

Story Generation from Sequence of Independent Short Descriptions

Parag Jain[†] Priyanka Agrawal[†] Abhijit Mishra[†] Mohak Sukhwani[†]
Anirban Laha[†] Karthik Sankaranarayanan[†]

[†]IBM Research India,
{pajain06,priyanka.agrawal,abhijimi,mosukhwa,anirlaha,kartsank}
@in.ibm.com

ABSTRACT

Existing Natural Language Generation (NLG) systems are weak AI systems and exhibit limited capabilities when language generation tasks demand higher levels of creativity, originality and brevity. Effective solutions or, at least evaluations of modern NLG paradigms for such creative tasks have been elusive, unfortunately. This paper introduces and addresses the **task of coherent story generation from independent descriptions**, describing a scene or an event. Towards this, we explore along two popular text-generation paradigms – (1) Statistical Machine Translation (SMT), posing story generation as a translation problem and (2) Deep Learning, posing story generation as a sequence to sequence learning problem. In SMT, we chose two popular methods such as phrase based SMT (PB-SMT) and syntax based SMT (SYNTAX-SMT) to ‘translate’ the incoherent input text into stories. We then implement a deep recurrent neural network (RNN) architecture that encodes sequence of variable length input descriptions to corresponding latent representations and decodes them to produce well formed comprehensive story like summaries. The efficacy of the suggested approaches is demonstrated on a publicly available dataset with the help of popular machine translation and summarization evaluation metrics. We believe, a system like ours has different interesting applications- for example, creating news articles from phrases of event information.

CCS CONCEPTS

•Computing methodologies → Natural language generation;
Supervised learning; Neural networks;

KEYWORDS

Story, Natural Language Generation, Sequential Learning

ACM Reference format:

Parag Jain[†] Priyanka Agrawal[†] Abhijit Mishra[†] Mohak Sukhwani[†]
Anirban Laha[†] Karthik Sankaranarayanan[†]

[†]IBM Research India,

{pajain06,priyanka.agrawal,abhijimi,mosukhwa,anirlaha,kartsank}

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SIGKDD’17, Halifax - Canada

© 2017 Copyright held by the owner/author(s). 123-4567-24-567/08/06...\$15.00

DOI: 10.475/123.4

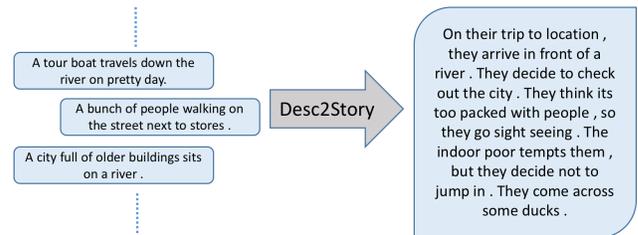


Figure 1: Our system generates variable length stories for variable length input descriptions. The standalone textual descriptions describing a scene or event are converted to human like coherent summaries.

@in.ibm.com. 2017. Story Generation from Sequence of Independent Short Descriptions. In *Proceedings of Workshop on Machine Learning for Creativity, Halifax - Canada, Aug 2017 (SIGKDD’17)*, 7 pages.
DOI: 10.475/123.4

1 INTRODUCTION

Recent advances in machine learning based approaches for natural language generation have led to exploration of many diverse but related text generation tasks. However, the existing systems/ approaches can be classified as weak AI systems. According to the classical definition [14], a strong AI based NLG system should perform language generation in the same manner, expressing similar levels of creativity, originality and brevity as humans. We proceed towards building such systems for the difficult task of automatic story generation that demands the above mentioned human qualities.

