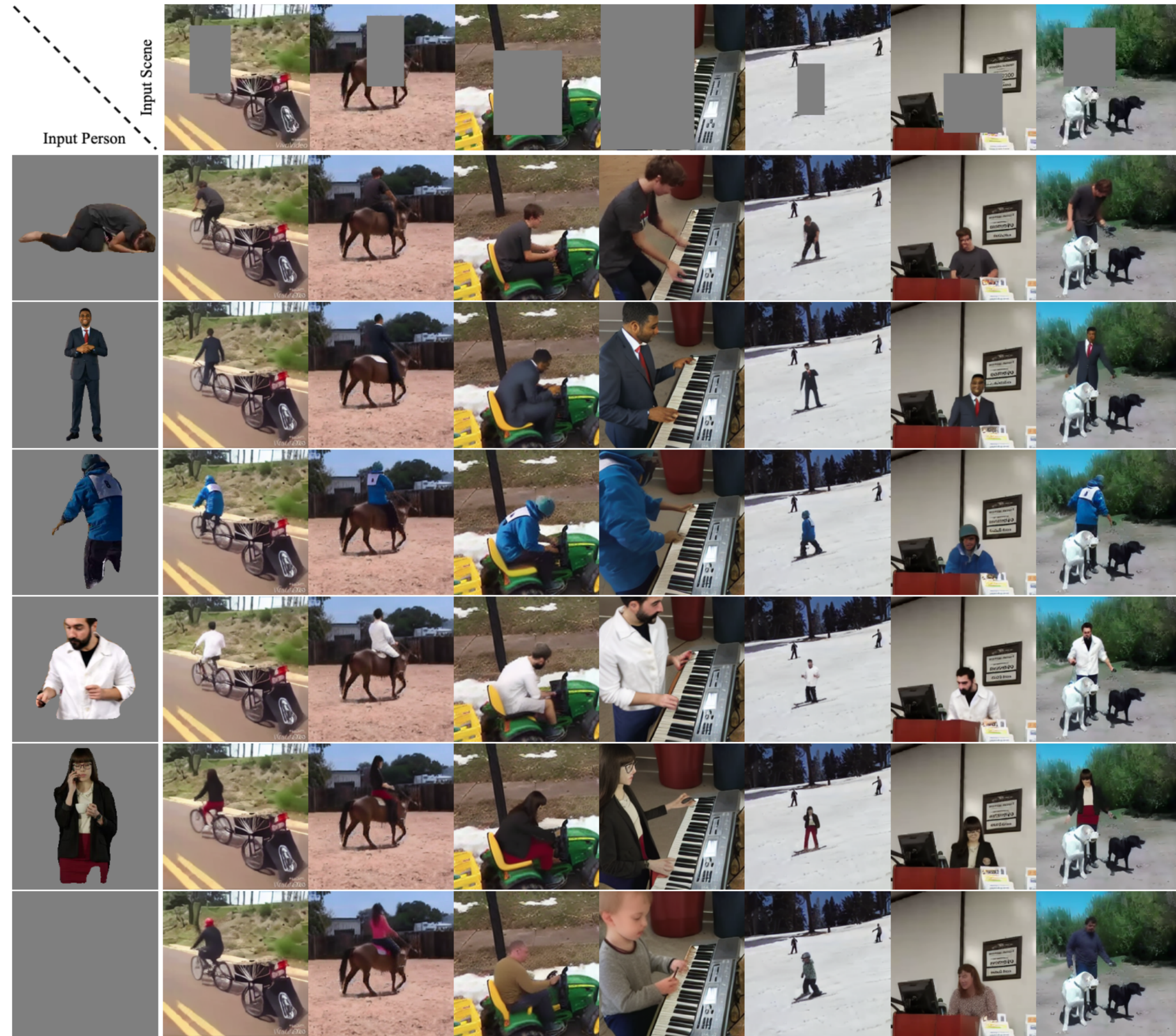


Putting People in Their Place: Affordance-Aware Human Insertion into Scenes

Sumith Kulal¹ Tim Brooks² Alex Aiken¹ Jiajun Wu¹
 Jimei Yang³ Jingwan Lu³ Alexei A. Efros² Krishna K. Singh³
¹Stanford University, ²UC Berkeley, ³Adobe Research

