

# Storyboarding of Recipes: Grounded Contextual Generation

Khyathi Raghavi Chandu    Eric Nyberg    Alan W Black

Language Technologies Institute, Carnegie Mellon University

{kchandu, eh, awb}@cs.cmu.edu

## Abstract

Information need of humans is essentially multimodal in nature, enabling maximum exploitation of situated context. We introduce a dataset for sequential procedural (*how-to*) text generation from images in cooking domain. The dataset consists of 16,441 cooking recipes with 160,479 photos associated with different steps. We setup a baseline motivated by the best performing model in terms of human evaluation for the Visual Story Telling (ViST) task. In addition, we introduce two models to incorporate high level structure learnt by a Finite State Machine (FSM) in neural sequential generation process by: (1) Scaffolding Structure in Decoder (SSiD) (2) Scaffolding Structure in Loss (SSiL). Our best performing model (SSiL) achieves a METEOR score of 0.31, which is an improvement of 0.6 over the baseline model. We also conducted human evaluation of the generated grounded recipes, which reveal that 61% found that our proposed (SSiL) model is better than the baseline model in terms of overall recipes. We also discuss analysis of the output highlighting key important NLP issues for prospective directions.

## 1 Introduction

Interpretation is heavily conditioned on context. Real world interactions provide this context in multiple modalities. In this paper, the context is derived from vision and language. The description of a picture changes drastically when seen in a sequential narrative context. Formally, this task is defined as: given a sequence of images  $\mathbb{I} = \{I_1, I_2, \dots, I_n\}$  and pairwise associated textual descriptions,  $\mathbb{T} = \{T_1, T_2, \dots, T_n\}$ ; for a new sequence  $\mathbb{I}'$ , our task is to generate the corresponding  $\mathbb{T}'$ . Figure 1 depicts an example for making *vegetable lasagna*, where the input is the first row and the output is the second row. We call this a ‘*storyboard*’, since it unravels the most important steps of a procedure associated with corresponding natural language text. The sequential context differ-

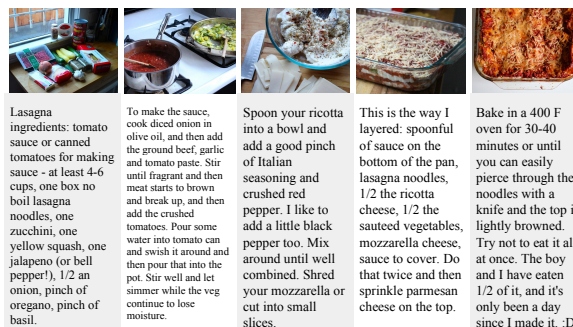


Figure 1: Storyboard for the recipe of vegetable lasagna

entiates this task from image captioning in isolation. The dataset is similar to that of ViST (Huang et al., 2016) with an apparent difference between stories and instructional in-domain text which is the clear transition in phases of the narrative. This task supplements the task of ViST with richer context of goal oriented procedure (*how-to*). Numerous online blogs and videos depict various categories of *how-to* guides for games, do-it-yourself (DIY) crafts, technology etc. This task lays initial foundations for full fledged storyboarding of a given video, by selecting the right junctions/clips to ground significant events and generate sequential textual descriptions. We are going to focus on the domain of cooking recipes in the rest of this paper. In this paper, we discuss our approach in generating more structural/coherent cooking recipes by explicitly modeling the state transitions between different stages of cooking (phases). We introduce a framework to apply traditional FSMs to incorporate more structure in neural generation.

The two main contributions of this paper are: (1) A dataset of 16k recipes targeted for sequential multimodal procedural text generation, (2) Two models (SSiD: Structural Scaffolding in Decoder, and SSiL: Structural Scaffolding in Loss) for incorporating high level structure learnt by an FSM into a neural text generation model to improve structure/coherence.

