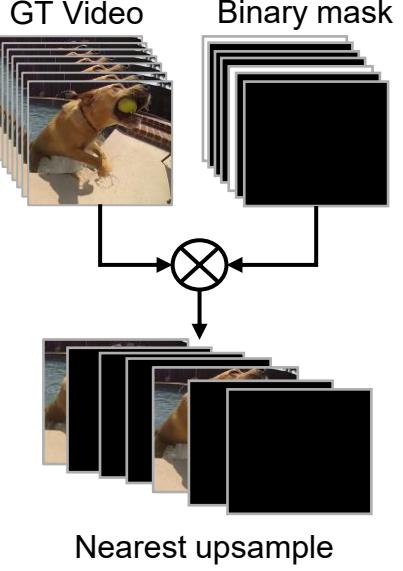
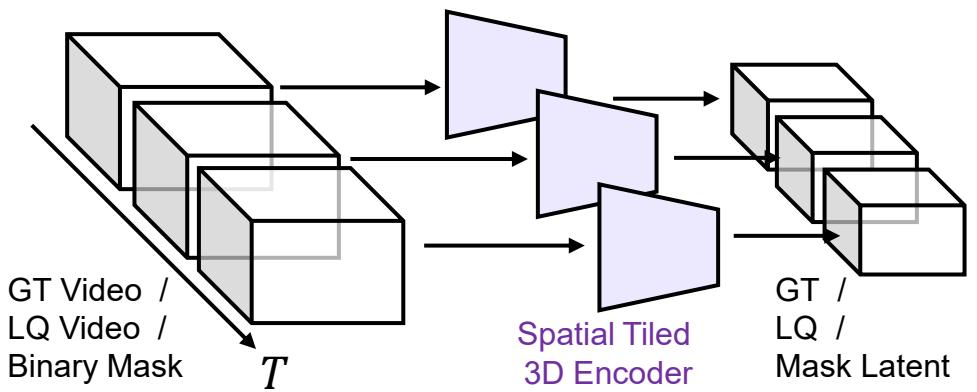


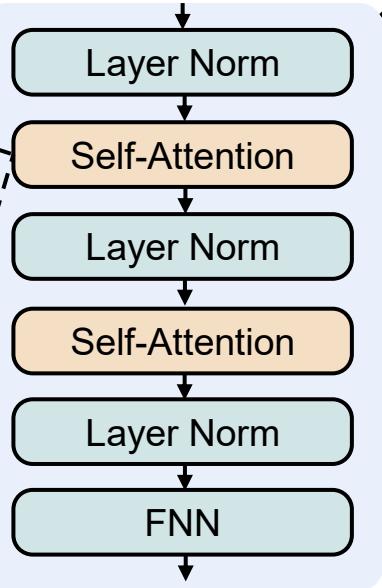
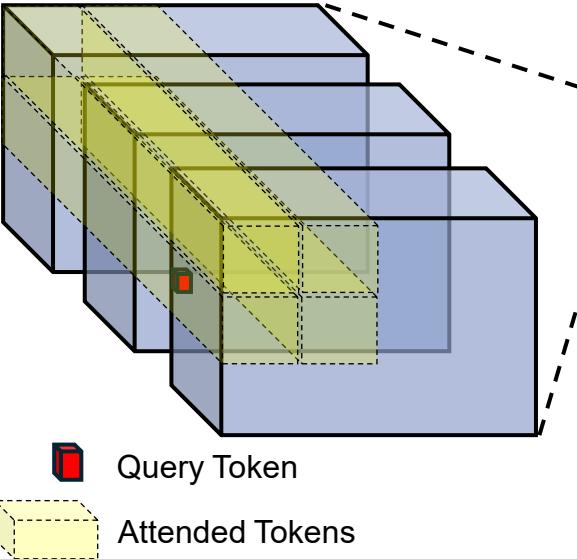
## Video degradation



## Spatial-tiled + temporal chunk encoding



## Sparse attention



GT latent

$\sigma_1 \epsilon_1$

$\oplus$

GT latent

$\sigma_2 \epsilon_2$

$\oplus$

GT latent

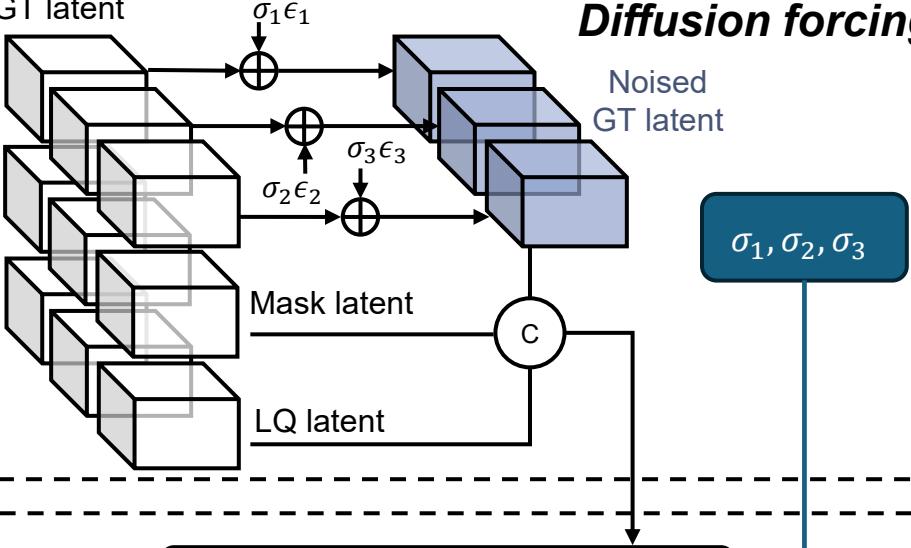
$\sigma_3 \epsilon_3$

$\oplus$

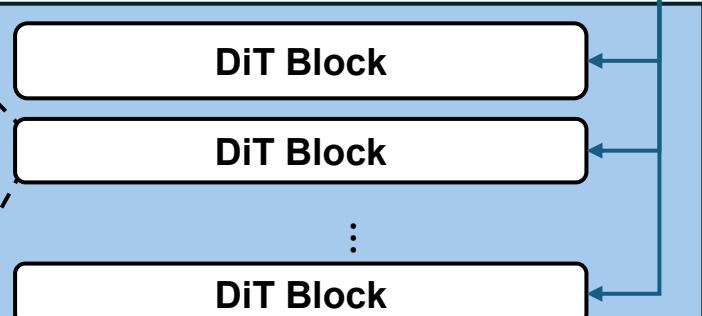
Noised GT latent

$\sigma_1, \sigma_2, \sigma_3$

## Diffusion forcing



## Patchify



## UnPatchify

## Predict velocity