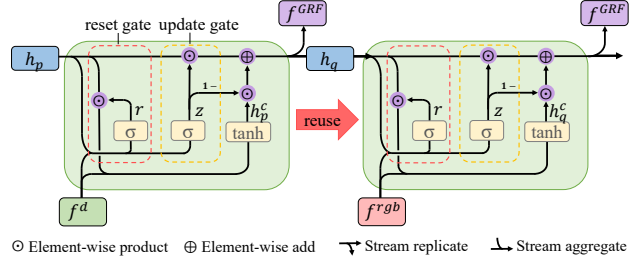


Figure 4. GRF fusion block. At step $p$, the input of GRF fusion block is one of the features from depth or RGB, and in the next step, input is the other. Both GRF fusion blocks share the same set of parameters.

to explain in detail the principle and work flow of the GRF fusion block.

In contrast to the commonly used GRU for encoding information in a temporal way, the way we use GRF to fuse RGB+Depth features is somehow different. Specifically, the GRU handles the fusion of RGB and depth information in a 'modality', rather than a 'sequential', way. And for the multi-stage fusion process, similar as other deep neural networks, it is more like the low-level multi-modal feature to guide the following high-level multi-modal feature to be merged.

**Hidden State** The hidden state $h_p$, along with the current input $f$ to control the reset and update gates. The output of the previous stage will be used as the hidden state of current stage. At the first step of fusion between depth feature $f^d$ and RGB feature $f^{rgb}$, that is $p = 1$, we use the sum fusion of two modal features to initialize the hidden state, as $h_0 = f^d + f^{rgb}$.

**Reset Gate** At step $p$, the hidden state $h_p$ and the current input $f^d$ together to decide the status of the reset gate $r$ by

$$r = \sigma \left( W_r \left( f^d, h_p \right) \right) \tag{1}$$

The two feature stream $f^d$ and $h_p$ are concatenated and fed into a convolution operation. $W_r$ represents the corresponding weights in the convolution. The sigmoid function $\sigma$ converts each value in the feature tensor into the range of $(0, 1)$ and acts as a gate signal.

**Update Gate** The update gate $z$ is also decided by the hidden state $h_p$ and input features $f^d$. Through another convolution operation with weight $W_z$ and the sigmoid function $\sigma$, we get $z$ as,

$$z = \sigma \left( W_z \left( f^d, h_p \right) \right) \tag{2}$$

Theoretically, reset and update gates essentially learn a set of weights that control the amount of information that is retained or discarded, experimental studies in section 4 show the effectiveness of the reset and update gates.

**Adaptive Memory** Through the element-wise product $\odot$, the reset gate $r$ determines how much information in the

| | scene completion | | | semantic scene completion | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | prec. | recall | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| Lin *et al.* ([27]) | 58.5 | 49.9 | 36.4 | 0.0 | 11.7 | 13.3 | 14.1 | 9.4 | 29.0 | 24.0 | 6.0 | 7.0 | 16.2 | 1.1 | 12.0 |
| Geiger *et al.* ([14]) | 65.7 | 58.0 | 44.4 | 10.2 | 62.5 | 19.1 | 5.8 | 8.5 | 40.6 | 27.7 | 7.0 | 6.0 | 22.6 | 5.9 | 19.6 |
| SSCNet [35] | 57.0 | **94.5** | 55.1 | 15.1 | **94.7** | 24.4 | 0.0 | 12.6 | 32.1 | 35.0 | 13.0 | 7.8 | 27.1 | 10.1 | 24.7 |
| EsscNet [43] | **71.9** | 71.9 | 56.2 | 17.5 | 75.4 | 25.8 | 6.7 | 15.3 | **53.8** | 42.4 | 11.2 | 0 | 33.4 | 11.8 | 26.7 |
| DDR-SSC [25] | 71.5 | 80.8 | 61.0 | 21.1 | 92.2 | **33.5** | 6.8 | 14.8 | 48.3 | 42.3 | 13.2 | **13.9** | 35.3 | 13.2 | 30.4 |
| GRFNet | 68.4 | 85.4 | **61.2** | **24.0** | 91.7 | 33.3 | **19.0** | **18.1** | 51.9 | **45.5** | **13.4** | 13.3 | **37.3** | **15.0** | **32.9** |

Table 1. Results on the NYU dataset [33]. Bold numbers represent the best scores.

| | scene completion | | | semantic scene completion | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | prec. | recall | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| Zheng *et al.*( [44]) | 60.1 | 46.7 | 34.6 | - | - | - | - | - | - | - | - | - | - | - | - |
| Firman *et al.*( [11]) | 66.5 | 69.7 | 50.8 | - | - | - | - | - | - | - | - | - | - | - | - |
| SSCNet [35] | 75.4 | **96.3** | 73.2 | 32.5 | 92.6 | 40.2 | 8.9 | 33.9 | 57.0 | **59.5** | 28.3 | 8.1 | 44.8 | 25.1 | 40.0 |
| TS3D [13] | 80.2 | 91.0 | 74.2 | 33.8 | **92.9** | 46.8 | **27.0** | 27.9 | **61.6** | 51.6 | 27.6 | **26.9** | 44.5 | 22.0 | 42.1 |
| DDR-SSC [25] | **88.7** | 88.5 | 79.4 | **54.1** | 91.5 | 56.4 | 14.9 | 37.0 | 55.7 | 51.0 | 28.8 | 9.2 | 44.1 | 27.8 | 42.8 |
| GRFNet | 87.2 | 91.0 | **80.1** | 50.3 | 91.8 | **58.1** | 18.4 | **42.7** | 60.6 | 52.8 | **34.6** | 11.5 | **46.6** | **30.8** | **45.3** |

Table 2. Results on the NYUCAD dataset [44]. Bold numbers represent the best scores.

past needs to be "memorized".

$$h'_p = r \odot h_p \qquad (3)$$

when the reset gate $r$ is close to 1, the "memorized" information $h'_p$ will be kept and then passed to the current fusion operation with current input feature $f^d$. The preserved "memory" $h'_p$ and feature $f^d$ are concatenated together to perform a linear transformation (convolution), and then activated by a $\tanh$ function.

$$h^c_p = \tanh\left(W_h\left(f^d, h'_p\right)\right) \qquad (4)$$

$h^c_p$ acts similarly to the memory cell in the LSTM and helps the GRF fusion block to remember long term information within the multi-stage fusion.