## 2 Related Work

**Why domain constraint?** Martin et al. (2017) and Khalifa et al. (2017) demonstrated that the predictive ability of a seq2seq model improves as the language corpus is reduced to a specialized domain with specific actions. Our choice of restricting domain to recipes is inspired from this, where the set of events are specialized (such as ‘cut’, ‘mix’, ‘add’) although we are not using event representations explicitly. These specialized set of events are correlated to phases of procedural text as described in the following sections.

**Planning while writing content:** A major challenge faced by neural text generation (Lu et al., 2018) while generating long sequences is the inability to maintain structure, contravening the coherence of the overall generated text. This aspect was also observed in various tasks like summarization (Liu et al., 2018), story generation (Fan et al., 2019). Pre-selecting content and planning to generate accordingly was explored by Puduppully et al. (2018) and Lukin et al. (2015) in contrast to generate as you proceed paradigm. Fan et al. (2018) adapt a hierarchical approach to generate a premise and then stories to improve coherence and fluency. Yao et al. (2018) experimented with static and dynamic schema to realize the entire storyline before generating. However, in this work we propose a hierarchical multi task approach to perform structure aware generation.

**Comprehending Food:** Recent times have seen large scale datasets in food, such as Recipe1M (Marin et al., 2018), Food-101 (Bossard et al., 2014). Food recognition (Arora et al., 2019) addresses understanding food from a vision perspective. Salvador et al. (2018) worked on generating cooking instructions by inferring ingredients from an image. Zhou et al. (2018) proposed a method to generate procedure segments for YouCook2 data. In NLP domain, this is studied as generating procedural text by including ingredients as checklists (Kiddon et al., 2016) or treating the recipe as a flow graph (Mori et al., 2014). Our work is at the intersection of two modalities (language and vision) by generating procedural text for recipes from a sequence of images. (Bosselut et al., 2017) worked on reasoning non-mentioned causal effects thereby improving the understanding and generation of procedural text for cooking recipes. This is done by dynamically tracking entities by modeling actions using state transformers.

**Visual Story Telling:** Research at the intersection of language and vision is accelerating with tasks like image captioning (Hossain et al., 2019),

visual question answering (Wu et al., 2017), visual dialog (Das et al., 2017; Mostafazadeh et al., 2017; De Vries et al., 2017; de Vries et al., 2018). ViST (Huang et al., 2016) is a sequential vision to language task demonstrating differences between descriptions in isolation and stories in sequences. Similarly, Gella et al. (2018) created VideoStory dataset from videos on social media with the task of generating a multi-sentence story captions for them. Smilevski et al. (2018) proposed a late fusion based model for ViST challenge. Kim et al. (2018) attained the highest scores on human readability in this task by attending to both global and local contexts. We use this as our baseline model and propose two techniques on top of this baseline to impose structure needed for procedural text.

## 3 Data Description

We identified two *how-to* blogs from: *instructables.com* and *snapguide.com*, comprising step-wise instructions (images and text) of various *how-to* activities like games, crafts etc,. We gathered 16,441 samples with 160,479 photos for food, dessert and recipe topics. We used 80% for training, 10% for validation and 10% for testing our models. In some cases, there are multiple images for the same step and we randomly select an image from the set of images. We indicate that there is a potential space for research here, in selecting most distinguishing/representative/meaningful image. Details of the datasets are presented in Table 1. The data and visualization of distribution of topics is here<sup>1</sup>. A trivial extension could be done on other domains like gardening, origami crafts, fixing guitar strings etc, which is left for future work.

## 4 Model Description

We first describe a baseline model for the task of storyboarding cooking recipes in this section. We then propose two models with incremental improvements to incorporate the structure of procedural text in the generated recipes : SSiD (Scaffolding Structure in Decoder) and SSiL (Scaffolding Structure in Loss). The architecture of *scaffolding* structure is presented in Figure 2, of which different aspects are described in the following subsections.

