

Appendix: Fine-Grained Scene Graph Generation with Data Transfer

Ao Zhang^{1*}, Yuan Yao^{2*}, Qianyu Chen², Wei Ji^{1†}, Zhiyuan Liu^{2†},
Maosong Sun², and Tat-Seng Chua¹

¹ Sea-NExT Joint Lab, Singapore

School of Computing, National University of Singapore, Singapore

² Department of Computer Science and Technology

Institute for Artificial Intelligence, Tsinghua University, Beijing, China

Beijing National Research Center for Information Science and Technology, China

aozhang@u.nus.edu, yaoyuanthu@163.com

1 VG-1800 Dataset

The VG-1800 dataset aims to provide reliable evaluation for the large-scale scene graph generation.

1.1 Dataset Construction

We construct the dataset based on original Visual Genome dataset [2] by the following steps: (1) **Filtration**. Instead of simply auto-filtering [7] and choosing the top frequent predicate categories [1], we manually filter out unreasonable predicate categories, including misspelling predicates (e.g., *i frot of*), adjectives (e.g., *white*), nouns (e.g., *car*), and relative clauses (e.g., *who has*). To provide enough relation instances for robust evaluation, we retain all object categories and predicate categories with over 5 samples. (2) **Split**. We split the VG dataset into 70% training and 30% test. Following VG-50 split [5], we further split out 5,000 images from the training set as the validation set, and ensure at least 5 samples on the test set and at least 1 samples on the training set for each predicate category.

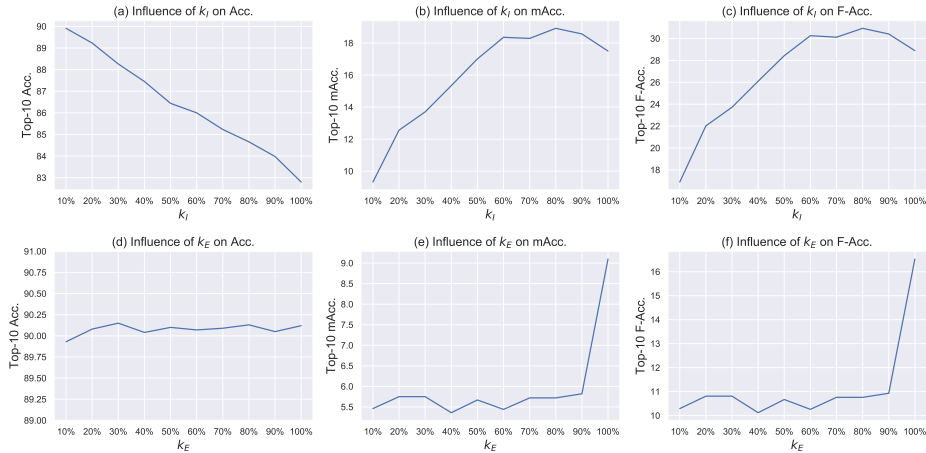
1.2 Dataset Statistics

Finally, the dataset contains 70,098 object categories, 1,807 predicate categories and 272,084 distinct relation triplets. It consists of 66,289, 4,995, and 32,893 images for training set, validation set and test set respectively. There are on average 19.5 objects and 16.0 relations for each image.

Comparison with Other VG Splits. We also compare our VG-1800 split with other splits based on Visual Genome [2] dataset, including a conventional VG-50

* indicates equal contribution.

† Corresponding author: jiwei@nus.edu.sg, liuzy@tsinghua.edu.cn

Fig. 1: Influence of k_I and k_E in different metrics.

and the other two large-scale SGG splits VG8K and VG8K-LT. Our VG-1800 can provide a more reliable evaluation for large-scale SGG. (1) When compared with VG8K, we provide a much cleaner dataset by manually cleaning the noise. For example, VG8K does not filter out nouns and adjectives, which will lead to an unreliable evaluation. (2) When compared with VG8K-LT, a cleaner version of VG8K, we provide a much stable evaluation for large amount of tail classes. More specifically, as shown in Table 1, our VG-1800 contains more test images. Meanwhile, VG-1800 also contains more samples of tail classes. As shown in Figure 2, our VG-1800 has 1,807 predicate classes with no less than 5 samples, while VG8K-LT has only 526 classes that have no less than 5 samples.

