

Generating Visually Descriptive Language from Object Layouts

Author: Xuwang Yin, Vicente Ordonez

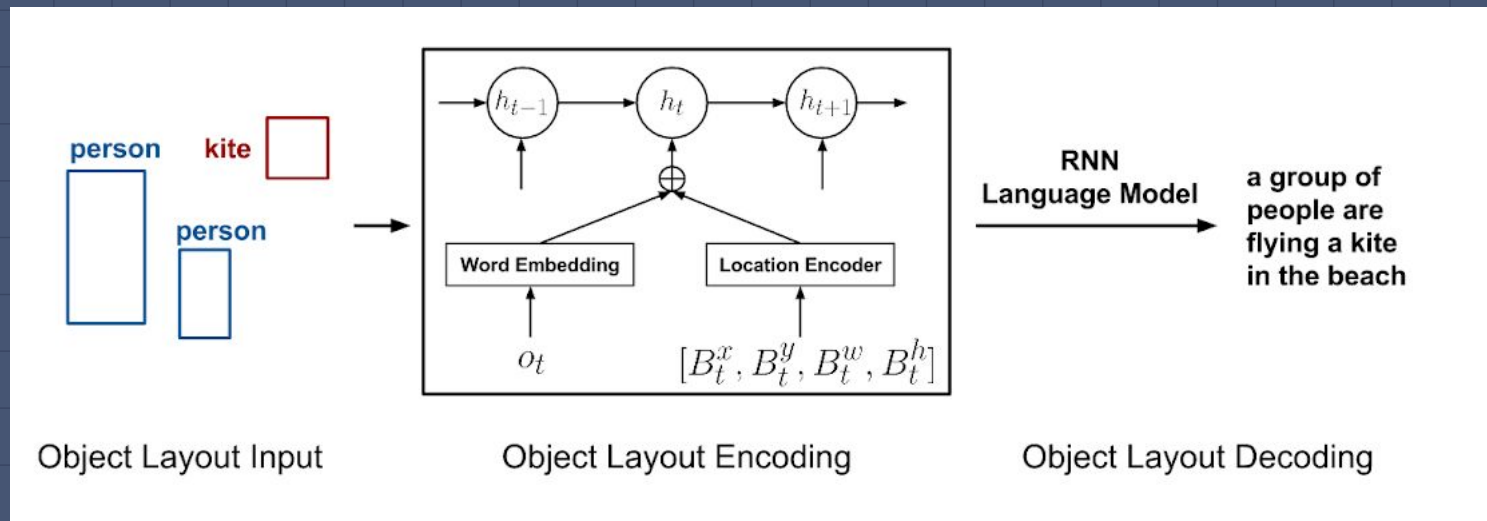
Presenter: Xing Zeng



Motivations

- ▣ Image Captioning is still a challenging problem
- ▣ Tackle a simpler problem instead: describing object layouts only
- ▣ This could be used as a middle stage to better image captioning models

Task



Encoder

Each time takes input pair

(o_t, l_t)

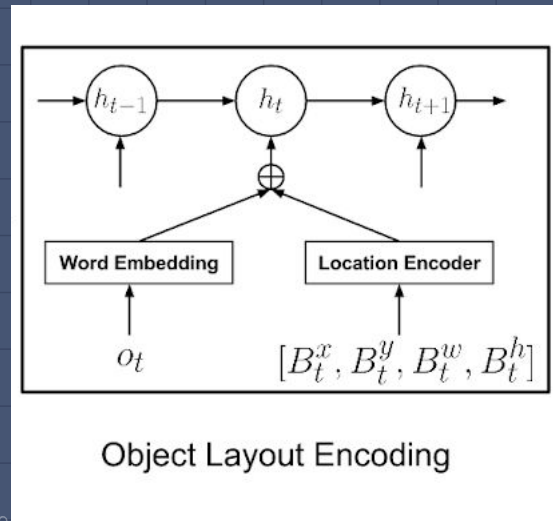
o_t = object category one-hot encoding

l_t = [left-most position (B_t^x),

top-most position (B_t^y),

width of the box (B_t^w),

height of the box (B_t^h)]



