



# *Look Twice as Much as You Say: Scene Graph Contrastive Learning for Self-Supervised Image Caption Generation*

**Chunhui Zhang<sup>1</sup>, Chao Huang<sup>2</sup>, Youhuan Li<sup>3</sup>, Xiangliang Zhang<sup>4</sup>, Yanfang Ye<sup>4</sup>, and Chuxu Zhang<sup>1</sup>**  
Brandeis University,<sup>1</sup> University of Hong Kong,<sup>2</sup> Hunan University,<sup>3</sup> University of Notre Dame<sup>4</sup>

# Motivation

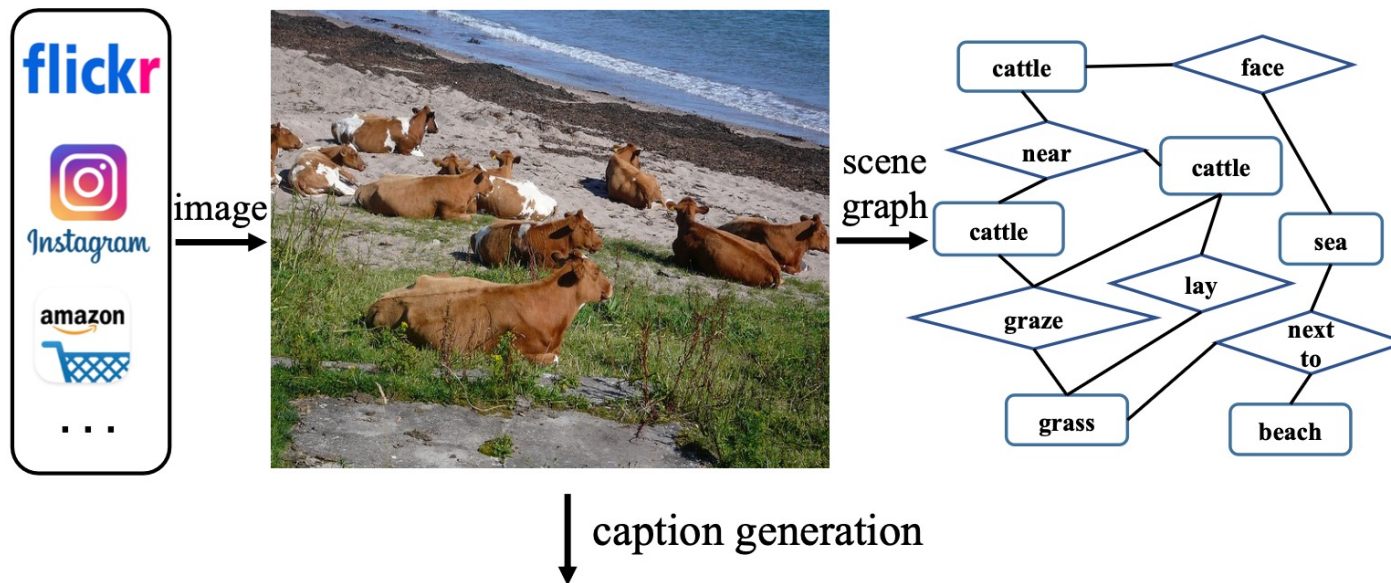
1. Background: Improving image captioning via leveraging unlabeled images from unpaired web sources.

2. Challenge:

- Cross-modal data - Unlike previous studies working on single- modal data (e.g., image, text, or graph), image caption generation is a cross-modal task on the intersection of image and text;
- Complex task - Image caption generation is a complex task that has to generate new content rather than simple classification or prediction task studied in previous work.

# Motivation

3. Target:



**GT: A herd of cattle laying on top of a sandy beach.**

---

**1% labels are used:**

C-GAT: A group standing a a a a.

