

TransForming Visual Scene Graphs to Image Captions

将场景图转换为图像文本

Xu Yang¹ Jiawei Peng¹ Zihua Wang¹ Haiyang Xu^{2*} Qinghao Ye²
Chenliang Li² Songfang Huang² Fei Huang² Zhangzikang Li¹ Yu Zhang^{1*}

1 东南大学计算机科学与工程学院、软件学院、人工智能学院-模式学习与挖掘实验室 (Palm)

2 阿里巴巴达摩院

Image captioning

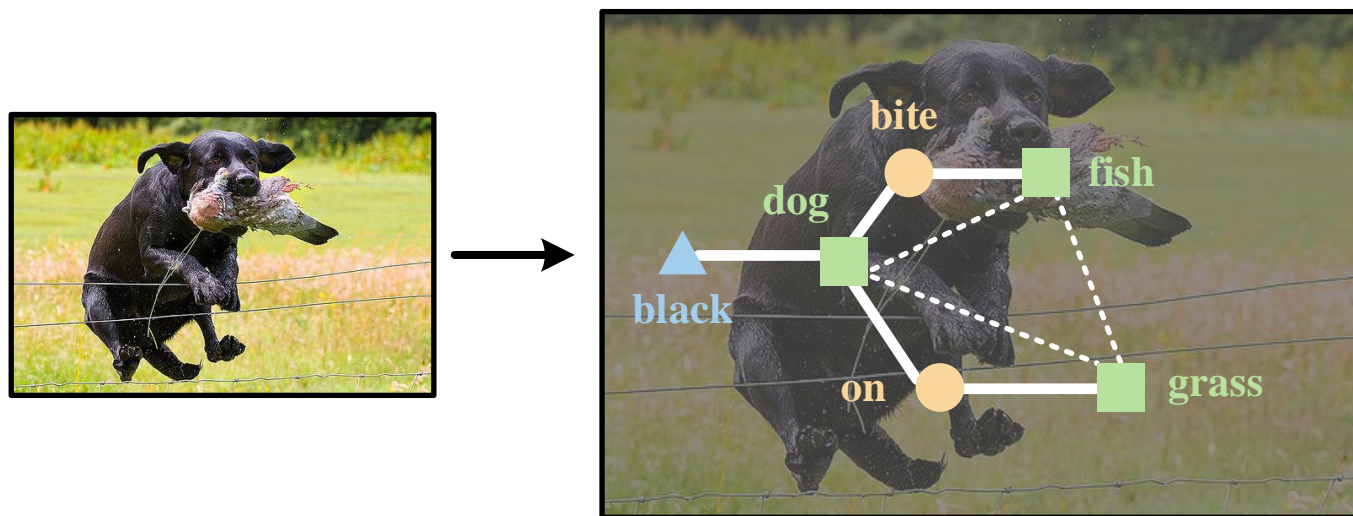
Generating a sentence to describe the image



A black dog biting a fish is running on the grass.

Image captioning with Scene Graphs

Structured information used to represent **objects**, **attributes**, and **relations** in an image



