## CSC2539 - Datasets and Metrics for Image Caption Generation

Kaustav Kundu

University of Toronto

- Conceptual
  - Specific: Identifying people and locations
  - Generic: Related to scene understanding

- Conceptual
  - Specific: Identifying people and locations
  - Generic: Related to scene understanding
- Non Visual



I don't chew up the couch and pee in the kitchen mama!

Source: SBU caption dataset



Patriots quarterback Tom Brady holds the Super Bowl MVP trophy for the third time during a news conference held a day after New England beat Seattle 28-24 in Glenate, Artz Brady won a pickup truck for his standout performance and its hoping to give it to tearmise Medicoff Buffer, three sites gener interception seated the victory. (Lemine Squirel@stift) (mages)

#### Source: CBC News Website

- Conceptual
  - Specific: Identifying people and locations
  - Generic: Related to scene understanding
- Non Visual



I don't chew up the couch and pee in the kitchen mama!

Patriots quarterback Tom Brady holds the Super Bowl MVP trophy for the third time during a news conference held a day after New England beat Seattle 28-24 in Glenate, Artz Brady won a pickup truck for his standout performance and its hoping to give it to tearmise Medicoff Buffer, three sites gener interception seated the victory. (Lemine Squirel@stift) (mages)

Source: SBU caption dataset

Source: CBC News Website

Perceptual

From a professional photographer's point of view

## Types of Image Descriptions

- Conceptual
  - Specific: Identifying people and locations
  - Generic: Related to scene understanding
     Focus of the today's topic
- Non Visual



I don't chew up the couch and pee in the kitchen mama!



Particle quarterback Tom Brady holds the Super Bowl MVP trophy for the third time during a news conference held a day after New England beat Seattle 28-24 in Glendale, Ariz. Brady won a pickup truck for his standout performance and is hoping to give it to tearmise Malcoim Butler, whose late-game interception sealed the victory. (Jamie Sogirie/Givi) trages)

Source: SBU caption dataset

Source: CBC News Website

#### Perceptual

From a professional photographer's point of view

#### • Datasets for image caption generation

- Single sentence generation
- Multiple sentence/paragraph generation

#### • Datasets for image caption generation

- Single sentence generation
- Multiple sentence/paragraph generation
- Datasets for video caption generation

- Datasets for image caption generation
  - Single sentence generation
  - Multiple sentence/paragraph generation
- Datasets for video caption generation
- Datasets for referring expressions task

#### • Datasets for image caption generation

- Single sentence generation
- Multiple sentence/paragraph generation
- Datasets for video caption generation
- Datasets for referring expressions task
- Metrics
  - Image measures
  - Text measures
    - Automatic measures
    - Human based measures

#### UIUC Pascal Sentence<sup>1</sup>



- A camouflaged plane sitting on the green grass.
- A plane painted in camouflage in a grassy field
- A small camouflaged airplane parked in the grass.
- Camouflage airplane sitting on grassy field.
- Parked camouflage high wing aircraft.

- 1000 images randomly sampled from PASCAL VOC 2008 training + validation data with 20 object categories.
- 5 generic conceptual descriptions per image.

<sup>&</sup>lt;sup>1</sup>Rashtchian et. al., *Collecting Image Annotations Using Amazon's Mechanical Turk*, 2010. [Dataset Link]

#### UIUC Pascal Sentence<sup>1</sup>



- A camouflaged plane sitting on the green grass.
- A plane painted in camouflage in a grassy field
- A small camouflaged airplane parked in the grass.
- Camouflage airplane sitting on grassy field.
- Parked camouflage high wing aircraft.

lssues:

- Only 1000 images to train and test models.
- Simple captions and images.
- 25% captions do not contain verbs. 15% contain static verbs like *sit, stand, wear, look*.

<sup>&</sup>lt;sup>1</sup>Rashtchian et. al., *Collecting Image Annotations Using Amazon's Mechanical Turk*, 2010. [Dataset Link]

## Flickr 8k, Flickr 30k



- A biker in red rides in the countryside.
- A biker on a dirt path.
- A person rides a bike off the top of a hill and is airborne.
- A person riding a bmx bike on a dirt course.
- The person on the bicycle is wearing red.

- 8k images in Flickr8k,<sup>2</sup> >30k images in Flickr30k,<sup>3</sup> with 5 descriptions per image.
- More image sentence pairs to train and test models.
- 21% images (vs 40% images in UIUC Pascal Sentence dataset) have static verbs like *sit, stand, wear, look* or no verbs.

