

1 Introducing Clockwork VAEs

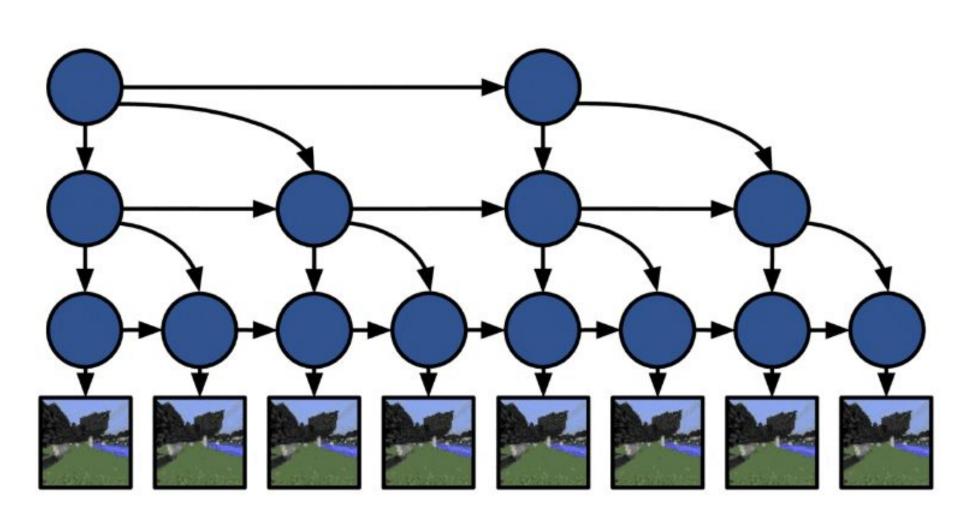
Clockwork VAE is a temporally-abstract recurrent latent variable models.

It predicts the future on 2 multiple time-scales, preserving long-term dependencies.

It can predict upto 1000 frames 3 on video datasets while preserving high-level details.

It can separate and store slow-changing semantics at higher levels of the hierarchy.

It adapts to the speed of the 5 sequence, shifting information across slow and fast moving recurrent chains.

