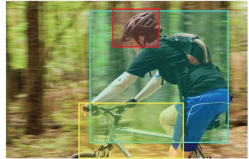


Introduction

Phrase Localization

- Given an image and its textual description, locate the image regions that correspond to the noun phrases in the description.



A woman wearing a black helmet riding on a bike.

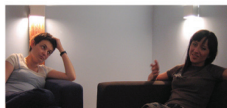


A man is working his horse on a racetrack. (examples from Flickr30Entities [1] dataset)

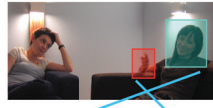
Our Contribution

- For the task of phrase localization, we propose a structured matching of phrases and regions that encourages the semantic relations between phrases to agree with the visual relations between regions.
- We formulate structured matching as a discrete optimization problem and relax it to a linear program to enable end-to-end training with neural networks.

Motivation



A woman is sitting down and leaning her head on her hand while another woman is smiling and sitting next to her.



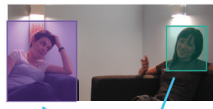
her head her hand leaning on



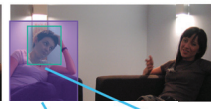
her head her hand leaning on



A woman is sitting down and leaning her head on her hand while another woman is smiling and sitting next to her.



a woman her head partial coreference



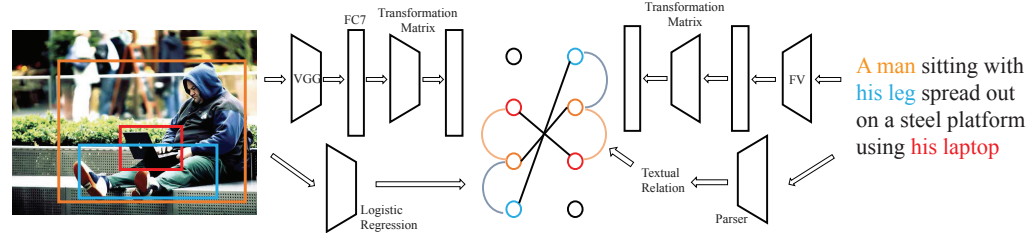
a woman her head partial coreference

- Phrase localization requires a deep understanding of semantic relations among phrases.
- This leads to the problem of structured matching of regions and phrases:
 - individual regions agree with their corresponding phrases.
 - visual relations among regions agree with textual relations among corresponding phrases.

Contact

Email: mzwang@umich.edu

Approach



Partial Coreference

Partial coreference: introduced by possessive pronouns "his", "her" or "its". For example:

A woman is dressed in Asian garb with a basket of goods on her hip.
An instructor is teaching his students how to escape a hold in a self-defense class.



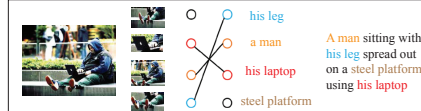
Bipartite Matching

- Denote y_{ij} as a matching configuration, $y_{ij} = 1$ is phrase p_i is matched with region r_j , $y_{ij} = 0$ otherwise.

- Denote w_{ij} as the weight of phrase p_i and region r_j .

- Solve the bipartite matching as a problem of Linear Programming:

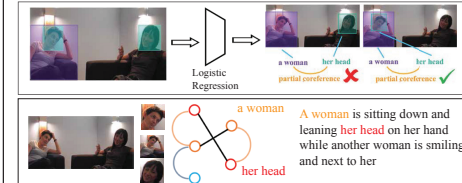
$$\begin{aligned} \max_y & \sum_{i=1}^n \sum_{j=1}^m w_{ij} y_{ij} \\ \text{s.t.} & \sum_{j=1}^m y_{ij} = 1, i = 1, 2, \dots, n \\ & \sum_{i=1}^n y_{ij} \leq 1, j = 1, 2, \dots, m \\ & 0 \leq y_{ij} \leq 1, i = 1, \dots, n, j = 1, \dots, m. \end{aligned}$$



