

# Modular Graph Transformer Networks for Multi-Label Image Classification

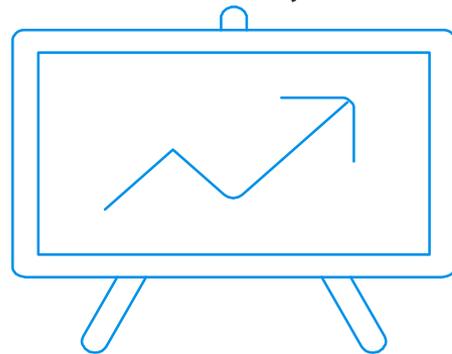
Hoang D.Nguyen, Xuan-Son Vu, Duc-Trong Le

@AAAI 2021

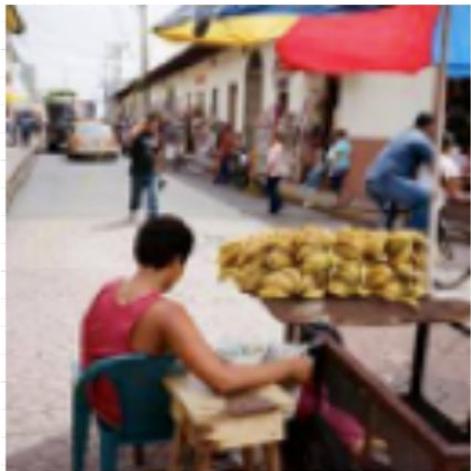


# Outline

- ① Multi-Label Image Classification
- ① Motivating Example
- ① Modular Graph Transformer Network (MGTN)
- ① Experiments



# Multi-Label Image Classification



**Predict**

**Person, Chair,  
Umbrella, Car**

# Motivating Example (1)



person, chair, umbrella,  
car



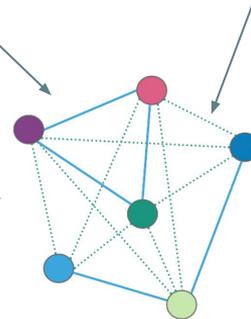
person, chair, dining table,  
cup



person, umbrella



dining table, cup, bowl



**Semantic Graph between Labels**

# Motivating Example (2)



person, chair, umbrella,  
car



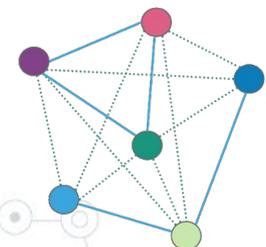
person, chair, dining table,  
cup



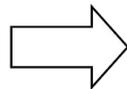
person, umbrella



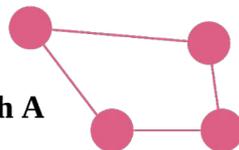
dining table, cup, bowl



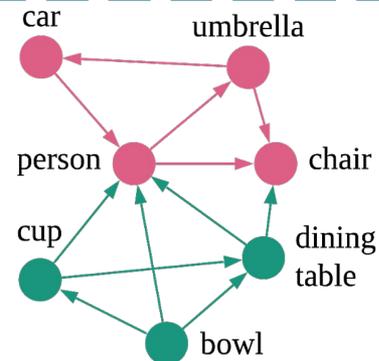
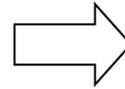
**Graph-Information**



