

# Adaptive Visual Scene Understanding: Incremental Scene Graph Generation

Naitik Khandelwal<sup>1,2</sup>, Xiao Liu<sup>1,2</sup>, Mengmi Zhang<sup>1,2,\*</sup>

<sup>1</sup>College of Computing and Data Science, Nanyang Technological University, Singapore

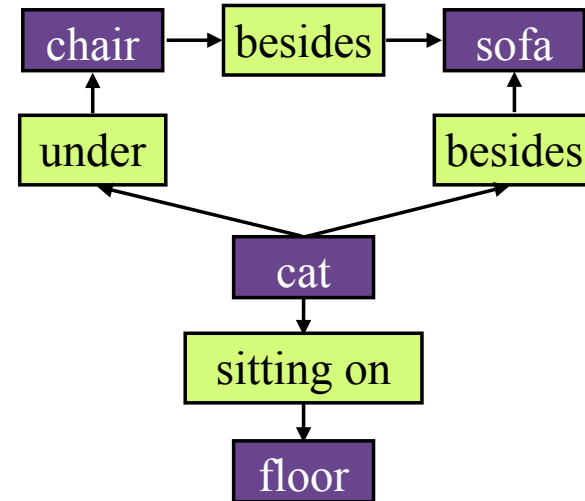
<sup>2</sup>Deep NeuroCognition Lab, I2R and CFAR, Agency for Science, Technology and Research, Singapore



arxiv

github

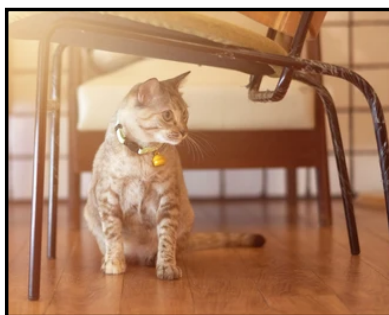
# Introduction



Scene Graph Generation (SGG)

# Introduction

**Home & Hotel**



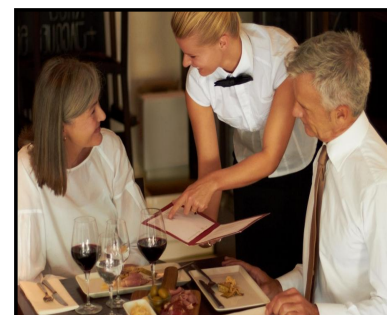
<cat, under, chair>  
<chair, on, floor>  
<sofa, besides, chair>

**Outdoor**



<man, besides, car>  
<car, parked on, street>  
<tree, besides, sidewalk>

**Shopping & Dine**



<waitress, holding, menu>  
<waitress, serving, customers>  
<menu, above, table>

.....



New **objects** and new **relations** emerge in new scenes

Time

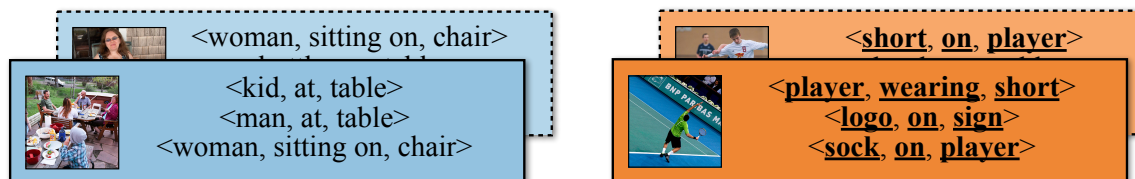
Continual Scene Graph Generation (CSEGG)