### 4.1 Baseline Model (Glocal):

The baseline model is inspired from the best performing system in ViST challenge with respect to

<sup>1</sup><https://storyboarding.github.io/story-boarding/>

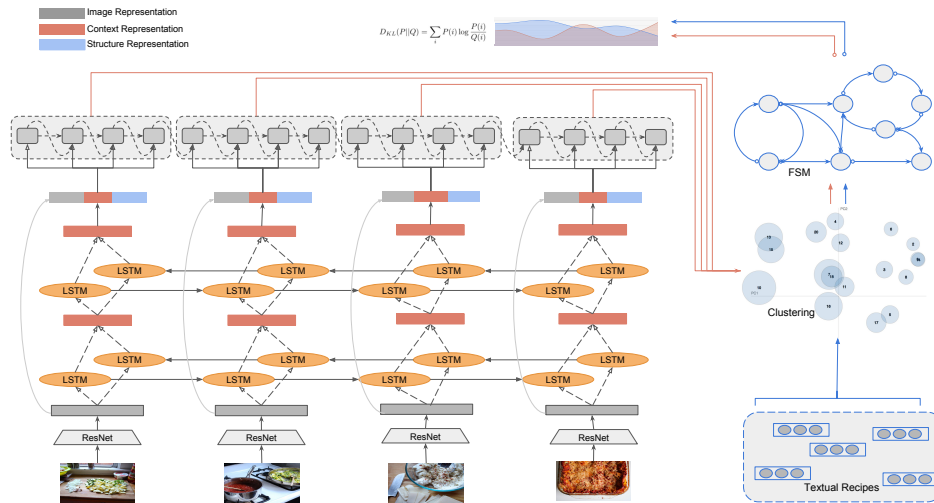


Figure 2: Architecture for incorporating high level structure in neural recipe generation

Data Sources	# Recipes	# Avg Steps
<i>instructables</i>	9,101	7.14
<i>snapguide</i>	7,340	13.01

Table 1: Details of dataset for *storyboarding* recipes

human evaluation (Kim et al., 2018). The images are first resized into 224 X 224. Image features for each step are extracted from the penultimate layer of pre-trained ResNet-152 (He et al., 2016). These features are then passed through an affinity layer to obtain an image feature of dimension 1024. To maintain the context of the entire recipe (global context), the sequence of these image features are passed through a two layered Bi-LSTM with a hidden size of 1024. To maintain specificity of the current image (local context), the image features for the current step are concatenated using a skip connection to the output of the Bi-LSTM to obtain global representation. Dropout of 0.5 is applied systematically at the affinity layer to obtain the image feature representation and after the Bi-LSTM layer. Batch normalization is applied with a momentum 0.01. This completes the encoder part of the sequence to sequence architecture. These global vectors are used for decoding each step. These features are passed through a fully connected layer to obtain a representation of 1024 dimension followed by a non-linear transformation using ReLU. These features are then passed through a decoder LSTM for each step in the recipe which are trained by teacher forcing. The overall coherence in generation is addressed by feeding the decoder state of the previous step to the next one. This is a seq2seq model translating one modality into another. The

model is optimized using Adam with a learning rate of 0.001 and weight decay of 1e-5.

The model described above does not explicitly cater to the structure of the narration of recipes in the generation process. However, we know that procedural text has a high level structure that carries a skeleton of the narrative. In the subsequent subsections, we present two models that impose this high level narrative structure as a scaffold. While this scaffold lies external to the baseline model, it functions on imposing the structure in decoder (SSiD) and in the loss term (SSiL).

#### 4.2 Scaffolding Structure in Decoder (SSiD):

There is a high level latent structure involved in a cooking recipe that adheres to transitions between steps, that we define as *phases*. Note that the steps and phases are different here. To be specific, according to our definition, one or more steps map to a phase (this work does not deal with multiple phases being a part of a single step). Phases may be ‘listing ingredients’, ‘baking’, ‘garnishing’ etc.,. The key idea of the SSiD model is to incorporate the sequence of phases in the decoder to impose structure during text generation

There are two sources of supervision to drive the model: (1) multimodal dataset  $\mathbb{M} = \{\mathbb{I}, \mathbb{T}\}$  from Section 3, (2) unimodal textual recipes<sup>2</sup>  $\mathbb{U}$  to learn phase sequences. Finer phases are learnt using clustering followed by an FSM.