Table 1: Comparison between different datasets’ predicate filtration and split ratio. Filtration denotes the method to remove noisy predicates. The Train, Val, and Test denote number of images in training set, validation set and test set.

Dataset	Filtration	Train	Val	Test
VG-50 [5]	-	57,723	5,000	26,446
VG8K [7]	Auto	97,961	2,000	4,871
VG8K-LT [1]	Auto	97,623	1,999	4,860
VG-1800	Manual	66,289	4,995	32,893

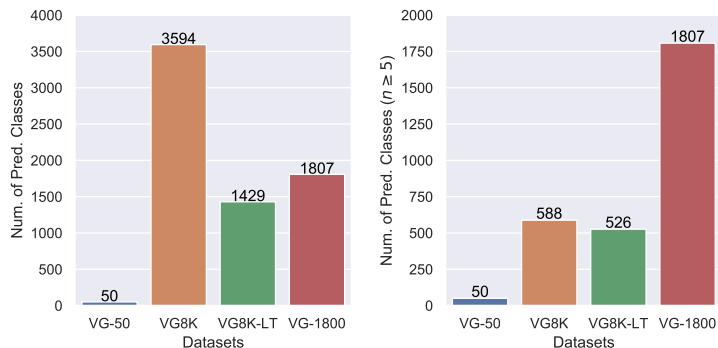


Fig. 2: Comparison of number of predicate classes between different splits on test set. The $n \geq 5$ denotes predicate classes having no less than 5 samples.

2 Implementation Details

2.1 VG-50

All normally trained baseline models including Motif, Transformer, VCTree, and GPS-Net are reproduced by us. For fair comparison, we equally remove all re-sampling and reweighting strategies. Moreover, to encourage informative SGG, we remove the frequency bias in the training and inference of base models, which may make the result tend to have higher mR@K, F@K and lower R@K than results in their original papers.

For Transformer, some implementation details are different. The batch size can be enlarged to 16 on 2 GPUs. The learning rate is reduced to 0.08 for PREDCLS and SGDET for training stability. For Transformer on SGCLS, where the training is even more unstable, we further lower the learning rate to 0.016. All experiments are done on RTX-2080ti GPUs.

2.2 VG-1800

Compared with VG-50, the same backbone, parameter fixation, learning rate, optimizer, and learning schedule are used on VG-1800 dataset. Specially, due to the significant increase of (*subject*, *predicate*, *object*) combinations, we equally remove all frequency bias items on VG-1800 to reduce machines' memory usage. For internal and external transfer, the k_I is set to 90% and k_E is set to 100%.

For baselines, we find that due to the significant difference between the number of head predicates and tail predicates, the bi-level resampling in BGNN [3] will make the model pay most of the attention on tail classes while ignoring head classes. The drop out rate of images that do not contain rare predicates are set to almost 100%. This lead to a bad convergence of BGNN. Thus, we remove bi-level resampling for BGNN results. For RelMix [1], we equip the proposed VilHub loss and predicate feature mixup to the Neural Motif model. Similar

with bi-level sampling, we find that the reweighting strategy also lead to worse results. Thus, we do not include a reweighting version like VG-50.

3 Supplementary Experiments

3.1 Influence of k_I and k_E

Influence of k_I . To provide a more detailed analysis on the influence of k_I , we report the performance on Acc and mAcc with different k_I . As shown in Figure 1, with the increase of internal transfer percentage k_I , Acc decreases linearly, while mAcc first increases when $k_I \leq 80\%$ and then decreases. The phenomenon shows that transferring more in internal transfer does not necessarily mean higher mAcc. For VG-1800, The first 80% internal data transfer is helpful to improve mAcc, while the last confident 20% will harm the overall performance. We guess the last 20% data may contain too noise, which will lower the data quality for model training.

Influence of k_E . As for k_E , external transfer shows almost no influence on Acc and mAcc when $k_E \leq 90\%$, while significantly boost mAcc when $k_E = 100\%$. Contrary to our observations for k_I , the last 10% samples which are believed to be unuseful by models, seem to bring the most profitable boost for mAcc. We guess the reason is that model can not distinguish well between tail classes and NA samples, while this part of the data is essential to provide more training samples for tail classes.

3.2 Adaptive Threshold.

In our IETrans, when determining how much data to transfer, we equally use a fixed percentage number for all relational triplets, which seems to be sub-optimal. Thus, we also tried an adaptive threshold by considering the prediction score of concrete relational triplet instances.

