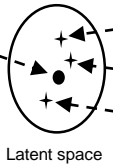


Introduction

Background: Image and text matching is an important task in the computer vision field. It remains challenging due to the heterogeneous representations.



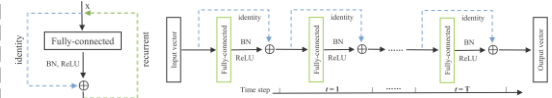
- (1) a woman dancing around with a umbrella in her hand.
- (2) a woman with an umbrella balances on one leg.
- (3) a young hipster woman holding a plaid umbrella.

Goal: How to better match images and texts, while retaining the parameters?

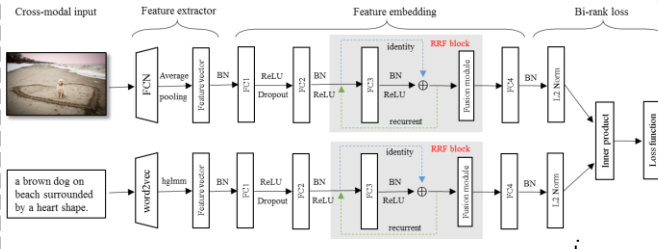
Key ideas:

- ❖ We propose a deep matching network using recurrent residual fusion (RRF) as building blocks for improving feature embeddings between image and text.
- ❖ We employ a bi-rank loss function to enforce separability of two modalities in the embedding space.

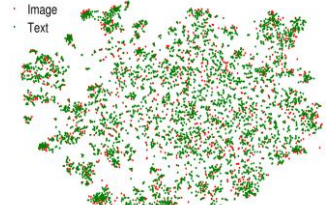
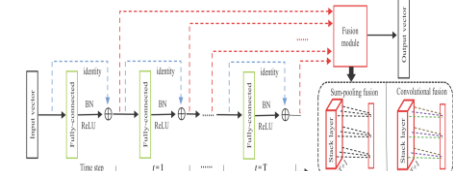
Method



RRF-Net architecture: it is composed of three components, including feature extractor, feature embedding, and bi-rank loss.



RRF building block: identity connections, recurrent connections and a fusion module.



Results

Ablation study: present comprehensive results on Flickr30K.

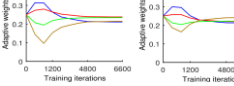
(1) Effect of recurrent steps.

(3) Effect of bi-rank loss.

Method	R@1	R@5	R@10	R@15	R@20
Baseline	45.0	75.5	83.6	86.5	88.5
RRF-Net, T=1	46.4	76.1	84.3	87.3	89.3
RRF-Net, T=2	46.9	76.8	84.8	87.7	89.7
RRF-Net, T=3	47.6	77.4	85.4	88.3	90.3
RRF-Net, T=4	46.2	76.6	85.1	87.6	89.6

Method	R@1	R@5	R@10	R@15	R@20
Baseline, bi-directional	45.4	73.8	82.5	85.4	87.4
Baseline, bi-rank	45.0	75.5	83.6	86.5	88.5
RRF-Net, bi-directional	46.4	76.5	84.1	87.4	89.4
RRF-Net, bi-rank	47.6	77.4	85.4	88.3	90.3

(2) Effect of fusion modules.



Model ensemble

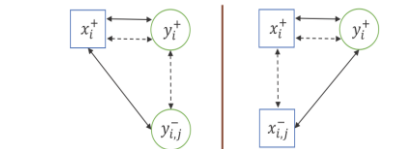
Method	R@1	R@5	R@10	R@15	R@20
RRF-Net, M = {1, 2, 3}	45.8	75.9	84.2	87.1	89.1
RRF-Net, M = {1, 2, 3, 4}	47.1	76.8	85.0	87.6	89.6
RRF-Net with core fusion	47.6	77.4	85.4	88.3	90.3

Cross-dataset	Method	R@1	R@5	R@10	R@15	R@20
Train: Flickr30K	Baseline	23.7	52.0	64.0	71.9	75.8
	RRF-Net, T=3	24.8	53.0	64.8	72.8	76.7
Train: MSCOCO	Baseline	26.4	54.3	66.2	74.1	78.0
	RRF-Net, T=3	27.2	55.8	67.1	75.0	78.9

Comparison with the state-of-the-art approaches.

Method	Flickr30K dataset					MSCOCO dataset						
	R@1	R@5	R@10	R@15	R@20	R@1	R@5	R@10	R@15	R@20		
DVSA [15]	22.2	48.2	61.4	15.2	37.7	50.5	38.4	69.9	80.5	27.4	60.2	74.8
SC-NLNet [17]	23.0	50.7	62.9	16.8	42.0	56.5	NA	NA	NA	NA	NA	NA
Mean vector [18]	24.8	52.5	64.3	20.5	46.3	59.3	33.2	61.8	75.1	24.2	56.4	72.4
Deep CCA [12]	27.9	56.9	68.2	26.8	52.9	66.9	NA	NA	NA	NA	NA	NA
CGM+HGLMM [14]	35.0	62.0	73.8	25.0	52.7	66.0	39.4	67.9	80.9	25.1	59.8	76.6
m-RNN [27]	35.4	63.8	73.7	22.8	50.7	63.1	41.0	73.0	83.5	29.0	42.2	77.0
RNN-FV [20]	35.6	62.5	74.2	27.4	55.9	70.0	41.5	72.0	82.9	29.2	64.7	80.4
RNN+Embed [25]	33.6	64.1	74.9	26.2	56.3	69.6	42.8	73.1	84.1	32.6	68.6	82.8
mC/NEmbedNet [19]	40.3	68.9	79.9	29.7	60.1	72.1	50.1	79.7	89.2	30.6	75.2	86.9
2WayNet [4]	49.8	67.5	NA	36.0	55.6	NA	55.8	75.2	NA	39.7	63.3	NA
RRF-Net	47.6	77.4	87.1	35.4	68.3	79.9	56.4	85.3	91.5	43.9	78.1	88.6

Bi-rank loss: it computes both image-to-text rank loss (Left) and text-to-image rank loss (Right).



$$l_{2t} = \sum_{j=1}^N (\alpha_1 \max[0, s(x_i^+, y_j^+) - s(x_i^+, y_{i,j}^-)] + \alpha_2 \max[0, s(x_i^+, y_i^+) - s(y_i^+, y_{i,j}^-)])$$

$$l_{2i} = \sum_{j=1}^N (\alpha_1 \max[0, s(y_i^+, x_i^+) - s(y_i^+, x_{i,j}^-)] + \alpha_2 \max[0, s(y_i^+, x_i^+) - s(x_{i,j}^-, y_i^+)])$$

$$l(x_i^+, y_i^+, x_i^-, y_i^-) = \frac{\beta_1 l_{2t} + \beta_2 l_{2i}}{N}$$

References

- [1] A. Eisenstein et al. Linking image and text with 2-way nets. In CVPR, 2017.
- [15] L. Ma et al. Multimodal convolutional neural networks for matching image and sentence. In ICCV, 2015.
- [16] B. Klein et al. Associating neural word embeddings with deep image representations using Fisher vectors. In CVPR, 2016.
- [3] R. Wang et al. Learning deep structure-preserving image-text embeddings. In CVPR, 2016.