We introduce the task of generating coherent narratives form a sequence of independent short descriptions, as shown in Figure 1. The standalone descriptions (in general), although are sufficiently informative enough to describe a scene or an event, they lack the flavor of temporal context and human touch, viz. describing a Christmas celebration scene as ‘*The church was all ready for the big day. The candles were lit and people were arriving. In came the alter boys. They began to read from the book. They put the lights out and the candle were light*’ is far more apt and intriguing than describing it by independent sentences - ‘*people are in a large candlelit room*

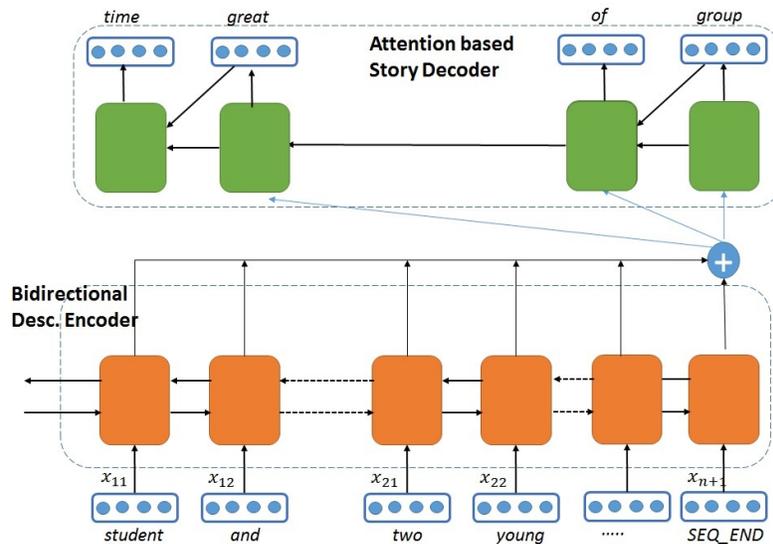


Figure 2: Sequence to Sequence Recurrent Neural Network Architecture for Desc2Story. Orange and green colored blocks represent GRU cells in encoder and decoder respectively.

in which there is a Christmas tree., *ritual candles are lit on stands in a darkened room.*, *a priest is holding a special ceremony.* etc.

The proposed system of story generation is framed as a sequence-to-sequence neural machine learning network that learns latent representations of the input descriptions to generate output summaries. The success of RNNs to model sequential language tasks have made it a de-facto choice of NLP community. From character level language models [28] to task of machine translations [35], RNNs have outperformed all of its present day contemporaries. We envision an end-to-end RNN system that understands the context holistically and generates cohesive textual stories. The input to our system is a set of short textual descriptions, each describing a scene or an event and the system produces semantically rich human like stories. Figure 1 demonstrates a working example.

Our work is motivated by two independent considerations – (1) Introduce the concept of ‘creativity’ for the machine generated text, i.e. a set of same input descriptions should have ability to generate varied story-line when parameterized over user choice of ‘theme’. (2) Leverage underlying ephemeral cues in input text to generate contextually relevant and coherent output stories. We take a step towards this by proposing a new computational regime of creative text content generation. We draw in our parallels from neural language generation modeling and more specifically sequence-to-sequence recurrent neural network, SEQ2SEQ. Unlike phrase Based [39] and syntax based [19] SMT which translate words and phrases independently within a sentence, SEQ2SEQ considers entire sentence as a unit to be translated. We leverage SEQ2SEQ architecture *i.e.*, deep recurrent neural network (RNN) with attention mechanism for the problem of coherent story generation. Additionally, we also pose story generation as a translation problem and report results for both phrase based and syntax based SMT towards this.

We provide an outline of all explored approaches and our network design choices in Section 3. This is followed by detailed Sections 4-5 describing our experimental setup, results and analysis respectively on publicly available dataset. Lastly we describe avenues where the proposed schema can be utilized along with conclusions and future directions in Section 6.

2 RELATED WORK

Text based story generation as a cognitive task has seen considerable traction over last many years. The onslaught of deep learning architectures and advancements in hardware infrastructure has made this task much more relevant in today’s time. Image captioning [24, 25, 38], video understanding [10, 31, 37], narratives in virtual environments [1] to novellas by computers [5, 21] are all manifestations of present day interest AI community has over text based story generating creative systems.

Computational models to mimic human creativity (here: *Story Generation*) are not new. In past storytelling has been approached as an analytical activity of discovering best fact to present to the listener [26]. In an attempt to add engagement, reflections and make stories user centric, more sensitive models modifying stories with human feedback were proposed [3, 27]. The traditional approaches were succeeded by much more sophisticated virtual environments with multi agents framework keeping tracks of plots, narratives, user interest etc. [1, 9, 34]. Text based interactive narratives are the other kind of storytelling variations that have been tried in past. Multitude of approaches ranging from evolutionary (genetic) algorithms [20] to mining of crowd sourced information [15, 33] from web have also been explored.

