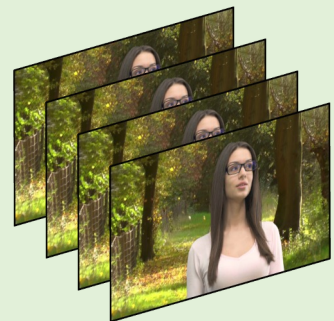


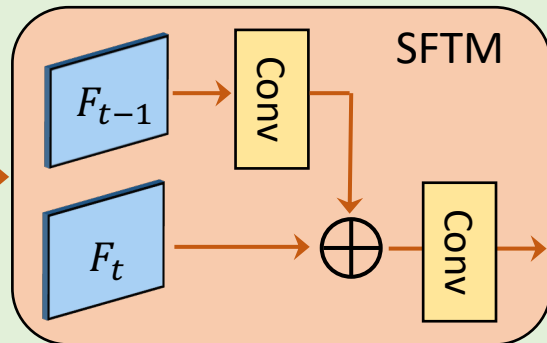
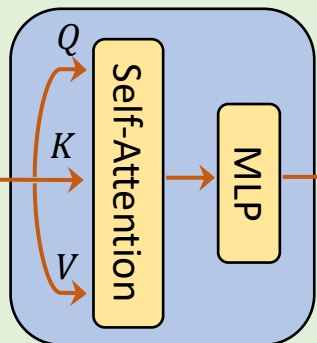
# Feature Modeling



$T \times H \times W \times 3$

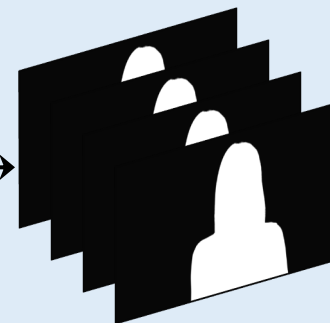
backbone

Transformer Encoder



FPN

# Prediction

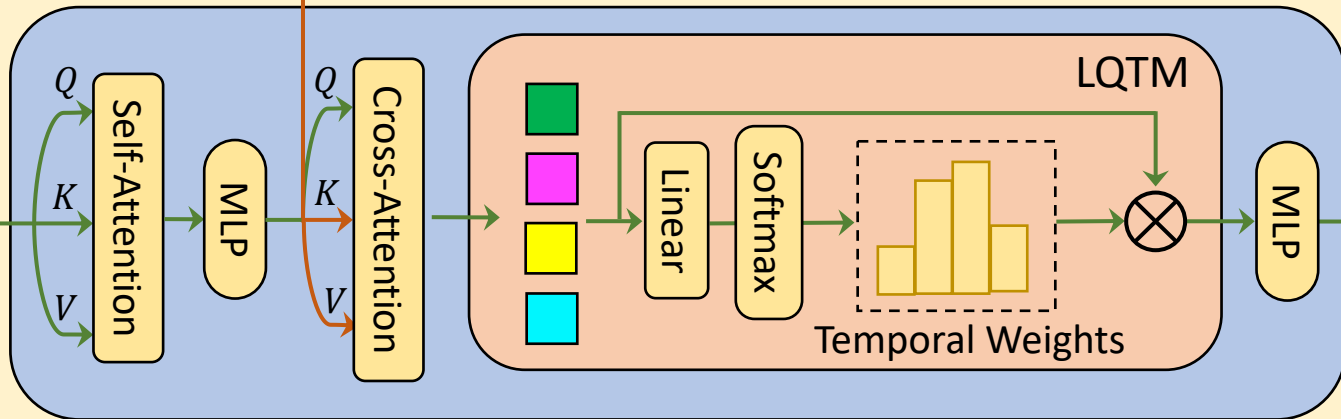


$T \times H \times W \times 1$

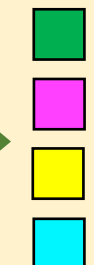
# Query Modeling



$T \times C$



Transformer Decoder



$T \times C$

