



# Generative Adversarial Graph Convolutional Networks for Human Action Synthesis

Bruno Degardin<sup>1,3,4</sup> João Neves<sup>2,3</sup> Vasco Lopes<sup>2,3,4</sup>  
<sup>1</sup>IT - Instituto de Telecomunicações <sup>2</sup>NOVA LINCS

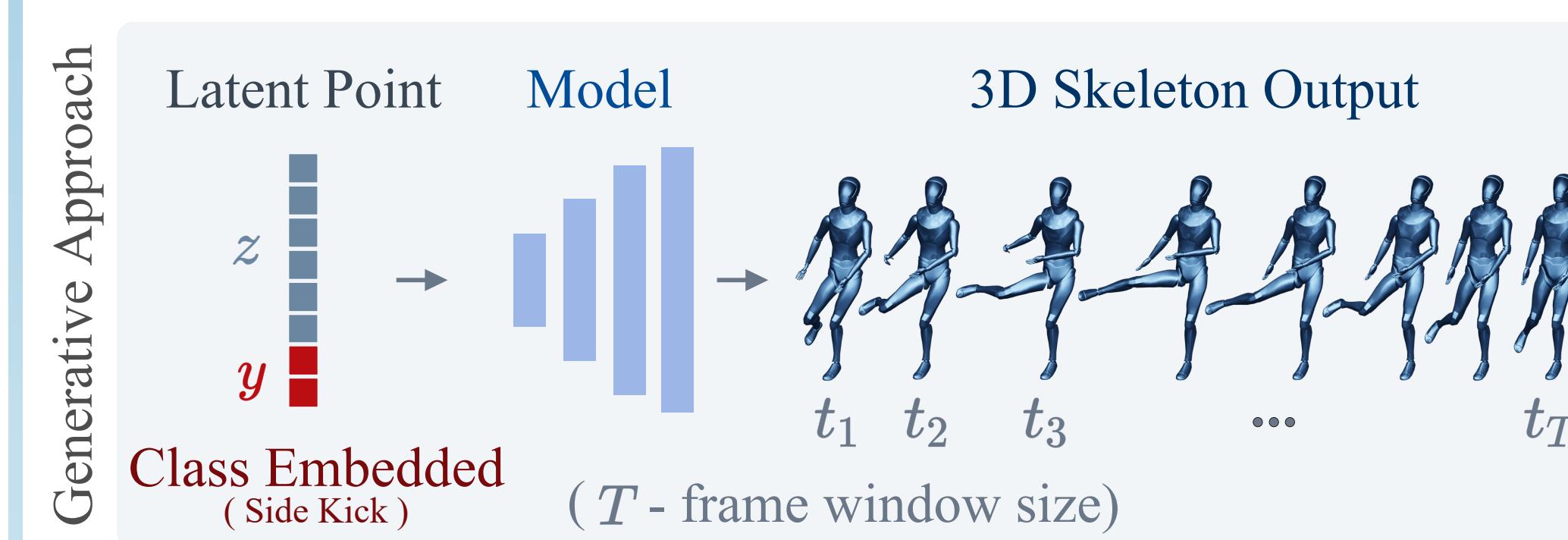
João Brito<sup>4</sup> Ehsan Yaghoubi<sup>3</sup> Hugo Proença<sup>1,3</sup>  
<sup>3</sup>Universidade da Beira Interior <sup>4</sup>DeepNeuronic

Paper, Code & Pretrained Models:  
<https://github.com/DegardinBruno/Kinetic-GAN>



## Problem Definition and Contribution

**Goal:** Synthesising human actions with realistic body movements (unconstrained environment).



### Motivation:

- Extracting 3D skeleton pose becomes difficult due to RGB and depth noises.
- Current approaches are limited in conditioning desirable actions and generate global movement.

### Key Contributions:

- A new scalable Generative Adversarial Graph Convolutional Network to synthesise human actions.
- Architecture extension to a conditional model, latent space disentanglement and stochastic variation to generate desired actions (120) with increased diversity.
- Kinetic-GAN achieves state-of-the-art performance on NTU RGB+D, NTU-120 RGB+D and Human3.6M.

## Problem Formulation

**Main idea:** Kinetic-GAN employs graph convolutions at each graph's resolution level  $l$  to overcome structural information loss from conventional convolutional kernels.