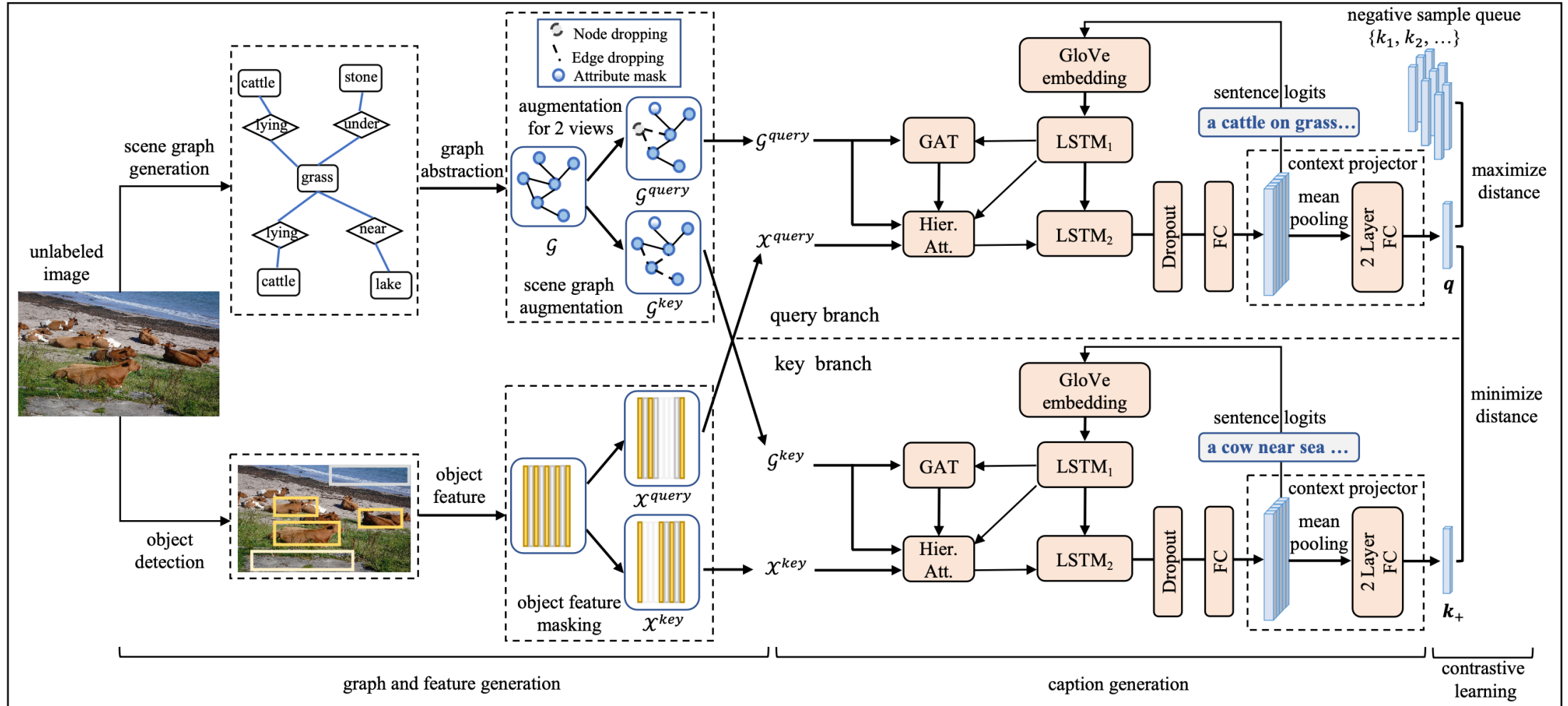
VSUA: A cattle a a a a.

SGAE: A cow is a a a a.