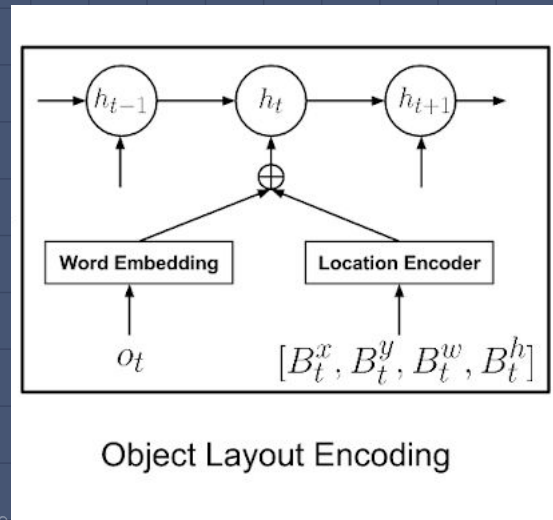
Encoder

Embedding at each time step

$$X_t = W_o o_t + W_l [B_t^x, B_t^y, B_t^w, B_t^h]$$

Hidden state at each time step

$$h_t^e = \text{LSTM}(h_{t-1}^e, x_t | W_{\text{encoder}})$$



Decoder

$$p(s|h^{\text{encoder}}) = \prod_t p(s_t|h^{\text{encoder}}, s_{<t})$$

$$p(s_t|h^{\text{encoder}}, s_{<t}) \\ = \text{softmax}(W_h h_{t-1}^d + b_h)$$

$$h_{t-1}^d \\ = \text{LSTM}(h_{t-2}^d, W_s s_{t-1} | W_{\text{decoder}})$$

RNN
Language Model

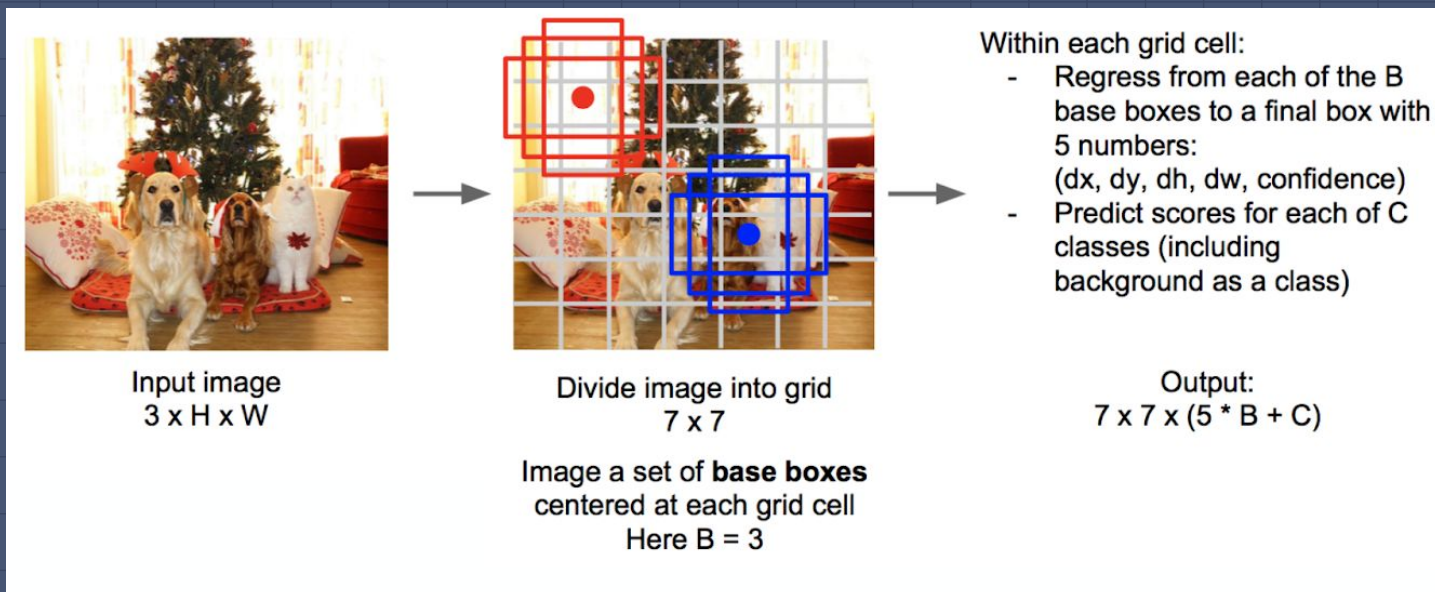
a group of
people are
flying a kite
in the beach

