

Generating Images from Captions with Attention

Elman Mansimov, Emilio Parisotto, Jimmy Ba and Ruslan Salakhutdinov

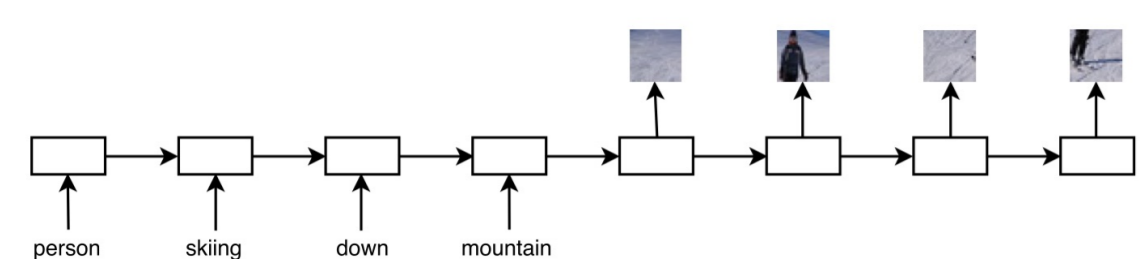
INTRODUCTION

Why to condition on captions ?

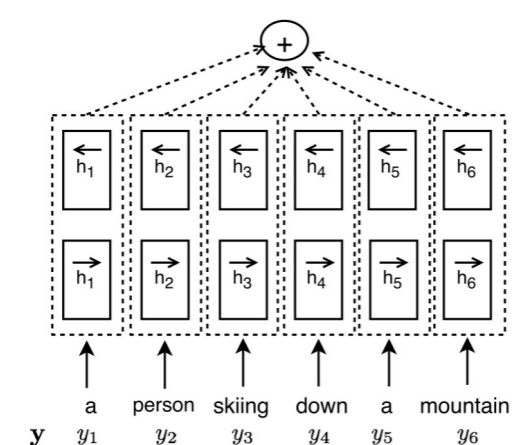
- Captions could be used to simplify image modelling task.
- Generating images conditioned on novel captions helps better understand its generalization.

Key Ideas

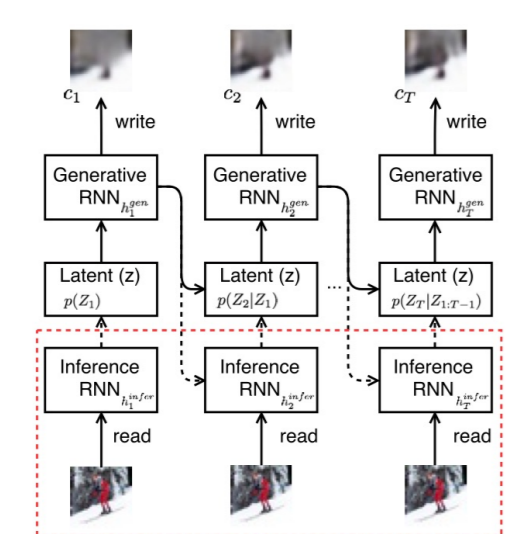
- Treat the problem as part of sequence-to-sequence framework [2, 9].



- Caption y is represented as sequence of words (y_1, y_2, \dots, y_N) , where N is the length of the sequence.



- Image x is represented as a sequence of $p \times p$ patches drawn on a $w \times h$ canvas c_t over time $t = 1, \dots, T$.

