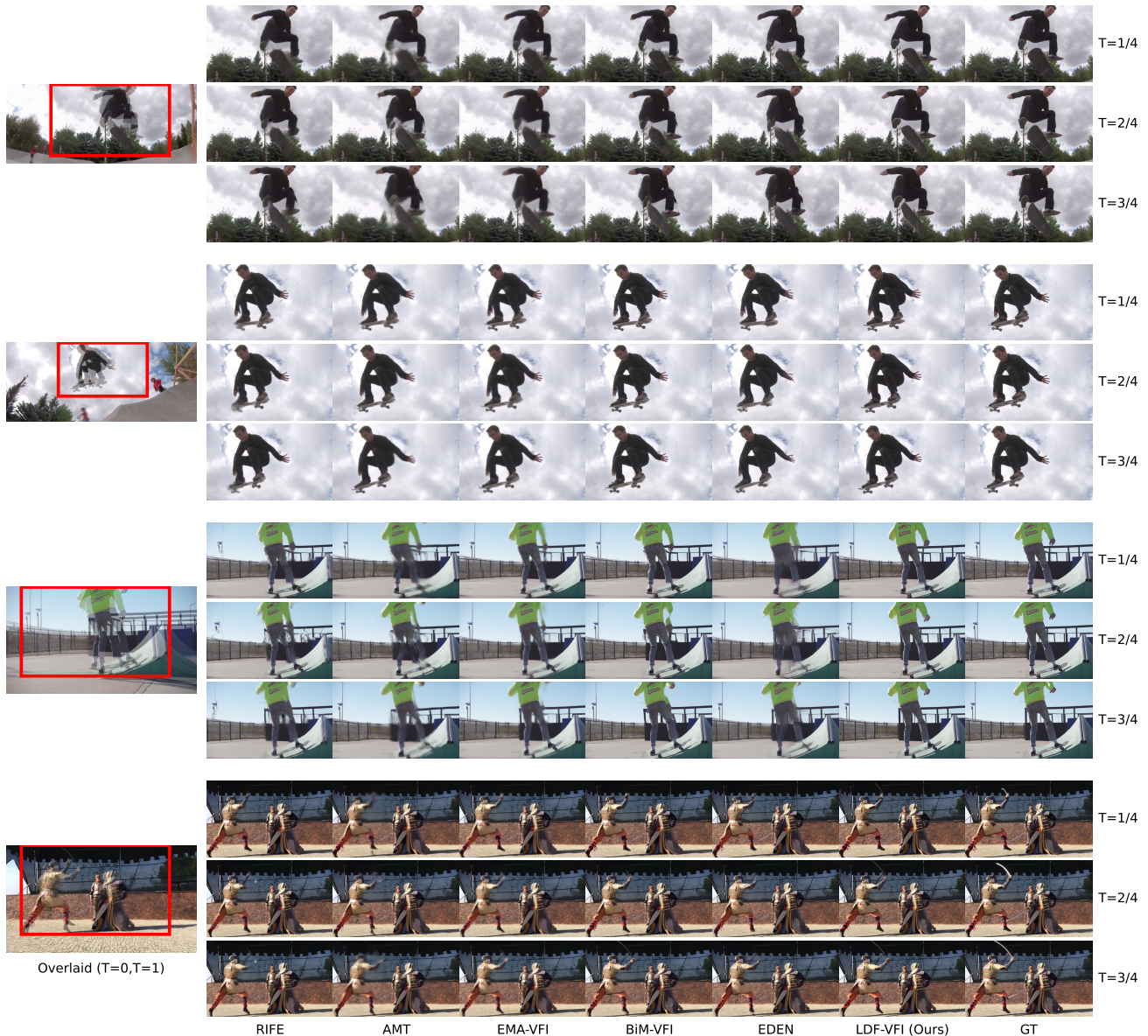


# Towards Holistic Modeling for Video Frame Interpolation with Auto-regressive Diffusion Transformers

## Supplementary Material



**Figure A1.** Qualitative comparisons of  $16\times$  interpolation on SNU-FILM-entire benchmark.

**More visualization** We provide additional visual comparisons in Figure A1.

**Impact of temporal chunk length** Table A1 compares chunk sizes under similar training computation (32K steps for length 5 vs. 8K steps for length 20). A longer chunk (20) improves temporal continuity, e.g., reducing X4K FVD from 510.63 to 121.29.

**Table A1.** Impact of temporal chunk length under similar total computation.

Temporal chunk	SNU-FILM-8 $\times$		X4K-16 $\times$	
	LPIPS $\downarrow$	FVD $\downarrow$	LPIPS $\downarrow$	FVD $\downarrow$
Length 5	<b>0.059</b>	42.15	0.198	510.63
<b>Length 20</b>	0.066	<b>24.87</b>	<b>0.126</b>	<b>121.29</b>