



DenseCap: Fully Convolutional Localization Networks for Dense Captioning

Justin Johnson, Andrej Karpathy, Li Fei-Fei (CVPR 2015)

Gaurav (13274)

Palash Chauhan (13455)

Shubham Agarwal (13674)

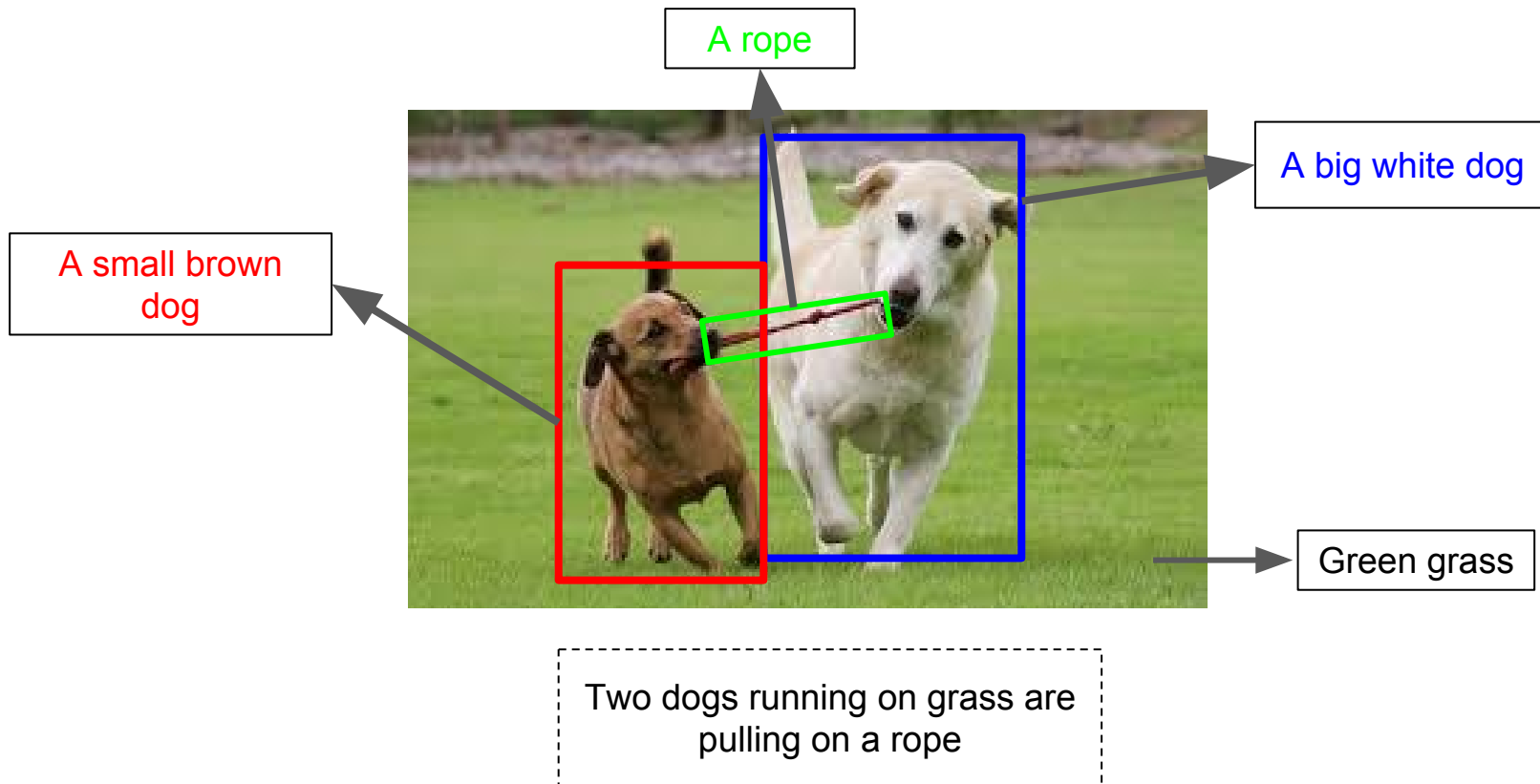
CS698: Recent Advances in Computer Vision

