



Multi-Modal Sarcasm Detection Based on Relationship Dependence of Knowledge Graph



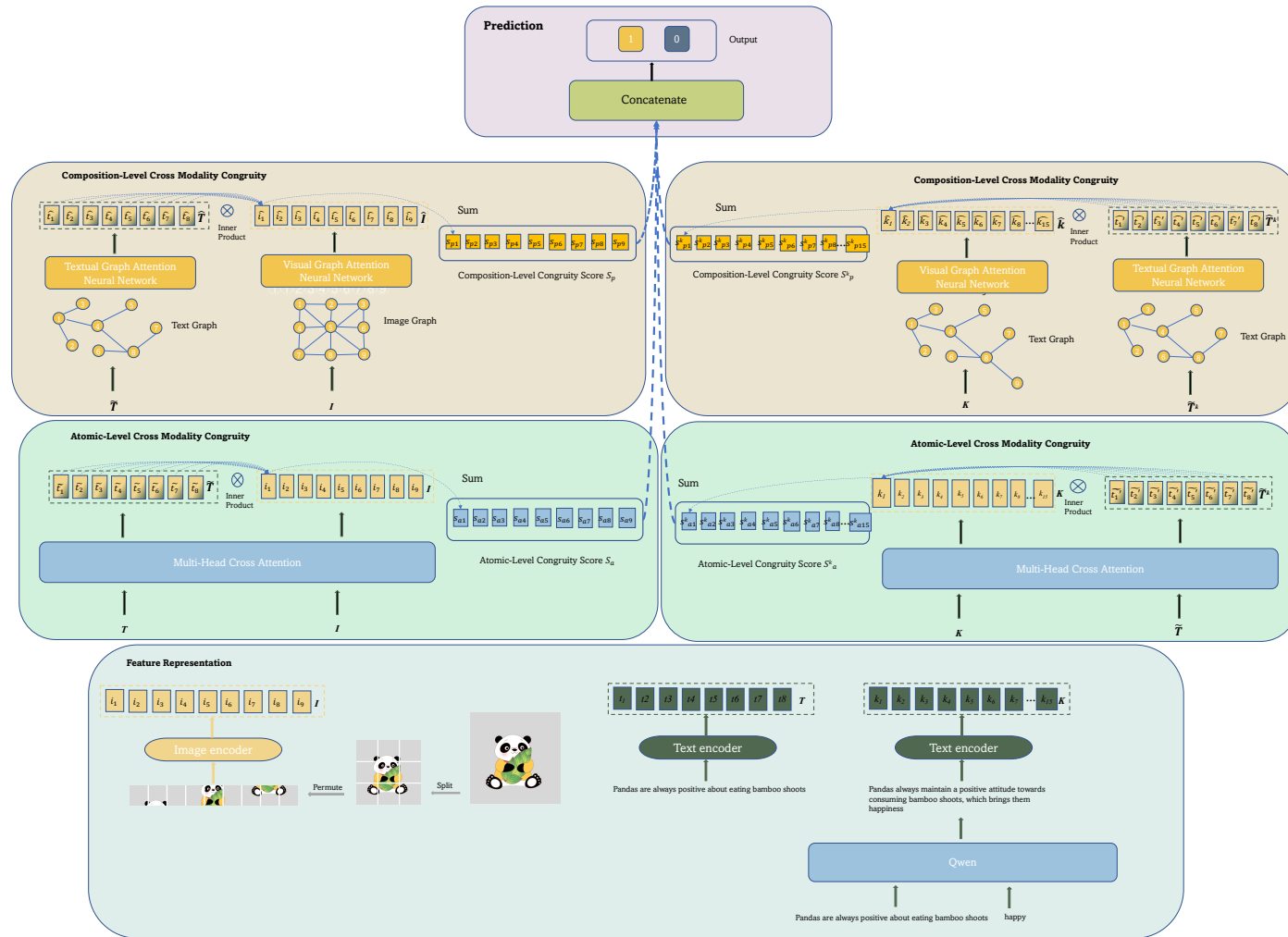
Reported by Danling.Wei

Motivation

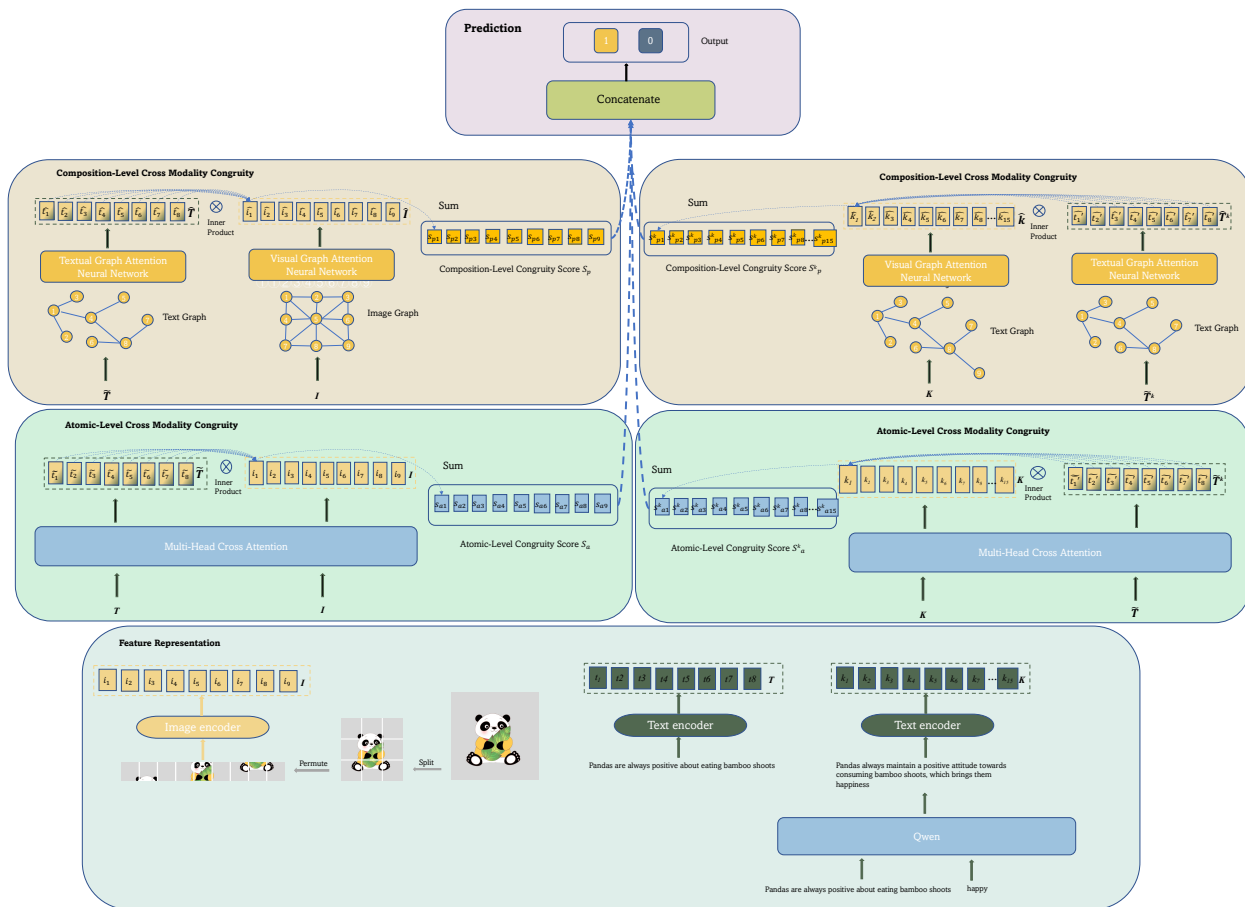


Figure 1: An example of sarcasm along with the corresponding image and different types of external knowledge extracted from the image. The sarcasm sentence represents the need for some good news. However, the image of the TV program is switched to bad news depicting severe storms (bad weather) which contradicts the sentence.

Overview



Method



$$\mathbf{head}_i = \text{softmax} \left(\frac{(\mathbf{T}\mathbf{W}_q^i)^\top}{\sqrt{d/h}} (\mathbf{I}\mathbf{W}_k^i) \right) (\mathbf{I}\mathbf{W}_v^i), \quad (1)$$

$$\tilde{\mathbf{T}} = \text{norm}(\mathbf{T} + \text{MLP}([\mathbf{head}_1 \parallel \mathbf{head}_2 \parallel \dots \parallel \mathbf{head}_h])), \quad (2)$$