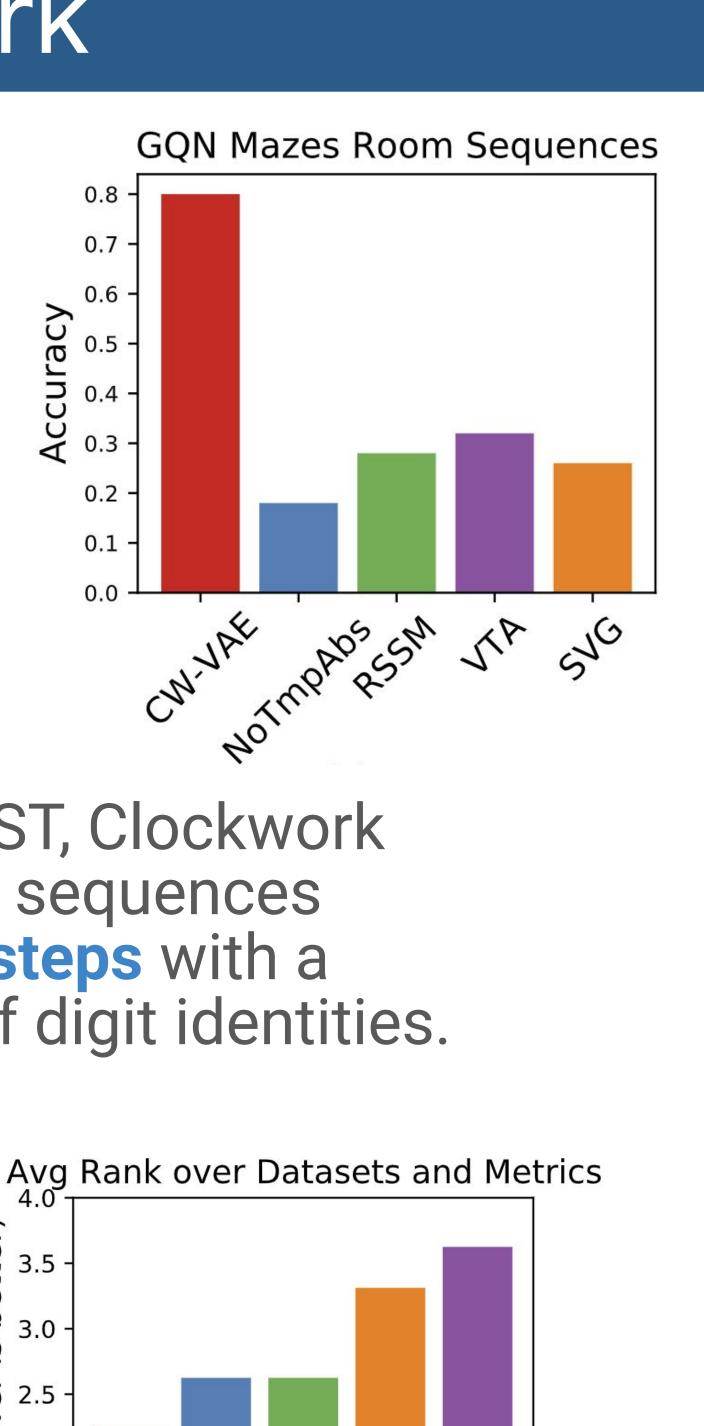


Temporally abstract recurrent chains in a Clockwork VAE with 3 levels and temporal abstraction factor 2.

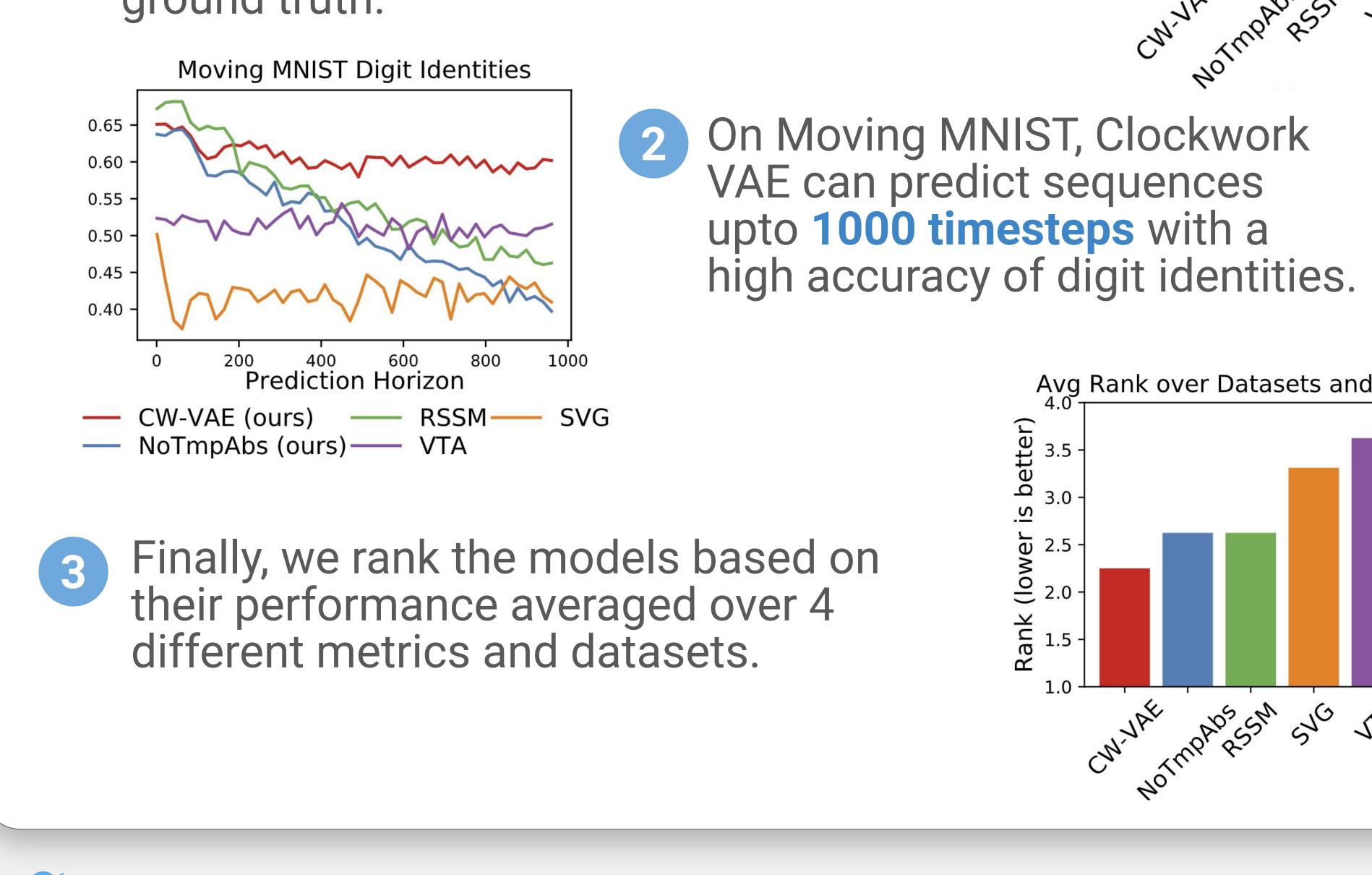
The amount of spatial abstraction determines the number of parameters in the model. Adding temporal abstraction does not affect the model size.

