# CSC2523 - Datasets and Metrics for Visual Recognition with Text

Kaustav Kundu

University of Toronto

# Topics that Involve Images and Text

- Detecting text in images
- Generating textual descriptions of images/videos
- Visual retrieval based on complex textual queries
- Word-sense disambiguation
- Text to image/video alignment
- Learning visual models via text
- Using text to improve visual parsing
- Questions and answers

[Slide: S. Fidler]

# Overview

- Datasets
  - 2D Images
    - UIUC Pascal Sentence
    - Flickr 8k, Flickr 30k
    - Microsoft CoCo
    - Abstract Scenes

# Overview

- Datasets
  - 2D Images
    - UIUC Pascal Sentence
    - Flickr 8k, Flickr 30k
    - Microsoft CoCo
    - Abstract Scenes
  - RGB-D images
    - Chen et. al.
    - DAQUAR

# Overview

- Datasets
  - 2D Images
    - UIUC Pascal Sentence
    - Flickr 8k, Flickr 30k
    - Microsoft CoCo
    - Abstract Scenes
  - RGB-D images
    - Chen et. al.
    - DAQUAR
  - Videos
    - YouCook (Jason Corso, Ror)
    - TACoS

# Overview

- Datasets
  - 2D Images
    - UIUC Pascal Sentence
    - Flickr 8k, Flickr 30k
    - Microsoft CoCo
    - Abstract Scenes
  - RGB-D images
    - Chen et. al.
    - DAQUAR
  - Videos
    - YouCook (Jason Corso, Ror)
    - TACoS

- Metrics
  - Image measures
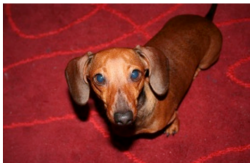
# Overview

- Datasets
  - 2D Images
    - UIUC Pascal Sentence
    - Flickr 8k, Flickr 30k
    - Microsoft CoCo
    - Abstract Scenes
  - RGB-D images
    - Chen et. al.
    - DAQUAR
  - Videos
    - YouCook (Jason Corso, Ror)
    - TACoS

- Metrics
  - Image measures
  - Text measures
    - Automatic measures
    - Human based measures

# Types of Image Descriptions

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

# Types of Image Descriptions

- 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!*
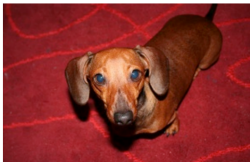
Source: *SBU caption dataset*



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 Glendale, Ariz. Brady won a pickup truck for his standout performance and is hoping to give it to teammate Malcolm Butler, whose late-game interception sealed the victory. (Jamie Squire/Getty Images)

Source: *CBC News Website*

# Types of Image Descriptions

- 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 Glendale, Ariz. Brady won a pickup truck for his standout performance and is hoping to give it to teammate Malcolm Butler, whose late-game interception sealed the victory. (Jamie Squire/Getty Images)

Source: *CBC News Website*

- Perceptual
  *From a professional photographer's point of view*

- 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.

---

[1]Rashtchian et. al., *Collecting Image Annotations Using Amazon's Mechanical Turk*, 2010. [Dataset Link]

- 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.

Issues:

- 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*.

---

[1]Rashtchian et. al., *Collecting Image Annotations Using Amazon's Mechanical Turk*, 2010. [Dataset Link]
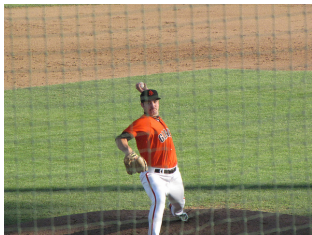
# 2D Image Datasets - Flickr 8k, Flickr 30k



- 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,[2] >30k images in Flickr30k,[3] with 5 descriptions per image.

- More image sentence pairs to train and test models, but no image based labels.

- 21% images (vs 40% images in UIUC Pascal Sentence dataset) have static verbs like *sit, stand, wear, look* or no verbs.

---

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

[3]Young et. al., *From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions*, 2014. [Datset Link]
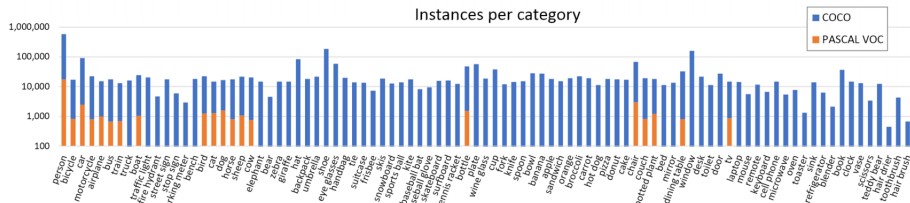
# 2D Image Datasets - Microsoft CoCo[4]



- 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.

