

Visual Storytelling

Ting-Hao (Kenneth) Huang^{1*}, Francis Ferraro^{2*}, Nasrin Mostafazadeh³, Ishan Misra¹, Aishwarya Agrawal⁴, Jacob Devlin⁶, Ross Girshick⁵, Xiaodong He⁶, Pushmeet Kohli⁶, Dhruv Batra⁴, C. Lawrence Zitnick⁵, Devi Parikh⁴, Lucy Vanderwende⁶, Michel Galley⁶, Margaret Mitchell⁶

Microsoft Research

1 Carnegie Mellon University, 2 Johns Hopkins University, 3 University of Rochester,

4 Virginia Tech, 5 Facebook AI Research

6 Corresponding authors: {jdevlin,lucyv,mgalley,memitc}@microsoft.com

Abstract

We introduce the first dataset for **sequential vision-to-language**, and explore how this data may be used for the task of *visual storytelling*. The first release of this dataset, SIND¹ v.1, includes 81,743 unique photos in 20,211 sequences, aligned to both descriptive (caption) and story language. We establish several strong baselines for the storytelling task, and motivate an automatic metric to benchmark progress. Modelling concrete description as well as figurative and social language, as provided in this dataset and the storytelling task, has the potential to move artificial intelligence from basic understandings of typical visual scenes towards more and more human-like understanding of grounded event structure and subjective expression.

1 Introduction

Beyond understanding simple objects and concrete scenes lies interpreting causal structure; making sense of visual input to tie disparate moments together as they give rise to a cohesive narrative of events through time. This requires moving from reasoning about single images – static moments, devoid of context – to sequences of images that depict events as they occur and change. On the vision side, progressing from single images to images in context allows us to begin to create an artificial intelligence (AI) that can reason about a visual moment given what it has already seen. On the language side, progressing from literal description to narrative helps to learn more evaluative, conversational, and abstract

*T.H. and F.F. contributed equally to this work.

¹Sequential Images Narrative Dataset. This and future releases are made available on www.sind.ai.

			
DII	A group of people that are sitting next to each other.	Adult male wearing sunglasses lying down on black pavement.	The sun is setting over the ocean and mountains.
SIS	Having a good time bonding and talking.	[M] got exhausted by the heat.	Sky illuminated with a brilliance of gold and orange hues.

Figure 1: Example language difference between descriptions for images in isolation (DII) vs. stories for images in sequence (SIS).

language. This is the difference between, for example, “sitting next to each other” versus “having a good time”, or “sun is setting” versus “sky illuminated with a brilliance...” (see Figure 1). The first descriptions capture image content that is literal and concrete; the second requires further inference about what a *good time* may look like, or what is special and worth sharing about a particular sunset.

We introduce the first dataset of sequential images with corresponding descriptions, which captures some of these subtle but important differences, and advance the task of visual storytelling. We release the data in three tiers of language for the same images: (1) **Descriptions of images-in-isolation (DII)**; (2) **Descriptions of images-in-sequence (DIS)**; and (3) **Stories for images-in-sequence (SIS)**. This tiered approach reveals the effect of temporal context and the effect of narrative language. As all the tiers are aligned to the same images, the dataset facilitates directly modeling the relationship between literal and more abstract visual concepts, including the relationship between visual imagery and typical event patterns. We additionally propose an automatic evaluation metric which is best

beach (684)	breaking up (350)	easter (259)
amusement park (525)	carnival (331)	church (243)
building a house (415)	visit (321)	graduation ceremony (236)
party (411)	market (311)	office (226)
birthday (399)	outdoor activity (267)	father's day (221)

Table 1: The number of albums in our tiered dataset for the 15 most frequent kinds of stories.

correlated with human judgments, and establish several strong baselines for the visual storytelling task.

2 Motivation and Related Work

Work in vision to language has exploded, with researchers examining image captioning (Lin et al., 2014; Karpathy and Fei-Fei, 2015; Vinyals et al., 2015; Xu et al., 2015; Chen et al., 2015; Young et al., 2014; Elliott and Keller, 2013), question answering (Antol et al., 2015; Ren et al., 2015; Gao et al., 2015; Malinowski and Fritz, 2014), visual phrases (Sadeghi and Farhadi, 2011), video understanding (Ramanathan et al., 2013), and visual concepts (Krishna et al., 2016; Fang et al., 2015).