**Clustering:** K-Means clustering is performed on the sentence embeddings with compositional n-gram features (Pagliardini et al., 2018) on each step of the recipe in  $\mathbb{U}$ . Aligning with our intu-

<sup>2</sup>[www.fts.com/recipes.htm](http://www.fts.com/recipes.htm)

ition, when  $k$  is 3, it is observed that these clusters roughly indicate categories of *desserts*, *drinks* and *main course foods* (*pizza*, *quesadilla* etc.). However, we need to find out finer categories of the phases corresponding to the phases in the recipes. We use  $k$ -means clustering to obtain the categories of these phases. We experimented with different number of phases  $\mathbf{P}$  as shown in Table 2. For example, let an example recipe comprise of 4 steps i.e, a sequence of 4 images. At this point, each recipe can be represented as a hard sequence of phases  $\mathbf{r} = \langle p_1, p_2, p_3, p_4 \rangle$ .

**FSM:** The phases learnt through clustering are not ground truth phases. We explore the usage of an FSM to individually model hard and a softer representation of the phase sequences by leveraging the states in an FSM. We first describe how the hard representation is modeled. The algorithm was originally developed for building language models for limited token sets in grapheme to phoneme prediction. The iterative algorithm starts with an ergodic state for all phase types and uses entropy to find the best state split that would maximize the prediction. As opposed to phase sequences, each recipe is now represented as a state sequence (decoded from FSM) i.e,  $\mathbf{r} = \langle s_1, s_2, s_3, s_4 \rangle$  (hard states). This is a hard representation of the sequence of states.

We next describe how a soft representation of these states is modeled. Since the phases are learnt in an unsupervised fashion and the ground truth of the phases is not available, we explored a softer representation of the states. We hypothesize that a soft representation of the states might smooth the irregularities of phases learnt. From the output of the FSM, we obtain the state transition probabilities from each state to every other state. Each state  $s_i$  can be represented as  $\langle q_{ij} \forall j \in \mathbf{S} \rangle$  (soft states), where  $q_{ij}$  is the state transition probability from  $s_i$  to  $s_j$  and  $\mathbf{S}$  is the total number of states. This is the soft representation of state sequences.

The structure in the recipe is learnt as a sequence of phases and/or states (hard or soft). This is the structural *scaffold* that we would like to incorporate in the baseline model. In SSiD model, for each step in the recipe, we identify which phase it is in using the clustering model and use the phase sequence to decode state transitions from the FSM. The state sequences are concatenated to the decoder in the hard version and the state transition probabilities are concatenated in the decoder in the soft version at every time step.

At this point, we have 2 dimensions, one is the complexity of the phases ( $\mathbf{P}$ ) and the other is the

FST Complexity	1	20	40	60	80	100	120
20 Phases	11.27	11.60	12.31	13.71	12.32	12.51	12.36
40 Phases	12.03	12.44	11.48	12.58	12.50	<b>13.91</b>	11.82
60 Phases	11.13	11.18	12.74	12.26	12.47	12.98	11.47

Table 2: BLEU Scores for different number of phases ( $\mathbf{P}$ ) and states( $\mathbf{S}$ )

complexity of the states in FSM ( $\mathbf{S}$ ). Comprehensive results of searching this space is presented in Table 2. We plan to explore the usage of hidden markov model in place of FSM in future.

