

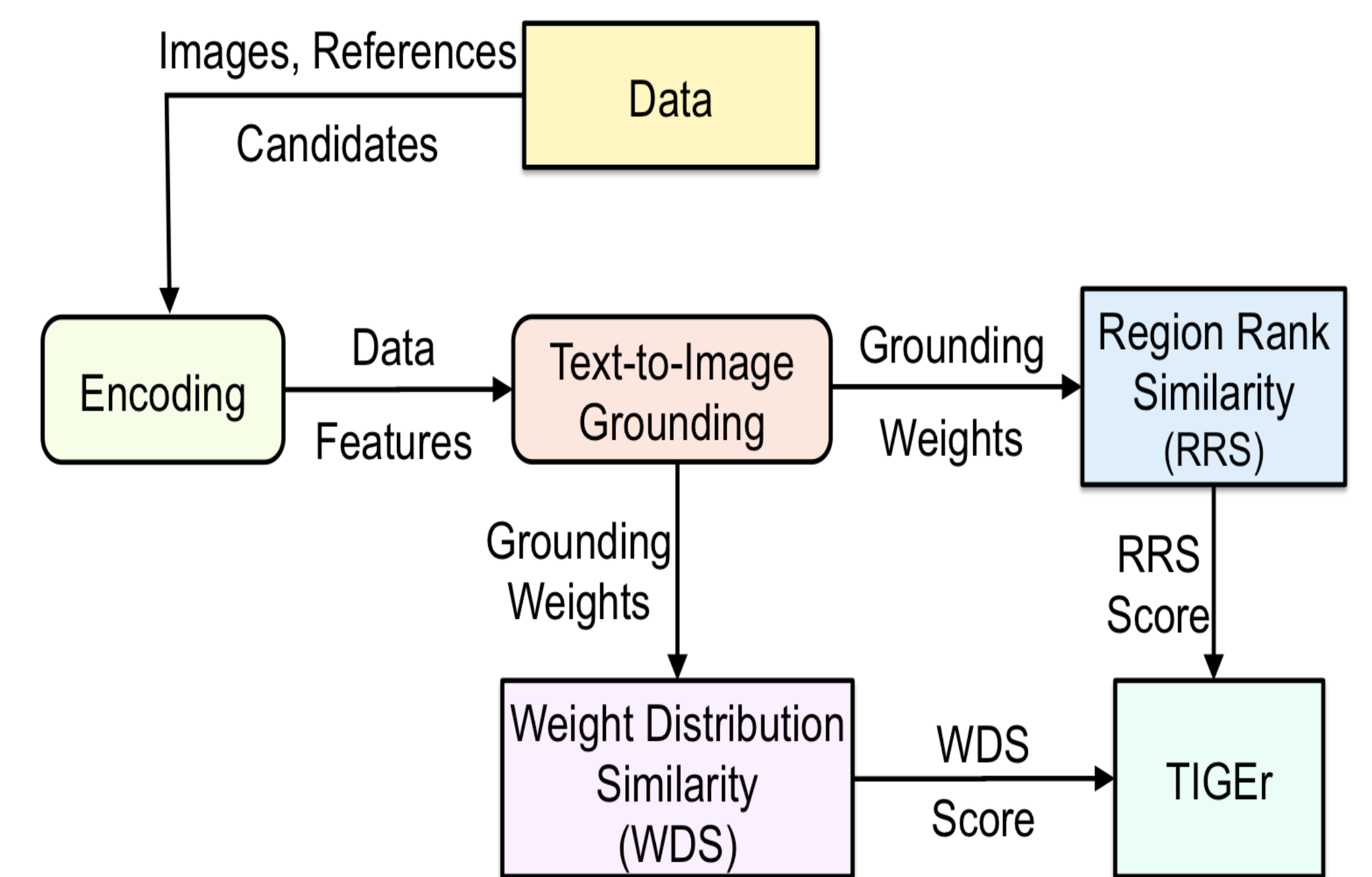
## Motivation and Contribution

References	
	A man with black curly hair is looking at a beer.
	A man is holding a bottle of beer in his hands.
	A man is reading the label on a beer bottle.
Candidate	
(a) A male with curling hair is staring at a drink with alcohol.	
BLEU-4: 0.00	ROUGE-L: 0.53
METEOR: 0.21	CIDEr: 0.33
SPICE: 0.07	TIGEr (ours): 0.88
(b) A woman with straight hair is holding a bottle of beer.	
BLEU-4: 0.54	ROUGE-L: 0.64
METEOR: 0.32	CIDEr: 1.63
SPICE: 0.27	TIGEr (ours): 0.69

- Metrics based on pure text-level comparison **lose image information** and face the challenge of **language ambiguity**.
- Propose a **novel automatic evaluation metric** called **TIGEr**.
  - Consider both image content and human-generated references.
