



1 Localization and Detection with Natural Language Queries

- Tradition object detection: person, cat, dog, car, motorbike, airplane, bed, television, ...
- Natural language query: a group of bikers, a building with lots of windows, ...

2 Discriminative Bimodal Networks (DBNet)

(a) Conditional captioning using RNNs

$$p(t|r) = p(t^1|\dots) \cdot p(t^2|\dots) \cdot p(t^3|\dots) \cdot p(t^4|\dots) \cdot p(t^5|\dots)$$

(b) Our discriminatively trained CNN

Typical previous work

- Caption generation model
- Output in a huge language space
- Using only positive training samples (matched text and regions), or a limited amount of negative samples

Our DBNet

- A fully discriminative model
- Learning a binary classifier is easier
- Able to use all possible negative samples across the training set

3 Model Architectures

- Visual features: Fast R-CNN (VGGNet/ResNet)
- Text features: Character-level CNN
 - Replicating a phrase until reaching the input length of the CNN (256 characters)
 - Non-saturating neurons: Leaky ReLU
 - Replaceable by other text embedding (e.g., skip-thoughts)
- The detection score of the image region r is given by a linear classifier dynamically generated according to the text phrase t :

$$f(x, r, t; \Theta) = \mathbf{w}(t)^\top \phi_{\text{rgn}}(x, r; \Theta_{\text{rgn}}) + b(t).$$
 where the classifier weights and bias are

$$\mathbf{w}(t) = \mathbf{A}_w^\top \phi_{\text{txt}}(t; \Theta_{\text{txt}}), \quad b(t) = \mathbf{a}_b^\top \phi_{\text{txt}}(t; \Theta_{\text{txt}}).$$
- An extra regularization term on the classifier parameter is beneficial for the SGD stability:

$$\Gamma_{\text{dynamic}} = \|\mathbf{w}(t)\|_2^2 + |b(t)|^2$$

4 Model Training: Labels, Objectives, and Optimization

- Any region-text pair (r, t) can be possibly used during training.
- r can be GT or proposed (EdgeBox), t can be an annotation on the same image as r , or it can be from the rest of the training set.
- Spatial overlapping based training labels:**
 - 1: r has large overlap with a ground truth region of t
 - $\langle \text{uncertain} \rangle$: (r, t) is not positive, and r has moderate overlap with a ground truth region of t . (excluded from training)
 - Text similarity based uncertainty augmentation:**
 - If t' is similar to t , t' should also have uncertain label with r
 - 0: otherwise (including all other text phrases).
- Given a region, the non-**$\langle \text{uncertain} \rangle$** phrases are categorized into: positive phrases (pos), negative phrases from the annotations on the same image (neg), and negative phrases from the rest images (rest).
- Training loss is normalized separately for the three categories:

$$L_i = \lambda_{\text{pos}} I_i^{\text{pos}} + \lambda_{\text{neg}} I_i^{\text{neg}} + \lambda_{\text{rest}} I_i^{\text{rest}}$$
 where we set $\lambda_{\text{pos}} = \lambda_{\text{neg}} + \lambda_{\text{rest}}$, and the sampling strategy in SGD determines $\lambda_{\text{neg}}, \lambda_{\text{rest}}$ (smaller).
- Optimization: (*initialization*) Initializing the visual pathway with the pretrained faster R-CNN; (*phase 1*) training the text pathway from scratch with the visual pathway unchanged; (*phase 2*) jointly finetuning the two pathways; (*phase 3*) decreasing the learning rate.