# CSEGG Benchmark

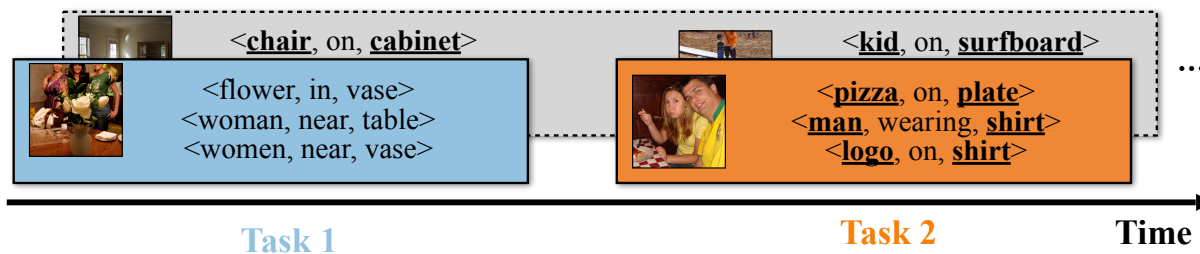
## Scenario 1: Relationship Incremental Learning



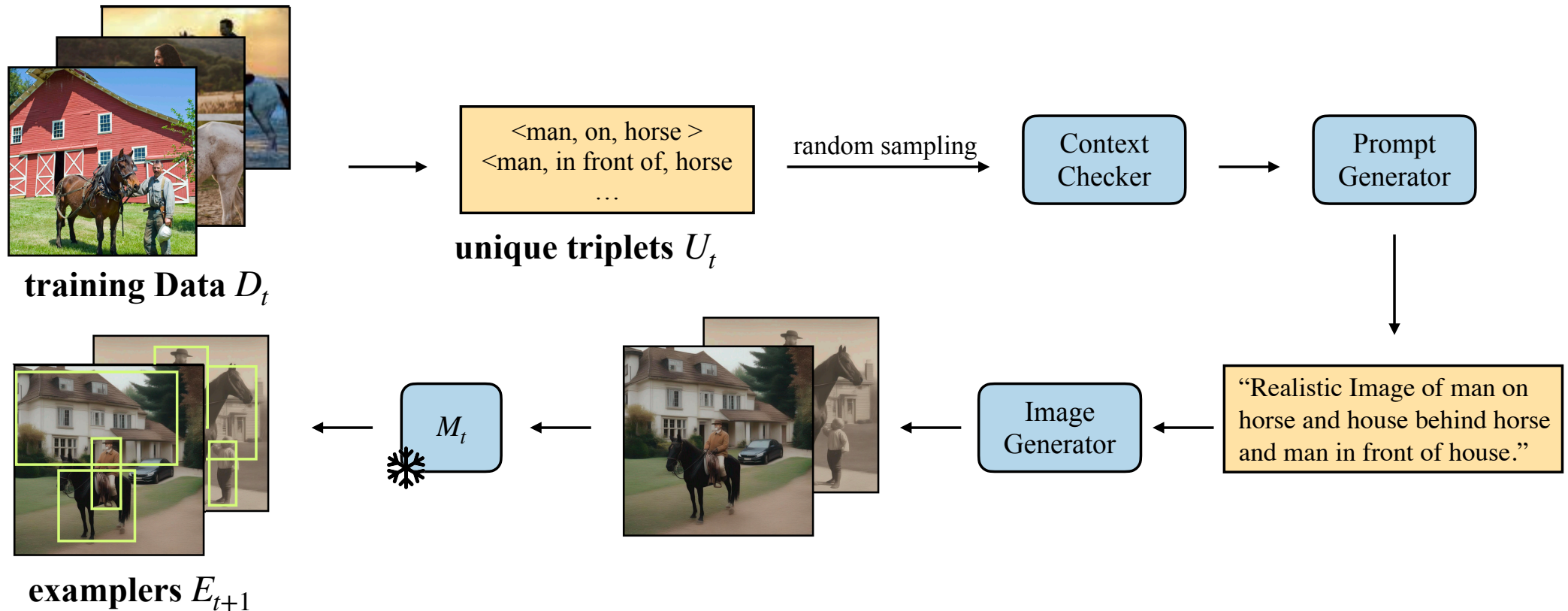
## Scenario 2: Scene Incremental Learning