For internal transfer, given a general relational triplet instance (o_s, p_G, o_o) , we are required to decide whether to transfer to its corresponding informative type p_I . We denote the model’s prediction score of object pair (o_s, o_o) on p_I as $s_{p_I}^*$. We denote the average and standard error of all (c_{o_s}, p_I, c_{o_o}) relational triplet instances’ prediction score on p_I as μ_I and σ_I . We conduct the transfer when $s_{p_I}^*$ satisfies:

$$s_{p_I}^* > \mu_I + k\sigma_I, \tag{1}$$

where k is a hyperparameter. The intuition is that if a general instance (o_s, p_G, o_o) ’s prediction score is over the average of all real (c_{o_s}, p_I, c_{o_o}) ’s prediction scores on p_I , (o_s, p_G, o_o) can probably be relabeled as an informative one. Meanwhile, the standard error is considered to further control the adaptive threshold.

By choosing different k including $\{-1.0, -0.5, 0.0, 0.5, 1.0\}$, we can get a curve with different Acc and mAcc trade-offs. As shown in Figure 3, the model with adaptive thresholds is overall worse than our fixed percentage.

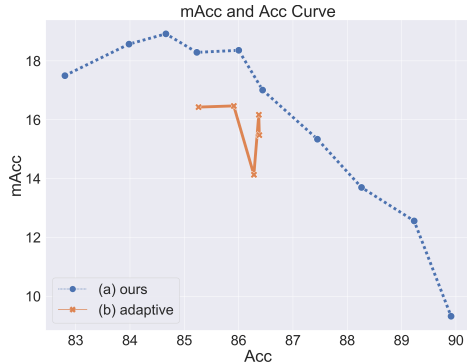


Fig. 3: Comparison of adaptive thresholds and fixed percentages for Internal Transfer.

A possible explanation is that the prediction score of an instance is non-linearly dependent on the number of its own instances and its similarity with different general classes, which results in inconsistency among different relational triplets. Especially when the number of an informative relational triplet is very small, the average of its prediction score is often near to zero, which will easily lead to an over-transfer problem. Thus, we leave the design of a more intelligent adaptive threshold for future work.

For external transfer, as shown in the paper, the prediction scores on NA of missed annotated samples are unreliable, i.e. the samples with the highest NA score bring maximum benefits for the model’s performance.

SGCLS Results on VG-1800. We also evaluate our method on SGCLS task

Table 2: SGCLS triplet-level evaluation results on VG-1800 dataset.

Models	Top-1				Top-5				Top-10			
	Acc	mAcc	F-Acc	Non-Zero	Acc	mAcc	F-Acc	Non-Zero	Acc	mAcc	F-Acc	Non-Zero
BGNN [3]	16.29	0.18	0.36	22	22.99	0.86	1.65	159	24.15	1.48	2.78	221
Motif [6]	18.93	0.18	0.36	37	26.08	0.74	1.43	90	27.28	1.15	2.21	121
-Focal Loss	18.55	0.14	0.28	29	25.91	0.51	1.00	52	27.14	0.76	1.48	80
-TDE [4]	18.02	0.11	0.22	15	24.85	0.38	0.75	38	26.14	0.56	1.10	53
-RelMix [1]	18.27	0.22	0.43	47	25.57	0.83	1.60	100	26.71	1.26	2.42	130
-IETrans ($k_I = 10\%$) (ours)	18.24	1.68	3.16	212	25.80	1.68	3.16	212	27.25	2.54	4.65	264
-IETrans ($k_I = 90\%$) (ours)	4.91	1.78	2.62	298	20.66	4.72	7.68	538	24.85	6.54	10.36	637

on VG-1800 dataset. As shown in table 2, the comparison with other baselines is similar to the results on PREDCLS task. When compared with Motif, our IETrans ($k_I = 10\%$) can achieve significant improvement on top-1 mAcc and Non-Zero metrics (over 5 times of Motif) with negligible degeneration (less than

1 point) on top-1 Acc metric. Our IETrans ($k_I = 90\%$) can further boost the mAcc and Non-Zero metrics, which shows the ability of our IETrans to generate informative scene graphs. When compared with other baselines, our IETrans can achieve the best F-Acc metrics across top-1, top-5, and top-10 evaluations. However, there is a large gap between SGCLS results and PREDCLS results (e.g., 2.62% vs. 4.70% for top-1 F-Acc of IETrans ($k_I = 90\%$)), which indicates that further effort should be made to explore the joint optimization of both objects and predicates.