4 Video Prediction Benchmark

For GQN Maze videos, we compute the prediction accuracy for three different high-level categories of sequences of rooms: agent staying in the same room, going into the hallway, and then coming back to the original room. Video generated by Clockwork VAE greatly matches with the ground truth.



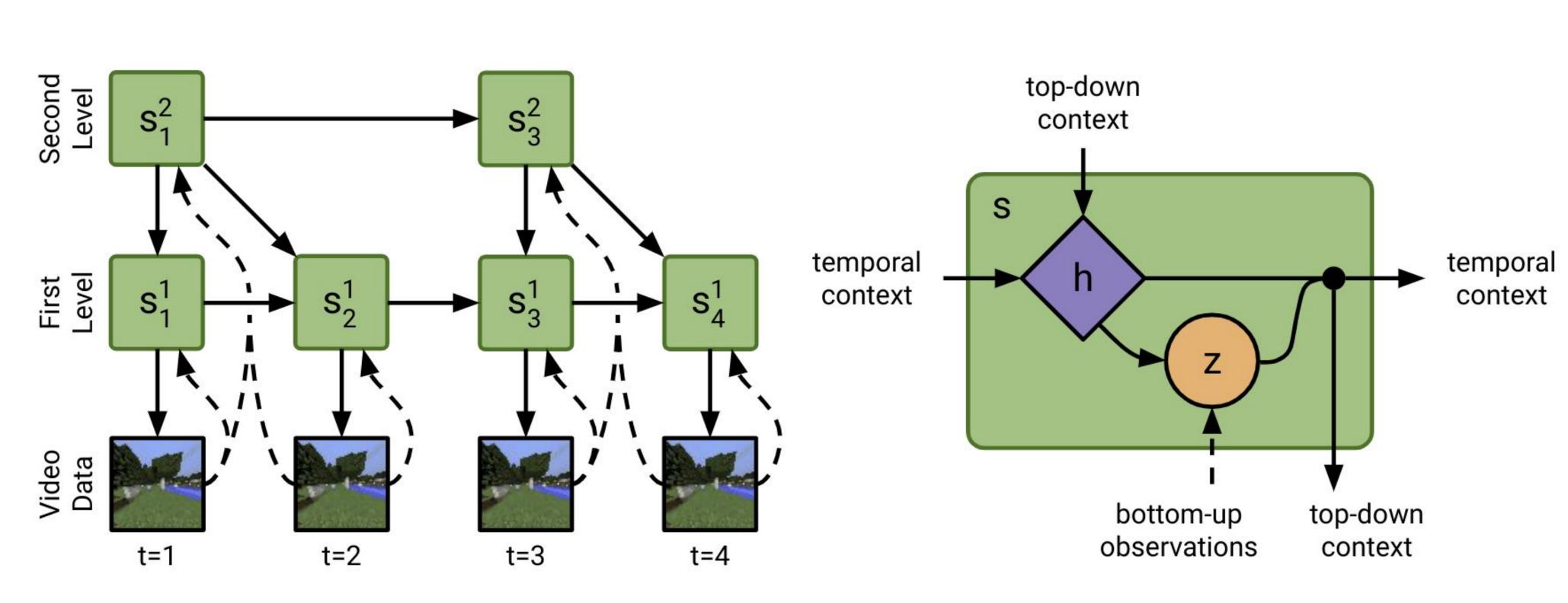
OADSSM SUG JA



Clockwork Variational Autoencoders

Vaibhav Saxena¹², Jimmy Ba¹², Danijar Hafner¹²³ ¹University of Toronto, ²Vector Institute, ³Google Brain

