

Dual Graph Convolutional Networks with Transformer and Curriculum Learning for Image Captioning

Xinzhi Dong¹, Chengjiang Long^{2*}, Wenju Xu³, Chunxia Xiao^{1*†}

School of Computer Science, Wuhan University¹, JD Finance America Corporation², OPPO US Research Center, InnoPeak Technology Inc.³

dongxz97@whu.edu.cn, cjfykx@gmail.com, xuwenju123@gmail.com, cxxiao@whu.edu.cn

Background & Problem



A stop sign is on a road.

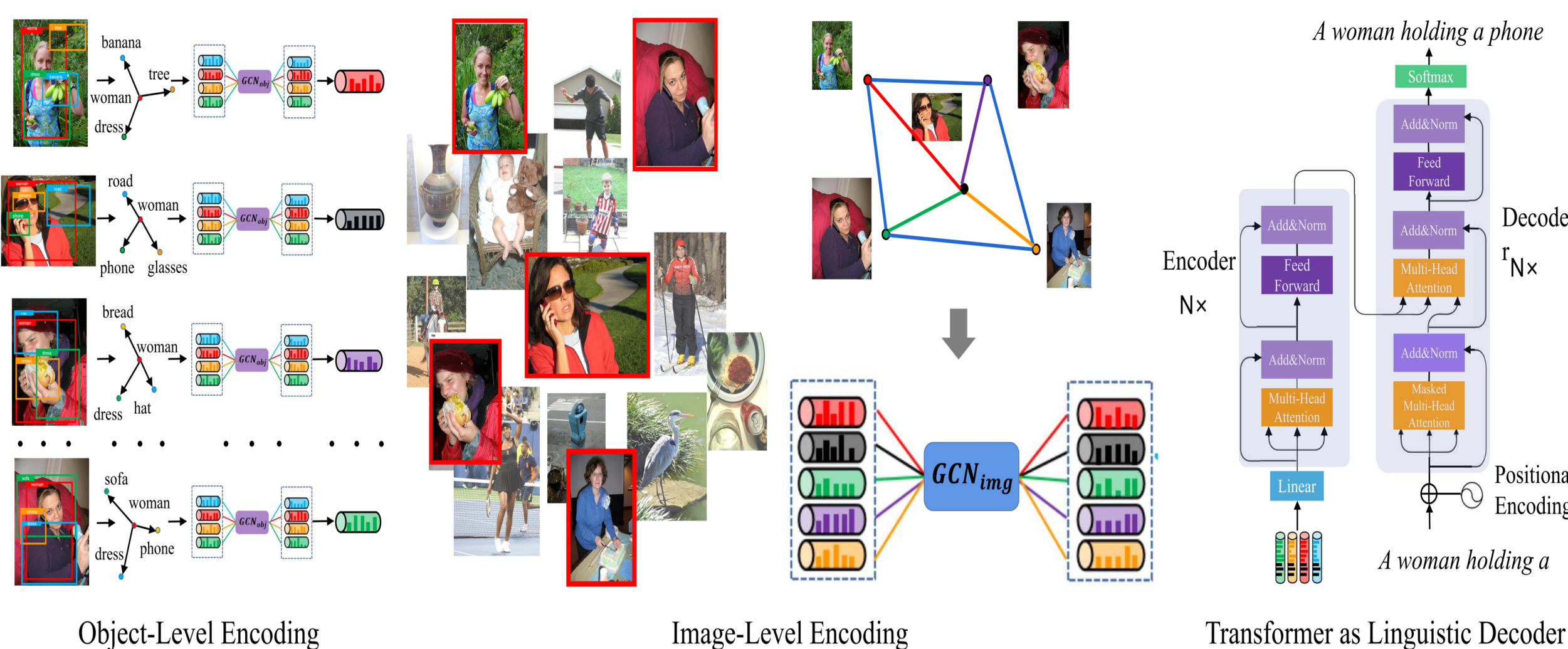
- Image captioning aims to generate a caption from a single image.
- It has excellent development prospects in unmanned driving, unmanned supermarket and other applications.
- Most of the difficulties lie in the accuracy, diversity and complexity of the generated caption.

Motivation

- Transformer has the advantage of parallelization, which is especially critical when generating text.
- Graph neural network (GCN) has excellent performance in extracting the relationship between objects.
- Curriculum learning can make the model train and learn from simple to difficult, and then improve the performance.