<sup>2</sup>Hodosh et. al., *Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics*, 2013. [Datset Link]

<sup>3</sup>Young et. al., From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions, 2014. [Datset Link]

# Microsoft CoCo<sup>4</sup>



- A baseball winds up to pitch the ball.
- A pitcher throwing the ball in a baseball game.
- A pitcher throwing a baseball on the mound.
- A baseball player pitching a ball on the mound.
- A left-handed pitcher throwing for the San Francisco giants.

- 120k train + validation images [vs 1k(Pascal), 31k(Flikr)].
- Instance level segmentations labels with 91 object classes and 2.5M labelled instances.
- Standard benchmark for image caption generation task.

<sup>4</sup>Lin et. al., Microsoft COCO: Common Objects in Context, 2014.[Dataset Link]

## Microsoft CoCo<sup>4</sup>



- 120k train + validation images [vs 1k(Pascal), 31k(Flikr)].
- Instance level segmentations labels with 91 object classes and 2.5M labelled instances.
- Standard benchmark for image caption generation task.

<sup>4</sup>Lin et. al., Microsoft COCO: Common Objects in Context, 2014.[Dataset Link] Kaustav Kundu (UofT) Datasets and Metrics

#### Abstract Scenes Dataset<sup>5</sup>



Jenny loves to play soccer but she is worried that Mike will kick the ball too hard.



Mike and Jenny play outside in the sandbox. Mike is afraid of an owl that is in the tree.

Source: L. Zitnick

- 1002 sets of scenes with 10 images in each.
- Reduced variability (hence complexity) than real word scenes.
- Descriptions have non-visual attributes.
- Clip-arts provide segmentation labels.

<sup>5</sup>Zitnick et.al., Bringing Semantics Into Focus Using Visual Abstraction, 2013. [Dataset Link]

Kaustav Kundu (UofT)

7 / 32

### Abstract Scenes Dataset<sup>5</sup>

Mike fights off a bear by giving him a hotdog while jenny runs away.



Source: L. Zitnick

- 1002 sets of scenes with 10 images in each.
- Reduced variability (hence complexity) than real word scenes.
- Descriptions have non-visual attributes.
- Clip-arts provide segmentation labels.

<sup>5</sup>Zitnick et.al., Bringing Semantics Into Focus Using Visual Abstraction, 2013. [Dataset Link]

Kaustav Kundu (UofT)

7 / 32

- Datasets for image caption generation
  - Single sentence generation
  - Multiple sentence/paragraph generation
- Datasets for video caption generation
- Datasets for referring expressions task
- Metrics
  - Image measures
  - Text measures
    - Automatic measures
    - Human based measures

#### Visual Genome Dataset<sup>6</sup>

Objects

Attributes

Relationships



<sup>6</sup>Krishna et. al., *Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations*, 2016. [Dataset Link]

### Visual Genome Dataset<sup>6</sup>

Objects





<sup>6</sup>Krishna et. al., Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations, 2016. [Dataset Link]

#### Visual Genome Dataset<sup>6</sup>



Source: Dataset Paper

Num.	Num.	Region			
images	categories	desc./image	Objs./image	Attr./image	Rel./image
$\sim$ 108k	$\sim$ 18k	$\sim 42$	$\sim 21$	$\sim$ 16	$\sim$ 18
Min desc	Max.	Word			
length	desc. length	count/desc.	Objs./region	Attr./region	Rel./region
1	16	$\sim 5$	$\sim$ 0.43	$\sim$ 0.41	$\sim 0.45$

<sup>6</sup>Krishna et. al., *Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations*, 2016. [Dataset Link]

### Krause et al<sup>7</sup>



#### Sentences

- 1) A girl is eating donuts with a boy in a restaurant
- 2) A boy and girl sitting at a table with doughnuts.
- 3) Two kids sitting a coffee shop eating some frosted donuts
- 4) Two children sitting at a table eating donuts.
- 5) Two children eat doughnuts at a restaurant table.

#### Paragraph

Two children are sitting at a table in a restaurant. The children are one little girl and one little boy. The little girl is eating a pink frosted donut with white icing lines on top of it. The girl has blonde hair and is wearing a green jacket with a black long sleeve shit underneath. The little boy is wearing a black zip up jacket and is holding his finger to his lip but is not eating. A metal napkin dispenser is in between them at the table. The wall next to them is white brick. Two adults are on the other side of the short white brick wall. The room has white circular lights on the ceiling and a large window in the front of the restaurant. It is daylight outside.