### 4.3 Scaffolding Structure in Loss (SSiL):

In addition to imposing structure via SSiD, we explored measuring the deviation of the structure learnt through phase/state sequences from the original structure. This leads to our next model where the deviation of the structure in the generated output from that of the original structure is reflected in the loss. The decoded steps are passed through the clustering model to get phase sequences and then state transition probabilities are decoded from FSM for the generated output. We go a step further to investigate the divergence between the phases of generated and original steps. This can also be viewed as hierarchical multi-task learning (Sanh et al., 2018). The first task is to decode each step in the recipe (which uses a cross entropy criterion,  $\mathbf{L}_1$ ). The second task uses KL divergence between phase sequences of decoded and original steps to penalize the model (say,  $\mathbf{L}_2$ ). When there are  $\tau$  steps in a recipe, we obtain  $o(s_1^\tau)$  and  $g(s_1^\tau)$  as the distributions of phases comprising of soft states for the original and generated recipes respectively. We measure the KL divergence( $D_{KL}$ ) between these distributions:

$$D_{KL}(o(s_1^\tau)||g(s_1^\tau)) = \sum_{i=1}^{\tau} \sum_{j=1}^S o(s_i[j]) \log \frac{o(s_i[j])}{g(s_i[j])}$$

Each task optimizes different functions and we minimize the combination of the two losses.

$$\sum_{I,T \in \mathbb{I}, \mathbb{T}} \mathbf{L}_1(I, T) + \alpha \sum_{U \in \mathbb{U}} \mathbf{L}_2(U)$$

This combined loss is used to penalize the model. Here,  $\alpha$  is obtained from KL annealing (Bowman et al., 2015) function that gradually increases the weight of KL term from 0 to 1 during train time.

## 5 Results and Discussion

The two dimensions explored in FSM are  $\mathbf{P}$  and  $\mathbf{S}$  and exhaustive results are presented in Table





Models					Phenomena
<i>Glocal Model</i>	This is a simple recipe for making a delicious chicken salad.	You will need: a butter knife a plate of bread flour a little bit of salt a dash of pepper flakes a couple of tablespoons of olive oil a pinch of sugar.	Add butter evenly on the pan.	Put the chicken on the grill and set aside.	<ul style="list-style-type: none"> <li>- Ingredients phase wrongly identified.</li> <li>- Wrong ingredients.</li> <li>- Improper conclusion.</li> </ul>
<i>SSiD Model</i>	This is a simple recipe for making a delicious and easy dish.	Ingredients: 4 pounds chicken 2 tsp salt, ¼ tsp sugar, marinara sauce, mozzarella cheese ( I used provolone ). Tools: a knife, an oven for the chicken, tongs. Mix all ingredients in a bag.	Add butter evenly on the pan.	Serve the baked chicken wings and enjoy the evening!	<ul style="list-style-type: none"> <li>- Learnt majority structure (step 1)</li> <li>+ Got 'tongs' right because of separate tools mention.</li> <li>- The action of baking is not explicitly mentioned (before 'baked' wings).</li> </ul>
<i>SSiL Model</i>	You will need: 5 pounds of chicken wings, ½ cup all purpose flour, ½ tsp salt, 2 tsp of paprika, melted butter, silicon mat, baking pan.	Preheat oven to 450 F. Mix dry ingredients in the dry ziplock bag.	Place a mat on the baking pan and spread butter evenly on it.	Spread the chicken pieces on butter on the baking pan. Bake until crispy for 30 minutes. Serve and enjoy!	<ul style="list-style-type: none"> <li>+ Global context of baking maintained in preheating.</li> <li>+ Non-repetitive ingredients phase.</li> <li>+ Referring expressions (baking pan -&gt; it).</li> <li>- Not mentioned tools (tongs).</li> </ul>

Figure 3: Comparison of generated storyboards for *Easy Oven Baked Crispy Chicken Wings*

Models	BLEU	METEOR	ROUGE-L
Glocal	10.74	0.25	0.31
SSiD (hard phases)	11.49	0.24	0.31
SSiD (hard states)	11.93	0.25	0.31
SSiD (soft phases)	13.91	0.29	0.32
<b>SSiL (soft phases)</b>	<b>16.38</b>	<b>0.31</b>	<b>0.34</b>

Table 3: Evaluation of storyboarding recipes

2. The BLEU score (Papineni et al., 2002) is the highest when  $P$  is 40 and  $S$  is 100. Fixing these values, we compare the models proposed in Table 3. The models with hard phases and hard states are not as stable as the one with soft phases since backprop affects the impact of the scaffolded phases. Upon manual inspection, a key observation is that for SSiD model, most of the recipes followed a similar structure. It seemed to be conditioned on a global structure learnt from all recipes rather than the current input. However, SSiL model seems to generate recipe that is conditioned on the structure of that particular example.