### Spatiotemporal Graph Convolution:

- The graph convolution, in one frame, is computed as:

$$\mathcal{S}(\mathbf{X}_l) = \sum_{i=1}^p \Lambda_{l_i}^{-\frac{1}{2}} (\mathbf{A}_{l_i} \odot \mathbf{M}_l) \Lambda_{l_i}^{-\frac{1}{2}} \mathbf{X}_l \mathbf{W}_{l_i}, \quad (1)$$

where  $\mathbf{W}_{l_i}$  are the stacked weight vectors.

- One-dimensional kernels on the temporal axis:

$$\mathcal{T}(\mathcal{S}(\mathbf{X}_l)) = \mathcal{S}(\mathbf{X}_l) * \mathbf{w}_l \quad (2)$$

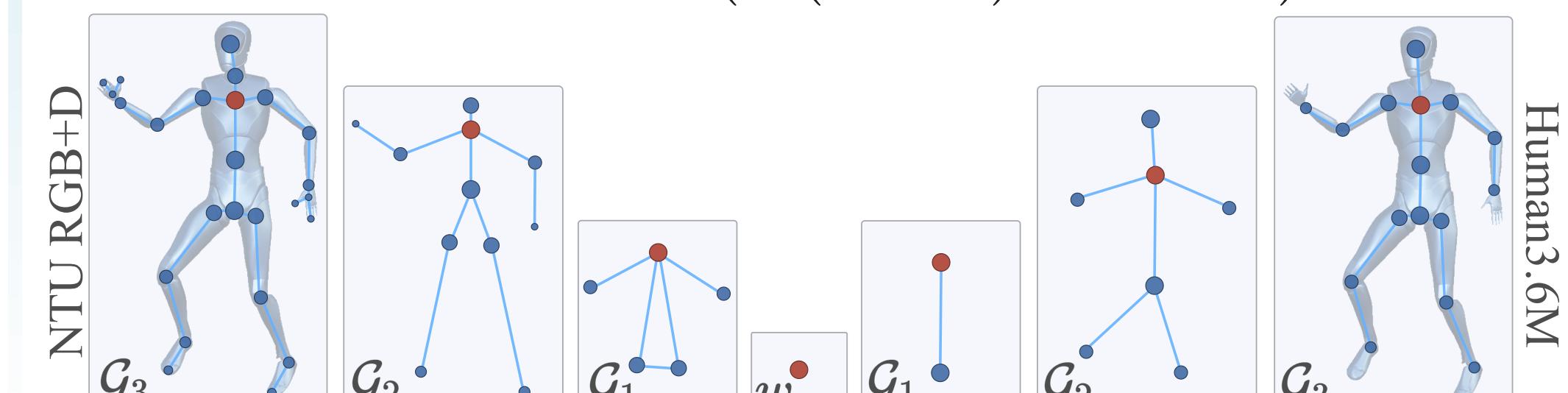
with temporal kernel  $\mathbf{w}_l \in \mathbb{R}^{1 \times t \times C}$  with  $t$  frames.