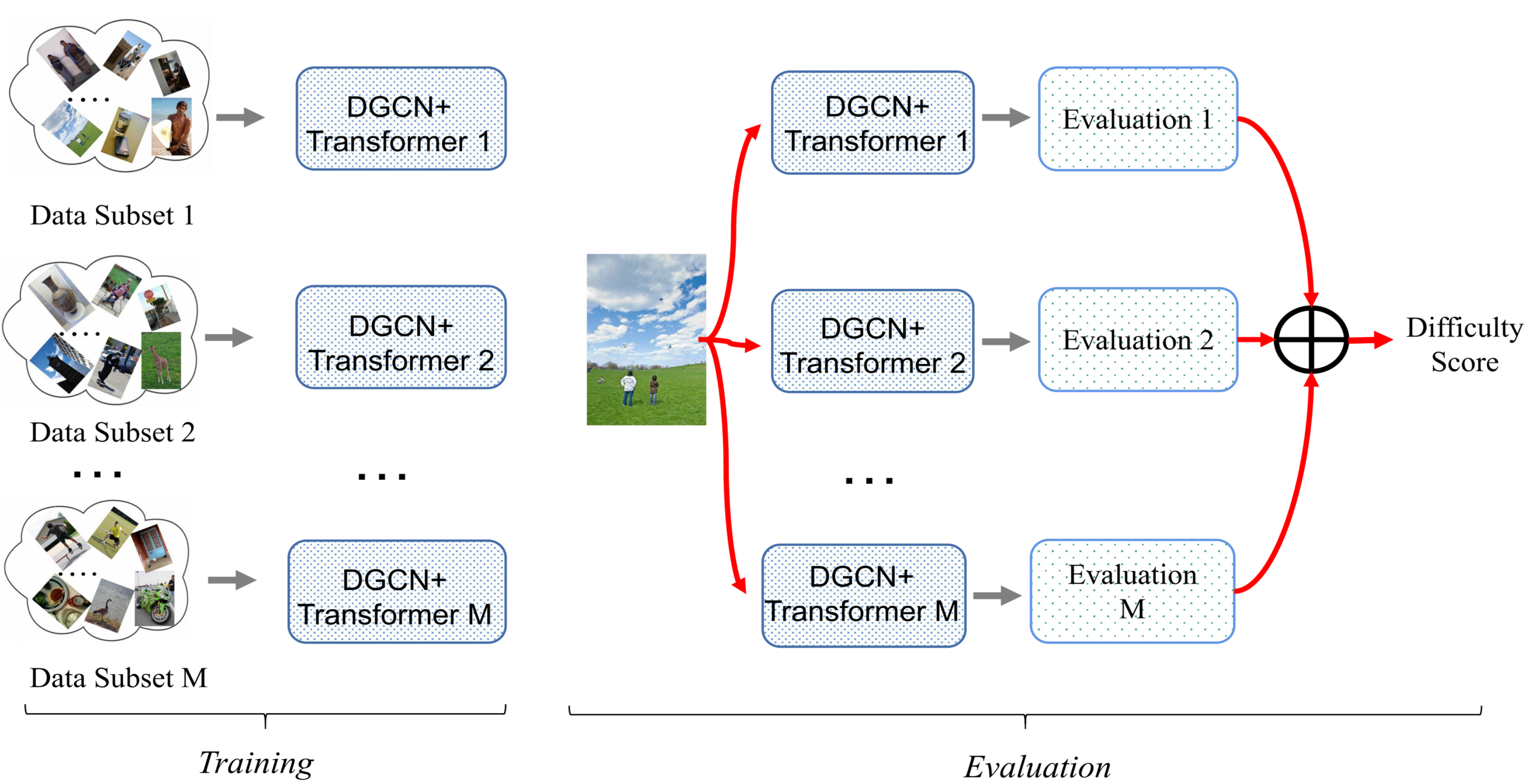
Proposed Method

Dual-GCN+Transformer Framework



- Our framework comprises three main parts: Dual-GCN (object-level and image-level) encoder, Transformer decoder, and Curriculum learning.

Curriculum Learning



- We use the cross-review mechanism used to determine the difficulty level for training examples:

$$DS(I, S) = \frac{1}{M-1} \sum_{(I, t) \in D_i, k \neq i} \epsilon_k(I, S)$$

Image I	Ground-truth Caption S	DS(I, S)
	<p>A cheesy pizza sitting on top of a table.</p> <p>This large pizza has a lot of cheese and tomato sauce.</p> <p>A prepared pizza being served at a table.</p> <p>A pizza with olive, cheese, tomato and onion.</p> <p>Baked pizza with herbs displayed on serving tray at table.</p>	0.214
	<p>A group of Army soldiers visit a memorial.</p> <p>A group of soldiers standing next to a bench.</p> <p>Soldiers performing a ceremony in a city park.</p> <p>Members of the military are in a park.</p> <p>Military officers standing on a public park in uniform.</p>	0.268

Dataset & Metrics & Setting

- Dataset: the Microsoft COCO Dataset.
- Metrics: BLEU-n, METEOR, ROUGE, and CIDEr.
- K=6, M=8 in our experiments.



