

3D Attention-Driven Depth Acquisition for Object Identification

Kai Xu, Yifei Shi, Lintao Zheng, Junyu Zhang, Min Liu, Hui Huang,
Hao Su, Daniel Cohen-Or and Baoquan Chen

National University of Defense Technology Shandong University
Shenzhen University SIAT Stanford University Tel-Aviv University



SA2016.SIGGRAPH.ORG



SIGGRAPH
ASIA 2016
MACAO

Background & motivation

- Robotic indoor scene modeling



Perception on object



SIGGRAPH
ASIA 2016
MACAO

Background & motivation

- Indoor environments acquisition and modeling

Dense Reconstruction



[Nießner et al. 2013]

Object Extraction



[Xu et al. 2015]



SIGGRAPH
ASIA 2016
MACAO

Background & motivation

A white thought bubble with a blue outline and three smaller bubbles leading to it, containing the text "What are these objects?".

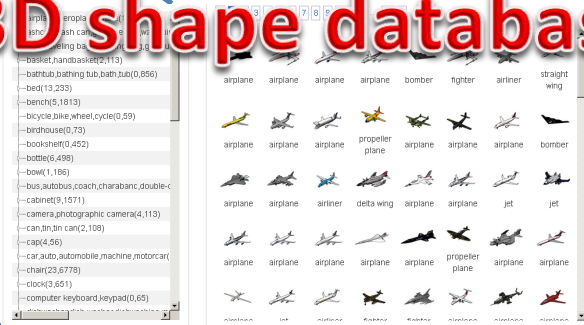
**What are
these objects?**





SIGGRAPH
ASIA 2016
MACAO

3D shape database

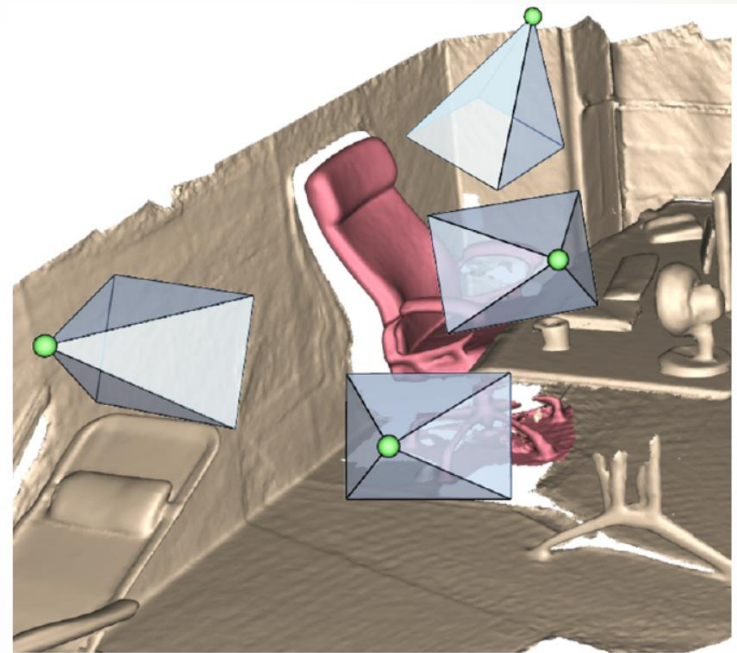


Object recognition



SIGGRAPH
ASIA 2016
MACAO

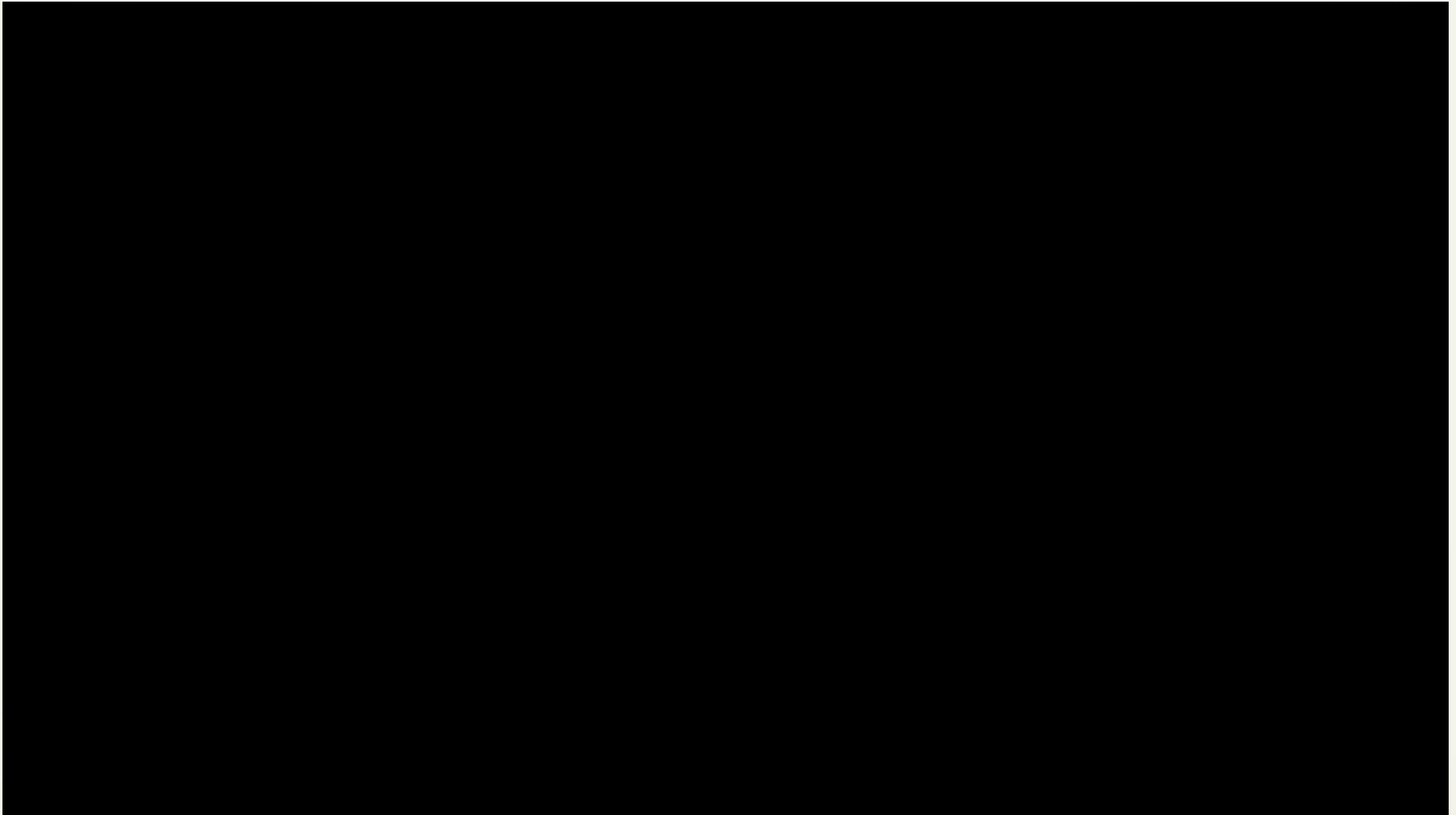
Active object recognition





SIGGRAPH
ASIA 2016
MACAO

Active object recognition



Problem setting

- A robot actively acquires new observations to gradually increase the confidence of object recognition
- Two key components:

Object classification

Estimate object class
based on so far acquired
observations

View planning

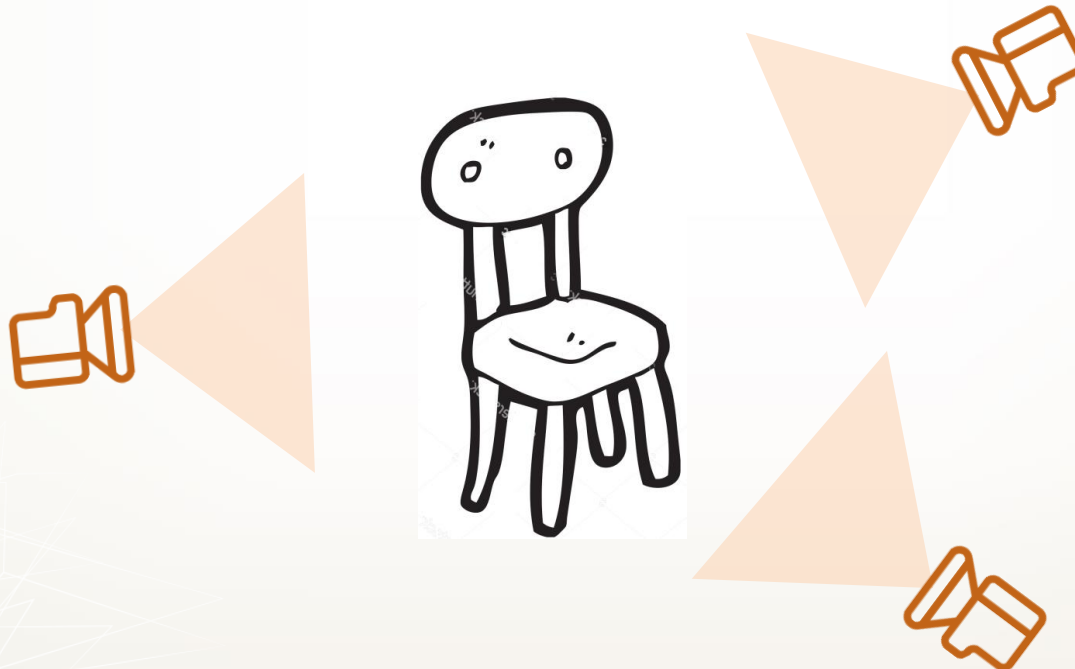
Predict the Next-Best-
View to maximize its
information gain