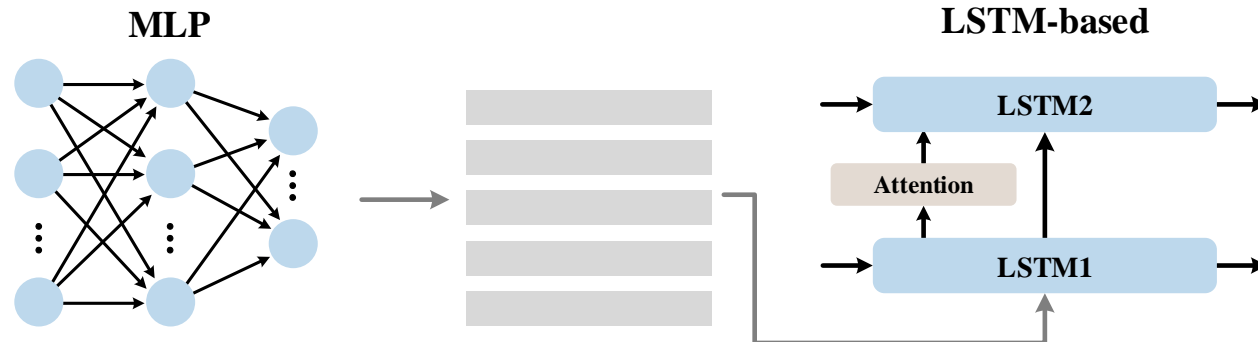
GNN-LSTM framework

Complex Training Strategies

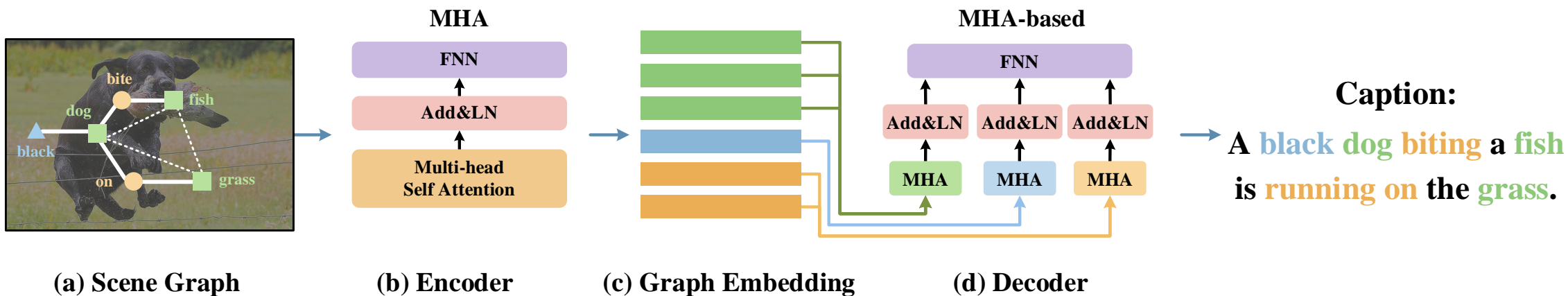
Heterogeneous structure Requires well-chosen training strategies.

Less descriptive

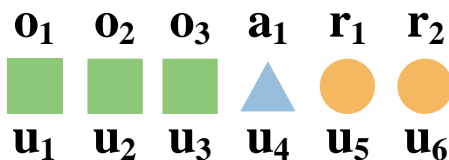
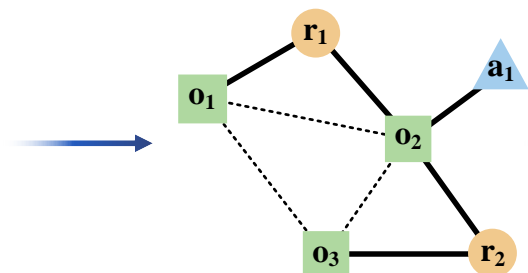
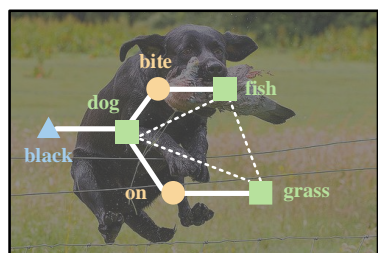
Embeddings are indiscriminately selected
Less descriptive captions



TransForming Scene Graphs to Captions (Ours)



