

Inputs ($t-1, t$)

Boxes + MV + DCT

Boxes
($t-1$)
 $N \times 4$

MV_t
60 × 60 × 2

DCT
(t)
80 × 80 × C

Box
Enc
 $N \times 32$

MV
ROI
Stats

DCT
ROI
Stats

ROI Extraction

32K params (box-aligned)

$N \times 64$

$N \times 32$

Concat
 $N \times 128$

BiLSTM
 $N \times 128$

Δpos
 $N \times 2$

$\Delta size$
 $N \times 2$

Fusion & Temporal

196K params

Boxes
(t)
 $N \times 4$