Source: Dataset paper

<sup>7</sup>Krause et.al., A Hierarchical Approach for Generating Descriptive Image Paragraphs, 2016.

Datasets and Metrics

## Krause et al<sup>7</sup>



#### Sentences

- 1) A girl is eating donuts with a boy in a restaurant
- 2) A boy and girl sitting at a table with doughnuts.
- 3) Two kids sitting a coffee shop eating some frosted donuts
- 4) Two children sitting at a table eating donuts.
- 5) Two children eat doughnuts at a restaurant table.

#### Paragraph

Two children are sitting at a table in a restaurant. The children are one little girl and one little boy. The little girl is eating a pink frosted donut with white icing lines on top of it. The girl has blonde hair and is wearing a green jacket with a black long sleeve shirt underneath. The little boy is wearing a black zip up jacket and is holding his finger to his lip but is not eating. A metal napkin dispenser is in between them at the table. The wall next to them is white brick. Two adults are on the other side of the short white brick wall. The room has white circular lights on the ceiling and a large window in the front of the restaurant. It is daylight outside.

#### Source: Dataset paper

#### • $\sim 20k$ images with following statistics (dataset to be public soon)

Dataset	Desc. Length	Sentence Length	Diversity*	Nouns	Adj.	Verbs	Pro- nouns
MS COCO	11.30	11.30	19.01	33.45	27.23	10.72	1.23
Krause et al	67.50	11.91	70.49	25.81	27.64	15.21	2.45

\* Diversity = 100 - Avg. CIDER similarity among sentences for each image

<sup>7</sup>Krause et.al., A Hierarchical Approach for Generating Descriptive Image Paragraphs, 2016.



**Description:** A big office desk is in the middle of the room. A Mac laptop is on top of the desk. There are a few bottles on top of the desk, on the right of the laptop. In front of the bottles there is a blue mug.

Source: S Fidler

<sup>8</sup>Kong et.al., What are you talking about? Text-to-Image Coreference, 2014. [Dataset Link] Kaustav Kundu (UoFT) Datasets and Metrics 11 / 32

# Kong et al<sup>8</sup>



**Description:** This room is filled with different types of furniture and home goods. The lights on the ceiling are strung across the room, they are circular and bright. At the back of the room, there are shelves filled with an assortment of pillows and blankets. There are a few couches facing away from those shelves. The couches have many pillows on top of them. On the second couch, which is dark green, sits a man in a plaid shirt. Another black couch faces the second couch. In front of the black couch is a shelf containing large brown bowls on the bottom shelf, towels on the second shelf, and vases on the top shelf. In front of the shelf is a dining table with brown wooden chairs, pink placemats, white dinnerware, and a brown glass bottle.

<sup>8</sup>Kong et.al., What are you talking about? Text-to-Image Coreference, 2014. [Dataset Link] Kaustav Kundu (UofT) Datasets and Metrics 11 / 32

# sent	# words	min # sent	max sent	min words	max words
3.2	39.1	1	10	6	144

<pre># nouns of interest</pre>	# pronouns	# scene mentioned	scene correct
3.4	0.53	0.48	83%

Table: Statistics per description.

- 1449 RGB-D images with 20 object categories.
- Long and complex descriptions.
- Significant co-reference.
- Deceiving information (object and scene mis-classification).

Source: S Fidler

<sup>8</sup>Kong et.al., What are you talking about? Text-to-Image Coreference, 2014. [Dataset Link] Kaustav Kundu (UofT) Datasets and Metrics <u>11 / 32</u>

- Datasets for image caption generation
  - Single sentence generation
  - Multiple sentence/paragraph generation
- Datasets for video caption generation
- Datasets for referring expressions task
- Metrics
  - Image measures
  - Text measures
    - Automatic measures
    - Human based measures



Source: Michaela Regneri

- 127 cooking videos with 20 different text descriptions/video.
- Time stamp labeling of textual descriptions with each description describing an activity label like wash, slicing, trash.
- Time stamp labelings of low level activity and participants(involving tool, patient, source, and target).
- Similarity scores of object activity pairs are available.