  - Measure the consistence with human attention distribution among image regions.

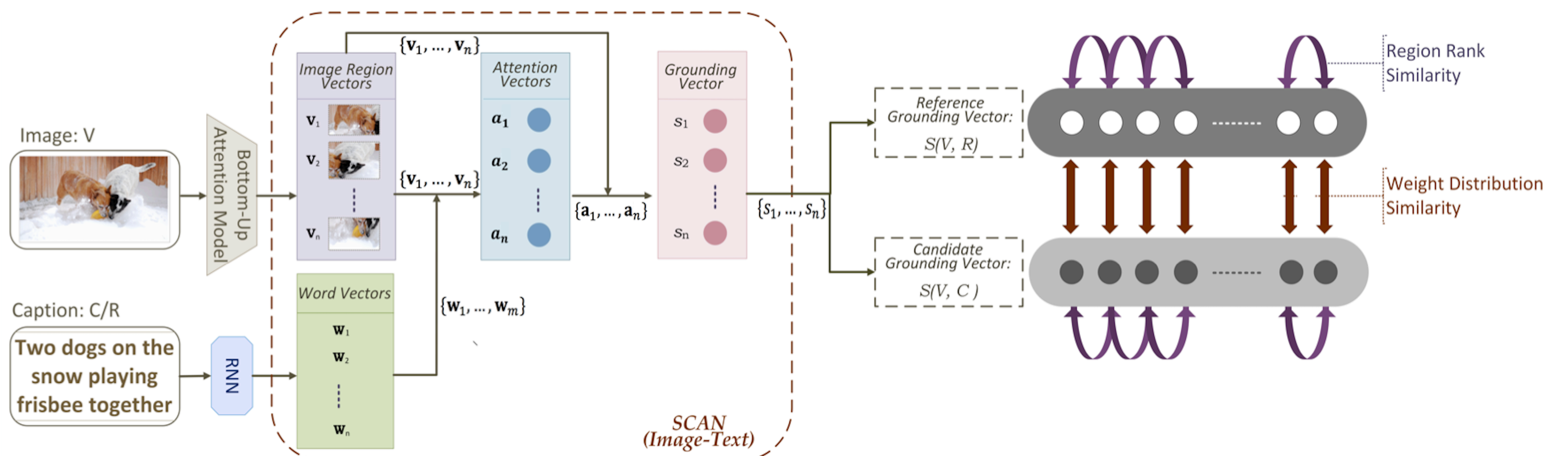
## TIGEr Framework

- Data Encoding**
  - Region-level & Word-level embedding vectors
- Text-to-Image Grounding**
  - Grounding a caption into each image region.
- [(Reference vs. Candidate) | Image]**
  - RRS: how similar is the order of image regions based on grounding weights?
  - WDS: how similar is the attention distributed by a caption among image regions?
- TIGEr**
  - Average value of RRS and WDS



## TIGEr Workflow

- Encoding images and texts by a pre-trained Bottom-Up Attention and a RNN model.
- Grounding texts and images by a pre-trained SCAN model.
- Calculating RRS based on Normalized Discounted Cumulative Gain (NDCG).
- Measuring WDS based on KL Divergence.



