



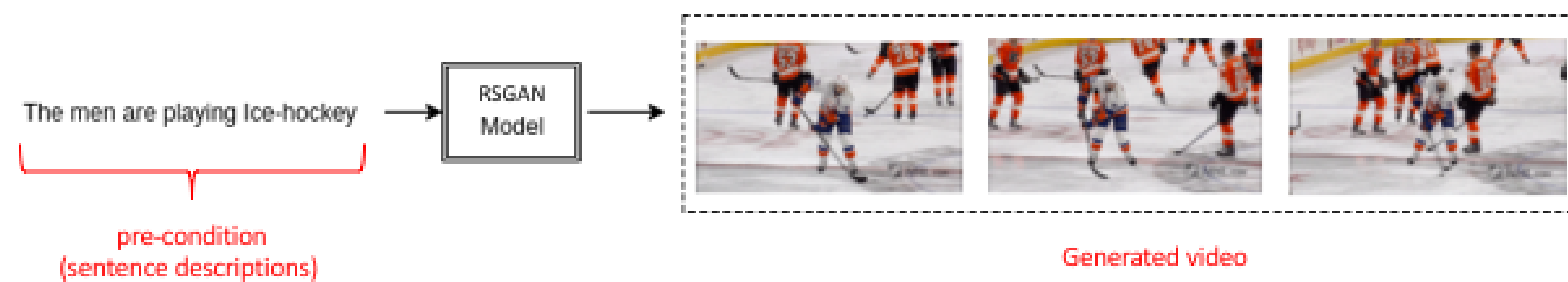
RSGAN: RECURRENT STACKED GENERATIVE ADVERSARIAL NETWORK FOR CONDITIONAL VIDEO GENERATION