## Scenario 3: Relationship Generalization



# Replays via Analysis by Synthesis (RAS)



# Main Results

| SGTR[32]     |                          |               |            |              |             |               |                          |               |             |              |              |               |
|--------------|--------------------------|---------------|------------|--------------|-------------|---------------|--------------------------|---------------|-------------|--------------|--------------|---------------|
| Methods      | Learning Scenario 1 (S1) |               |            |              |             |               | Learning Scenario 2 (S2) |               |             |              |              |               |
|              | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑        | BWT↑          | Avg.R↑                   | F↑            | mR↑         | mF↑          | FWT↑         | BWT↑          |
| Joint        | 20.15                    | 0             | 4.6        | 0            | -           | -             | 12.64                    | 0             | 9.84        | 0            | -            | -             |
| Replay@100%  | 16.17                    | -12.24        | 3.32       | -1.34        | -1.77       | -11.72        | 4.56                     | -4.13         | 4.56        | -5.61        | -1.045       | -30.25        |
| Naive        | 1.33                     | -28.7         | 0.86       | -1.74        | -2.03       | -60.67        | 0.51                     | -23.22        | 0.05        | -11.31       | -3.77        | -62.34        |
| EWC[24]      | 1.89                     | -28.4         | 0.96       | -1.72        | -1.17       | -52.45        | 0                        | -23.22        | 0           | -11.31       | -2.65        | -50.12        |
| RAS_GT       | 5.78                     | -26.51        | 1.43       | -1.54        | -1.2        | -44.27        | 0.98                     | -23.11        | 0.76        | -10.86       | -1.6         | -43.25        |
| PackNet[44]  | 7.19                     | -25.67        | 1.35       | -1.64        | -1.03       | -42.35        | 1.67                     | -22.77        | 0.9         | -10.33       | -1.4         | -42.45        |
| Replay@10%   | 8.55                     | -22.21        | 4.33       | -1.44        | <b>4.29</b> | -38.35        | 1.81                     | -20.72        | 1.15        | -9.64        | -0.9         | -40.67        |
| Replay@20%   | 9.25                     | -20.35        | 4.78       | -1.42        | 3.21        | -31.98        | 2.57                     | -17.17        | 1.56        | -8.07        | -0.67        | -38.27        |
| <b>Ours*</b> | <b>10.78</b>             | <b>-18.92</b> | <b>5.6</b> | <b>-1.39</b> | <b>2.3</b>  | <b>-25.56</b> | <b>3.45</b>              | <b>-10.23</b> | <b>2.75</b> | <b>-6.57</b> | <b>-0.54</b> | <b>-35.67</b> |

| TCNN[72]     |                          |               |            |              |             |               |                          |               |            |              |              |               |
|--------------|--------------------------|---------------|------------|--------------|-------------|---------------|--------------------------|---------------|------------|--------------|--------------|---------------|
| Methods      | Learning Scenario 1 (S1) |               |            |              |             |               | Learning Scenario 2 (S2) |               |            |              |              |               |
|              | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑        | BWT↑          | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑         | BWT↑          |
| Joint        | 19.53                    | 0             | 3.9        | 0            | -           | -             | 4.3                      | 0             | 3.7        | 0            | -            | -             |
| Replay@100%  | 13.45                    | -8.83         | 3.6        | -0.35        | -1.5        | -10.45        | 12.45                    | -4.13         | 3.2        | -0.56        | -2.1         | -20.34        |
| Naive        | 0.98                     | -21.2         | 0.74       | -1.35        | -3.45       | -43.87        | 0                        | -18.22        | 0.45       | -2.67        | -4.12        | -53.12        |
| EWC[24]      | 2.36                     | -21.05        | 0.67       | -1.34        | -2.34       | -39.89        | 0                        | -18.22        | 0.03       | 0            | -3.77        | -51.67        |
| PackNet[44]  | 3.2                      | -19.7         | 1.1        | -1.13        | -1.3        | -32.45        | 1.1                      | -17.82        | 0.84       | -1.97        | -2.84        | -40.34        |
| Replay@10%   | 5.67                     | -18.9         | 3.21       | -1.05        | <b>1.45</b> | -28.34        | 1.81                     | -16.72        | 1.03       | -1.74        | -1.4         | -43.56        |
| Replay@20%   | 6.23                     | -17.45        | 3.5        | -1.01        | 1.01        | -24.32        | 2.37                     | -15.17        | 1.45       | -1.53        | -1.1         | -38.56        |
| <b>Ours*</b> | <b>7.8</b>               | <b>-15.67</b> | <b>3.9</b> | <b>-0.95</b> | <b>0.5</b>  | <b>-19.83</b> | <b>4.67</b>              | <b>-11.31</b> | <b>2.2</b> | <b>-0.89</b> | <b>-0.97</b> | <b>-29.65</b> |

