



<http://dmirlab.com>

Improved Conditional VRNNs for Video Prediction

Lluís Castrejon*

Mila, Université de Montréal

lluis.castrejon@gmail.com

Nicolas Ballas

Facebook AI Research

nballas@fb.com

Aaron Courville

CIFAR, Mila, Université de Montréal

aaron.courville@umontreal.ca

Hao Zhang

2020/02/25

Outline



<http://dmirlab.com>

- Short-term prediction
- Long-term prediction

Improved VRNN



<http://dmirlab.com>

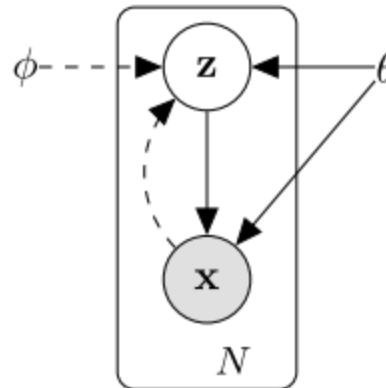
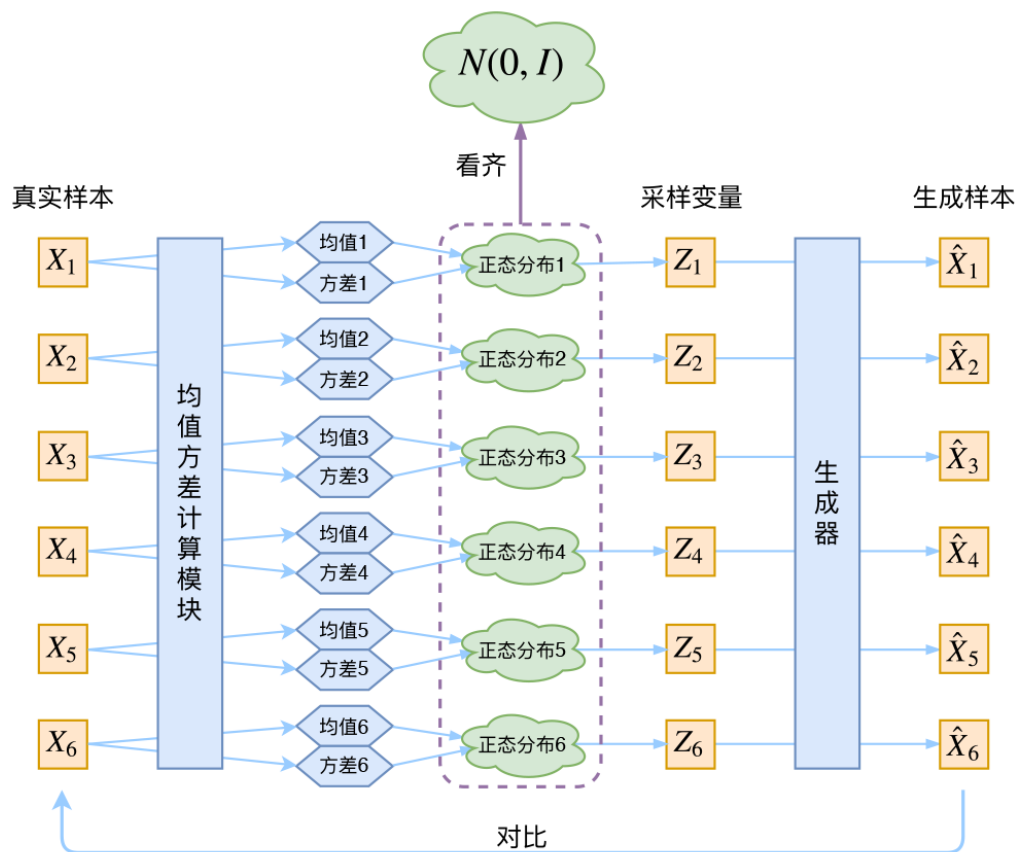


Figure 1: The type of directed graphical model under consideration. Solid lines denote the generative model $p_{\theta}(z)p_{\theta}(x|z)$, dashed lines denote the variational approximation $q_{\phi}(z|x)$ to the intractable posterior $p_{\theta}(z|x)$. The variational parameters ϕ are learned jointly with the generative model parameters θ .

