

Window

MFCC

MobileNetV2
& CA

Audio
Features

FC layers

Voting

\hat{Y}^n (*Filling Type*)

Frames

MobileNet2
& CA

RGB
Features

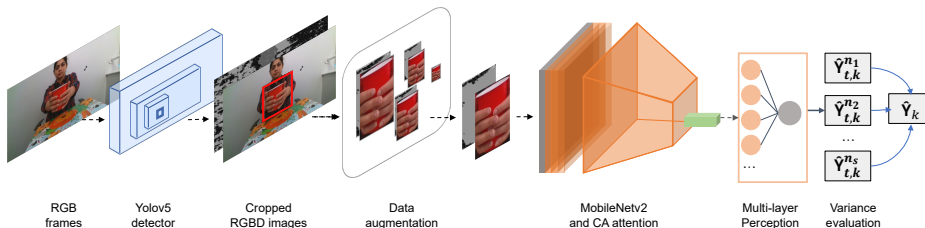
Concatenated
Features

Concatenate

LSTM

FC layers

\hat{Y}^n (*Filling Level*)



Pre-trained without fine-tuned

Trained from scratch

After average pooling

Pre-trained and fine-tuned