RAS outperforms all CSEGG baselines in S1 and S2

# Main Results

| SGTR[32]      |                          |               |             |              |             |               |                          |               |             |               |              |               |
|---------------|--------------------------|---------------|-------------|--------------|-------------|---------------|--------------------------|---------------|-------------|---------------|--------------|---------------|
| Methods       | Learning Scenario 1 (S1) |               |             |              |             |               | Learning Scenario 2 (S2) |               |             |               |              |               |
|               | Avg.R↑                   | F↑            | mR↑         | mF↑          | FWT↑        | BWT↑          | Avg.R↑                   | F↑            | mR↑         | mF↑           | FWT↑         | BWT↑          |
| Joint         | 20.15                    | 0             | 4.6         | 0            | -           | -             | 12.64                    | 0             | 9.84        | 0             | -            | -             |
| Replay@100%   | 16.17                    | -12.24        | 3.32        | -1.34        | -1.77       | -11.72        | 4.56                     | -4.13         | 4.56        | -5.61         | -1.045       | -30.25        |
| Naive         | 1.33                     | -28.7         | 0.86        | -1.74        | -2.03       | -60.67        | 0.51                     | -23.22        | 0.05        | -11.31        | -3.77        | -62.34        |
| EWC[24]       | 1.89                     | -28.4         | 0.96        | -1.72        | -1.17       | -52.45        | 0                        | -23.22        | 0           | -11.31        | -2.65        | -50.12        |
| <b>RAS GT</b> | <b>5.78</b>              | <b>-26.51</b> | <b>1.43</b> | <b>-1.54</b> | <b>-1.2</b> | <b>-44.27</b> | <b>0.98</b>              | <b>-23.11</b> | <b>0.76</b> | <b>-10.86</b> | <b>-1.6</b>  | <b>-43.25</b> |
| PackNet[44]   | 7.19                     | -25.67        | 1.35        | -1.64        | -1.03       | -42.35        | 1.67                     | -22.77        | 0.9         | -10.33        | -1.4         | -42.45        |
| Replay@10%    | 8.55                     | -22.21        | 4.33        | -1.44        | <b>4.29</b> | -38.35        | 1.81                     | -20.72        | 1.15        | -9.64         | -0.9         | -40.67        |
| Replay@20%    | 9.25                     | -20.35        | 4.78        | -1.42        | 3.21        | -31.98        | 2.57                     | -17.17        | 1.56        | -8.07         | -0.67        | -38.27        |
| <b>Ours*</b>  | <b>10.78</b>             | <b>-18.92</b> | <b>5.6</b>  | <b>-1.39</b> | <b>2.3</b>  | <b>-25.56</b> | <b>3.45</b>              | <b>-10.23</b> | <b>2.75</b> | <b>-6.57</b>  | <b>-0.54</b> | <b>-35.67</b> |

