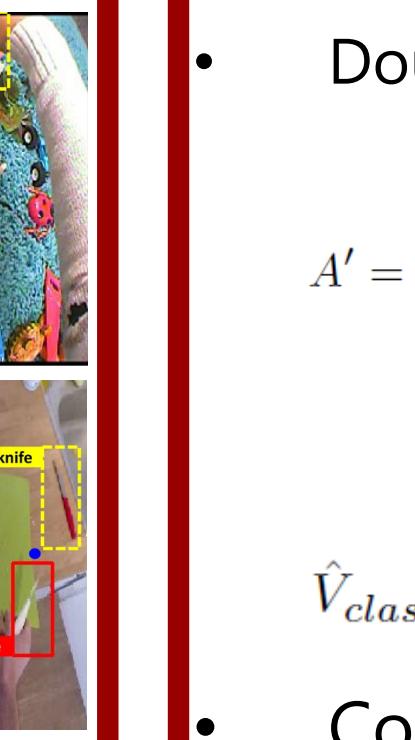


# A Self-Validation Network for Object-Level Human Attention Estimation

Zehua Zhang, Chen Yu, David Crandall  
Indiana University Bloomington  
<http://vision.soic.indiana.edu/mindreader/>

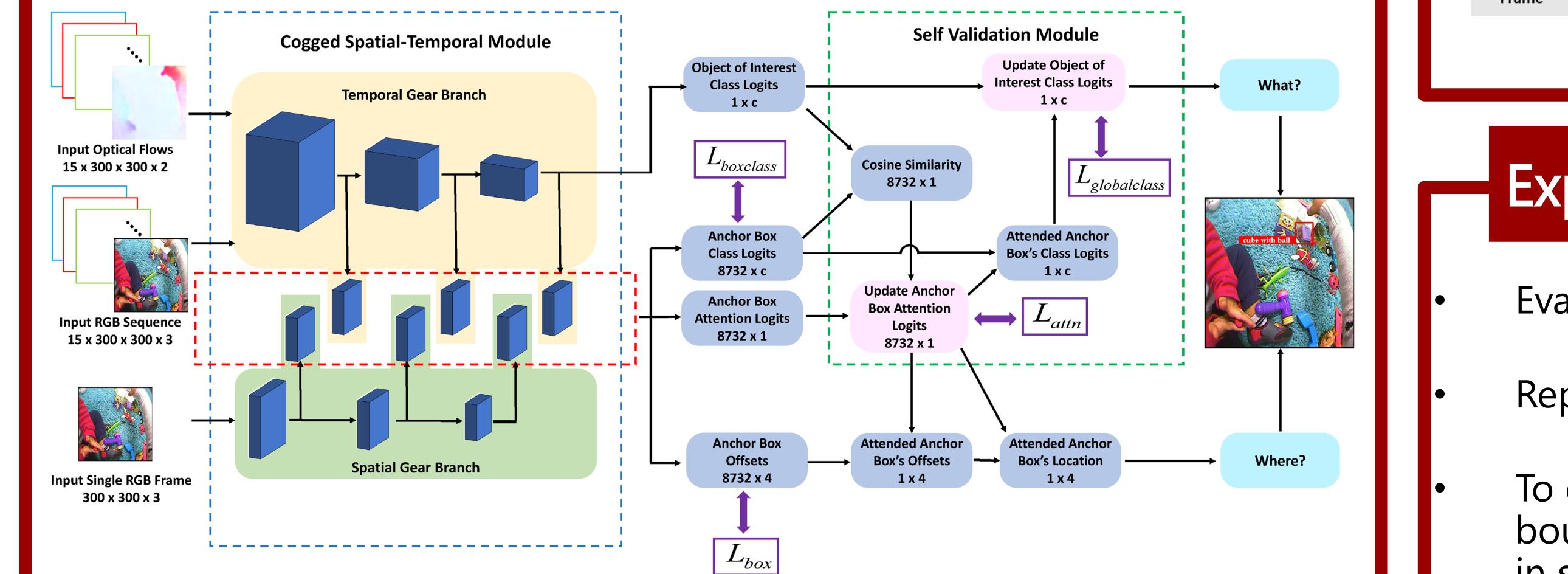
## Introduction

- Due to the foveated nature of the human vision system, people can focus their visual attention on a small region of their visual field at a time.
- Our goal is to jointly identify and locate the attended object among many others in first-person videos.
- Various applications:
  - Applied in VR/AR devices, driver assistance systems
  - Auxiliary information for other computer vision tasks