**Selective Fusion** $z \odot h_p$ : Indicates how much of the previous features should be preserved. $(1 - z) \odot h^c_p$: Indicates how much of the current information $h^c_p$ should be added. Similar to the former, here $(1-z)$ forgets some unimportant information in $h^c_p$. Or, it can be viewed as a choice of some information in $h^c_p$.

Combined with $f^d$ and $h_p$, the fusion result at step $p$ is,

$$h_q = z \odot h_p + (1 - z) \odot h^c_p \qquad (5)$$

This operation ignores some information in previous hidden state $h_p$, and adds some information from the current step. Update gate $z$ is equivalent to the *forget gate* in LSTM, and $1-z$ is equivalent to the *input gate* in LSTM. In this way, the *forget gate* $z$ and the *input gate* $(1-z)$ are linked. That is, if the previous information is ignored with a weight of $z$, then the information for the current input $h^c_p$ would be selected with a weight of $(1 - z)$. In our case, if the information in previous stage is depth feature and current input is RGB feature, this enables the complementary information to be

effectively merged. Accordingly, the output of the current step $f^{GRF}_q = h_p$ will also be passed to the next step.

**Single-stage Fusion Module** The bimodal information passes through multiple DDRs for feature extraction, and each process of the DDR corresponds to a stage. That is, single-stage fusion module has only one layer of DDR from RGB and depth branch, and the fusion block is performed after the DDR. In specific, the GRF module has two input features, $f^d$ and $f^{rgb}$. At step 1, the hidden state is initialised as mentioned above. We feed $f^d$ into GRF fusion block, and get the output $h_1$. Then at step 2, we reuse the same structure and the same parameters in the GRF fusion block. The input hidden state is replaced by $h_1$, and the input is the features $f^{rgb}$ extracted from the RGB image. As shown in Figure 4, we use the red line with an arrow to indicate the reuse of GRF fusion block at step 2.

**Multi-stage GRF Fusion Module** Features extracted by the earlier-stage DDR are at relatively low-level, while those by the later-stage DDR are at relatively high-level regarding to the semantic meaning representation. For multi-stage fusion, the features of the two modal data at each stage will be formed as a sequence. Taking the $N$-stage RGB-D fusion as an example, the feature sequence is $F = \left(f^d_1, f^{rgb}_1, f^d_2, f^{rgb}_2, \cdots, f^d_N, f^{rgb}_N\right)$. Each feature tensor in $F$ will be fed into the GRF fusion block serially. The GRU fusion block will be reused $2N$ times and all these fusion stages share the same group of parameters. Different with the single-stage fusion module which performs the multi-modal feature fusion at only one stage of the network, multi-modal features are fused in multi-stages which covers both of the high-level and low-level features. It is not only helpful to recover the details of the scene, but also important to propagate information among different stages.

Using the low-level feature to guide high-level feature is a common and reasonable approach in computer vision

5

| NYU | scene completion | | | semantic scene completion | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| method | prec. | recall | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| single-stage GRFNet | 66.5 | **85.9** | 60.1 | **27.5** | **92.9** | 28.1 | 10.7 | 14.9 | **60.1** | 33.8 | **17.3** | 10.1 | 30.4 | 14.7 | 31.0 |
| multi-stage GRFNet | **68.4** | 85.4 | **61.2** | 24.0 | 91.7 | **33.3** | **19.0** | **18.1** | 51.9 | **45.5** | 13.4 | **13.3** | **37.3** | **15.0** | **32.9** |
| NYUCAD | scene completion | | | semantic scene completion | | | | | | | | | | | |
| single-stage GRFNet | **88.4** | 89.1 | 79.7 | 50.0 | 91.4 | 56.4 | **18.7** | 41.3 | 56.8 | 52.7 | 33.5 | **16.3** | 45.2 | 30.0 | 44.8 |
| multi-stage GRFNet | 87.2 | **91.0** | **80.1** | **50.3** | **91.8** | **58.1** | 18.4 | **42.7** | **60.6** | **52.8** | **34.6** | 11.5 | **46.6** | **30.8** | **45.3** |

Table 3. Results of single-stage GRFNet and multi-stage GRFNet on both NYU and NYUCAD dataset.



(a) RGB and Depth images    (b) Ground truth    (c) GRFNet    (d) DDR-SSC    (e) SSCNet
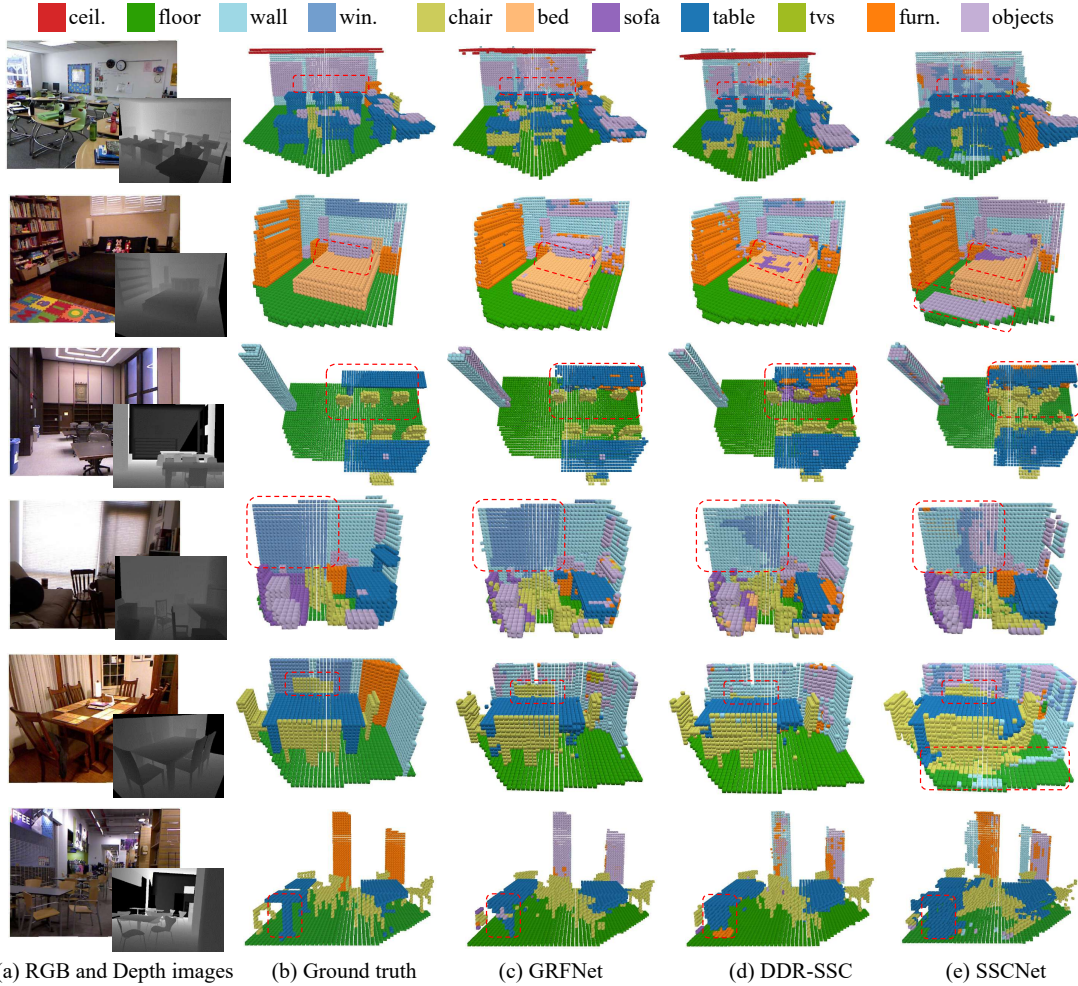
Figure 5. Qualitative results on NYUCAD. From left to right: Input RGB-D image, ground truth, results generated by our GRFNet, DDR-SSC [25], and SSCNet [35]. Overall, our completed semantic 3D scenes are less cluttered and show a higher voxel class accuracy compared to the others.