Attempts to generate cohesive comprehensible stories have not gone beyond short text snippets. Recent works on Visual storytelling [12, 18] aim at generating story by harnessing rich cognitive

	# Docs	Avg. # Sents		Avg. # Words		Avg. # non-overlapping Words
		Caption	Story	Caption	Story	(excluding stop words)
Train	32987	5	6	52	56	23
Val	4168	5	6	51	57	22
Test	4145	5	6	51	56	23

Table 1: Statistics of the dataset used for experimentation

information from a sequence of images. In this paper, we introduce and address a much tougher setting – starting from incoherent description, this is the first known work that discuss the prospects of converting it to human like stories. Our setting is deprived of the rich contextual information that is available in the form of images in traditional Visual storytelling [12]. Such a system, however, has many interesting applications- for example, creating news articles from phrases of event information.

This paper is a preliminary study on story generation from textual input and explores the effectiveness of two traditional language generation paradigms. Our experiment section presents both qualitative and quantitative analyses of the different approaches in the respective paradigms.

3 EXPLORED APPROACHES FOR STORY GENERATION

Story generation from independent textual-descriptions is the process of transforming one form of text into another. Towards this, we explore along two popular text-generation paradigms -(1) Statistical Machine Translation (SMT), posing story generation as a translation problem and (2) Deep Learning, posing story generation as a sequence to sequence learning problem. In SMT, we chose two popular methods such as phrase based SMT and syntax based SMT where as in sequence to sequence learning, we implement different variants of recurrent neural networks (RNN), empowered with attention mechanism.

3.1 Phrase based Statistical Machine Translation PB-SMT

Phrase based Statistical Machine Translation (PB-SMT) deals with (a) finding best possible target language phrase-map for phrases in the given source text (referred to as prediction of alignment) and (b) combining the phrase-maps together to synthesize the target language text with the objective to maximize fluency (referred to as

Method	BLEU-4	METEOR	TER	ROUGE-L
PB-SMT	3.50	10.30	102.95	0.179
SYNTAX-SMT	3.40	10.06	102.03	0.180
SEQ2SEQ (50)	1.63	0.07	89.38	0.160
SEQ2SEQ (128)	1.84	0.07	89.35	0.163
SEQ2SEQ (256)	1.98	0.07	89.23	0.166

Table 2: Evaluation results for SMT and SEQ2SEQ methods. For each SEQ2SEQ method, embedding dimensionality is mentioned in brackets.

decoding). In our setting incoherent captions are treated as *source* and the corresponding generated stories in the used data-sets are used as *target*. In PB-SMT, for learning phrase alignment, we use multi-threaded GIZA++ [8] and decoding is done using the MOSES decoder [13]. The development-set is used to tune the PB-SMT system using MERT mechanism [22].

3.2 Syntax based Statistical Machine Translation SYNTAX-SMT

Syntax-based translation (SYNTAX-SMT) is based on the idea of translating syntactic units, rather than single words or strings of words (as in phrase-based MT), i.e. (partial) parse trees of sentences/utterances [36]. For the task of story-generation, we use similar training configuration as PB-SMT to learn syntactic unit correspondences; for decoding the *chart-decoder* of moses [13] is used for syntax-based (tree-to-tree based) target-sentence generation.

3.3 Sequence-to-sequence Recurrent Neural Network

As an initial approach for story generation from independent descriptions, we use sequence-to-sequence recurrent neural network (SEQ2SEQ) architecture [32] popularly used for machine translation [35], question answering, summarization and other text generation tasks. We use a bidirectional encoder that processes the set of independent descriptions separated by a delimiter and SEQ_END marking the end of set. This encoding of descriptions is given to an attention enabled decoder that generates the story sequence. Refer to Figure 2 for an overview.

4 EXPERIMENT SETUP

We now describe the experimental setup below, sharing details on the dataset and configuration details of the story generation models.