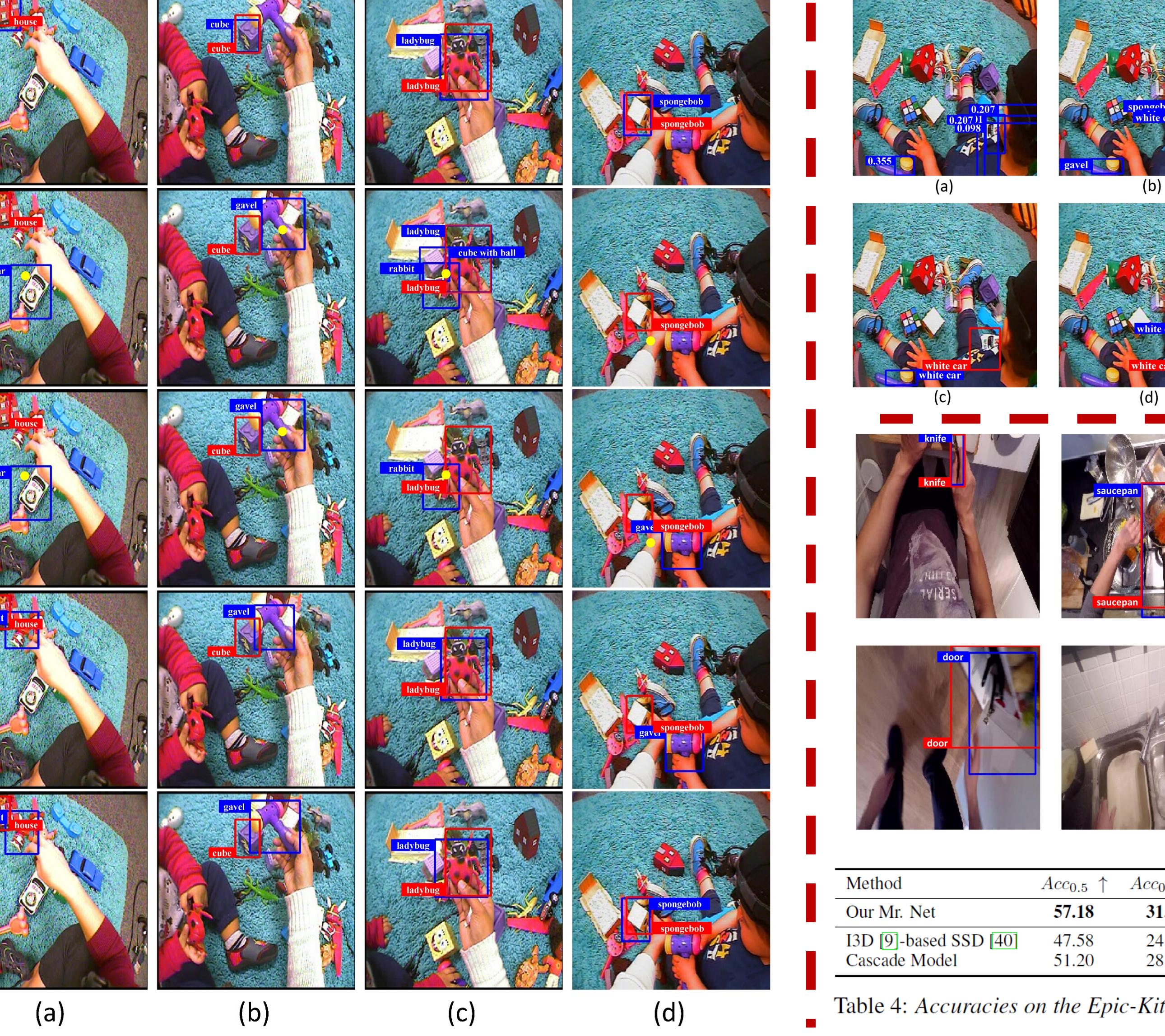
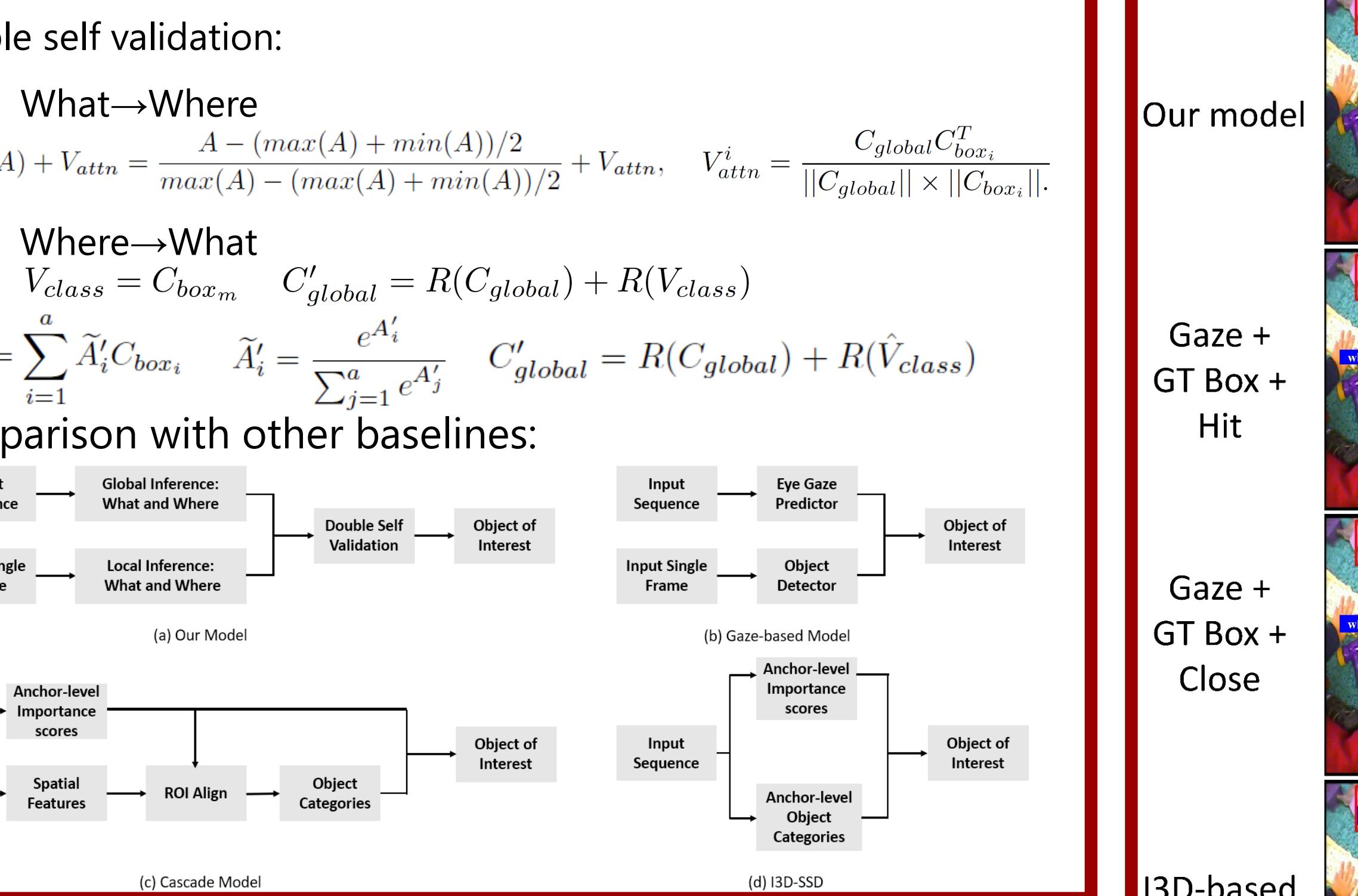
## Model

