

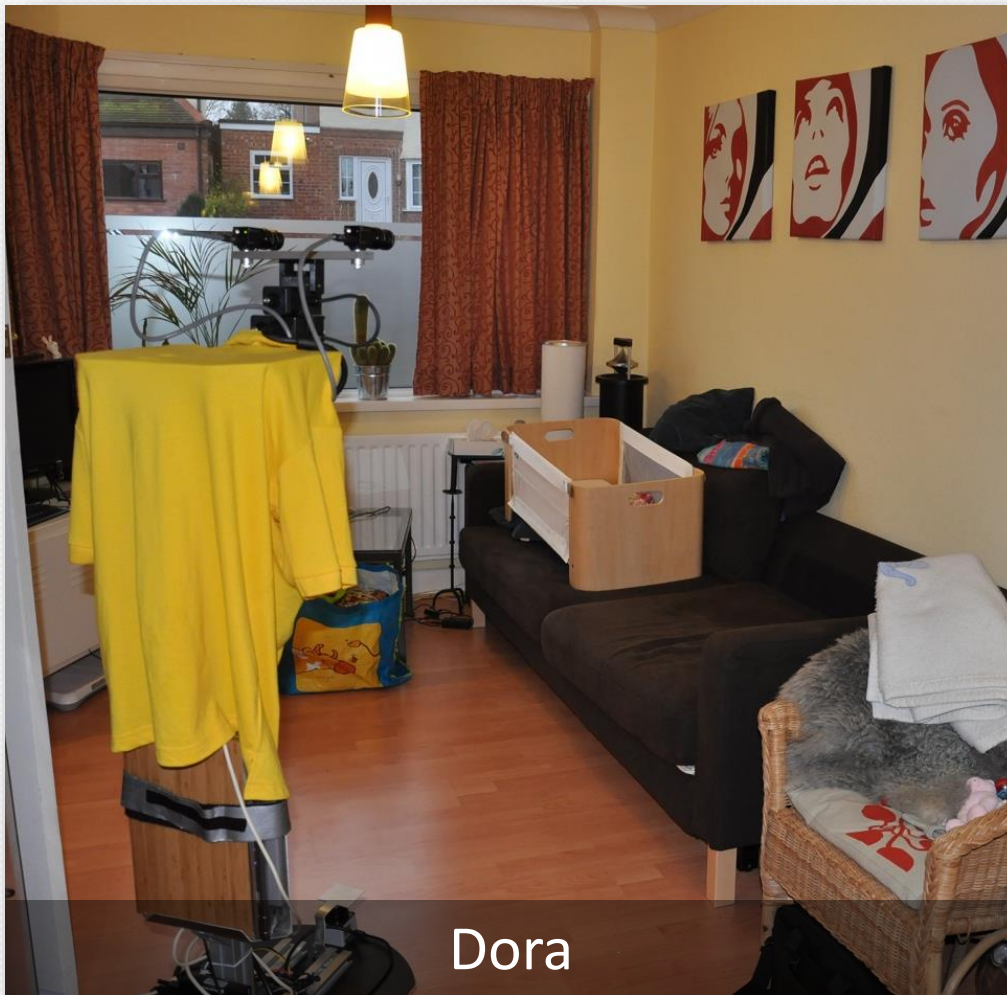
FROM SEMANTIC WORLD UNDERSTANDING TO COLLABORATION WITH DEEP REPRESENTATIONS

Andrzej PRONOBIS

University of Washington
KTH Royal Institute of Technology

www.pronobis.pro

SCENARIOS AND SYSTEMS



Dora

[Hanheide et al., AI'17]



InfoBot

[Chung*, Pronobis* et al.,
IROS'15, IROS'16]

