# VD3D: Taming Large Video Diffusion Transformers for 3D Camera Control

Sherwin Bahmani[1,2,3]  Ivan Skorokhodov[3]  Aliaksandr Siarohin[3]  Willi Menapace[3]  Guocheng Qian[3]  Michael Vasilkovsky[3]

Hsin-Ying Lee[3]  Chaoyang Wang[3]  Jiaxu Zou[3]  Andrea Tagliasacchi[1,4]  David B. Lindell[1,2]  Sergey Tulyakov[3]

[1] University of Toronto  [2] Vector Institute  [3] Snap Inc.  [4] SFU
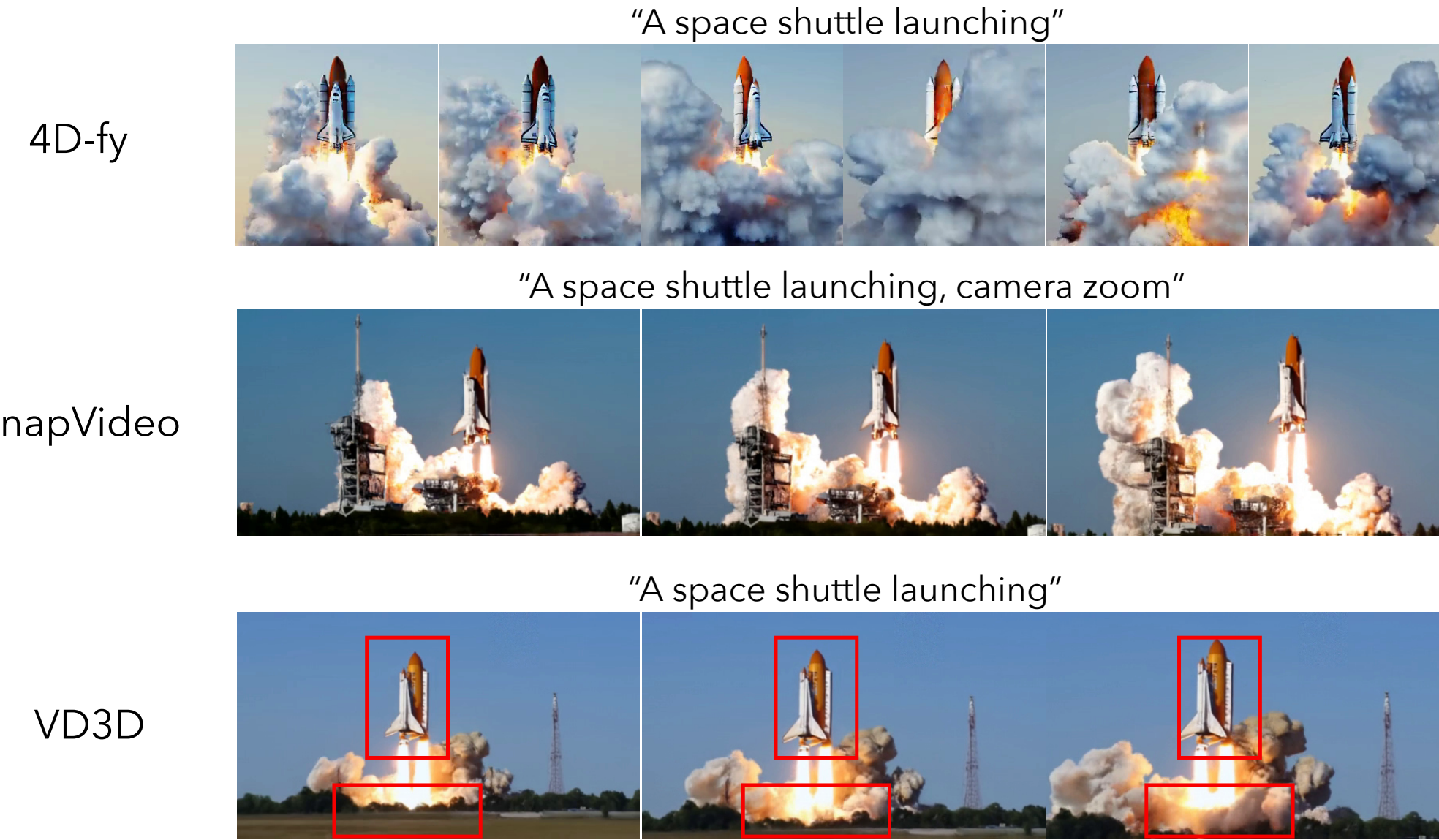
Website

## Motivation

- Text-to-4D generation approaches (e.g., **4D-fy**) lack photorealism
- Previous text-to-video generation methods (e.g., **SnapVideo**) can only control camera with text

4D-fy

"A space shuttle launching"



SnapVideo

"A space shuttle launching, camera zoom"



VD3D

"A space shuttle launching"



**Task:** We tackle camera-controllability for text-to-video diffusion transformers

**Inputs:**

- Text and/or image describing the scene content
- Sequence of camera matrices (extrinsics and intrinsics) describing the camera motion

**Output:** Video following the text, image, and camera motion conditioning

## Comparisons

MotionCtrl and CameraCtrl were designed for U-Net models and not transformers

"3 sheep enjoying spaghetti together"   "Melting ice cream dripping down the cone"

MotionCtrl

CameraCtrl

VD3D



## Method

Our approach injects camera control through Plücker conditioning into a pre-trained video diffusion transformer

"running horse"

text enc.

noise std. $\sigma$
framerate $\nu$
resolution $r$

compression
joint s.t. block
decompression

$\tilde{x}$

time

video pixels to patches

$(\boldsymbol{C}_f, \boldsymbol{K}_f)_{f=1}^{F}$

Plücker coords.