- 328k train + validation images [vs 1k(Pascal), 31k(Flikr)].
- Instance level segmentations labels with 91 object classes and 2.5M labelled instances.
- Benchmark to be released soon and more images to be added this year.

---

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

# 2D Image Datasets - Microsoft CoCo[4]



Source: *Dataset Paper*

- 328k train + validation images [vs 1k(Pascal), 31k(Flikr)].
- Instance level segmentations labels with 91 object classes and 2.5M labelled instances.
- Benchmark to be released soon and more images to be added this year.

---

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

# 2D Image Datasets - Abstract Scenes Dataset[5]



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.

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

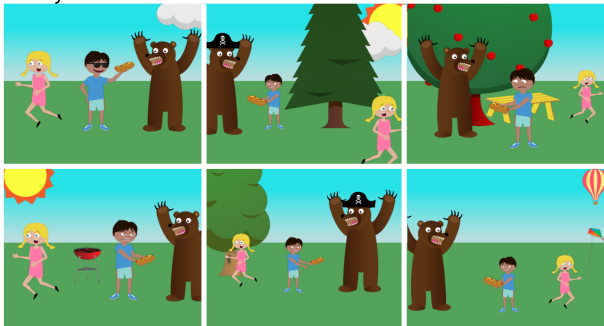Mike fights off a bear by giving him a hotdog while jenny runs away.



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.

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

# RGB-D Image Datasets - Kong et al



**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*

**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.

Source: *S Fidler*

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

| # nouns of interest | # 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*

[6]Kong et.al., What are you talking about? Text-to-Image Coreference, 2014. [Dataset Link]

**Q.** What is between the the two white and black garbage bins?
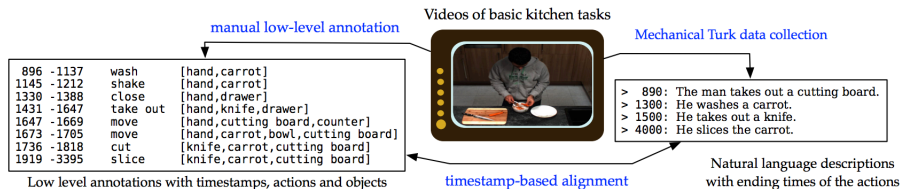**A.** Chair

**Q.** How many objects are between the fire extinguisher and the white oven on the floor?
**A.** 3

- 1449 RGB-D images with ∼9 Q&A pairs/image.
- 37 object classes.

---

[7]Malinowski et. al., A Multi-World Approach to Question Answering about Real-World Scenes based on Uncertain Input, 2014. [Dataset Link]

Videos of basic kitchen tasks

manual low-level annotation

Mechanical Turk data collection

```
 896 -1137    wash       [hand,carrot]
1145 -1212    shake      [hand,carrot]
1330 -1388    close      [hand,drawer]
1431 -1647    take out   [hand,knife,drawer]
1647 -1669    move       [hand,cutting board,counter]
1673 -1705    move       [hand,carrot,bowl,cutting board]
1736 -1818    cut        [knife,carrot,cutting board]
1919 -3395    slice      [knife,carrot,cutting board]
```

Low level annotations with timestamps, actions and objects

> 890: The man takes out a cutting board.
> 1300: He washes a carrot.
> 1500: He takes out a knife.
> 4000: He slices the carrot.

Natural language descriptions with ending times of the actions

timestamp-based alignment

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.

[8]Regneri et. al., Grounding Action Descriptions in Videos, 2013. [Dataset Link]

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 ∼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.

[9]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]

# Datasets - Summary

| | UIUC Sentence | Flickr 30k | MS CoCo | Abstract Scenes | Chen et al | DAQUAR | YouCook | TACoS |
|---|---|---|---|---|---|---|---|---|
| Generation | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Retrieval | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Word Sense Disambiguation | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ |
| Image-text Alignment | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Learning Visual Models | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ |
| Visual Parsing | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Q&A | | | | | | ✓ | | |

ICDAR Robust Reading Competition[10]

- Challeneges
  - **Born Digital**: Difficulties due to low resolution images, with compression artefacts, anti-aliasing.
  - **Scene Text**: Difficulties due to high illumination variability, perspective.
  - **Text in Videos**
  - **Accidental Scene Text**: Similar to *Scene Text*, but not centered.
- Tasks
  - **Localization/Tracking**
  - **Segmentation**
  - **Recognition**
  - **End-to-end**: Simultaneous localization and recognition