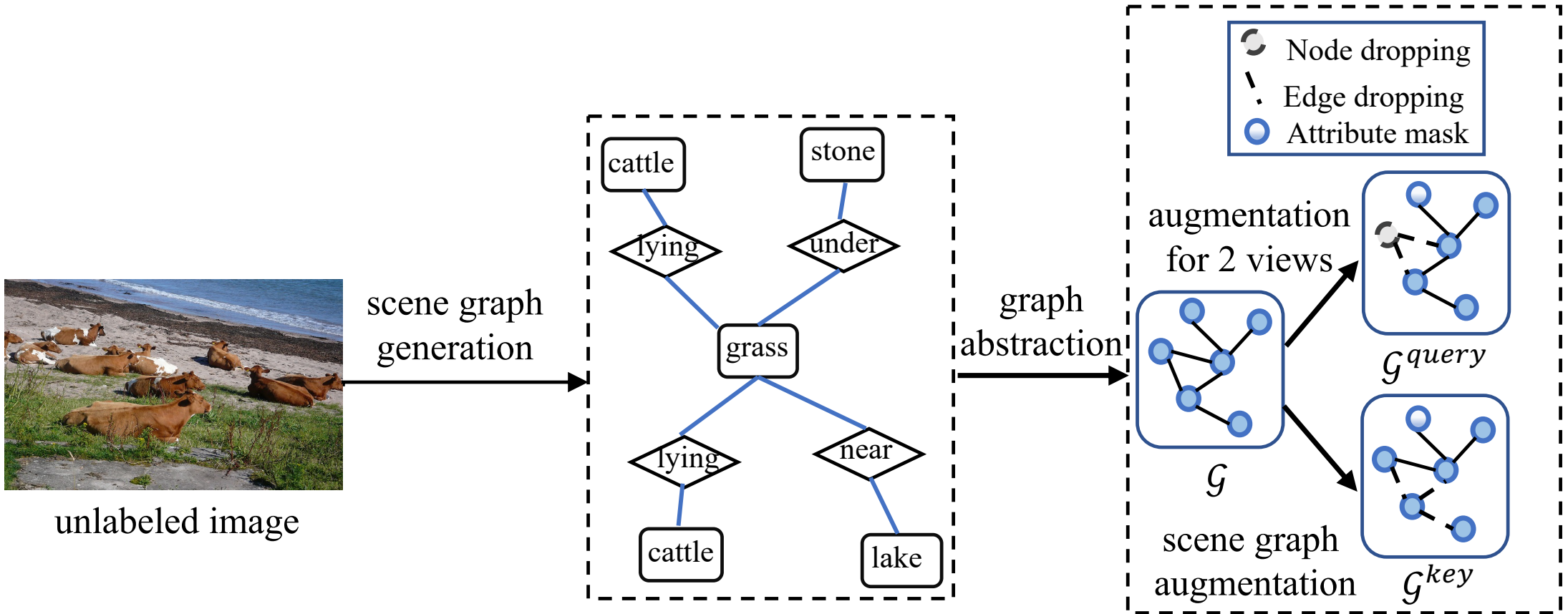
$M^2$  -T: A herd of sheep standing a a a a a a.