4 Discovered Visual Hierarchy Analysis

Visual Hierarchy Evaluation. A key element of conducting correct internal transfer is to find reasonable general-informative relation pairs. To evaluate the precision, we randomly choose 50 pairs with over 3 samples being transferred, so as to avoid involving too many noise-to-noise pairs. Then, human evaluation is conducted. The ratio of reasonable general-informative pairs is **76%** for VG-50, and **74%** for VG-1800.

Visualization. In the following, we show 100 discovered general-informative pairs for both VG-50 and VG-1800. The pairs are ranked by the number of samples which are transferred.

Table 3: Examples of discovered visual hierarchy in VG-50

<i>(window, on, building)</i>	\rightarrow	<i>(window, part of, building)</i>
<i>(man, wearing, arm)</i>	\rightarrow	<i>(man, wears, arm)</i>
<i>(boy, wearing, boy)</i>	\rightarrow	<i>(boy, wears, boy)</i>
<i>(pillow, on, bed)</i>	\rightarrow	<i>(pillow, lying on, bed)</i>
<i>(building, has, building)</i>	\rightarrow	<i>(building, made of, building)</i>
<i>(sign, on, building)</i>	\rightarrow	<i>(sign, mounted on, building)</i>
<i>(arm, of, arm)</i>	\rightarrow	<i>(arm, belonging to, arm)</i>
<i>(sign, on, building)</i>	\rightarrow	<i>(sign, hanging from, building)</i>
<i>(man, wearing, bag)</i>	\rightarrow	<i>(man, wears, bag)</i>
<i>(window, on, building)</i>	\rightarrow	<i>(window, belonging to, building)</i>
<i>(car, on, building)</i>	\rightarrow	<i>(car, parked on, building)</i>
<i>(clock, on, building)</i>	\rightarrow	<i>(clock, mounted on, building)</i>
<i>(man, wearing, building)</i>	\rightarrow	<i>(man, wears, building)</i>
<i>(man, has, arm)</i>	\rightarrow	<i>(man, wears, arm)</i>
<i>(man, wearing, boot)</i>	\rightarrow	<i>(man, wears, boot)</i>
<i>(bottle, on, bottle)</i>	\rightarrow	<i>(bottle, sitting on, bottle)</i>
<i>(window, on, bus)</i>	\rightarrow	<i>(window, belonging to, bus)</i>
<i>(book, on, book)</i>	\rightarrow	<i>(book, above, book)</i>
<i>(man, has, arm)</i>	\rightarrow	<i>(man, with, arm)</i>
<i>(ear, of, ear)</i>	\rightarrow	<i>(ear, belonging to, ear)</i>
<i>(hand, of, arm)</i>	\rightarrow	<i>(hand, belonging to, arm)</i>