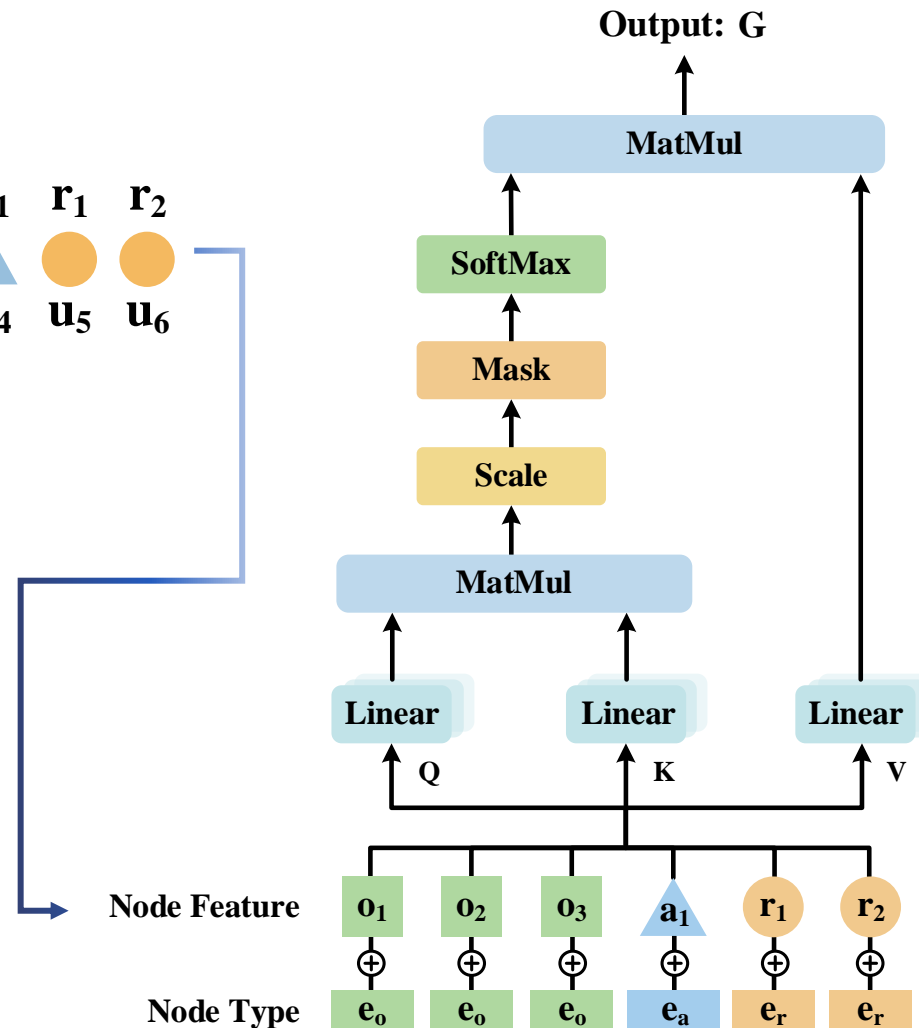
Input



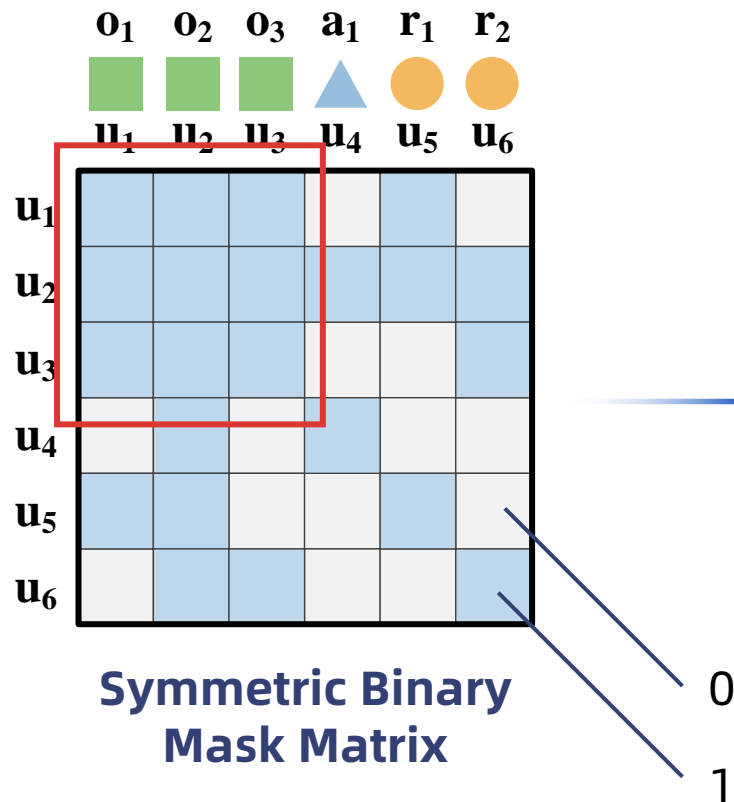
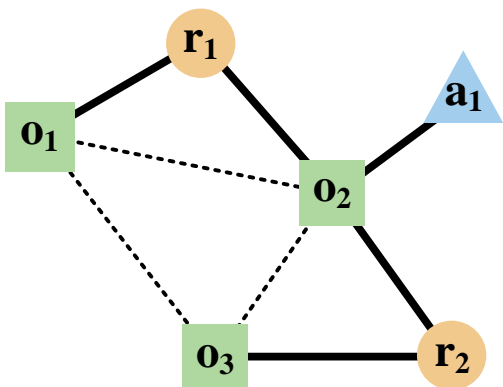
Object: $u_i = o_i + e_o, \quad 1 \leq i \leq N_o,$