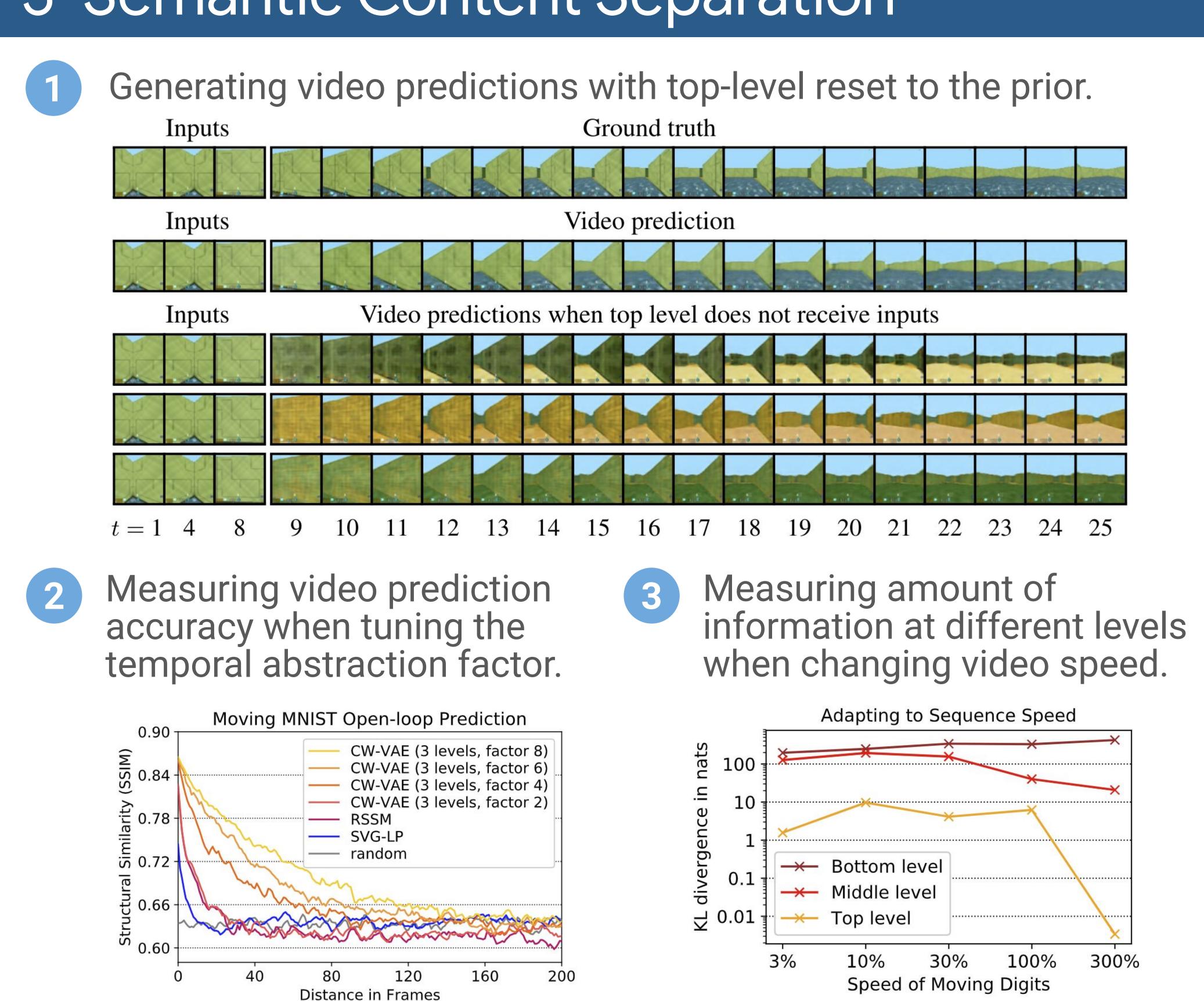
2 Hierarchy of Latent Sequences