<sup>9</sup>Regneri et. al., Grounding Action Descriptions in Videos, 2013. [Dataset Link] Kaustav Kundu (UofT) Datasets and Metrics

## YouCook<sup>10</sup>



She chops the egg with an egg chopper and put the egg chopper in a glass container. Then she takes the egg mixture in the steel bowl and the bread pieces and butter which are kept in plates on the kitchen counter top. Then she places it near the sink. Then she applies butter on the frying pan and takes the chopped egg kept in the steel bowl.

- 88 videos with  ${\sim}8$  descriptions/video.
- Each video annotated with human descriptions, tracks for 48 different objects (belonging to 7 categories), and time intervals of 7 different actions.

<sup>10</sup>Das et. al., A Thousand Frames in Just a Few Words: Lingual Description of Videos through Latent Topics and Sparse Object Stitching, 2013. [Dataset Link]

Kaustav Kundu (UofT)

Datasets and Metrics

# LSMDC<sup>11</sup>



He lights a match.



Spinning wearily in his chair, the director faces SOMEONE.



A missile explodes nearby. A few blocks away, the boys run past houses and across more lawns.



On the street, a helmeted soldier talks into a transceiver.

Source: Anna Rohrbach

• Audio descriptions and script data aligned with movie videos.

<sup>11</sup>Rohrbach et. al., Movie Description, 2017. [Dataset Link] Kaustav Kundu (UofT) Datasets and Metrics

- Audio descriptions (ADs) are descriptions for the visually impaired, prepared by trained describers and professional narrators.
- Usually ADs have better visual descriptions and accurate alignment than script data.
- Movies have more diversity and realistic than cooking videos.
- Statistics:

Num. movies	Num. clips	Num. sentences	Avg. length of clip	Avg. sentences / clip
200	128k	128k	4.1s	>1

<sup>11</sup>Rohrbach et. al., Movie Description, 2017. [Dataset Link]

Datasets and Metrics

# LSMDC<sup>12</sup>

#### Tasks

- Movie description
  - Generate a single sentence to describe a given clip.
- Movie annotation and retrieval
  - There are two tracks (Multiple Choice Test and Movie Retrieval)
  - In the Multiple Choice test, a video clip with 5 possible captions is given. And the correct caption needs to be determined.
  - In the Retrieval task, given a text query, the nearest video needs to be retrieved.
  - MovieQA<sup>11</sup> dataset has a similar Multiple Choice task, but more information (movie clips, plots, subtitles and ADs) can be used to determine the correct caption.
- Fill in the blanks
  - Given a clip and a sentence with a blank, the task is to fill that blank.

<sup>12</sup>Rohrbach et. al., Movie Description, 2017. [Dataset Link]

<sup>&</sup>lt;sup>11</sup>Tapaswi et. al., MovieQA: Understanding Stories in Movies through Question-Answering, 2016. [Dataset Link]

- Datasets for image caption generation
  - Single sentence generation
  - Multiple sentence/paragraph generation
- Datasets for video caption generation
- Datasets for referring expressions task
- Metrics
  - Image measures
  - Text measures
    - Automatic measures
    - Human based measures

#### Referring Expressions Dataset



- This task involves referring to the particular objects described in natural language.
- Based on the ReferIt task, with more descriptive language expressions.
- Several datasets<sup>13,14</sup> have been concurrently developed for this task.

<sup>13</sup>Yu et. al., Modeling Context in Referring Expressions, 2016. [Dataset Link]
 <sup>14</sup>Mao et. al., Generation and Comprehension of Unambiguous Object Descriptions, 2016.
 [Dataset Link]

- Datasets for image caption generation
  - Single sentence generation
  - Multiple sentence/paragraph generation
- Datasets for video caption generation
- Datasets for referring expressions task
- Metrics
  - Image measures
  - Text measures
    - Automatic measures
    - Human based measures

• IoU<sup>15</sup>(or Jaccard Index)

$$\mathsf{IoU}(A,B) = \frac{A \cap B}{A \cup B}$$

• Precision, Recall, F1 measure

$$P = \frac{TP}{TP + FP}$$
$$R = \frac{TP}{TP + FN}$$
$$F1 = \frac{2 \cdot P \cdot R}{P + R}$$

<sup>15</sup>Everingham et al

## BLEU<sup>17</sup> (BiLingual Evaluation Understudy)

*a* : candidate sentence, *b* : set of reference sentences,  $w_n$  : n-gram  $c_x(y_n)$  : count of n-gram  $y_n$  in sentence *x*.

• Based on n-gram based precision.

