



國立臺灣大學
National Taiwan University



Dual-MTGAN: Stochastic and Deterministic Motion Transfer for Image-to-Video Synthesis

ICPR 2020

^{1,2}Fu-En Yang*, ¹Jing-Cheng Chang*, ¹Yuan-Hao Lee, ^{1,2}Yu-Chiang Frank Wang

¹Graduate Institute of Communication Engineering, National Taiwan University

²ASUS Intelligent Cloud Services (AICS)

What is image-to-video synthesis?

Synthesize videos from an input image with the motion of interest.

Deterministic Motion Transfer

- Transfer motion pattern across different videos

Source
Video

Synthesized Videos with Different Identities
(Preserved Facial Expressions/Motions)

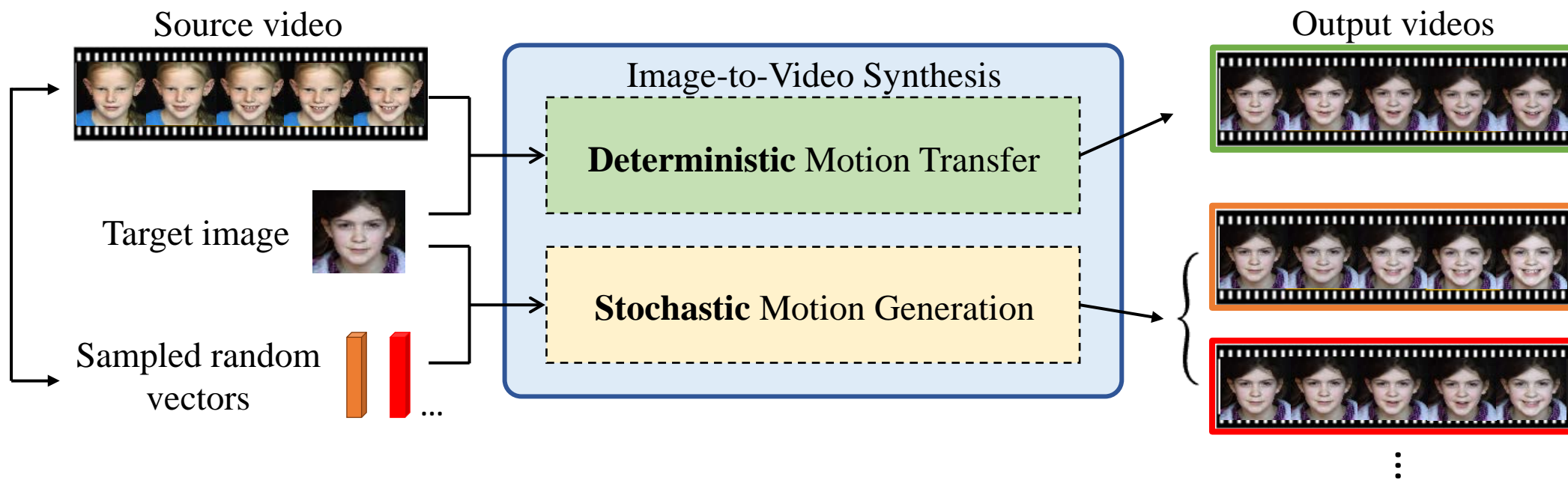


Stochastic Motion Generation

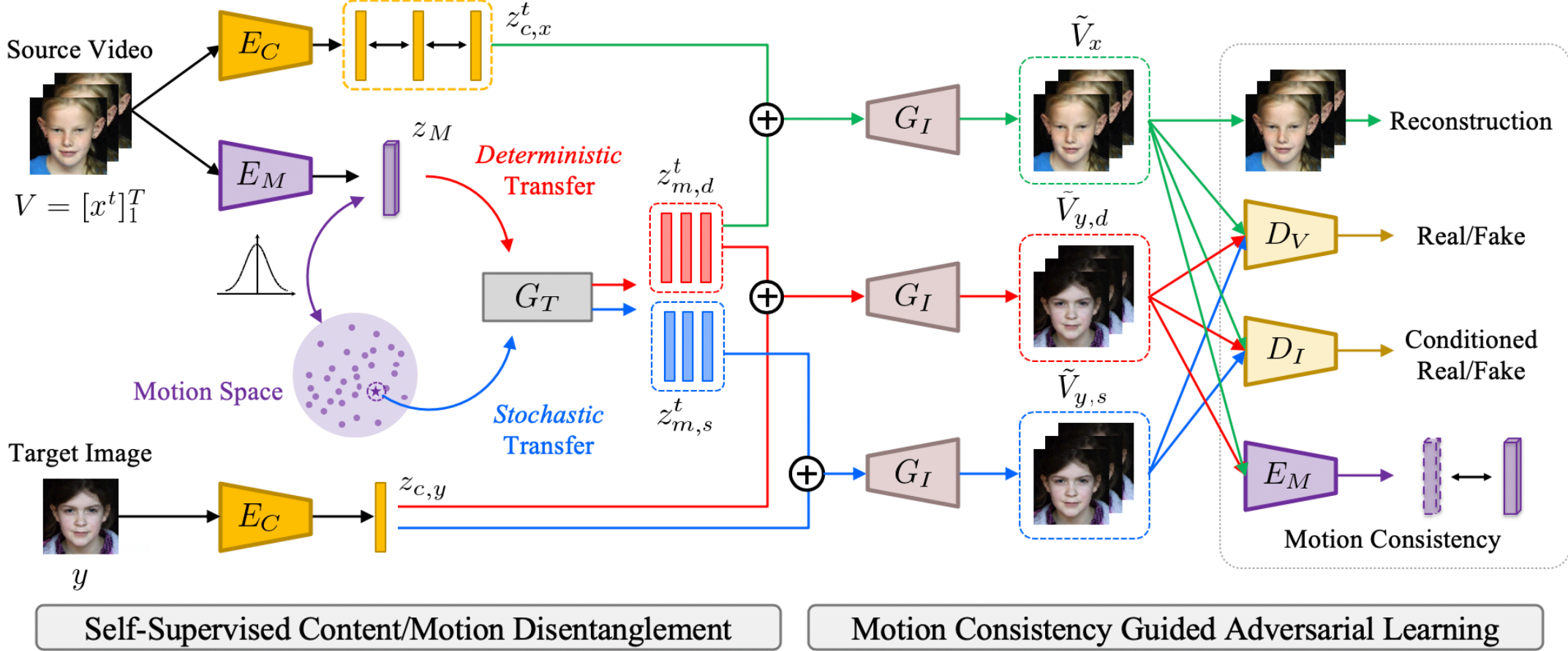
- Generate videos from an image or few frames



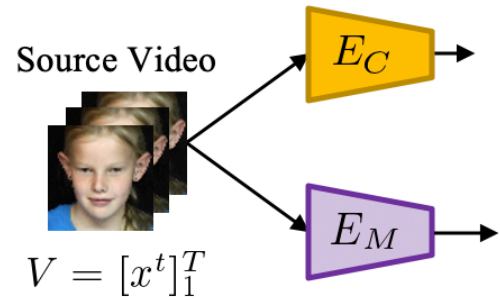
Can we jointly perform *deterministic* and *stochastic* motion transfer in a *unified* framework?



Method – Dual-Motion Transfer GAN (Dual-MTGAN)

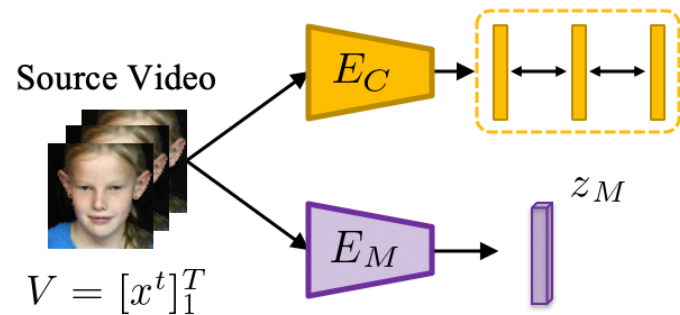


Method



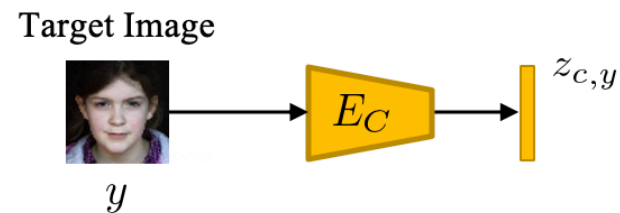
Method

– Self-Supervised Content/Motion Disentanglement



Temporal coherence across frames:

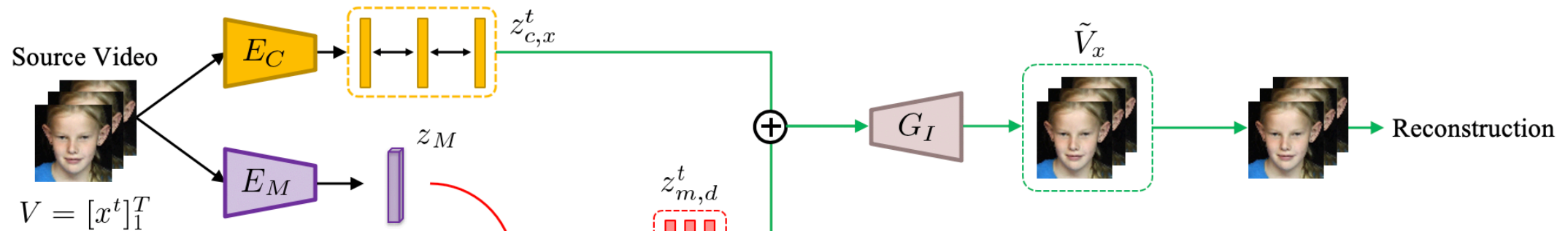
$$\mathcal{L}_C = \|E_C(x^t) - E_C(x^{t+1})\|_1$$