community. The proposed GRF module can preserves the complementary information and compensates the missing details. Particularly, the color and texture details in RGB image as well as the geometry and distance information in depth are complementary to each other. The "gate" structure in the GRF fusion block controls the feature fusion by learning weights (between 0 and 1). Moreover, for multi-stage fusion, the GRF fusion block has the potential to manage the interleaved modalities. To be specific, the bi-modal information has gradually transformed into abstract semantic information through the network, and the differ-

ence between their distributions is gradually reduced. Due the above merits, we employ the multi-stage GRF fusion module in our network. The effectiveness of multi-stage fusion module are supported and reflected by our experiments in section 4.

## 3.4. Training Protocol

The loss function used in our training process is the soft-max cross-entropy loss, and it is performed on the unnor-

| Method (NYU) | GS | MM | prec. | recall | IoU | mIoU |
|---|---|---|---|---|---|---|
| Concatenate Fusion | | | 70.6 | 76.2 | 57.6 | 25.9 |
| Sum Fusion | | | 67.6 | 79.4 | 57.6 | 25.7 |
| Max Fusion | | | 67.6 | 79.4 | 57.5 | 25.6 |
| Gated Fusion | ✓ | | **70.8** | 77.5 | 58.6 | 27.6 |
| LSTM Fusion | ✓ | ✓ | 68.0 | 82.3 | 59.6 | 28.3 |
| GRF Fusion | ✓ | ✓ | 66.5 | **85.9** | **60.1** | **31.0** |

| Method (NYUCAD) | GS | MM | prec. | recall | IoU | mIoU |
|---|---|---|---|---|---|---|
| Concatenate Fusion | | | 87.3 | 83.5 | 74.3 | 37.8 |
| Sum Fusion | | | 81.4 | 89.3 | 74.3 | 37.7 |
| Max Fusion | | | 81.9 | 87.8 | 73.3 | 36.5 |
| Gated Fusion | ✓ | | 82.1 | **91.3** | 76.0 | 40.2 |
| LSTM Fusion | ✓ | ✓ | 83.5 | **91.3** | 77.5 | 41.4 |
| GRF Fusion | ✓ | ✓ | **88.4** | 89.1 | **79.7** | **44.8** |

Table 4. Results of different single-stage fusion methods on the NYU and NYUCAD dataset. GS denotes Gate Structure, and MM represents Memory Mechanism. With IoU denotes the accuracy of semantic completion and mIoU denotes the accuracy of semantic scene completion.

malized network outputs $y$:

$$\mathcal{L} = -\sum_{c=0}^{N} w_c \hat{y}_{i,c} \log \left( \frac{e^{y_{ic}}}{\sum_{c'}^{N} e^{y_{ic'}}} \right) \tag{6}$$

where $\hat{y}_{i,c}$ are the one-hot ground truth vectors, *i.e.* $\hat{y}_{i,c} = 1$ if voxel $i$ is labeled by class $c$, otherwise $\hat{y}_{i,c} = 0$. $N$ is the number of classes, and $w_c$ is the loss weight, for balancing different classes, and the setting following SSCNet [35]. To compute the loss function, we ignore all voxels outside the field of view but include all voxels inside the view (empty, non-empty and occluded voxels).

We train the network from scratch with the initial learning rate 0.01 which is reduced by a factor of 0.1 after every ten epochs. We set the weight of empty voxels $w_0$ to 0.05 for data balancing and increase it by 0.05 for every 40 training epochs. Our model is trained using the SGD optimizer with a momentum of 0.9, weight decay of $10^{-4}$ and batch size is 4.

Please note, the order in which the modalities are fed into GRF fusion block has been always fixed (fist Depth, then RGB) for all of the experiments, however, according to our preliminary experiments, using different input order (first RGB, then Depth) has a minor impact on the performance.

## 4. Experiments

### 4.1. Datasets and Metrics

**Datasets** We evaluate the proposed method and compare it with the state-of-the-art methods on NYU [33] and NYU-CAD [11] datasets. NYU consists of 1449 indoor scenes, including 795 training samples and 654 testing samples. The RGBD images of NYU are captured via a Kinect RGBD sensor, and the 3D semantic scene completion labels are from Rock *et al.* ([32]). The annotations are fitted into the scenes by CAD models. NYUCAD uses the depth

maps generated from the projections of the 3D annotations to reduce the misalignment of depths and the annotations.

**Metrics** The primary evaluation metric is the voxel-level intersection over union (IoU) between the predicted labels and ground-truth labels. For semantic scene completion, the IoU is calculated for each category on both the observed and occluded voxels. For scene completion, all non-empty classes are treated as one category, IoU, precision, and recall of the binary predictions are evaluated on the occupied voxels.