4.1 Dataset

In the absence of availability of the any dataset for specific task, we use ‘Text Annotations’ from recently introduced Visual Storytelling Dataset [12] (VIST) for our experiments. Barring the images, we use both ‘Descriptions of Images-in-Isolation (DII)’ and ‘Stories of Images-in-Sequence (SIS)’ for our experiments. In all, there are 41300 image sequences aligned with caption and story pairs. The text data is crowd sourced using Amazon Mechanical Turk – participants were shown image sequences from the dataset to generate text descriptions for captions and stories. We believe these pairs act as a good (input-output) proxy for our setup as both of them describe a ‘scenes’ represented by image sequences.

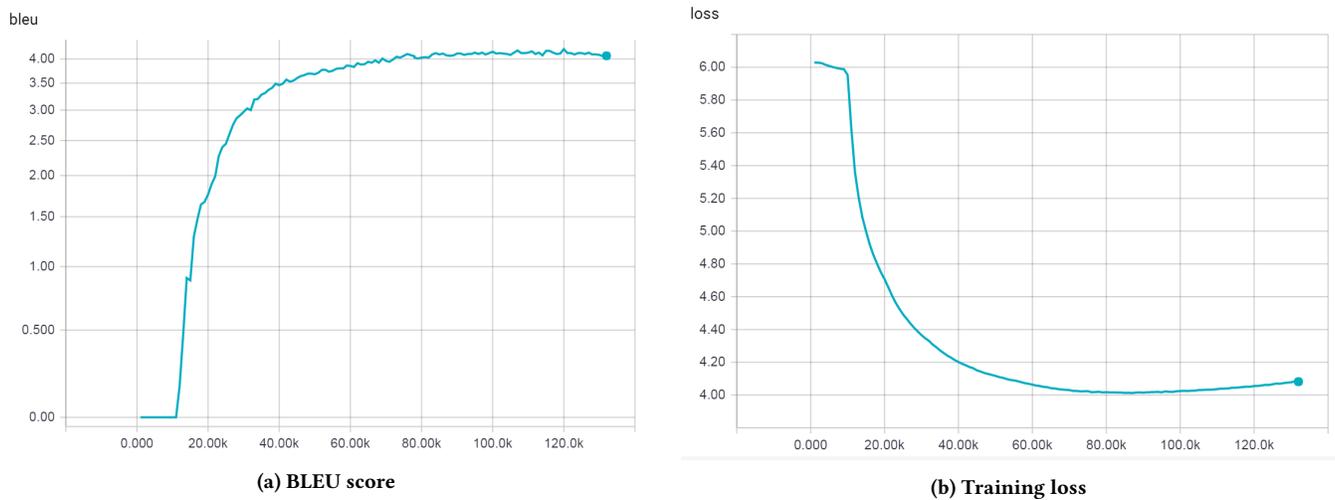


Figure 3: Graphs depicting changes in BLEU score and Loss over the iterations during training.

Independent captions describe each image as a separate event and stories are a coherent version of text describing the image sequence as a whole. Few sample input sets of descriptions and output stories are listed as a part of Table 3.

VIST dataset is split into training, validation and test sets and Table 1 summarizes various statistics of these sets like average number of words and sentences per caption and story. Usually the stories are longer and on an average have one sentence more compared to set of captions. More than 50% of words in a story are different from that in the captions. In fact, on an average, 41% of the story are non-stop words unseen in the input, thus making the task of generating the precise story tougher given limited set of training inputs.

4.2 Model Details

For SMT related experiments, we use the MOSES toolkit [13]. For phrase based translation model learning, *grow-diag-final-and* heuristic is opted and to tackle lexicalized reordering, we use the *msd-bidirectional-fe* model. GIZA++ is configured to apply the principles of IBM model 4 and 5 [6] and the HMM alignment model [23] for alignment learning. The assumption behind selecting the higher order IBM models is that we expect the possibility of addition of foreign words as well as dropping of words when a given text is transformed into a story, and higher order IBM models are good at tackling such nuances.