Self-Supervised Content/Motion Disentanglement

Method

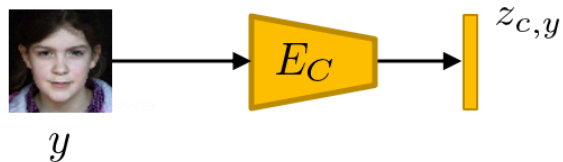
– Source Video Reconstruction



Reconstruction Loss:

$$\mathcal{L}_{rec} = \|\tilde{V}_x - V\|_1$$

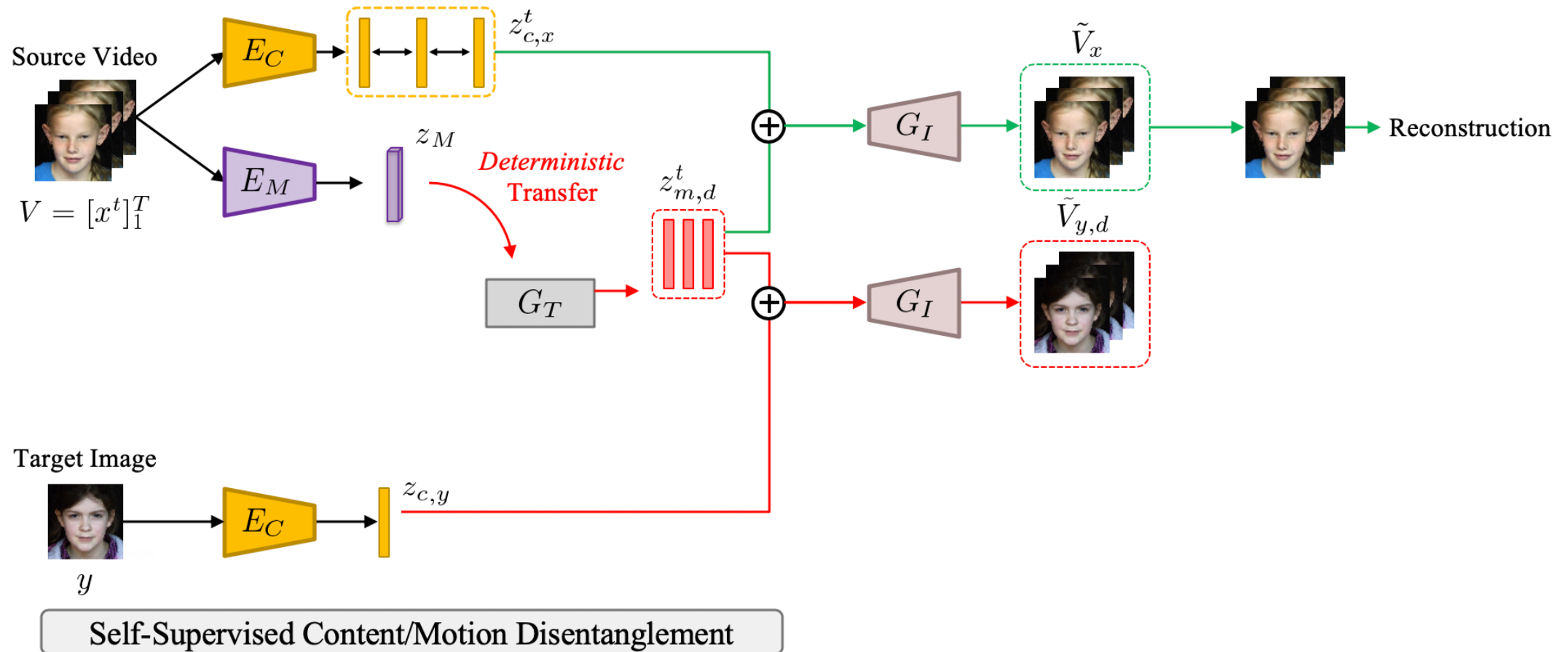
Target Image



Self-Supervised Content/Motion Disentanglement

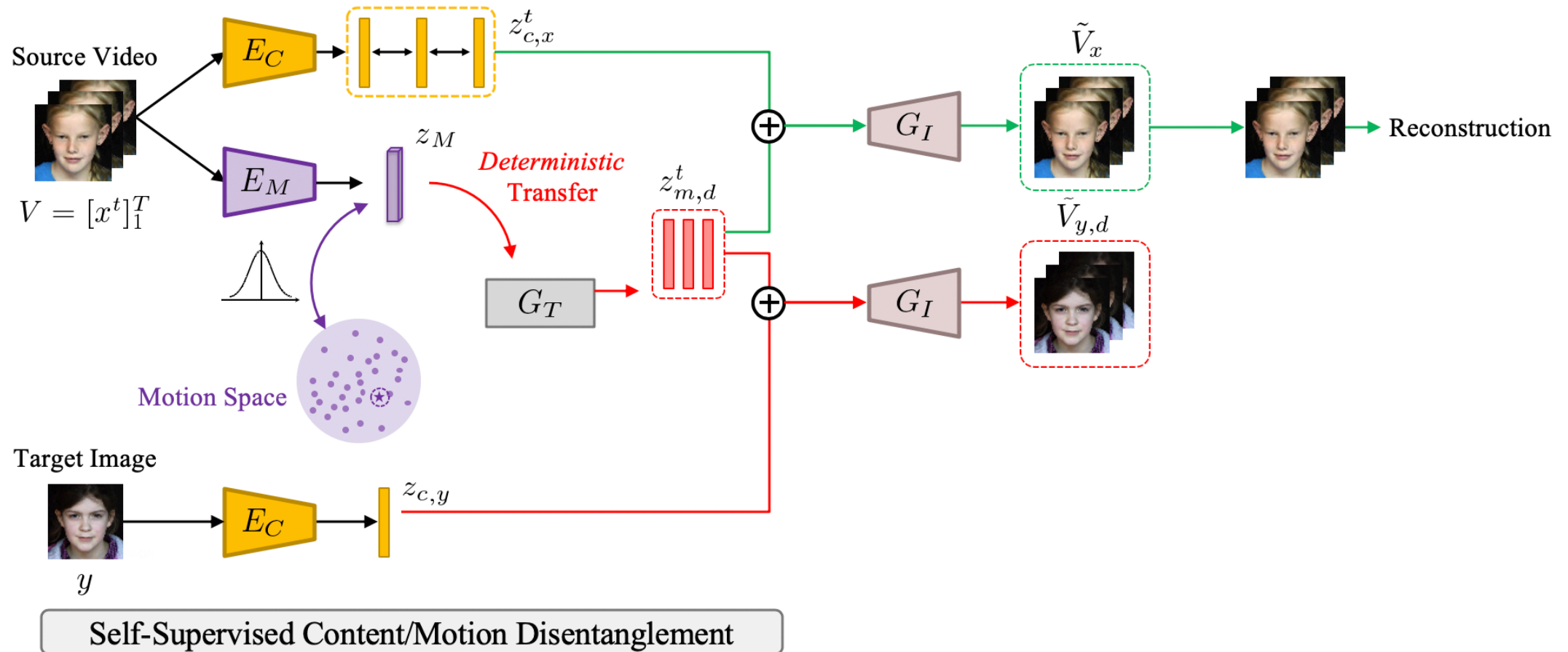
Method

– Deterministic Motion Transfer



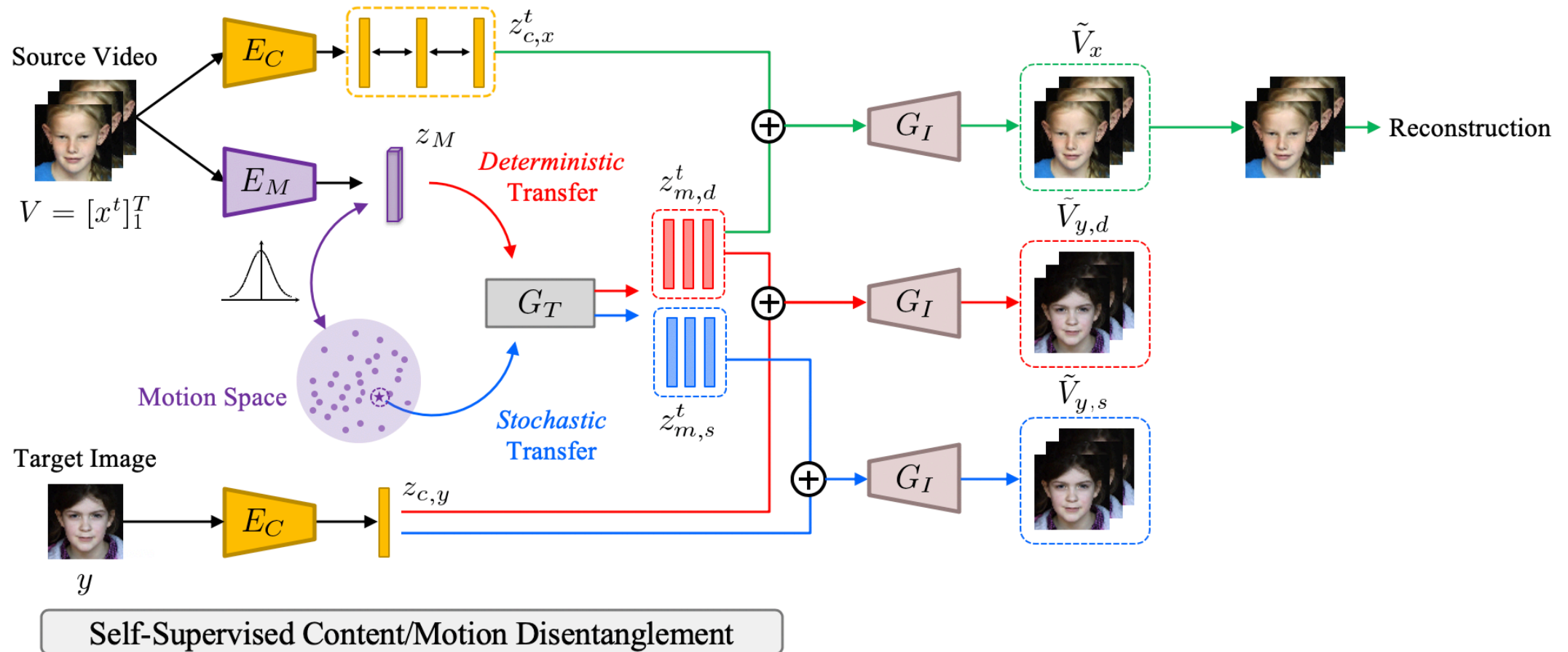