### 4.2. Comparisons with the State-of-the-art Methods

In the task of semantic scene completion, our GRFNet outperforms all existing methods and achieves the start-of-the-art accuracy. The results on NYU [33] and NYU-CAD [11] are shown in Table 1 and Table 2, respectively. The GRFNet improves the average IoU over DDR-SSC [25] by 2.5% on both NYU and NYUCAD datasets.

The experiments demonstrate that the proposed fusion block effectively utilizes multi-modality information. Since we focus on presenting a practical multi-modal data fusion approach (GRF), we maintain the consistency of the network structure for a fair comparison to prove that the improvement in accuracy comes from the GRF fusion block. In specific, our network framework is the same as DDR-SSC except for the fusion block.

### 4.3. Quantitative Analysis

Table 1 shows the quantitative results on NYU dataset [33] acquired by our method and other state-of-the-art methods. Approaches of Lin *et al.* ([27]) and Geiger *et al.* ([14]) are traditional methods. SSCNet [35], Essc-Net [43], and DDR-SSC [25] are CNN-based approaches. Compared to the classical approach SSCNet, the IoUs of GRFNet increase 6.1% and 8.2% for SC and SSC tasks, respectively. In the SSC task, our GRFNet gets 6.2% higher accuracy than EsscNet and achieves higher IoU in almost every category. SSCNet and EsscNet only use depth information, while DDR-SSC uses a multi-stage fusion structure to take advantage of the depth and RGB images. Our GRFNet also uses RGB-D information and achieves a 2.5% higher average IoU than DDR-SSC. Regarding the individual class accuracy, the IoUs for each category are also listed out in Table 1.

As shown in Table 2, GRFNet achieves outstanding performance on NYUCAD dataset as well. Specifically, compared to Zheng *et al.*'s( [44]) and Firman *et al.*'s( [11]) methods, GRFNet significantly improves the accuracy in two metrics. Since SSCNet only employs depth as input, the proposed GRFNet which use RGB and depth information achieves much more accurate results. Although TS3D [13] and DDR-SSC use both RGB and depth information, these methods only adopt simple fusion strategy. On the contrary, GRFNet benefited from the novel fusion block, to

| Method | Prec. | Recall | IoU | mIoU |
|---|---|---|---|---|
| Sum Fusion(DDR-SSC) | **71.5** | 80.8 | 61.0 | 30.4 |
| LSTM Fusion | 68.0 | 83.2 | 60.2 | 30.2 |
| GRF Fusion | 68.4 | **85.4** | **61.2** | **32.9** |

Table 5. Results of different multi-stage fusion methods on NYU dataset. With IoU represents the accuracy of scene completion, and mIoU denotes the accuracy of semantic scene completion.

| Fusion stages | Params [k] | FLOPs [G] |
|---|---|---|
| 1 | 794.59 | 193.47 |
| 2 | 803.39 | 366.65 |
| 3 | 812.19 | 539.84 |
| 4 | 820.99 | 713.02 |

Table 6. Params and FLOPs of multi-stage GRFNets with different number of fusion stages.

obtain 0.7% and 2.5% improvements compared to DDR-SSC, and 5.9% and 3.2% improvements compared to TS3D for SC and SSC tasks respectively. In summary, our approach achieves higher accuracy on most indicators than previous methods, especially the average IoU, which reflects the overall performance.

### 4.4. Qualitative Analysis

Figure 5 visualizes results of the semantic scene completion generated by the proposed GRFNet (c), DDR-SSC (d) and SSCNet (e). We mark the difference in visual quality with a red dotted box for reference. As can be seen, compared with both SSCNet and DDR-SSC, the scene completion results of our GRFNet are much more abundant in detail and less error-prone. More visualization results and analyses are provided in the supplemental materials.

### 4.5. Ablation Study

To study the effect of different components and design choices we perform an ablation study. We choose DDR-SSC [25] as the baseline, which is the most relevant work to the proposed GRFNet. Since our focus is the fusion strategy, the GRF module will be analyzed in detail below.

**Single-stage Fusion** To verify the effectiveness of our GRF fusion module, we compare the single-stage GRFNet with a variety of conventional fusion methods, including Concatenate Fusion, Sum Fusion, Max Fusion, and Gated Fusion. For better comparison, we replace the fusion block in the framework (as shown in Figure 3) by the compared single-stage fusion methods. The results of the comparison are shown in Table 4. We have two findings as following: 1) The fusion strategy using the gate structure is better than the one without the gate structure; 2) The memory mechanism can further enhance fusion effects.

As shown in Table 4, Sum Fusion, Max Fusion, and Concatenate Fusion achieve similar performance. And they are significantly lower than the other three modules in which contain adaptive selection mechanism. LSTM fusion and GRF Fusion have a memory mechanism, but Gated Fusion does not; therefore, the accuracy of the first two methods is better. LSTM is more complex and has more parameters than GRF. However, GRF Fusion is still 2.7% and 3.4% more accurate than LSTM on NYU and NYUCAD regarding to SSC accuracy, respectively.

**Multi-stage Fusion**

1). **Multi-stage Strategy** In Table 3, we compare the performance of single-stage GRFNet and the multi-stage GRFNet. Single stage-fusion can only fuse information at one of the network stages. While multi-stage GRFNet can use both the low-level and high-level information and get higher accuracy than the single-stage version on both datasets.

2). **Fusion Strategy** In DDR-SSC [25], sum fusion is used to fuse the four stages of features separately. The fusion results are cascaded and handed over to subsequent networks for semantic label prediction. GRF module employs the recurrent structure that uses the previous fusion results as the input for the next fusion stage, hence the information for each stage can be combined without additional cascading operations. As can been in Table 5, multi-stage GRFNet gets 0.2% higher average IoU than DDR-SSC on SC and 2.5% higher on SSC. And GRF fusion is 1% higher than LSTM fusion on SC and 2.7% higher on SSC.