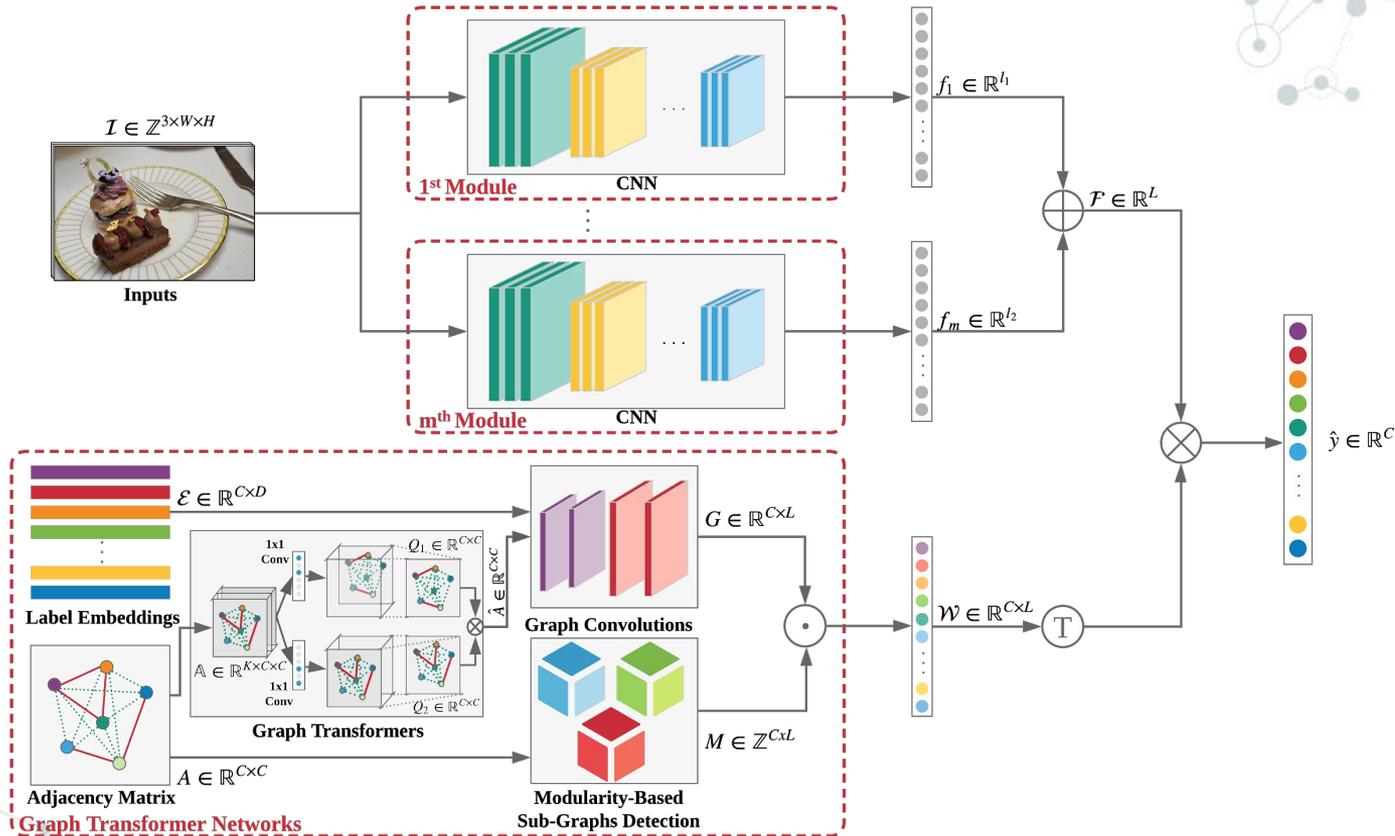
**Sub-graph A**



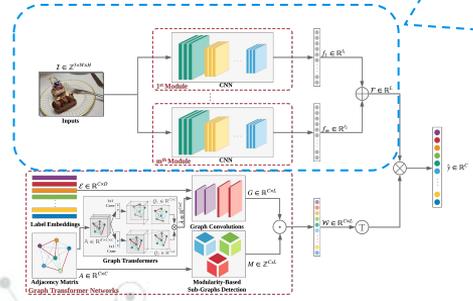
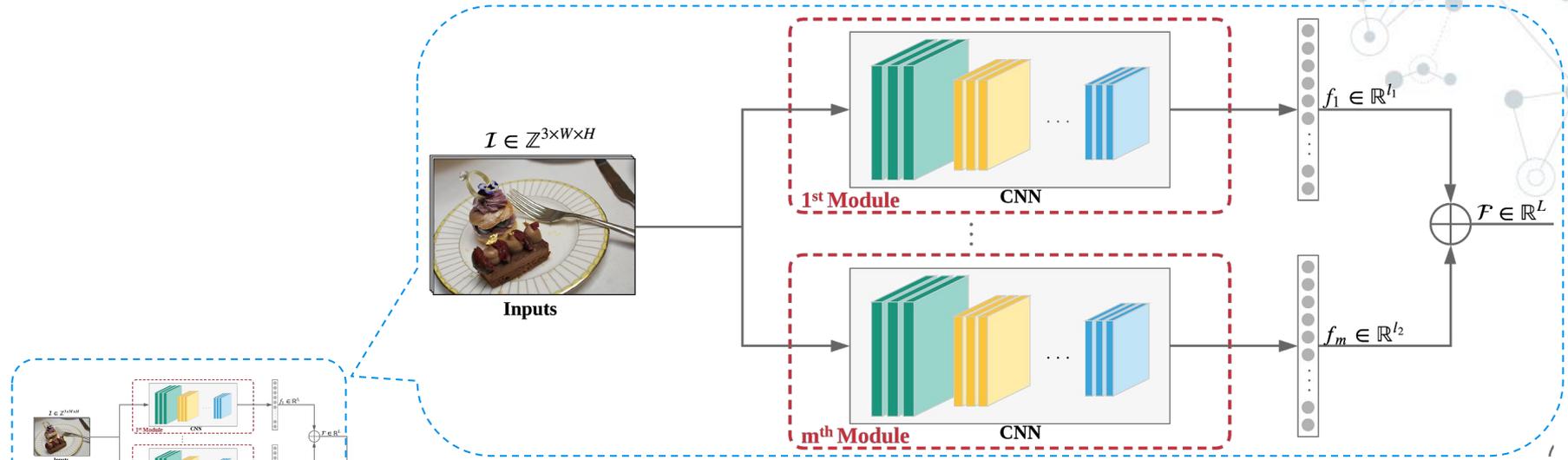
**Sub-graph B**