$$\mathbf{Q}_a = \frac{1}{\sqrt{d}} (\tilde{\mathbf{T}}\mathbf{I}^\top)$$

$$\mathbf{s}_a = \text{softmax}(\tilde{\mathbf{T}}\mathbf{W}_a + \mathbf{b}_a)^\top \mathbf{Q}_a, \quad (3)$$

$$\alpha_{i,j}^l = \frac{\exp(\text{LeakyReLU}(\mathbf{v}_l^\top [\Theta_l \mathbf{t}_i^l \parallel \Theta_l \mathbf{t}_j^l]))}{\sum_k \exp(\text{LeakyReLU}(\mathbf{v}_l^\top [\Theta_l \mathbf{t}_i^l \parallel \Theta_l \mathbf{t}_k^l]))}, \quad (4)$$

$$\mathbf{t}_i^{l+1} = \alpha_{i,i}^l \Theta_l \mathbf{t}_i^l + \sum_{j \in \mathcal{N}(i)} \alpha_{i,j}^l \Theta_l \mathbf{t}_j^l, \quad (5)$$

$$\mathbf{c} = \text{softmax}(\mathbf{T}\mathbf{W}_c + \mathbf{b}_c)^\top \tilde{\mathbf{T}}, \quad (6)$$

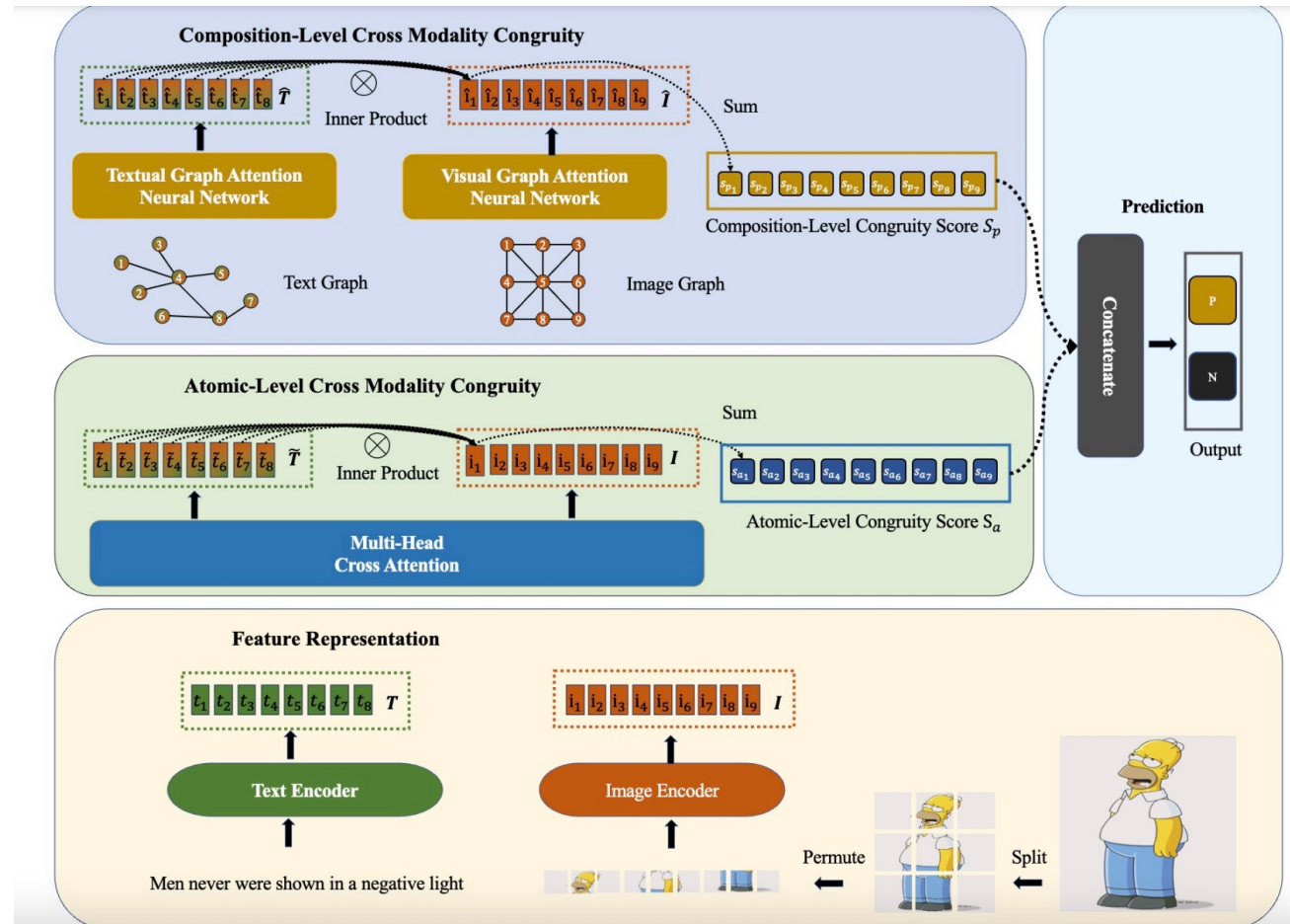
$$\mathbf{Q}_p = \frac{1}{\sqrt{d}} ([\hat{\mathbf{T}} \parallel \mathbf{c}] \hat{\mathbf{I}}^\top)$$