**Parameters and Flops of Different Fusion Stages**

In Figure 6, parameters and FLOPs of our network with different fusion stages are listed out. As can be seen, with the increasing of fusion stages, parameters increase slightly, which mainly due to the reuse of GRF fusion block, the only source for more parameters is the new added DDR block (for both Depth and RGB channel). On the contrast, FLOPs increase dramatically, which mainly come from the GRF fusion block and small portion come from DDR blocks. In our implementation, GRF fusion block still employ 3D convolutions, thus bring in relatively high compution costs. As we point out before, our focus of this work is to provide a new strategy for fusing the two-modal data in SSC, which can be improved by light-weight operations.

## 5. Conclusion

In this paper, we propose GRFNet with a novel gated recurrent fusion module to fuse RGB and depth information. Different from the existing fusion strategies, we emphasize the importance of the adaptive selectivity of information and the memory mechanism within the fusion block. Moreover, we further extend the single-stage GRFNet to a multi-stage version, which can fuse both low-level and high-level feature at different stages. Our approach has significant advantages over previous methods in multi-modal data fusion and achieves the state-of-the-art performance in semantic scene completion. Extensive comparison exper-

iments and ablation studies verify the effectiveness of the proposed method. In the future, one of our research interests would be to consider making the proposed GRFNet light-weight, for instance, replacing the 3D convolution of GRF fusion block with DDR.

## References

[1] Nicolas Ballas, Li Yao, Chris Pal, and Aaron Courville. Delving deeper into convolutional networks for learning video representations. *arXiv:1511.06432*, 2015.

[2] Long Chen, Karl Francis, and Wen Tang. Semantic augmented reality environment with material-aware physical interactions. In *ISMAR*, pages 135–136. IEEE, 2017.

[3] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *CVPR*, pages 1907–1915, 2017.

[4] Yanhua Cheng, Rui Cai, Zhiwei Li, Xin Zhao, and Kaiqi Huang. Locality-sensitive deconvolution networks with gated fusion for rgb-d indoor semantic segmentation. In *CVPR*, pages 3029–3037, 2017.

[5] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv:1406.1078*, 2014.

[6] Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. Attention-based models for speech recognition. In *NeurIPS*, pages 577–585, 2015.

[7] Camille Couprie, Clément Farabet, Laurent Najman, and Yann LeCun. Indoor semantic segmentation using depth information. In *ICLR*, 2013.

[8] Angela Dai, Daniel Ritchie, Martin Bokeloh, Scott Reed, Jürgen Sturm, and Matthias Nießner. Scancomplete: Large-scale scene completion and semantic segmentation for 3d scans. In *CVPR*, pages 4578–4587, 2018.

[9] Anh-Dzung Doan, Yasir Latif, Tat-Jun Chin, Yu Liu, Thanh-Toan Do, and Ian Reid. Scalable place recognition under appearance change for autonomous driving. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9319–9328, 2019.

[10] Ekrem Emre Yurdakul and Yucel Yemez. Semantic segmentation of rgbd videos with recurrent fully convolutional neural networks. In *ICCV*, pages 367–374, 2017.

[11] Michael Firman, Oisin Mac Aodha, Simon Julier, and Gabriel J Brostow. Structured prediction of unobserved voxels from a single depth image. In *CVPR*, pages 5431–5440, 2016.

[12] Huazhu Fu, Dong Xu, and Stephen Lin. Object-based multiple foreground segmentation in rgbd video. *TIP*, 26(3):1418–1427, 2017.

[13] Martin Garbade, Johann Sawatzky, Alexander Richard, and Juergen Gall. Two stream 3d semantic scene completion. *arXiv:1804.03550*, 2018.

[14] Andreas Geiger and Chaohui Wang. Joint 3d object and layout inference from a single rgb-d image. In *GCPR*, pages 183–195, 2015.

[15] Yanrong Guo and Tao Chen. Semantic segmentation of rgbd images based on deep depth regression. *Pattern Recognition Letters*, 109:55–64, 2018.

[16] Saurabh Gupta, Ross Girshick, Pablo Arbeláez, and Jitendra Malik. Learning rich features from rgb-d images for object detection and segmentation. In *ECCV*, pages 345–360. Springer, 2014.

[17] Xiaoguang Han, Zhen Li, Haibin Huang, Evangelos Kalogerakis, and Yizhou Yu. High-resolution shape completion using deep neural networks for global structure and local geometry inference. In *ICCV*, pages 85–93, 2017.

[18] Caner Hazirbas, Lingni Ma, Csaba Domokos, and Daniel Cremers. Fusenet: Incorporating depth into semantic segmentation via fusion-based cnn architecture. In *ACCV*, pages 213–228. Springer, 2016.

[19] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[20] Jian-Fang Hu, Wei-Shi Zheng, Jiahui Pan, Jianhuang Lai, and Jianguo Zhang. Deep bilinear learning for rgb-d action recognition. In *ECCV*, pages 335–351, 2018.

[21] Earnest Paul Ijjina and Krishna Mohan Chalavadi. Human action recognition in rgb-d videos using motion sequence information and deep learning. *Pattern Recognition*, 72:504–516, 2017.

[22] Le Kang, Peng Ye, Yi Li, and David Doermann. Convolutional neural networks for no-reference image quality assessment. In *CVPR*, pages 1733–1740, 2014.

[23] Christian Kerl, Jürgen Sturm, and Daniel Cremers. Dense visual slam for rgb-d cameras. In *IROS*, pages 2100–2106. IEEE, 2013.

[24] Yoon Kim, Yacine Jernite, David Sontag, and Alexander M Rush. Character-aware neural language models. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

[25] Jie Li, Yu Liu, Dong Gong, Qinfeng Shi, Xia Yuan, Chunxia Zhao, and Ian Reid. Rgbd based dimensional decomposition residual network for 3d semantic scene completion. In *CVPR*, pages 7693–7702, 2019.

[26] Zhen Li, Yukang Gan, Xiaodan Liang, Yizhou Yu, Hui Cheng, and Liang Lin. Lstm-cf: Unifying context modeling and fusion with lstms for rgb-d scene labeling. In *ECCV*, pages 541–557. Springer, 2016.

[27] Dahua Lin, Sanja Fidler, and Raquel Urtasun. Holistic scene understanding for 3d object detection with rgbd cameras. In *ICCV*, pages 1417–1424, 2013.

[28] Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear cnn models for fine-grained visual recognition. In *ICCV*, pages 1449–1457, 2015.

[29] Yan Lu and Dezhen Song. Robust rgb-d odometry using point and line features. In *ICCV*, pages 3934–3942, 2015.