SIGGRAPH
ASIA 2016
MACAO

The main challenge

- **Observation is partial and progressive**
 - Shape description/matching with partial data is hard
 - Observations from varying views

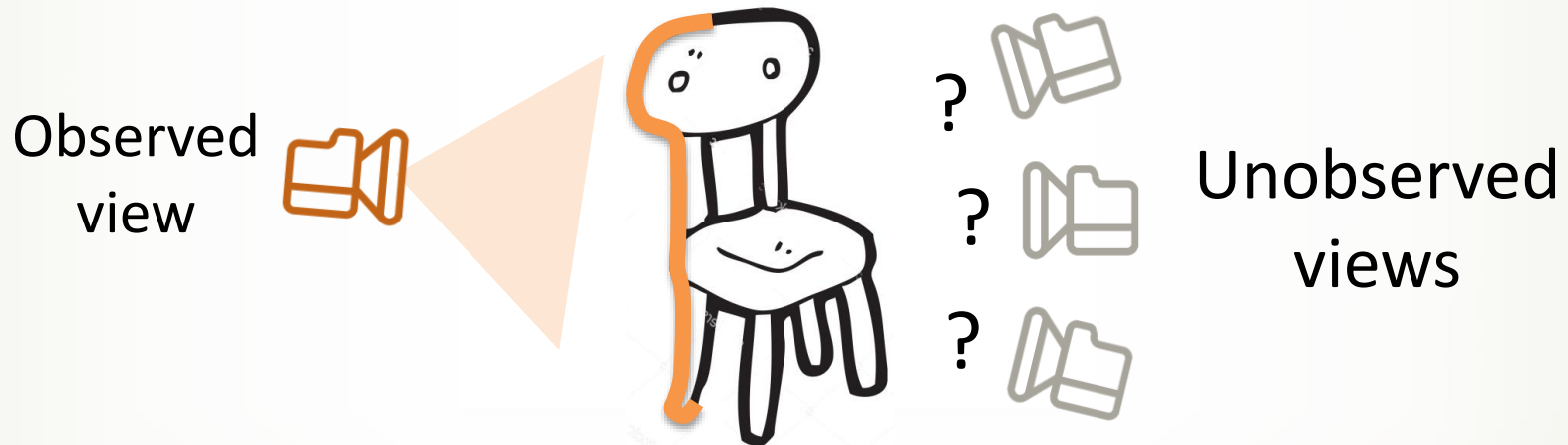


Sponsored by



The main challenge

- Observation is partial and progressive
 - View planning



How can you know which view is better without knowing its observation?

The main challenge

- Real indoor scenes are often cluttered
 - Degrade recognition accuracy
 - Invalidate the off-line learned viewing policy





SIGGRAPH
ASIA 2016
MACAO

Related work

SA2016.SIGGRAPH.ORG

CONFERENCE: 5 - 8 DECEMBER 2016 • EXHIBITION: 6 - 8 DECEMBER 2016 • THE VENETIAN MACAO, MACAO

Sponsored by

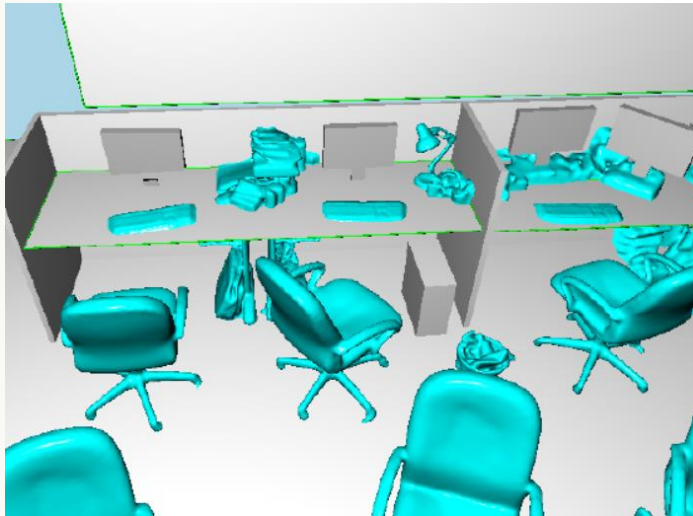




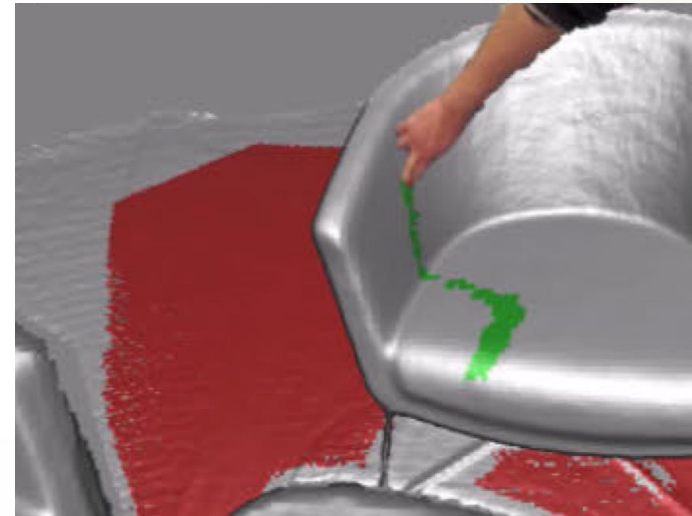
SIGGRAPH
ASIA 2016
MACAO

Related work

- Online scene analysis and modeling



Plane/Object Extraction
[Zhang et al. 2014]



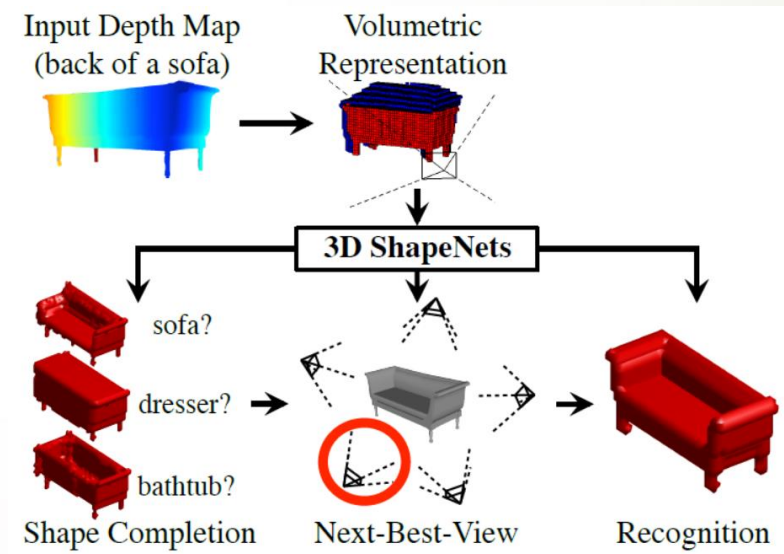
SemanticPaint
[Valentin et al. 2015]

Related work

- Active reconstruction and recognition



Next-best-view for reconstruction
[Wu et al. 2014]



Next-best-view for recognition
[Wu et al. 2015]