We tuned the trained SMT models using Minimum Error Rate Training (MERT) with default parameters (100 best list, max 25 iterations). We trained a 5-gram language model using the target side (stories) of the training data. For this, we use Kneser-Ney smoothing algorithm implemented in KenLM toolkit [11]. Batch training of multiple SMT systems was done using the Moses Job Scripts¹ experiment management system. We use similar training and tuning configuration for setting up the syntax based SMT system. The default rule-learning mechanism of MOSES is opted for

training by setting the hierarchical flag in the mooses training script. Decoding during tuning and testing is carried out using the mooses chart decoder.

For sequence to sequence learning, we implement Desc2Story framework used in our experiments. Fig 2 provides an overview of the model, our framework is based on sequence to sequence encoder decoder architecture including attention mechanism at the decoder [2]. Encoder is a single layer bidirectional Gated Recurrent Unit (GRU) [7] including dropout. Decoder is a 2 layer GRU using attention mechanism at each time step. We have used Adam optimizer to optimize the cross entropy loss function. Both encoder and decoder has a dropout [30] with probability $p = 0.8$. We tried 50, 128 and 256 dimensional RNN at encoder and decoder. During testing for SEQ2SEQ model, we use beam search with beam-width set to 5. Batch size of 32 was used while training.

4.3 Evaluation Metrics

For evaluation, we rely on popular Machine Translation evaluation metrics such as BLEU-4 [17], METEOR [4], and TER [29] and ROUGE-L [16], frequently used for evaluating summarization output. While these evaluation metrics may not adequately capture all necessary aspects of story generations, in the absence of alternative evaluation methods specific to story generation, these metrics could provide first a level insight into the applicability of popular NLG paradigms towards a difficult task like story generation.

5 RESULTS AND DISCUSSION

In this work, we attempted to try a few off-the-shelf techniques to solve a complex task as generating a story from simple textual descriptions in a sequence, without the taking another input or modality. The evaluation results are reported in Table 2, clearly suggesting that the overall scores are not very high. This indicates the limitations of basic off-the shelf methods for a highly creative task such as story generation even though they work well for simpler machine translation, paraphrase generation or document summarization tasks. This is a very complex task as stories

¹https://bitbucket.org/anoopk/moses_job_scripts

Input Descriptions	a student and an older man sit down to talk about technology. two young women speaking with an older man wearing a suit. man with gray hair and eye-glasses talking to woman in cap sitting at a computer. a young man and an older man are sitting in front of a monitor while referencing a book. group of students standing for picture in classroom.	group of older people gathered for bingo night in hall. a girl sitting at a table eating pancakes with fruit. a group of historical reenactors torment tourists who had hoped to have a fun vacation . a line of people walking through the woods.	the leaf has begun to change it 's color . close up of a bush with small pink flower clusters. purple and white flowers growing outdoors together .
PB-SMT	the family is so proud of him . the man is talking to her friends . this one has improved flash . i gave the speech . he is welcome to the graduation . we had a great time and we love him for a book . at the end of the day , and we had a great time . one of the class .	everyone arrived for the bingo hall . at the end of the night . we had a great dinner . more fruits . all of the history here . our reenactors torment was lovely . we took a walk through the forest .	in the woods . this begun to it . it is a clusters . the flowers were beautiful flowers and plants .
SYNTAX-SMT	the family is so proud of him . he has improved flash . to get to the fun . he and his . to my eyes . i turned on and my dad . we played in the world . the students love to read a book . and we had a great time .	the family gathered at the end of the night . there was a bingo hall . at t he end of the night . we ate breakfast at the end of the ruins . the fruits of our torment reenactors our lovely . we took a walk through the woods .	this begun in the woods . it was to support . there was a nice . the flowers and plants . the flowers were in clusters photos .
SEQ2SEQ (256)	a group of friends met for a meeting . they talked about the plans . they talked about the company 's company . they had a great time at the meeting . everyone had a great time .	a group of friends went on a trip to location . they went to a historical site . they saw many interesting buildings . they also had a map to eat . they had a great time .	we decided to take a trip to the botanical gardens . there was a beautiful view of the water . we also saw some lily pads . there was a lot of flowers on the ground . we also saw some lily pads .
Ground Truth	the gentleman sat with the boy to discuss the diagrams . he then asked the young ladies if they needed help with anything . he spoke to the man about his interest in technology . he then met with this fellow , to discuss his future plans in engineering . the students took a picture to remember the special day .	we went to the syrup festival. the kids got to eat pancakes. there was also a demonstration on how syrup is made. we got a tour of the woods. and it ended in the production lab.	i have been outside doing micro photography for a class project . love how it is possible to get a good blur in pictures like this . these flowers wer e so beautiful . one of my class subjects is nature . nothing beats getting out and taking pictures of sites like this . most people never experience this . last but not least a single red flower . this day of shooting turned out very good .