Attribute: $u_{N_o+i} = a_i + e_a, \quad 1 \leq i \leq N_a,$

Relation: $u_{N_o+N_a+i} = o_i + e_r, \quad 1 \leq i \leq N_r,$

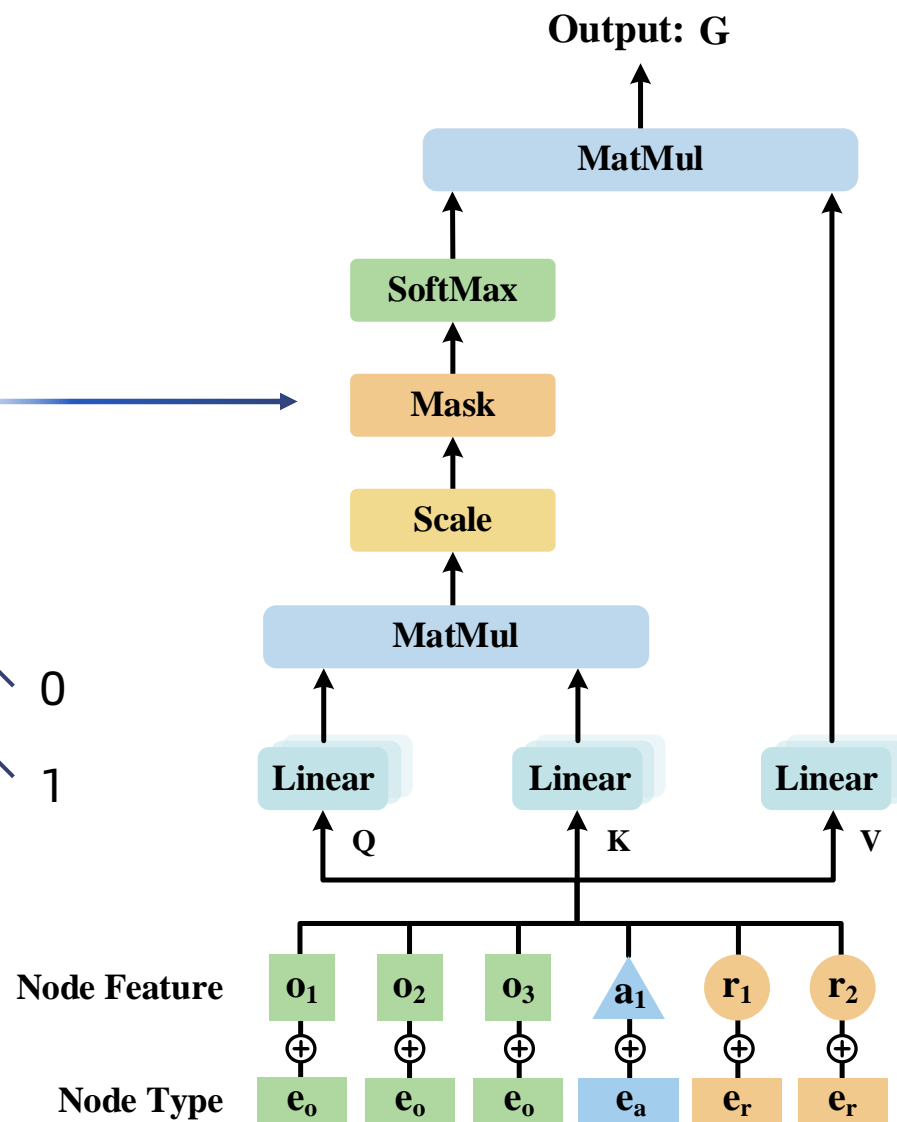