- We propose a novel unified model to enforce and leverage the visual consistency of our foveated visual system.



## Experiments

- Evaluate on ATT dataset and Epic-Kitchens dataset
- Report Acc\_0.5, Acc\_0.75 and mAcc
- To create stronger baselines, we directly use ground truth object bounding boxes and classes as well as ground truth objects in hands in several implementations of the baselines.



Method	Acc <sub>0.5</sub> ↑	Acc <sub>0.75</sub> ↑	mAcc ↑
Our Mr. Net	<b>74.27</b>	<b>46.78</b>	<b>44.78</b>
Gaze [68] + GT Box + Hit	25.26	25.26	25.26
Gaze [68] + GT Box + Closest	35.86	35.86	35.86
I3D [9]-based SSD [40]	<b>70.11</b>	42.10	40.85
Cascade Model	66.97	<b>45.10</b>	41.93
OIH Detectors + WH Classifier	37.16	37.16	37.16
Left Handed Model	38.31	38.31	38.31
Right Handed Model	39.00	39.00	39.00
OIH GT + WH Classifier	40.83	40.83	40.83
Either Handed Model	42.94	42.94	<b>42.94</b>
Center GT Box	23.97	23.97	23.97

Table 1: Accuracy of our method compared to others, on the ATT dataset. OIH represents Object-in-Hand, while WH means Which-Hand.

Streams	Self validation?		
	Training	Testing	Acc <sub>0.5</sub> ↑
Two	yes	yes	<b>74.27</b>
Two	yes	half	—
Two	yes	no	68.19
Two	no	yes	67.18
Two	no	half	—
Two	no	no	62.33
RGB	yes	yes	74.59
Flow	yes	yes	64.30
Flow	no	yes	—
Flow	no	no	18.40

Table 2: Ablation results. Testing with half means that the model is tested with only what→where validation.

Model	Acc <sub>0.5</sub> ↑	Acc <sub>0.75</sub> ↑	mAcc ↑
Mr. Net	<b>71.34</b>	<b>38.26</b>	<b>39.04</b>
Gaze [68] + GT Boxes Hit	26.46	26.46	26.46
Gaze [68] + GT Boxes Closest	36.81	36.81	36.81
I3D [9]-based SSD [40]	47.58	24.38	25.42
Cascade Model	51.20	28.18	28.36

Table 3: Results of online detection.

Method	Acc <sub>0.5</sub> ↑	Acc <sub>0.75</sub> ↑	mAcc ↑
Mr. Net	<b>71.34</b>	<b>38.26</b>	<b>39.04</b>
Gaze [68] + GT Boxes Hit	26.46	26.46	26.46
Gaze [68] + GT Boxes Closest	36.81	36.81	36.81
I3D [9]-based SSD [40]	67.43	37.90	37.22
Cascade Model	65.96	38.01	37.93

Table 4: Accuracies on the Epic-Kitchen dataset.