**SGCL: A group of cattle grazing in the beach in front of the water.**

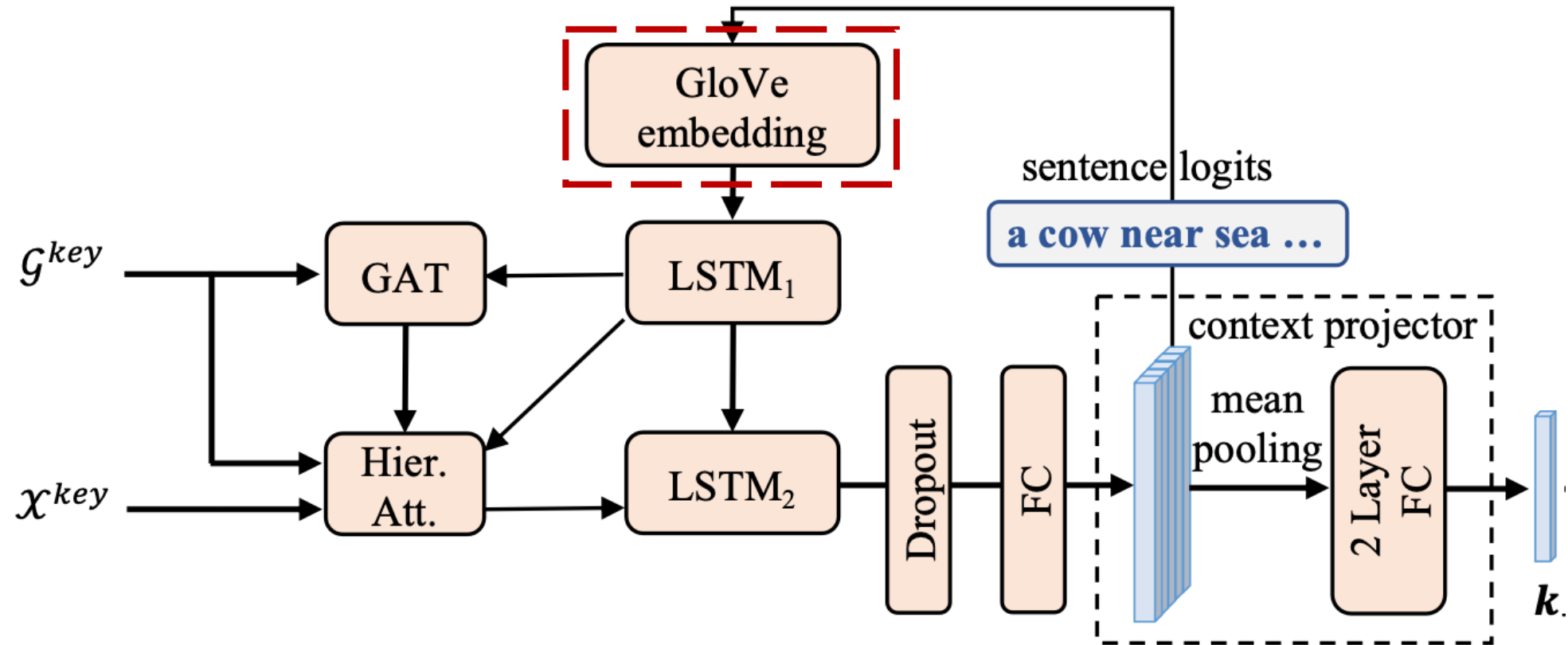
# The pipeline of Our Model



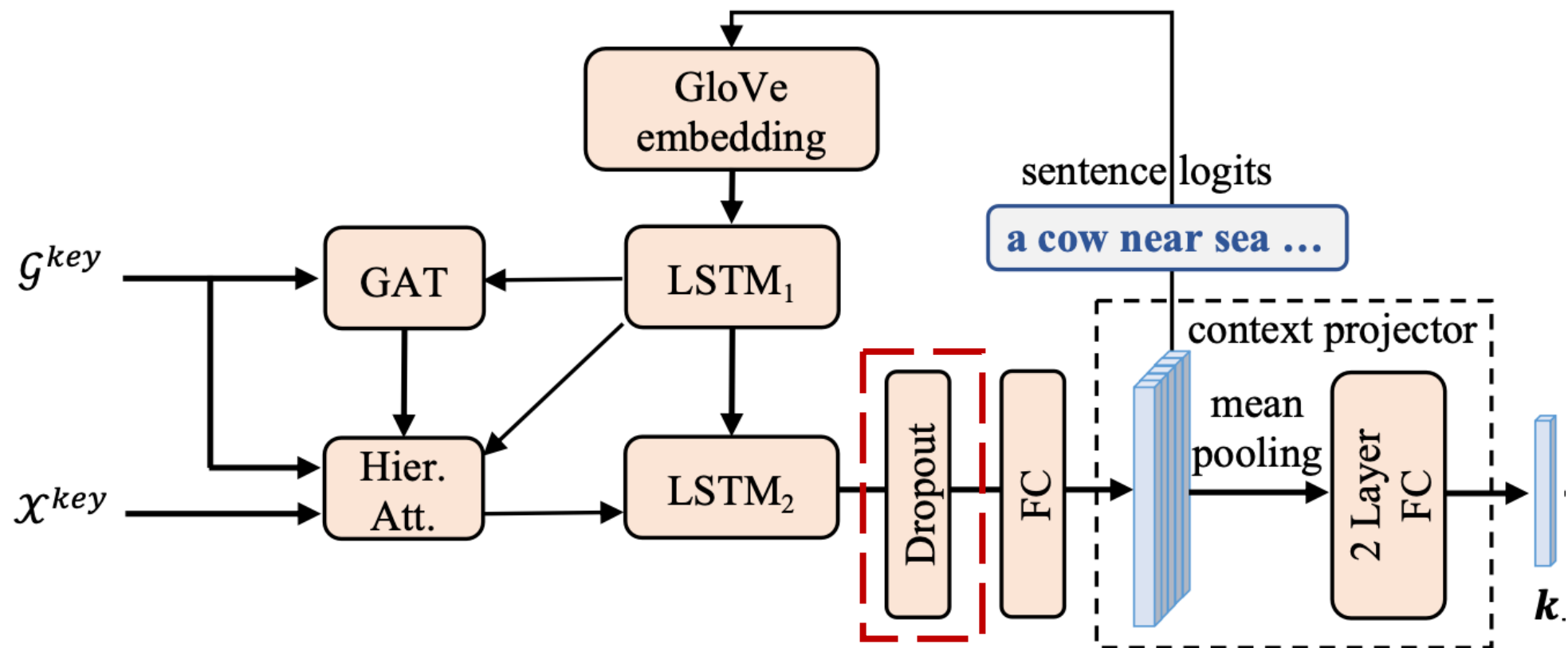
# Scene Graph Augmentation



# Pretrained Word Embedding *for NLP Information*

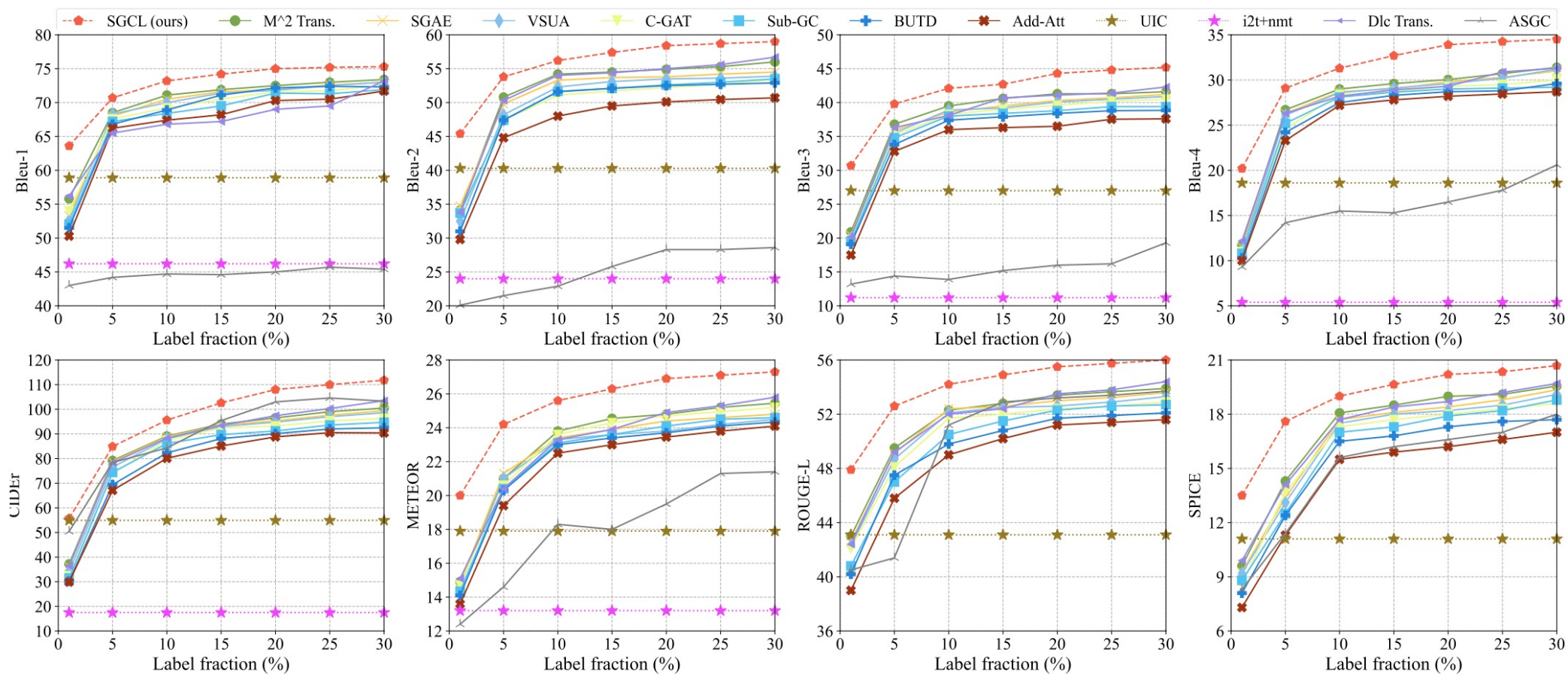


# Big Dropout Rate as Semantics Augmentations





# Experiment



**Figure 3: Performances of all models with limited labels (Note that ROUGE-L and SPICE of i2t+nmt are not shown due to missing values in the original work).**