State of Art Presentation

Contents

1. Introduction
2. Related Work
3. Model Architecture
4. Dataset
5. Evaluation and Results
6. Conclusion

Introduction



Whole Image

Image Regions

label density

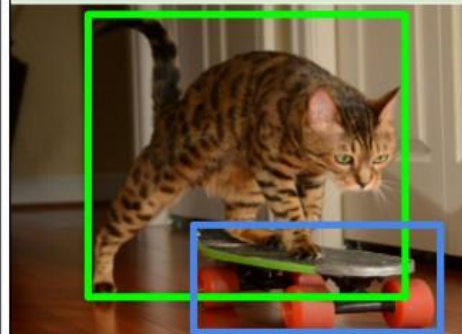
Single
Label

Classification



Cat

Detection



Cat

Skateboard

Sequence

Captioning



A cat
riding a
skateboard

Dense Captioning



Orange spotted cat

Skateboard with
red wheels

Cat riding a
skateboard

Brown hardwood
flooring

label
complexity

Related Work

Object Detection

- [O. Russakovsky et. al. 2015]
- [S. Ren et. al. 2015]
- [R. Girshick et. al. 2014]
- [A. Krizhevsky et. al. 2012]

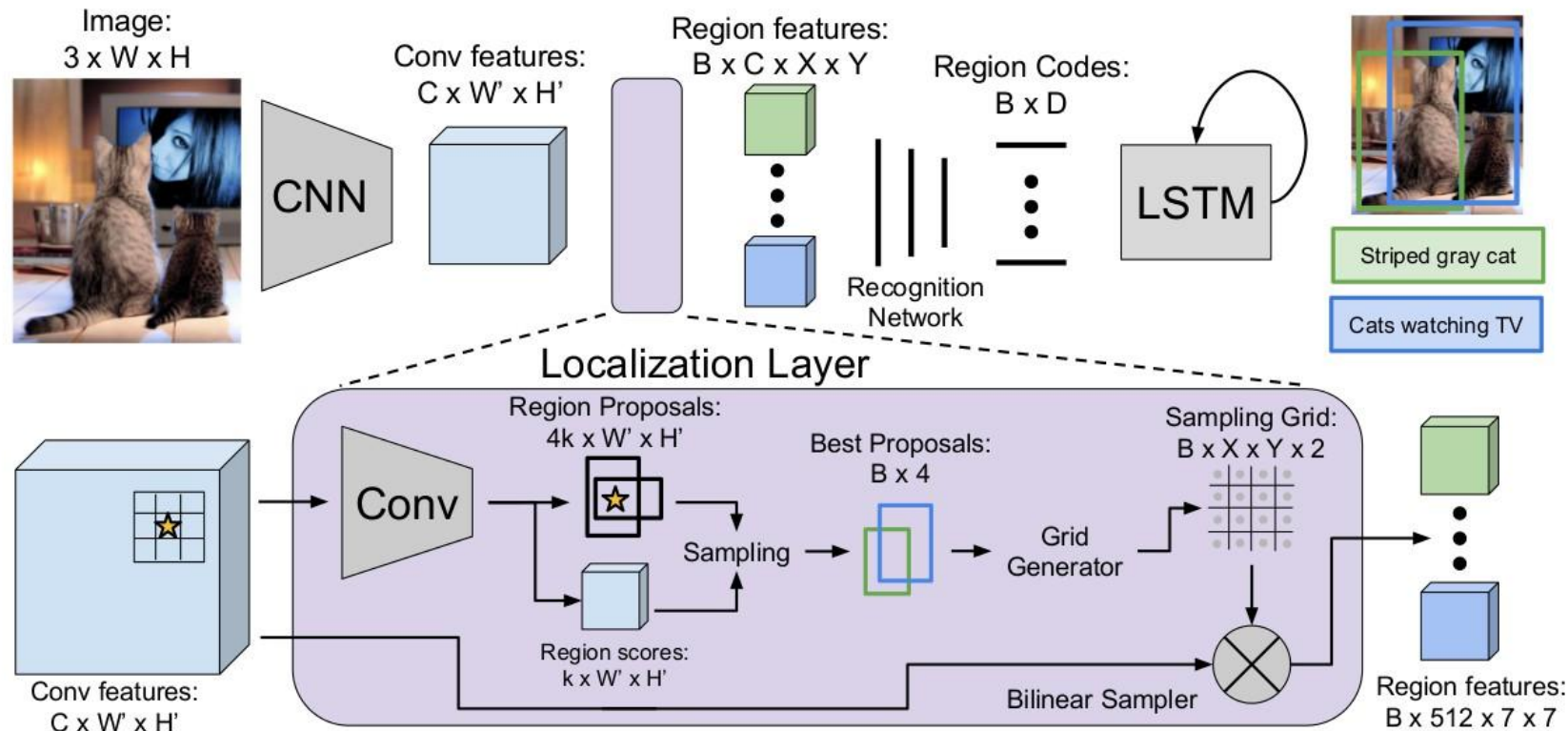
Dense Captioning

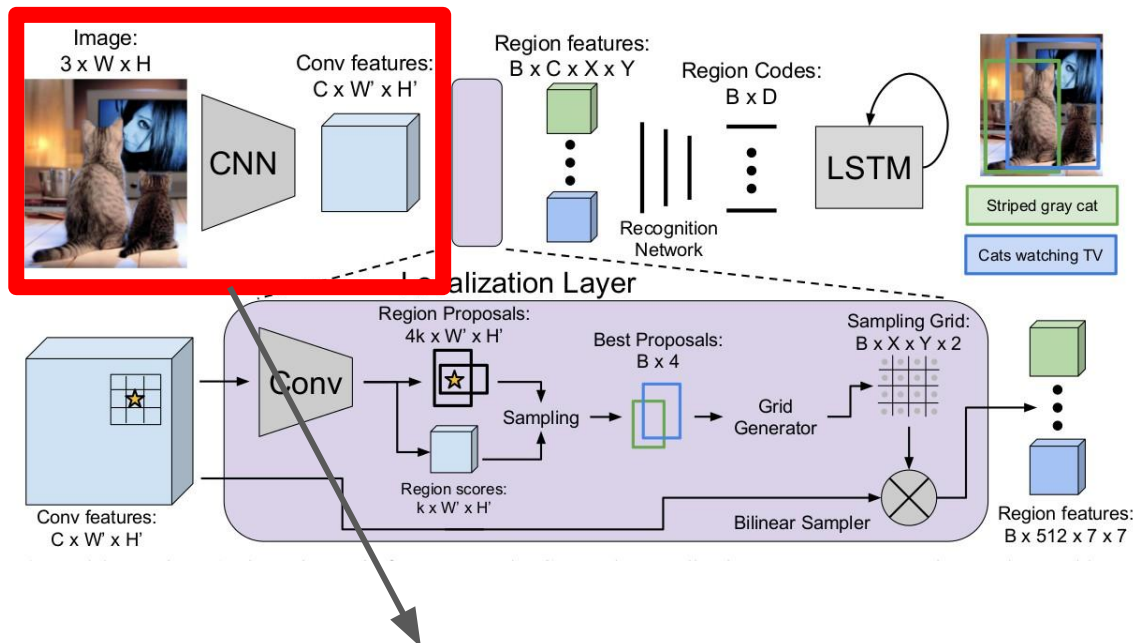
- [Karpathy et. al. 2014]

Image Captioning

- [Xu et. al. 2015]
- [Vinyals et. al. 2014]
- [Vedantam et. al. 2014]
- [A. Graves et. al. 2013]
- [G. Kulkarni et. al. 2011]
- [A. Farhadi et. al. 2010]