5 Benchmarking Protocol

- Based on the Visual Genome dataset
 - Bounding boxes with text phrase descriptions
 - Additional spell checking and auto-correction
- Localization:** find a known-to-be-existing entity
- Detection:** localize all matched entities if exists any
 - Should have negative images (no matched entity)
 - Impractical to test all queries on each image
- We proposed the first benchmarking protocol** for visual entity detection with language queries. 3 difficulty levels with increasing number of randomly chosen negative images per query phrase:
 - Level 0: no negative image
 - Level 1: the same number as the positive image
 - Level 2: 5 times as the number of positive images or 20 (whichever is greater) for each test phrase
- Detection metric:** average precision (AP)
 - mean AP (mAP): averaging APs over all phrases; each phrase has its own decision threshold.
 - global AP (gAP): a single AP for any phrases; all phrases share the same decision threshold. (requires scores to be calibrated over phrases; practical for zero-shot settings)

Scenario	# unique phrases	
	Training set	Test set
Training set	2,113,688	180,363
Test set (unseen on training)	119,315	
Per image (level 0)	~43	
Per image (level 1)	~92	
Per image (level 2)	~775	

6 Experiments

We released our implementation in both Caffe+MATLAB (original results) and TensorFlow

- Support large batch size for stable training.
- Independent APIs for development and evaluation.

a. Localization performance

Method	Top-1 recall/% for IoU@			Median IoU
	0.3	0.5	0.7	
DenseCap	25.7	10.1	2.4	0.092
SCRC	27.8	11.0	2.5	0.115
DBNet w/o bias	36.3	22.4	9.4	0.124
DBNet	38.3	23.7	9.9	0.152
DBNet (ResNet)	42.3	26.4	11.2	0.205

b.1. Detection APs for IoU@0.3

Difficulty level:	0		1		2	
	AP / %	mAP	gAP	mAP	gAP	mAP
DenseCap	36.2	1.8	22.9	1.0	4.1	0.1
SCRC	38.5	2.2	37.5	1.7	29.3	0.7
DBNet	48.1	23.1	45.5	21.0	26.7	8.0
DBNet (ResNet)	51.1	24.2	48.3	22.2	29.7	9.0

b.2. Detection APs for IoU@0.5

Difficulty level:	0		1		2	
	AP / %	mAP	gAP	mAP	gAP	mAP
DenseCap	15.7	0.5	10.0	0.3	1.7	0.0
SCRC	16.5	0.5	16.3	0.4	12.8	0.2
DBNet	30.0	10.8	28.8	9.9	17.7	3.9
DBNet (ResNet)	32.6	11.5	31.2	10.7	19.8	4.3

b.3. Detection APs for IoU@0.7

Difficulty level:	0		1		2	
	AP / %	mAP	gAP	mAP	gAP	mAP
DenseCap	3.4	0.0	2.1	0.0	0.3	0.0
SCRC	3.4	0.0	3.4	0.0	2.7	0.0
DBNet	11.6	2.1	11.4	2.0	7.6	0.9
DBNet (ResNet)	12.9	2.2	12.6	2.1	8.5	0.9

- DBNet outperforms captioning models.
- The dynamic bias term helps.
- DBNet shows higher per-phrase performance in term of mAP.
- DBNet's scores are better "calibrated" over different phrases according to gAP.
- DBNet shows more robustness to negative images and rare text phrases.

c. Quantitative comparison

Green: Ground truth; Red: DenseCap; Yellow: DBNet

Green: Ground truth; Red: SCRC; Yellow: DBNet

Text phrases: a man jumping a skateboard, a man wearing a red shirt, a red white and blue baseball cap, three people hanging out in the background

GT: semitransparent & dashed
Detected: boxes of solid edges

e. Ablation study on the major components of DBNet

Prune ambiguous phrases	Phrases from other images	Finetune visual pathway	Localization				Detection (Level-1)						
			Recall / % for IoU@			Median IoU	mAP / % for IoU@			gAP / % for IoU@			
			0.3	0.5	0.7		0.3	0.5	0.7	0.3	0.5	0.7	
No	No	No	30.6	17.5	7.8	0.066	0.211	35.5	22.0	8.6	8.3	3.1	0.4
Yes	No	No	34.5	21.2	9.0	0.113	0.237	39.0	24.6	9.7	15.5	7.4	1.6
Yes	Yes	No	34.7	21.1	8.8	0.119	0.238	41.3	25.6	10.0	17.2	7.9	1.6
Yes	Yes	Yes	38.3	23.7	9.9	0.152	0.258	45.5	28.8	11.4	21.0	9.9	2.0