## Metric Performance

- TIGEr achieved a noticeable improvement in the assessment of caption quality on three benchmark datasets.
- Identifying irrelevant human-written captions in HI is relatively easy for all metrics, while judging the quality of two correct human-annotated captions in HC is more difficult than other comparison groups.
- Given the change of reference sizes, TIGEr achieves a higher judgment accuracy and more stable performance.

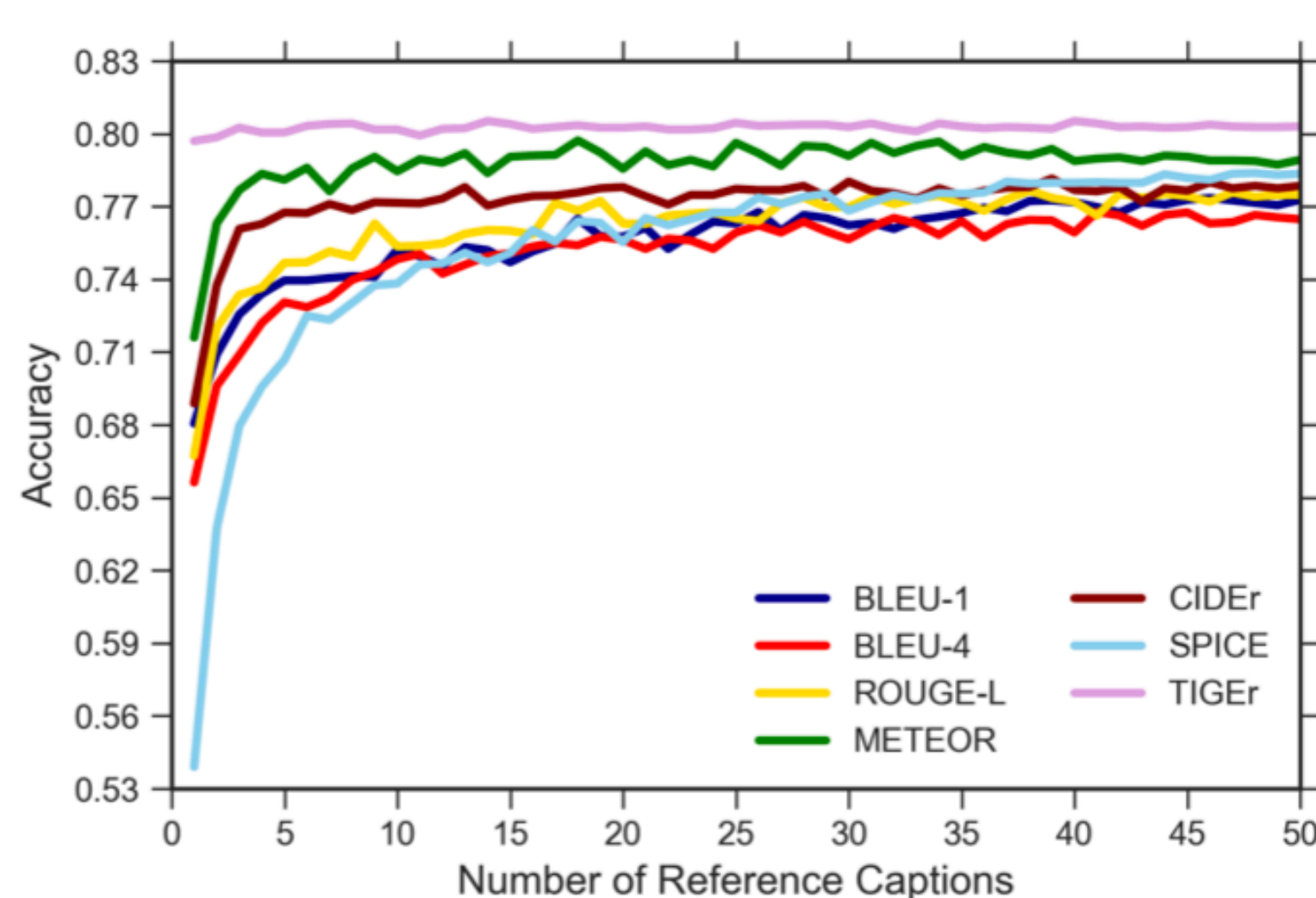
	Composite		Flickr8k	
	$\tau$	$\rho$	$\tau$	$\rho$
BLEU-1	0.280	0.353	0.323	0.404
BLEU-4	0.205	0.352	0.138	0.387
ROUGE-L	0.307	0.383	0.323	0.404
METEOR	0.379	0.469	0.418	0.519
CIDEr	0.378	0.472	0.439	0.542
SPICE	0.419	0.514	0.449	0.596
<b>Ours</b>				
RRS	0.388	0.479	0.418	0.521
WDS	0.433	0.526	0.464	0.572
<b>TIGEr</b>	<b>0.454</b>	<b>0.553</b>	<b>0.493</b>	<b>0.606</b>