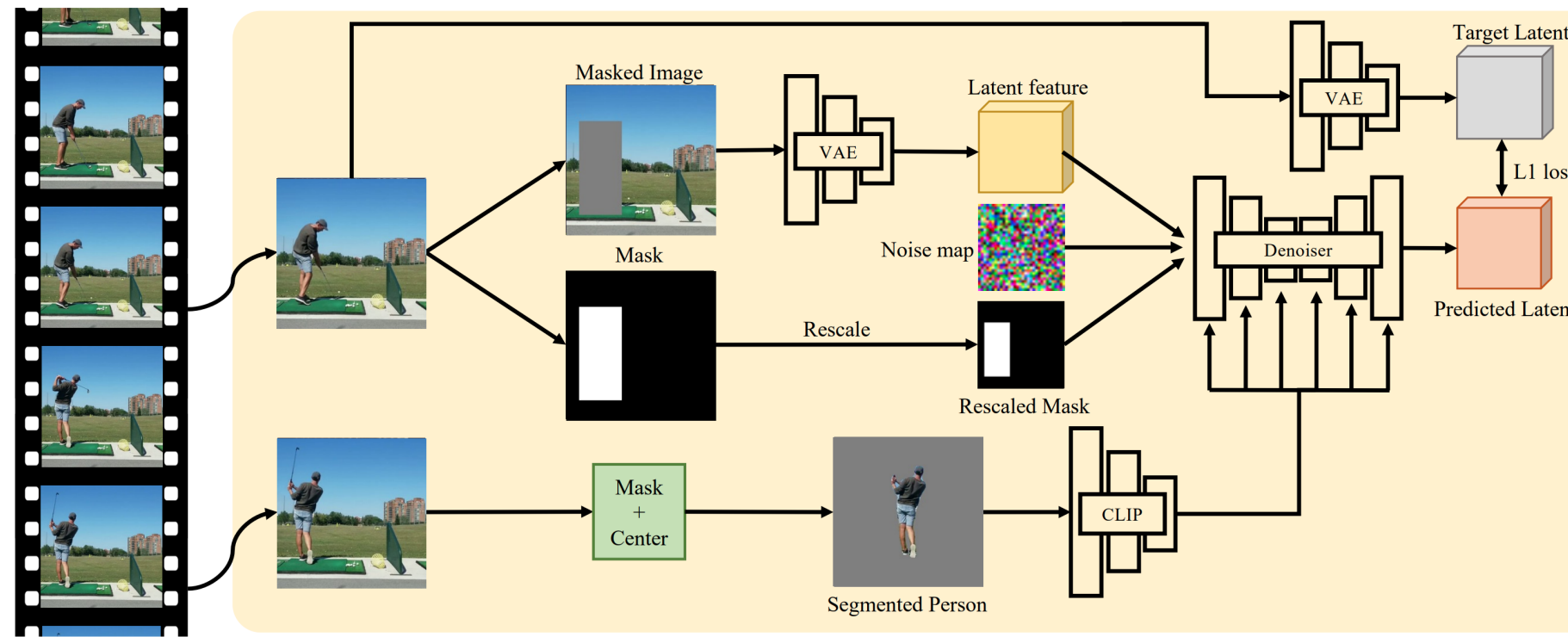


Photo-realistic Affordance-Aware Human Insertion into Scenes



Inputs: person image (left) and scene image with marked region (top)
Outputs: realistic insertion of the person into the scene image
Method: large-scale diffusion model trained in a self-supervised fashion on videos
Data: 2.4 million videos of humans moving around in scenes.
Highlights: self-supervised learning, affordances, image synthesis and editing

Learning Architecture Overview



Quantitative Results

	FID	PCKh
Image (w/o aug)	13.17	8.32
Image (w/ aug)	13.01	10.66
Video (w/o aug)	12.10	15.80
Video (w/ aug)	10.08	17.60