# Modular Graph Transformer Network (MGTN)

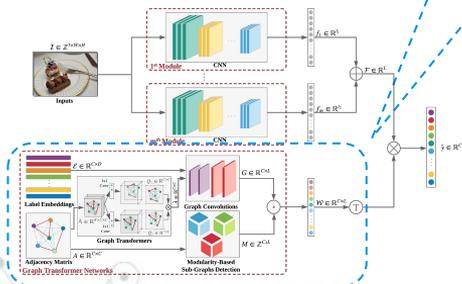
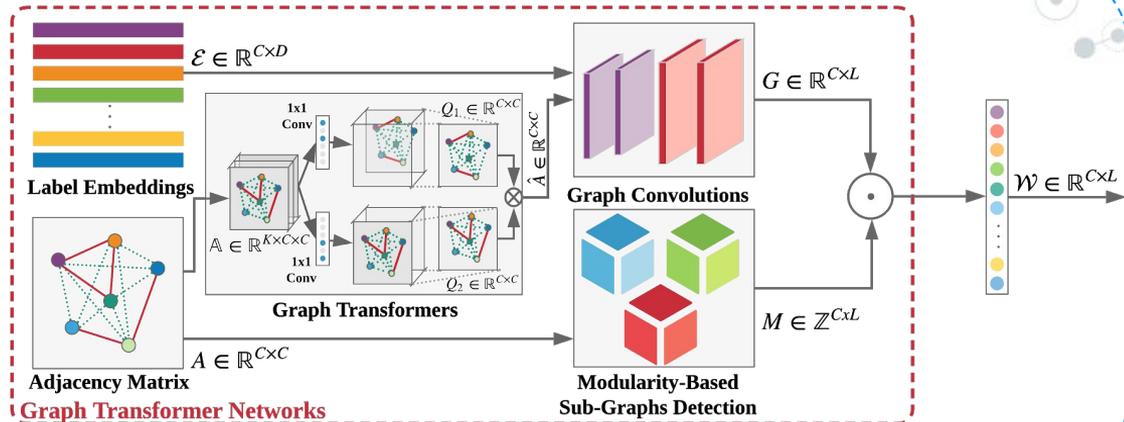


# Convolutional Neural Networks



**Idea:** Exploit multiple modules of CNNs to infer the image-label representation  $F \in \mathbb{R}^L$  for each image input  $I \in \mathbb{Z}^{3 \times W \times H}$

# Graph Transformer Networks



**Idea:** Exploit semantic information ( $\mathcal{E} \in \mathbb{R}^{C \times D}$ ) & topological information ( $A \in \mathbb{R}^{C \times C}$ ) between the set of  $C$  labels

