

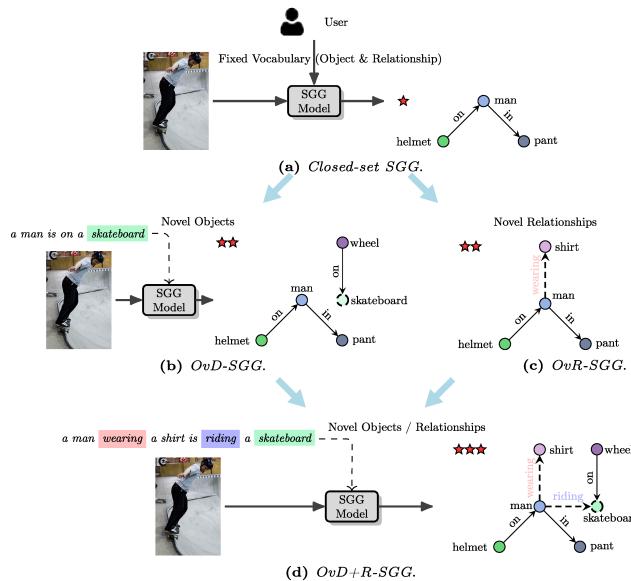
Zuyao (Joseph) Chen<sup>1,2</sup> Jinlin Wu<sup>2,3</sup> Zhen Lei<sup>2,3,4</sup> Zhaoxiang Zhang<sup>2,3,4</sup> Chang Wen Chen<sup>1</sup>

<sup>1</sup>The Hong Kong Polytechnic University <sup>2</sup>Centre for Artificial Intelligence and Robotics, HKISI-CAS

<sup>3</sup>Institute of Automation, CAS <sup>4</sup>University of Chinese Academy of Sciences

## Motivation

Scene Graph Generation (SGG) predominantly operates within a **closed-set** setup, constraining its recognition abilities to a predefined set of objects and relationships.



- Can the model predict unseen objects or relationships?
- What if BOTH unseen objects AND unseen relationships are present?

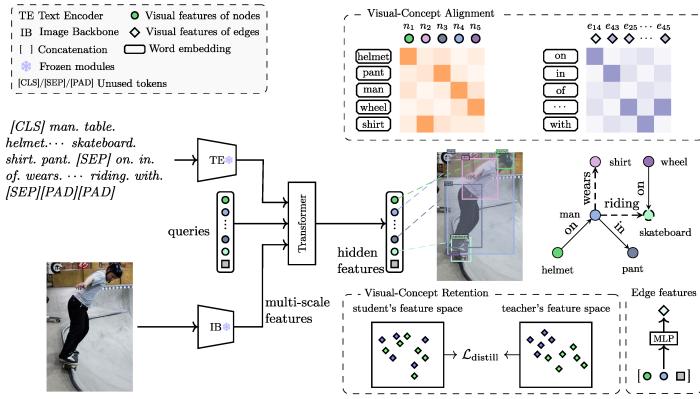
## Challenges

- no relation-aware pre-trained models
- distinguish unseen relationships from “background” (OvR-SGG, OvD+R-SGG)

## Highlights

1. Introduced a unified framework for open-vocabulary SGG, addressing limitations of traditional methods
2. Proposed novel visual-concept alignment and retention mechanisms for improved generalization to unseen concepts
3. Achieved state-of-the-art results, demonstrating the effectiveness of OvSGTR

## Unified framework: OvSGTR



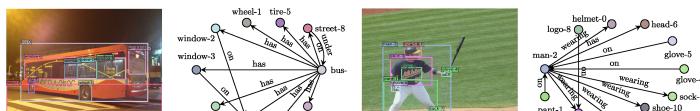
- DETR-like structure instead of R-CNN
- Visual-concept alignment
- Visual-concept retention
- Relation-aware pre-training