### Graph Upsampling and Downsampling Paths:

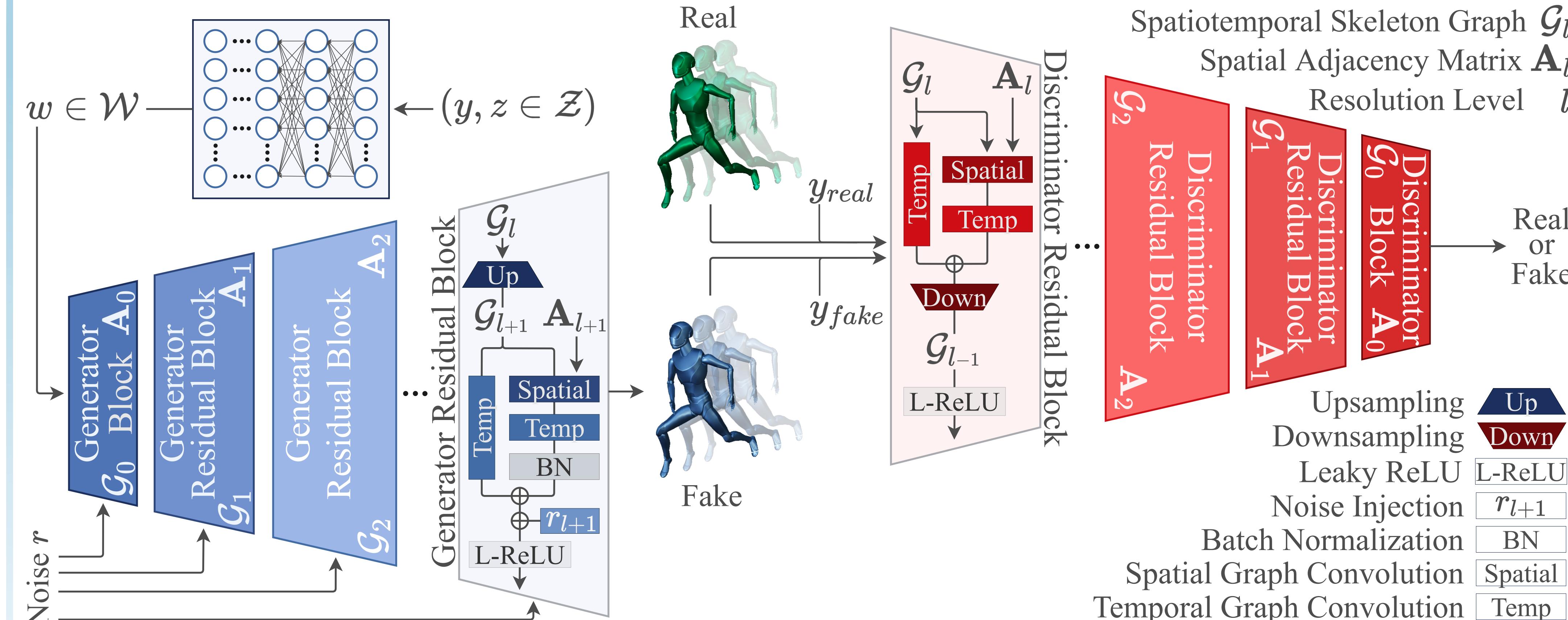
- Graph upsampling: Introducing new vertices and one-dimensional interpolation over the temporal axis.
- Temporal skip connections bring training stability.

$$\mathbf{X}_{l+1} = \mathcal{T}(\mathcal{S}(Up(\mathbf{X}_l))) + \mathcal{T}(Up(\mathbf{X}_l)) \quad (3)$$

$$\mathbf{X}_{l-1} = Down(\mathcal{T}(\mathcal{S}(\mathbf{X}_l)) + \mathcal{T}(\mathbf{X}_l))$$