Object Layout Decoding

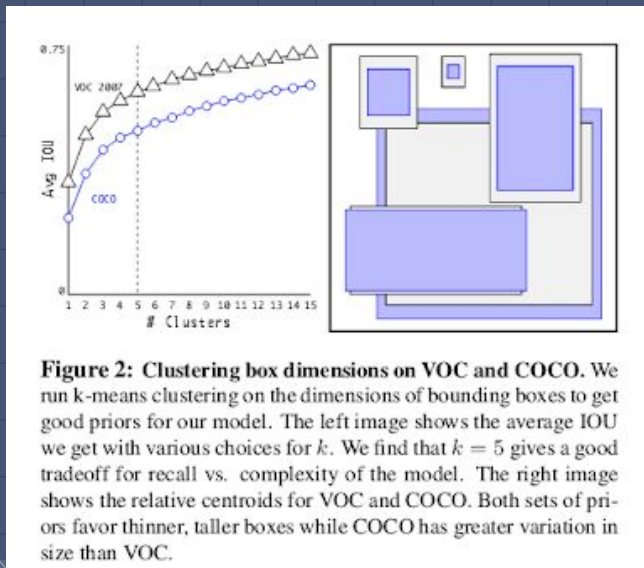
Variants

- ▣ OBJ2TEXT
 - Basic Variant
- ▣ OBJ2TEXT-YOLO
 - Object layout are generated from model YOLO instead of taking the ground truth
- ▣ OBJ2TEXT-YOLO + CNN-RNN
 - In addition to YOLO , extract visual feature using VGG-16
 - Feed encoded object layouts plus visual feature to the decoder