SIGGRAPH
ASIA 2016
MACAO

Method

SA2016.SIGGRAPH.ORG

CONFERENCE: 5 - 8 DECEMBER 2016 • EXHIBITION: 6 - 8 DECEMBER 2016 • THE VENETIAN MACAO, MACAO

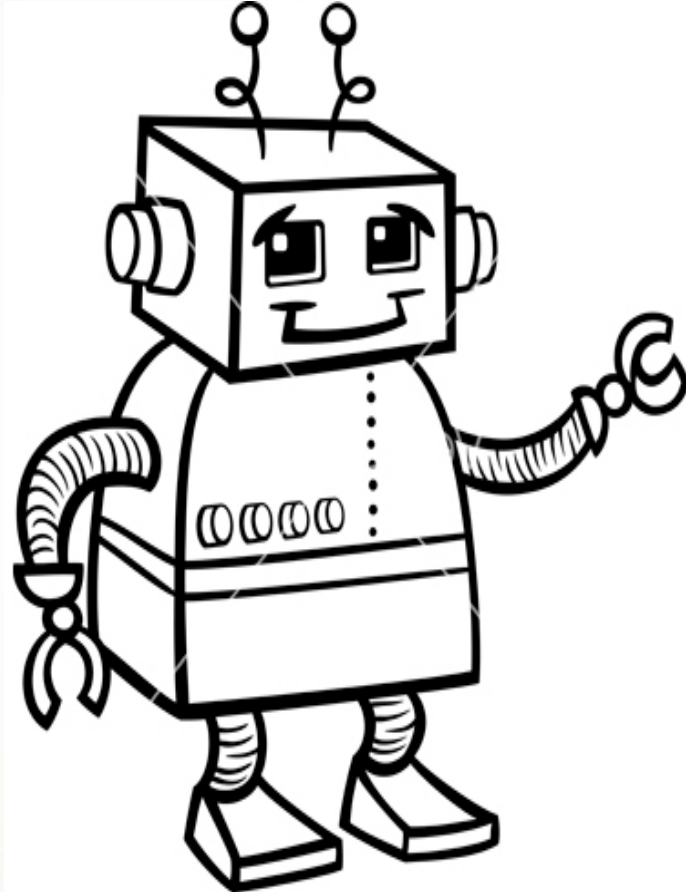
Sponsored by





SIGGRAPH
ASIA 2016
MACAO

The general framework

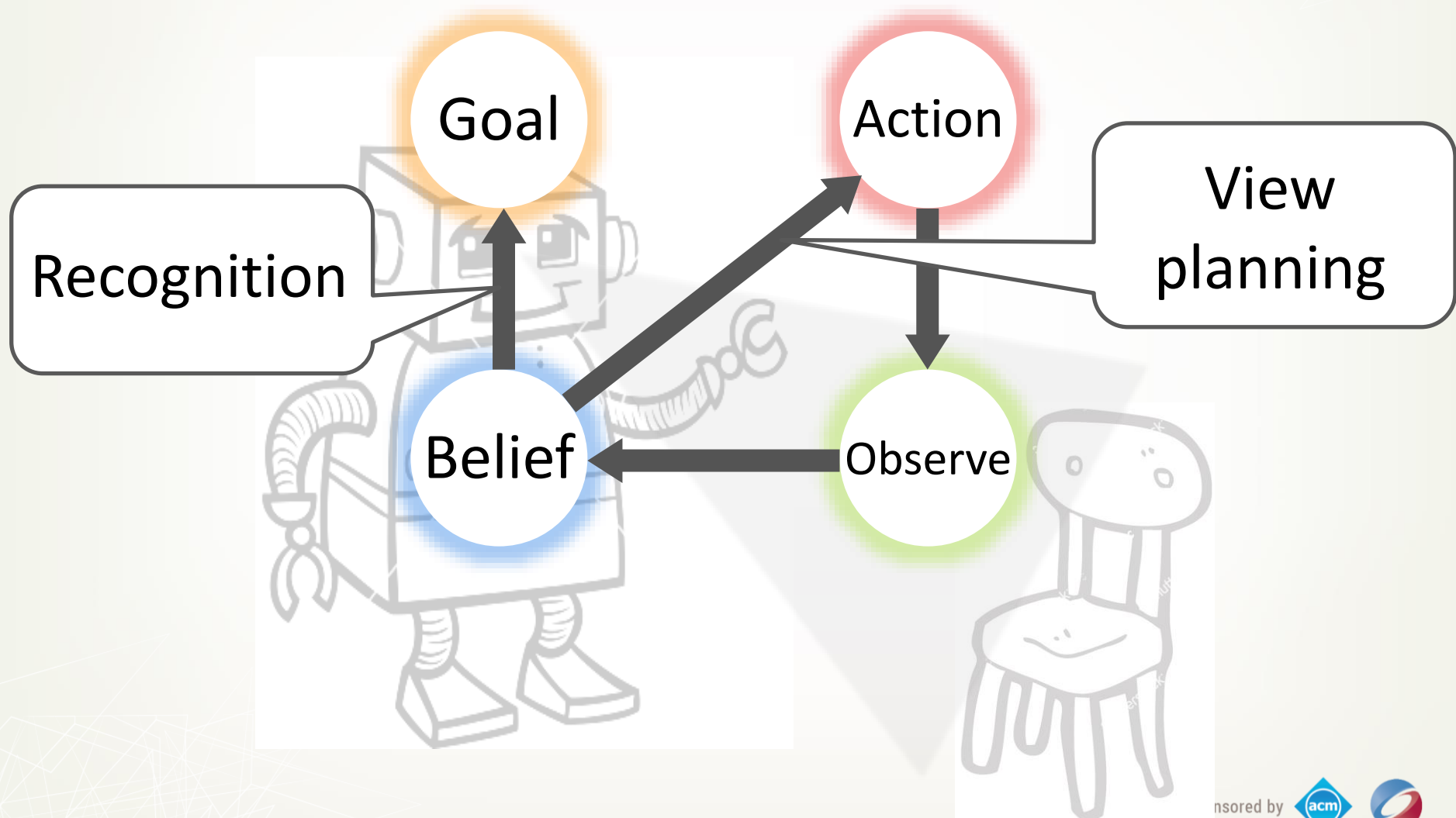


sponsored by





The general framework

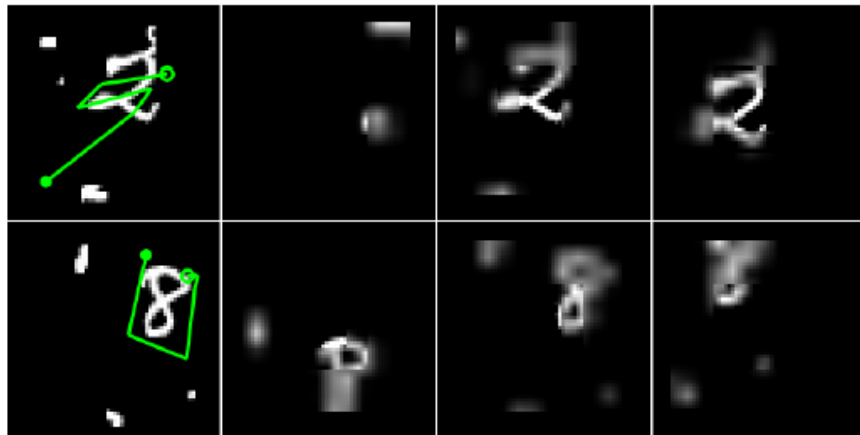


An attentional formulation

“Humans *focus attention selectively on parts* of the visual space to acquire information when and where it is needed, and combine information from different fixations over time to build up an *internal representation* of the scene”

Internal representation

Ronald Rensink



Hand-writing recognition
[Mnih et al. 2014]

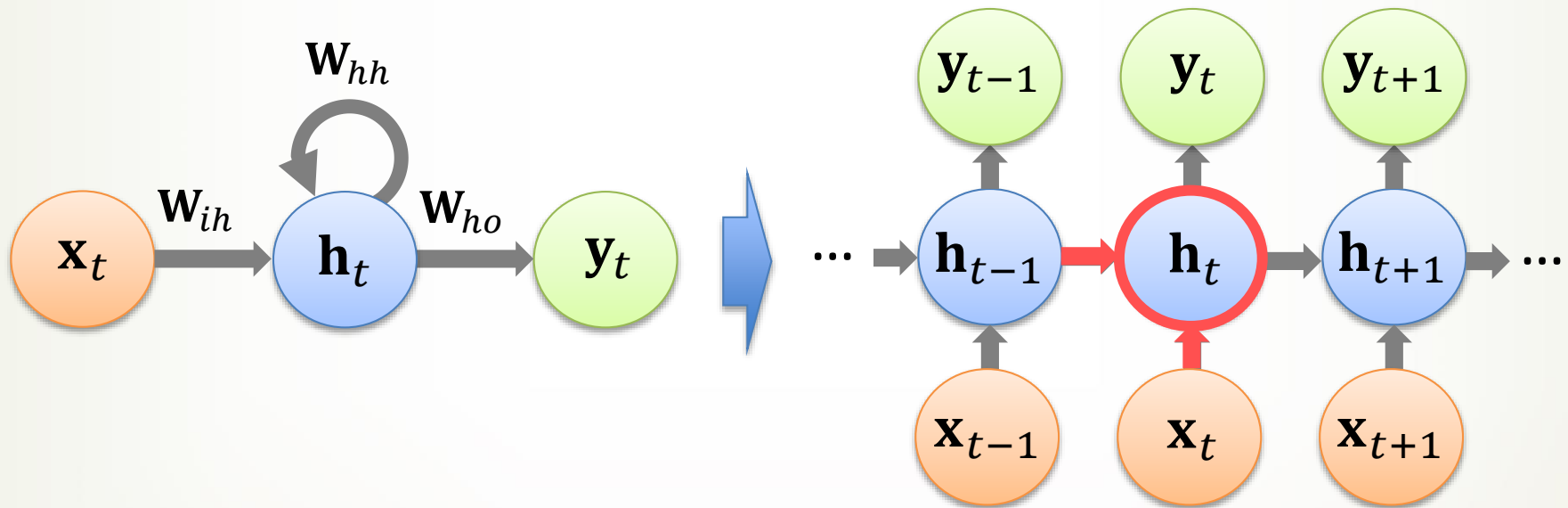


A woman is throwing a frisbee in a park.