# Graph Transformers

**Idea:** Dynamically select important topological connections of the label graph.

Adjacency Tensor

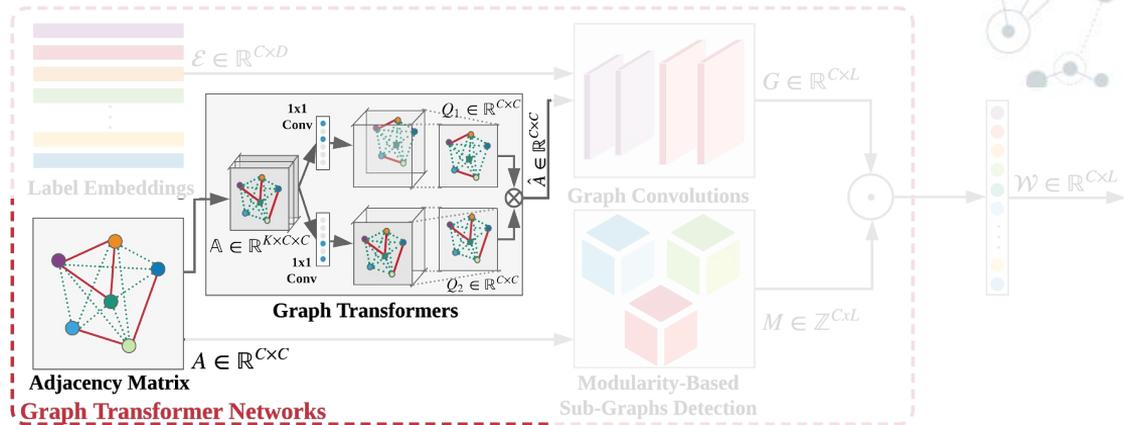
$$\mathbb{A}_{kij} = \begin{cases} 1 & \text{if } P_{ij} \in [t_{k-1}, t_k), i \neq j \\ 0 & \text{otherwise} \end{cases}$$

where  $\mathbb{A}_1 \equiv I$

$$P_{ij} = \varrho * A_{ij} / d_i$$

$$d_i = \sum_k A_{i,k}$$

$$t_k \in [0, 1], k \in \{1, \dots, K\}$$



Two  $1 \times 1$  convolutions  $Q_1, Q_2 \in \mathbb{R}^{C \times C}$

(Yun et al. 2019)

$$Q_1 = \phi(\hat{A}, \text{softmax}(W_\phi^1))$$

$$W_\phi^1, W_\phi^2 \in \mathbb{R}^{1 \times 1 \times K}$$

$$Q_2 = \phi(\hat{A}, \text{softmax}(W_\phi^2))$$

$$\hat{A} \in \mathbb{R}^{C \times C}$$

$$\hat{A} = \eta(Q_1 Q_2 + I)$$

where

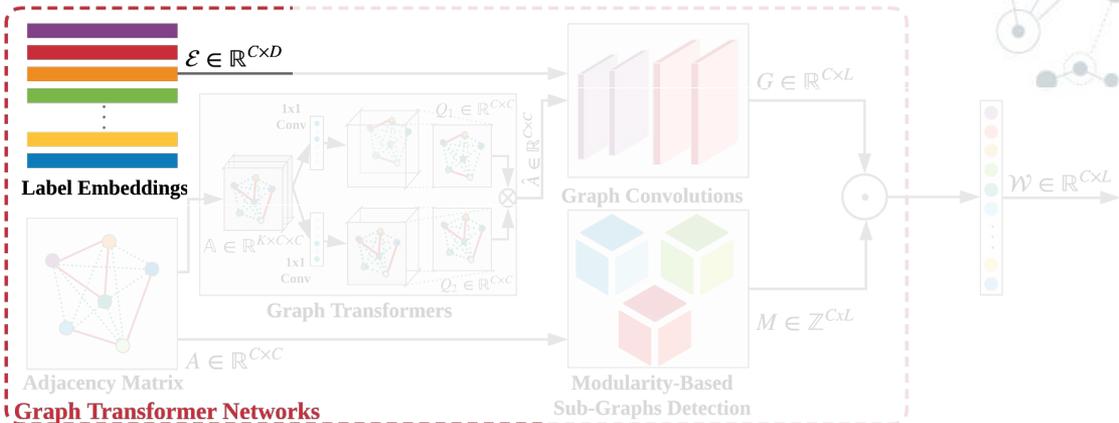
$$\eta(A) = d^{-\frac{1}{2}} A d^{-\frac{1}{2}}$$

# Eigenvector-based Embedding Transformation

**Idea:** Exploit semantic information from label embeddings

Pre-Trained Embeddings:

- 1) **Char2Vec** (Kim et al. 2015): character-level,  $D = 300$
- 2) **BERT\_Base** (Devlin et al. 2018): word-level, averaging the last layer,  $D = 768$
- 3) **RoBERTa** (Liu et al. 2019): word-level, averaging the last layer,  $D = 768$