# Ablation Study

**Table 1: Performances of different model variants with various graph augmentation strategies (Note: N - node dropping, E - edge dropping, A - node attribute masking, O - object feature masking).**

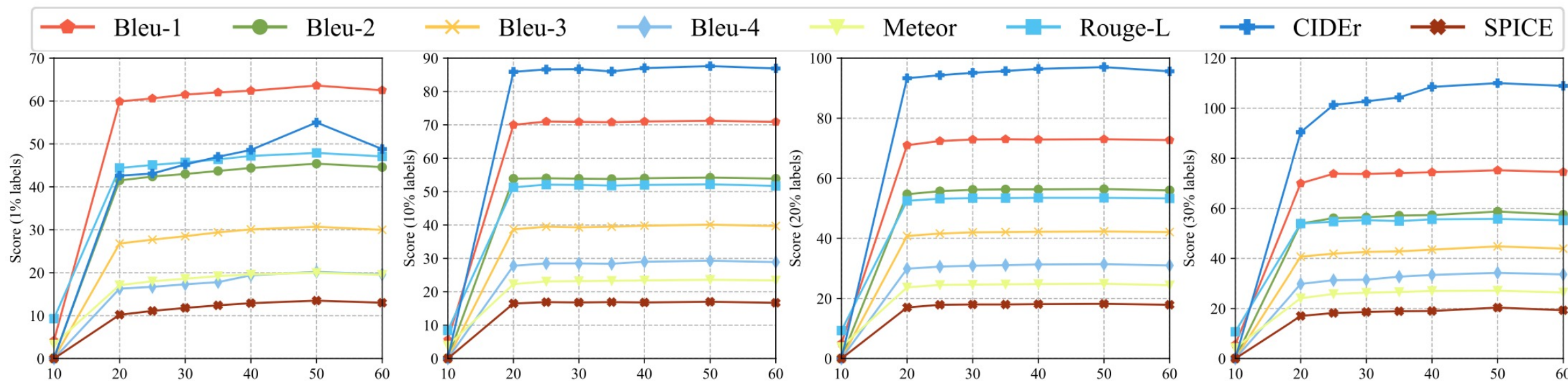
Label	N	E	A	O	B-1	B-2	B-3	B-4	C.	M.	R-L	S.
1%					61.8	43.8	28.9	18.2	47.7	18.5	46.3	11.9
				✓	62.5	44.6	30.0	19.1	49.2	19.9	47.0	13.1
	✓			✓	62.5	44.5	29.0	18.5	52.9	19.1	47.3	13.2
		✓		✓	63.1	44.3	28.8	18.6	52.2	19.3	47.2	13.0
			✓	✓	63.0	45.1	29.9	19.3	53.3	19.6	47.5	13.3
	✓	✓	✓	✓	<b>63.6</b>	<b>45.4</b>	<b>30.7</b>	<b>20.2</b>	<b>55.0</b>	<b>20.0</b>	<b>47.9</b>	<b>13.5</b>
5%					69.4	51.7	36.6	26.2	75.9	22.2	49.4	16.3
				✓	70.3	53.0	38.6	27.9	79.4	23.9	51.9	17.3
	✓			✓	62.5	52.8	38.9	28.5	81.4	19.1	51.9	17.2
		✓		✓	63.1	53.5	38.7	28.6	82.2	19.3	52.0	17.1
			✓	✓	70.3	53.3	39.2	28.1	82.3	24.1	52.2	17.4
	✓	✓	✓	✓	<b>70.7</b>	<b>53.8</b>	<b>39.8</b>	<b>29.1</b>	<b>84.9</b>	<b>24.2</b>	<b>52.6</b>	<b>17.6</b>