Encoder

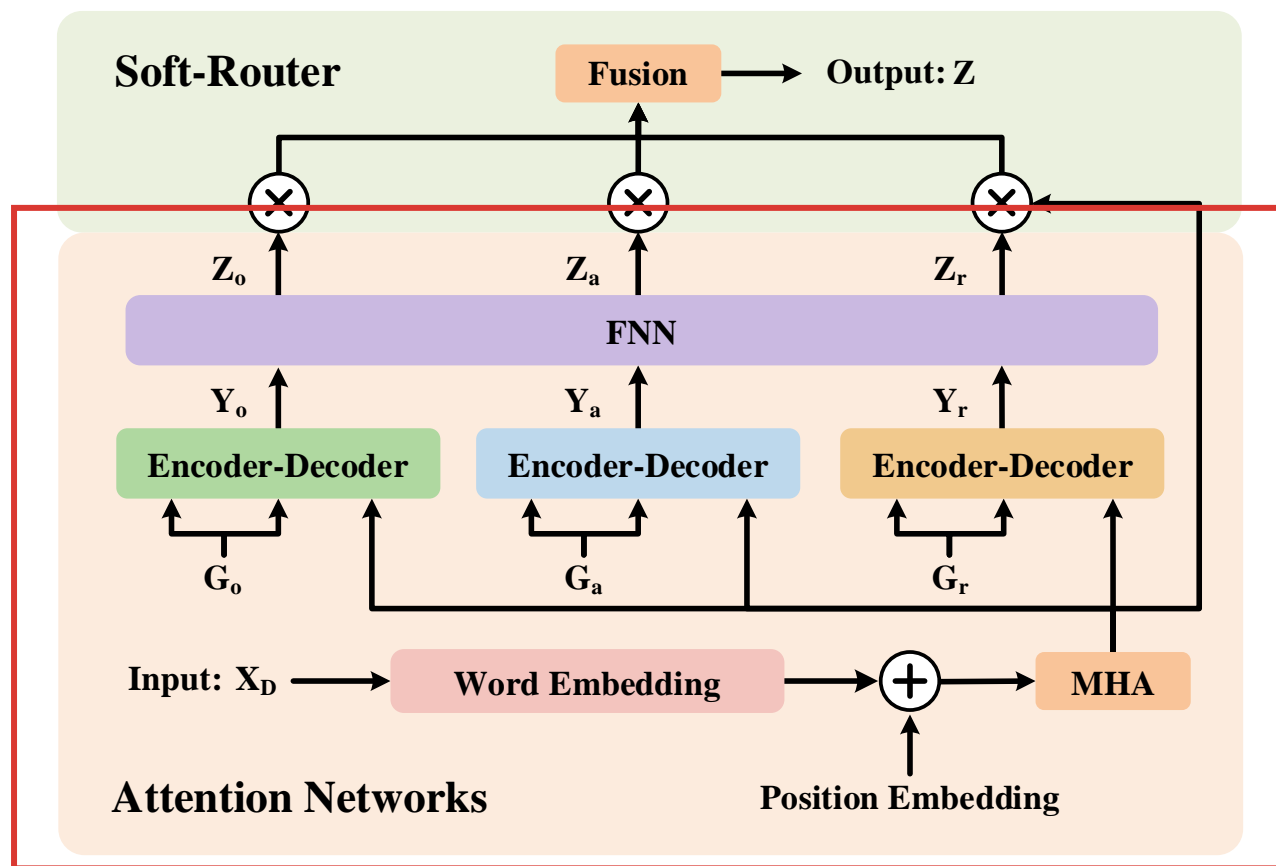


Masked:

- Control whether two nodes are connected or not
- Node can be updated by aggregating more neighbour embeddings