[30] Seong-Jin Park, Ki-Sang Hong, and Seungyong Lee. Rdfnet: Rgb-d multi-level residual feature fusion for indoor semantic segmentation. In *ICCV*, pages 4980–4989, 2017.

[31] Xiaofeng Ren, Liefeng Bo, and Dieter Fox. Rgb-(d) scene labeling: Features and algorithms. In *CVPR*, pages 2759–2766. IEEE, 2012.

[32] Jason Rock, Tanmay Gupta, Justin Thorsen, JunYoung Gwak, Daeyun Shin, and Derek Hoiem. Completing 3d ob-

ject shape from one depth image. In *CVPR*, pages 2484–2493. IEEE, 2015.

[33] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, pages 746–760. Springer, 2012.

[34] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NeurIPS*, pages 568–576, 2014.

[35] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *CVPR*, pages 190–198, 2017.

[36] Maryam Sultana, Arif Mahmood, Sajid Javed, and Soon Ki Jung. Unsupervised rgbd video object segmentation using gans. *arXiv:1811.01526*, 2018.

[37] Jacob Varley, Chad DeChant, Adam Richardson, Joaquín Ruales, and Peter Allen. Shape completion enabled robotic grasping. In *IROS*, pages 2442–2447, 2017.

[38] Jinghua Wang, Zhenhua Wang, Dacheng Tao, Simon See, and Gang Wang. Learning common and specific features for rgb-d semantic segmentation with deconvolutional networks. In *ECCV*, pages 664–679. Springer, 2016.

[39] Ningning Wang and Xiaojin Gong. Adaptive fusion for rgb-d salient object detection. *arXiv:1901.01369*, 2019.

[40] Thomas Whelan, Hordur Johannsson, Michael Kaess, John J Leonard, and John McDonald. Robust real-time visual odometry for dense rgb-d mapping. In *ICRA*, 2013.

[41] Joe Yue-Hei Ng, Matthew Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat Monga, and George Toderici. Beyond short snippets: Deep networks for video classification. In *CVPR*, pages 4694–4702, 2015.

[42] Haokui Zhang, Ying Li, Peng Wang, Yu Liu, and Chunhua Shen. Rgb-d based action recognition with light-weight 3d convolutional networks. *arXiv preprint arXiv:1811.09908*, 2018.

[43] Jiahui Zhang, Hao Zhao, Anbang YaoE, Yurong Chen, Li Zhang, and Hongen LiaoE. Efficient semantic scene completion network with spatial group convolution. In *ECCV*, pages 733–749, 2018.

[44] Bo Zheng, Yibiao Zhao, Joey C Yu, Katsushi Ikeuchi, and Song-Chun Zhu. Beyond point clouds: Scene understanding by reasoning geometry and physics. In *CVPR*, pages 3127–3134, 2013.

## 6. More Details of GRFNet

### 6.1. Detailed Architectures

The details of the proposed network structure are shown in Table 7. PWConv represents the point-wise convolution, and it is used to adjust the number of channels of the feature map. The down-sample layer in our network is composed of a max-pooling layer and a convolution layer with stride set as 2. The outputs of the two layers are concatenated before fed into the subsequent layers.

### 6.2. Dimensional Decomposition Residual Block

In Table 9, we show the details of the Dimensional Decomposition Residual (DDR) [25] block. DDR $(k, w, s, d)$ denotes the DDR block with the kernel size $k$, the output channels of feature maps $w$, the stride $s$ and the dilation rate $d$. DDRConv represents the proposed DDR convolution within a DDR block.

The most significant advantage of using DDR in 3D tasks is the reduction of the number of parameters and the amount of calculation. In a DDR block, the 3D convolution with the kernel size $k \times k \times k$ is decomposed into three consecutive layers with filter size $1 \times 1 \times k$, $1 \times k \times 1$ and $k \times 1 \times 1$. The most common value for $k$ is 3. The computational costs of the original block and the DDR block are proportional to $c^{in} \times c^{out} \times k \times k \times k$ and $c^{in} \times c^{out} \times (k+k+k)$, where $c^{in}$ and $c^{out}$ are the numbers of input and output channels.

We assume $c^{in} = c^{out} = w$ and ignore bias, then the parameter quantity changes from $w^2 \times k^3$ to $w^2 \times 3k$. The bottleneck structure within the DDR block further reduces its parameter amount and calculation cost.

### 6.3. 2D to 3D Projection

Each point in depth can be projected to a position in the 3D space. We voxelize this entire 3D space with meshed grids to obtain a 3D volume. In the projection layer, every feature tensor is projected into the 3D volume at the location corresponding to its position in depth. With the feature projection layer, the 2D feature maps extracted by the 2D CNN are converted to a view-independent 3D feature volume.

## 7. Recurrent Property of GRFNet

In Figure 6 and Figure 7, network structures with different fusion strategies are listed out. Specifically, Figure 6 is the network structure with GRF fusion block, which means different stages will share the same GRF fusion block, and it has the 'recurrent' property between different stages, that is the low-level fusion could be part of guidance and contribute the following the high-level stages. On the contrary, Figure 7 given the network workflow which use other fusion blocks that does not have the 'recurrent' property, including *Sum Fusion*, *Average Fusion*, *Bilinear Fusion*, *Concatenation Fusion*, *Max Fusion*, *Transformation Fusion*. And as denoted with the blue arrow, different fusion blocks are separated with each other, until they are concated together to be feded to the light-weight ASPP module.

