

Choosing Linguistics over Vision to Describe Images



INPUT



An ocean boat is travelling in narrow water

OUTPUT

Ankush Gupta, Yashaswi Verma, C.V. Jawahar
IIIT Hyderabad, India

What is present in an Image?



Keywords : { boeing 747, airport, parked }

Phrase : “boeing 747 at the airport”

Sentence : “A boeing 747 is parked at the airport.”

Advantages

Useful in Image Indexing and Retrieval

Limitations



{ black, dog, car }

“A **black dog** is sitting inside a car.”



{ black, dog, car }

“A dog is sitting inside a **black car**.”

Problem Statement

TRAINING DATASET



A cyclist relaxes on a bench and gazes toward the ocean.



A man on a bicycle with a racing suit.



A silver bicycle is parked in a living room.



Two bicyclists pass spectators on the road near a field.

UNSEEN IMAGE



?

Key Challenges

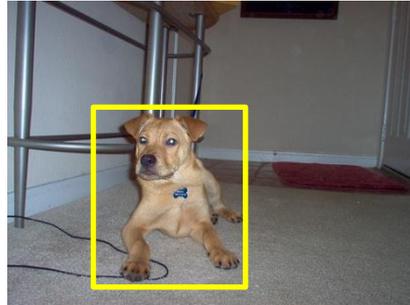


- Detecting Image content (objects, attributes, scene, action) is tough
- How they interact with each other?
- Determining correct order of words



This is a photograph of one **person** and one **sky**. The white **person** is **by** the blue **sky**. (Kulkarni et al. CVPR 2011)

- Use trained object detectors and scene/attribute classifiers



Object : Dog

Attribute : Brown

Scene : Room

- Use corpus statistics to smooth noisy vision predictions
- Predict associated verb/preposition
- Sentence generation
 - Language Model based
 - Template based

Hypothesis

- Image inherits characteristics of similar images



(white, cow), (cow, with, ear),
(cow, in, field), (grassy, field)



(white, cow), (bull, standing),
(bull, in, field), (grassy, field)

TEST IMAGE



(brown, cow), (white, cow),
(young, cow), (cow, in, field),
(grassy, field), (cow, with, ear)



(brown, cow), (young, calf),
(cow, on, grass), (cow, in, field)



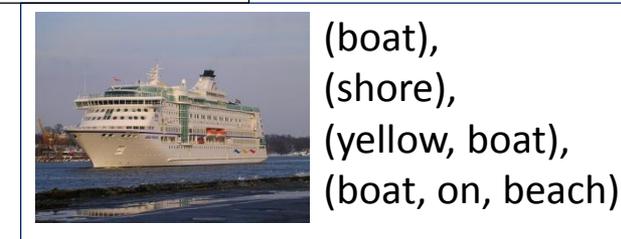
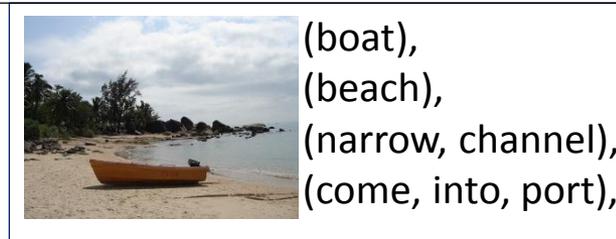
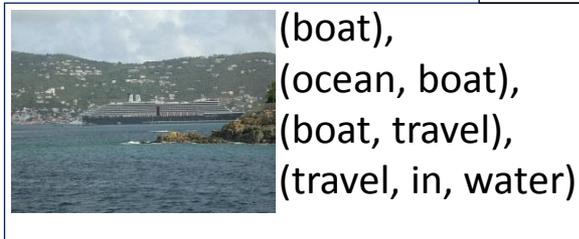
(brown, cow), (cow on farm),
(cow, with, tag), (tag, in, ear)

Overview of our Approach



Input Image

Neighbouring Images with Extracted Phrases



Ranked List of Predicted Triples

((ocean, boat), travel), (travel, in, (narrow, water)), ((ocean, boat), in, (narrow, water))
((ocean, boat), beach), (beach, on, (ocean, shore)), ((ocean, boat), on, (ocean, shore))
((ocean, boat), dock), (dock, at, (large, port)), ((ocean, boat), at, (large, port))

An ocean boat is travelling in narrow water.

