

# Latent Alignment of Procedural Concepts in Multimodal Recipes

Hossein Rajaby Faghihi, Roshanak Mirzaee, Sudarshan Paliwal and Parisa Kordjamshidi

Michigan State University

{rajabyfa, mirzaeem, paliwal, kordjams}@msu.edu

## Abstract

We propose a novel alignment mechanism to deal with procedural reasoning on a newly released multimodal QA dataset, named RecipeQA. Our model is solving the textual cloze task which is a reading comprehension on a recipe containing images and instructions. We exploit the power of attention networks, cross-modal representations, and a latent alignment space between instructions and candidate answers to solve the problem. We introduce constrained max-pooling which refines the max-pooling operation on the alignment matrix to impose disjoint constraints among the outputs of the model. Our evaluation result indicates a 19% improvement over the baselines.

## 1 Introduction

Procedural reasoning by following several steps to achieve a goal is an essential part of our daily tasks. However, this is challenging for machines due to the complexity of instructions and commonsense reasoning required for understanding the procedure (Dalvi et al., 2018; Yagcioglu et al., 2018; Bosselut et al., 2017).

In this paper, we tackle the task of procedural reasoning in a multimodal setting for understanding cooking recipes. The RecipeQA dataset (Yagcioglu et al., 2018) contains recipes from internet users. Thus, understanding the text is challenging due to the different language usage and informal nature of user-generated texts. The recipes are along with images provided by users which are taken in an unconstrained environment. This exposes a level of difficulty similar to real-world problems.

The tasks proposed with the dataset include textual cloze, visual cloze, visual ordering, and visual coherence. Here, we focus on textual cloze. An example of this task is shown in Figure 1. The input to the task is a set of multimodal instructions,

three textual items from the question and a placeholder to be filled by the answer. The answer has to be chosen from four options. The three question items and the correct answer make a sequence which correctly describes the steps of the recipe.

To design our model, we rely on the intuition that given question items, each answer describes exactly one step of the recipe. Hence, we design a model to make explicit alignments between the candidate answers and each step and use those alignment results, given the question information. This alignment space is latent due to not having any direct supervision based on provided annotations.

Using multimodal information and representations by making a joint space for comparison has been broadly investigated in the recent research (Hessel et al., 2019; Wu et al., 2019; Li et al., 2019; Su et al., 2020; Yu et al., 2019; Fan and Zhou, 2018; Tan and Bansal, 2019; Nam et al., 2017). Our work differs from those as we do not have direct supervision on multimodal alignments. Moreover, the task we are solving uses the sequential nature of visual and textual modality as a weak source of supervision to build a neural model to compare the textual representation of context and the answers for a given question representation.

Procedural reasoning has been investigated on different tasks (Amac et al., 2019; Park et al., 2017). While PRN (Amac et al., 2019) is proposed on RecipeQA, their model does not apply to the textual cloze task. (Park et al., 2017) is using procedural reasoning on multimodal information to generate a story from a sequence of images. However, the textual cloze task is about filling a blank in a sequence given a set of textual options.

Our model exploits the latent alignment space and the positional encoding of questions and answers while applying a novel approach for constraining the output space of the latent alignment. Moreover, we exploit cross-modality representa-

**Pizza Pancakes**





Step1: You need the following ingredients ...  
400 gr. flour 3 eggs ...

Step2: Take a bowl and add the flour and ...

Step3: Take a cutting board and knife ...

Step4: Bake the veggies in separate pieces...

Step5: Heat up the pan and poor a little ...

Step 1
Step 2
Step 3
Step 5

---

Question: Choose the best title for the missing blank to correctly complete the recipe.

\_\_\_\_\_ Making the Dough.      Preparing Veggies.      Baking.

Answers:       A. Preparation      B. Pizza Cones      C. Fillings      D. Cut the Portrait

Figure 1: A sample of textual cloze task

tions based on cross attention to investigate the benefits from information flow between images and instructions. We compare our results to the provided baselines in (Yagcioglu et al., 2018) and achieve the state-of-the-art by improving over 19%.