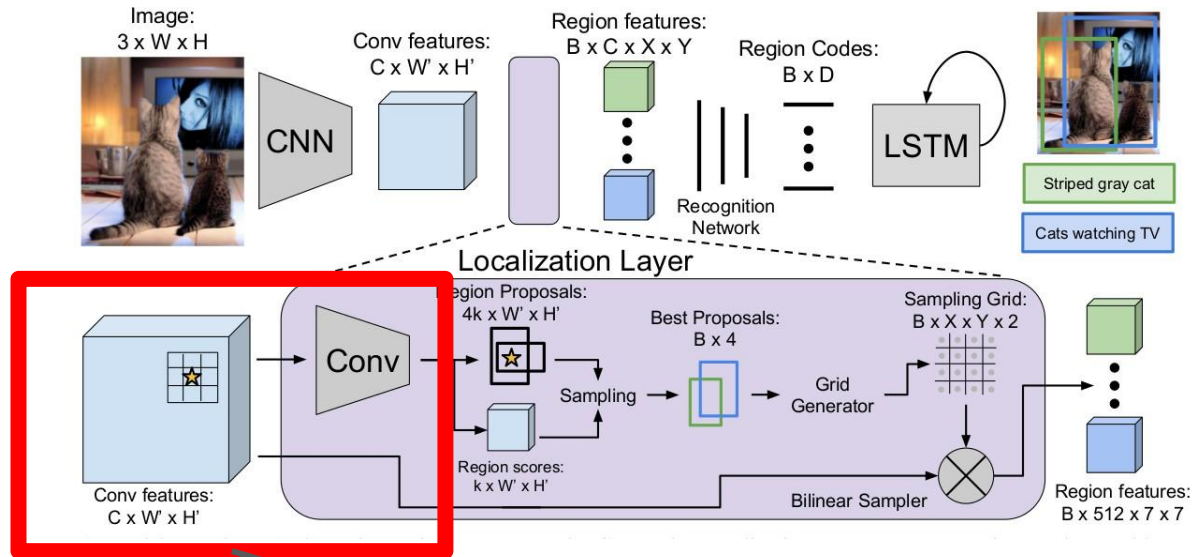
Model Architecture





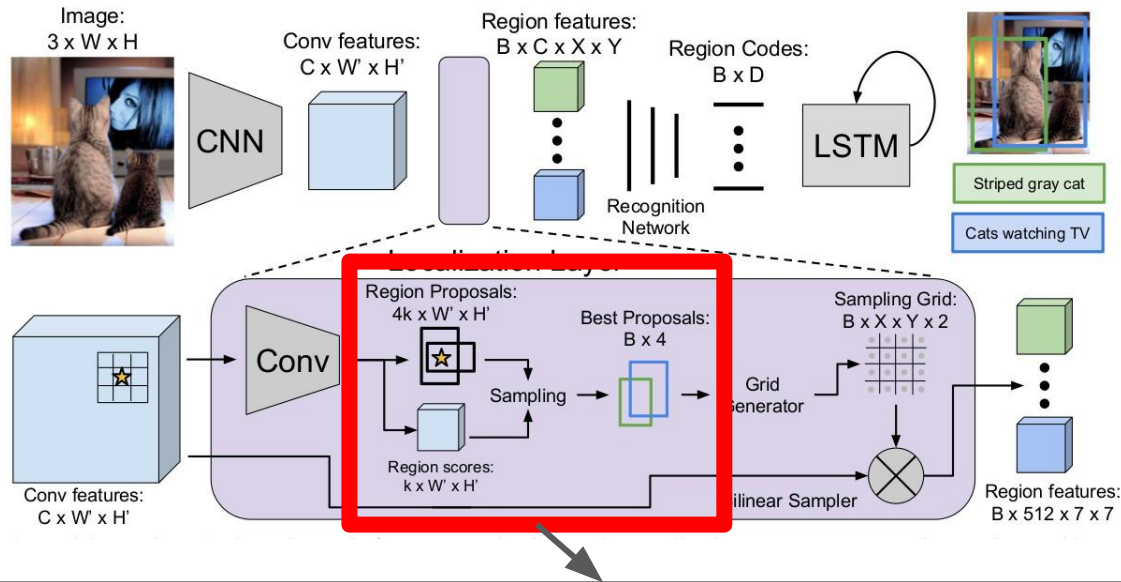
Convolutional Network

- VGG-16 architecture
- Final pooling layer removed
- Input: Image, $3 \times W \times H$
- Output: Feature Tensor, $C \times W' \times H'$