$$\mathbf{s}_p = \text{softmax}([\hat{\mathbf{T}} \parallel \mathbf{c}] \mathbf{W}_p + \mathbf{b}_p)^\top \mathbf{Q}_p, \quad (7)$$

$$\mathbf{p}_v = \text{softmax}(\mathbf{I}\mathbf{W}_v + \mathbf{b}_v), \quad (9)$$

$$\mathbf{y}' = \text{softmax}(\mathbf{W}_y [\mathbf{p}_v \odot \mathbf{s}_a \parallel \mathbf{p}_v \odot \mathbf{s}_p] + \mathbf{b}_y), \quad (10)$$

Overview



Experiments

Table 3.2 Comparison results for sarcasm detection (%).

| 模型 | | 准确率 | 精确率 | 召回率 | F1 |
|-----|-----------------|--------------|--------------|--------------|--------------|
| 文本 | TextCNN | 78.06 | 73.30 | 77.20 | 73.54 |
| | Bi-LSTM | 80.20 | 75.66 | 75.40 | 76.21 |
| | BERT | 81.75 | 76.62 | 81.82 | 78.20 |
| 图像 | ResNet | 64.01 | 53.90 | 71.23 | 60.79 |
| | ViT | 68.32 | 54.91 | 71.05 | 62.23 |
| 多模态 | HFM | 83.44 | 76.57 | 84.15 | 80.18 |
| | InCrossMGs | 85.57 | 81.26 | 84.32 | 82.14 |
| | CMGCN | 85.63 | - | - | 82.65 |
| | HKE | 85.92 | 81.40 | 84.93 | 83.13 |
| | RDKG (Owen) | 86.68 | 82.81 | 85.04 | 83.91 |
| | RDKG (MiniGPT4) | 86.39 | 82.49 | 84.62 | 83.55 |

Table 2: Comparison results for sarcasm detection. † indicates ResNet backbone and ‡ indicates ViT backbone.

| Model | | Acc(%) | P(%) | R(%) | F1(%) |
|-------------|--------------|--------|--------------|--------------|-------|
| Text | TextCNN | 80.03 | 74.29 | 76.39 | 75.32 |
| | Bi-LSTM | 81.90 | 76.66 | 78.42 | 77.53 |
| | SMSD | 80.90 | 76.46 | 75.18 | 75.82 |
| | BERT | 83.85 | 78.72 | 82.27 | 80.22 |
| Image | Image | 64.76 | 54.41 | 70.80 | 61.53 |
| | ViT | 67.83 | 57.93 | 70.07 | 63.43 |
| Multi-Modal | HFM† | 83.44 | 76.57 | 84.15 | 80.18 |
| | D&R Net† | 84.02 | 77.97 | 83.42 | 80.60 |
| | Att-BERT† | 86.05 | 80.87 | 85.08 | 82.92 |
| | InCrossMGs‡ | 86.10 | 81.38 | 84.36 | 82.84 |
| | CMGCN‡ | 86.54 | - | - | 82.73 |
| | Ours† | 87.02 | 82.97 | 84.90 | 83.92 |
| Ours‡ | 87.36 | 81.84 | 86.48 | 84.09 | |

Table 3: Results of different knowledge types.

| Knowledge Type | Acc(%) | F1(%) |
|----------------------------|--------------|--------------|
| w/o external knowledge | 87.36 | 84.09 |
| Image Attributes | 86.43 | 83.30 |
| ANPs | 86.35 | 83.54 |
| Image Captions | 88.26 | 84.84 |
| Image Captions (w/o image) | 86.60 | 83.28 |



Thanks!