

Audio



Audio encoder

BYOL-A

Projection

Video frames



Visual encoder

TimeSformer

Projection

Face encoder

Inception ResNet

LSTM

Projection

F_U

F_A

F_V

F_F

Cross-Attention Fusion (CAF)

CA

CA

CA

F_S

SA

F_{CAF}

Linear

L_{cls}

L_{ss}