<i>(bottle, on, bottle)</i>	<i>→</i>	<i>(bottle, above, bottle)</i>
<i>(window, in, building)</i>	<i>→</i>	<i>(window, part of, building)</i>
<i>(light, on, building)</i>	<i>→</i>	<i>(light, mounted on, building)</i>
<i>(door, on, building)</i>	<i>→</i>	<i>(door, to, building)</i>
<i>(food, on, food)</i>	<i>→</i>	<i>(food, lying on, food)</i>
<i>(bowl, on, bowl)</i>	<i>→</i>	<i>(bowl, above, bowl)</i>
<i>(man, wearing, man)</i>	<i>→</i>	<i>(man, wears, man)</i>
<i>(flower, in, flower)</i>	<i>→</i>	<i>(flower, painted on, flower)</i>
<i>(woman, wearing, bag)</i>	<i>→</i>	<i>(woman, wears, bag)</i>
<i>(tree, near, building)</i>	<i>→</i>	<i>(tree, in front of, building)</i>
<i>(woman, wearing, arm)</i>	<i>→</i>	<i>(woman, wears, arm)</i>
<i>(window, of, building)</i>	<i>→</i>	<i>(window, part of, building)</i>
<i>(clock, on, building)</i>	<i>→</i>	<i>(clock, part of, building)</i>
<i>(bus, on, building)</i>	<i>→</i>	<i>(bus, parked on, building)</i>
<i>(man, wearing, coat)</i>	<i>→</i>	<i>(man, wears, coat)</i>
<i>(building, has, building)</i>	<i>→</i>	<i>(building, with, building)</i>
<i>(boy, wearing, arm)</i>	<i>→</i>	<i>(boy, wears, arm)</i>
<i>(man, on, arm)</i>	<i>→</i>	<i>(man, riding, arm)</i>
<i>(tree, has, branch)</i>	<i>→</i>	<i>(tree, with, branch)</i>
<i>(woman, wearing, boot)</i>	<i>→</i>	<i>(woman, wears, boot)</i>
<i>(pillow, on, bed)</i>	<i>→</i>	<i>(pillow, above, bed)</i>
<i>(woman, has, arm)</i>	<i>→</i>	<i>(woman, with, arm)</i>
<i>(window, on, building)</i>	<i>→</i>	<i>(window, to, building)</i>
<i>(glass, on, bottle)</i>	<i>→</i>	<i>(glass, sitting on, bottle)</i>
<i>(sign, on, building)</i>	<i>→</i>	<i>(sign, says, building)</i>
<i>(man, wearing, bike)</i>	<i>→</i>	<i>(man, wears, bike)</i>
<i>(tire, on, building)</i>	<i>→</i>	<i>(tire, on back of, building)</i>
<i>(branch, on, branch)</i>	<i>→</i>	<i>(branch, growing on, branch)</i>
<i>(book, on, book)</i>	<i>→</i>	<i>(book, laying on, book)</i>
<i>(car, on, car)</i>	<i>→</i>	<i>(car, parked on, car)</i>
<i>(man, wearing, bench)</i>	<i>→</i>	<i>(man, wears, bench)</i>
<i>(elephant, has, ear)</i>	<i>→</i>	<i>(elephant, using, ear)</i>
<i>(person, on, beach)</i>	<i>→</i>	<i>(person, standing on, beach)</i>
<i>(wing, on, plane)</i>	<i>→</i>	<i>(wing, attached to, plane)</i>
<i>(windshield, on, building)</i>	<i>→</i>	<i>(windshield, of, building)</i>
<i>(arm, on, arm)</i>	<i>→</i>	<i>(arm, belonging to, arm)</i>
<i>(woman, holding, bag)</i>	<i>→</i>	<i>(woman, carrying, bag)</i>
<i>(window, on, bike)</i>	<i>→</i>	<i>(window, part of, bike)</i>
<i>(ear, on, ear)</i>	<i>→</i>	<i>(ear, belonging to, ear)</i>
<i>(man, on, beach)</i>	<i>→</i>	<i>(man, walking on, beach)</i>
<i>(boy, has, boy)</i>	<i>→</i>	<i>(boy, wears, boy)</i>
<i>(roof, on, building)</i>	<i>→</i>	<i>(roof, covering, building)</i>
<i>(leaf, on, branch)</i>	<i>→</i>	<i>(leaf, growing on, branch)</i>