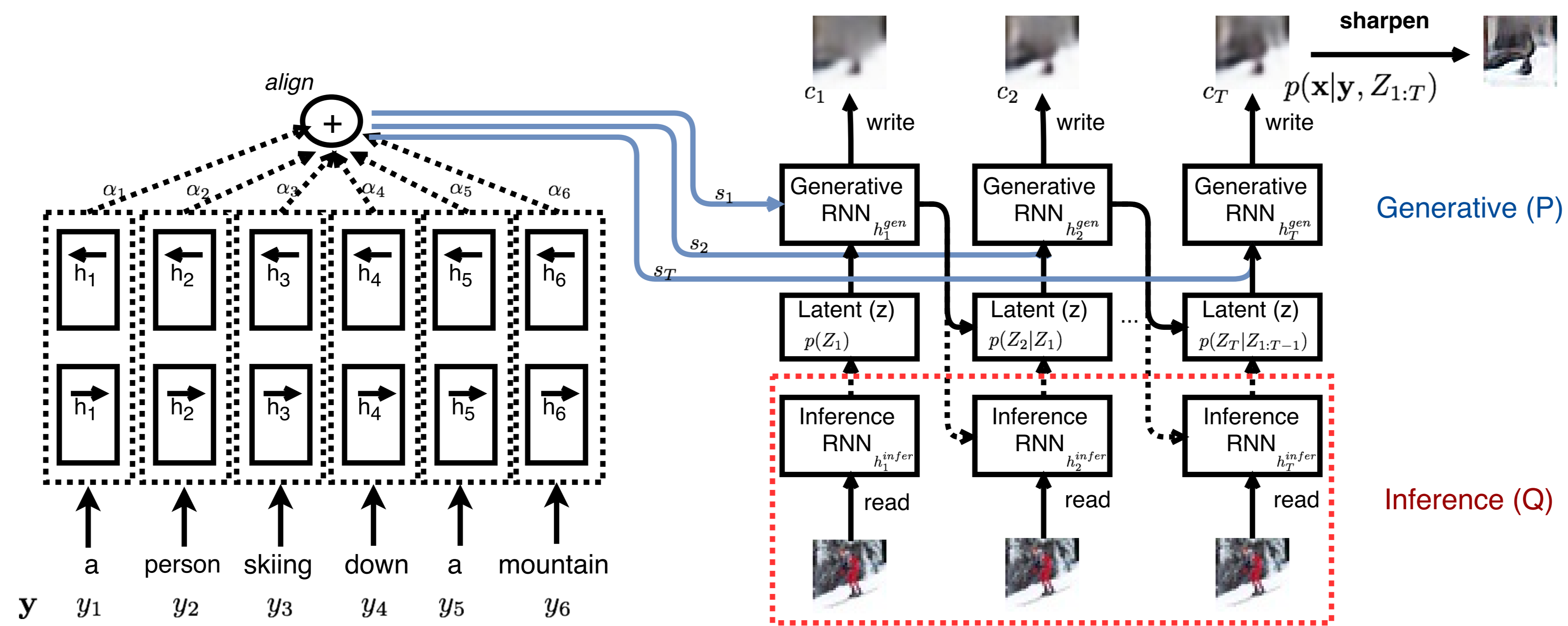


REFERENCES

- [1] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015.
- [2] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *EMNLP*, 2014.
- [3] E. Denton, S. Chintala, A. Szlam, and R. Fergus. Deep generative image models using a laplacian pyramid of adversarial networks. In *NIPS*, 2015.
- [4] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [5] K. Gregor, I. Danihelka, A. Graves, and D. Wierstra. DRAW: A recurrent neural network for image generation. In *ICML*, 2015.
- [6] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 1997.
- [7] D. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014.
- [8] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.
- [9] I. Sutskever, O. Vinyals, and Q. Le. Sequence to sequence learning with neural networks. In *NIPS*, 2014.

MODEL DESCRIPTION

The proposed model consists of four main parts: language model, image model, alignment model and post-processing model.



- Bidirectional LSTM [6] computes a sequence of forward $\vec{h}_{1..N}^{lang}$ and backward $\overleftarrow{h}_{1..N}^{lang}$ hidden states respectively, which are concatenated together into the sentence representation $h^{lang} = [\vec{h}_{1..N}^{lang}, \overleftarrow{h}_{1..N}^{lang}]$.

- The image model [5] iteratively computes the following set of equations over time $t = 1, \dots, T$:

$$\begin{aligned} \hat{x}_t &= x - \sigma(c_{t-1}), & z_t &\sim P(Z_t | Z_{1:t-1}) = \mathcal{N}(\mu(h_{t-1}^{gen}), \sigma(h_{t-1}^{gen})), \\ r_t &= read(x_t, \hat{x}_t, h_{t-1}^{gen}), & s_t &= align(h_{t-1}^{gen}, h^{lang}), \\ h_t^{infer} &= LSTM^{infer}(h_{t-1}^{infer}, [r_t, h_{t-1}^{gen}]), & h_t^{gen} &= LSTM^{gen}(h_{t-1}^{gen}, [z_t, s_t]), \\ Q(Z_t | x, y, Z_{1:t-1}) &= \mathcal{N}(\mu(h_t^{infer}), \sigma(h_t^{infer})), & c_t &= c_{t-1} + write(h_t^{gen}), \\ & & \tilde{x} &\sim P(x | y, Z_{1:T}) = \prod_i P(x_i | y, Z_{1:T}) = \prod_i \text{Bern}(\sigma(c_{T,i})). \end{aligned}$$

- $align$ operator [1] outputs a dynamic sentence representation s_t at each timestep by computing a weighted sum of hidden states of words using alignment probabilities $\alpha_{1..N}^t$:

$$s_t = align(h_{t-1}^{gen}, h^{lang}) = \alpha_1^t h_1^{lang} + \alpha_2^t h_2^{lang} + \dots + \alpha_N^t h_N^{lang} \quad \alpha_k^t \propto \exp(v^T \tanh(Uh_k^{lang} + Wh_{t-1}^{gen} + b))$$