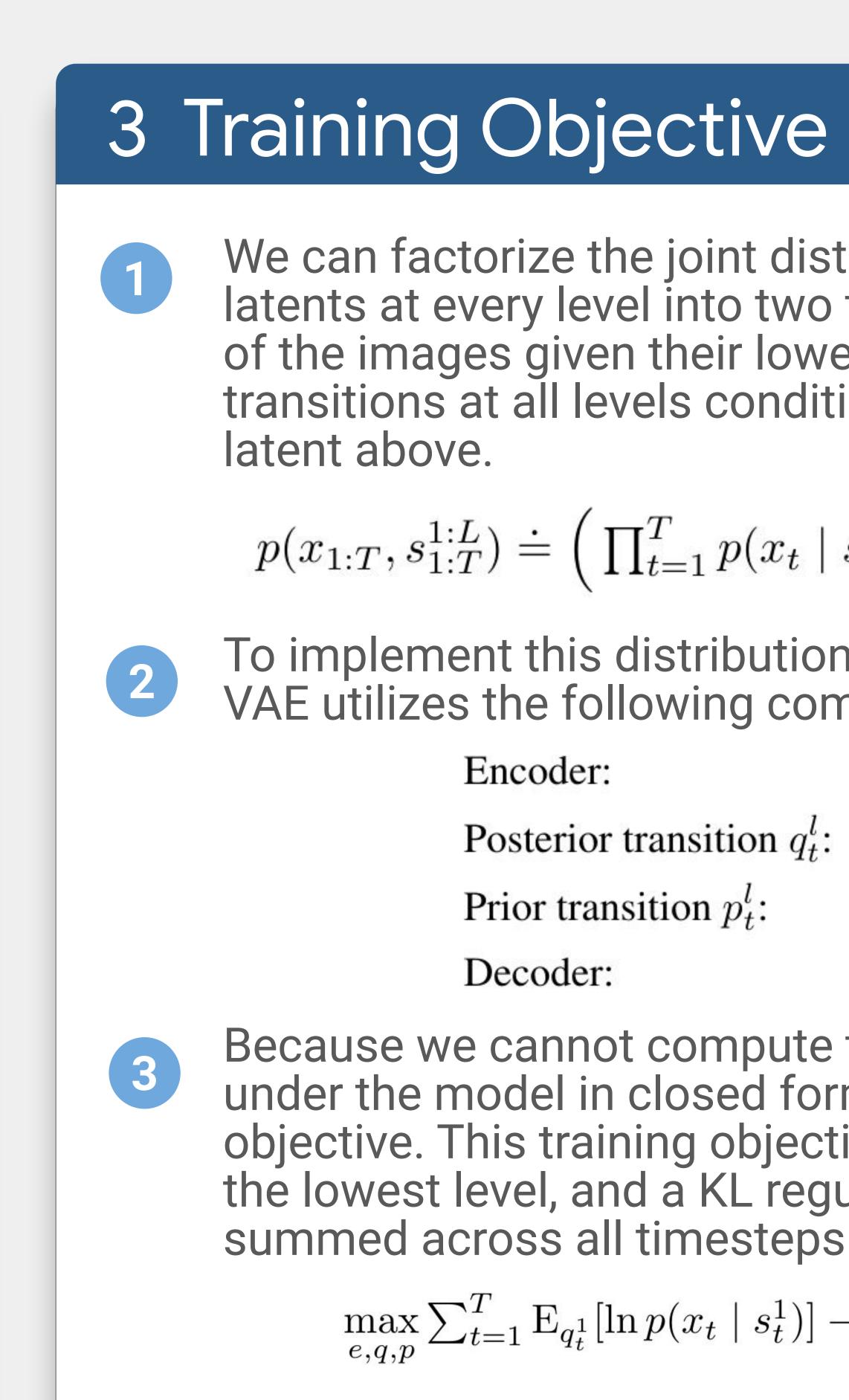
Connection used for bot prior and posterior