Table 3: Example Case studies on VIST dataset

can be generated with a greater diversity in content compared to the above tasks. While the narrative generation is itself a grand challenge, metrics to evaluate such a task is also inadequate. It is well-known that BLEU and other metrics have their shortcomings even in simpler natural language generation settings as they are mostly based on n-gram match. This is another reason why scores

of our baseline models are low even though we found coherent enough outputs as shown in Table 3.

From Table 2, it is clear that SMT systems produce higher BLEU/METEOR/ROUGE values with PB-SMT being the best. SEQ2SEQ scores, on the other hand, suffer as SEQ2SEQ models are embedding-based approaches which can lead to less scores on n-gram based

match evaluation even though they are able to produce better readable and coherent stories. This is evident from the anecdotal examples in Table 3. This is supported also by the fact that TER, which is not based on n-gram match, is lower in case of SEQ2SEQ indicating better translation quality by lowering the translation edit rate. Figures 3a and 3b show how the training of SEQ2SEQ model saturates with the iterations (each iteration is a new batch), suggesting that the model has indeed converged.

Referring to the examples as in Table 3, we find that all methods are doing well to generate a story which is readable and coherent. This shows the models are capable enough to capture the grammatical aspects of a story based on the training data used. However, all models fall short of generating a story which is semantically related to the input descriptions, therefore, highlighting the need for more sophisticated models for creative story generation. We believe, our current models are too simplistic to realize beyond the co-occurrence statistics in the output, thus, unable to figure out the semantic relatedness between the input and output. This could possibly be tackled to some extent by *pretraining* the encoder and decoder in the SEQ2SEQ model using widely available data in the form of essays and novels. Another promising alternative is modeling the same problem with hierarchical recurrent neural networks. However, in the long run, this kind of task also calls for novel generation of story at test time not just relying on mapping between input and output words/phrases learnt at training time. Thus, to achieve the distant goal of novel story generation which is somewhat semantically connected with the input, we have to look beyond just translation based techniques, which essentially maps input words to output words in sequence.

6 CONCLUSION AND FUTURE WORK

We introduced the task of automatically generating stories from independent one-liners. The scope of this task can eventually be expanded to generating stories from fewer input parameters like theme, actors, etc. This task is challenging not only because of insufficiency of information compared to the task attempted in VIST dataset but also because of the fact that the possible output space of stories is large. This is aggravated by the unavailability of good metrics to evaluate the approaches tried out. Here we explored off-the-shelf sequence generation methods like statistical machine translation and sequence to sequence neural networks as preliminary approaches towards solving this problem.

A future direction is to design trainable metrics for evaluating stories holistically to include aspects on creativity, coherency, novelty and other parameters compared to current score computation which is based on exact match. Creating an appropriate dataset for textual story generation is another important direction that we would like pursue.