Structured Matching

- Denote z_{ijst} as the joint configuration of phrase p_i , p_s with r_j , r_t .
- Relaxation: Refer y_{ij} as the probability of p_i is matched with r_j . Then z_{ijst} is the joint probability of p_i is matched with r_j , p_s is matched with r_t . With the rule of marginalization:

$$\sum_{t=1}^m z_{ijst} = \sum_{t=1}^m \Pr(R(p_i) = r_j, R(p_s) = r_t) = \Pr(R(p_i) = r_j) = y_{ij}$$



$$\begin{aligned} \max_{y \in \mathcal{Y}} & \sum_{i=1}^n \sum_{j=1}^m w_{ij} y_{ij} + \lambda \sum_{(i,s) \in Q} \sum_{j,t} z_{ijst} g(r_j, r_t) \\ \text{s.t.} & \sum_{j=1}^m y_{ij} = 1, \text{ for } i = 1, 2, \dots, n \\ & \sum_{i=1}^n y_{ij} \leq 1, \text{ for } j = 1, 2, \dots, m \\ & \sum_{t=1}^m z_{ijst} = y_{ij} \text{ for any } i, j, s \\ & \sum_{j=1}^m z_{ijst} = y_{st} \text{ for any } i, s, t \\ & 0 \leq y_{ij} \leq 1, \text{ for all } i, j \\ & 0 \leq z_{ijst} \leq 1, \text{ for all } i, j, s, t. \end{aligned}$$

Experiments

Experiment Setup

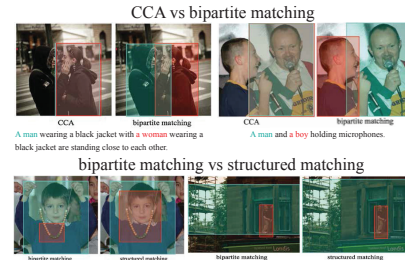
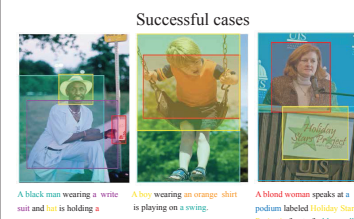
- Dataset: Flickr30K Entities [1].
 - 31783 images and 500k regions.
 - 500k noun phrases and 70k unique phrases.
- Evaluate with Recall@1 across all phrases.
- A region is true if it overlaps with the ground truth in terms of IoU > 0.5.

Results

Methods	Accuracy (Recall@1)
CCA [1]	25.30
NonlinearSP [2]	26.70 (43.89)
SCRC [9]	27.80
GroundR [3]	29.02 (47.70)
MCB [28]	(48.69)
CCA [29]	(50.89)
Ours: CCA+Fast-RCNN	39.44
Ours: Matching	41.78
Ours: Structured Matching	42.08

Methods	accuracy (Recall@1) on PC phrases only							
	person	cloth	body	anim	vehic	instru	scene	other
Bipartite Matching								
Structured Matching								
CCA[1]	29.58	24.20	10.52	33.40	34.75	35.80	20.20	20.75
GroundR[3]	44.24	9.93	1.91	45.17	46.00	20.99	30.20	16.12
CCA[29]	(53.80)	(34.04)	(7.27)	(49.23)	(58.75)	(22.84)	(52.07)	(24.13)
Ours: CCA+FRON	(64.73)	(46.88)	(17.21)	(65.83)	(68.75)	(37.65)	(51.39)	(31.77)
Ours: Bipartite	55.39	32.78	16.25	53.80	48.50	19.14	28.97	23.56
Ours: Structured	57.94	34.43	16.44	56.56	51.50	27.16	33.42	26.23
Upperbound	89.36	66.48	39.39	84.56	91.00	69.75	75.05	67.40

Qualitative Results



References

- Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image to sentence models. ICCV 2015
- Wang, L., Li, Y., Lazebnik, S.: Learning deep structure-preserving image-text embeddings. CVPR 2016
- Rohrbach, A., Rohrbach, M., Hu, R., Darrell, T., Schiele, B.: Grounding of textual phrases in images by reconstruction. ECCV 2016
- Hu, R., Xu, H., Rohrbach, M., Feng, J., Saenko, K., Darrell, T.: Natural language object retrieval. CVPR 2016
- Fukui, A., Park, D.H., Yang, D., Rohrbach, A., Darrell, T., Rohrbach, M.: Multimodal compact bilinear pooling for visual question answering and visual grounding. arXiv 2016
- Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. IJCV 2016

Code

<https://github.com/mingzhe/structured-matching>