| TCNN[72]     |                          |               |            |              |             |               |                          |               |            |              |              |               |
|--------------|--------------------------|---------------|------------|--------------|-------------|---------------|--------------------------|---------------|------------|--------------|--------------|---------------|
| Methods      | Learning Scenario 1 (S1) |               |            |              |             |               | Learning Scenario 2 (S2) |               |            |              |              |               |
|              | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑        | BWT↑          | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑         | BWT↑          |
| Joint        | 19.53                    | 0             | 3.9        | 0            | -           | -             | 4.3                      | 0             | 3.7        | 0            | -            | -             |
| Replay@100%  | 13.45                    | -8.83         | 3.6        | -0.35        | -1.5        | -10.45        | 12.45                    | -4.13         | 3.2        | -0.56        | -2.1         | -20.34        |
| Naive        | 0.98                     | -21.2         | 0.74       | -1.35        | -3.45       | -43.87        | 0                        | -18.22        | 0.45       | -2.67        | -4.12        | -53.12        |
| EWC[24]      | 2.36                     | -21.05        | 0.67       | -1.34        | -2.34       | -39.89        | 0                        | -18.22        | 0.03       | 0            | -3.77        | -51.67        |
| PackNet[44]  | 3.2                      | -19.7         | 1.1        | -1.13        | -1.3        | -32.45        | 1.1                      | -17.82        | 0.84       | -1.97        | -2.84        | -40.34        |
| Replay@10%   | 5.67                     | -18.9         | 3.21       | -1.05        | <b>1.45</b> | -28.34        | 1.81                     | -16.72        | 1.03       | -1.74        | -1.4         | -43.56        |
| Replay@20%   | 6.23                     | -17.45        | 3.5        | -1.01        | 1.01        | -24.32        | 2.37                     | -15.17        | 1.45       | -1.53        | -1.1         | -38.56        |
| <b>Ours*</b> | <b>7.8</b>               | <b>-15.67</b> | <b>3.9</b> | <b>-0.95</b> | 0.5         | <b>-19.83</b> | <b>4.67</b>              | <b>-11.31</b> | <b>2.2</b> | <b>-0.89</b> | <b>-0.97</b> | <b>-29.65</b> |

**Decomposing** scene graphs into **smaller, diverse** components with clear prompts is **more effective** than directly storing and using ground truth scene graphs for image generation

# Main Results

| SGTR[32]     |                          |               |            |              |             |               |                          |               |             |              |              |               |
|--------------|--------------------------|---------------|------------|--------------|-------------|---------------|--------------------------|---------------|-------------|--------------|--------------|---------------|
| Methods      | Learning Scenario 1 (S1) |               |            |              |             |               | Learning Scenario 2 (S2) |               |             |              |              |               |
|              | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑        | BWT↑          | Avg.R↑                   | F↑            | mR↑         | mF↑          | FWT↑         | BWT↑          |
| Joint        | 20.15                    | 0             | 4.6        | 0            | -           | -             | 12.64                    | 0             | 9.84        | 0            | -            | -             |
| Replay@100%  | 16.17                    | -12.24        | 3.32       | -1.34        | -1.77       | -11.72        | 4.56                     | -4.13         | 4.56        | -5.61        | -1.045       | -30.25        |
| Naive        | 1.33                     | -28.7         | 0.86       | -1.74        | -2.03       | -60.67        | 0.51                     | -23.22        | 0.05        | -11.31       | -3.77        | -62.34        |
| EWC[24]      | 1.89                     | -28.4         | 0.96       | -1.72        | -1.17       | -52.45        | 0                        | -23.22        | 0           | -11.31       | -2.65        | -50.12        |
| RAS_GT       | 5.78                     | -26.51        | 1.43       | -1.54        | -1.2        | -44.27        | 0.98                     | -23.11        | 0.76        | -10.86       | -1.6         | -43.25        |
| PackNet[44]  | 7.19                     | -25.67        | 1.35       | -1.64        | -1.03       | -42.35        | 1.67                     | -22.77        | 0.9         | -10.33       | -1.4         | -42.45        |
| Replay@10%   | 8.55                     | -22.21        | 4.33       | -1.44        | <b>4.29</b> | -38.35        | 1.81                     | -20.72        | 1.15        | -9.64        | -0.9         | -40.67        |
| Replay@20%   | 9.25                     | -20.35        | 4.78       | -1.42        | 3.21        | -31.98        | 2.57                     | -17.17        | 1.56        | -8.07        | -0.67        | -38.27        |
| <b>Ours*</b> | <b>10.78</b>             | <b>-18.92</b> | <b>5.6</b> | <b>-1.39</b> | 2.3         | <b>-25.56</b> | <b>3.45</b>              | <b>-10.23</b> | <b>2.75</b> | <b>-6.57</b> | <b>-0.54</b> | <b>-35.67</b> |