Such work focuses on direct, literal description of image content. While this is an encouraging first step in connecting vision and language, it is far from the capabilities needed by intelligent agents for naturalistic interactions. There is a significant difference, yet unexplored, between remarking that a visual scene shows “sitting in a room” – typical of most image captioning work – and that the same visual scene shows “bonding”. The latter description is grounded in the visual signal, yet it brings to bear information about social relations and emotions that can be additionally inferred in context (Figure 1). Visually-grounded stories facilitate more evaluative and figurative language than has previously been seen in vision-to-language research: If a system can recognize that colleagues look *bored*, it can remark and act on this information directly.

Storytelling itself is one of the oldest known human activities (Wiessner, 2014), providing a way to educate, preserve culture, instill morals, and share advice; focusing AI research towards this task therefore has the potential to bring about more human-like intelligence and understanding.

3 Dataset Construction

Extracting Photos We begin by generating a list of “storyable” event types. We leverage the idea that

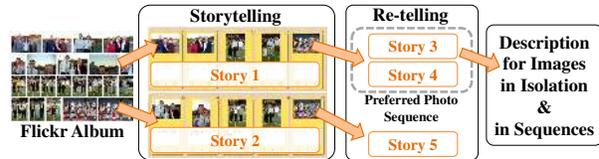


Figure 2: Dataset crowdsourcing workflow.



Figure 3: Interface for the *Storytelling* task, which contains: 1) the photo album, and 2) the storyboard.

“storyable” events tend to involve some form of possession, e.g., “John’s birthday party,” or “Shabnam’s visit.” Using the Flickr data release (Thomee et al., 2015), we aggregate 5-grams of photo titles and descriptions, using Stanford CoreNLP (Manning et al., 2014) to extract possessive dependency patterns. We keep the heads of possessive phrases if they can be classified as an *EVENT* in WordNet3.0, relying on manual winnowing to target our collection efforts.² These terms are then used to collect albums using the Flickr API.³ We only include albums with 10 to 50 photos where all album photos are taken within a 48-hour span and CC-licensed. See Table 1 for the query terms with the most albums returned.

The photos returned from this stage are then presented to crowd workers using Amazon’s Mechanical Turk to collect the corresponding stories and descriptions. The crowdsourcing workflow of developing the complete dataset is shown in Figure 2.

Crowdsourcing Stories In Sequence We develop a 2-stage crowdsourcing workflow to collect naturalistic stories with text aligned to images. The first stage is *storytelling*, where the crowd worker selects a subset of photos from a given album to form a photo sequence and writes a story about it (see Figure 3). The second stage is *re-telling*, in which the worker writes a story based on one photo sequence

²We simultaneously supplemented this data-driven effort by a small hand-constructed gazetteer.

³<https://www.flickr.com/services/api/>

					
DII	A black frisbee is sitting on top of a roof.	A man playing soccer outside of a white house with a red door.	The boy is throwing a soccer ball by the red door.	A soccer ball is over a roof by a frisbee in a rain gutter.	Two balls and a frisbee are on top of a roof.
DIS	A roof top with a black frisbee laying on the top of the edge of it.	A man is standing in the grass in front of the house kicking a soccer ball.	A man is in the front of the house throwing a soccer ball up	A blue and white soccer ball and black Frisbee are sitting on top of the roof top.	Two soccer balls and a Frisbee are sitting on top of the roof top.
SIS	A discuss got stuck up on the roof.	Why not try getting it down with a soccer ball?	Up the soccer ball goes.	It didn't work so we tried a volley ball.	Now the discuss, soccer ball, and volleyball are all stuck on the roof.

Figure 4: Example descriptions of images in isolation (DII); descriptions of images in sequence (DIS); and stories of images in sequence (SIS).

generated by workers in the first stage.

In both stages, all album photos are displayed in the order of the time that the photos were taken, with a “storyboard” underneath. In *storytelling*, by clicking a photo in the album, a “story card” of the photo appears on the storyboard. The worker is instructed to pick at least five photos, arrange the order of selected photos, and then write a sentence or a phrase on each card to form a story; this appears as a full story underneath the text aligned to each image. Additionally, this interface captures the alignments between text and photos. Workers may skip an album if it does not seem storyable (e.g., a collection of coins). Albums skipped by two workers are discarded. The interface of *re-telling* is similar, but it displays the two photo sequences already created in the first stage, which the worker chooses from to write the story. For each album, 2 workers perform *storytelling* (at \$0.3/HIT), and 3 workers perform *re-telling* (at \$0.25/HIT), yielding a total of 1,907 workers. All HITs use quality controls to ensure varied text at least 15 words long.