Phrase Extraction

- 9 distinct types of phrases are extracted from training image descriptions
 - E.g. *(subject)*, *(attribute, subject)*, *(subject, verb)*, *(verb, prep, object)*, *(attribute, object)* etc.
- Each noun is expanded up to at most 3 hyponym levels using WordNet

Sentence	Synonym	A black and white pug is looking at the camera.
Phrases	No	<p>(pug_(subj)), (black_(attr), pug_(subj)), (white_(attr), pug_(subj)), (pug_(subj), look_(verb)), (camera_(obj)), (look_(verb), at_(prep), camera_(obj))</p>
Phrases	Yes	<p>(dog_(subj)), (black_(attr), dog_(subj)), (white_(attr), dog_(subj)), (dog_(subj), look_(verb)), (camera_(obj)), (look_(verb), at_(prep), camera_(obj))</p>

Step-1 : Predicting Phrase Relevance



- Given an unseen image I , its k most similar images (τ^k) are computed based on weighted distance between image features
 - Color, Texture, Scene (GIST) and Shape (SIFT) Descriptors
- Joint Probability of associating phrase y_i with test image I is determined based on
 - Presence / Absence of y_i in neighboring image $J \in \tau^k$
 - Corpus count of y_i
 - Feature similarity between image I and J

Step-2 : Parameter Learning



- Two types of parameters in phrase prediction model
 - For Optimal combination of different image features
 - For smoothing of phrase presence/absence probability using corpus
- Error function is defined such that
 - Probability of predicting any phrase not present in image I should be minimized
 - Probability of predicting any phrase present in image I should be greater than any other phrase
- Use Gradient Descent method to estimate parameters

Step-3 : Phrase Integration

(object ₁)	(attribute ₁ , object ₁)	(object ₁ , verb)	(verb, prep, object ₂)	(object ₂)	(attribute ₂ , object ₂)	(object ₁ , prep, object ₂)
<p>boat</p> <p>ocean channel riverboat ship</p>	<p>ocean ship</p> <p>ocean boat</p> <p>lone boat canal boat vintage ship</p>	<p>boat come</p> <p>boat travel</p> <p>ship position boat beach boat sail</p>	<p>travel in water</p> <p>sit on water come into port sail in water park at harbor</p>	<p>port sea</p> <p>water</p> <p>ocean coast</p>	<p>mountain village tropical beach ocean shore</p> <p>narrow water</p> <p>dirty shore</p>	<p>boat at port ship near coast</p> <p>boat in water</p> <p>boat in city boat in river</p>



((ocean, boat), travel), (travel, in, (narrow, water)), ((ocean, boat), in, (narrow, water))



$t = \{ ((attr_1, obj_1), verb), (verb, prep, (attr_2, obj_2)), (obj_1, prep, obj_2) \}$

- Natural Language Generation
 - Content Selection
 - Kulkarni et al. CVPR 2011 : <person, under, road>
 - Our Approach : <person, on, road>
 - Aggregation
 - *Syntactic Aggregation* : E.g. “A dark-skinned person is climbing with a pick-axis **and** posing with a green chile.”
 - Surface Realization
 - SimpleNLG (Gatt and Reiter, ENLG 2009)

1. UIUC PASCAL Sentence dataset

- 1,000 images



A girl riding a brown horse.

A girl wearing a red blouse riding a brown horse.

A woman riding a brown horse.

A young girl riding a brown horse.

A young girl wearing a helmet riding a pony.

2. IAPR TC-12 Benchmark

- 20,000 images



A grey and brown terraced house with a large, ornate entrance and two cars that are parked in front of it.

Qualitative Results (PASCAL)



A black ferrari is parked in front of a green tree.

A sporty car is parked on a concrete driveway.



An adult hound is laying on an orange couch.

A sweet cat is curling on a pink blanket.



A blond woman is posing with an elvis impersonator.

An orange fixture is hanging in a messy kitchen.