---

[10]Dataset Link

# Image Measures

- IoU[11] (or Jaccard Index)

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

- Precision, Recall, F1 measure

$$P = \frac{TP}{TP + FP}$$

$$R = \frac{TP}{TP + FN}$$

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

---

[11] Everingham et al

# BLEU[13] (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.

- $$\text{BLEU}_n(a, b) = \frac{\sum\limits_{w_n \in a} \min\left(c_a(w_n), \max\limits_{j=1,\dots,|b|} c_{b_j}(w_n)\right)}{\sum\limits_{w_n \in a} c_a(w_n)}$$

- BLEU or $\text{BLEU}_{\text{Overall}}$ is a geometric mean of n-gram scores from 1 to 4.

---

[12]Detailed results in: Callison-Burch et. al., 2006; Reiter et. al., 2008; Hodosh et. al., 2013
[13]Papineni et. al., BLEU: A Method for Automatic Evaluation of Machine Translation, 2002

# BLEU[13] (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.

- $$\text{BLEU}_n(a, b) = \frac{\sum\limits_{w_n \in a} \min\left(c_a(w_n), \max\limits_{j=1,\ldots,|b|} c_{b_j}(w_n)\right)}{\sum\limits_{w_n \in a} c_a(w_n)}$$

- BLEU or $\text{BLEU}_{\text{Overall}}$ is a geometric mean of n-gram scores from 1 to 4.
- Strength
  - Automatic, easy to compute
- Weakness[12]
  - 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 $b$.

---

[12]Detailed results in: Callison-Burch et. al., 2006; Reiter et. al., 2008; Hodosh et. al., 2013
[13]Papineni et. al., BLEU: A Method for Automatic Evaluation of Machine Translation, 2002

# Rouge[14] (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.

- $\text{ROUGE}_n (a, b) = \dfrac{\sum\limits_{j=1}^{|b|} \sum\limits_{w_n \in b_j} \min \left( c_a (w_n), c_{b_j} (w_n) \right)}{\sum\limits_{j=1}^{|b|} \sum\limits_{w_n \in b_j} c_{b_j} (w_n)}$

---

[14]Lin et. al., ROUGE: A Package for Automatic Evaluation of Summaries, 2004

# Rouge[14] (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.

- $$\text{ROUGE}_n(a, b) = \frac{\sum\limits_{j=1}^{|b|} \sum\limits_{w_n \in b_j} \min\left(c_a(w_n), c_{b_j}(w_n)\right)}{\sum\limits_{j=1}^{|b|} \sum\limits_{w_n \in b_j} c_{b_j}(w_n)}$$

- There are other variants like $\text{ROUGE}_S$, $\text{ROUGE}_L$

---

[14]Lin et. al., ROUGE: A Package for Automatic Evaluation of Summaries, 2004

# Rouge[14] (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.

- $$\text{ROUGE}_n(a, b) = \frac{\sum\limits_{j=1}^{|b|} \sum\limits_{w_n \in b_j} \min\left(c_a(w_n), c_{b_j}(w_n)\right)}{\sum\limits_{j=1}^{|b|} \sum\limits_{w_n \in b_j} c_{b_j}(w_n)}$$

- There are other variants like $\text{ROUGE}_S$, $\text{ROUGE}_L$

- Similar strengths and weaknesses as BLEU.

---

[14]Lin et. al., ROUGE: A Package for Automatic Evaluation of Summaries, 2004

# METEOR[15] (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.



the cat sat on the mat     the cat sat on the mat

on the mat sat the cat     on the mat sat the cat

\# criss-crosses of left alignment is less

Source: *Wikipedia*

---

[15]Banerjee et. al., METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments, 2005

# METEOR[15] (Metric for Evaluation of Translation with Explicit ORdering)

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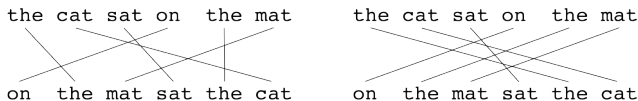
- An alignment between $a$ and $b$ is first computed.



```
the cat sat on  the mat        the cat sat on   the mat



on  the mat sat the cat        on  the mat sat the cat
```
  # criss-crosses of left alignment is less

Source: *Wikipedia*

- $\texttt{METEOR} = \max\limits_{j=1,\dots,|b|} \left( \dfrac{10PR}{R + 9P} \right) \left( 1 - \dfrac{1}{2} \left( \dfrac{\#\text{chunks}}{\#\text{matched unigrams}} \right)^3 \right)$
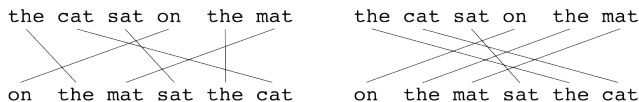
  $P$ = unigram precision, $R$ = unigram recall,
  chunks = set of unigrams adjacent in $a$ and $b_j$ (Example on the right has 3 chunks).