	HC	HI	HM	MM	All
BLEU-1	51.20	95.70	91.20	58.20	74.08
BLEU-4	53.00	92.40	86.70	59.40	72.88
ROUGE-L	51.50	94.50	92.50	57.70	74.05
METEOR	<b>56.70</b>	97.60	<b>94.20</b>	63.40	77.98
CIDEr	53.00	98.00	91.50	64.50	76.75
SPICE	52.60	93.90	83.60	48.10	69.55
<b>TIGEr (ours)</b>	56.00	<b>99.80</b>	92.80	<b>74.20</b>	<b>80.70</b>

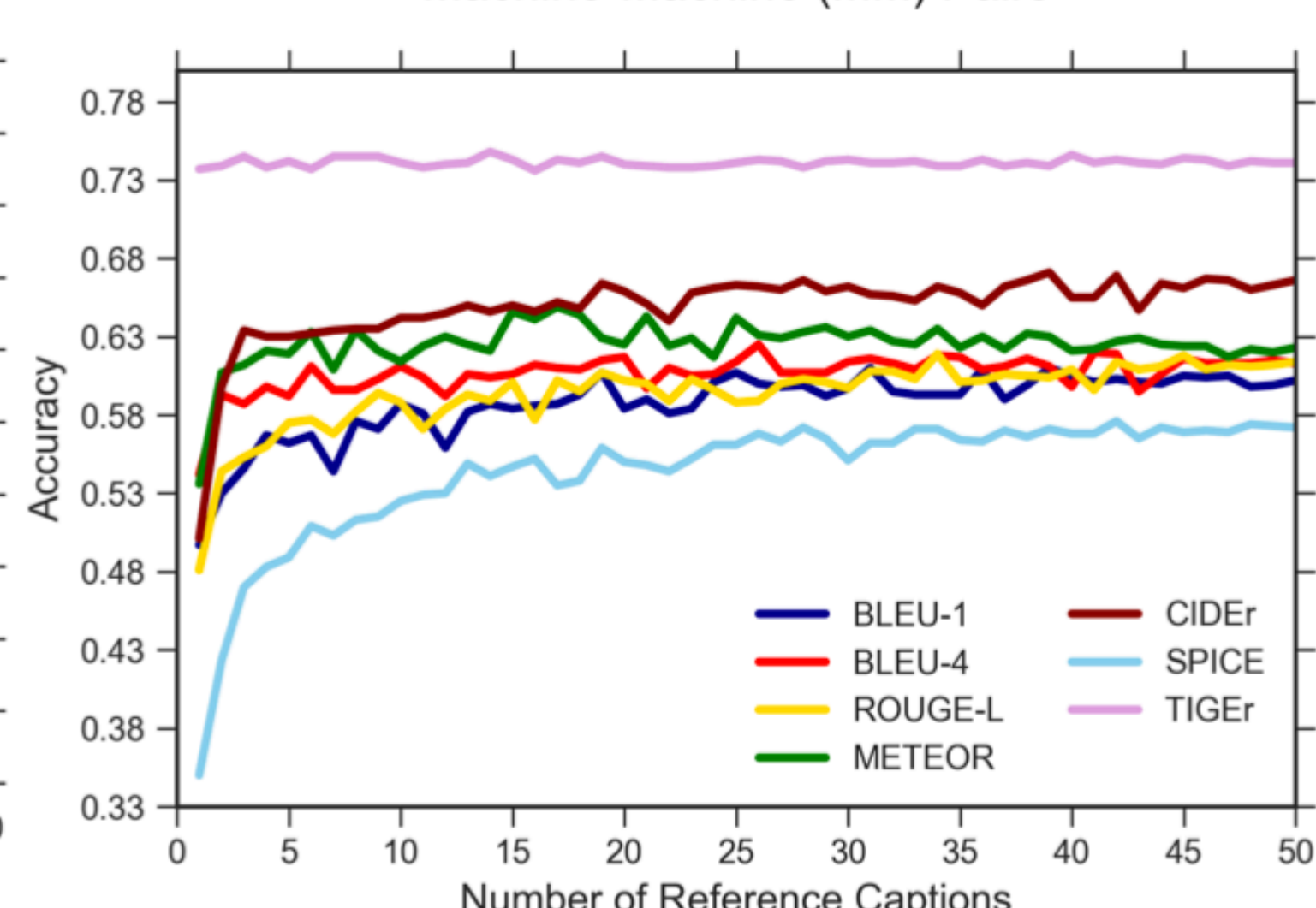
Accuracy of metrics at matching human judgments on PASCAL-50S with 5 reference captions. The highest accuracy per pair type is shown in bold font. HC: human-human correct, HI: human-human incorrect, HM: human-machine, MM: machine-machine, ALL: all pairs.