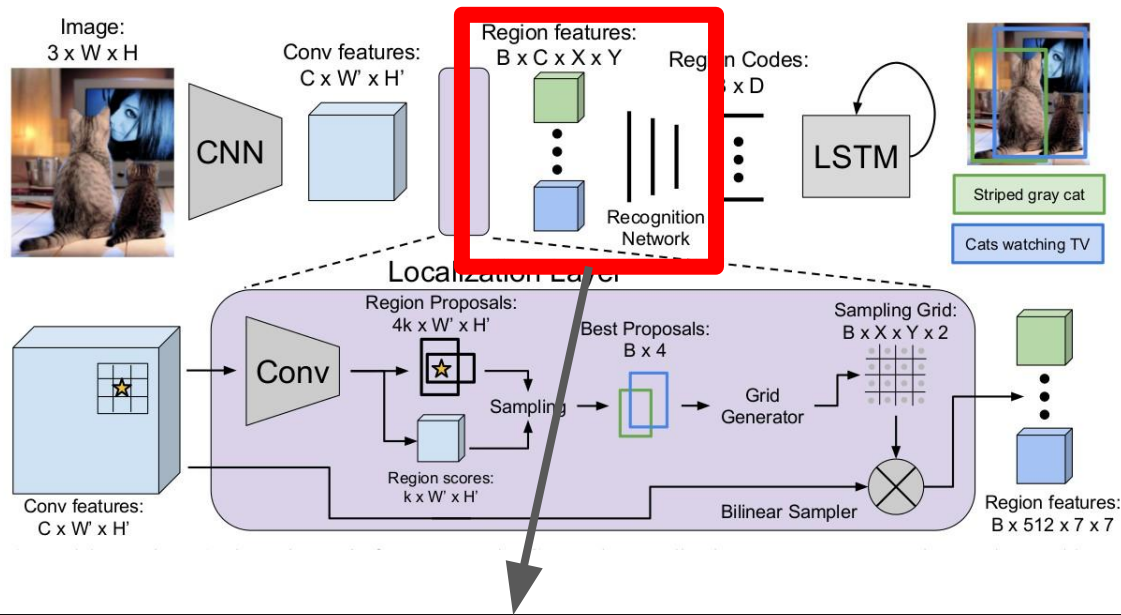
Convolutional Anchors, Box Regression

- Approach based on Faster R-CNN
- Predict region proposals by regressing offsets from a set of translation-invariant anchors
- k anchor boxes at every point
- Regress from anchors to the region proposals