| TCNN[72]     |                          |               |            |              |             |               |                          |               |            |              |              |               |
|--------------|--------------------------|---------------|------------|--------------|-------------|---------------|--------------------------|---------------|------------|--------------|--------------|---------------|
| Methods      | Learning Scenario 1 (S1) |               |            |              |             |               | Learning Scenario 2 (S2) |               |            |              |              |               |
|              | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑        | BWT↑          | Avg.R↑                   | F↑            | mR↑        | mF↑          | FWT↑         | BWT↑          |
| Joint        | 19.53                    | 0             | 3.9        | 0            | -           | -             | 4.3                      | 0             | 3.7        | 0            | -            | -             |
| Replay@100%  | 13.45                    | -8.83         | 3.6        | -0.35        | -1.5        | -10.45        | 12.45                    | -4.13         | 3.2        | -0.56        | -2.1         | -20.34        |
| Naive        | 0.98                     | -21.2         | 0.74       | -1.35        | -3.45       | -43.87        | 0                        | -18.22        | 0.45       | -2.67        | -4.12        | -53.12        |
| EWC[24]      | 2.36                     | -21.05        | 0.67       | -1.34        | -2.34       | -39.89        | 0                        | -18.22        | 0.03       | 0            | -3.77        | -51.67        |
| PackNet[44]  | 3.2                      | -19.7         | 1.1        | -1.13        | -1.3        | -32.45        | 1.1                      | -17.82        | 0.84       | -1.97        | -2.84        | -40.34        |
| Replay@10%   | 5.67                     | -18.9         | 3.21       | -1.05        | <b>1.45</b> | -28.34        | 1.81                     | -16.72        | 1.03       | -1.74        | -1.4         | -43.56        |
| Replay@20%   | 6.23                     | -17.45        | 3.5        | -1.01        | 1.01        | -24.32        | 2.37                     | -15.17        | 1.45       | -1.53        | -1.1         | -38.56        |
| <b>Ours*</b> | <b>7.8</b>               | <b>-15.67</b> | <b>3.9</b> | <b>-0.95</b> | 0.5         | <b>-19.83</b> | <b>4.67</b>              | <b>-11.31</b> | <b>2.2</b> | <b>-0.89</b> | <b>-0.97</b> | <b>-29.65</b> |

RAS outperforms Replay@20% (~2 Gb) while requiring **less storage** (~1.2 Gb)