# Ablation Study

**Table 2: Effectiveness of loading pre-trained word embedding (P) and freezing word embedding (F).**

Label	P	F	B-1	B-2	B-3	B-4	C.	M.	R.-L	S.
1%			62.9	44.8	30.1	19.8	54.1	19.3	47.0	13.0
	✓		63.0	45.1	30.2	19.9	54.4	19.6	47.3	13.0
	✓	✓	<b>63.6</b>	<b>45.4</b>	<b>30.7</b>	<b>20.2</b>	<b>55.0</b>	<b>20.0</b>	<b>47.9</b>	<b>13.5</b>
5%			69.8	53.2	38.8	28.3	80.0	23.6	52.1	16.9
	✓		70.1	53.6	39.0	28.3	79.8	23.8	52.3	17.1
	✓	✓	<b>70.7</b>	<b>53.8</b>	<b>39.8</b>	<b>29.1</b>	<b>84.9</b>	<b>24.2</b>	<b>52.6</b>	<b>17.6</b>
10%			71.8	58.3	41.1	33.5	91.7	23.5	52.9	17.5
	✓		72.4	58.3	41.5	33.7	92.6	23.6	53.4	18.3
	✓	✓	<b>73.2</b>	<b>56.2</b>	<b>42.1</b>	<b>31.3</b>	<b>94.6</b>	<b>25.6</b>	<b>54.2</b>	<b>19.0</b>
20%			74.1	57.8	43.2	33.1	103.9	26.4	54.5	19.6
	✓		74.6	58.0	43.5	33.5	105.5	26.4	55.1	20.0
	✓	✓	<b>75.0</b>	<b>58.4</b>	<b>44.3</b>	<b>33.9</b>	<b>108.0</b>	<b>26.9</b>	<b>55.5</b>	<b>20.2</b>

# Ablation Study



**Figure 4: Impact of dropout rate at the output layer on model performance.**

# Case Show

image:



## 1% labels are used:

$M^2$ -T: A sheep standing in a a a.	A man girl a a a a a.	A girl tennis a a tennis tennis	A elephant of in a a a.	A cow standing standing a a a.
SGAE: A sheep of standing a a a.	A man is a a a a a.	A man girl a a a a a.	A man is a a a a a.	A sheep of in a a a a a.
VSUA: A sheep of in a a a.	A man standing a a a a a.	A girl girl a a a a a.	A people of a a a a a a.	A cow cow cow a a a a a.
C-GAT: A group of a a a a.	A man is a a a a a.	A man girl a a a a a.	A street of a a a a.	A sheep of a a a.
<b>SGCL: A couple of sheep standing in the grass in a field.</b>	<b>A group of people playing a frisbee standing by the sea.</b>	<b>A woman is holding a tennis racket on a tennis ball.</b>	<b>A group of people standing in a street with a building.</b>	<b>A herd of cows walking down a road in the grass.</b>

## 30% labels are used:

$M^2$ -T: A group of sheep standing in a fenced area.	Two men playing frisbee on a dirt field.	A woman is holding a tennis racket in her hand.	A man riding an elephant in front of a building.	A group of cows standing next to each other on a field.
SGAE: A white sheep is standing in the grass.	A group of men playing a game of frisbee.	A man hitting a tennis ball on a tennis court.	A group of people standing next to an elephant.	A cow standing on top of a lush green field.
VSUA: A couple of sheep standing next to each other.	A man holding a frisbee in his hand.	Two men playing frisbee on a dirt field.	An elephant standing in front of a building.	A group of cows are standing in the grass.
C-GAT: A group of sheep grazing in a grassy field.	A man holding a frisbee in his hand.	A woman is playing tennis on the court.	A man riding on the back of an elephant.	A brown cow standing next to a brown cow.
<b>SGCL: A couple of sheep standing on a lush green field near a fence.</b>	<b>A man is jumping in the air to catch a frisbee on a sea beach.</b>	<b>A woman is trying to hitting a tennis ball on a tennis court.</b>	<b>An elephant walking down a street with people in the background.</b>	<b>A herd of black cows standing next to each other on a lush green field.</b>



*Thank you for your listening.*