S. NAHA, K. M. NAYEM AND M. L. ISLAM SCHOOL OF INFORMATICS AND COMPUTING, INDIANA UNIVERSITY

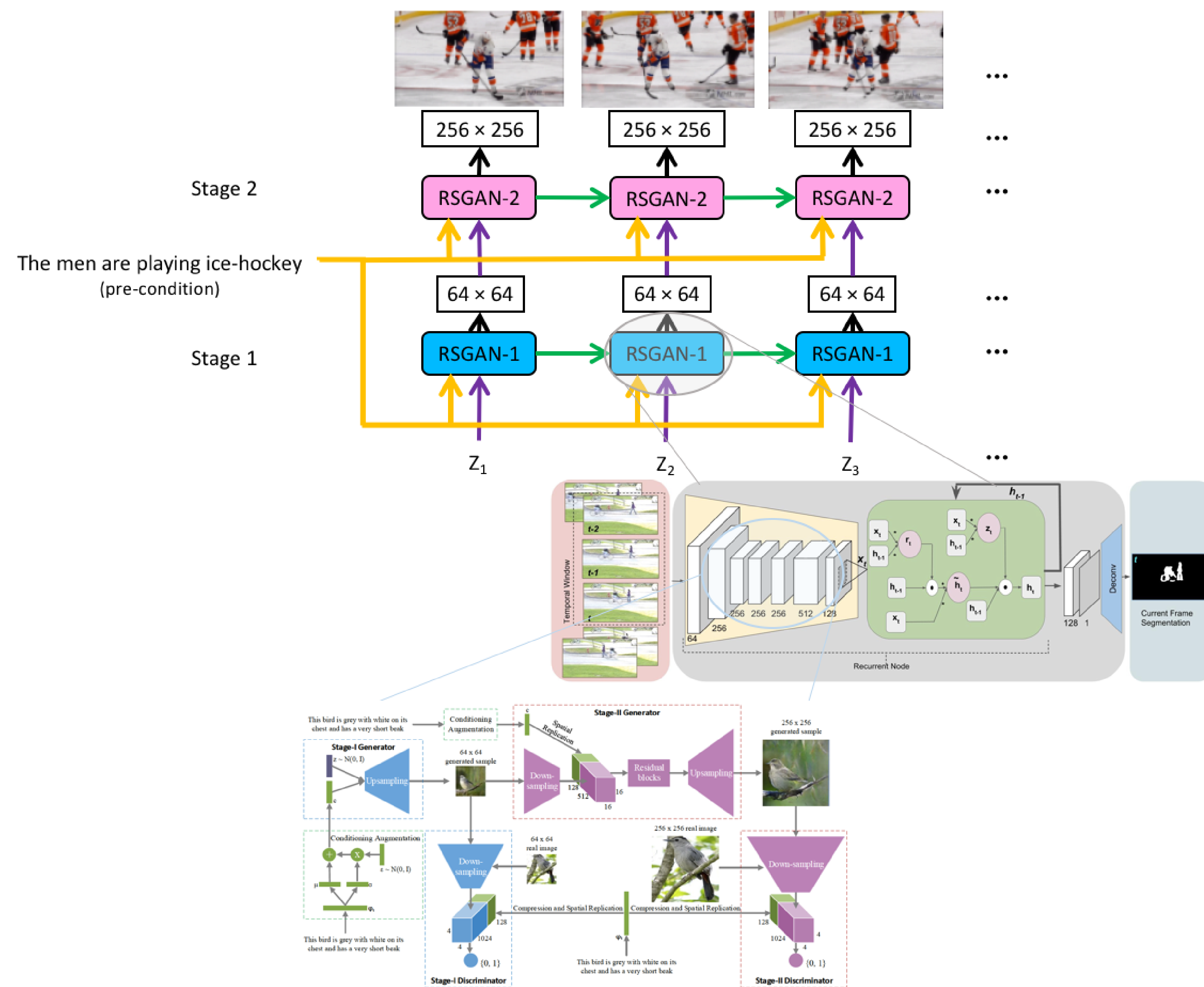
INTRODUCTION

Generating video frames based on a pre-condition is a challenging problem and requires understanding of per frame contents and visual dynamics and their relevacies to the pre-condition. In this project, we propose a novel Recurrent Stacked Generative Adversarial Network (RSGAN) based model to generate video frames based on a given pre-condition. In our knowledge, this is the first work to address the problem of conditional video generation using adversarial network. We can address the problem of generating videos based on pre-conditions such as,

1. action classes
2. fMRI signals
3. sentence descriptions



OUR APPROACH



Each RSGAN module is - consist of a StackGAN model.
- connected by a Fully Convolutional LSTM Network.

OBJECTIVE FUNCTION

Conditioned on Gaussian latent variables c_0 , Stage-I RSGAN trains discriminator D_0 and generator G_0 by alternatively maximizing \mathcal{L}_{D_0} and minimizing \mathcal{L}_{G_0} .

Stage-I RSGAN:

$$\mathcal{L}_{D_0} = \mathbb{E}_{(I_0, t) \sim p_{data}} [\log D_0(I_0, \varphi_t)] + \mathbb{E}_{z \sim p_z, t \sim p_{data}} [\log(1 - D_0(G_0(z, c_0), \varphi_t))]$$

$$\mathcal{L}_{G_0} = \mathbb{E}_{z \sim p_z, t \sim p_{data}} [\log(1 - D_0(G_0(z, c_0), \varphi_t))] + \lambda D_{KL}(\mathcal{N}(\mu_0(\varphi_t), \Sigma_0(\varphi_t)) || \mathcal{N}(0, I))$$

Conditioned on the low resolution sample s_0 and Gaussian latent variables c , discriminator D and generator G in Stage-II RSGAN is trained by alternatively maximizing \mathcal{L}_D and minimizing \mathcal{L}_G .