Caption-level correlation between metrics and human grading scores in Composite and Flickr 8K dataset by using Kendall tau and Spearman rho. All p-values < 0.01.

All Pairs




Machine-Machine (MM) Pairs



## Analysis

- Image region has a higher grounding weight with the corresponding caption than other unrelated regions.
- Text-to-image grounding is more challengeable at action-level compared to object-level.
- Reference captions may not fully cover visual information and TIGEr can measure a caption quality by considering the semantic information of image contents.
- Human interpretation inspired by the image is hard to be judged by an automatic evaluation metric.

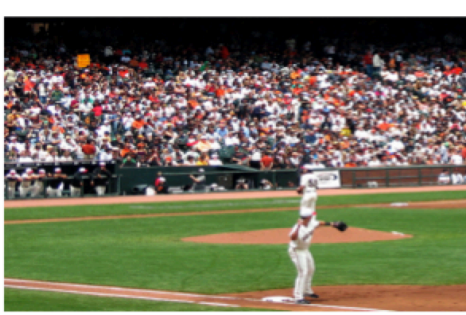


**R:** A curly-haired child is blowing away dandelion seeds while laying in a field of lush green grass.

**Human: 5**

**TIGEr: 1**

**C:** As he lays in the grass, he plucks a spent dandelion, and makes a wish as he blows the little pinwheels into the air.