Qualitative Results (IAPR TC-12)



In this image, a dark-skinned person is climbing with a pick-axis and posing with a green chile. A green garden is surrounded with a striped chair.



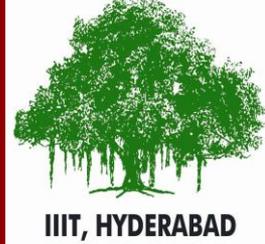
In this image, a child and a woman are sitting at an orange table. A child is sitting with tangerine.



In this image, a king-size bed is made with a white bedcover. A black cabinet is made with a brown telephone and standing on left.

Application-1

(Multiple Interesting Descriptions)



- A **white toddler** is playing with a condiment package.
- A **wide-eyed baby** is holding in a highchair.
- A **female child** is smiling in a floral shirt.



- A **panama railway** is sitting on an off-track.
- An **idle train** is heading toward a double-decker station.
- A young teenager is dancing near a **long bus**.



- A **black lamb** is standing in a dry bush.
- A **shaggy sheep** is walking through a wild grass.
- A **gray wolf** is standing in a wild grass.

Application-2 (Phrase Annotation)

“aeroplane at airport”



“sit on sofa”



“graffiti-covered bus”



- Automatic Evaluation
 - BLEU score
 - Rouge score
- Human Evaluation
 - Readability
 - Relevance



Quantitative Results(Automatic Evaluation)

DATASET	SYNONYM	BLEU-1	BLEU-2	BLEU-3	ROUGE-1
PASCAL	Yes	0.41	0.11	0.02	0.28
PASCAL	No	0.36	0.09	0.01	0.21
PASCAL Human (std.)	-	0.64	0.42	0.24	0.50
IAPR TC-12	Yes	0.21	0.07	0.01	0.14
IAPR TC-12	No	0.15	0.06	0.01	0.11

Higher Score means Better Performance

Quantitative Results (Human Evaluation)

- Likert scale of {1,2,3} where 1 is good, 2 is ok and 3 is bad

DATASET	SYNONYM	READABILITY	RELEVANCE
PASCAL	Yes	1.19	1.57
PASCAL	No	1.24	1.76
IAPR TC-12	Yes	1.38	2.32
IAPR TC-12	No	1.41	2.55

Lower Score means Better Performance

Qualitative Results (Comparison)



Kulkarni et al. 2011 : “This is a picture of one tree, one **road** and one **person**. The rusty tree is under the red road. The colorful **person** is near the rusty tree, and **under** the **road**.”

Yang et al. 2011 : “The **person** is showing the bird on the street.”

Li et al. 2011 : “Black **women** hanging from a black tree. Colored **man in** the tree.”

Ours : “**An American eagle is perching on a thick rope.**”

Quantitative Results (Comparison)



APPROACH	BLEU-1	BLEU-2	BLEU-3	ROUGE-1
Baby Talk (Kulkarni et al. 2011)	0.30	-	-	-
Ours	0.47	0.19	0.06	0.33
Corpus Guided (Yang et al. 2011)	0.41	0.13	0.03	0.31
Ours	0.54	0.23	0.07	0.41

Higher Score means Better Performance

Advantages over Previous Approaches



- Use of available image descriptions
- Unlimited Vocabulary
- Grammatically correct sentences
- Scalable
 - No use of trained detectors/classifiers
 - Performance likely to improve as the collection of available annotated images grows
 - Domain-Independent sentence generation approach

- Requires a collection of images with their corresponding descriptions
- Needs sufficient number of images of different categories
 - Trained detectors/classifiers might work well for categories with very low frequency in a dataset

- A generic method which simultaneously benefits from all 3 diverse sources of information
 - Visual clues
 - Corpus statistics
 - Available descriptions
- Superior performance on the UIUC PASCAL sentence data set and baseline for IAPR TC-12 benchmark.
- Future Work : Incorporating more visual clues from detectors and classifiers

Acknowledgment



- Prashanth Mannem from LTRC, IIIT Hyderabad for helpful discussions and pointers
- Thanks to Indian Association for Research in Computing Science (IARCS) and Microsoft Research India for travel support

THANK YOU