<i>(head, of, arm)</i>	<i>→ (head, belonging to, arm)</i>
<i>(wheel, on, building)</i>	<i>→ (wheel, on back of, building)</i>
<i>(tree, near, building)</i>	<i>→ (tree, along, building)</i>
<i>(bird, on, bird)</i>	<i>→ (bird, sitting on, bird)</i>
<i>(door, on, door)</i>	<i>→ (door, to, door)</i>
<i>(woman, has, bag)</i>	<i>→ (woman, with, bag)</i>
<i>(man, wearing, hat)</i>	<i>→ (man, wears, hat)</i>
<i>(man, on, arm)</i>	<i>→ (man, standing on, arm)</i>
<i>(sign, on, building)</i>	<i>→ (sign, attached to, building)</i>
<i>(letter, on, building)</i>	<i>→ (letter, painted on, building)</i>
<i>(bird, on, bird)</i>	<i>→ (bird, standing on, bird)</i>
<i>(ear, of, cat)</i>	<i>→ (ear, belonging to, cat)</i>
<i>(window, on, bench)</i>	<i>→ (window, part of, bench)</i>
<i>(window, near, building)</i>	<i>→ (window, part of, building)</i>
<i>(wheel, on, bike)</i>	<i>→ (wheel, on back of, bike)</i>
<i>(building, near, building)</i>	<i>→ (building, across, building)</i>
<i>(elephant, has, ear)</i>	<i>→ (elephant, between, ear)</i>
<i>(man, has, bag)</i>	<i>→ (man, with, bag)</i>
<i>(engine, on, plane)</i>	<i>→ (engine, mounted on, plane)</i>
<i>(man, wearing, chair)</i>	<i>→ (man, wears, chair)</i>
<i>(woman, wearing, building)</i>	<i>→ (woman, wears, building)</i>
<i>(sign, on, sign)</i>	<i>→ (sign, mounted on, sign)</i>
<i>(plate, on, bowl)</i>	<i>→ (plate, above, bowl)</i>
<i>(man, wearing, face)</i>	<i>→ (man, wears, face)</i>
<i>(leg, of, giraffe)</i>	<i>→ (leg, part of, giraffe)</i>
<i>(pillow, above, bed)</i>	<i>→ (pillow, lying on, bed)</i>
<i>(tire, on, bike)</i>	<i>→ (tire, on back of, bike)</i>
<i>(leg, of, arm)</i>	<i>→ (leg, belonging to, arm)</i>
<i>(man, wearing, cap)</i>	<i>→ (man, wears, cap)</i>
<i>(bird, has, bird)</i>	<i>→ (bird, with, bird)</i>
<i>(trunk, of, ear)</i>	<i>→ (trunk, belonging to, ear)</i>
<i>(roof, of, building)</i>	<i>→ (roof, covering, building)</i>
<i>(plate, on, bottle)</i>	<i>→ (plate, above, bottle)</i>
<i>(man, wearing, ear)</i>	<i>→ (man, wears, ear)</i>
<i>(man, has, building)</i>	<i>→ (man, wears, building)</i>
<i>(window, on, arm)</i>	<i>→ (window, part of, arm)</i>

Table 4: Examples of discovered visual hierarchy in VG-1800

<i>(window, on, building)</i>	<i>→ (window, on exterior of, building)</i>
<i>(man, wearing, arm)</i>	<i>→ (man, wearing striped, arm)</i>
<i>(arm, of, arm)</i>	<i>→ (arm, stretched out on, arm)</i>

<i>(man, has, arm)</i>	→	<i>(man, stretching out, arm)</i>
<i>(cloud, in, cloud)</i>	→	<i>(cloud, floating through, cloud)</i>
<i>(pillow, on, bed)</i>	→	<i>(pillow, propped up on, bed)</i>
<i>(tree, in, background)</i>	→	<i>(tree, visible in, background)</i>
<i>(leg, of, arm)</i>	→	<i>(leg, belonging to, arm)</i>
<i>(boat, in, boat)</i>	→	<i>(boat, sailing on, boat)</i>
<i>(building, has, building)</i>	→	<i>(building, seen outside, building)</i>
<i>(cloud, in, building)</i>	→	<i>(cloud, floating through, building)</i>
<i>(hand, of, arm)</i>	→	<i>(hand, hand of, arm)</i>
<i>(boat, on, boat)</i>	→	<i>(boat, sailing on, boat)</i>
<i>(cloud, in, sky)</i>	→	<i>(cloud, floating through, sky)</i>
<i>(building, in, background)</i>	→	<i>(building, visible in, background)</i>
<i>(man, has, arm)</i>	→	<i>(man, sheltering, arm)</i>
<i>(man, wearing, arm)</i>	→	<i>(man, dressed in, arm)</i>
<i>(head, of, arm)</i>	→	<i>(head, turning, arm)</i>
<i>(man, has, arm)</i>	→	<i>(man, losing, arm)</i>
<i>(cloud, in, airplane)</i>	→	<i>(cloud, floating through, airplane)</i>
<i>(cloud, in, background)</i>	→	<i>(cloud, floating through, background)</i>
<i>(window, on, awning)</i>	→	<i>(window, on exterior of, awning)</i>
<i>(man, has, arm)</i>	→	<i>(man, pointing with, arm)</i>
<i>(man, has, arm)</i>	→	<i>(man, spreading, arm)</i>
<i>(window, on, bus)</i>	→	<i>(window, lining side of, bus)</i>
<i>(window, on, balcony)</i>	→	<i>(window, on exterior of, balcony)</i>
<i>(cloud, in, beach)</i>	→	<i>(cloud, floating through, beach)</i>
<i>(car, on, building)</i>	→	<i>(car, driving alongside, building)</i>
<i>(window, on, building)</i>	→	<i>(window, lining side of, building)</i>
<i>(man, wearing, arm)</i>	→	<i>(man, lifting up, arm)</i>
<i>(shirt, on, arm)</i>	→	<i>(shirt, worn by, arm)</i>
<i>(man, wearing, bag)</i>	→	<i>(man, wearing striped, bag)</i>
<i>(bottle, on, bottle)</i>	→	<i>(bottle, kept on, bottle)</i>
<i>(man, wearing, background)</i>	→	<i>(man, wearing striped, background)</i>
<i>(boy, wearing, boy)</i>	→	<i>(boy, striped, boy)</i>
<i>(cloud, in, air)</i>	→	<i>(cloud, floating through, air)</i>
<i>(woman, has, arm)</i>	→	<i>(woman, raising, arm)</i>
<i>(car, on, building)</i>	→	<i>(car, moving down, building)</i>
<i>(airplane, in, airplane)</i>	→	<i>(airplane, flying under, airplane)</i>
<i>(tile, on, bathroom)</i>	→	<i>(tile, fixed to, bathroom)</i>
<i>(arm, on, arm)</i>	→	<i>(arm, stretched out on, arm)</i>
<i>(cloud, in, arm)</i>	→	<i>(cloud, floating through, arm)</i>
<i>(head, of, arm)</i>	→	<i>(head, belonging to, arm)</i>
<i>(man, wearing, arm)</i>	→	<i>(man, kicking up, arm)</i>
<i>(airplane, on, airplane)</i>	→	<i>(airplane, taking off from, airplane)</i>
<i>(sign, on, building)</i>	→	<i>(sign, strapped, building)</i>