## 2 Proposed Model

We design a model to solve a structured output prediction on the textual cloze task. The intuition of our model is that the correct answer option should describe precisely one instruction, and this instruction should not be already described with other items in the question. Hence, our model assumes the instruction and question as the context and candidate answers as an additional input to the alignment process. Moreover, to incorporate the order of the sequence in question items and the placeholder, we utilize a one-hot encoding vector of positions to be concatenated with the candidate answers and question items’ representations.

We give the instructions to a sentence splitter using Stanford Core NLP library (Manning et al., 2014). The output is then tokenized by Flair data structure (Akbiik et al., 2018) and embedded with BERT (Devlin et al., 2019). The words’ embeddings are passed to an LSTM layer and the last layer is used as the instruction representation. We propose two different approaches to include images representations. These proposals are described in Section 3.3. An overview of our approach is shown in Figure 2.

Question representation is the last layer of an LSTM on question items. The representation of each question item is the concatenated vector of a one-hot position encoding and word embedding obtained from BERT. The candidate answers’ representations are computed using the same approach.

We concatenate the question representation to each instruction. Then, the similarity of each candidate answer and instruction is computed using the cosine similarity and form a similarity matrix. We use  $S$  to denote the similarity matrix. The rows of this matrix are candidate answers and the columns represent the recipe steps. The value of  $S_{ij}$  indicates the similarity score of candidate  $i$  and step  $j$ .

For training the model, we define two different objectives directly applied to the similarity matrix. The textual cloze task does not have the direct supervision required for the alignment between candidates and steps, and our objective is designed to use the answer of the question to train this latent space of alignments. For imposing the constraint of the alignment to be disjoint between steps and candidates, one way is to simply compute the maximum of each row in the similarity matrix and use that as the aligned step for each candidate answer; However, we introduce constrained max-pooling which is a more sophisticated approach as shown in Figure 3. We compare these two alternatives in the experimental results. We apply an iterative process to select the most related pair of instruction (a column) and answer candidate (a row) while removing the related column and row each time until all candidate answers find their aligned instruction. We denote the final selected maximum scores by  $m = (S_{1i_1}, S_{2i_2}, S_{3i_3}, S_{4i_4})$ , where  $i_c \in [1, number\_of\_steps]$  is the index of the step with maximum alignment score with candidate  $cand$  for all pairs of candidates  $c$  and  $d$ ,  $c \neq d \implies i_c \neq i_d$ .

Respectively, we define two following objectives. The first objective maximizes the distance between the maximum score of the correct answer and the

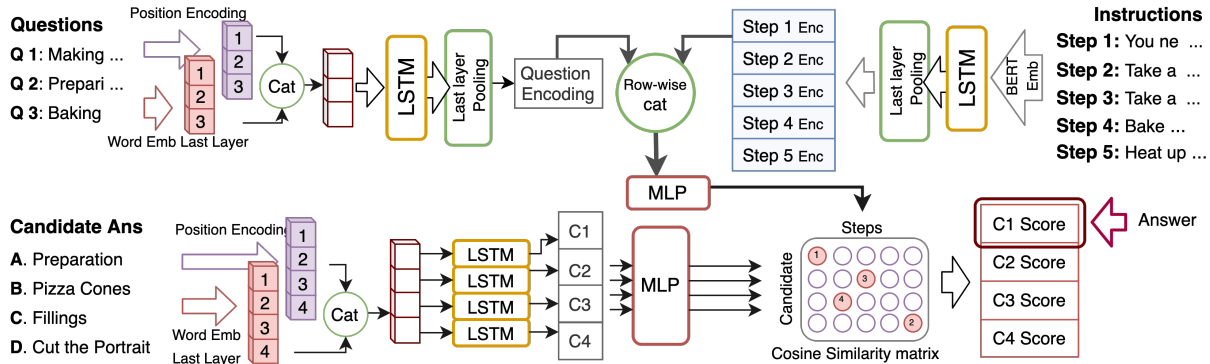


Figure 2: An overview of proposed model

maximum score of another random wrong answer candidate. Furthermore, by fixing the instruction with the maximum alignment with the correct answer, it decreases the score of the other candidates alignments with that instruction. The second objective, increases the maximum similarity score of the answer to approach to 1 while decreasing the other maximum scores to be lower than 0.1.

