

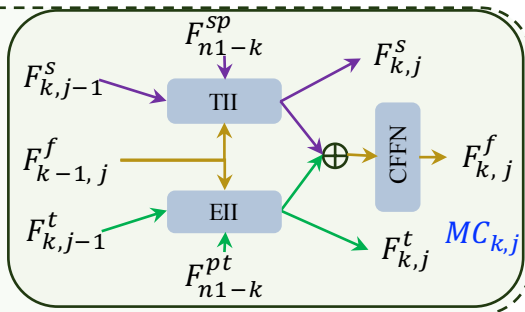
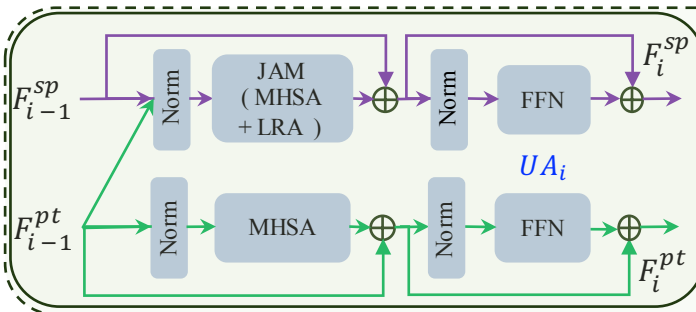
F^{sp} Feature from superpixel embedding

F^{pt} Feature from patch embedding

$F_{k,j}^s$ Features fused with semantics

$F_{k,j}^t$ Features fused with structure

$F_{k,j}^f$ Feature fused with structure and semantics



UA_i i-th Transformer cell

$MC_{k,j}$ Structure cell at k-th level and j-th scale

CRT Collaborative relation Transformer

MDN Compensative decoding network

\oplus Addition operation

\odot Concatenation

SPE Superpixel embedding MHSA Multi-head self-attention CFFN Convolution feedforward layer

PTE Patch embedding JAM Joint Attention Module LRA Locality-related cross attention

EII Semantic information injection module

TII Structure information injection module

FFN Feedforward layer