<i>(man, wearing, air)</i>	<i>→</i>	<i>(man, wearing striped, air)</i>
<i>(sign, on, arrow)</i>	<i>→</i>	<i>(sign, strapped, arrow)</i>
<i>(window, on, building)</i>	<i>→</i>	<i>(window, adorning, building)</i>
<i>(cat, has, cat)</i>	<i>→</i>	<i>(cat, possesses, cat)</i>
<i>(cloud, in, boat)</i>	<i>→</i>	<i>(cloud, floating through, boat)</i>
<i>(person, has, arm)</i>	<i>→</i>	<i>(person, stretching out, arm)</i>
<i>(clock, on, building)</i>	<i>→</i>	<i>(clock, attached to side of, building)</i>
<i>(train, on, building)</i>	<i>→</i>	<i>(train, switching, building)</i>
<i>(woman, has, arm)</i>	<i>→</i>	<i>(woman, combing, arm)</i>
<i>(boy, wearing, arm)</i>	<i>→</i>	<i>(boy, striped, arm)</i>
<i>(window, on, building)</i>	<i>→</i>	<i>(window, on the side of, building)</i>
<i>(window, of, building)</i>	<i>→</i>	<i>(window, on exterior of, building)</i>
<i>(arrow, on, arrow)</i>	<i>→</i>	<i>(arrow, printed, arrow)</i>
<i>(woman, wearing, arm)</i>	<i>→</i>	<i>(woman, wearing striped, arm)</i>
<i>(head, of, arm)</i>	<i>→</i>	<i>(head, turned to, arm)</i>
<i>(branch, on, branch)</i>	<i>→</i>	<i>(branch, sticking up on, branch)</i>
<i>(man, on, arm)</i>	<i>→</i>	<i>(man, swimming with, arm)</i>
<i>(wing, on, airplane)</i>	<i>→</i>	<i>(wing, on left side of, airplane)</i>
<i>(mountain, in, background)</i>	<i>→</i>	<i>(mountain, visible in, background)</i>
<i>(bowl, on, bowl)</i>	<i>→</i>	<i>(bowl, placed on, bowl)</i>
<i>(cloud, in, blue sky)</i>	<i>→</i>	<i>(cloud, floating through, blue sky)</i>
<i>(window, on, arrow)</i>	<i>→</i>	<i>(window, on exterior of, arrow)</i>
<i>(foot, of, arm)</i>	<i>→</i>	<i>(foot, belonging to, arm)</i>
<i>(toilet, in, bathroom)</i>	<i>→</i>	<i>(toilet, installed in, bathroom)</i>
<i>(sign, on, building)</i>	<i>→</i>	<i>(sign, anchored to, building)</i>
<i>(wall, on, building)</i>	<i>→</i>	<i>(wall, making up, building)</i>
<i>(kite, in, air)</i>	<i>→</i>	<i>(kite, flying through, air)</i>
<i>(leaf, on, building)</i>	<i>→</i>	<i>(leaf, growing on, building)</i>
<i>(man, wearing, arm)</i>	<i>→</i>	<i>(man, adjusting, arm)</i>
<i>(person, wearing, arm)</i>	<i>→</i>	<i>(person, striped, arm)</i>
<i>(tile, on, bathroom)</i>	<i>→</i>	<i>(tile, installed on, bathroom)</i>
<i>(bag, on, bag)</i>	<i>→</i>	<i>(bag, kept in, bag)</i>
<i>(sky, in, sky)</i>	<i>→</i>	<i>(sky, stretched across, sky)</i>
<i>(bear, has, bear)</i>	<i>→</i>	<i>(bear, scratching, bear)</i>
<i>(line, on, building)</i>	<i>→</i>	<i>(line, painted in, building)</i>
<i>(man, wearing, building)</i>	<i>→</i>	<i>(man, wearing striped, building)</i>
<i>(book, on, book)</i>	<i>→</i>	<i>(book, arranged on, book)</i>
<i>(wall, near, building)</i>	<i>→</i>	<i>(wall, making up, building)</i>
<i>(window, on, advertisement)</i>	<i>→</i>	<i>(window, lining side of, advertisement)</i>
<i>(sink, in, bathroom)</i>	<i>→</i>	<i>(sink, mounted in, bathroom)</i>
<i>(wave, in, arm)</i>	<i>→</i>	<i>(wave, cresting in, arm)</i>
<i>(window, in, building)</i>	<i>→</i>	<i>(window, on exterior of, building)</i>
<i>(blanket, on, bed)</i>	<i>→</i>	<i>(blanket, laying over, bed)</i>