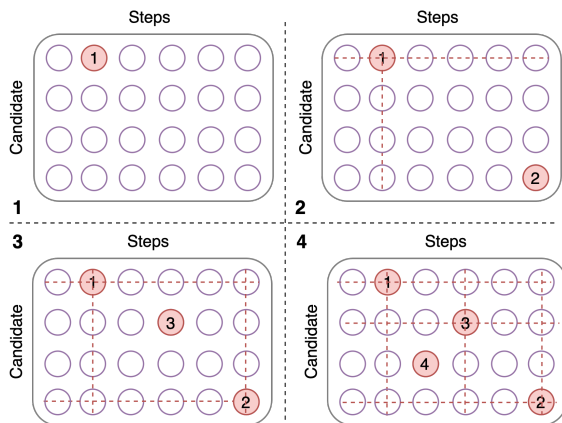


Figure 3: The matrix operation for constrained max-pooling

$$Loss = \max(0, S_{r_i} - S_{a_i} + 0.1) + \sum_{c \neq a} \max(0, S_{c_i} - S_{a_i} + 0.1) \quad (1)$$

$$Loss = (1 - S_{a_i}) + \sum_{c \neq a} \max(0, S_{c_i} - 0.1) \quad (2)$$

Where  $a \in \{1, 2, 3, 4\}$  is the correct answer number and  $r$  is a random index from  $\{1, 2, 3, 4\} - \{a\}$ . The main difference in objective 1 and objective 2 is the regularization term

on the selected instruction column in the alignment matrix.

### 3 Experiment

#### 3.1 Baselines

**Hasty Student** (Tapaswi et al., 2016) is a simple approach considering only the similarity between elements in question and candidate answers. This baseline fails to get good results due to the intrinsic of the task.

**Impatient Reader** (Hermann et al., 2015) computes attention from answers to the recipe for each candidate and despite being a complicated approach, yet it fails to get good results on the task. Moreover, multimodal Impatient reader approach uses both instructions and corresponding images.

#### 3.2 Results

The RecipeQA textual cloze task contains 7837 training, 961 validation, and 963 test examples. A learning rate of  $4 - e1$  is used for the first half and then  $8 - e2$  for the second half of training iterations. We use the momentum of 0.9 for all variations of our model. We train for 30 iterations with a batch size of 1 and optimize the weights using an SGD optimizer. For word embedding, the pre-trained BERT embedding in Flair framework is used. For the image representations, ResNet50 (He et al., 2016) pre-trained on Imagenet (Russakovsky et al., 2015) using PyTorch library (Paszke et al., 2019) is applied.

Table 1 presents the experimental results. We call the model variations which use the loss objective in Equation (1) as Model-obj 1 and the ones that use the loss in Equation (2) as Model-obj 2. Using the objective in Formula (1) yields better results in all experiments. This indicates the benefit

of using the column-wise disjoint constraint on the similarity matrix. Also, using multimodal information yields 1.12% improvement. We elaborate further on the comparison between multimodal and unimodal results in Section 4.

We provide our Pytorch implementation publicly available on Github <sup>1</sup>.

### 3.3 Multimodal Results

In order to investigate the usefulness of the images in solving the textual cloze task, we propose two different models that incorporate the image representation in addition to the textual information of recipe steps. The first variation receives ResNet50 representations of the images and, after applying an LSTM layer, pulls the last layer as image representation. Finally, it concatenates the image representation to the question and instruction representation in the main architecture before applying the MLP and computing the cosine similarities.

The second variation as shown in figure 4, uses a more complex architecture introduced in LXMERT (Tan and Bansal, 2019). We modify the architecture of LXMERT and apply it to the word embedding and image representations to flow the information from each to another. The updated word embedding and image representations are passed to an LSTM, and its last layer is used to represent the visual and textual information of a step. In the end, these representations are concatenated to each other and the question representation to build the instruction vector representation. We report the results of these model variations in Table 1. Using the cross modality representations based on LXMERT provided extensive way to flow the information from text and image to each other and yields the best results.