• BLEU<sub>n</sub>(a, b) = 
$$\frac{\sum_{w_n \in a} \min\left(c_a(w_n), \max_{j=1,\dots,|b|} c_{b_j}(w_n)\right)}{\sum_{w_n \in a} c_a(w_n)}$$

• BLEU or BLEU<sub>Overall</sub> is a geometric mean of n-gram scores from 1 to 4.

<sup>16</sup>Detailed results in: Callison-Burch et. al., 2006; Reiter et. al., 2008; Hodosh et. al., 2013
 <sup>17</sup>Papineni et. al., BLEU: A Method for Automatic Evaluation of Machine Translation, 2002

## BLEU<sup>17</sup> (BiLingual Evaluation Understudy)

*a* : candidate sentence, *b* : set of reference sentences,  $w_n$  : n-gram  $c_x(y_n)$  : count of n-gram  $y_n$  in sentence *x*.

• Based on n-gram based precision.

• BLEU<sub>n</sub>(a, b) = 
$$\frac{\sum_{w_n \in a} \min\left(c_a(w_n), \max_{j=1,\dots,|b|} c_{b_j}(w_n)\right)}{\sum_{w_n \in a} c_a(w_n)}$$

- BLEU or BLEU<sub>Overall</sub> is a geometric mean of n-gram scores from 1 to 4.
- Strength
  - Automatic, easy to compute
- Weakness<sup>16</sup>
  - No constraints on the ordering of n-grams.
  - Each n-gram is treated equally.
  - A measure of fluency rather than semantic similarity between a and b.

<sup>16</sup>Detailed results in: Callison-Burch et. al., 2006; Reiter et. al., 2008; Hodosh et. al., 2013
 <sup>17</sup>Papineni et. al., BLEU: A Method for Automatic Evaluation of Machine Translation, 2002

## Rouge<sup>18</sup> (Recall Oriented Understudy of Gisting Evaluation)

*a* : candidate sentence, *b* : set of reference sentences,  $w_n$  : n-gram  $c_x(y_n)$  : count of n-gram  $y_n$  in sentence *x*.

• Based on n-gram based recall.

• ROUGE<sub>n</sub>(a, b) = 
$$\frac{\sum_{j=1}^{|b|} \sum_{w_n \in b_j} \min\left(c_a(w_n), c_{b_j}(w_n)\right)}{\sum_{j=1}^{|b|} \sum_{w_n \in b_j} c_{b_j}(w_n)}$$

<sup>18</sup>Lin et. al., ROUGE: A Package for Automatic Evaluation of Summaries, 2004

## Rouge<sup>18</sup> (Recall Oriented Understudy of Gisting Evaluation)

*a* : candidate sentence, *b* : set of reference sentences,  $w_n$  : n-gram  $c_x(y_n)$  : count of n-gram  $y_n$  in sentence *x*.

• Based on n-gram based recall.

• ROUGE<sub>n</sub>(a, b) = 
$$\frac{\sum_{j=1}^{|b|} \sum_{w_n \in b_j} \min\left(c_a(w_n), c_{b_j}(w_n)\right)}{\sum_{j=1}^{|b|} \sum_{w_n \in b_j} c_{b_j}(w_n)}$$

There are other variants like ROUGE<sub>S</sub>, ROUGE<sub>L</sub>

<sup>18</sup>Lin et. al., ROUGE: A Package for Automatic Evaluation of Summaries, 2004

*a* : candidate sentence, *b* : set of reference sentences,  $w_n$  : n-gram  $c_x(y_n)$  : count of n-gram  $y_n$  in sentence *x*.

• Based on n-gram based recall.

• ROUGE<sub>n</sub>(a, b) = 
$$\frac{\sum_{j=1}^{|b|} \sum_{w_n \in b_j} \min\left(c_a(w_n), c_{b_j}(w_n)\right)}{\sum_{j=1}^{|b|} \sum_{w_n \in b_j} c_{b_j}(w_n)}$$

- There are other variants like ROUGE<sub>S</sub>, ROUGE<sub>L</sub>
- Similar strengths and weaknesses as BLEU.

<sup>18</sup>Lin et. al., ROUGE: A Package for Automatic Evaluation of Summaries, 2004

# METEOR<sup>19</sup> (Metric for Evaluation of Translation with Explicit ORdering)

- a : candidate sentence, b : set of reference sentences
  - An alignment between *a* and *b* is first computed.