OBJ2TEXT-YOLO Variant



OBJ2TEXT-YOLO Variant



B^x, B^y remains the same

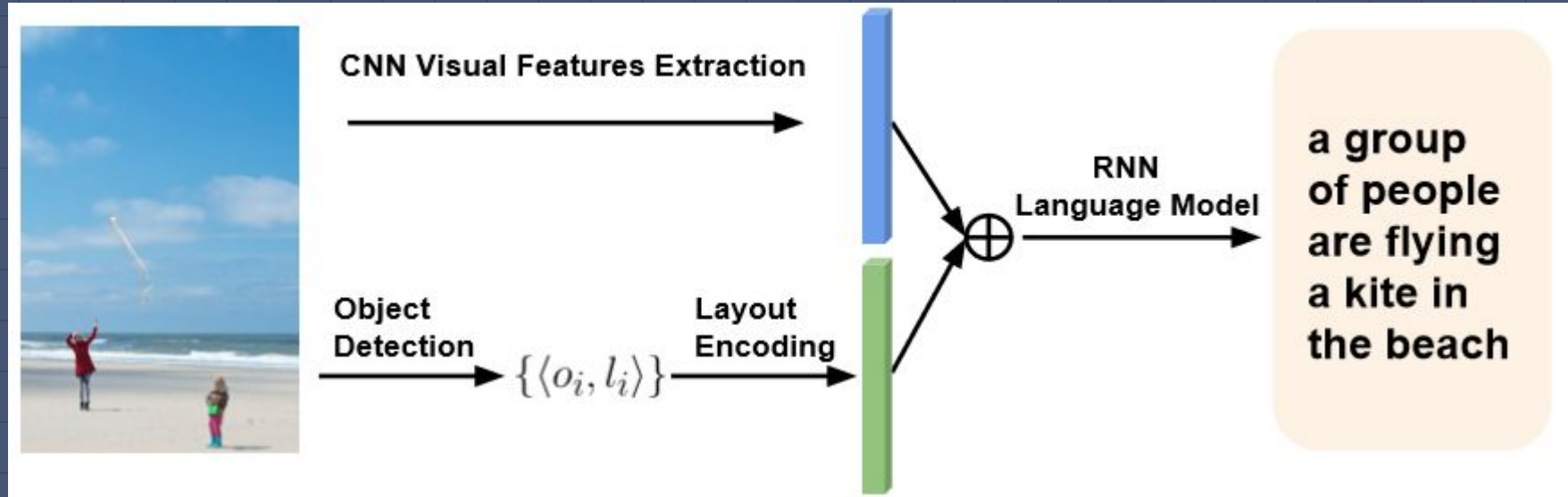
$$B^w = p_w e^{t1}$$

$$B^h = p_h e^{t2}$$

p_w, p_h are priors

$t1, t2$ are output from NN

OBJ2TEXT-YOLO + CNN-RNN Variant



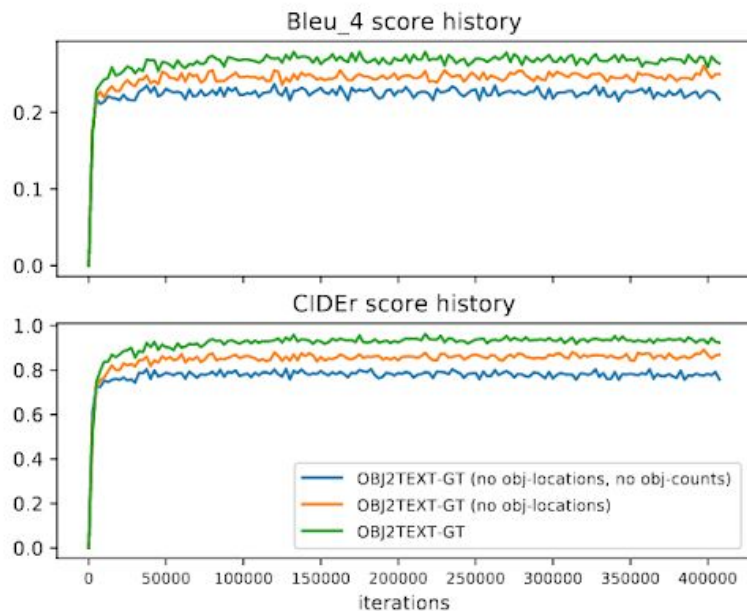
Evaluation

Train & Validation on the MS-COCO training Dataset

Test on the MS-COCO official test set

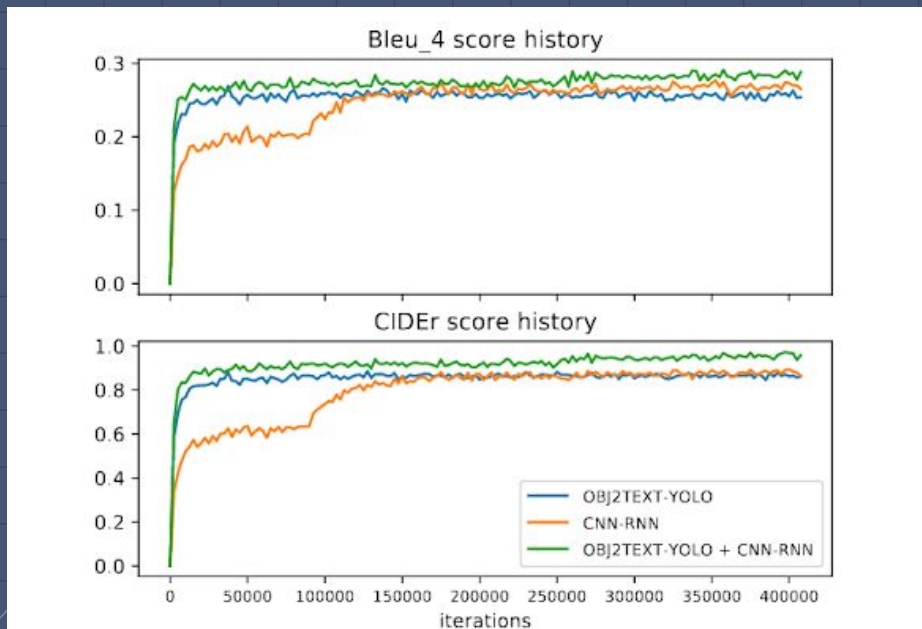


Evaluation: Ablation on OBJ2TEXT



(a) Score histories of lesioned versions of the proposed model for the task of object layout captioning.

Evaluation: YOLO-based variants



(b) Score histories of image captioning models. Performance boosts of CNN-RNN and combined model around iteration 100K and 250K are due to fine-tuning of the image CNN model.

Evaluation: Human based

"Two Alternative Forced-Choice Evaluation (2AFC)":

User are presented with one image and two alternatives,
choose the best one that describe it.

Done on Amazon Mechanical Turk.



Evaluation: Human based

Alternatives	Choice-all	Choice-agreement	Agreement
OBJ2TEXT-GT vs. OBJ2TEXT-GT (no obj-locations)	54.1%	62.1%	40.6%
OBJ2TEXT-YOLO vs. CNN+RNN	45.6%	40.6%	54.7%
OBJ2TEXT-YOLO + CNN-RNN vs. CNN-RNN	58.1%	65.3%	49.5%
OBJ2TEXT-GT vs. HUMAN	23.6%	9.9%	58.8%

Table 2: Human evaluation results using two-alternative forced choice evaluation. Choice-all is percentage the first alternative was chosen. Choice-agreement is percentage the first alternative was chosen only when all annotators agreed. Agreement is percentage where all annotators agreed (random is 25%).

Choice-All: Percentage of times A was picked over B

Agreement: Percentage of times all 3 user select the same method

Choice-Agreement: Percentage of times A was picked over B and was agreed by all 3 users among all Agreement

Some Sample Output

<http://www.cs.virginia.edu/~xy4cm/obj2text/samples/>



Potential Extensions

- ▣ A better combination with visual features
- ▣ Better Seq2Seq model
- ▣ Different Training mechanism
- ▣ Can this be done in reverse i.e. TEXT2OBJ?