Eigenvector-based transformation for pre-trained embedding  $E$

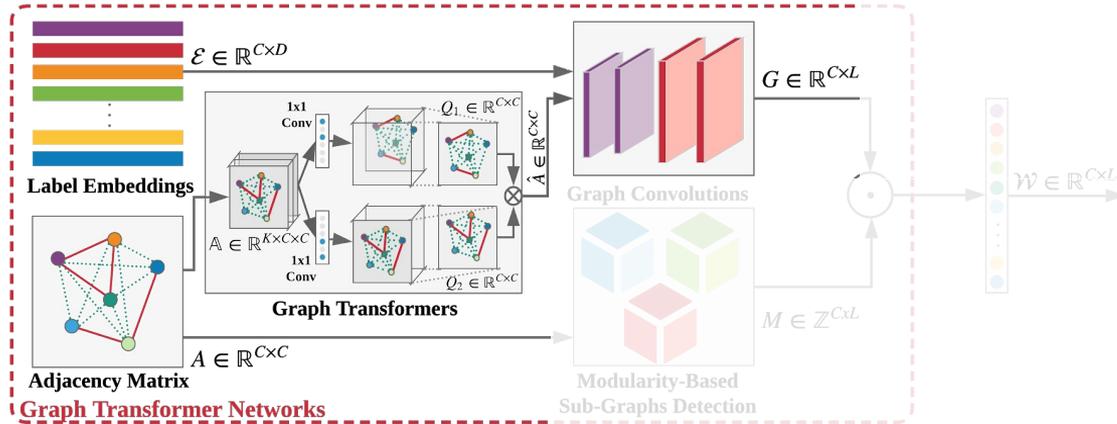
$$\mathcal{E} = E \cdot C^T$$

where  $C_i$  the eigenvector centrality of the label  $i$ -th

$$C_i = \frac{1}{\lambda} \sum_k a_{k,i} C_k$$

$\lambda \neq 0$  is the largest eigenvalue

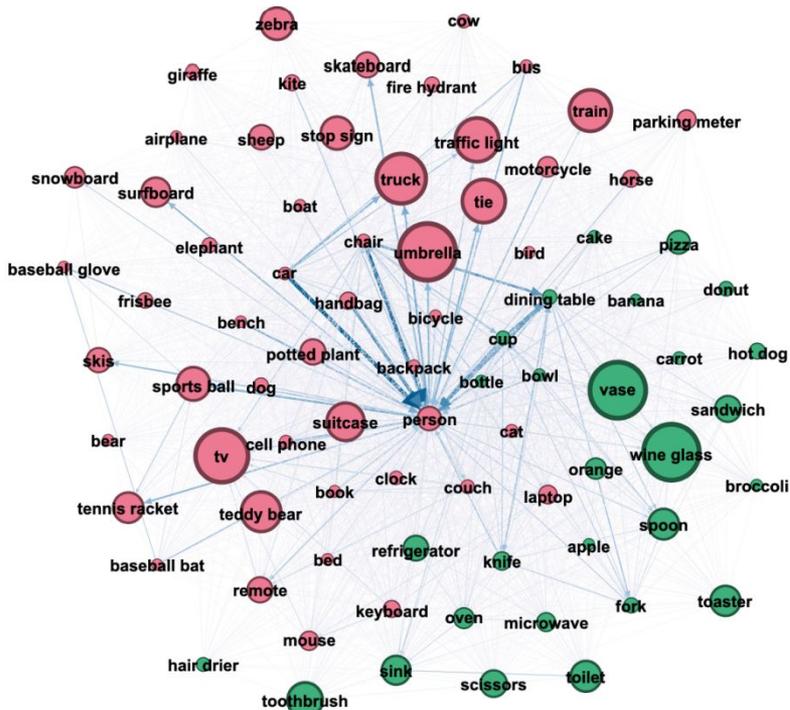
# Graph Convolutional Network



**Idea:** Jointly exploit label-level word embedding and topological information via graph convolutional networks (GCN, Kipf et al. 2017)

$$G = \text{GCN}(\mathcal{E}, \hat{A}) \quad G \in \mathbb{R}^{C \times L}$$

# Modularity-based Sub-Graphs Detection



MS-COCO

The Clauset-Newman-Moore agglomeration algorithm (Clauset et al. 2004)

# Modularity-based Enhancement

**Idea:** Exploit different highly inter-connected sets of objects among sub-graphs

The modularity of a graph :