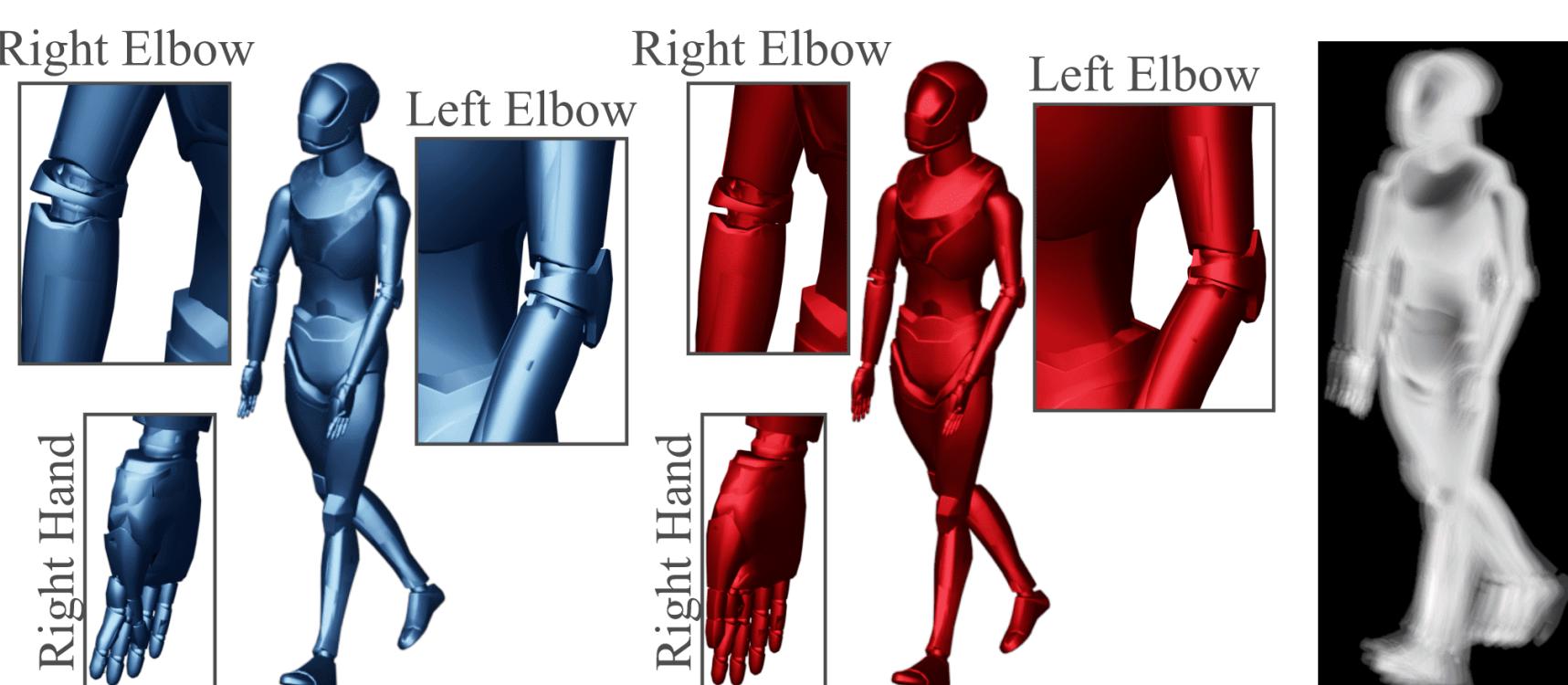
## Method



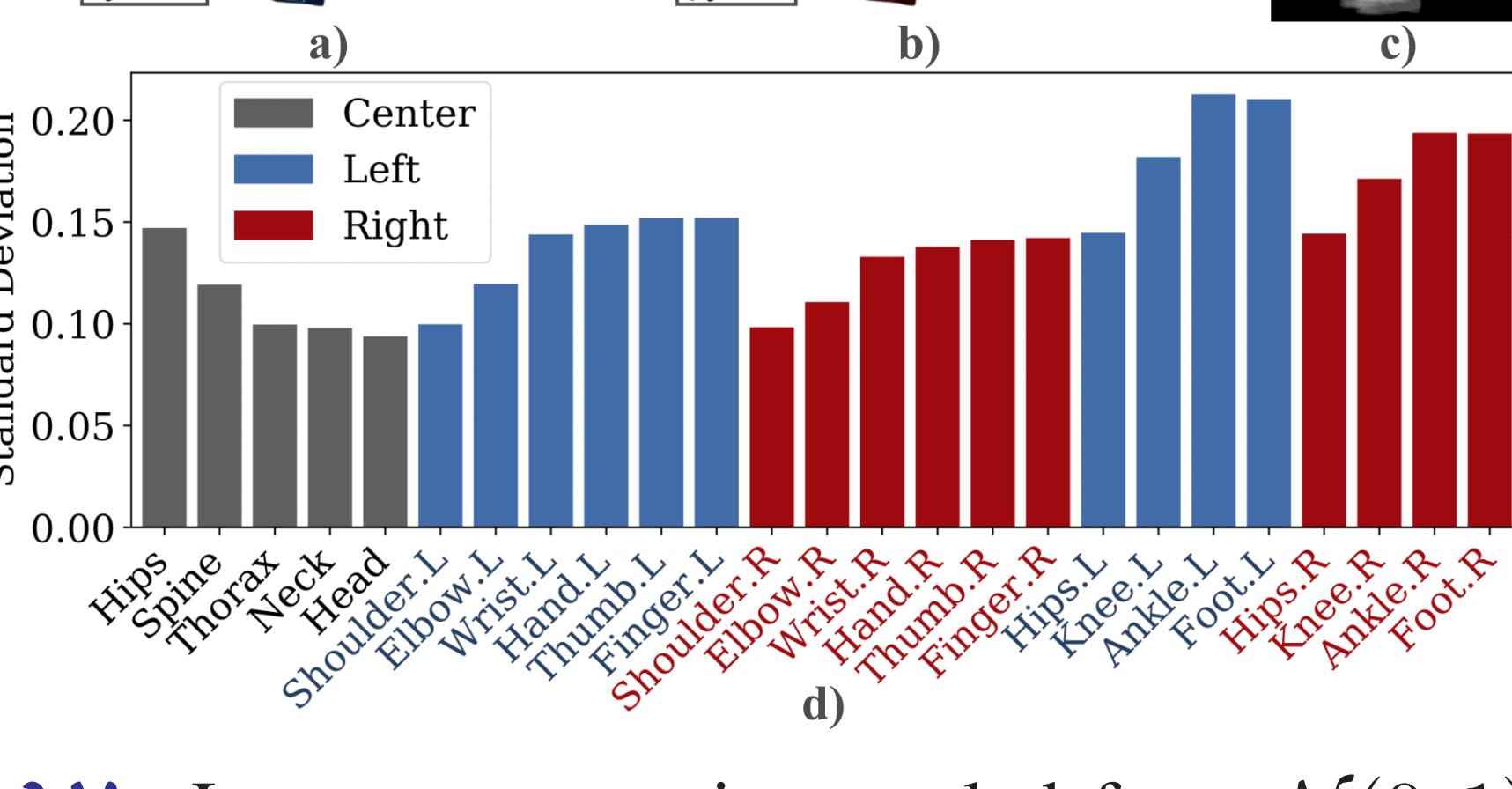
- Embedded class representation  $y$  is concatenated to  $z$  in  $\mathcal{G}$  (channel-wise) in the discriminator.
- Noise injector provides a second input to the generator to produce variations between samples (same latent point).
- Non-linear mapping network to produce an intermediate latent space  $\mathcal{W}$  to allow less entangled latent factors.
- Wasserstein distance with gradient penalty loss is used to drive the learning process of Kinetic-GAN.