$\ddot{\boldsymbol{p}} \in \mathbb{R}^6$

Plücker coords to patches

XAttn + FF

FIT Block (repeated $B\times$)

latent tok. | SAttn + FF | latent tok. | SAttn + FF | latent tok. | SAttn + FF | latent tok.

$z$   $z$

Plücker conditioning

XAttn + FF

FF

FF$_{\ddot{p}}$  $\ddot{c}$ Plücker tokens   $v$ video tokens

patches to pixels

$\hat{x}$

### Plücker conditioning

$\oplus$   Conv$_{res}$

FF$_{cam}$

XAttn+FF  XAttn$_{cam}$

$\oplus$

$z$  $v$  $\ddot{c}$

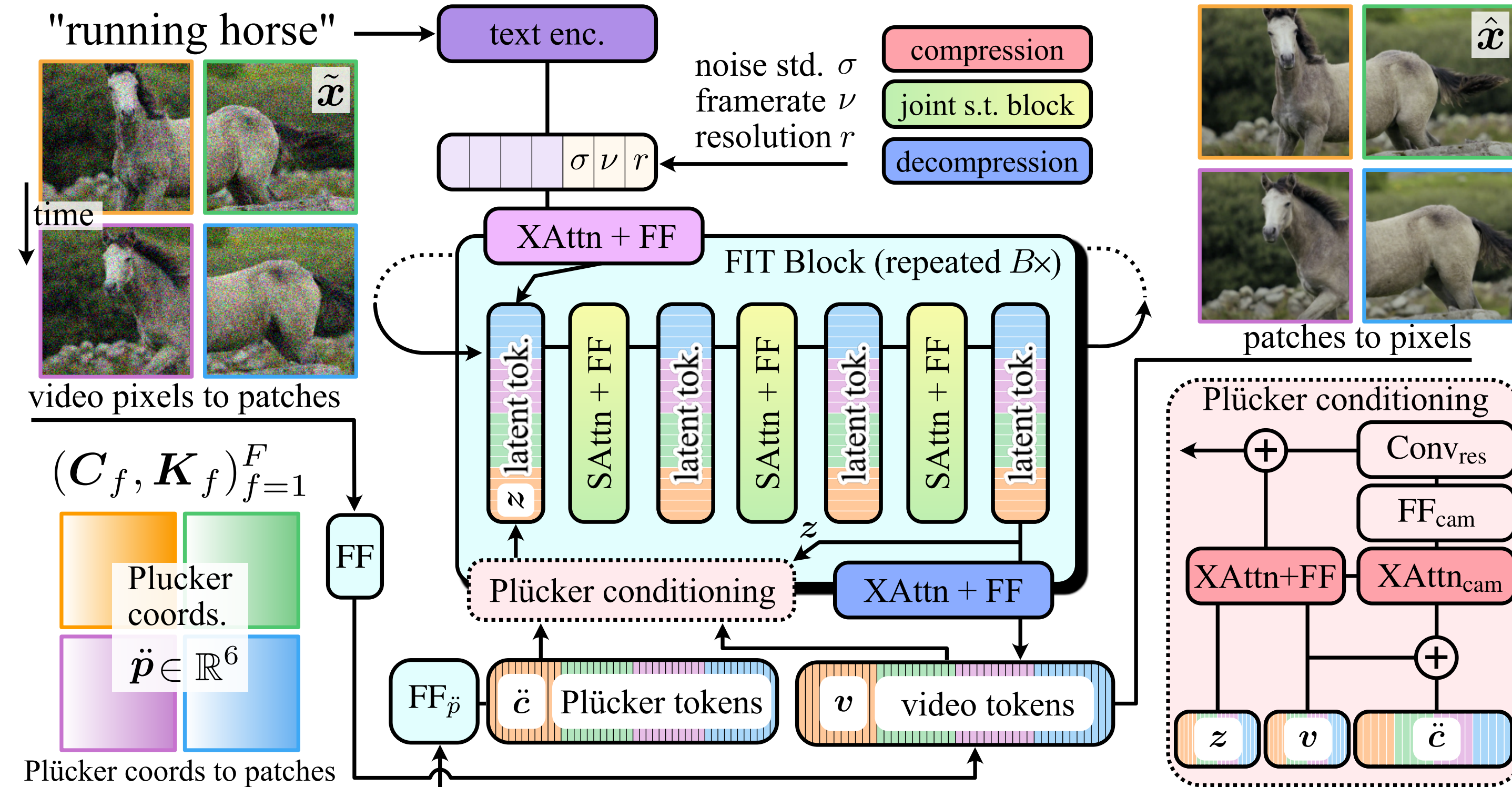## Image-to-Multiview Generation

VD3D can generate multiple viewpoints of a scene given a single real image

Camera   "A bedroom with a bed, lamps and a window"   "A house sitting in the middle of a grassy field"



## Out-of-Distribution Cameras

Rotate Counterclockwise

Rotate Clockwise

Pan Down

Pan Right