Improved VRNN



<http://dmirlab.com>



Improved Conditional VRNNs for Video Prediction (Lluis et al, ICCV2019)

<http://dmirlab.com>

Improved VRNN



<http://dmirlab.com>

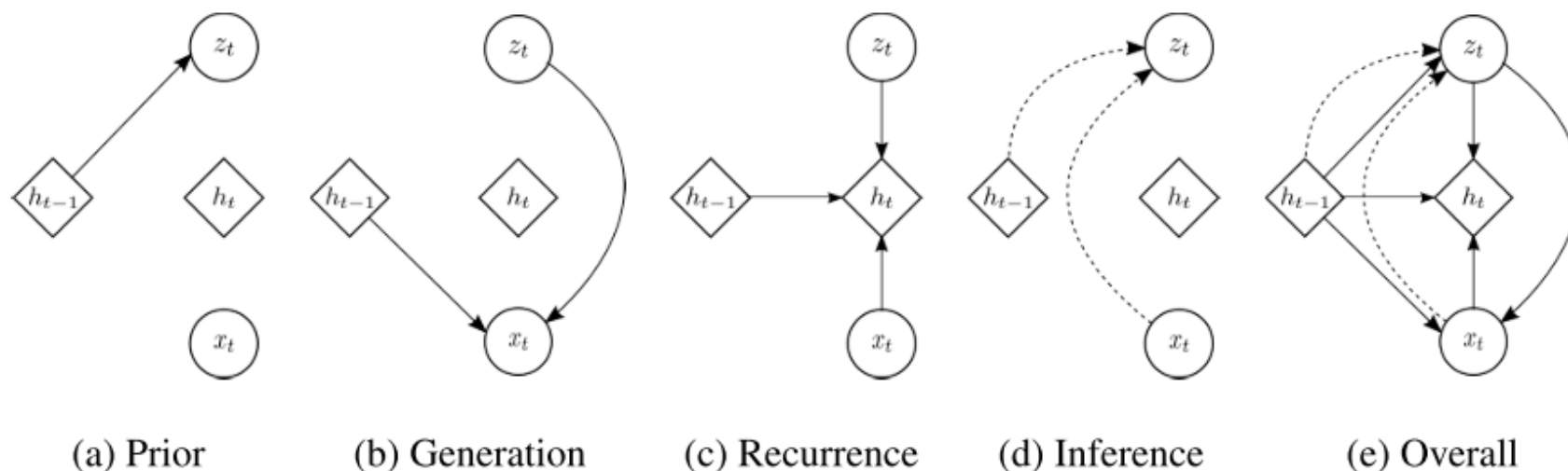


Figure 1: Graphical illustrations of each operation of the VRNN: (a) computing the conditional prior using Eq. (5); (b) generating function using Eq. (6); (c) updating the RNN hidden state using Eq. (7); (d) inference of the approximate posterior using Eq. (9); (e) overall computational paths of the VRNN.

Improved VRNN



<http://dmirlab.com>

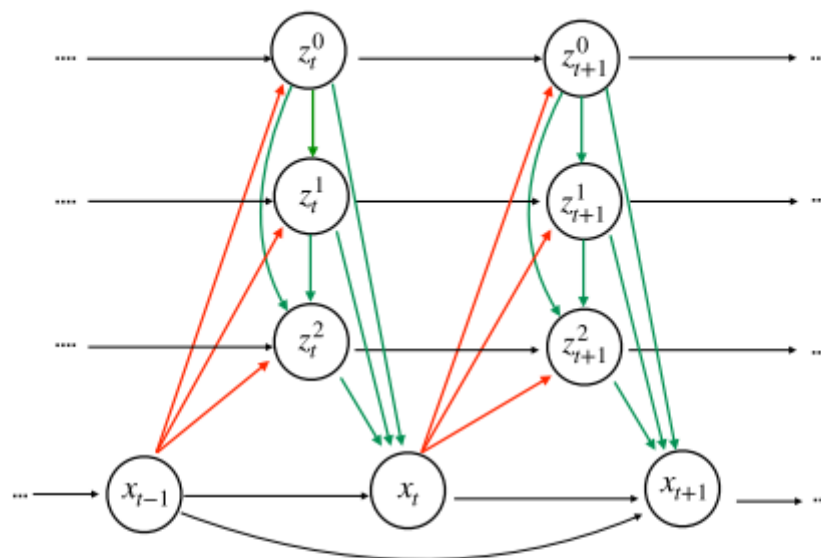


Figure 2: **Graphical model for the learned prior with the dense latent connectivity pattern.** Arrows in **red** show the connections from the input at the previous timestep to current latent variables. Arrows in **green** highlight **skip connections** between latent variables and connections to outputs. Arrows in **black** indicate recurrent temporal connections. We empirically observe that this dense-connectivity pattern **eases** the training of **latent hierarchies**.