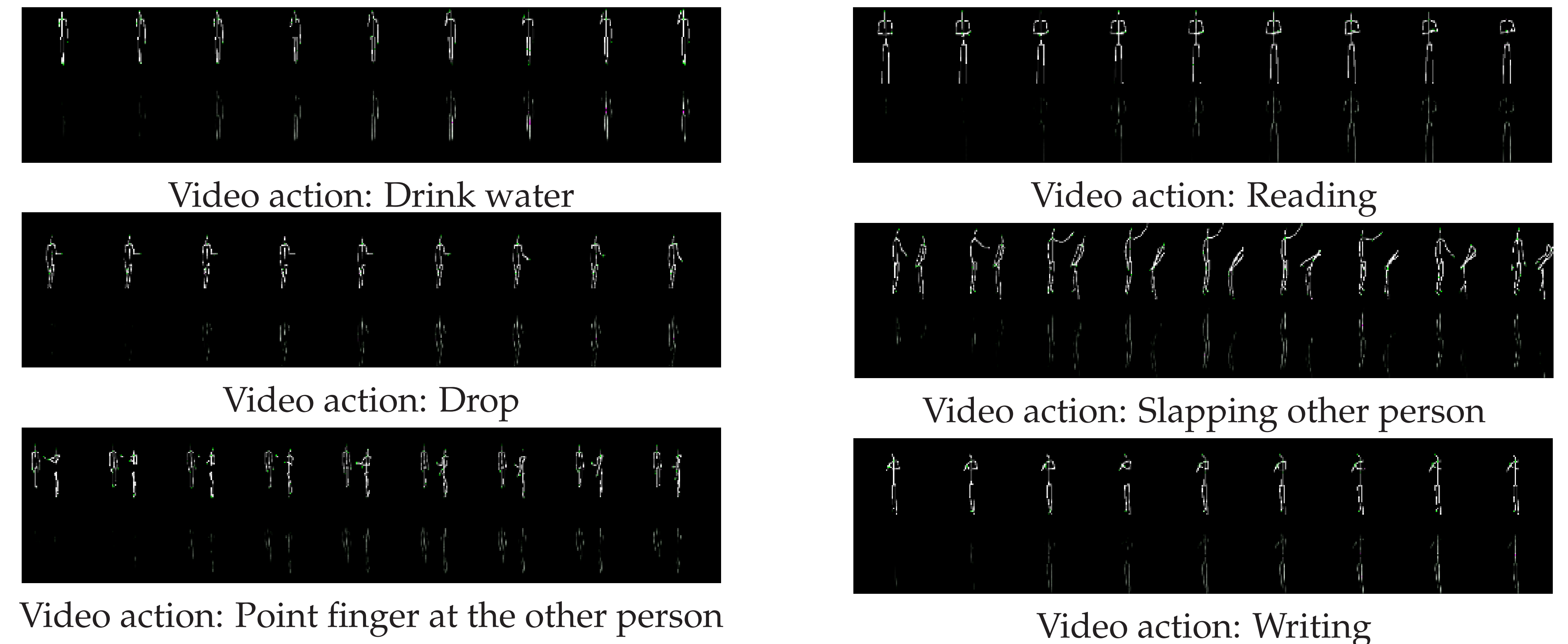
Stage-II RSGAN:

$$\mathcal{L}_D = \mathbb{E}_{(I, t) \sim p_{data}} [\log D(I, \varphi_t)] + \mathbb{E}_{s_0 \sim p_{G_0}, t \sim p_{data}} [\log(1 - D(G(s_0, c), \varphi_t))]$$

$$\mathcal{L}_G = \mathbb{E}_{s_0 \sim p_{G_0}, t \sim p_{data}} [\log(1 - D(G(s_0, c), \varphi_t))] + \lambda D_{KL}(\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t)) || \mathcal{N}(0, I))$$

RESULT

Ground truth (1st row) and the partial result on Convolutional LSTM (2nd row).



• Right now, we are trying to generate video with simple details. That's why we are using NTU RGB+D Action Recognition Dataset (skeletal data).

FUTURE RESEARCH

- Generate video with complex details and multiple moving objects.
- Use fMRI dataset of human brain, to generate video.

REFERENCES

- [1] H. Zhang, T. Xu, H. Li, S. Zhang, X. Huang, X. Wang, and D. Metaxas. StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks. *arXiv preprint arXiv:1612.03242*, 2016.
- [2] S. Valipour, M. Siam, M. Jagersand, and N. Ray. Recurrent Fully Convolutional Networks for Video Segmentation. *arXiv preprint arXiv:1611.09904*, 2016.
- [3] O. Mogren. C-RNN-GAN: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:1611.09904*, 2016.