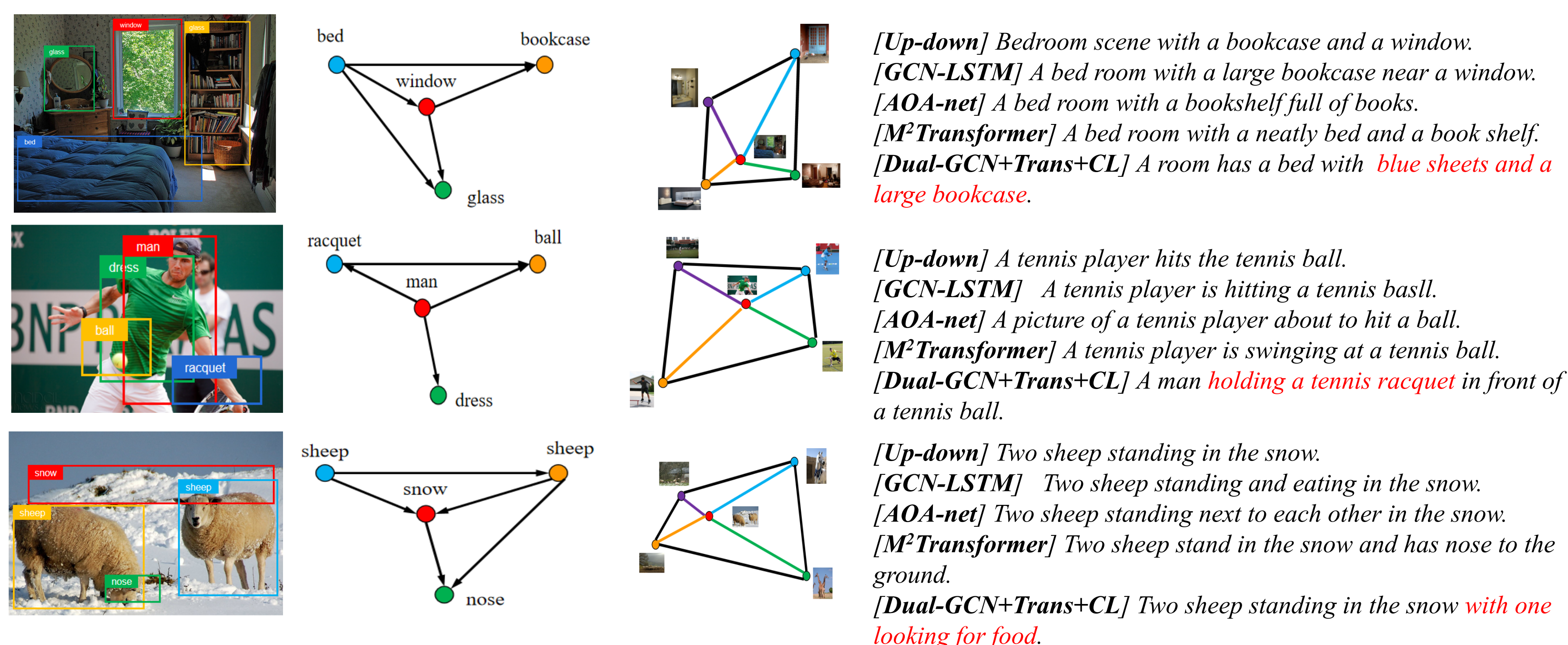
Our source code is available here. Just scan it!

Comparisons with State-of-the-Art

Quantitative Results

	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE	CIDEr
Up-down [1]	80.2	64.1	49.1	36.9	27.6	57.1	117.9
AOA-NET [21]	81.0	65.8	51.4	39.4	29.1	58.9	126.9
GCN-LSTM [53]	80.8	65.5	50.8	38.7	28.5	58.5	125.3
GCN-LSTM+HIP [54]	81.6	66.2	51.5	39.3	28.8	59.0	127.9
SGAE [50]	81.0	65.6	50.7	38.5	28.2	58.6	123.8
M ² Transformer [5]	81.6	66.4	51.8	39.7	29.4	59.2	127.9
Dual-GCN+Transformer+CL	82.2	67.6	52.4	39.7	29.7	59.7	129.2