# Main Results

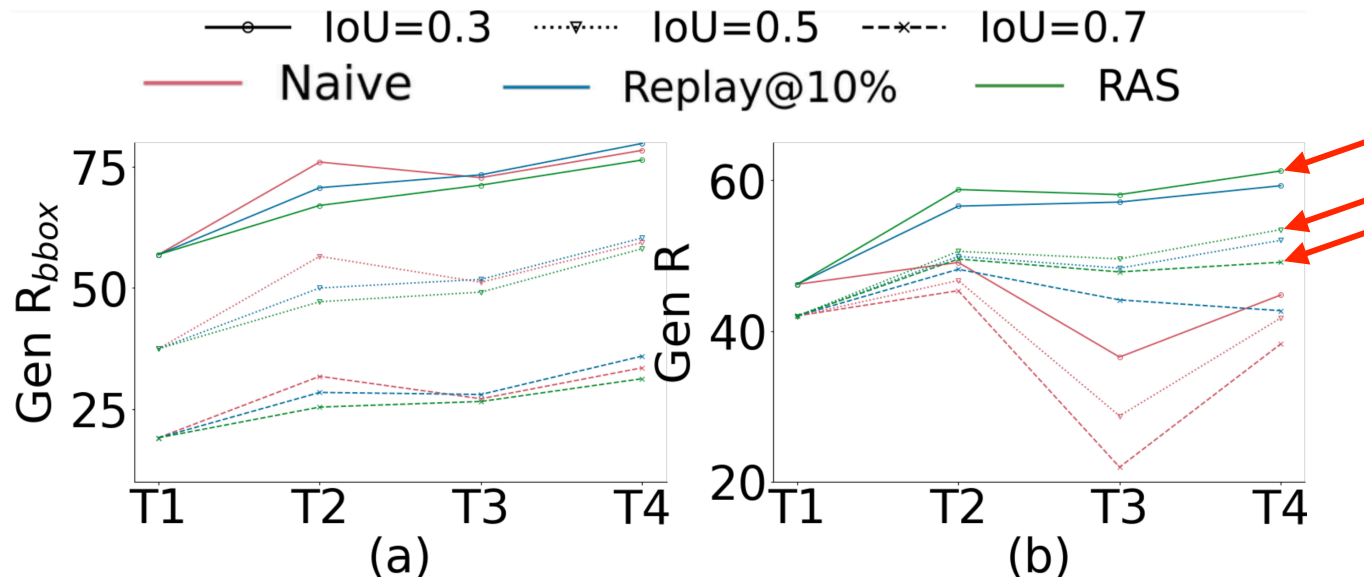
| SGTR[32]     |                          |               |               |               |                |                |                          |               |               |               |                |                |
|--------------|--------------------------|---------------|---------------|---------------|----------------|----------------|--------------------------|---------------|---------------|---------------|----------------|----------------|
| Methods      | Learning Scenario 1 (S1) |               |               |               |                |                | Learning Scenario 2 (S2) |               |               |               |                |                |
|              | Avg.R $\uparrow$         | F $\uparrow$  | mR $\uparrow$ | mF $\uparrow$ | FWT $\uparrow$ | BWT $\uparrow$ | Avg.R $\uparrow$         | F $\uparrow$  | mR $\uparrow$ | mF $\uparrow$ | FWT $\uparrow$ | BWT $\uparrow$ |
| Joint        | 20.15                    | 0             | 4.6           | 0             | -              | -              | 12.64                    | 0             | 9.84          | 0             | -              | -              |
| Replay@100%  | 16.17                    | -12.24        | 3.32          | -1.34         | -1.77          | -11.72         | 4.56                     | -4.13         | 4.56          | -5.61         | -1.045         | -30.25         |
| Naive        | 1.33                     | -28.7         | 0.86          | -1.74         | -2.03          | -60.67         | 0.51                     | -23.22        | 0.05          | -11.31        | -3.77          | -62.34         |
| EWC[24]      | 1.89                     | -28.4         | 0.96          | -1.72         | -1.17          | -52.45         | 0                        | -23.22        | 0             | -11.31        | -2.65          | -50.12         |
| RAS_GT       | 5.78                     | -26.51        | 1.43          | -1.54         | -1.2           | -44.27         | 0.98                     | -23.11        | 0.76          | -10.86        | -1.6           | -43.25         |
| PackNet[44]  | 7.19                     | -25.67        | 1.35          | -1.64         | -1.03          | -42.35         | 1.67                     | -22.77        | 0.9           | -10.33        | -1.4           | -42.45         |
| Replay@10%   | 8.55                     | -22.21        | 4.33          | -1.44         | <b>4.29</b>    | -38.35         | 1.81                     | -20.72        | 1.15          | -9.64         | -0.9           | -40.67         |
| Replay@20%   | 9.25                     | -20.35        | 4.78          | -1.42         | 3.21           | -31.98         | 2.57                     | -17.17        | 1.56          | -8.07         | -0.67          | -38.27         |
| <b>Ours*</b> | <b>10.78</b>             | <b>-18.92</b> | <b>5.6</b>    | <b>-1.39</b>  | 2.3            | <b>-25.56</b>  | <b>3.45</b>              | <b>-10.23</b> | <b>2.75</b>   | <b>-6.57</b>  | <b>-0.54</b>   | <b>-35.67</b>  |