Image caption generation
[Xu et al. 2015]

Recurrent Attention Model

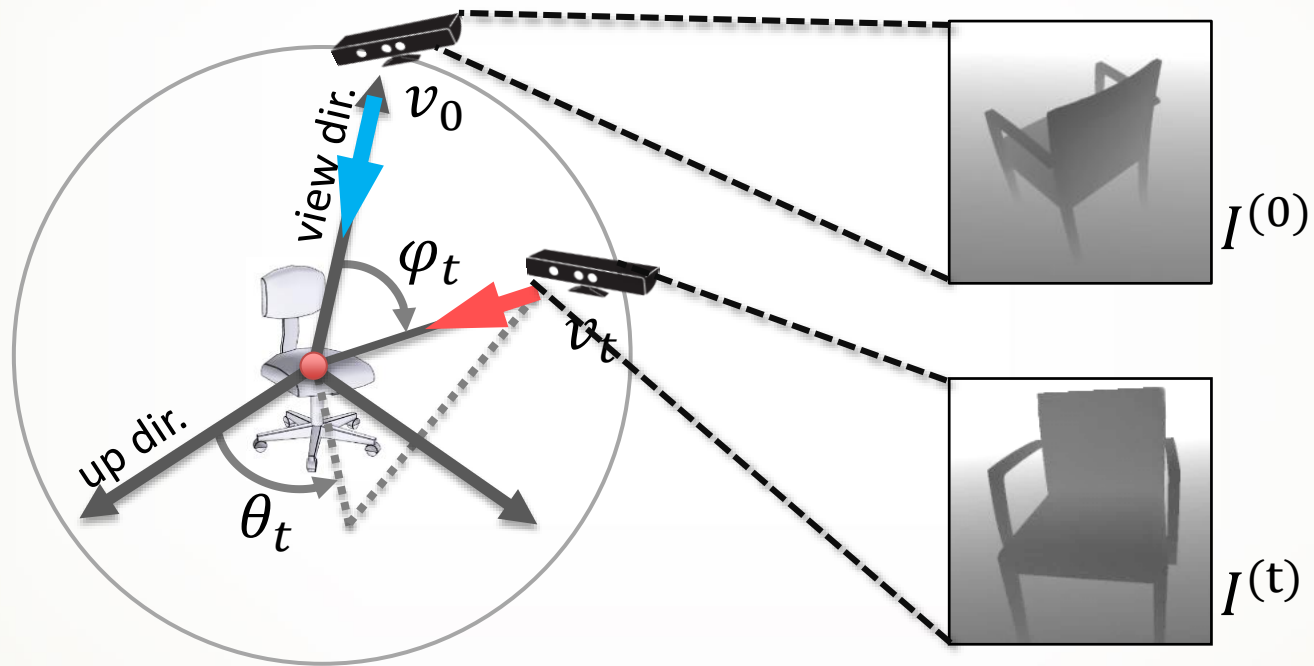
- Recurrent Neural Networks (RNN)