Source: Wikipedia

<sup>&</sup>lt;sup>19</sup>Banerjee et. al., METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments, 2005

# METEOR<sup>19</sup> (Metric for Evaluation of Translation with Explicit ORdering)

- a : candidate sentence, b : set of reference sentences
  - An alignment between *a* and *b* is first computed.



<sup>&</sup>lt;sup>19</sup>Banerjee et. al., METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments, 2005

# METEOR<sup>19</sup> (Metric for Evaluation of Translation with Explicit ORdering)

- a : candidate sentence, b : set of reference sentences
  - An alignment between *a* and *b* is first computed.



- Smoother penalization of different ordering of chunks.
- Higher correlation with human consensus scores.

<sup>19</sup>Banerjee et. al., METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments, 2005

a : candidate sentence, b : set of reference sentences

• 
$$CIDEr_n(a, b) = \frac{1}{|b|} \sum_{j=1}^{|b|} \frac{\mathbf{g}^n(a) \cdot \mathbf{g}^n(b_j)}{\|\mathbf{g}^n(a)\|\|\mathbf{g}^n(b_j)\|}$$
  
 $\mathbf{g}^n(x)$ : vector formed by TF-IDF scores of all n-grams in  $x$ .  
 $CIDEr(a, b) = \sum_{n=1}^{N} w_n CIDEr_n(a, b)$ 

<sup>20</sup>Vedantam et. al., CIDEr: Consensus-based Image Description Evaluation, 2014

a : candidate sentence, b : set of reference sentences

• 
$$CIDEr_n(a, b) = \frac{1}{|b|} \sum_{j=1}^{|b|} \frac{\mathbf{g}^n(a) \cdot \mathbf{g}^n(b_j)}{\|\mathbf{g}^n(a)\|\|\mathbf{g}^n(b_j)\|}$$
  
 $\mathbf{g}^n(x)$ : vector formed by TF-IDF scores of all n-grams in x.  
 $CIDEr(a, b) = \sum_{n=1}^{N} w_n CIDEr_n(a, b)$ 

- Gives more weight-age to *important* n-grams.
- Higher correlation with human consensus scores compared to above metrics.

<sup>20</sup>Vedantam et. al., CIDEr: Consensus-based Image Description Evaluation, 2014

# SPICE<sup>21</sup>

#### Motivation

#### 'False positive' (High n-gram similarity)

A young girl standing on top of a tennis court.



A shiny metal pot filled with some diced veggies.



A giraffe standing on top of a green field.



The pan on the stove has chopped vegetables in it.

....n-gram overlap is not necessary or sufficient for two sentences to mean the same ....SPICE primarily addresses false positives

Source: Peter Anderson

<sup>21</sup>Anderson et. al., SPICE: Semantic Propositional Image Caption Evaluation, 2016

#### 'False negative' (Low n-gram similarity)

# "A young girl standing on top of a basketball court"



Semantic propositions:

 There is girl
 The girl is young
 The girl is standing
 There is court
 The court is used for basketball
 The girl is on the court

Source: Peter Anderson

<sup>22</sup>Anderson et. al., SPICE: Semantic Propositional Image Caption Evaluation, 2016

## Key Idea – scene graphs<sup>1</sup>



<sup>2</sup> Klein & Manning: Accurate Unlexicalized Parsing, ACL 2003

<sup>3</sup>Schuster et. al: Generating semantically precise scene graphs from textual descriptions for improved image retrieval, EMNLP 2015

Source: Peter Anderson

<sup>23</sup>Anderson et. al., SPICE: Semantic Propositional Image Caption Evaluation, 2016

Kaustav Kundu (UofT)

Datasets and Metrics

26 / 32

## SPICE<sup>24</sup>

SPICE calculated as an F-score over tuples, with:

- Merging of synonymous nodes, and
- Wordnet synsets used for tuple matching and merging.

Given candidate caption c, a set of reference captions S, and the mapping T from captions to tuples:

$$P(c,S) = \frac{|T(c) \otimes T(S)|}{|T(c)|}$$
$$R(c,S) = \frac{|T(c) \otimes T(S)|}{|T(S)|}$$
$$SPICE(c,S) = F_1(c,S) = \frac{2 \cdot P(c,S) \cdot R(c,S)}{P(c,S) + R(c,S)}$$

Source: Peter Anderson

<sup>24</sup>Anderson et. al., SPICE: Semantic Propositional Image Caption Evaluation, 2016

Kaustav Kundu (UofT)

Datasets and Metrics

27 / 32

# SPICE<sup>25</sup>

#### Example

#### **Reference** captions

"People playing with kites outside in the desert." "A group of people at a park flying a kite." "A group of people flying a kite on a sandy beach" "People on the beach flying kites in the wind." "A couple people out flying a kite on some sand."