Method

– Learning Motion Latent Space



Method

– Stochastic Motion Transfer



Method

– Adversarial Learning

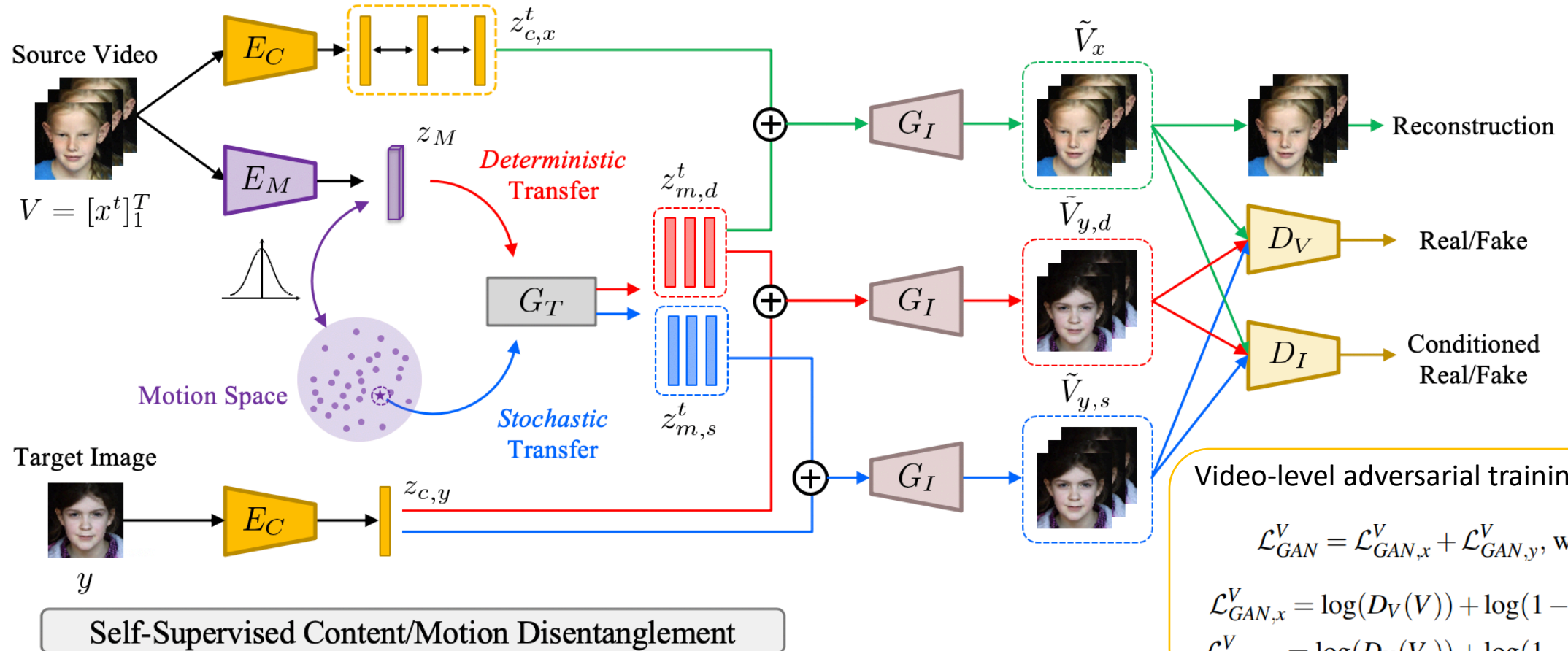


Image-level conditional adversarial training:

$$\mathcal{L}_{GAN}^I = \mathcal{L}_{GAN,x}^I + \mathcal{L}_{GAN,y}^I, \text{ where}$$

$$\mathcal{L}_{GAN,x}^I = \log(D_I(x^1, S_I(V))) + \frac{1}{2}[\log(1 - D_I(x^1, S_I(\tilde{V}_x))) + \log(1 - D_I(y, S_I(V)))]$$

$$\mathcal{L}_{GAN,y}^I = \log(D_I(y, y)) + \frac{1}{2}[\log(1 - D_I(y, S_I(\tilde{V}_y))) + \log(1 - D_I(x^1, S_I(V_y)))],$$

Video-level adversarial training:

$$\mathcal{L}_{GAN}^V = \mathcal{L}_{GAN,x}^V + \mathcal{L}_{GAN,y}^V, \text{ where}$$

$$\mathcal{L}_{GAN,x}^V = \log(D_V(V)) + \log(1 - D_V(\tilde{V}_x)),$$

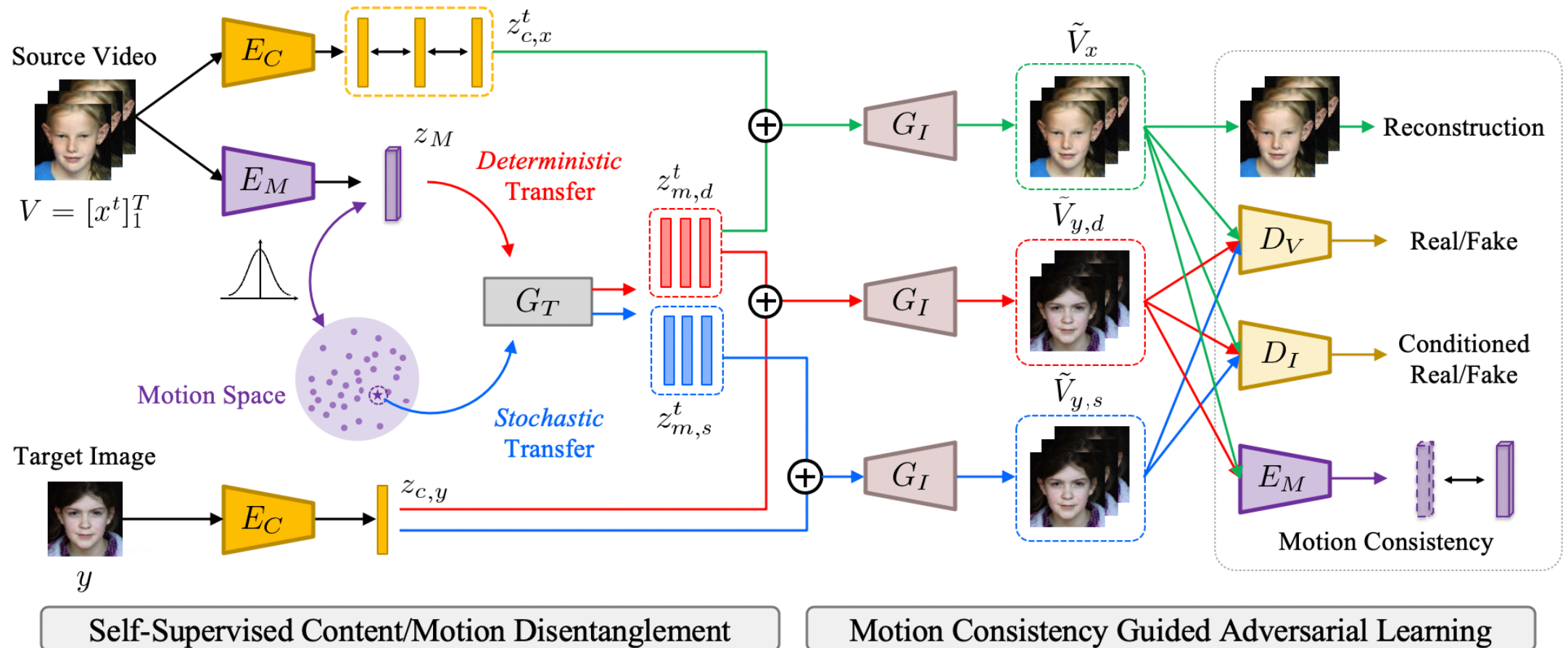
$$\mathcal{L}_{GAN,y}^V = \log(D_V(V_y)) + \log(1 - D_V(\tilde{V}_y)).$$