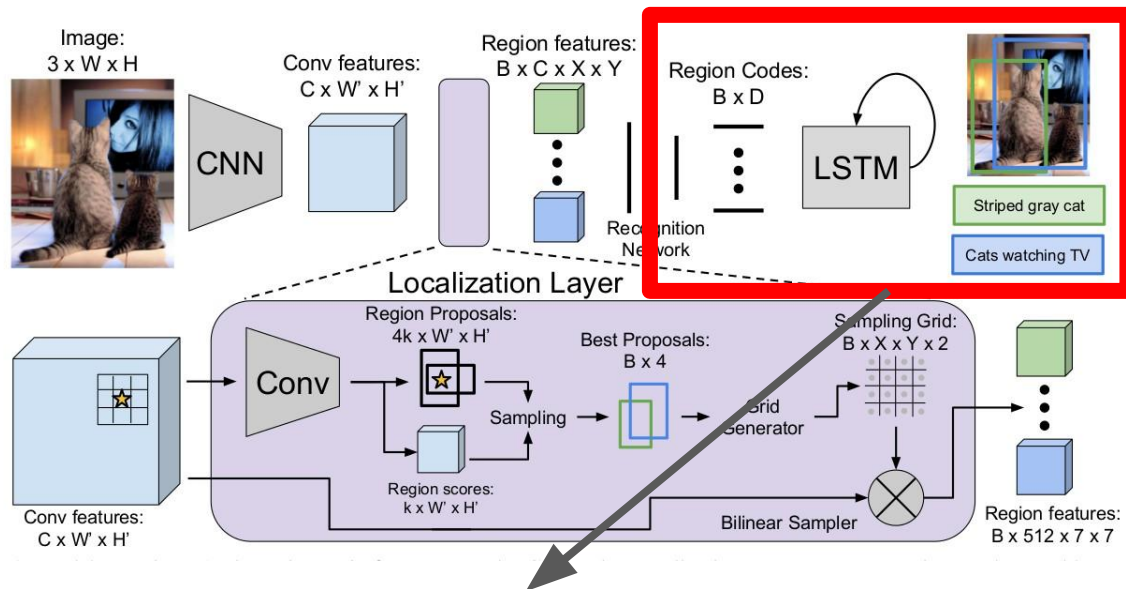
Box Sampling

- $5k \times W' \times H'$: Coordinates and scores of regional proposals
- $W = 720, H = 540, k = 12 \rightarrow 17280$ proposals \rightarrow expensive
- During training, sample a minibatch of $B=256$ regions with atmost $B/2$ positives and rest negatives.
- Positive, if $\text{IoU} > 0.7$ with some ground-truth region
- Negative, if $\text{IoU} < 0.3$ with all ground-truth regions
- Also, the predicted region of maximal IoU with each ground-truth region is positive



Recognition Network

- Fully-connected neural network
- Input (from bilinear interpolation): tensor, $B \times C \times X \times Y$
- Produces a code for every region
- Hence, compactly encodes visual appearance
- $D = 4096$

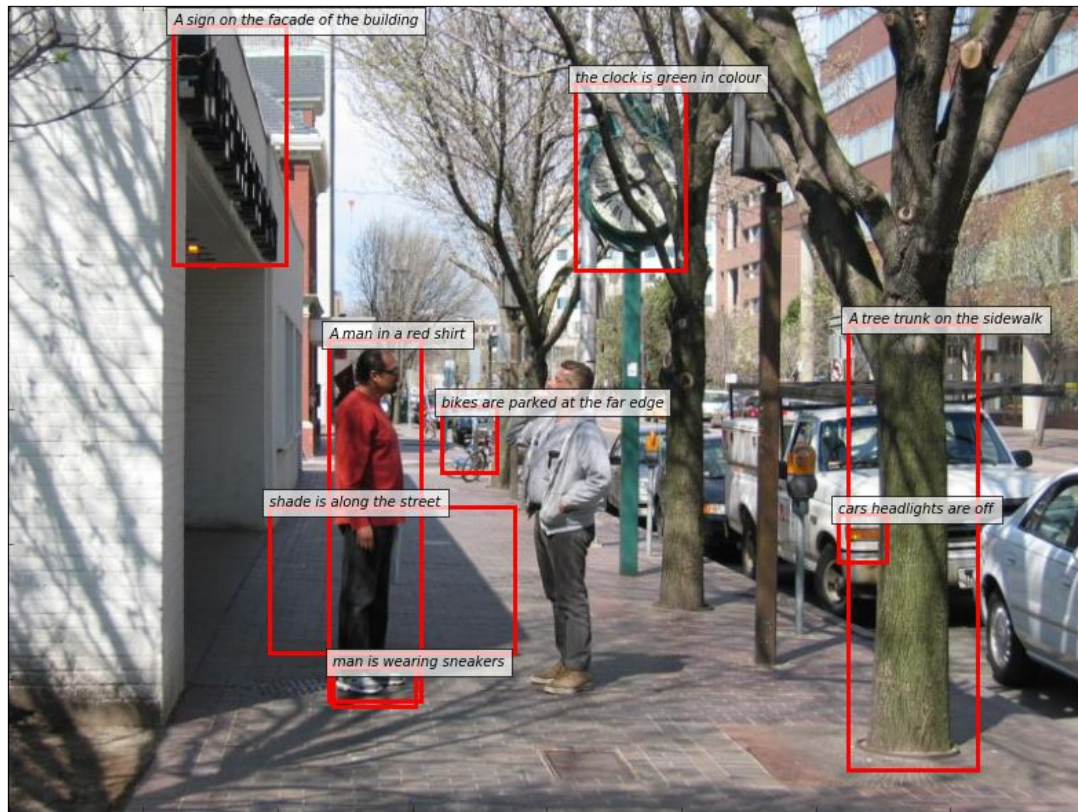


RNN Language Model

- use the region codes to condition an RNN language model
- Training: Feed word vectors appended with a START token and the region code. Obtain output vectors of length $|V| + 1$. $|V|$ = token vocabulary
- Loss function on output vectors \rightarrow average cross entropy
- Testing: Feed the region code in the beginning. Sample most likely next token and feed it to RNN in the next time step