**Human Evaluation:** We have also performed human evaluation by conducting user preference study to compare the baseline with our best performing SSiL model. We randomly sampled generated outputs of 20 recipes and asked 10 users to answer two preferences: (1) overall recipe based on images, (2) structurally coherent recipe. Our SSiL model was preferred 61% and 72.5% for overall and structural preferences respectively. This shows that while there is a viable space to improve structure, generating an edible recipe needs to be explored to improve the overall preference.

### 5.1 Qualitative Analysis:

Figure 3 presents the generated text from the three models with an analysis described below.

**Coherence of Referring Expressions:** Introducing referring expressions is a key aspect of co-

herence (Dale, 2006, 1992), as seen in the case of ‘baking pan’ being referred as ‘it’ in the SSiL model.

**Context Maintenance:** Maintaining overall context explicitly affects generating each step. This is seen in SSiL model where ‘preheating’ in the second step is learnt from *baking* step that appears later although the image does not show an oven.

**Schema for Procedural Text:** Explicit modeling of structure has enabled SSiD and SSiL models to conclude the recipe by generating words like ‘serve’ and ‘enjoy’. Lacking this structure, glocal model talks about ‘setting aside’ at the end.

**Precision of Entities and Actions:** SSiD model introduces ‘sugar’ in ingredients after generating ‘salt’. A brief manual examination revealed that this co-occurrence is a common phenomenon. SSiL model misses ‘tongs’ in the first step.

## 6 Conclusions

Our main focus in this paper is instilling structure learnt from FSMs in neural models for sequential procedural text generation with multimodal data. We gather a dataset of 16k recipes where each step has text and associated images. We setup a baseline inspired from the best performing model in ViST. We propose two ways of imposing structure from phases and states of a recipe derived from FSM. The first model imposes structure on the decoder and the second model imposes structure on the loss function by modeling it as a hierarchical multi-task learning problem. We show that our proposed approach improves upon the baseline and achieves a METEOR score of 0.31. We plan to explore explicit evaluation of the latent structure learnt. We plan on exploring backpropable variants as a scaffold for structure and also extend the techniques to other *how-to* domains in future.