As can be seen in Table 8, results of networks with different fusion strategies are listed out. The first two rows are the results of networks which employ the MaxFusion and SumFusion, and last row are the results of GRFNet. For both of the scene completion and semantic scene completion tasks, GRFNet boost the performance significantly. We

| Module | Operation | Output Size<br>2D: $Height \times Width \times Channels$<br>3D: $Depth \times Height \times Width \times Channels$ | Kernel | Stride | Dilation |
|---|---|---|---|---|---|
| Feature Extractor | PWConv | $640 \times 480 \times 8$ | 1 | 1 | 1 |
| | 2D DDR | $640 \times 480 \times 8$ | 3 | 1 | 1 |
| | 2D DDR | $640 \times 480 \times 8$ | 3 | 1 | 1 |
| | 2D - 3D Projection | $240 \times 144 \times 240 \times 8$ | - | - | - |
| | Down-sample | $120 \times 72 \times 120 \times 16$ | 3 | 2 | 1 |
| | 3D DDR | $120 \times 72 \times 120 \times 16$ | 3 | 1 | 1 |
| | Down-sample | $60 \times 36 \times 60 \times 64$ | 3 | 2 | 1 |
| | 3D DDR | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 1 |
| Feature Fusion | GRF stage 1 | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 1 |
| | 3D DDR | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 2 |
| | GRF stage 2 | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 1 |
| | 3D DDR | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 3 |
| | GRF stage 3 | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 1 |
| | 3D DDR | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 5 |
| | GRF stage 4 | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 1 |
| LW-ASPP | PWConv | $60 \times 36 \times 60 \times 64$ | 1 | 1 | 1 |
| | 3D DDR | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 3 |
| | 3D DDR | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 6 |
| | 3D DDR | $60 \times 36 \times 60 \times 64$ | 3 | 1 | 9 |
| | GlobalAvgPool | $60 \times 36 \times 60 \times 64$ | - | - | - |
| | Concatenate | $60 \times 36 \times 60 \times 320$ | - | - | - |
| Output | PWConv | $60 \times 36 \times 60 \times 160$ | 1 | 1 | 1 |
| | PWConv | $60 \times 36 \times 60 \times 12$ | 1 | 1 | 1 |
| | ArgMax | $60 \times 36 \times 60 \times 12$ | - | - | - |

Table 7. The details of the proposed (GRFNet) network architecture. Including module name, layer operation, output size, kernel size, stride and dilation.

| | scene completion | | | semantic scene completion | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| method | prec. | recall | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| multi-stage MaxFusion | 82.5 | **91.8** | 76.8 | 48.2 | 91.4 | 53.3 | 16.5 | 36.5 | 54.6 | 50.5 | 29.3 | **11.8** | 41.2 | 25.0 | 41.7 |
| multi-stage SumFusion | **88.7** | 88.5 | 79.4 | **54.1** | 91.5 | 56.4 | 14.9 | 37.0 | 55.7 | 51.0 | 28.8 | 9.2 | 44.1 | 27.8 | 42.8 |
| multi-stage GRFNet | 87.2 | 91.0 | **80.1** | 50.3 | **91.8** | **58.1** | **18.4** | **42.7** | **60.6** | **52.8** | **34.6** | 11.5 | **46.6** | **30.8** | **45.3** |

Table 8. Results of multi-stage networks with different fusion blocks on NYUCAD dataset. As can be seen, compared with the results based MaxFusion and SumFusion blocks, GRFNet (with GRF fusion block) achieves better accuracy both in scene completion and semantic scene completion tasks. Which may due to the advantage of 'recurrent' property of the proposed GRF fusion block.
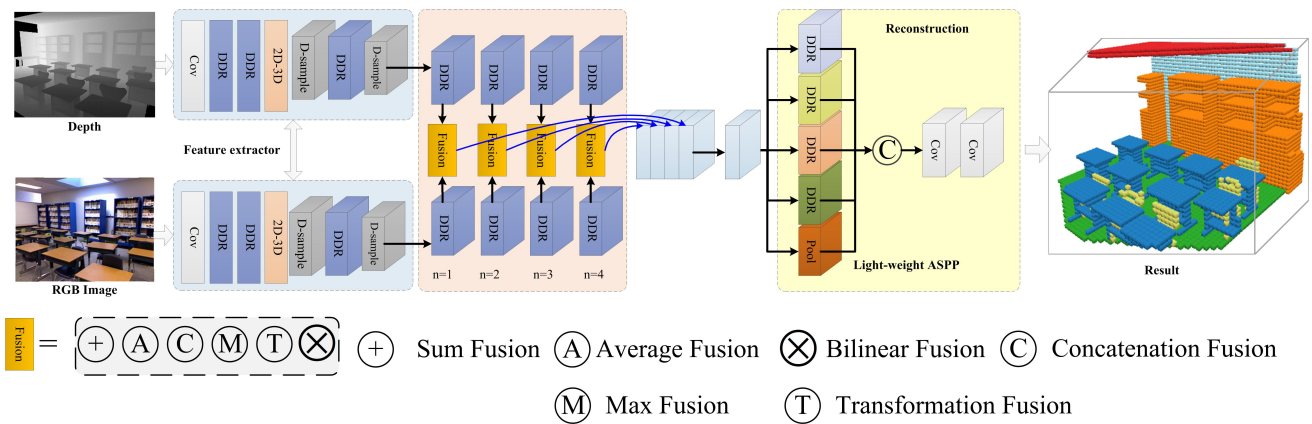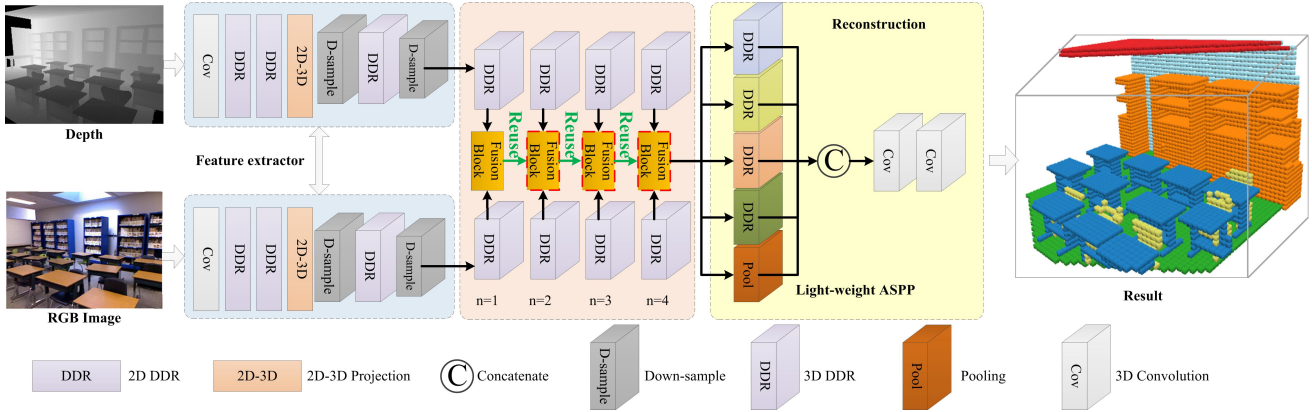
| Operation | Kernel | Channels | Stride | Dilation |
|---|---|---|---|---|
| PWConv | $1 \times 1 \times 1$ | $w/4$ | 1 | 1 |
| DDRConv | $1 \times 1 \times k$ | $w/4$ | $s$ | $d$ |
| DDRConv | $1 \times k \times 1$ | $w/4$ | $s$ | $d$ |
| DDRConv | $k \times 1 \times 1$ | $w/4$ | $s$ | $d$ |
| PWConv | $1 \times 1 \times 1$ | $w$ | 1 | 1 |

Table 9. Details of the DDR $(k, w, s, d)$ block. $k$ is the kernel size, $w$ is the output channels of the feature map, $s$ is the stride and $d$ is the dilation rate of the convolution.

supect that is because the 'recurrent' property of the proposed GRF fusion block.