Aggregate information

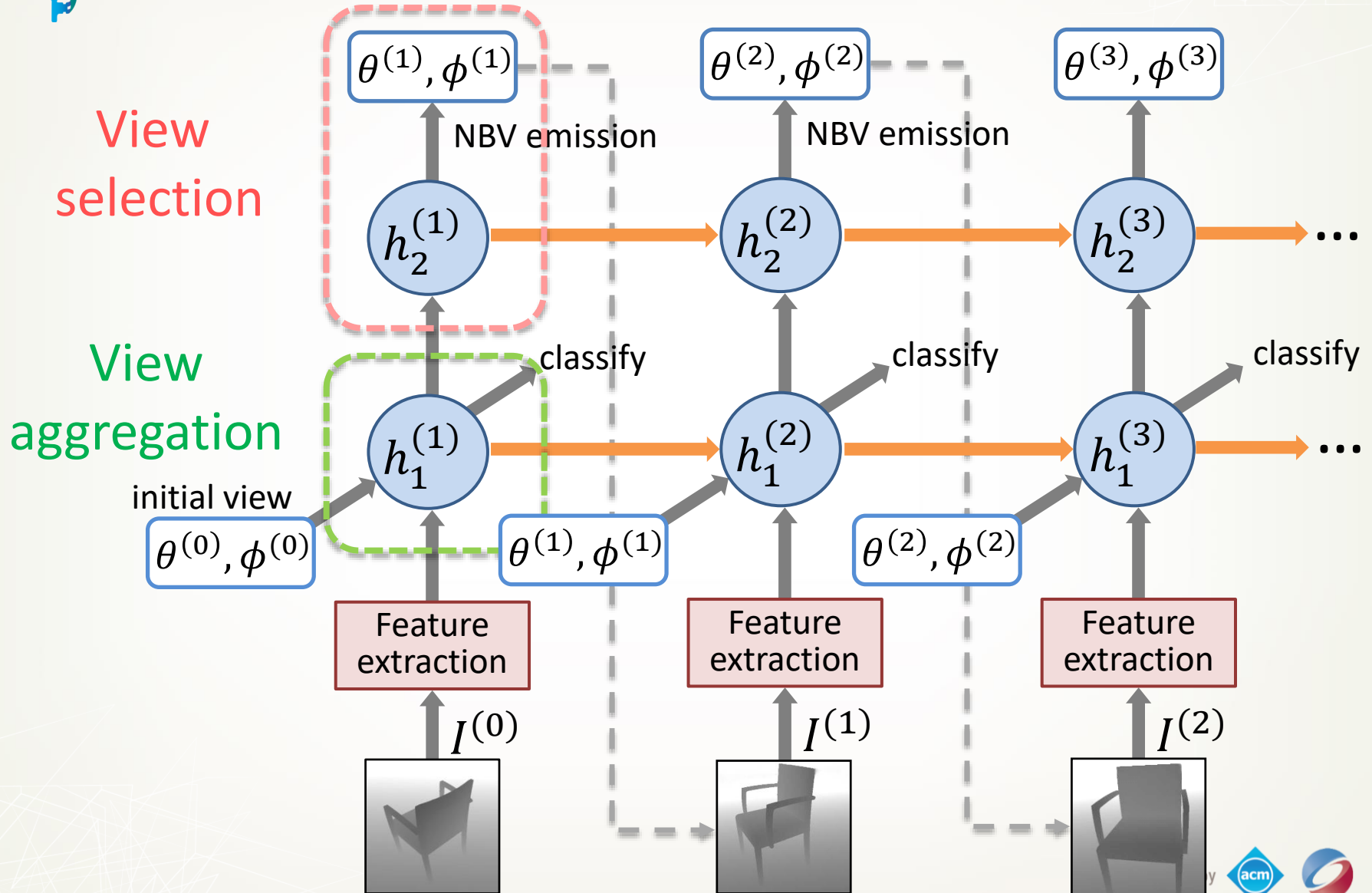


View-based observation



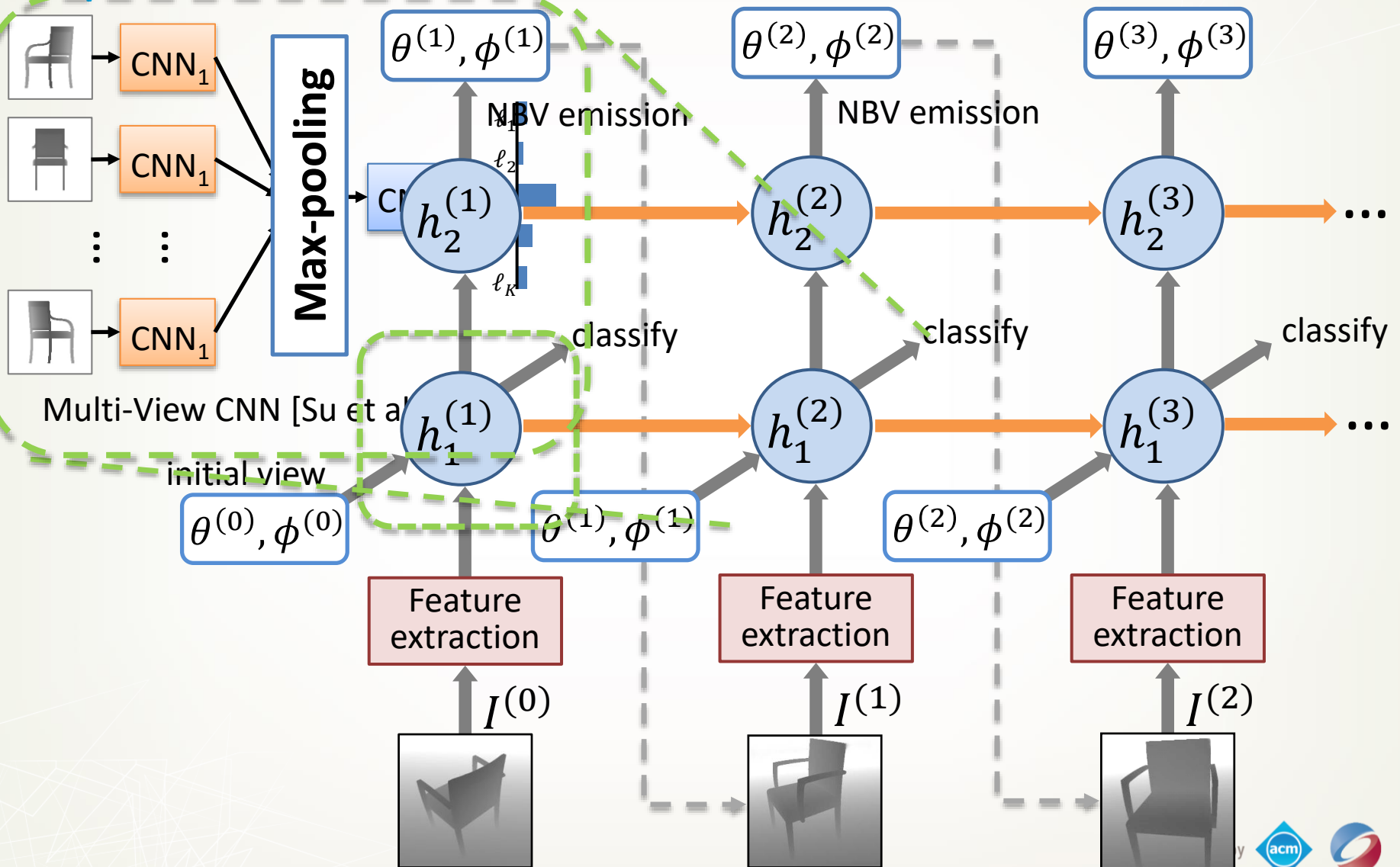


3D Recurrent Attention Model





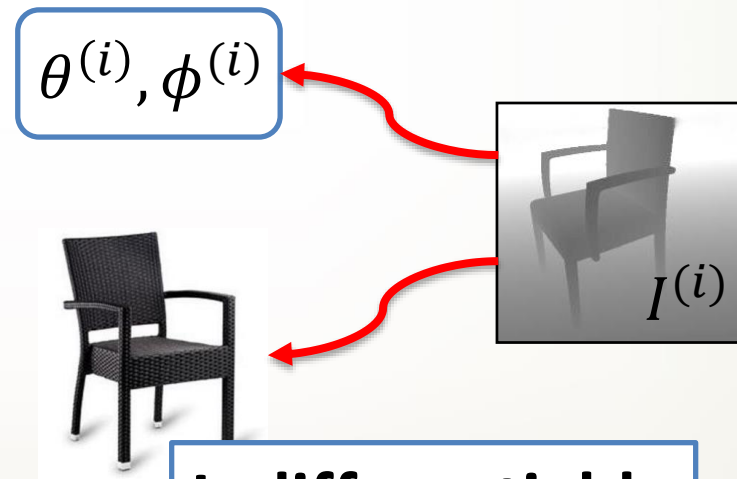
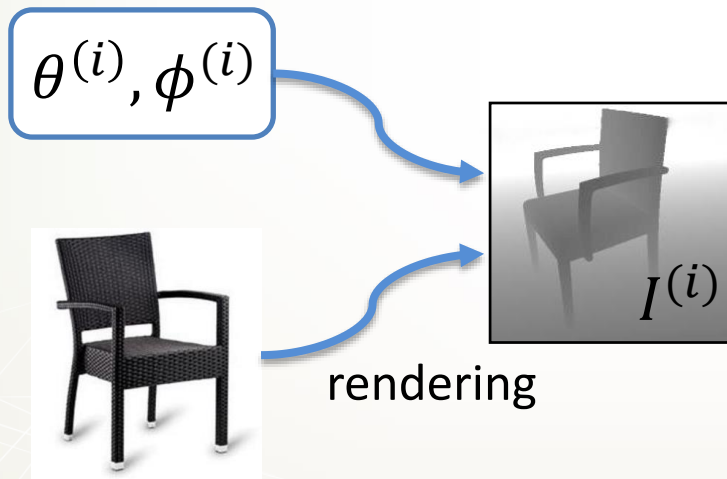
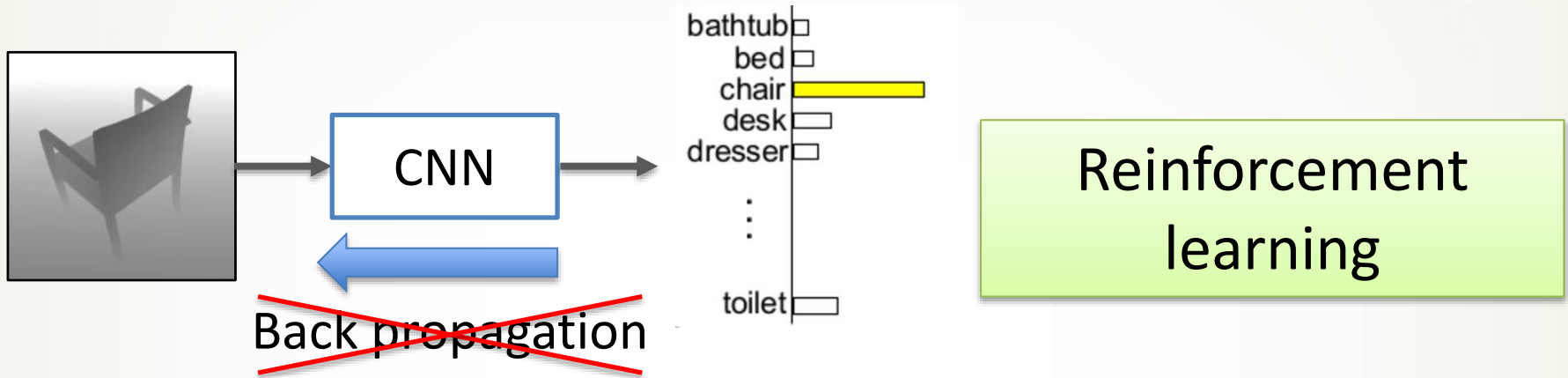
3D Recurrent Attention Model





