x T timesteps

K Frames

Noised latent input vidéo

Adaptive multi-frame sampling

prompt

Pre-trained T2I diffusion model

with token merging for self-attention

Output vidéo

*Input video*

*Adaptive stochastic permutation*

permuted indices

indices

*Timestep*

*Frame grouping*

t

t'

t''