Method

– Motion Consistency

Motion consistency:

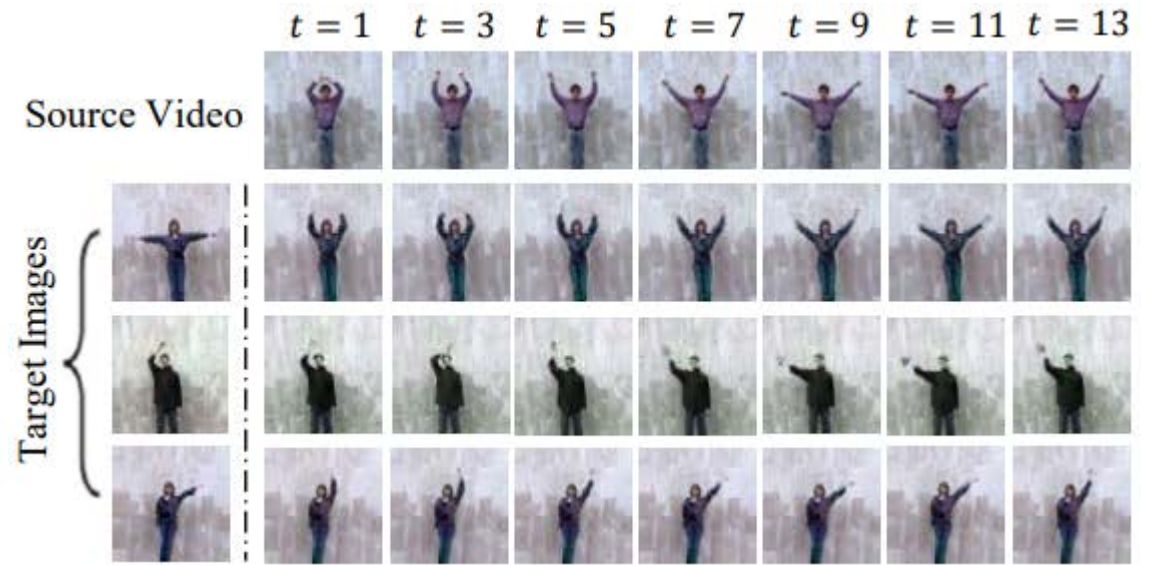
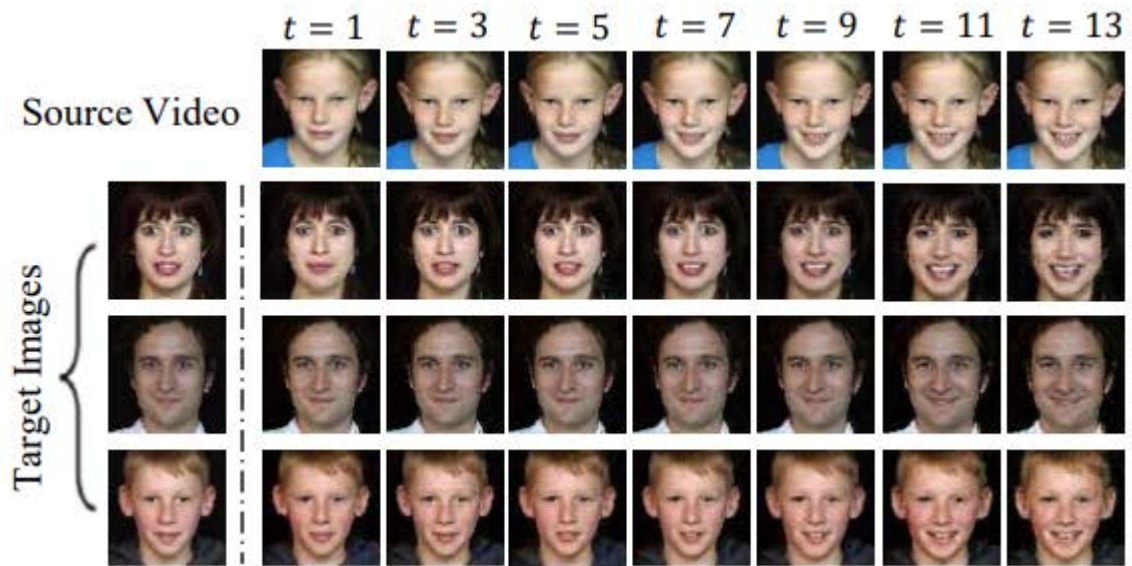
$$\mathcal{L}_M = \|E_M(\tilde{V}_x) - z_M\|_1 + \|E_M(\tilde{V}_{y,d}) - z_M\|_1$$