Crowdsourcing Descriptions of Images In Isolation & Images In Sequence We also use crowdsourcing to collect descriptions of images-in-isolation (DII) and descriptions of images-in-sequence (DIS), for the photo sequences with stories from a majority of workers in the first task (as Figure 2). In both DII and DIS tasks, workers are asked to follow the instructions for image captioning proposed in MS COCO (Lin et al., 2014) such as *describe all the important parts*. In DII, we use

Data Set	#(Txt, Img) Pairs (k)	Vocab Size (k)	Avg. #Tok	%Abs	Frazier	Yngve	Ppl
Brown	52.1	47.7	20.8	15.2%	18.5	77.2	194.0
DII	151.8	13.8	11.0	21.3%	10.3	27.4	147.0
DIS	151.8	5.0	9.8	24.8%	9.2	23.7	146.8
SIS	252.9	18.2	10.2	22.1%	10.5	27.5	116.0

Table 2: A summary of our dataset, following the proposed analyses of Ferraro et al. (2015), including the Frazier and Yngve measures of syntactic complexity. The balanced Brown corpus (Marcus et al., 1999), provided for comparison, contains only text. Perplexity (Ppl) is calculated against a 5-gram language model learned on a generic 30B English words dataset scraped from the web.

Desc.-in-Iso.			Desc.-in-Seq.			Story-in-Seq.		
man	sitting	black	chatting	amount	trunk	went	[female]	see
woman	white	large	gentleman	goers	facing	got	today	saw
standing	two	front	enjoys	sofa	bench	[male]	decided	came
holding	young	group	folks	egg	enjoying	took	really	started
wearing	image		shoreline	female		great	time	

Table 3: Top words ranked by normalized PMI.

the MS COCO image captioning interface.⁴ In DIS, we use the storyboard and story cards of our *storytelling* interface to display a photo sequence, with MS COCO instructions adapted for sequences. We recruit 3 workers for DII (at \$0.05/HIT) and 3 workers for DIS (at \$0.07/HIT).

Data Post-processing We tokenize all storylets and descriptions with the CoreNLP tokenizer, and replace all people names with generic MALE/FEMALE tokens,⁵ and all identified named entities with their entity type (e.g., location). The data is released as *training*, *validation*, and *test* following an 80%/10%/10% split on the stories-in-sequence albums. Example language from each tier is shown in Figure 4.

4 Data Analysis

Our dataset includes 10,117 Flickr albums with 210,819 unique photos. Each album on average has 20.8 photos ($\sigma = 9.0$). The average time span of each album is 7.9 hours ($\sigma = 11.4$). Further details of each tier of the dataset are shown in Table 2.⁶

We use normalized pointwise mutual information to identify the words most closely associated with each tier (Table 3). Top words for descriptions-

⁴<https://github.com/tylin/coco-ui>

⁵We use those names occurring at least 10,000 times. <https://ssa.gov/oact/babynames/names.zip>

⁶We exclude words seen only once.

	METEOR	BLEU	Skip-Thoughts
r	0.22 (2.8e-28)	0.08 (1.0e-06)	0.18 (5.0e-27)
ρ	0.20 (3.0e-31)	0.08 (8.9e-06)	0.16 (6.4e-22)
τ	0.14 (1.0e-33)	0.06 (8.7e-08)	0.11 (7.7e-24)

Table 4: Correlations of automatic scores against human judgements, with p-values in parentheses.

in-isolation reflect an impoverished disambiguating context: References to people often lack social specificity, as people are referred to as simply “man” or “woman”. Single images often do not convey much information about underlying events or actions, which leads to the abundant use of posture verbs (“standing”, “sitting”, etc.). As we turn to descriptions-in-sequence, these relatively uninformative words are much less represented. Finally, top story-in-sequence words include more storytelling elements, such as names (*[male]*), temporal references (*today*) and words that are more dynamic and abstract (*went, decided*).

5 Automatic Evaluation Metric

Given the nature of the complex storytelling task, the best and most reliable evaluation for assessing the quality of generated stories is human judgment. However, automatic evaluation metrics are useful to quickly benchmark progress. To better understand which metric could serve as a proxy for human evaluation, we compute pairwise correlation coefficients between automatic metrics and human judgments on 3,000 stories sampled from the SIS training set.