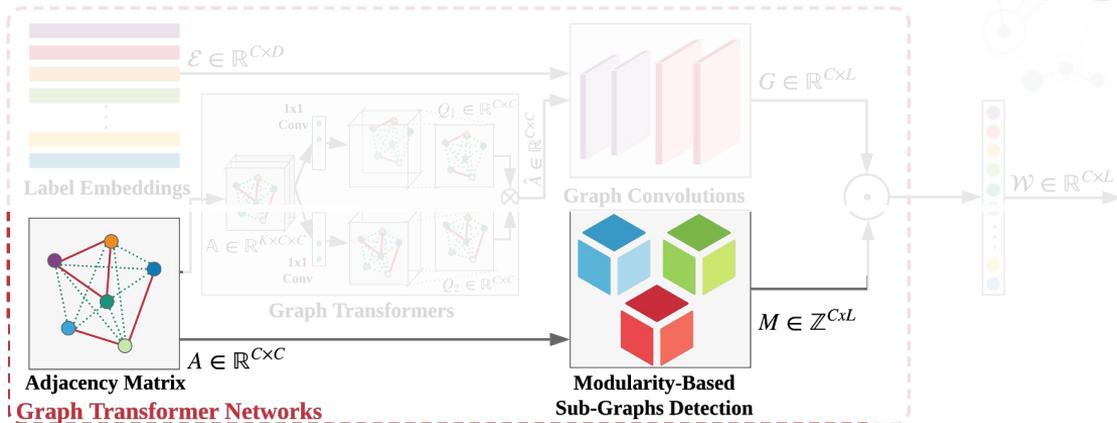
$$\Omega = \frac{1}{2m} \sum_{i,j} \left[ A_{i,j} - \frac{d_i d_j}{2m} \right] \delta(c_i, c_j)$$

where

$$m = \frac{1}{2} \sum_{i,j} A_{i,j}$$

$$d_i = \sum_k A_{i,k}$$

$\delta(u, v)$  is 1 if  $u = v$  otherwise 0

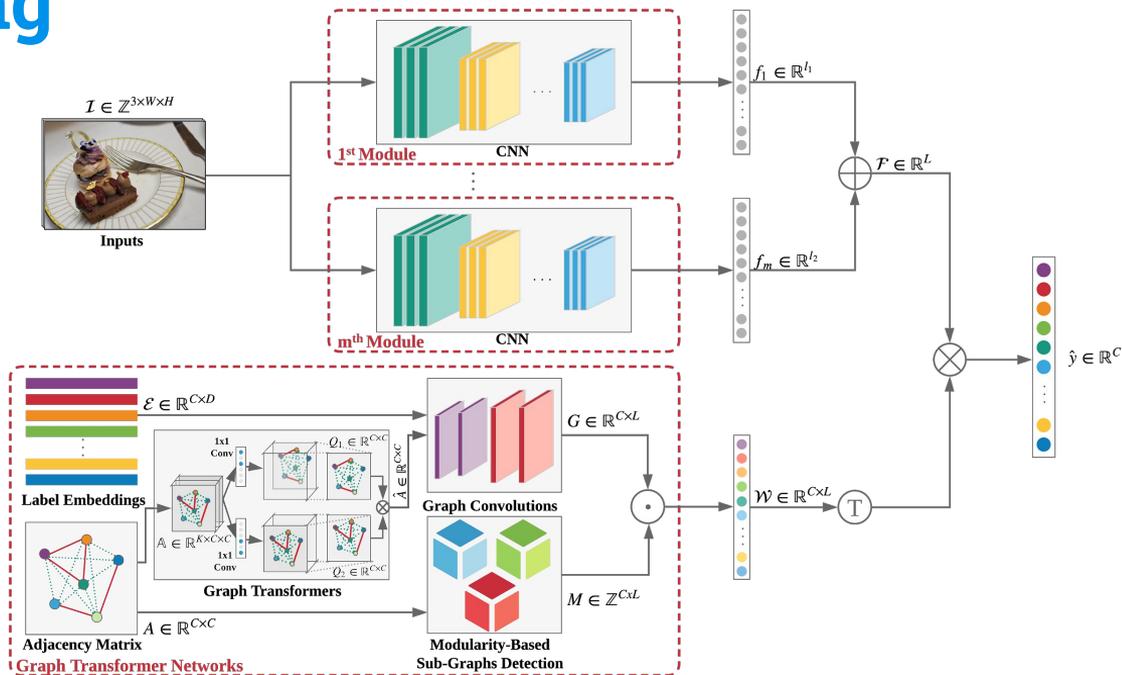


Suppose that  $m$  sub-graphs are discovered. Multiple CNNs with configurable backbones are employed using a control tensor  $M$  with a threshold  $\tau$  :

$$M_{iv} = \begin{cases} \tau & \text{if } S_i = p \text{ and } v \in f_p \\ \frac{1-\tau}{m-1} & \text{otherwise} \end{cases}$$

where  $S$  is the sub-graph assignment,  $f_p \in \mathbb{R}^{l_p}$  is the image-level representation