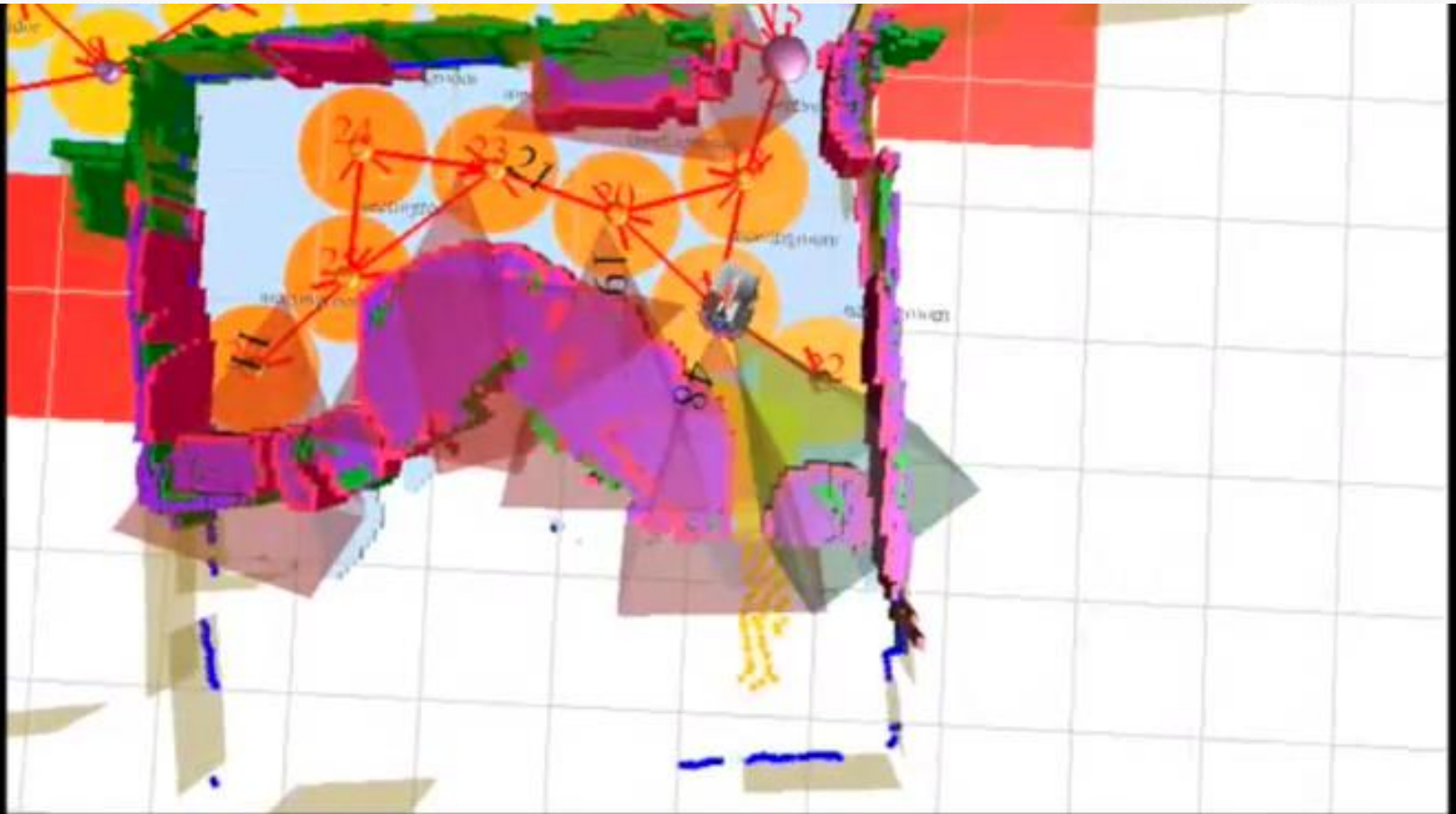
DORA: ACTIVE PLACE/OBJECT SEARCH [Hanheide et al., AI'17]

A photograph of a kitchen interior. In the foreground, a bright yellow cloth hangs from a metal stand. To the left is a dark wooden cabinet with glassware. In the background, there is a window with orange curtains, a white countertop with various items, and a blue patterned bag on the floor. A hanging light fixture is visible above the window.

Dora, find me some cornflakes!

DORA: APPROACH

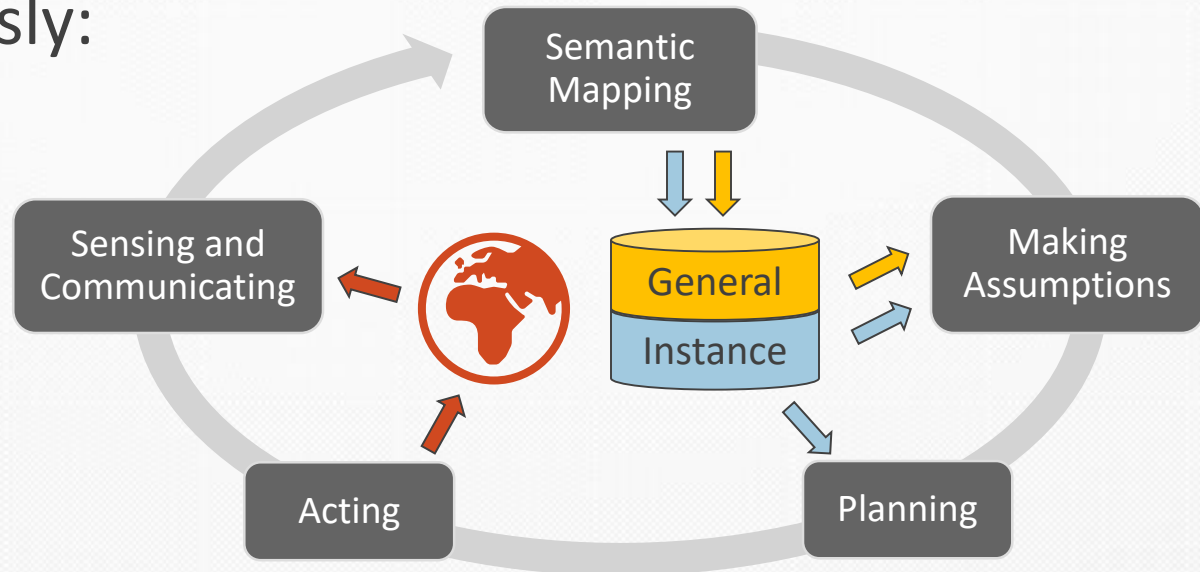
[Hanheide et al., AI'17]



of dialogue and sensing actions

DORA: CONTINUAL PLANNING WITH SEMANTICS

- Continual planning paradigm
- Semantic world knowledge is key
 - Instance – about current environment
 - General/common-sense – about human environments
- Continuously:



- Trades exploration vs. exploitation in a principled way

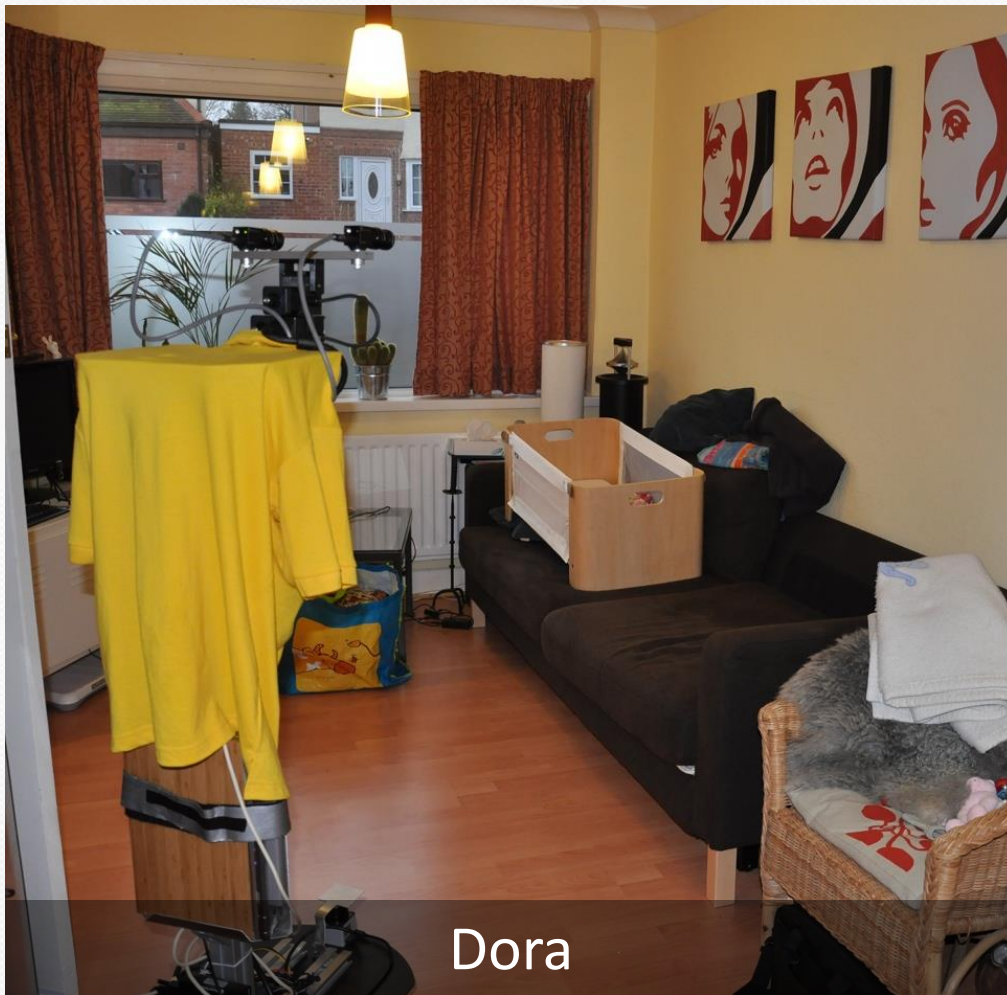
DORA: LANGUAGE-BASED INTERACTION [Hanheide et al., AI'17]

- Receiving goals in terms of semantic attributes
- Verbalizing explanations of failures
- Gathering unobservable instance knowledge
 - Polar and open questions
- Gathering common-sense knowledge

Robot: Starting plan execution.
Robot: let's see if there is someone here
Robot: hello human
Human: hello
Robot: ok
Robot: is there a container in this room ?
Human: yes
Robot: ok

◀ and general knowledge questions

SCENARIOS AND SYSTEMS



Dora

[Hanheide et al., AI'17]

