

Stage I: Pose Generation



Initial & Target
HOI Pose

Render



Depth Map

Seg Map

Hand Keypoints

Feature Extraction

Action Network



HOI Pose Sequence

Render



Depth Seq

Seg Seq

Hand Keypoints Seq

Stage II: Appearance Generation

Control Block I

Control Block II



First Frame

Enc

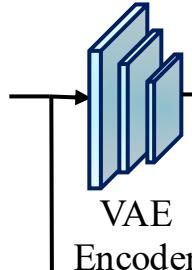
Noise

Base Block I

Base Block II

Dec

Stage III: Motion Generation

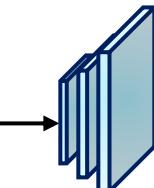


VAE Encoder

[First Frame,
Noise]

DiT Block

DiT Block



VAE Decoder



Generated Videos

Controllable Video
Generation