

Input Image



Feature
Extraction



- Human Appearance Stream
- Object Appearance Stream
- Pose Feature Stream
- Spatial Feature Stream
- PAM
- Hold book Language Priors

Human Appearance



H^{ws} Block

Object Appearance



O^{ws} Block

Augmentation

Pose Feature

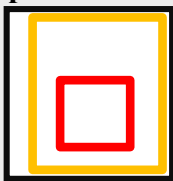


Hold
book

P^{ws} Block

Augmentation

Spatial Feature



Hold
book

S^{ws} Block

H^{sh} Block

O^{sh} Block

P^{sh} Block

S^{sh} Block

A^{sh} Block

SH-VCM

Hold ?



Hold book

A^{ws} Block

H^{sp} Block

O^{sp} Block

P^{sp} Block

S^{sp} Block

A^{sp} Block

SP-VCM

Hold book?