Visualization Results

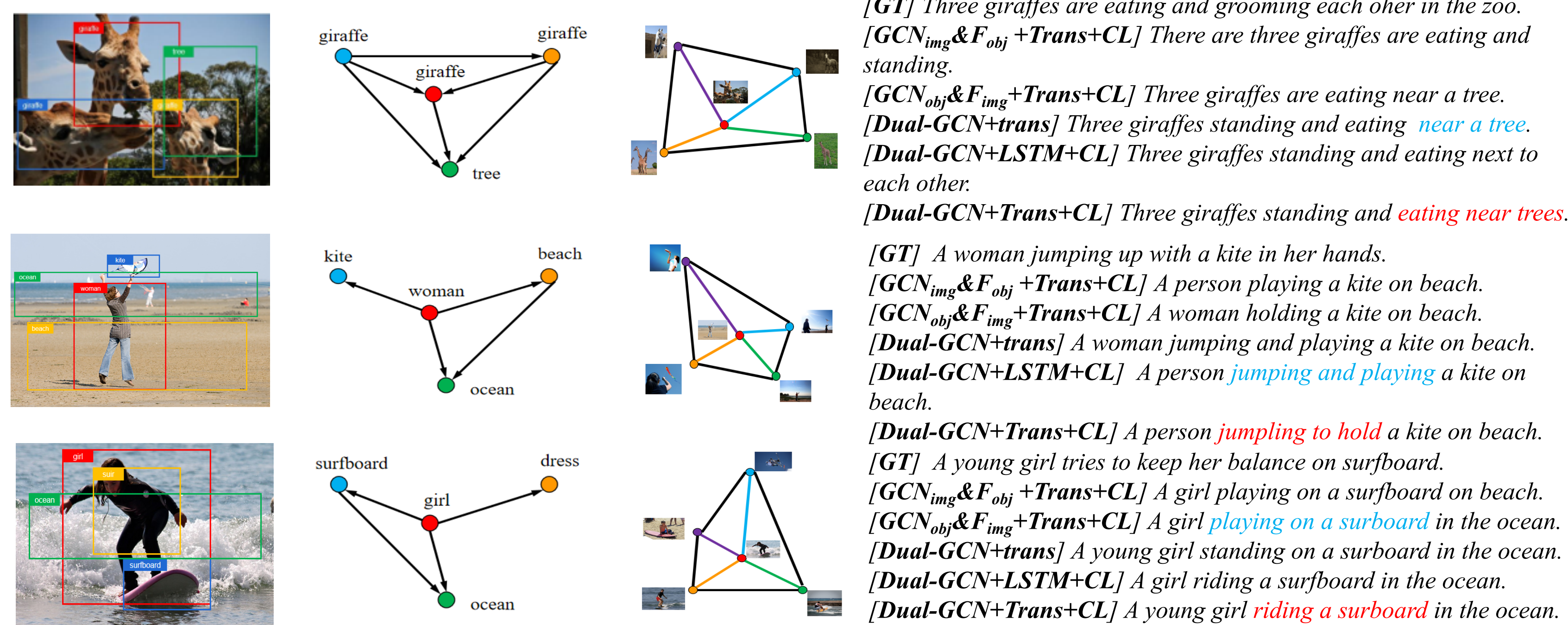


Abation Study

Quantitative Results

	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE	CIDEr
F _{obj} +Transformer	71.1	55.9	45.3	28.0	25.2	52.8	114.3
GCN _{obj} +Transformer	80.5	64.7	49.2	37.6	27.9	57.4	125.8
F _{obj} +Transformer+CL	72.3	56.8	46.1	28.4	25.4	53.0	114.7
GCN _{obj} +Transformer+CL	81.6	66.8	51.6	39.3	28.3	58.4	129.4
F _{img} +Transformer+CL	80.2	64.1	49.0	37.4	27.4	57.0	124.0
GCN _{img} +Transformer+CL	82.0	66.9	52.2	39.8	29.6	59.4	129.0
GCN _{obj} &F _{img} +Transformer+CL	81.9	67.3	52.0	39.5	29.4	59.6	129.4
GCN _{img} &F _{obj} +Transformer+CL	82.0	67.3	52.3	39.7	29.8	59.6	129.0
Dual-GCN+Transformer	81.8	66.2	51.6	39.2	28.7	58.9	128.0
Dual-GCN+LSTM+CL	81.4	65.8	51.3	38.8	28.0	58.4	127.6
Dual-GCN+Transformer+CL	82.2	67.6	52.4	39.7	29.7	59.7	129.2

Visualization Results



Failure Cases



Conclusion

- A novel framework that takes Dual-GCN as a powerful visual encoder and transformer as the linguistic decoder for image captioning.
- To our best knowledge, we are the first one to apply the Curriculum Learning as the training strategy on image captioning model in an easy-to-hard manner.
- We conduct extensive experiments on the MS COCO dataset. The results strongly demonstrate that our proposed approach outperforms state-of-the-art methods.