

⊗ Cross Attention

⊕ Residual Update

Slow
Video DiT
(World Planner)

Fast
Action DiT
(Action Expert)

Cross Attention & FFN

AdaLN

Layer-wise
Joint Attention

Cross Attention & FFN

AdaLN

Self Attn

K / V

Video Context Routing

Joint Attention

Q K V

AdaLN

Obs-guided query

Q K V

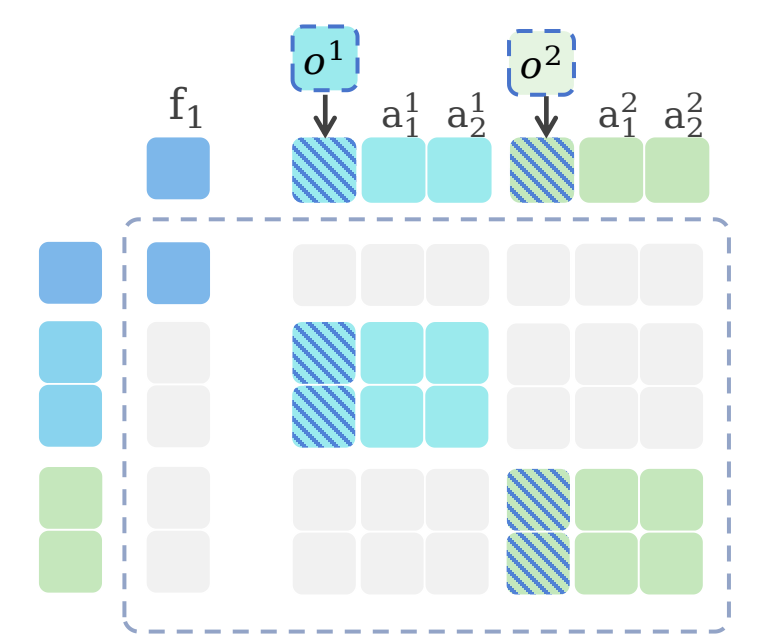
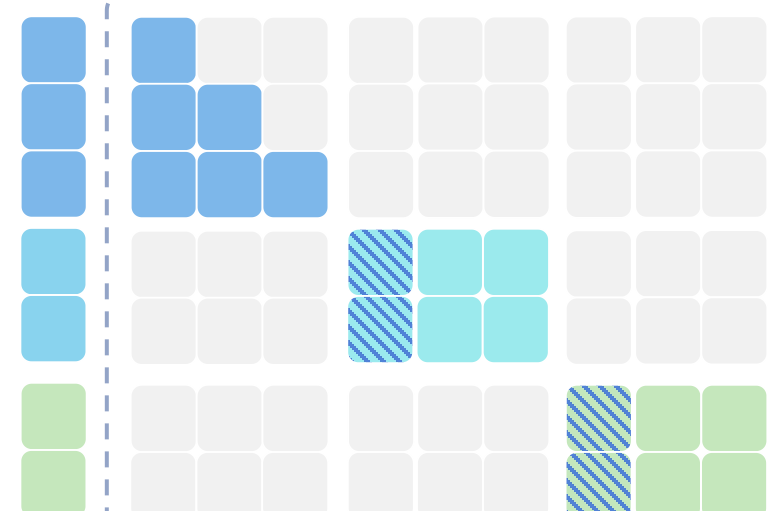
AdaLN

Training

Inference

Video Branch Action Chunk1 Action Chunk2

f_1 f_2 f_n o^1 a_1^1 a_2^1 o^2 a_1^2 a_2^2



Obs-guided query Updated K / V



VAE Encoder

Rolling K/V Memory



Vision Encoder

State Encoder