# 8. More Qualitative Results

Figure 8 shows some visualized results on NYU-CAD [11] dataset. As shown in Figure 8, the proposed GRFNet achieves better results than DDR-SSC [25] and SS-CNet [35], and is much more accurate in shape completion and semantic segmentation. The color information is beneficial for the prediction of semantic labeling. In Figure 8, the prediction of furniture in the second, third, and fifth

Figure 6. Network workflow with GRF fusion block



Figure 7. Network workflow with other fusion blocks

rows is more accurate than the method using only depth.

When an object consists of several parts with various appearances, the RGB information may result in inconsistencies in the semantics of the local areas. For example, in the first row of Figure 8, most of the walls are white, while the right part of the wall is brown. This makes it difficult for the network to predict the semantics of that brown wall. In this case, the geometric information contained in the depth image can effectively eliminate ambiguity and provide a reasonable inference.

Besides, when different objects are in similar colors, depth can provide adequate information for distinguishing objects. In the fourth row of Figure 8, the chair is very similar to the background in the RGB image, but it can be easily distinguished in depth. Therefore, in general, the fusion of RGB-D is critical, and our GRFNet can effectively improve the accuracy of SSC.

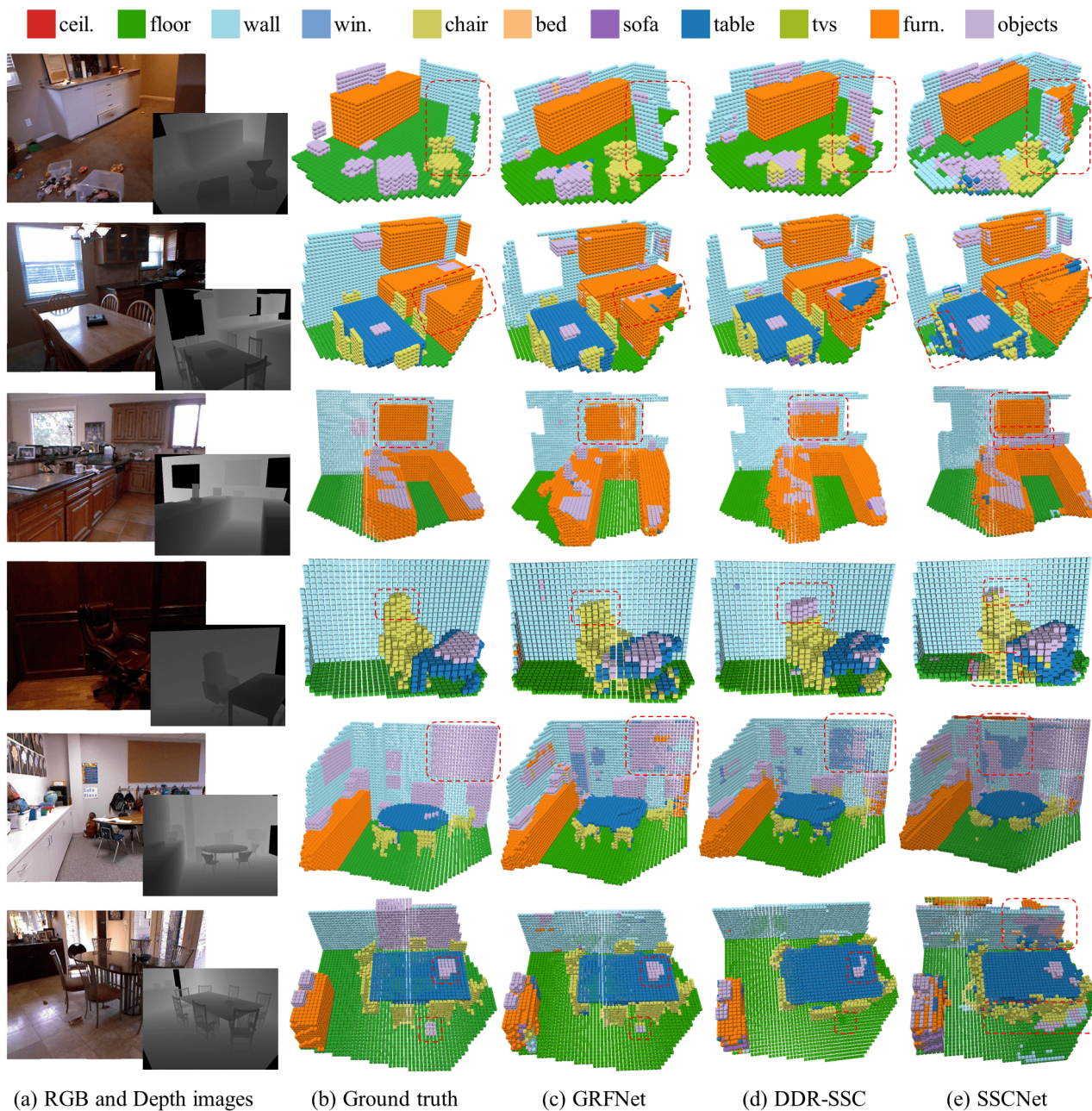(a) RGB and Depth images     (b) Ground truth     (c) GRFNet     (d) DDR-SSC     (e) SSCNet

Figure 8. Qualitative results on NYUCAD. From left to right: Input RGB-D image, ground truth, results obtained by our proposed GRFNet, results obtained by DDR-SSC [25] and SSCNet [35]. Overall, our completed semantic 3D scenes are less cluttered and show a higher voxel-wise accuracy compared to DDR-SSC and SSCNet.