## Improving Quality and Diversity

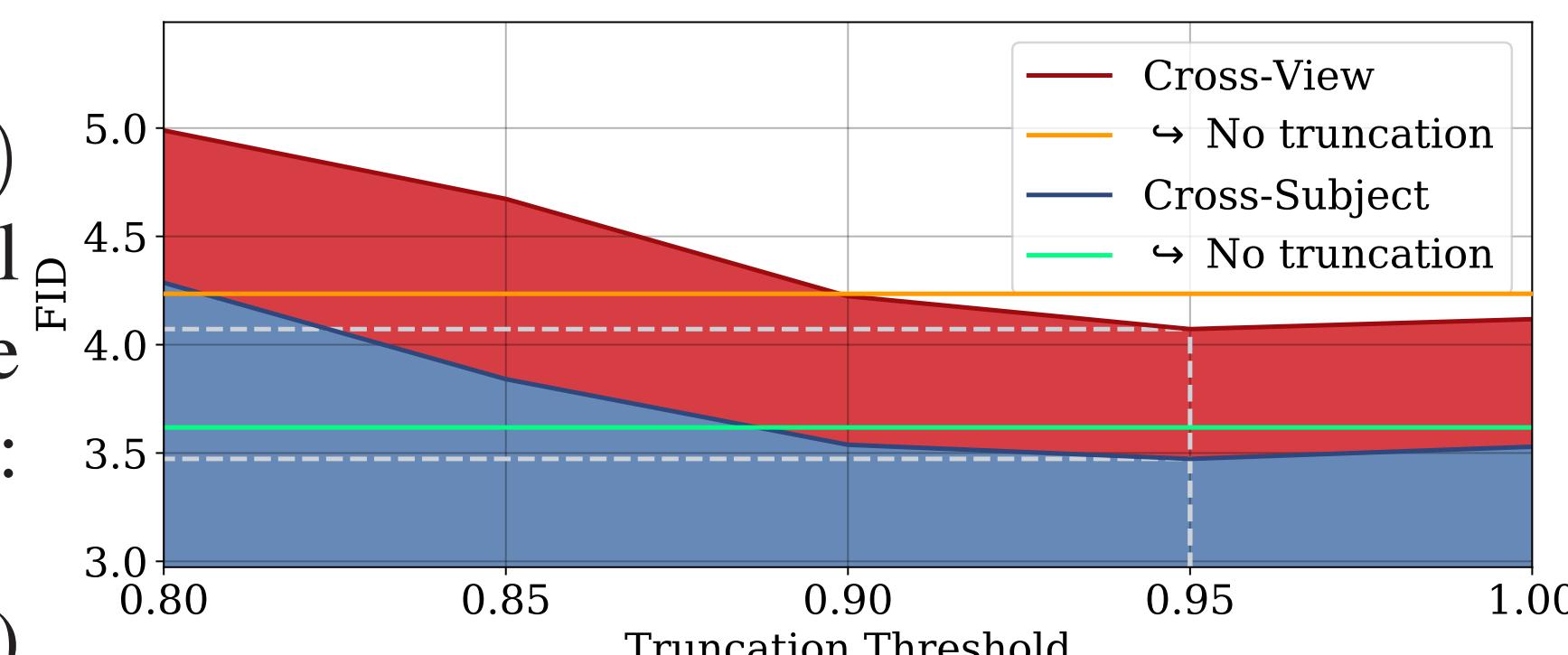
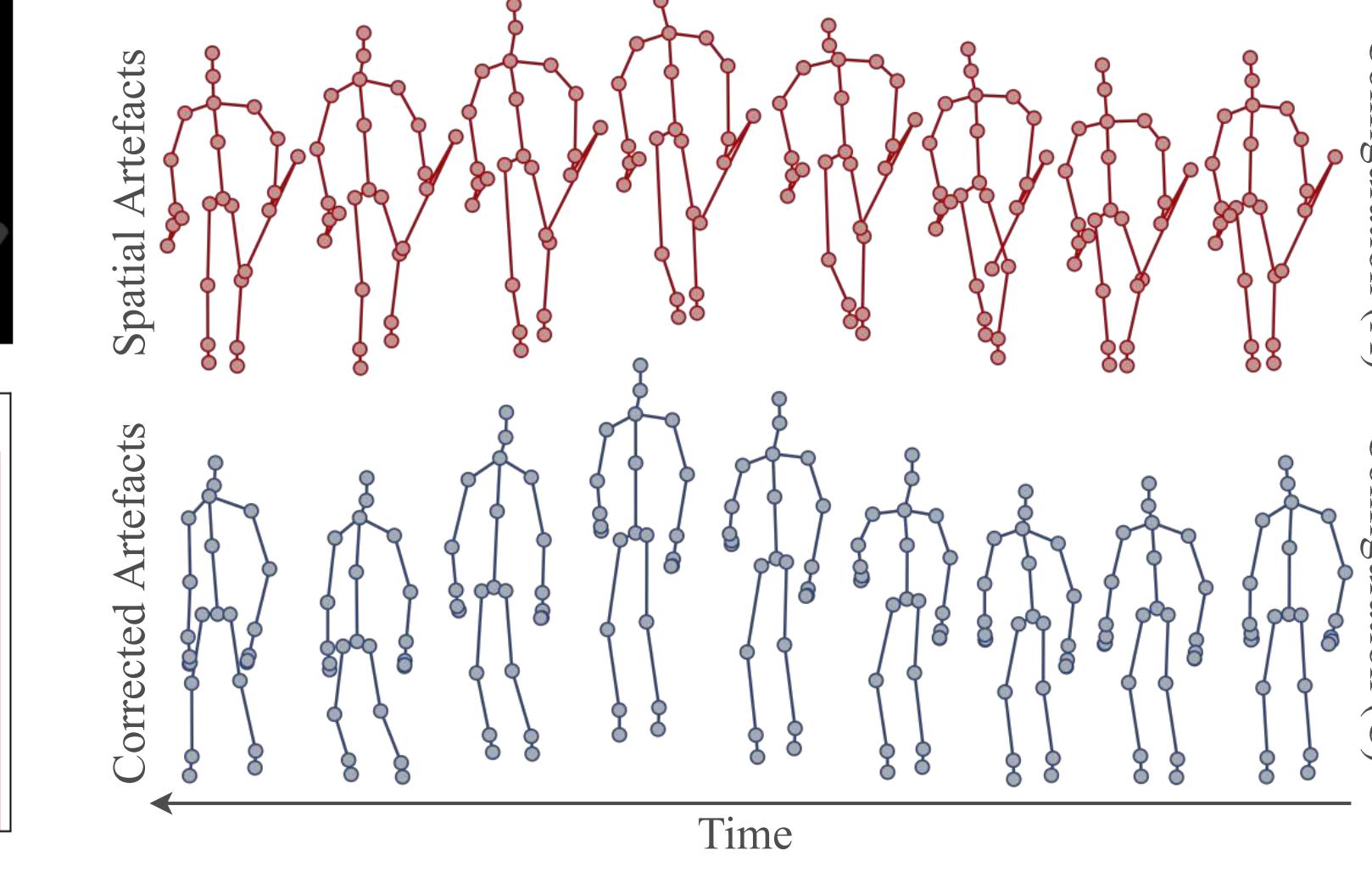
**Stochastic Variation:** Current approaches are attached to autoregressive techniques and Gaussian process reducing the ability of variation.