Visual Gnome Dataset

- 94,313 images
- 4,100,413 snippets of text (43.5 per image)



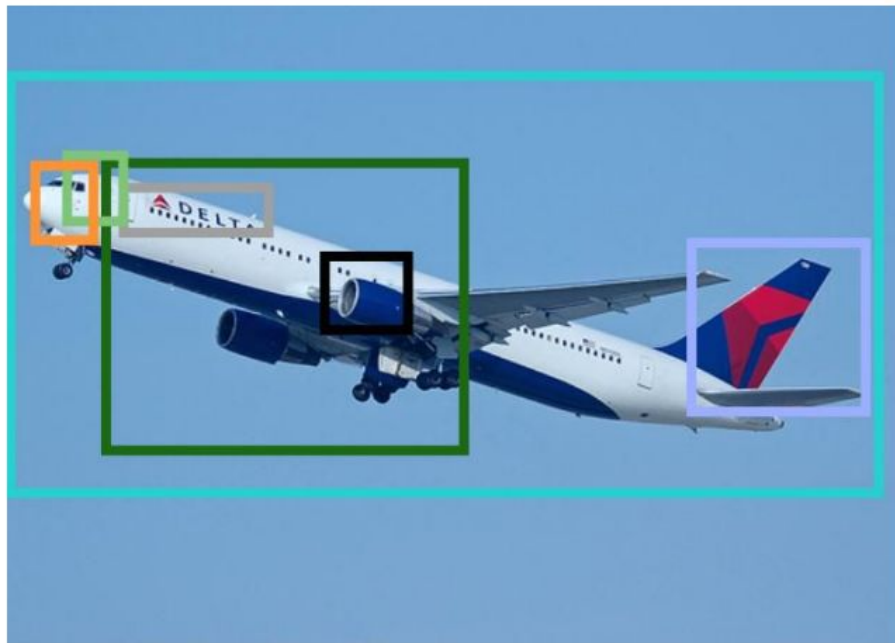
Evaluation Metric

- Mean Average Precision
 - Localization
 - Intersection over union
 - Language accuracy
 - METEOR score

Baseline Models

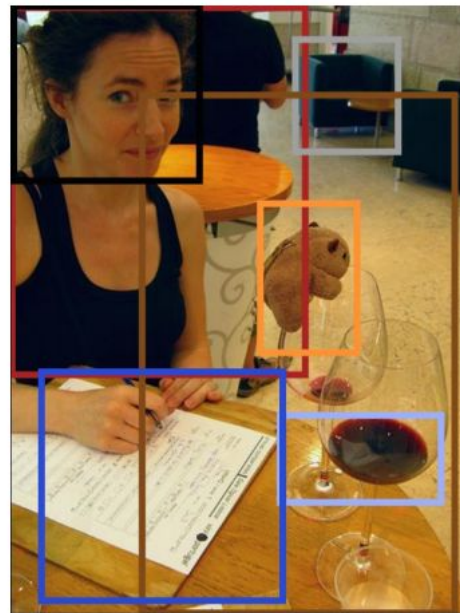
- Baseline Models
 - Region RNN model
 - Full Image RNN model
 - FCLN on EB
- Region Proposals during testing
 - Ground truth boxes(GT)
 - Edge Boxes(EB)
 - Region Proposal Network(RPN)

Results - Qualitative



plane is flying. tail of the plane. red and white plane.
plane is white. engine on the plane. windows on the
plane. nose of the plane.

A large jetliner flying through a blue sky.



woman wearing a black shirt. teddy
bear is brown. chair is black. glass of
wine. table is brown. woman with
brown hair. paper on the table.