# Learning



$$W = G \odot M \quad \hat{y} = W^T F$$

$$\mathcal{L} = -\frac{1}{C} \sum_{c=1}^C y_c \log(\sigma(\hat{y}_c)) + (1 - y_c) \log(1 - \sigma(\hat{y}_c))$$

# Experimental Setups

- ◎ **Dataset:** MS-COCO (81 labels), Fashion550K (66 labels)
- ◎ **Evaluation metrics:** mAP, per-class (CP, CR, CF1), overall (OP, OR, OF1)
- ◎ **Preprocessing:** Resize images 512x512 to 448x448 with random horizontal flips
- ◎ **Implementation:** Dual ResNeXt-50-32x4d backbones, 2-layer GCN,  $\tau = 0.999$ , adjacency thresholds (MS-COCO = [0.1, 0.3, 1.0] , Fashion550K = [0.2, 0.4, 1.0]) learning rate decays by a factor of 10 for every 20 epochs.

# Experimental Results - MS-COCO

METHOD	MAP	CP	CR	CF1	OP	OR	OF1
CNN-RNN (WANG ET AL. 2016)	61.2	-	-	-	-	-	-
SRN (ZHU ET AL. 2017)	77.1	81.6	65.4	71.2	82.7	69.9	75.8
BASELINE(RESNET101) (HE ET AL. 2016)	77.3	80.2	66.7	72.8	83.9	70.8	76.8
MULTI-EVIDENCE (GE, YANG, AND YU 2018)	-	80.4	70.2	74.9	85.2	72.5	78.4
ML-GCN (CHEN ET AL. 2019B)	82.4	84.4	71.4	77.4	85.8	74.5	79.8
ML-GCN (RESNEXT50 WITH IMAGENET)	86.2	85.8	77.3	81.3	86.2	79.7	82.8
A-GCN (LI ET AL. 2019)	83.1	84.7	72.3	78.0	85.6	75.5	80.3
KSSNET (WANG ET AL. 2020B)	83.7	84.6	73.2	77.2	87.8	76.2	81.5
SGTN (OUR) (VU ET AL. 2020)	86.6	77.2	<b>82.2</b>	79.6	76.0	<b>82.6</b>	79.2
MGTN(BASE)	86.9	<b>89.4</b>	74.5	81.3	<b>90.9</b>	76.3	83.0
MGTN(FINAL)	<b>87.0</b>	86.1	77.9	<b>81.8</b>	87.7	79.4	<b>83.4</b>

# Model Analysis - Embeddings

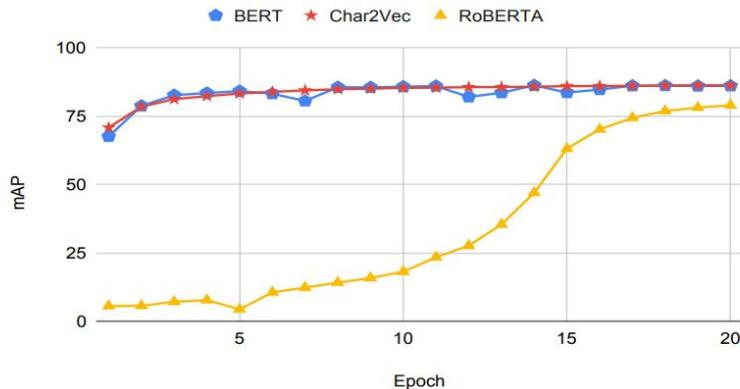


Figure 4: Learning patterns of MGTN with different label embeddings in 20 epochs. The MGTN model with the setting using RoBERTA<sub>avg\_12</sub> label embedding shows a slow learning speed in comparison to others.