- Random noise is added joint-wise after each generator's graph convolution.



**Reducing Spatial Artefacts:** Batch normalization omits information concerning individual feature's magnitude (magnified by the generator). Hence, we regularize the use of batch normalization on the generator.



**Truncation Trick on  $\mathcal{W}$ :** Latent space  $z$  is sampled from  $\mathcal{N}(0, 1)$  and mapped to  $\mathcal{W}$ . Ranges of low density (training data) are not well represented, becoming difficult to learn for the generator. Despite some variation losses, during inference time, we scale the deviation of  $w$  as:

$$\mathbf{w}' = \mathbb{E}_{z \sim \mathbb{P}_z} [f(z)] + \psi(\mathbf{w} - \mathbb{E}_{z \sim \mathbb{P}_z} [f(z)]), \quad (4)$$

## Quantitative and Qualitative Results

Method	FID	MMD <sub>a</sub>	MMD <sub>s</sub>	FID	MMD <sub>a</sub>	MMD <sub>s</sub>
<i>Cross-Subject</i>						
c-GAN	27.480	0.919	0.975	31.875	0.993	1.088
CSGN	6.030	0.873	0.954	7.114	0.910	0.991
A Proposed	5.621	0.836	0.927	6.528	0.883	0.953
B + No Residual	19.723	0.892	0.961	21.331	0.971	1.030
C + Regular BN	4.751	0.815	0.917	5.328	0.867	0.940
D + Noise Inject	4.698	0.811	0.895	5.102	0.851	0.933
E + Mapping Net	<b>3.618</b>	<b>0.772</b>	<b>0.871</b>	<b>4.235</b>	<b>0.824</b>	<b>0.913</b>
<i>Cross-View</i>						

Method	FID	MMD <sub>a</sub>	MMD <sub>s</sub>	FID	MMD <sub>a</sub>	MMD <sub>s</sub>
<i>NTU RGB+D</i>						
c-GAN	54.403	1.037	1.104	58.531	1.082	1.141
<b>Kinetic-GAN</b>	<b>5.967</b>	<b>0.819</b>	<b>0.906</b>	<b>6.751</b>	<b>0.847</b>	<b>0.934</b>

Method	MMD <sub>a</sub>	MMD <sub>s</sub>	MMD <sub>a</sub>	MMD <sub>s</sub>
E2E	0.991	0.805		
EPVA	0.996	0.806		
adv-EPVA	0.977	0.792		
SkeletonVAE	0.452	0.467		
SkeletonGAN	0.419	0.419		
c-SkeletonGAN	0.195	0.218		
c-GAN	0.161	0.187		
SA-GCN	0.146	0.134		
<b>Kinetic-GAN</b>	<b>0.071</b>	<b>0.082</b>		

Method	Cross-Subject	Cross-View
SkeletonVAE	0.992	1.079
SkeletonGAN	0.698	0.999
c-SkeletonGAN	0.338	0.371
c-GAN	0.334	0.365
SA-GCN	0.285	0.316
<b>Kinetic-GAN</b>	<b>0.256</b>	<b>0.295</b>