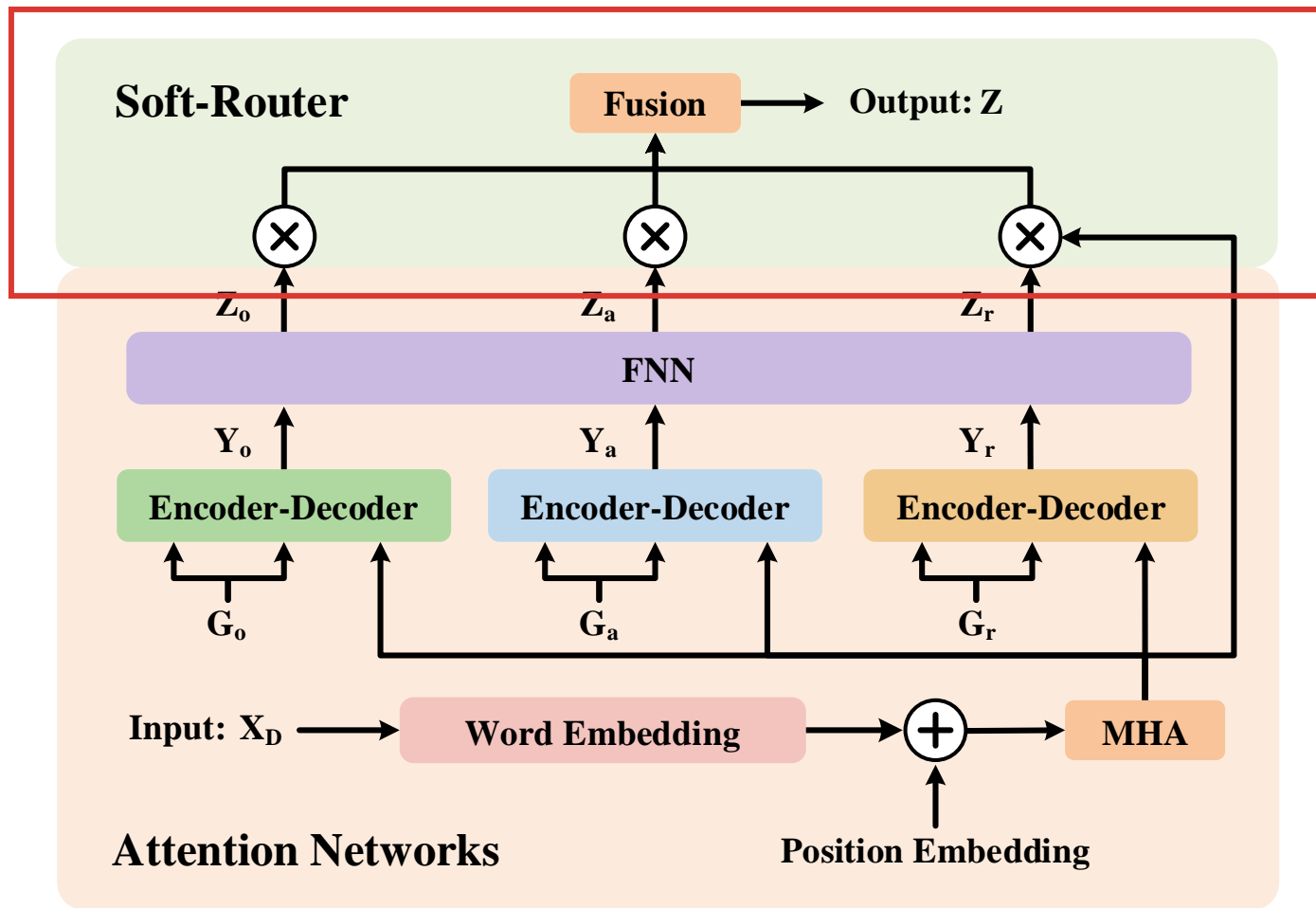
Decoder



Input: X_D, G_o, G_a, G_r
 SA: $X = \text{MHA}(Q = K = V = X_D)$,
 EXP_O: $Y_o = \text{MHA}(Q = X, K = V = G_o)$,
 EXP_A: $Y_a = \text{MHA}(Q = X, K = V = G_a)$,
 EXP_R: $Y_r = \text{MHA}(Q = X, K = V = G_r)$,
 FFN: $Z_o, Z_a, Z_r = \text{FFN}(Y_o, Y_a, Y_r)$,

The sketch of the MOE-decoder, we use different colors to denote different experts: green/blue/yellow correspond to object/attribute/relation experts.

Output



Input: $x, z_o, z_a, z_r,$

ATT: $\alpha = \{\alpha_o, \alpha_a, \alpha_r\}$
 $= \text{Softmax}(\{x^T z_o, x^T z_a, x^T z_r\})$

Output: $z = \alpha_o z_o + \alpha_a z_a + \alpha_r z_r,$

e.g.:

If the last word of generated caption is an **adjective** “black” , the next word is more like to be a **noun** and thus α_o should be a large value for using more **object embeddings** instead of the other embeddings.

TransForming Visual Scene Graphs to Image Captions

将场景图转换为图像文本

Thank you for listening!

Xu Yang¹ Jiawei Peng¹ Zihua Wang¹ Haiyang Xu^{2*} Qinghao Ye²
Chenliang Li² Songfang Huang² Fei Huang² Zhangzikang Li¹ Yu Zhang^{1*}