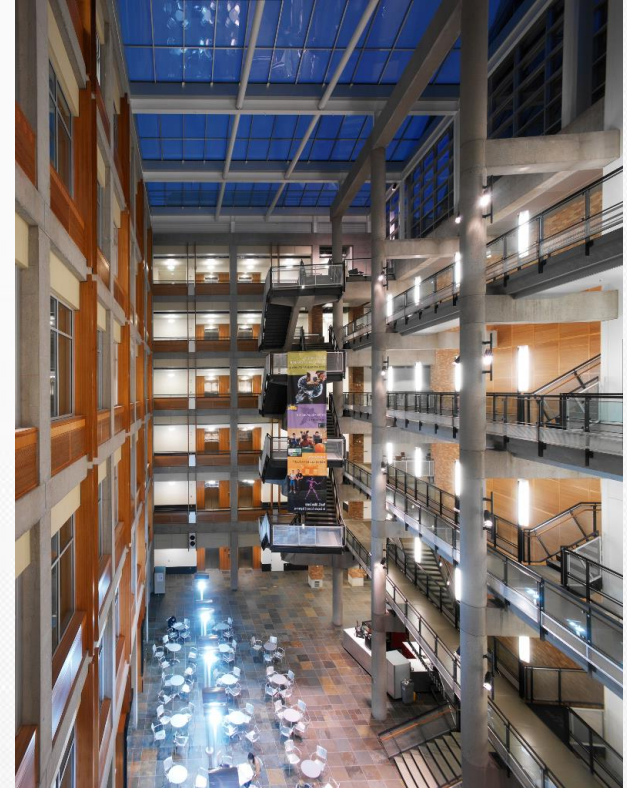


InfoBot

[Chung*, Pronobis* et al.,
IROS'15, IROS'16]

INFOBOT: CONCEPT

- Real-world environments are dynamic
- Humans in large environments spend time collecting information about current state:
- New type of service robots
 - Assistants answering questions about up-to-date state
 - Assist human co-inhabitants in daily tasks



End-to-end system with web interface

The screenshot displays the INFOBOT web interface. At the top, there are two tabs: "Your Questions" (active) and "Everyone's Questions". Below the tabs is a text input field labeled "Type question here". Underneath the input field are three controls: "Make Public" with a green toggle switch, "Email Notification" with a grey toggle switch, and "Deadline" with a text field showing "8:46 PM" and a calendar icon. A blue "Submit" button is to the right of these controls. Below the submission form, a confirmation message is shown: "Is there any free food in the lunchroom?" followed by "Received your question." and "In Queue". The message is attributed to "DUB-E Today at 8:39 am". At the bottom, there is a "Write a comment..." text field, a "Post" button, and a "Thank You DUB-E ❤️ 0" button.

Your Questions | Everyone's Questions

Type question here

Make Public ☒ Email Notification ☐ Deadline 8:46 PM **Submit**

██████████ Today at 8:39 am **Cancel**

private email deadline today at 9:30 am

Is there any free food in the lunchroom?

DUB-E Today at 8:39 am
Received your question.


In Queue

Write a comment... **Post** Thank You DUB-E ❤️ 0

End-to-end system with web interface

Your Questions | Everyone's Questions

Type question here

Make Public ☒ **Email Notification** ☐ **Deadline** 8:46 PM 

Submit

██████████ Today at 8:39 am Cancel

private email deadline today at 9:30 am

Is there any free food in the lunchroom?

DUB-E Today at 8:39 am
Working on your question!

Running

Write a comment... Post Thank You DUB-E ❤️ 0



End-to-end system with web interface

Your Questions | Everyone's Questions

Type question here

Make Public ☒ Email Notification ☐ Deadline

██████████ Today at 8:39 am

Cancel

private email deadline today at 9:30 am

Is there any free food in the lunchroom?

DUB-E Today at 8:39 am

Working on your question!

Running

Write a comment... Thank You DUB-E ❤️ 0



End-to-end system with web interface

Your Questions | Everyone's Questions

Type question here

Make Public ☒ Email Notification ☐ Deadline

Today at 8:39 am

Cancel

private email deadline today at 9:30 am

Is there any free food in the lunchroom?

Yes

Success :)

Write a comment... Post Thank You DUB-E ❤️ 0



- Goal: practical usage, typical questions
- Comprehensive survey (2 buildings, 111 responses)
- Wizard-of-Oz deployment (4 days, 45 unique users)
 - Users use the web interface
 - Operator tele-operates the robot and posts answers

- Goal: practical usage, typical questions
- Comprehensive survey (2 buildings, 111 responses)
- Wizard-of-Oz deployment (4 days, 45 unique users)
 - Users use the web interface
 - Operator tele-operates the robot and posts answers
- Questions:

Is there anyone in {location}?

Is {person} in his/her office?

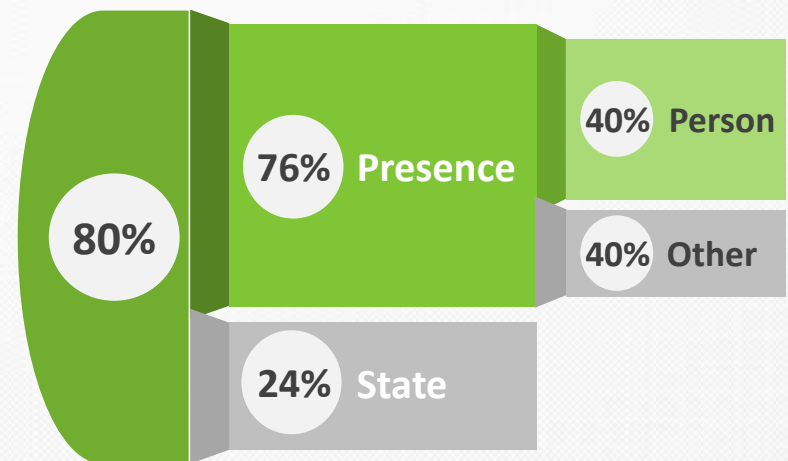
Is there any food in the downstairs kitchen?