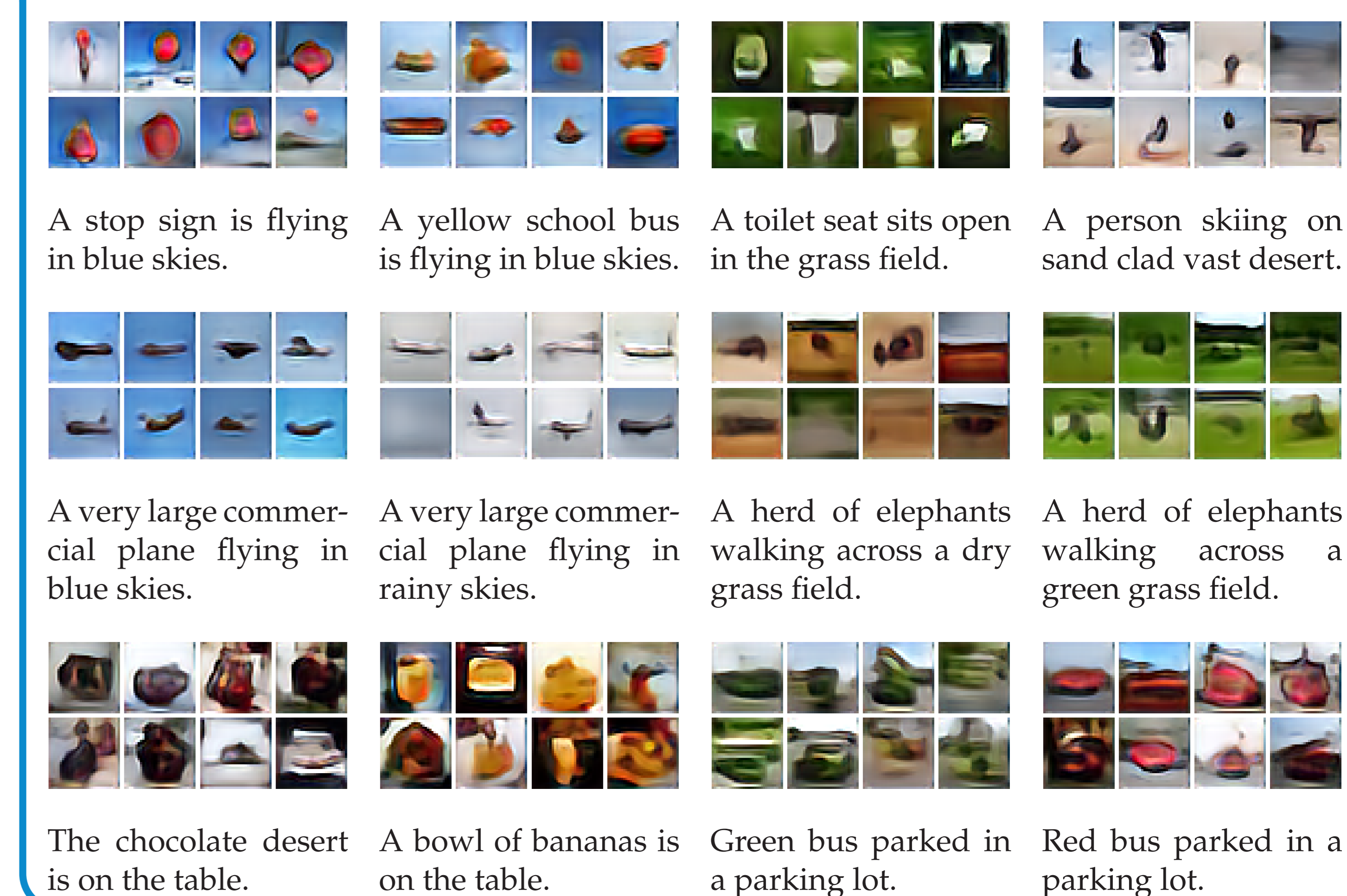
- Images are generated by discarding the inference network and by sampling latent variables $Z_{1:t}$ from prior distribution.
- Finally, images are sharpened using a deterministic adversarial network [3, 4] trained on residuals of a Laplacian pyramid.

LEARNING

The model is trained to optimize the variational lower bound \mathcal{L} of image x given caption y using the SGVB [7] algorithm.

$$\mathcal{L} = \mathbb{E}_{Q(Z_{1:T} | y, x)} \left[\log p(x | y, Z_{1:T}) - \sum_{t=2}^T D_{KL}(Q(Z_t | Z_{1:t-1}, y, x) \| P(Z_t | Z_{1:t-1}, y)) \right] - D_{KL}(Q(Z_1 | x) \| P(Z_1)).$$

GENERATED IMAGES



EXAMPLES OF ALIGNMENT



MICROSOFT COCO [8] RESULTS

| Model | Microsoft COCO (before post-processing) | | | | | Similarity SSI |
|-----------------|---|------|------|------|-------|----------------|
| | R@1 | R@5 | R@10 | R@50 | Med r | |
| LAPGAN | - | - | - | - | - | 0.08 |
| Fully-Conn VAE | 1.0 | 6.6 | 12.0 | 53.4 | 47 | 0.156 |
| Conv-Deconv VAE | 1.0 | 6.5 | 12.0 | 52.9 | 48 | 0.164 |
| skipthoughtDRAW | 2.0 | 11.2 | 18.9 | 63.3 | 36 | 0.157 |
| noalignDRAW | 2.8 | 14.1 | 23.1 | 68.0 | 31 | 0.155 |
| alignDRAW | 3.0 | 14.0 | 22.9 | 68.5 | 31 | 0.156 |

Top: Image retrieval and similarity results of different models. **R@K** is Recall@K (higher is better). **Med r** is the median rank (lower is better). **SSI** is Structural Similarity Index, which is between -1 and 1 (higher is better).