- Clockwork VAE consists of a hierarchy of recurrent latent variables, where each level transitions at a different clock speed.
- Transitions slow down exponentially as we go up in the hierarchy • with a factor we call the temporal abstraction factor.
- The posterior at each level is comprised of top-down and bottom-up information, while the prior (generation) is only top-down.
- Each latent state is a combination of a deterministic and a stochastic state, with the weights determining the deterministic state shared by the posterior and prior computations.

5 Semantic Content Separation



Connection only used for posterior (given data)



6 Find out more

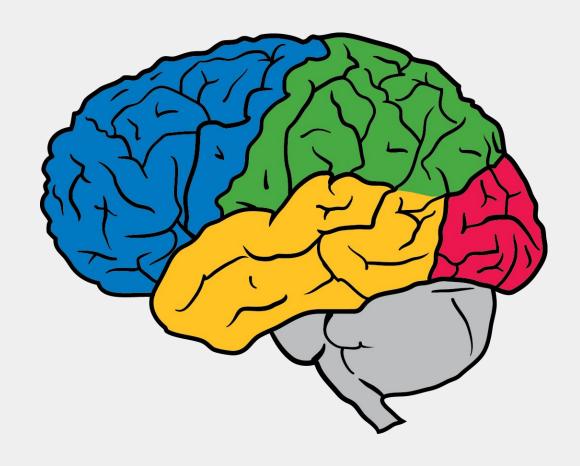
Try out Clockwork VAEs!

Useful resources are available on the project website:

danijar.com/cwvae

- this paper.
- Video prediction examples on Minecraft, KTH Action, GQN Mazes, and Moving MNIST.
- TensorFlow and a JAX reimplementation.
- Video generation code for level-resetting.

Website with videos, data, and code: danijar.com/cwvae



We can factorize the joint distribution of a sequence of images and latents at every level into two terms: (1) the reconstruction terms of the images given their lowest level latents, and (2) state transitions at all levels conditioned on the previous latent and the

$p(x_{1:T}, s_{1:T}^{1:L}) \doteq \left(\prod_{t=1}^{T} p(x_t \mid s_t^1)\right) \left(\prod_{l=1}^{L} \prod_{t \in \mathcal{T}_l} p(s_t^l \mid s_{t-1}^l, s_t^{l+1})\right)$

To implement this distribution and its inference model, Clockwork VAE utilizes the following components:

 $e_t^l = e(x_{t:t+k^{l-1}-1})$ Encoder: Posterior transition q_t^l : $q(s_t^l \mid s_{t-1}^l, s_t^{l+1}, e_t^l)$ $p(s_t^l \mid s_{t-1}^l, s_t^{l+1})$ Prior transition p_t^l : $p(x_t \mid s_t^1).$ Decoder:

Because we cannot compute the likelihood of the training data under the model in closed form, we use the ELBO as our training objective. This training objective optimizes a reconstruction loss at the lowest level, and a KL regularizer at every level in the hierarchy summed across all timesteps.

 $\max_{e,q,p} \sum_{t=1}^{T} \mathrm{E}_{q_t^1} [\ln p(x_t \mid s_t^1)] - \sum_{l=1}^{L} \sum_{t \in T_l} \mathrm{E}_{q_{t-1}^l q_t^{l+1}} \left[\mathrm{KL}[q_t^l \parallel p_t^l] \right]$

Minecraft dataset released with

Training and evaluation code in