---

[15]Banerjee et. al., METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments, 2005

# METEOR[15] (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.



the cat sat on   the mat       the cat sat on   the mat

on  the mat sat the cat        on  the mat sat the cat

# criss-crosses of left alignment is less

Source: *Wikipedia*

- $\texttt{METEOR} = \max_{j=1,\ldots,|b|} \left( \dfrac{10PR}{R + 9P} \right) \left( 1 - \dfrac{1}{2} \left( \dfrac{\#\text{chunks}}{\#\text{matched unigrams}} \right)^3 \right)$

  $P$ = unigram precision, $R$ = unigram recall,
  chunks = set of unigrams adjacent in $a$ and $b_j$ (Example on the right has 3 chunks).

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

[15]Banerjee et. al., METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments, 2005

# CIDEr[16]

$a$ : candidate sentence, $b$ : set of reference sentences

- $CIDEr_n(a, b) = \dfrac{1}{|b|} \sum\limits_{j=1}^{|b|} \dfrac{\mathbf{g^n}(a) \cdot \mathbf{g^n}(b_j)}{\left\|\mathbf{g^n}(a)\right\| \left\|\mathbf{g^n}(b_j)\right\|}$

  $\mathbf{g^n}(x)$ : vector formed by TF-IDF scores of all n-grams in $x$.

  $CIDEr(a, b) = \sum\limits_{n=1}^{N} w_n CIDEr_n(a, b)$

---

[16]Vedantam et. al., CIDEr: Consensus-based Image Description Evaluation, 2014

# CIDEr[16]

$a$ : candidate sentence, $b$ : set of reference sentences

- $CIDEr_n(a, b) = \dfrac{1}{|b|} \sum\limits_{j=1}^{|b|} \dfrac{\mathbf{g^n}(a) \cdot \mathbf{g^n}(b_j)}{\left\| \mathbf{g^n}(a) \right\| \left\| \mathbf{g^n}(b_j) \right\|}$

  $\mathbf{g^n}(x)$ : vector formed by TF-IDF scores of all n-grams in $x$.

  $CIDEr(a, b) = \sum\limits_{n=1}^{N} w_n CIDEr_n(a, b)$

- Gives more weight-age to *important* n-grams.

- Higher correlation with human consensus scores compared to above metrics.

---

[16]Vedantam et. al., CIDEr: Consensus-based Image Description Evaluation, 2014

- **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%.

---

[17]Hodosh et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

# Ranking based measures

- **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.

---

[17]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.[17] shows that both automatic ranking based measures are *more* robust than metrics that consider only the quality of the first result.

---

[17]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.[18]

---

[18]Hodosh et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

[19]Manning et. al., Introduction to Information Retrieval, 2008

- Measuring quality of a single best result
  - **Rating system of 1-4** from Hodosh et. al.[18]

- Measuring ranked candidates
  - **Success@k** = % image sentence pairs for which at least one relevant result is found in the top-k list.
  - **R-precision**[19] = average % of relevant items in the top-k list.

---

[18]Hodosh et. al., Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics, 2013

[19]Manning et. al., Introduction to Information Retrieval, 2008

# Challenges involving humans

- Datasets
  - Using humans to make binary/choosing decisions, rather than complex decisions. Helps in faster and quality annotation.[20]
  - Games to make the creation of datasets more interesting for annotators.[21]
  - Anyhow involves post-processing to remove spelling mistakes, and sometimes grammatical mistakes.

---

[20] Parikh et. al., 2011; Vedantam et. al., 2014
[21] Deng et. al., 2013; Kazemzadeh et. al., 2014
[22] More details in Reiter et. al., 2008; Hodosh et. al., 2013

# Challenges involving humans

- Datasets
  - Using humans to make binary/choosing decisions, rather than complex decisions. Helps in faster and quality annotation.[20]
  - Games to make the creation of datasets more interesting for annotators.[21]
  - 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.[22]
  - 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.

---

[20]Parikh et. al., 2011; Vedantam et. al., 2014
[21]Deng et. al., 2013; Kazemzadeh et. al., 2014
[22]More details in Reiter et. al., 2008; Hodosh et. al., 2013