Is there anything in my mailbox?

Is the door to the conference room open?

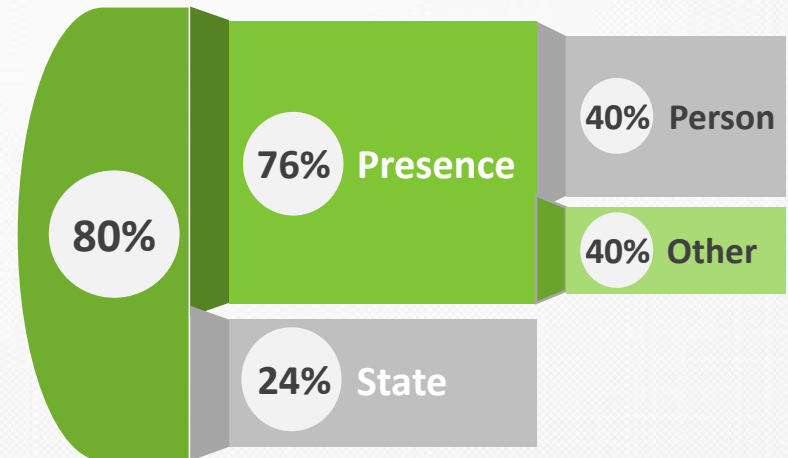
Is the reception still open?

How noisy is it in the atrium right now?



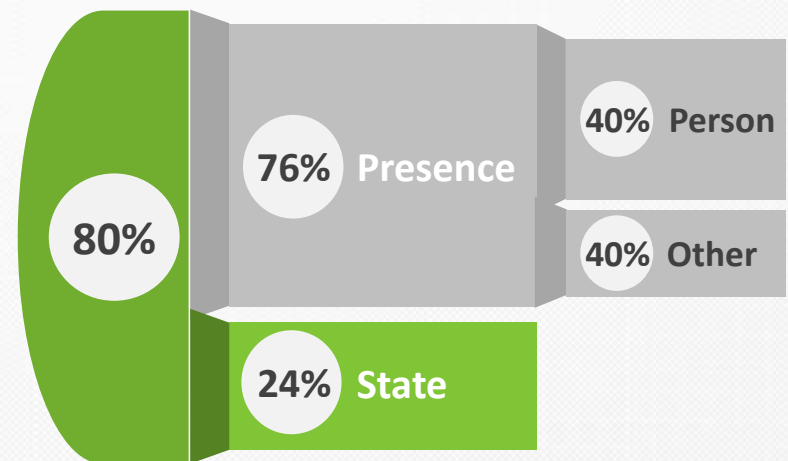
- Goal: practical usage, typical questions
- Comprehensive survey (2 buildings, 111 responses)
- Wizard-of-Oz deployment (4 days, 45 unique users)
 - Users use the web interface
 - Operator tele-operates the robot and posts answers
- Questions:

Is there anyone in {location}?
Is {person} in his/her office?
Is there any food in the downstairs kitchen?
Is there anything in my mailbox?
Is the door to the conference room open?
Is the reception still open?
How noisy is it in the atrium right now?