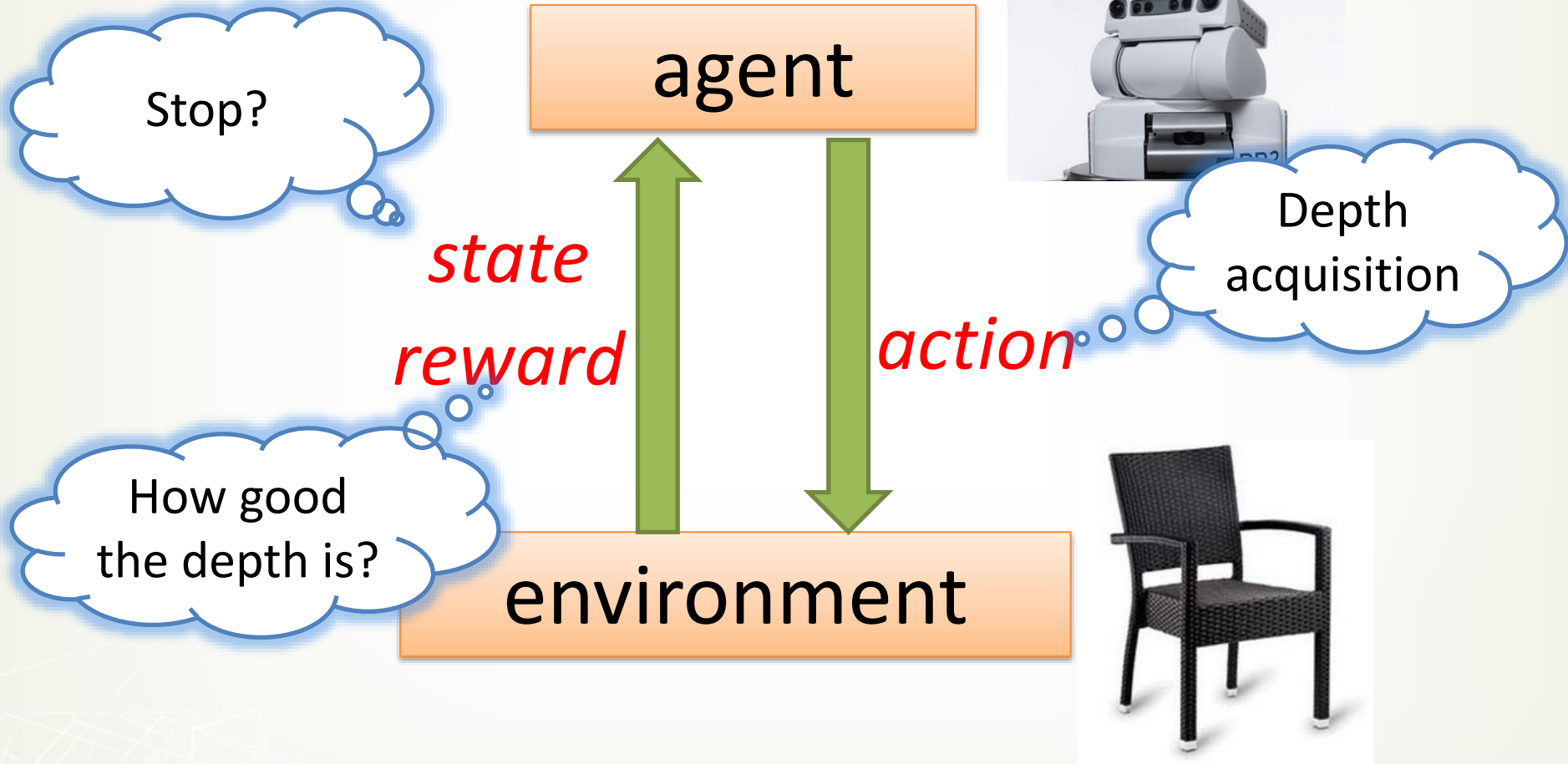
SIGGRAPH
ASIA 2016
MACAO

Network training





Reinforcement learning





Reward

$$r_t = H_t(p_t, \bar{p}) + I_t(p_t, p_{t-1}) - C_t$$

prediction
accuracy

information
gain

movement
cost



Part-level attention



occlusion

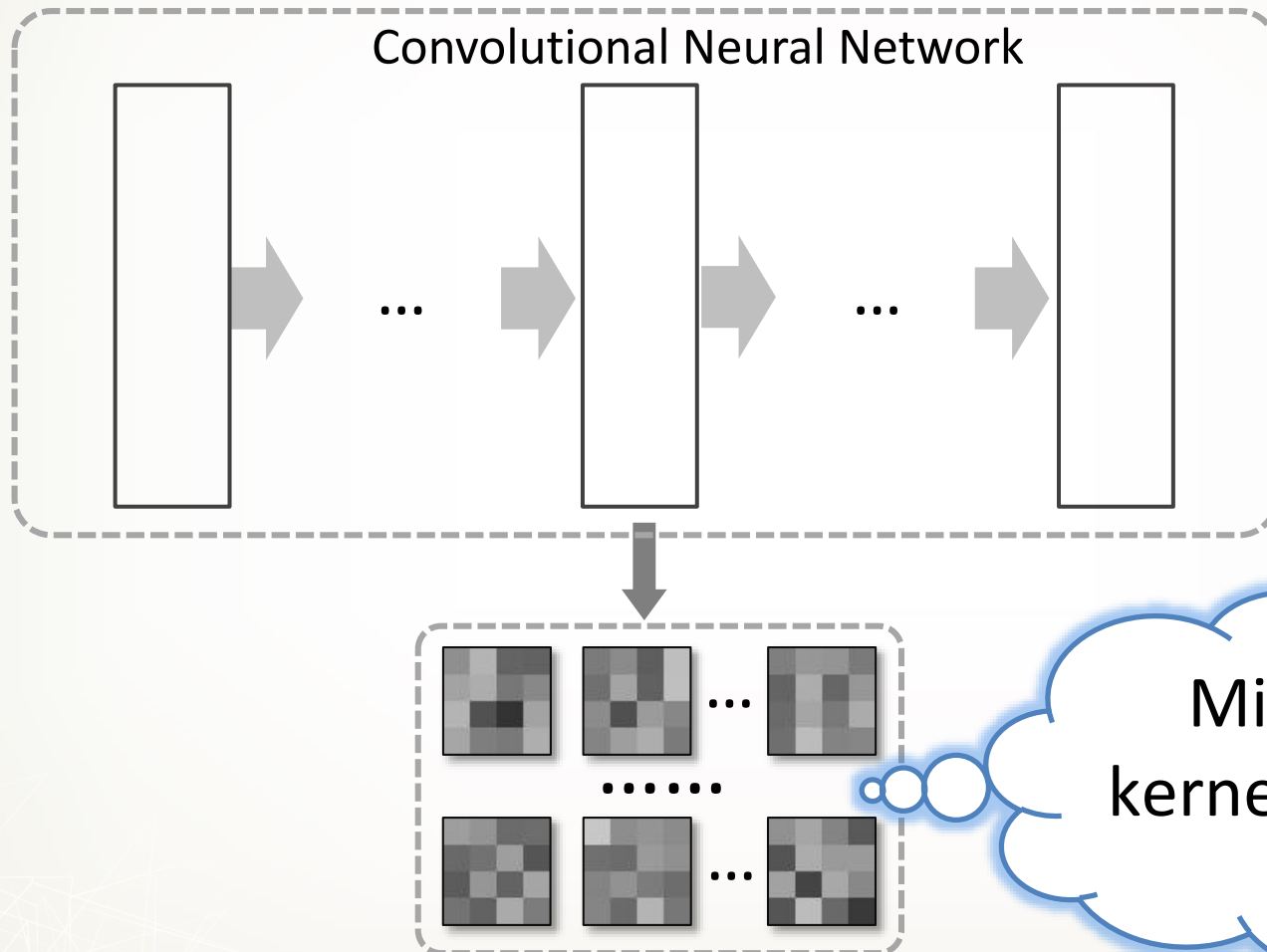


Informative parts

How to distinguish these two chairs?



Attention extraction

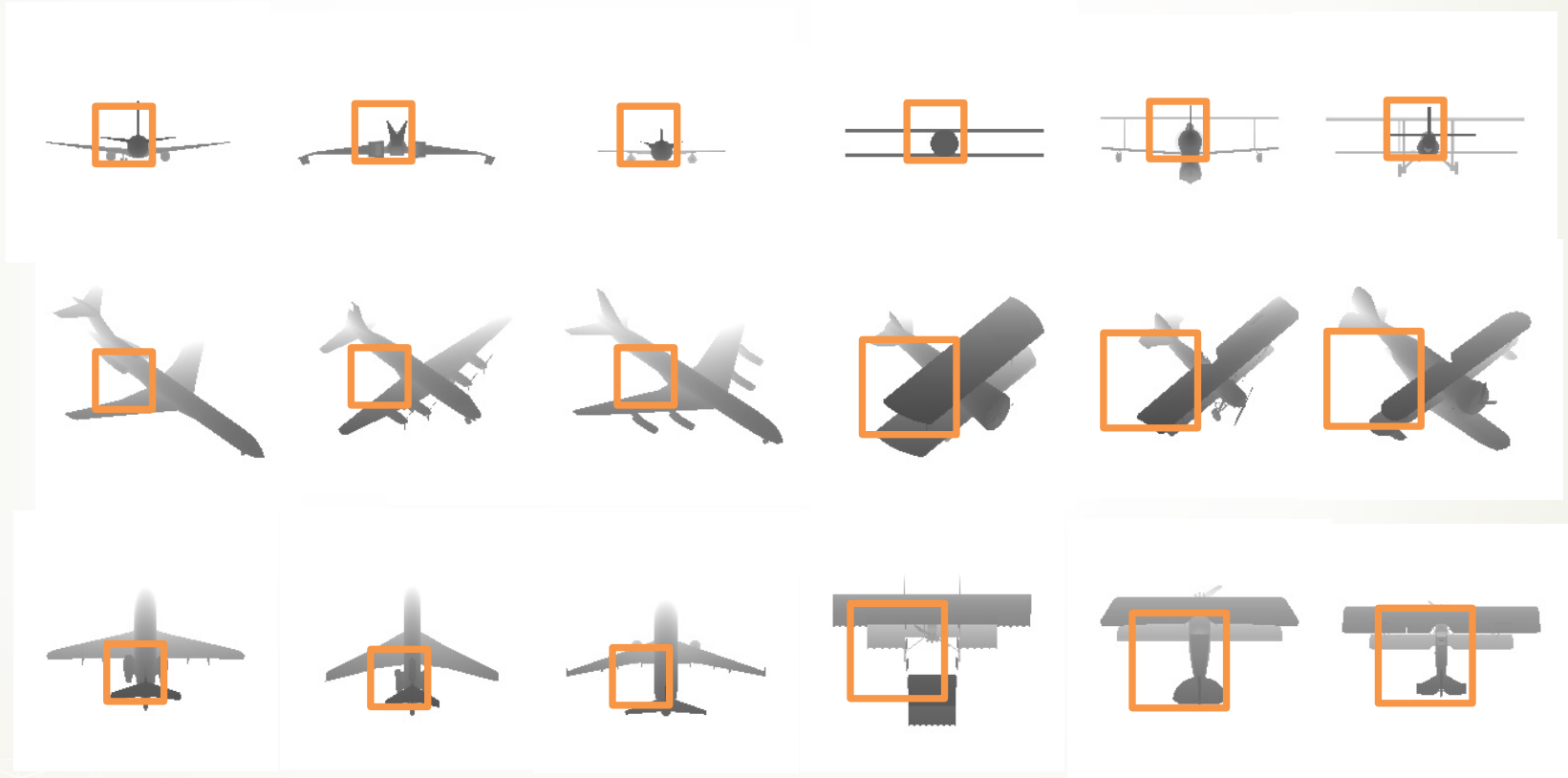




Attention extraction

One wing

Two wings





SIGGRAPH
ASIA 2016
MACAO

Results and evaluation

SA2016.SIGGRAPH.ORG

CONFERENCE: 5 - 8 DECEMBER 2016 • EXHIBITION: 6 - 8 DECEMBER 2016 • THE VENETIAN MACAO, MACAO

Sponsored by





SIGGRAPH
ASIA 2016
MACAO

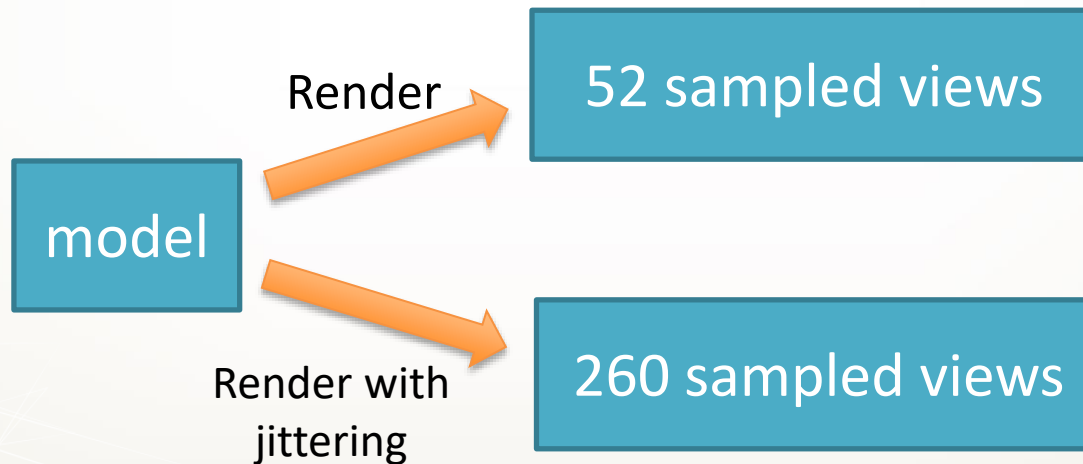
Database



57,452 models
57 categories



12,311 models
40 categories

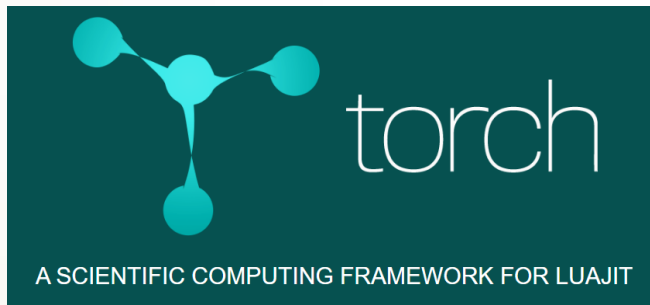




SIGGRAPH
ASIA 2016
MACAO

Timing

Database	MV-RNN train	MV-RNN test
ShapeNet	49 hr.	0.1 sec.
ModelNet40	22 hr.	0.1 sec.