REFERENCES

- [1] Ruth Aylett. 1999. Narrative in virtual environments-towards emergent narrative. In *AAAI*.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [3] Paul Bailey. 1999. Searching for storiness: Story-generation from a reader's perspective. In *Working notes of the Narrative Intelligence Symposium*.
- [4] Satyanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*.
- [5] J Barrie. 2014. Computers Are Writing Novels: Read A Few Samples Here. (2014).
- [6] Peter F Brown, Vincent J Della Pietra, Stephen A Della Pietra, and Robert L Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational linguistics* (1993).
- [7] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078* (2014).
- [8] Qin Gao and Stephan Vogel. 2008. Parallel implementations of word alignment tool. In *Software Engineering, Testing, and Quality Assurance for Natural Language Processing*.
- [9] Pablo Gervás, Belén Díaz-Agudo, Federico Peinado, and Raquel Hervás. 2005. Story plot generation based on CBR. *Knowledge-Based Systems* (2005).
- [10] Sergio Guadarrama, Niveda Krishnamoorthy, Girish Malkarnenkar, Subhashini Venugopalan, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2013. Youtube2text: Recognizing and describing arbitrary activities using semantic hierarchies and zero-shot recognition. In *ICCV*.
- [11] Kenneth Heafield. 2011. KenLM: Faster and smaller language model queries. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*.
- [12] Ting-Hao K. Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Jacob Devlin, Aishwarya Agrawal, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual Storytelling. In *North American Chapter of the Association for Computational Linguistics*.
- [13] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*.
- [14] Ray Kurzweil. 2005. *The singularity is near: When humans transcend biology*.
- [15] Boyang Li, Stephen Lee-Urban, George Johnston, and Mark Riedl. 2013. Story Generation with Crowdsourced Plot Graphs. In *AAAI*.
- [16] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out: Proceedings of the ACL-04 workshop*.
- [17] Chin-Yew Lin and Franz Josef Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Association for Computational Linguistics*.
- [18] Yu Liu, Jianlong Fu, Tao Mei, and Chang Wen Chen. 2016. Storytelling of Photo Stream with Bidirectional Multi-thread Recurrent Neural Network. *CoRR abs/1606.00625* (2016). <http://arxiv.org/abs/1606.00625>
- [19] Adam Lopez. 2008. Statistical machine translation. *ACM Computing Surveys (CSUR)* (2008).
- [20] Neil McIntyre and Mirella Lapata. 2009. Learning to tell tales: A data-driven approach to story generation. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*. Association for Computational Linguistics, 217–225.
- [21] James Richard Meehan. 1976. *The metanovel: writing stories by computer*. Technical Report. DTIC Document.
- [22] Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics-Volume 1*.
- [23] Franz Josef Och and Hermann Ney. 2000. Improved statistical alignment models. In *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*.
- [24] Vicente Ordonez, Girish Kulkarni, and Tamara L Berg. 2011. Im2text: Describing images using 1 million captioned photographs. In *NIPS*.
- [25] Jia-Yu Pan, Hyung-Jeong Yang, Pinar Duygulu, and Christos Faloutsos. 2004. Automatic image captioning. In *ICME*.
- [26] Lyn Pemberton. 1989. A modular approach to story generation. In *Proceedings of the fourth conference on European chapter of the Association for Computational Linguistics*.
- [27] Rafael Pérez Y Pérez and Mike Sharples. 2001. MEXICA: A computer model of a cognitive account of creative writing. *Journal of Experimental & Theoretical Artificial Intelligence* (2001).
- [28] Devendra Kumar Sahu and Mohak Sukhwani. 2015. Sequence to sequence learning for optical character recognition. *arXiv preprint arXiv:1511.04176* (2015).
- [29] Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Association for machine translation in the Americas*.
- [30] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* (2014).
- [31] Mohak Sukhwani and CV Jawahar. 2015. TennisVid2Text: Fine-grained Descriptions for Domain Specific Videos. In *BMVC*.

- [32] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. *CoRR* (2014).
- [33] Reid Swanson and Andrew Gordon. 2008. Say Anything: A massively collaborative open domain story writing companion. *Interactive Storytelling* (2008), 32–40.
- [34] Mariët Theune, Sander Faas, DKJ Heylen, and Anton Nijholt. 2003. The virtual storyteller: Story creation by intelligent agents. (2003).
- [35] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, ffukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *CoRR* (2016).
- [36] Kenji Yamada and Kevin Knight. 2001. A syntax-based statistical translation model. In *Proceedings of the 39th Annual Meeting on Association for Computational Linguistics*.
- [37] Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher Pal, Hugo Larochelle, and Aaron Courville. 2015. Describing Videos by Exploiting Temporal Structure. In *ICCV*.
- [38] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. 2016. Image captioning with semantic attention. In *CVPR*.
- [39] Richard Zens, Franz Josef Och, and Hermann Ney. 2002. Phrase-based statistical machine translation. In *Annual Conference on Artificial Intelligence*.