CLIP Image Encoder

MLPs

**Learned embeddings $f_1$** indoor on the floor dolly motion.

image prompt $I$

fine visual embeddings

coarse visual embeddings

$X^I_t$ $X^I_t$

concated keys

K img Projection

V img Projection

K Projection

V Projection

concated values

Q Projection

Cross-Frame Attention

K Projection

V Projection

Q Projection

Temporal Attention

$X_t$

$X_{t-1}$

Blocks in U-net