*A man and a woman sitting
at a table with a cake.*

Results - Quantitative

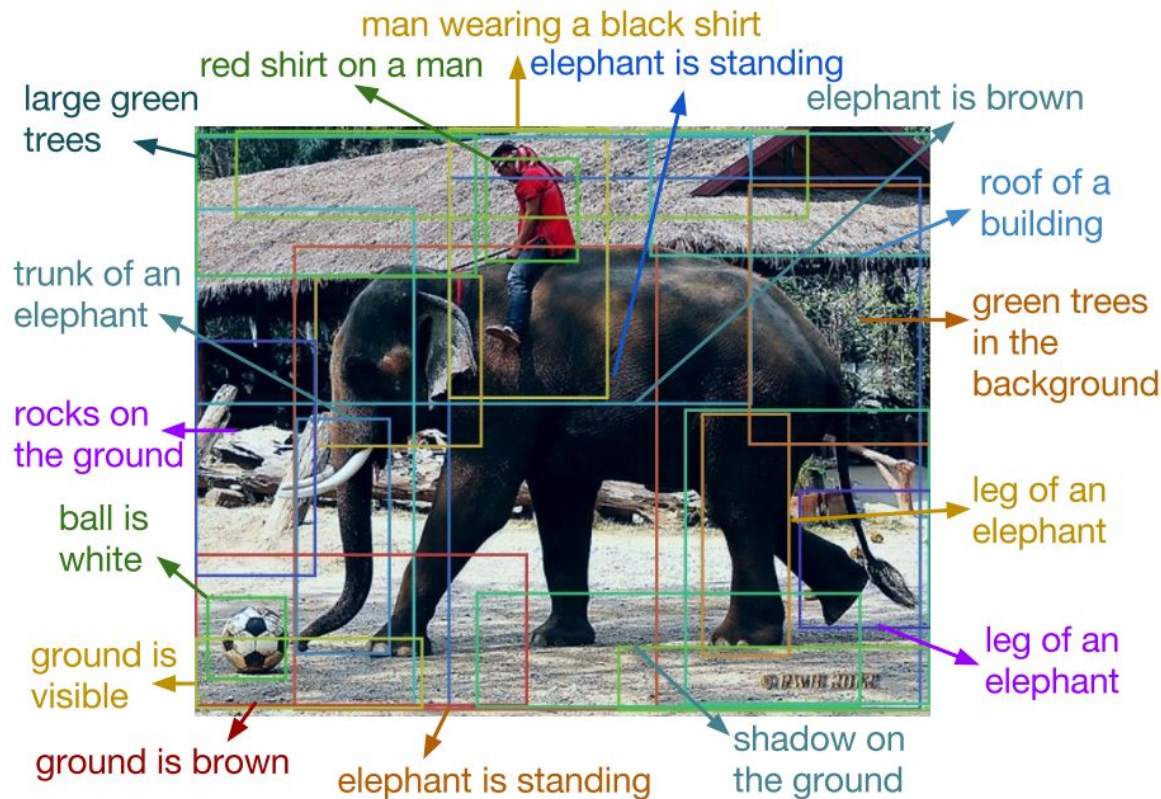
Region source	Language (METEOR)			Dense captioning (AP)			Test runtime (ms)			
	EB	RPN	GT	EB	RPN	GT	Proposals	CNN+Recog	RNN	Total
Full image RNN [21]	0.173	0.197	0.209	2.42	4.27	14.11	210ms	2950ms	10ms	3170ms
Region RNN [21]	0.221	0.244	0.272	1.07	4.26	21.90	210ms	2950ms	10ms	3170ms
FCLN on EB [13]	0.264	0.296	0.293	4.88	3.21	26.84	210ms	140ms	10ms	360ms
Our model (FCLN)	0.264	0.273	0.305	5.24	5.39	27.03	90ms	140ms	10ms	240ms

Further Applications

- Image retrieval using regions and captions

GT image	Query phrases	Retrieved Images				
	<div>man playing tennis outside</div> <div>logo with red letters</div> <div>pair of white shoes</div> <div>red and black tennis racket</div>					
	<div>black seat on bike</div> <div>chrome exhaust pipe</div> <div>white and black motorcycle</div> <div>woman in a store</div>					

Conclusion



Questions ?