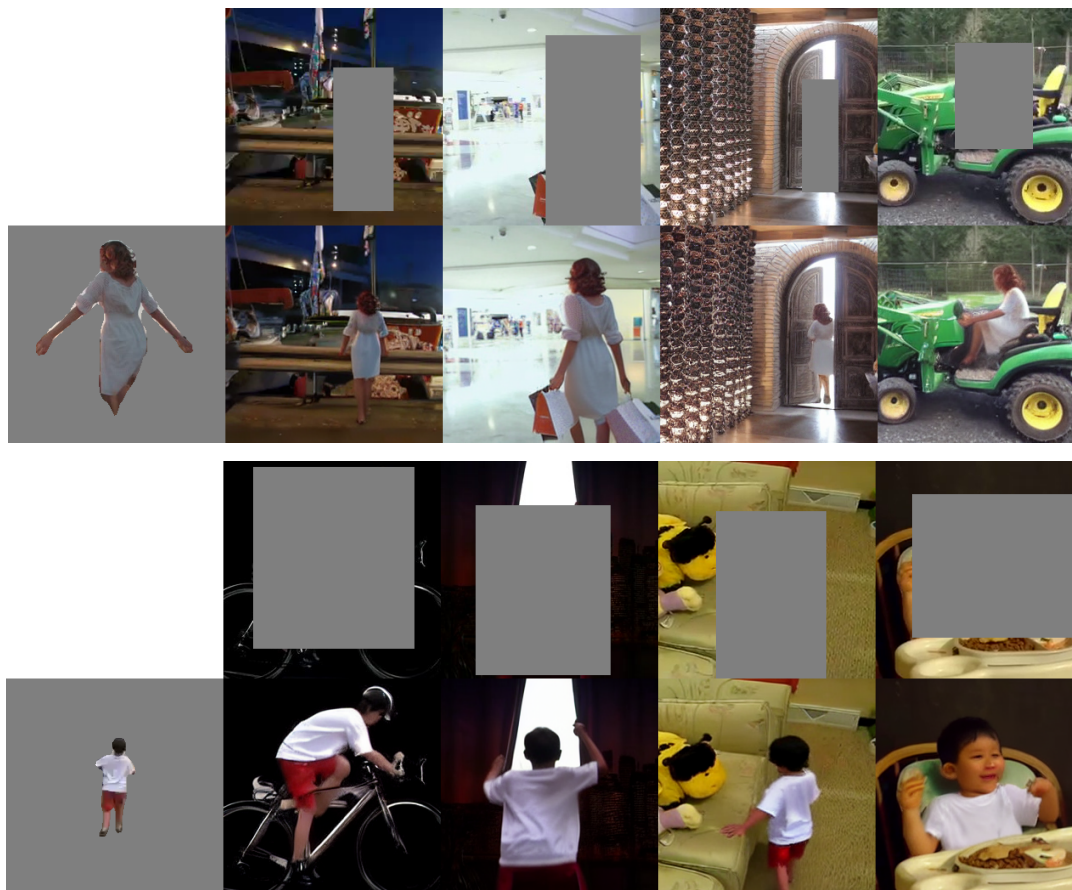
Video data is **critical** for this task, image only data even with aug performs poorly.

	FID	PCKh
Small (scratch)	12.37	15.10
Large (scratch)	11.23	15.87
Large (SD finetune)	10.08	17.60

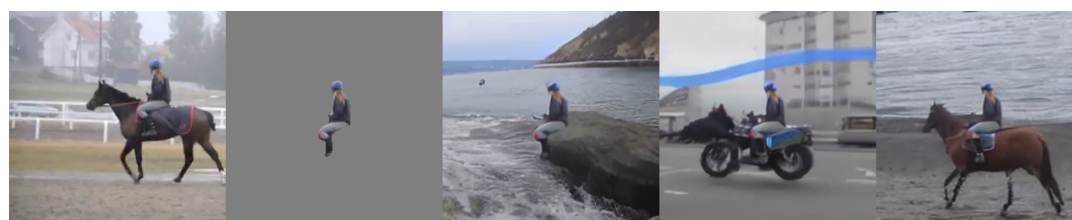
Large-scale models are also **critical**. Initializing with Stable-Diffusion helps.