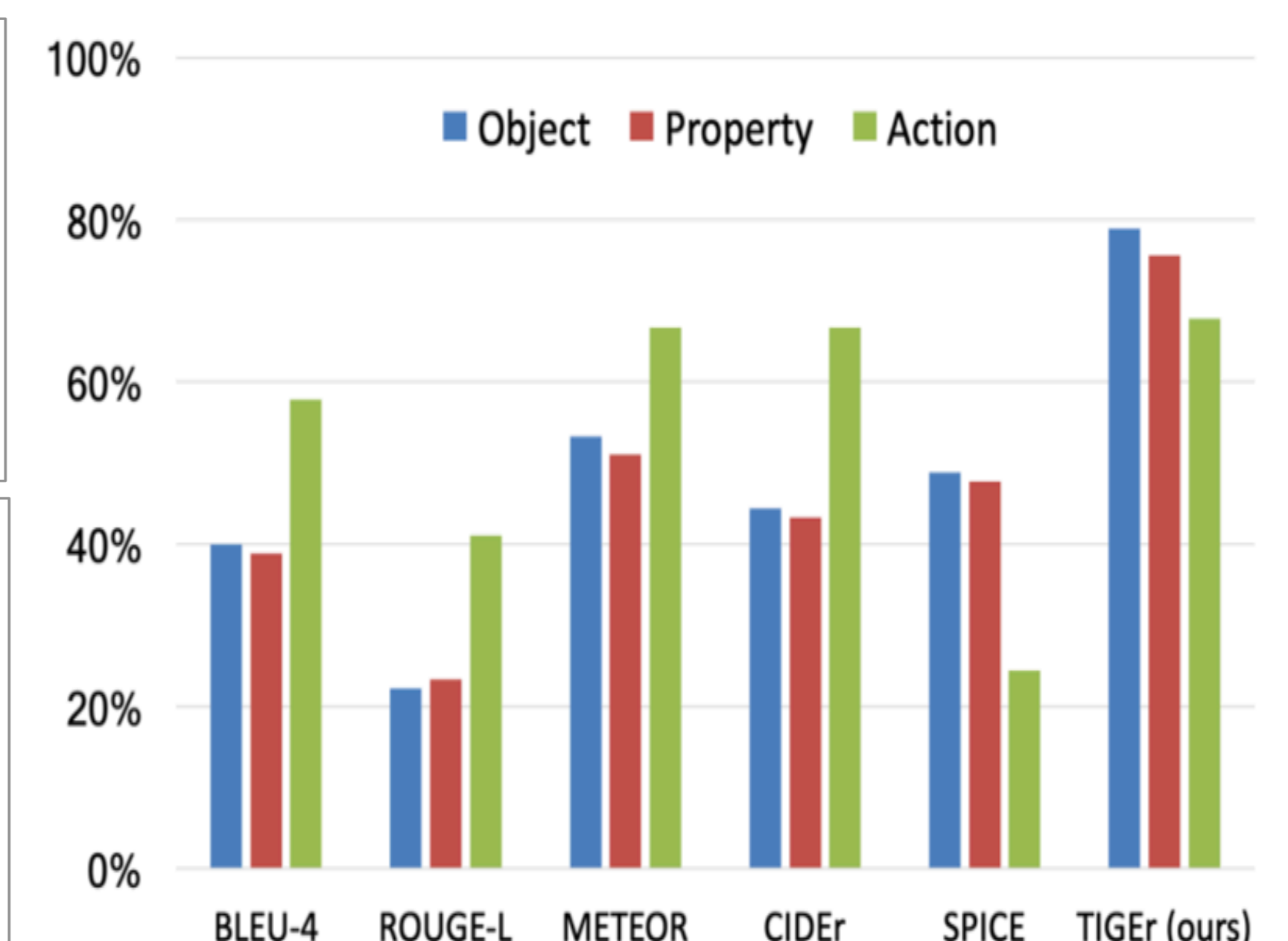



**R:** A baseball player getting ready to throw the ball from a base.

**Human: 1**

**TIGEr: 4**

**C:** A professional baseball game being played at night.






A bunch of people **sit in** an open court yard.

region a: 0.58    region a: 0.24  
region b: 0.12    region b: 0.37

A group of people **walking around** a parking lot.



A young kid hitting a baseball with a bat close to **some chairs**.

region a: 0.16    region a: 0.04  
region b: 0.12    region b: 0.31

A young kid hitting a baseball with a bat close to **a garage**.

## Related Resource

REO-Relevance, Extraness, Omission: A Fine-grained Evaluation for Image Captioning. In EMNLP-IJCNLP'19.

- A fine-grained evaluation on description adequacy
- Candidate vs. Image or (Image + References)

Github Link:

<https://github.com/SeleenaJM/CapEval>