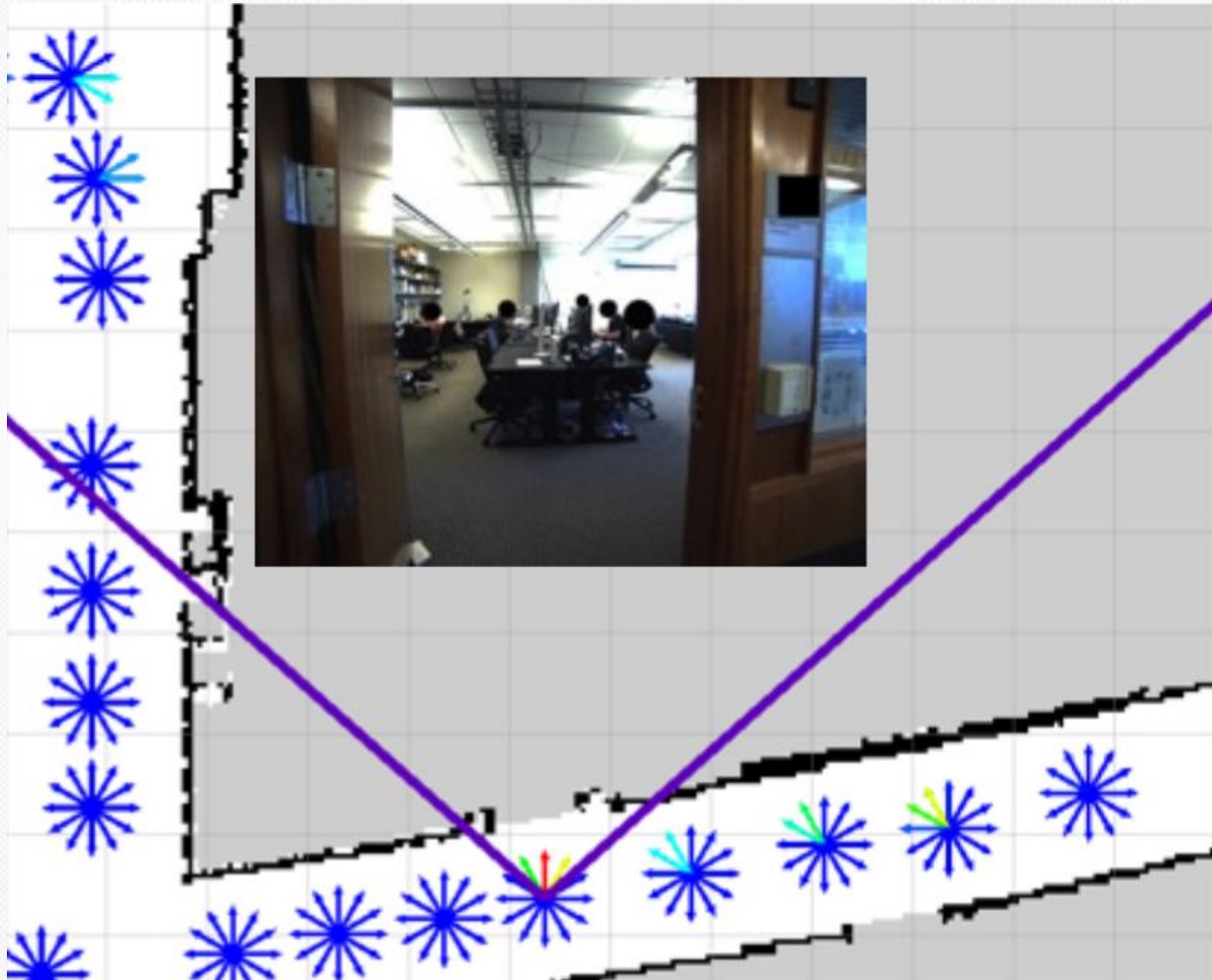
- Goal: practical usage, typical questions
- Comprehensive survey (2 buildings, 111 responses)
- Wizard-of-Oz deployment (4 days, 45 unique users)
 - Users use the web interface
 - Operator tele-operates the robot and posts answers
- Questions:

Is there anyone in {location}?
Is {person} in his/her office?
Is there any food in the downstairs kitchen?
Is there anything in my mailbox?
Is the door to the conference room open?
Is the reception still open?
How noisy is it in the atrium right now?



- User satisfaction: only 16% not satisfied with answers

Is there anyone in the mobile robotics lab?