## Experiments

### Experimental results of Closed-set SGG on VG150 test set.

SGG model	Backbone	Detector	Params.	R@20/50/100			mR@20/50/100	Time (s)		
				R@20	R@50	R@100				
IMP	RX-101	Faster R-CNN	146M/308M	17.7	25.5	30.7	2.7	4.1	5.3	0.25
MOTIFS			205M/367M	25.5	32.8	37.2	5.0	6.8	7.9	0.27
VCTREE			197M/358M	24.7	31.5	36.2	-	-	-	0.38
HL-Net			220M/382M	26.0	33.7	38.1	-	-	-	0.10
FCSGG	HRNetW48	-	87M/87M	13.6	18.6	22.5	2.3	3.2	3.9	0.13
SGTR	R-101	DETR	36M/96M	-	24.6	28.4	-	-	-	0.21
VS <sup>3</sup>	Swin-T	DETR	93M/233M	26.1	34.5	39.2	-	-	-	0.16
VS <sup>3</sup>	Swin-L	DETR	124M/432M	27.3	36.0	40.9	4.4	6.5	7.8	0.24
OvSGTR	Swin-T	DETR	41M/178M	27.0	35.8	41.3	5.0	7.2	8.8	0.13
OvSGTR	Swin-B	DETR	41M/238M	27.8	36.4	42.4	5.2	7.4	9.0	0.19

state-of-the-art performance  
fewer parameters & lower inference latency

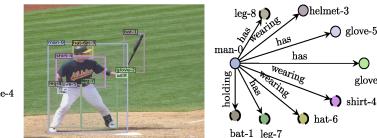


Trained on Closed-set SGG

## Experiments



Trained on OvD+R-SGG



Real-world example

### Experimental results of OvD-SGG.

Method	Base+Novel (Object)		Novel (Object)	
	PREDCLS	SGDET	PREDCLS	SGDET
VS <sup>3</sup> -T	50.10 / 52.05	15.07 / 18.73	46.91 / 49.13	10.08 / 13.65
OvSGTR-T	<b>60.58</b> / 62.10	<b>18.14</b> / 23.20	<b>59.01</b> / 60.65	<b>12.06</b> / 16.49
OvSGTR-B	<b>60.83</b> / 62.33	<b>21.35</b> / 26.22	<b>59.30</b> / 60.95	<b>15.58</b> / 19.96

### Experimental results of OvR-SGG.

Method	Base+Novel (Relation)		Novel (Relation)	
	R@50	R@100	R@50	R@100
VS <sup>3</sup> -T	15.60	17.30	0.00	0.00
OvSGTR-T <sup>†</sup>	<b>17.71</b>	20.00	0.34	0.41
OvSGTR-T	<b>20.46</b>	<b>23.86</b>	<b>13.45</b>	<b>16.19</b>
OvSGTR-B <sup>†</sup>	18.58	20.84	0.08	0.10
OvSGTR-B	<b>22.89</b>	<b>26.65</b>	<b>16.39</b>	<b>19.72</b>

### Experimental results of OvD+R-SGG.

Method	Joint Base+Novel		Novel (Object)		Novel (Relation)	
	R@50	R@100	R@50	R@100	R@50	R@100
VS <sup>3</sup> -T	5.88	7.20	6.00	7.51	0.00	0.00
OvSGTR-T <sup>†</sup>	7.88	10.06	6.82	9.23	0.00	0.00
OvSGTR-T	<b>13.53</b>	<b>16.36</b>	<b>14.37</b>	<b>17.44</b>	<b>9.20</b>	<b>11.19</b>
OvSGTR-B <sup>†</sup>	11.23	14.21	13.27	16.83	1.78	2.57
OvSGTR-B	<b>17.11</b>	<b>21.02</b>	<b>17.58</b>	<b>21.72</b>	<b>14.56</b>	<b>18.20</b>

<sup>†</sup> refers to w.o. distillation

Models and Code are available:

<https://github.com/gpt4vision/OvSGTR>