## 4 Discussion and Analysis

We did qualitative analysis using some examples and their results to better understand the behaviour of the proposed model. Our model is almost able to detect all matched candidates with the instructions (in case that there exist multiple matches) but fails to choose the one that completes the sequence of the question items. This indicates the shortage of procedural hints inside our architecture while the latent alignment is proven to be practical. By analysing the results, we found interesting cases

<sup>1</sup><https://github.com/HLR/LatentAlignmentProcedural>

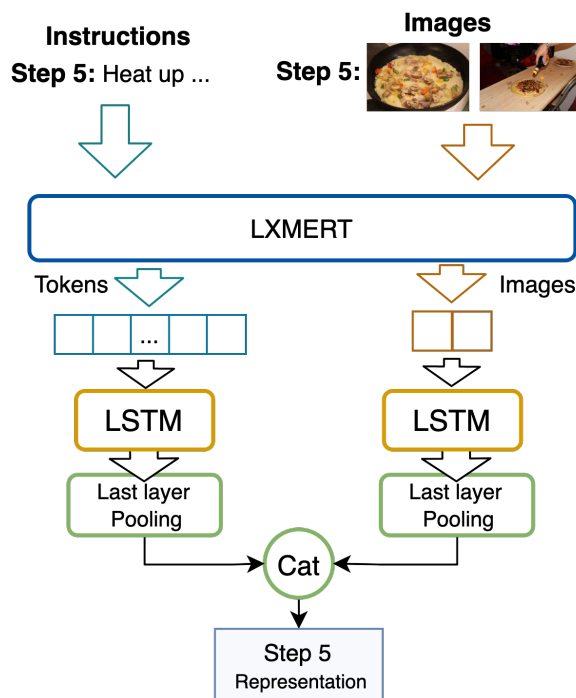


Figure 4: Using LXMERT for integrating multimodal information on steps

where either multimodal or unimodal architectures could yield more accurate predictions.

#### Multimodal -, Unimodal +:

- Images contain misleading information (see example in Figure 5).
- Image quality is low.
- Images are not showing the steps correctly.
- Text contains direct mentions of candidate answers.

Questions  
Materials Needed \_\_\_\_\_? Sprinkling Microwave

Candidate answers:  
A) Cutting B. Halloween Baked Apples  
C. Turn Apple Slices D. Pile It Up

Uni-modal scores A: -0.520 B: -0.688 C: -0.607 D: -0.625  
Multi-modal Scores A: -0.129 B: -0.249 C: -0.088 D: -0.153

Figure 5: The image is misleading the multimodal setting to choose apple slices rather than cutting option

#### Multimodal +, Unimodal -:

- The sequence of the images provide detailed steps and good quality.
- The entities in candidates answers are shown in the pictures but not in the text.
- The recipes instructions are very short and the images provide more information.

Models	Accuracy	p@2
Human	73.6	-
Hasty Student	26.89	-
Impatient Reader	28.03	-
Impatient Reader (multimodal)	29.07	-
Model-Obj 1	46.35	<b>78.7</b>
Model-Obj 2	43.36	-
Model-Obj 1 (multimodal)	45.41	
Model-Obj 1 (multimodal) + LXMERT	<b>47.5</b>	77.5
Model-Obj 1 (multimodal) + LXMERT - ConstrainedMaxPooling	46.9	76.3

Table 1: Evaluation on the test set

In some cases, the multimodal information can fix the errors resulted from not considering the order of events in the proposed architecture. Our intuition is that, although, the textual model does not contain information from previous steps, the images carry useful information on what has been already done. An example of this is shown in Figure 6, where co-reference resolution is required to answer the question correctly.

