

MoCoGan: Decomposing Motion and Content for Video Generation

As the title says, the main point of this work is the separation of a video in two subspaces: Content and Motion. This paper is proposing a GAN based framework that generates **video frames** from random vectors.

Each vector consists of a content part and a motion part, and because the content doesn't change much in a short term video, that part remains the same for the whole generation (extracted from a Gaussian distribution), while the motion part is stochastic generated (using a RNN).

Motion and Content Decomposed GAN

To be able to generate videos of different lengths and videos with the same length but executed with different speeds, in this paper is used a latent space Z_i of images. A video of length K is represented by a path $[z_1, \dots, z_K]$ in this latent space. This space is further decomposed into the content and motion subspaces.

The content subspace is modelled using a Gaussian (normal) distribution, and the z_C variable is generated only once for the whole video.

$$\mathbf{z}_C \sim p_{Z_C} \equiv \mathcal{N}(\mathbf{z}|0, I_{d_C}) ,$$

The motion subspace is modelled by a path in the Z_m subspace. Since not all the paths in this subspace correspond to plausible motion, a recurrent neural network is used to generate valid paths. At each time step, this network takes a vector sampled from a Gaussian distribution and outputs a vector in the Z_m subspace. The RNN consists of one GRU layer.

Networks

This work uses 4 neural networks: the recurrent network that generates **motion paths**, the **image generator** G_i , the **image discriminator** D_i and the **video discriminator** D_v .

G_i generates a video clip by sequentially mapping vectors from Z_i to images, from a sequence of vectors like :

$$\left[\begin{bmatrix} \tilde{\mathbf{z}}_C \\ \mathbf{z}_M^{(1)} \end{bmatrix}, \dots, \begin{bmatrix} \mathbf{z}_C \\ \mathbf{z}_M^{(K)} \end{bmatrix} \right]$$

to a sequence of images :

$$\tilde{\mathbf{v}} = [\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(K)}],$$

where

z_M is from the recurrent network. Both D_i and D_v play the judge role, providing criticisms to G_i and R_M . D_i is specialised in criticising G_i based on individual images, while D_v provides criticism based on the generated video clip. D_v is also based on a spatio-temporal CNN architecture and evaluates the generated motion. Because G_i has no concept of motion, the criticism from D_v goes to the recurrent network.

The authors note that D_v should be enough because it provides feedback on both static image appearance and video dynamics but D_i improves the convergence of the adversarial training.

Objective function

The learning problem of MoCoGAN is:

$$\max_{G_I, R_M} \min_{D_I, D_V} \mathcal{F}_V(D_I, D_V, G_I, R_M)$$

The objective function is:

$$\mathbb{E}_{\mathbf{v}}[-\log D_I(S_1(\mathbf{v}))] + \mathbb{E}_{\tilde{\mathbf{v}}}[-\log(1 - D_I(S_1(\tilde{\mathbf{v}})))] + \\ \mathbb{E}_{\mathbf{v}}[-\log D_V(S_T(\mathbf{v}))] + \mathbb{E}_{\tilde{\mathbf{v}}}[-\log(1 - D_V(S_T(\tilde{\mathbf{v}})))] ,$$

In this formula, the first and second terms encourage D_I to output 1 for a video frame from a real video and 0 for a generated one. Similarly, the third and fourth terms encourage D_V to output 1 or 0 depending if the frame sequence is real or generated.

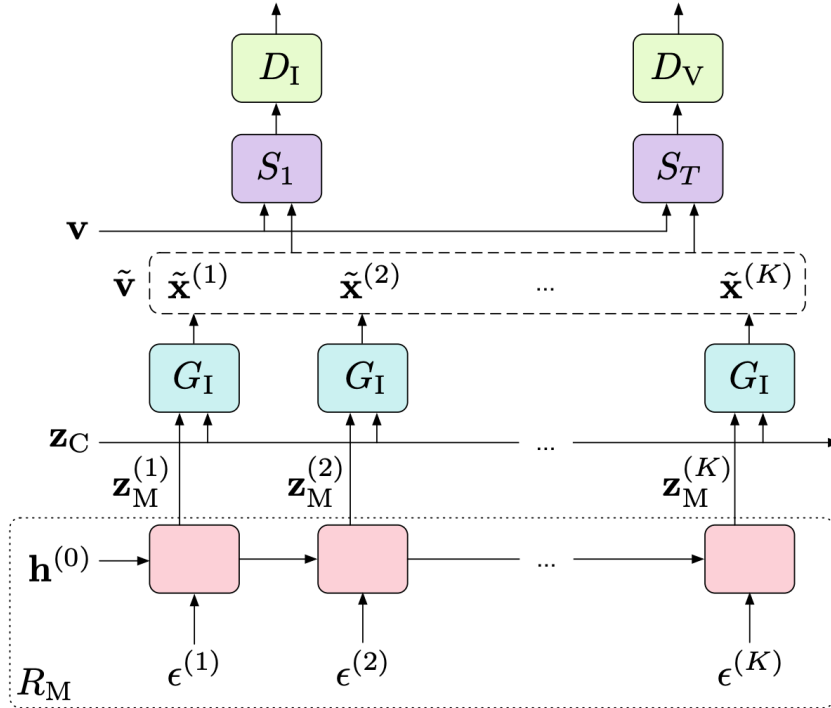


Figure 2: The MoCoGAN framework for video generation. For a video, the content vector, \mathbf{z}_C , is sampled once and fixed. Then, a series of random variables $[\epsilon^{(1)}, \dots, \epsilon^{(K)}]$ is sampled and mapped to a series of motion codes $[\mathbf{z}_M^{(1)}, \dots, \mathbf{z}_M^{(K)}]$ via the recurrent neural network R_M . A generator G_I produces a frame, $\tilde{\mathbf{x}}^{(k)}$, using the content and the motion vectors $\{\mathbf{z}_C, \mathbf{z}_M^{(k)}\}$. The discriminators, D_I and D_V , are trained on real and fake images and videos, respectively, sampled from the training set \mathbf{v} and the generated set $\tilde{\mathbf{v}}$. The function S_1 samples a single frame from a video, S_T samples T consecutive frames.