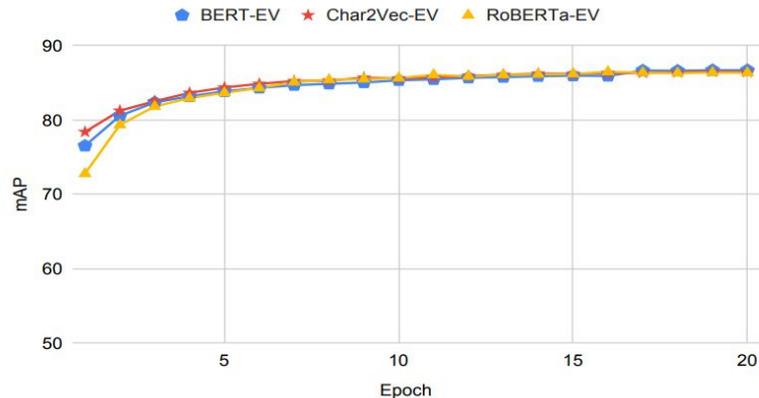
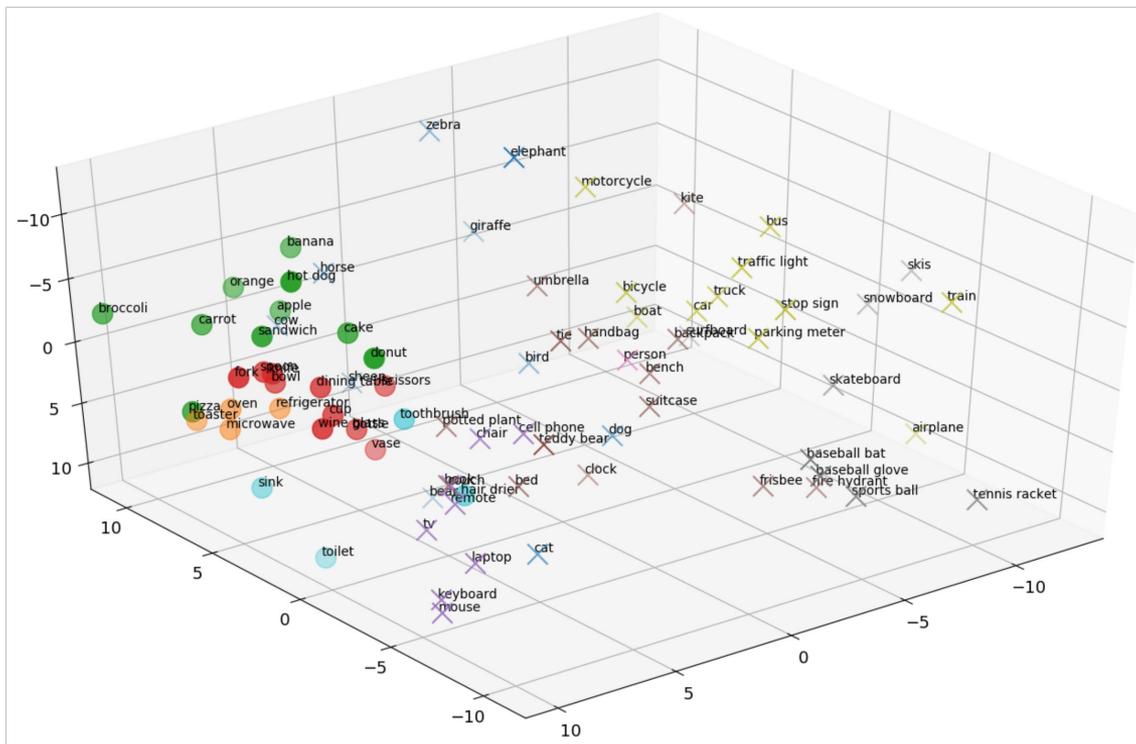


Figure 5: The EV-enhancement for label embedding helps the MGTN's model learn faster, even MGTN with the setting using the RoBERTA<sub>avg\_12</sub> now learns faster. Note: y-axis here is ranged in [50, 100] for visibility.

# 3D t-SNE Visualization on MS-COCO



# Conclusion

- ◎ Introduce **Modular Graph Transformer Network (MGTN)** to solve multi-label image classification
  - Employ **multiple CNN backbones** on unfolded **sub-graphs**
  - Exploit **topological and semantic information** via the graph transformer and EV-based embedding transformation.
- ◎ **MGTN** shows significant improvements against SOTA methods on MS-COCO and Fashion550K

# Thanks!

**Any questions?**