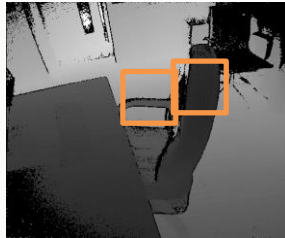
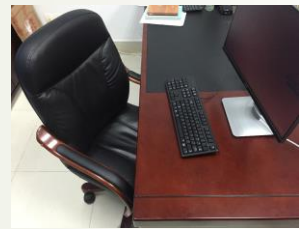




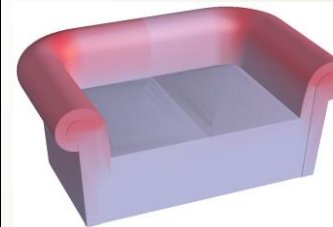
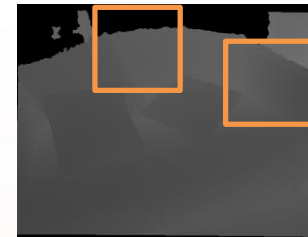
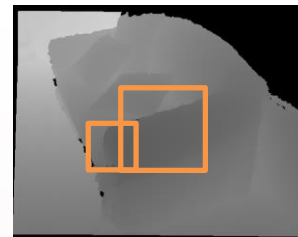
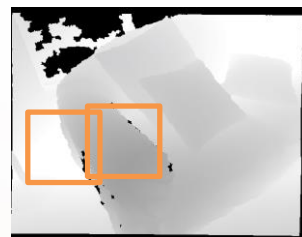
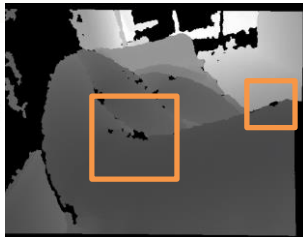
SIGGRAPH
ASIA 2016
MACAO