Qualitative Results

Same Person in Different Scenes



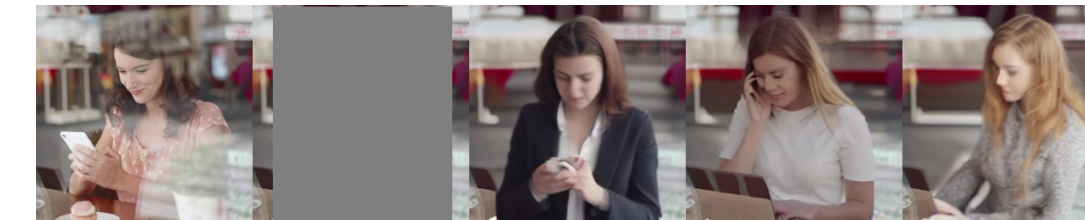
Scene Hallucination



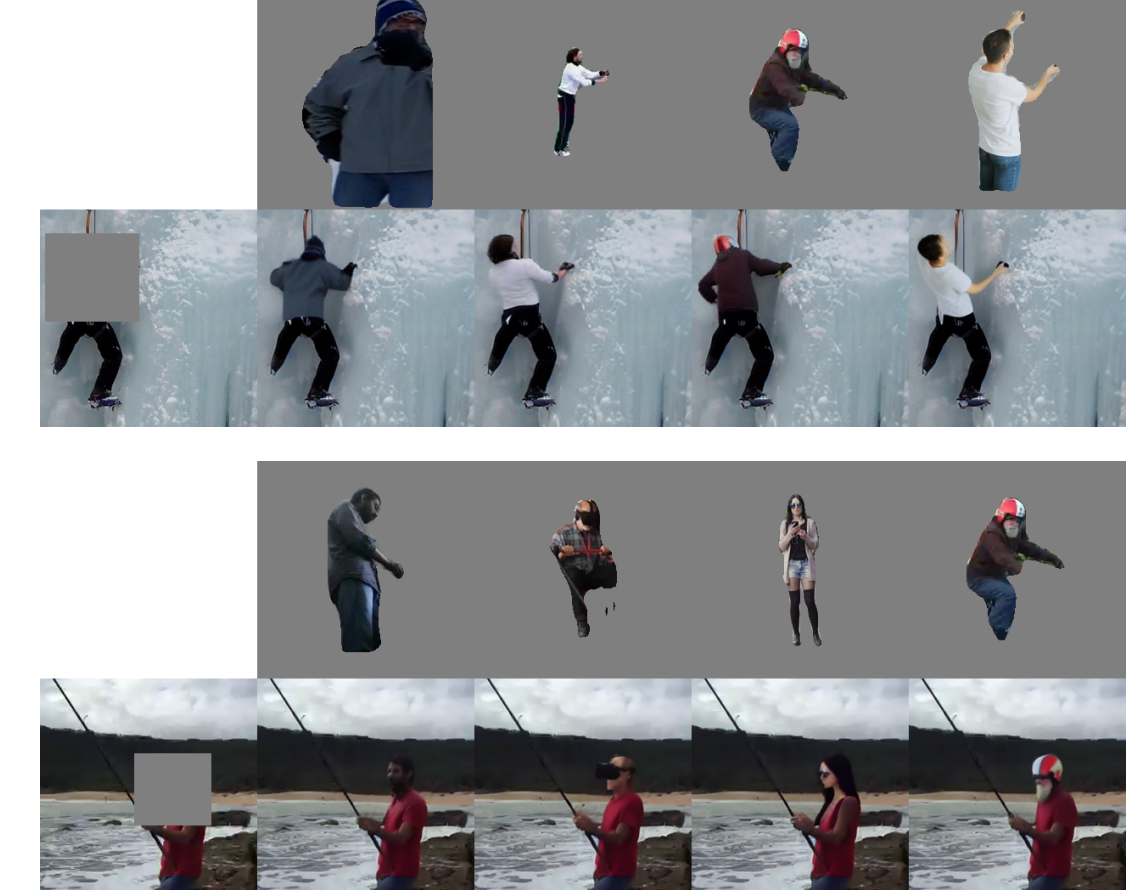
Different People in Same Scene



Person Hallucination



Partial Body Completion



Cloth Swapping