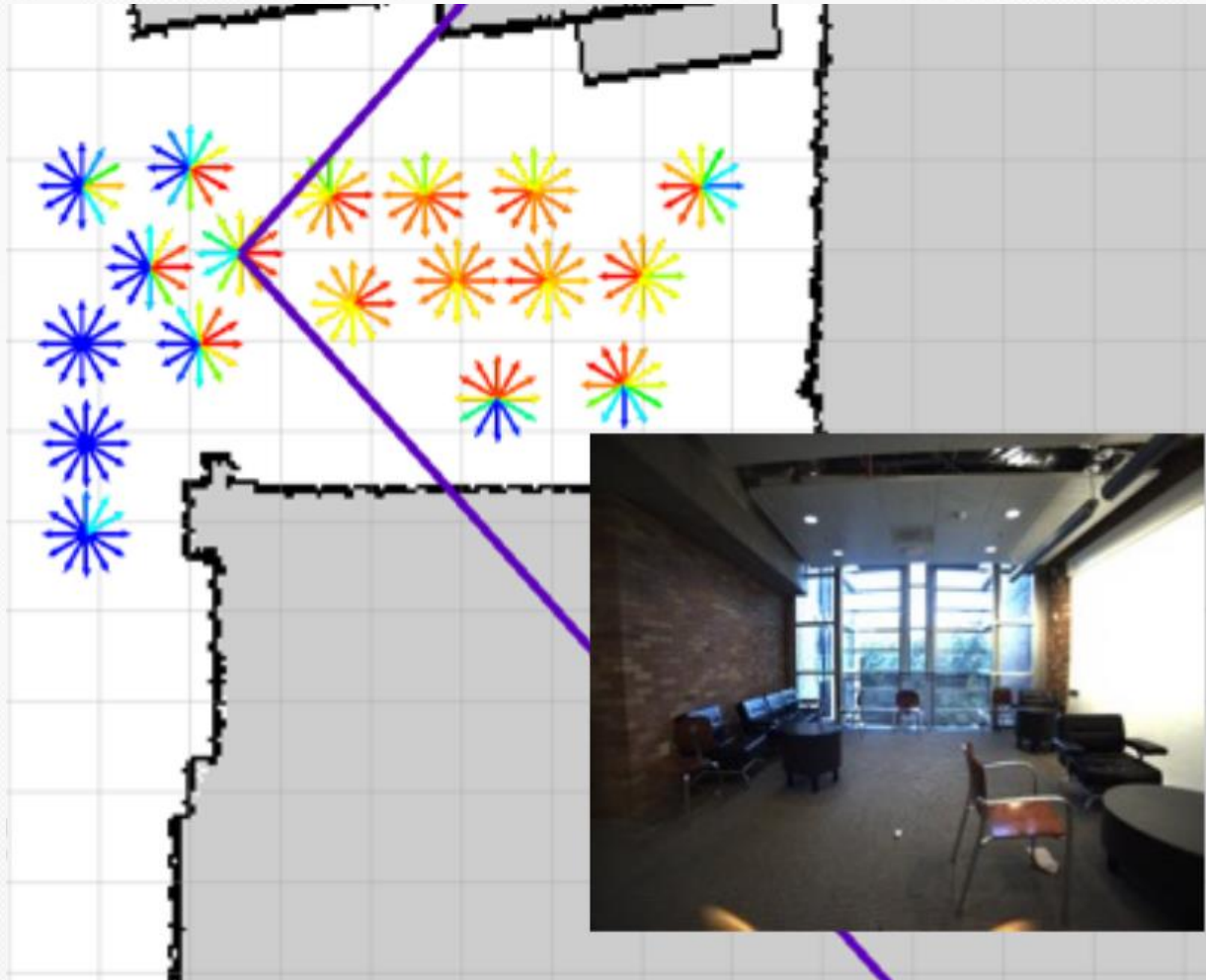
Ground truth: **Yes**

Is there anyone in the mobile robotics lab?



Ground truth: **No**

Is the breakout area occupied?

