Input Video:  Input Audio:

$N$ blocks

$T \times H \times W \times C$    $T \times C$

$T \times (H \times W + 1) \times C$

Spatial Fusion

$T \times H \times W \times C$    $T \times C$

$T \times T \times C$

Temporal A-V Fusion    Temporal V-A Fusion

: Video features at different stages    ◎ : Concatenation

: Audio features    ⬡ : Multiplication and summation

🔷🔶🟢🔺🟣 : Audio-constrained query    ⊗ : Element-wise multiplication    ⊕ : Element-wise addition

$(N - 1)$ blocks

$T \times H \times W \times C$

$T \times C \times 1 \times 1$    ⊗ gate

Pooling

Conv

feature from $i$-th block    feature from $(i + 1)$-th block

Video features from backbone    Prediction masks $M_{pos}$

FPN    $F_{seg}$    ⬡    $M_i$    Merge

Dynamic Kernel Head $G_{Kernel}$    Reference Head $G_{Ref}$    $R_i$

Audio-Queried Decoder

$T \times H \times W \times C$    $N \times T \times C$

🔷🔶🟢🔺🟣

Audio features from backbone

Audio-Constrained Query

(a) Decoupled Audio-Visual Transformer    (b) Blockwise-Encoded Gate    (c) Audio-Queried Decoding