For the human judgements, we again use crowdsourcing on MTurk, asking five judges per story to rate how strongly they agreed with the statement “If these were my photos, I would like using a story like this to share my experience with my friends”.⁷ We take the average of the five judgments as the final score for the story. For the automatic metrics, we use METEOR,⁸ smoothed-BLEU (Lin and Och, 2004), and Skip-Thoughts (Kiros et al., 2015) to compute similarity between each story for a given sequence. Skip-thoughts provide a Sentence2Vec embedding which models the semantic space of novels.

As Table 4 shows, METEOR correlates best with human judgment according to all the correlation co-

⁷Scale presented ranged from “Strongly disagree” to “Strongly agree”, which we convert to a scale of 1 to 5.

⁸We use METEOR version 1.5 with `hter` weights.

Beam=10	Greedy	-Dups	+Grounded
23.55	19.10	19.21	–

Table 6: Captions generated per-image with METEOR scores.

Beam=10	Greedy	-Dups	+Grounded
23.13	27.76	30.11	31.42

Table 7: Stories baselines with METEOR scores.

efficients. This signals that a metric such as METEOR which incorporates paraphrasing correlates best with human judgement on this task. A more detailed study of automatic evaluation of stories is an area of interest for a future work.

6 Baseline Experiments

We report baseline experiments on the storytelling task in Table 7, training on the SIS tier and testing on half the SIS validation set (valtest). Example output from each system is presented in Table 5. To highlight some differences between story and caption generation, we also train on the DII tier in isolation, and produce captions per-image, rather than in sequence. These results are shown in Table 7.

To train the story generation model, we use a sequence-to-sequence recurrent neural net (RNN) approach, which naturally extends the single-image captioning technique of Devlin et al. (2015) and Vinyals et al. (2014) to multiple images. Here, we encode an image *sequence* by running an RNN over the `fc7` vectors of each image, in reverse order. This is used as the initial hidden state to the story decoder model, which learns to produce the story one word at a time using softmax loss over the training data vocabulary. We use Gated Recurrent Units (GRUs) (Cho et al., 2014) for both the image encoder and story decoder.

In the baseline system, we generate the story using a simple beam search (size=10), which has been successful in image captioning previously (Devlin et al., 2015). However, for story generation, the results of this model subjectively appear to be very poor – the system produces generic, repetitive, high-level descriptions (e.g., “This is a picture of a dog”). This is a predictable result given the label bias problem inherent in maximum likelihood training; recent work has looked at ways to address this issue directly (Li et al., 2016).



+ <i>Viterbi</i>	This is a picture of a family. This is a picture of a cake. This is a picture of a dog. This is a picture of a beach. This is a picture of a beach.
+ <i>Greedy</i>	The family gathered together for a meal. The food was delicious. The dog was excited to be there. The dog was enjoying the water. The dog was happy to be in the water.
- <i>Dups</i>	The family gathered together for a meal. The food was delicious. The dog was excited to be there. The kids were playing in the water. The boat was a little too much to drink.
+ <i>Grounded</i>	The family got together for a cookout. They had a lot of delicious food. The dog was happy to be there. They had a great time on the beach. They even had a swim in the water.

Table 5: Example stories generated by baselines.

To establish a stronger baseline, we explore several decode-time heuristics to improve the quality of the generated story. The first heuristic is to lower the decoder beam size substantially. We find that using a beam size of 1 (greedy search) significantly increases the story quality, resulting in a 4.6 gain in METEOR score. However, the same effect is not seen for caption generation, with the greedy caption model obtaining worse quality than the beam search model. This highlights a key difference in generating stories versus generating captions.

Although the stories produced using a greedy search result in significant gains, they include many repeated words and phrases, e.g., “The kids had a great time. And the kids had a great time.” We introduce a very simple heuristic to avoid this, where the same content word cannot be produced more than once within a given story. This improves METEOR by another 2.3 points.

An advantage of comparing captioning to storytelling side-by-side is that the captioning output may be used to help inform the storytelling output. To this end, we include an additional baseline where “visually grounded” words may only be produced if they are licensed by the caption model. We define the set of visually grounded words to be those which occurred at higher frequency in the caption training than the story training:

$$\frac{P(w|T_{caption})}{P(w|T_{story})} > 1.0 \quad (1)$$

We train a separate model using the caption annotations, and produce an n-best list of captions for each image in the valtest set. Words seen in at

least 10 sentences in the 100-best list are marked as ‘licensed’ by the caption model. Greedy decoding without duplication proceeds with the additional constraint that if a word is visually grounded, it can only be generated by the story model if it is licensed by the caption model for the same photo set. This results in a further 1.3 METEOR improvement.