| TCNN[72]     |                          |               |               |               |                |                |                          |               |               |               |                |                |
|--------------|--------------------------|---------------|---------------|---------------|----------------|----------------|--------------------------|---------------|---------------|---------------|----------------|----------------|
| Methods      | Learning Scenario 1 (S1) |               |               |               |                |                | Learning Scenario 2 (S2) |               |               |               |                |                |
|              | Avg.R $\uparrow$         | F $\uparrow$  | mR $\uparrow$ | mF $\uparrow$ | FWT $\uparrow$ | BWT $\uparrow$ | Avg.R $\uparrow$         | F $\uparrow$  | mR $\uparrow$ | mF $\uparrow$ | FWT $\uparrow$ | BWT $\uparrow$ |
| Joint        | 19.53                    | 0             | 3.9           | 0             | -              | -              | 4.3                      | 0             | 3.7           | 0             | -              | -              |
| Replay@100%  | 13.45                    | -8.83         | 3.6           | -0.35         | -1.5           | -10.45         | 12.45                    | -4.13         | 3.2           | -0.56         | -2.1           | -20.34         |
| Naive        | 0.98                     | -21.2         | 0.74          | -1.35         | -3.45          | -43.87         | 0                        | -18.22        | 0.45          | -2.67         | -4.12          | -53.12         |
| EWC[24]      | 2.36                     | -21.05        | 0.67          | -1.34         | -2.34          | -39.89         | 0                        | -18.22        | 0.03          | 0             | -3.77          | -51.67         |
| PackNet[44]  | 3.2                      | -19.7         | 1.1           | -1.13         | -1.3           | -32.45         | 1.1                      | -17.82        | 0.84          | -1.97         | -2.84          | -40.34         |
| Replay@10%   | 5.67                     | -18.9         | 3.21          | -1.05         | <b>1.45</b>    | -28.34         | 1.81                     | -16.72        | 1.03          | -1.74         | -1.4           | -43.56         |
| Replay@20%   | 6.23                     | -17.45        | 3.5           | -1.01         | 1.01           | -24.32         | 2.37                     | -15.17        | 1.45          | -1.53         | -1.1           | -38.56         |
| <b>Ours*</b> | <b>7.8</b>               | <b>-15.67</b> | <b>3.9</b>    | <b>-0.95</b>  | 0.5            | <b>-19.83</b>  | <b>4.67</b>              | <b>-11.31</b> | <b>2.2</b>    | <b>-0.89</b>  | <b>-0.97</b>   | <b>-29.65</b>  |

CSEGG models still find it **challenging** to learn **new scenes incrementally**

# Main Results



RAS is more proficient in generalizing to **classify relationships** among **unknown objects**

# Future Work

- Explore CSEGG on video-based datasets.
- Develop a synthetic CSEGG dataset to analyze continual learning under controlled conditions.
- Integrate a generative model with fine-grained control signals



We support Open Science!  
Scan QR code for  
Papers, code, data

PhD students and postdocs wanted! Join us!