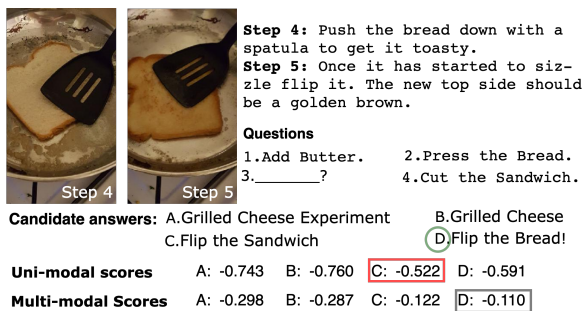


Figure 6: The images lead the model to understand that "it" refers to bread rather than sandwich

Furthermore, we have tested our multimodal architecture with representations of ResNet101 and the results dropped. We confirmed this experiment by re-implementing Hasty Student approach on visual coherence task (that has 68% accuracy with ResNet50) and obtained 35% lower than ResNet50. This can be due to the lack of quality of images resulting in extra noise when using a more complicated network. Thus, ResNet50 achieves better accuracy by producing more abstract representations of the images.

## 5 Conclusion and Future Work

We proposed a model for RecipeQA textual cloze task which exploits the latent alignment of question items with instructions. Moreover, we investigated

the benefit of using multimodal information in this task by comparing three different architectures and provided qualitative analysis on some examples to justify the results. Our model exceeded the baselines and improved the SOTA by over 19%. As a future direction, we will investigate the usage of the latent alignment in other tasks. We will apply more complex methods on textual abstractions and attention mechanisms to link the candidate answers with the recipe instructions. Investigating how to incorporate the question order in the architecture is another direction.

## Acknowledgments

This work is (partially) supported by the Office of Naval Research grant N00014-19-1-2308.

## References

- Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649.
- Mustafa Sercan Amac, Semih Yagcioglu, Aykut Erdem, and Erkut Erdem. 2019. [Procedural reasoning networks for understanding multimodal procedures](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 441–451, Hong Kong, China. Association for Computational Linguistics.
- Antoine Bosselut, Omer Levy, Ari Holtzman, Corin Ennis, Dieter Fox, and Yejin Choi. 2017. Simulating action dynamics with neural process networks. In *Sixth International Conference on Learning Representations (ICLR)*.
- Bhavana Dalvi, Lifu Huang, Niket Tandon, Wen-tau Yih, and Peter Clark. 2018. [Tracking state changes in procedural text: a challenge dataset and models](#)

- for process paragraph comprehension. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1595–1604, New Orleans, Louisiana. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of deep bidirectional transformers for language understanding**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Haoqi Fan and Jiatong Zhou. 2018. Stacked latent attention for multimodal reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1072–1080.
- 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.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in neural information processing systems*, pages 1693–1701.
- Jack Hessel, Lillian Lee, and David Mimno. 2019. Unsupervised discovery of multimodal links in multi-image, multi-sentence documents. In *EMNLP*.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven 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.
- Hyeonseob Nam, Jung-Woo Ha, and Jeonghee Kim. 2017. Dual attention networks for multimodal reasoning and matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 299–307.
- Cesc Chunseong Park, Youngjin Kim, and Gunhee Kim. 2017. Retrieval of sentence sequences for an image stream via coherence recurrent convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):945–957.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. **PyTorch: An imperative style, high-performance deep learning library**. In H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. **Vi-bert: Pre-training of generic visual-linguistic representations**. In *Eighth International Conference on Learning Representations (ICLR)*.
- Hao Tan and Mohit Bansal. 2019. **LXMERT: Learning cross-modality encoder representations from transformers**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111, Hong Kong, China. Association for Computational Linguistics.
- Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2016. Movieqa: Understanding stories in movies through question-answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4631–4640.
- Hao Wu, Jiayuan Mao, Yufeng Zhang, Yuning Jiang, Lei Li, Weiwei Sun, and Wei-Ying Ma. 2019. Unified visual-semantic embeddings: Bridging vision and language with structured meaning representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6609–6618.
- Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. 2018. **RecipeQA: A challenge dataset for multimodal comprehension of cooking recipes**. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1358–1368, Brussels, Belgium. Association for Computational Linguistics.
- Zhou Yu, Yuhao Cui, Jun Yu, Dacheng Tao, and Qi Tian. 2019. Multimodal unified attention networks for vision-and-language interactions. *arXiv preprint arXiv:1908.04107*.