Improved VRNN



<http://dmirlab.com>

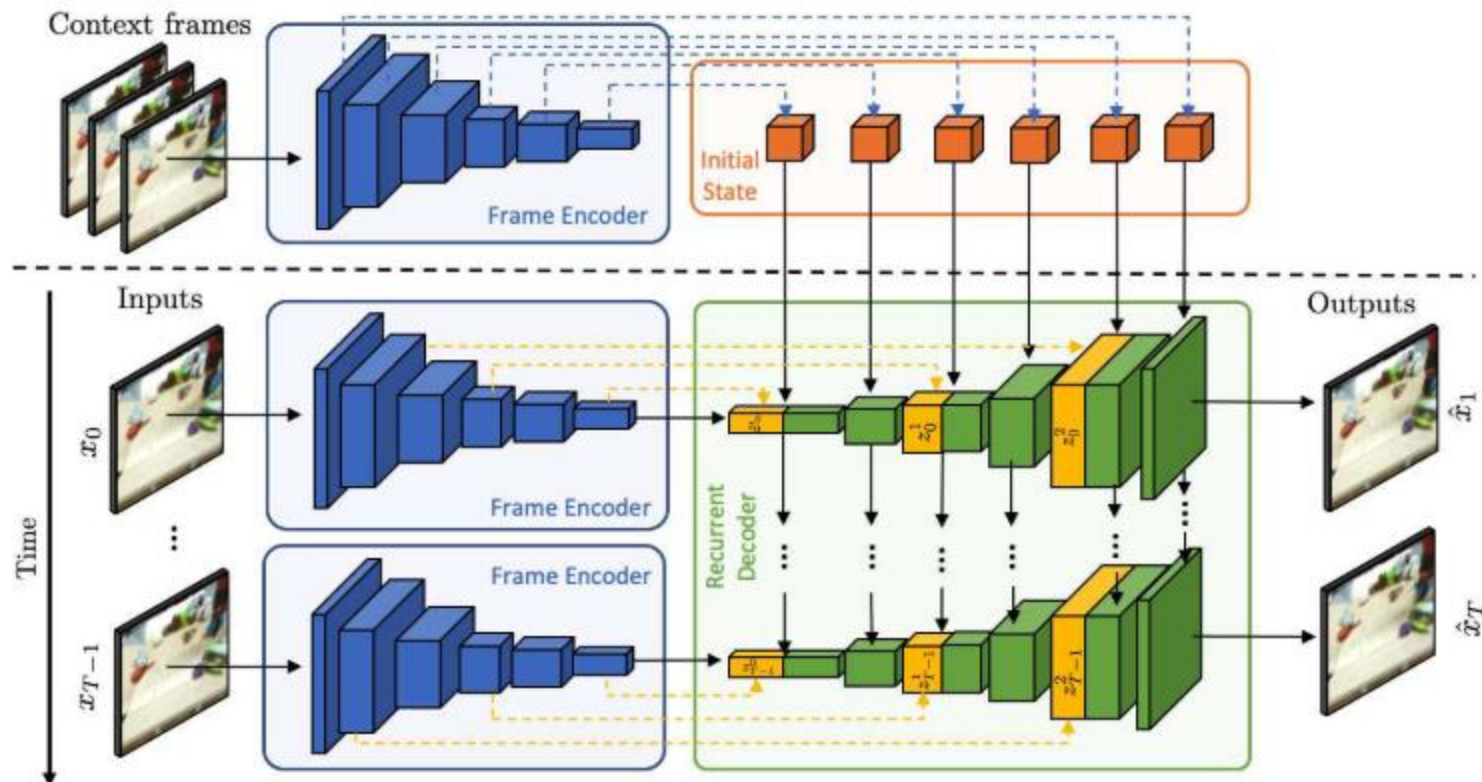


Figure 3: **Model Parametrization.** Our model uses a CNN to encode frames individually. The representation of the context frames is used to initialize the states of the prior, posterior and likelihood networks, all of which use recurrent networks. At each timestep, the decoder receives an encoding of the previous frame, a set of latent variables (either from the prior or the posterior) and its previous hidden state and predicts the next frame in the sequence.

Improved Conditional VRNNs for Video Prediction (Lluis et al, ICCV2019)

<http://dmirlab.com>



<http://dmirlab.com>

Thank you !
Q&A

<http://dmirlab.com>