Candidate caption & scene graph "a group of people flying kites on a beach"



#### Source: Peter Anderson

<sup>25</sup>Anderson et. al., SPICE: Semantic Propositional Image Caption Evaluation, 2016

Kaustav Kundu (UofT)

#### Datasets and Metrics

28 / 32

Reference scene graph

- Pros:
  - Places importance on capturing details about objects, attributes and relationships.
  - Higher correlation with humans compared to n-gram based metrics.

- Cons:
  - This metric does not check whether the grammar is correct.
  - Depends on semantic parsers, which might not always be correct.
  - Equal weighting of different nouns, attributes, relationships.

<sup>26</sup>Anderson et. al., SPICE: Semantic Propositional Image Caption Evaluation, 2016

- Recall@k = % image sentence pairs for which the ground truth sentence was present in the top-k list.
- Median rank = k at which the system has a recall of 50%.

<sup>&</sup>lt;sup>27</sup>Hodosh et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

- Recall@k = % image sentence pairs for which the ground truth sentence was present in the top-k list.
- Median rank = k at which the system has a recall of 50%.
- Such measures can be used for retrieval based systems.

<sup>&</sup>lt;sup>27</sup>Hodosh et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

- Recall@k = % image sentence pairs for which the ground truth sentence was present in the top-k list.
- Median rank = k at which the system has a recall of 50%.
- Such measures can be used for retrieval based systems.
- Hodosh et. al.<sup>27</sup> shows that both automatic ranking based measures are *more* robust than metrics that consider only the quality of the first result.

<sup>&</sup>lt;sup>27</sup>Hodosh et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

- Measuring quality of a single best result
  - Rating system of 1-4 from Hodosh et. al.<sup>28</sup>

 $^{28}\mbox{Hodosh}$  et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

<sup>29</sup>Manning et. al., Introduction to Information Retrieval, 2008

Kaustav Kundu (UofT)

Datasets and Metrics

- Measuring quality of a single best result
  - Rating system of 1-4 from Hodosh et. al.<sup>28</sup>
- Measuring ranked candidates
  - Success@k = % image sentence pairs for which at least one relevant result is found in the top-k list.
  - **R-precision**<sup>29</sup> = average % of relevant items in the top-k list.

 $^{28}\mbox{Hodosh}$  et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

<sup>29</sup>Manning et. al., Introduction to Information Retrieval, 2008

Kaustav Kundu (UofT)

Datasets and Metrics

- Datasets
  - Using humans to make binary/choosing decisions, rather than complex decisions. Helps in faster and quality annotation.<sup>30</sup>
  - Games to make the creation of datasets more interesting for annotators.  $^{\rm 31}$
  - Anyhow involves post-processing to remove spelling mistakes, and sometimes grammatical mistakes.

<sup>30</sup>Parikh et. al., 2011; Vedantam et. al., 2014
<sup>31</sup>Deng et. al., 2013; Kazemzadeh et. al., 2014
<sup>32</sup>More details in Reiter et. al., 2008; Hodosh et. al., 2013

- Datasets
  - Using humans to make binary/choosing decisions, rather than complex decisions. Helps in faster and quality annotation.<sup>30</sup>
  - Games to make the creation of datasets more interesting for annotators.  $^{\rm 31}$
  - Anyhow involves post-processing to remove spelling mistakes, and sometimes grammatical mistakes.
- Metrics
  - Hodosh et. al. used qualification tests to get *experts* to compare correlation between human based measures and automatic measures.<sup>32</sup>
  - Common practice to use averaged responses from humans rather than single responses. Vedantam et. al.(2014) uses as many as 50 human responses per image sentence pair to ensure the quality of responses.

<sup>&</sup>lt;sup>30</sup>Parikh et. al., 2011; Vedantam et. al., 2014

<sup>&</sup>lt;sup>31</sup>Deng et. al., 2013; Kazemzadeh et. al., 2014

<sup>&</sup>lt;sup>32</sup>More details in Reiter et. al., 2008; Hodosh et. al., 2013