Visualization of attentions

Part-level attention



View sequence



View sequence



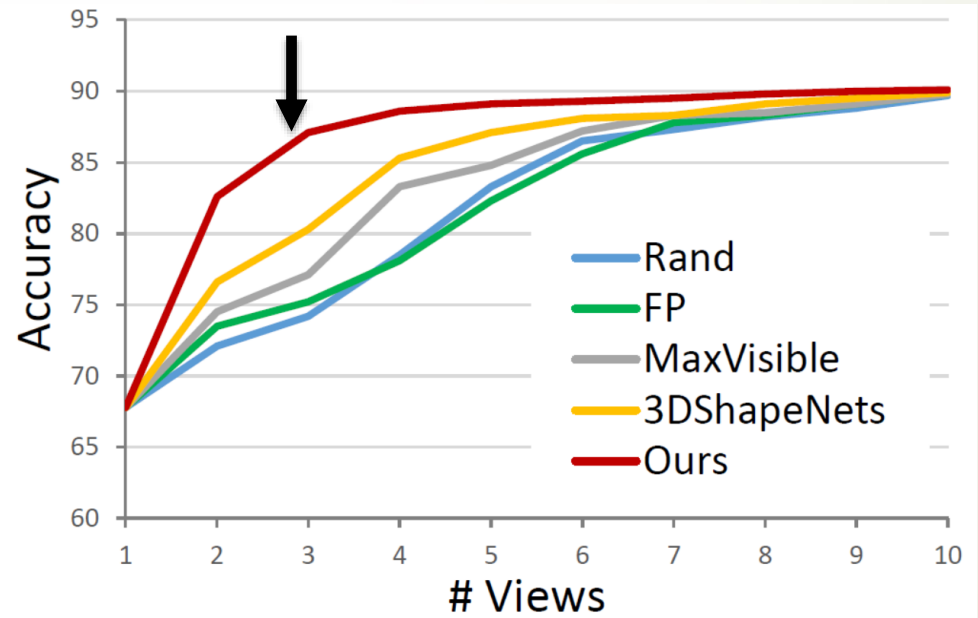
SIGGRAPH
ASIA 2016
MACAO

NBV estimation



PRINCETON
MODELNET

40 classes



Classification Accuracy

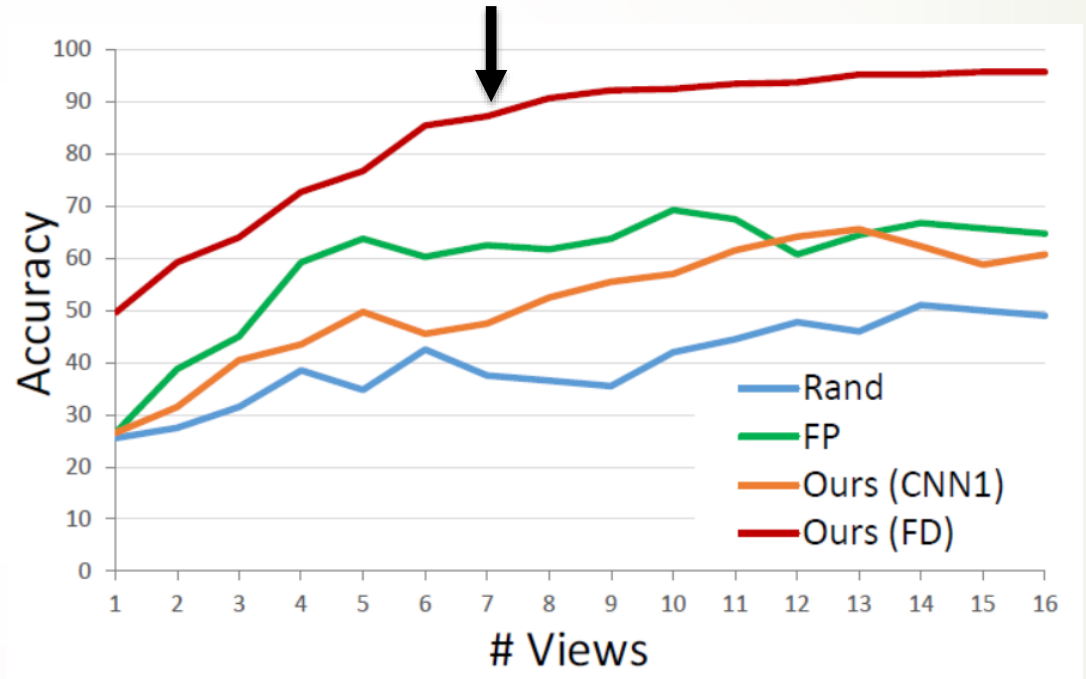


SIGGRAPH
ASIA 2016
MACAO

NBV estimation under occlusion



⋮



Classification Accuracy

Results on real scenes





SIGGRAPH
ASIA 2016
MACAO

Results on real scenes





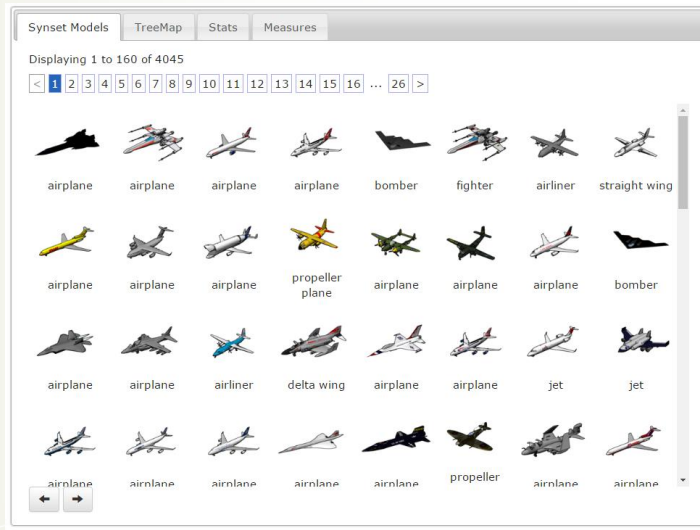
SIGGRAPH
ASIA 2016
MACAO

Results on real scenes



Limitations

- Recognizable objects
- No contextual information





SIGGRAPH
ASIA 2016
MACAO

Future works: Multi-modal recognition

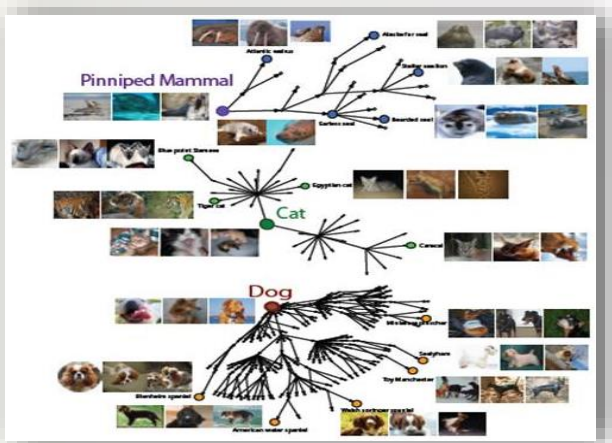
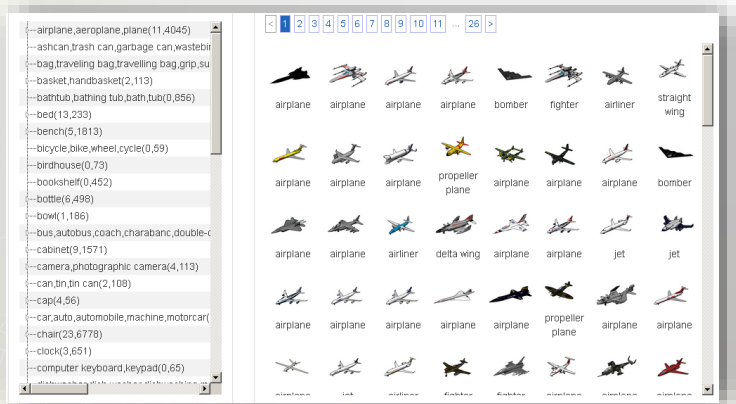
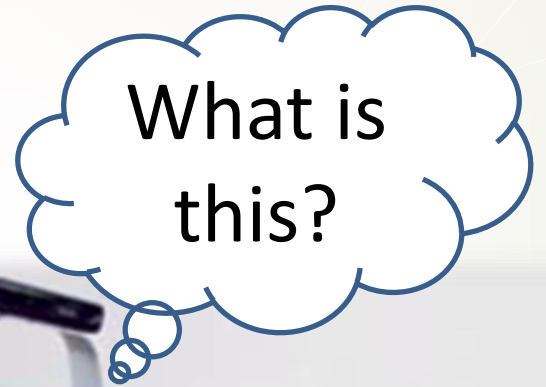


Image database



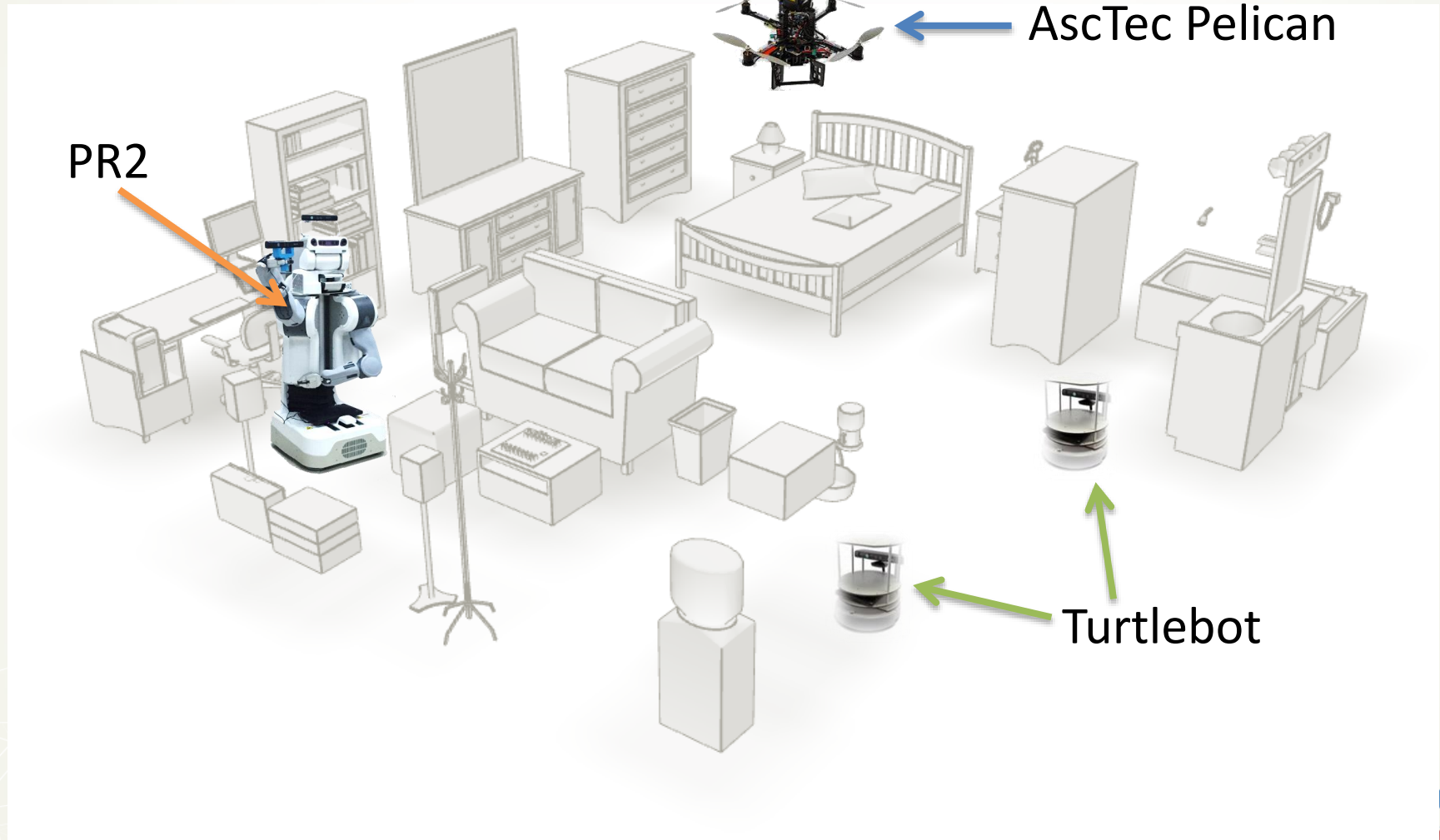
Shape database





SIGGRAPH
ASIA 2016
MACAO

Future: Multi-robot scene reconstruction & understanding

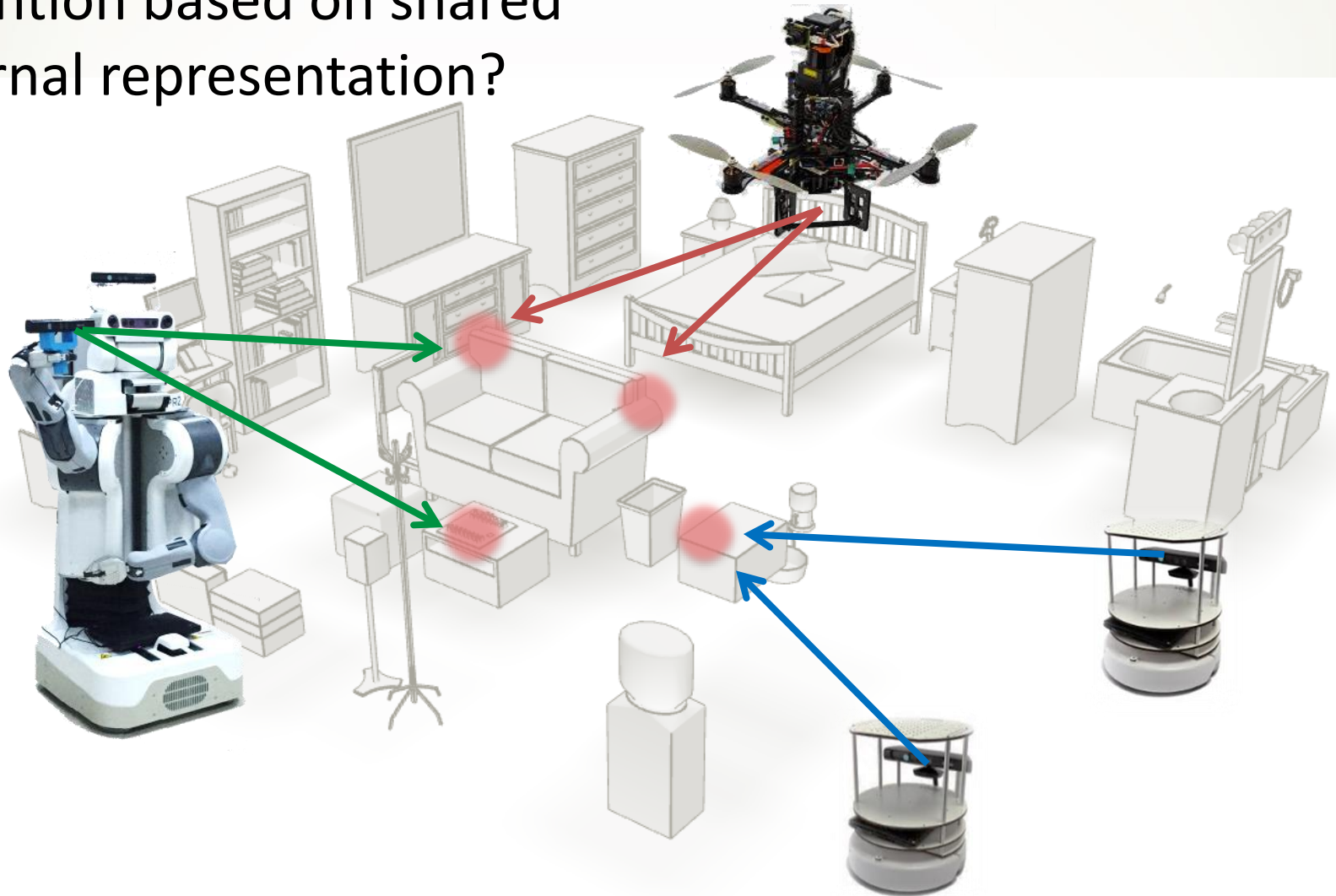




SIGGRAPH
ASIA 2016
MACAO

Future: Multi-robot attention model

Attention based on shared internal representation?



Thank you

Q & A



More details: kevinkaixu.net & yifeishi.net

SA2016.SIGGRAPH.ORG