## References

- Sandhya Arora, Gauri Chaware, Devangi Chinchankar, Eesha Dixit, and Shevi Jain. 2019. Survey of different approaches used for food recognition. In *Information and Communication Technology for Competitive Strategies*, pages 551–560. Springer.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 2014. Food-101—mining discriminative components with random forests. In *European Conference on Computer Vision*, pages 446–461. Springer.
- Antoine Bosselut, Omer Levy, Ari Holtzman, Corin Ennis, Dieter Fox, and Yejin Choi. 2017. Simulating action dynamics with neural process networks. *arXiv preprint arXiv:1711.05313*.
- Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. 2015. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*.
- Robert Dale. 1992. *Generating referring expressions: Constructing descriptions in a domain of objects and processes*. The MIT Press.
- Robert Dale. 2006. Generating referring expressions.
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. 2017. Visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 2.
- Harm De Vries, Florian Strub, Sarath Chandar, Olivier Pietquin, Hugo Larochelle, and Aaron C Courville. 2017. Guesswhat?! visual object discovery through multi-modal dialogue. In *CVPR*, volume 1, page 3.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. *arXiv preprint arXiv:1805.04833*.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. *arXiv preprint arXiv:1902.01109*.
- Spandana Gella, Mike Lewis, and Marcus Rohrbach. 2018. A dataset for telling the stories of social media videos. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 968–974.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- MD Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. 2019. A comprehensive survey of deep learning for image captioning. *ACM Computing Surveys (CSUR)*, 51(6):118.
- Ting-Hao Kenneth Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual storytelling. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1233–1239.
- Ahmed Khalifa, Gabriella AB Barros, and Julian Togelius. 2017. Deeptingle. *arXiv preprint arXiv:1705.03557*.
- Chloé Kiddon, Luke Zettlemoyer, and Yejin Choi. 2016. Globally coherent text generation with neural checklist models. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 329–339.
- Taehyeong Kim, Min-Oh Heo, Seonil Son, Kyoung-Wha Park, and Byoung-Tak Zhang. 2018. Glacnet: Glocal attention cascading networks for multi-image cued story generation. *arXiv preprint arXiv:1805.10973*.
- Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. *arXiv preprint arXiv:1801.10198*.
- Sidi Lu, Yaoming Zhu, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Neural text generation: past, present and beyond. *arXiv preprint arXiv:1803.07133*.
- Stephanie M Lukin, Lena I Reed, and Marilyn A Walker. 2015. Generating sentence planning variations for story telling. In *16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, page 188.
- Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2018. Recipe1m: A dataset for learning cross-modal embeddings for cooking recipes and food images. *arXiv preprint arXiv:1810.06553*.
- Lara J Martin, Prithviraj Ammanabrolu, Xinyu Wang, William Hancock, Shruti Singh, Brent Harrison, and Mark O Riedl. 2017. Event representations for automated story generation with deep neural nets. *arXiv preprint arXiv:1706.01331*.
- Shinsuke Mori, Hirokuni Maeta, Yoko Yamakata, and Tetsuro Sasada. 2014. Flow graph corpus from recipe texts. In *LREC*, pages 2370–2377.
- Nasrin Mostafazadeh, Chris Brockett, Bill Dolan, Michel Galley, Jianfeng Gao, Georgios P Spithourakis, and Lucy Vanderwende. 2017. Image-grounded conversations: Multimodal context for natural question and response generation. *arXiv preprint arXiv:1701.08251*.

- Matteo Pagliardini, Prakhar Gupta, and Martin Jaggi. 2018. Unsupervised learning of sentence embeddings using compositional n-gram features. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, volume 1, pages 528–540.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Ratish Puduppully, Li Dong, and Mirella Lapata. 2018. Data-to-text generation with content selection and planning. *arXiv preprint arXiv:1809.00582*.
- Amaia Salvador, Michal Drozdal, Xavier Giro-i Nieto, and Adriana Romero. 2018. Inverse cooking: Recipe generation from food images. *arXiv preprint arXiv:1812.06164*.
- Victor Sanh, Thomas Wolf, and Sebastian Ruder. 2018. A hierarchical multi-task approach for learning embeddings from semantic tasks. *arXiv preprint arXiv:1811.06031*.
- Marko Smilevski, Ilija Lalkovski, and Gjorgji Madzarov. 2018. Stories for images-in-sequence by using visual and narrative components. *arXiv preprint arXiv:1805.05622*.
- Harm de Vries, Kurt Shuster, Dhruv Batra, Devi Parikh, Jason Weston, and Douwe Kiela. 2018. Talk the walk: Navigating new york city through grounded dialogue. *arXiv preprint arXiv:1807.03367*.
- Qi Wu, Damien Teney, Peng Wang, Chunhua Shen, Anthony Dick, and Anton van den Hengel. 2017. Visual question answering: A survey of methods and datasets. *Computer Vision and Image Understanding*, 163:21–40.
- Lili Yao, Nanyun Peng, Weischedel Ralph, Kevin Knight, Dongyan Zhao, and Rui Yan. 2018. Plan-and-write: Towards better automatic storytelling. *arXiv preprint arXiv:1811.05701*.
- Luowei Zhou, Chenliang Xu, and Jason J Corso. 2018. Towards automatic learning of procedures from web instructional videos. In *Thirty-Second AAAI Conference on Artificial Intelligence*.