<i>(man, wearing, backpack)</i>	→	<i>(man, carrying, backpack)</i>
<i>(hair, of, arm)</i>	→	<i>(hair, on head of, arm)</i>
<i>(kite, in, beach)</i>	→	<i>(kite, kite in, beach)</i>
<i>(man, wearing, arm)</i>	→	<i>(man, dressed, arm)</i>
<i>(boy, has, arm)</i>	→	<i>(boy, outstretched, arm)</i>
<i>(picture, on, bed)</i>	→	<i>(picture, framed on, bed)</i>
<i>(window, on, banner)</i>	→	<i>(window, on exterior of, banner)</i>
<i>(arm, of, arm)</i>	→	<i>(arm, belonging to, arm)</i>
<i>(man, on, board)</i>	→	<i>(man, going off, board)</i>
<i>(shadow, on, arm)</i>	→	<i>(shadow, cast over, arm)</i>
<i>(arm, of, arm)</i>	→	<i>(arm, around neck of, arm)</i>

References

1. Abdelkarim, S., Agarwal, A., Achlioptas, P., Chen, J., Huang, J., Li, B., Church, K., Elhoseiny, M.: Exploring long tail visual relationship recognition with large vocabulary. In: Proceedings of ICCV. pp. 15921–15930 (2021)
2. Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L.J., Shamma, D.A., et al.: Visual Genome: Connecting language and vision using crowdsourced dense image annotations. IJCV pp. 32–73 (2017)
3. Li, R., Zhang, S., Wan, B., He, X.: Bipartite graph network with adaptive message passing for unbiased scene graph generation. In: Proceedings of CVPR. pp. 11109–11119 (2021)
4. Tang, K., Niu, Y., Huang, J., Shi, J., Zhang, H.: Unbiased scene graph generation from biased training. In: Proceedings of CVPR. pp. 3716–3725 (2020)
5. Xu, D., Zhu, Y., Choy, C.B., Fei-Fei, L.: Scene graph generation by iterative message passing. In: Proceedings of CVPR. pp. 5410–5419 (2017)
6. Zellers, R., Yatskar, M., Thomson, S., Choi, Y.: Neural Motifs: Scene graph parsing with global context. In: Proceedings of CVPR. pp. 5831–5840 (2018)
7. Zhang, J., Kalantidis, Y., Rohrbach, M., Paluri, M., Elgammal, A., Elhoseiny, M.: Large-scale visual relationship understanding. In: Proceedings of the AAAI. pp. 9185–9194 (2019)