It is interesting to note what a strong effect relatively simple heuristics have on the generated stories. We do not intend to suggest that these heuristics are the *right* way to approach story generation. Instead, the main purpose is to provide clear baselines that demonstrate that story generation has fundamentally different challenges from caption generation; and the space is wide open to explore for training and decoding methods to generate fluent stories.

7 Conclusion and Future Work

We have introduced the first dataset for **sequential vision-to-language**, which incrementally moves from images-in-isolation to stories-in-sequence. We argue that modelling the more figurative and social language captured in this dataset is essential for evolving AI towards more human-like understanding. We have established several strong baselines for the task of visual storytelling, and have motivated METEOR as an automatic metric to evaluate progress on this task moving forward.

References

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and

- Devi Parikh. 2015. Vqa: Visual question answering. In *International Conference on Computer Vision (ICCV)*.
- Jianfu Chen, Polina Kuznetsova, David Warren, and Yejin Choi. 2015. Déjà image-captions: A corpus of expressive descriptions in repetition. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 504–514, Denver, Colorado, May–June. Association for Computational Linguistics.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *CoRR*.
- Jacob Devlin, Hao Cheng, Hao Fang, Saurabh Gupta, Li Deng, Xiaodong He, Geoffrey Zweig, and Margaret Mitchell. 2015. Language models for image captioning: The quirks and what works. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 100–105, Beijing, China, July. Association for Computational Linguistics.
- Desmond Elliott and Frank Keller. 2013. Image description using visual dependency representations. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1292–1302, Seattle, Washington, USA, October. Association for Computational Linguistics.
- Hao Fang, Saurabh Gupta, Forrest N. Iandola, Rupesh Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C. Platt, C. Lawrence Zitnick, and Geoffrey Zweig. 2015. From captions to visual concepts and back. In *Computer Vision and Pattern Recognition (CVPR)*.
- Francis Ferraro, Nasrin Mostafazadeh, Ting-Hao K. Huang, Lucy Vanderwende, Jacob Devlin, Michel Galley, and Margaret Mitchell. 2015. A survey of current datasets for vision and language research. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 207–213, Lisbon, Portugal, September. Association for Computational Linguistics.
- Haoyuan Gao, Junhua Mao, Jie Zhou, Zhiheng Huang, Lei Wang, and Wei Xu. 2015. Are you talking to a machine? dataset and methods for multilingual image question. In C. Cortes, N.D. Lawrence, D.D. Lee, M. Sugiyama, R. Garnett, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2287–2295. Curran Associates, Inc.
- Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June.
- Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors. In C. Cortes, N.D. Lawrence, D.D. Lee, M. Sugiyama, R. Garnett, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 3276–3284. Curran Associates, Inc.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalanditis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. 2016. Visual genome: Connecting language and vision using crowdsourced dense image annotations.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. *NAACL HLT 2016*.
- Chin-Yew Lin and Franz Josef Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of the 42Nd Annual Meeting on Association for Computational Linguistics, ACL '04*, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014*, pages 740–755. Springer.
- Mateusz Malinowski and Mario Fritz. 2014. A multi-world approach to question answering about real-world scenes based on uncertain input. In Z. Ghahramani, M. Welling, C. Cortes, N.D. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1682–1690. Curran Associates, Inc.
- Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60.
- Mitchell Marcus, Beatrice Santorini, Mary Ann Marcinkiewicz, and Ann Taylor. 1999. Brown corpus, treebank-3.
- Vignesh Ramanathan, Percy Liang, and Li Fei-Fei. 2013. Video event understanding using natural language descriptions. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, pages 905–912. IEEE.
- Mengye Ren, Ryan Kiros, and Richard Zemel. 2015. Exploring models and data for image question answering.

- In C. Cortes, N.D. Lawrence, D.D. Lee, M. Sugiyama, R. Garnett, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2935–2943. Curran Associates, Inc.
- Mohammad Amin Sadeghi and Ali Farhadi. 2011. Recognition using visual phrases. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1745–1752. IEEE.
- Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. 2015. The new data and new challenges in multimedia research. *arXiv preprint arXiv:1503.01817*.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2014. Show and tell: a neural image caption generator. In *CVPR*.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Computer Vision and Pattern Recognition*.
- Polly W Wiessner. 2014. Embers of society: Firelight talk among the ju/hoansi bushmen. *Proceedings of the National Academy of Sciences*, 111(39):14027–14035.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. *arXiv preprint arXiv:1502.03044*.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78.