Result

– Deterministic Face Reenactment and Motion Retargeting

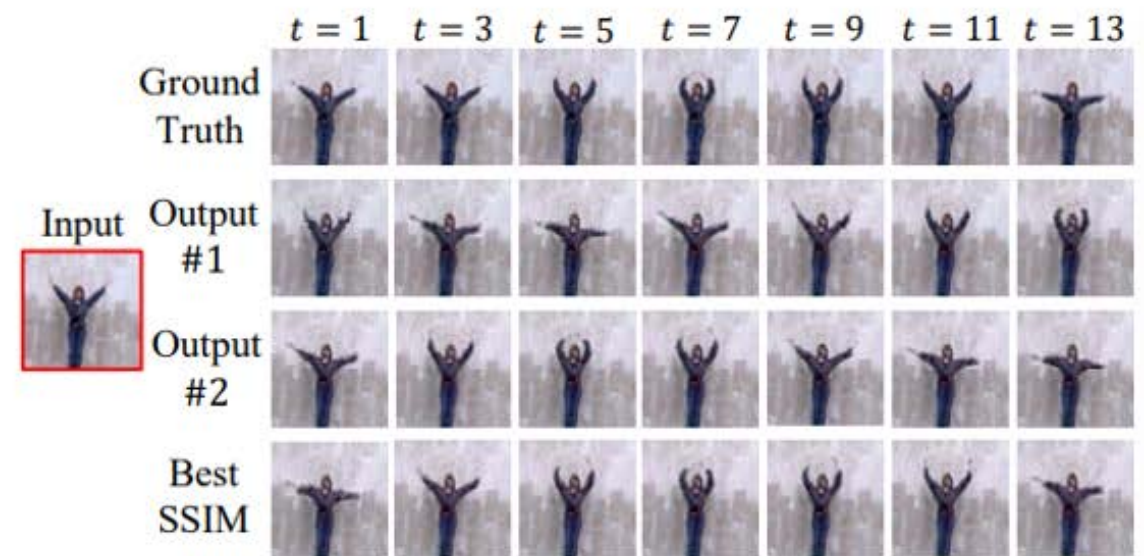
- Facial expression & Human actions



Result

– Stochastic Robot Movement and Action Generation

- Robot pushing & Human actions



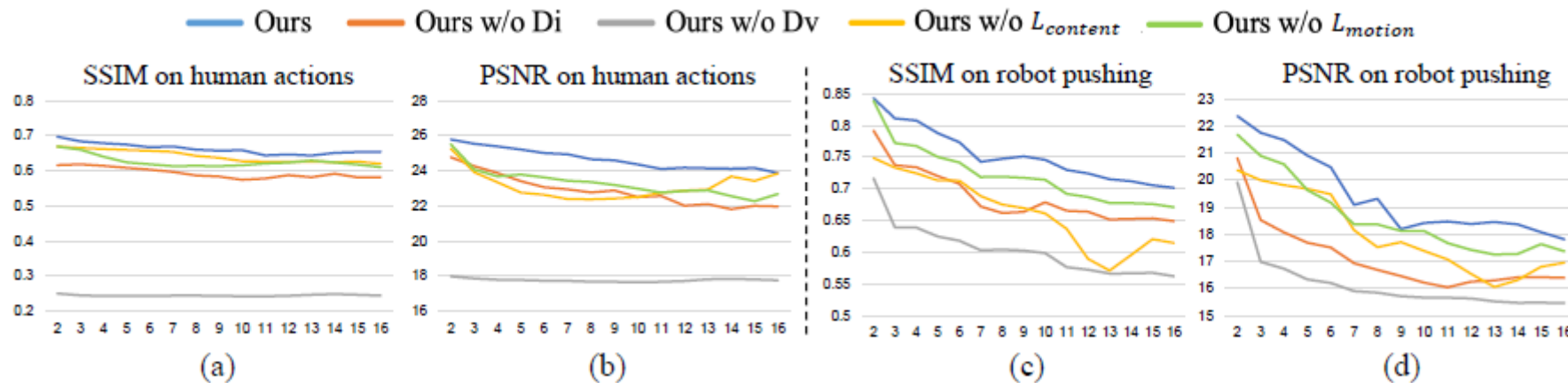
Result

– Quantitative Comparison & Ablation Study

- Quantitative Comparison

Method	<i>Robot pushing</i>	
	SSIM (\uparrow)	LPIPS (\uparrow)
SVG	0.815 ± 0.006	0.0398 ± 0.0005
Monkey-Net	0.783 ± 0.008	N/A
Ours	0.827 ± 0.007	0.0422 ± 0.0003

- Ablation Study



Conclusion

- Given an input image, our proposed model allows transfer of motion patterns from video data, or synthesis of video sequences with motion diversity.
- By enforcing appearance coherence and motion consistency, our Dual-MTGAN factorizes visual latent representations into disjoint features describing content and motion information in a self-supervised manner.

Thanks for listening!