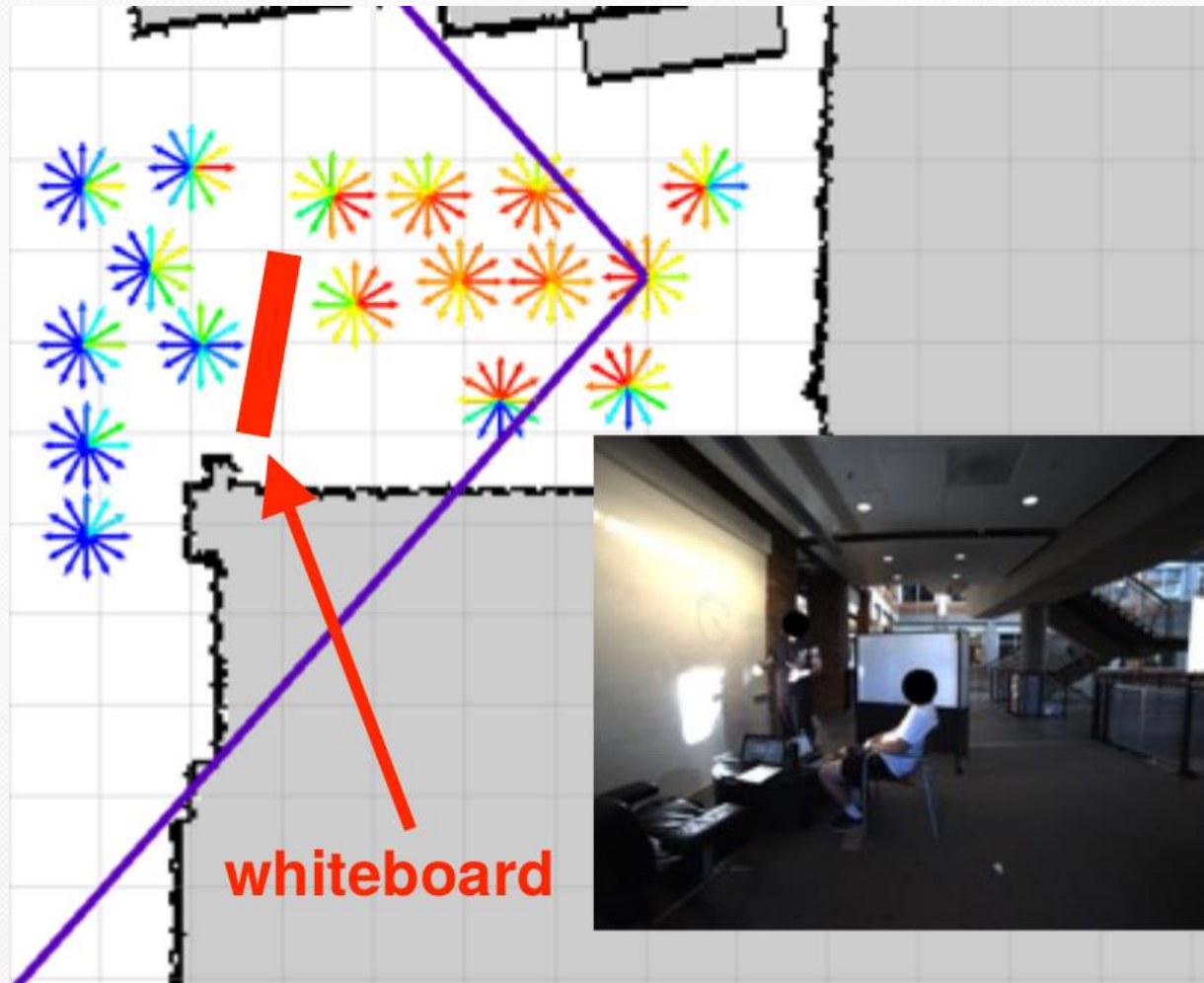


Ground truth: **No**

INFOBOT: AUTONOMOUS SYSTEM

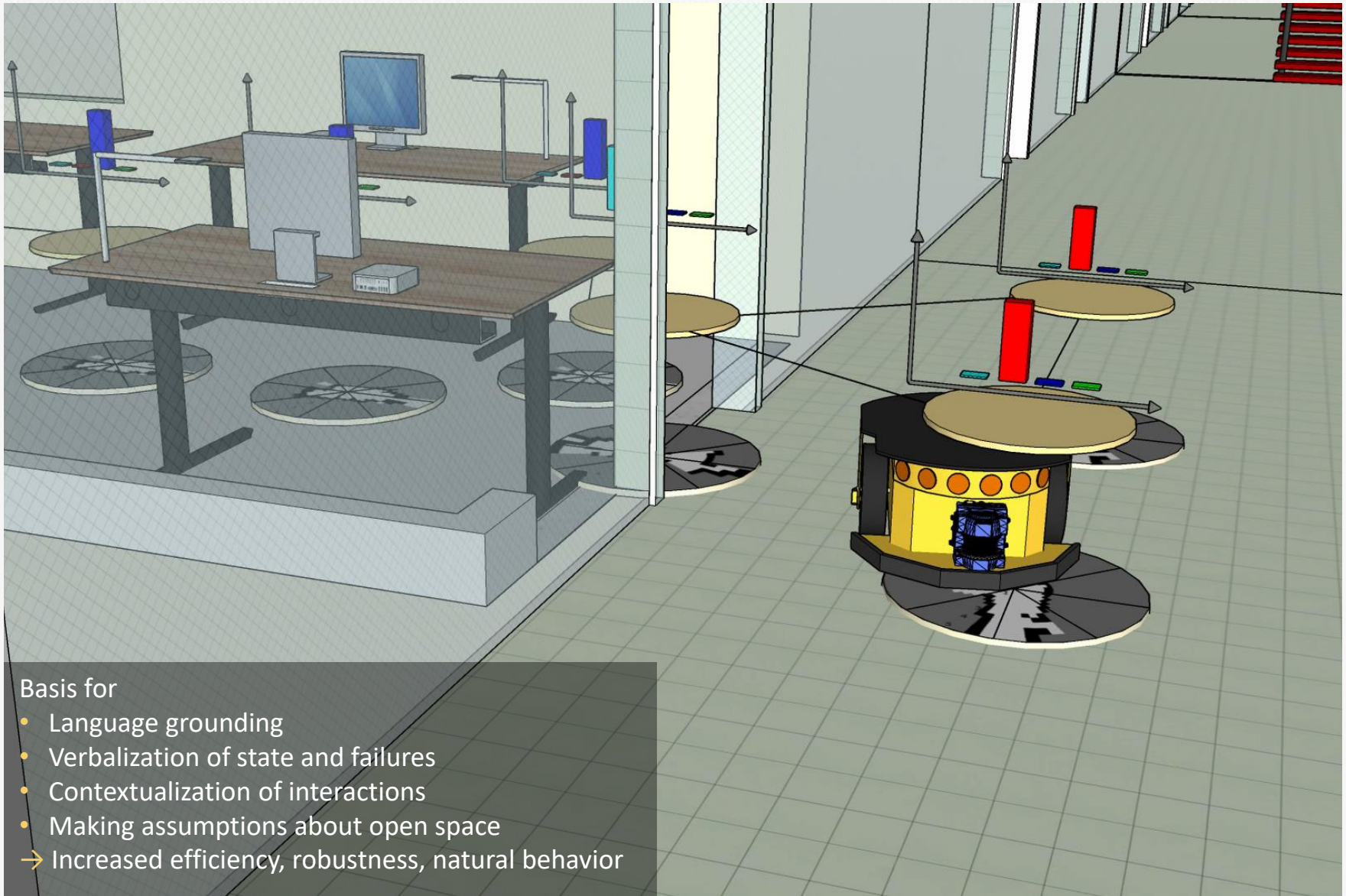
[Chung*, Pronobis*
et al., IROS'16]

Is the breakout area occupied?

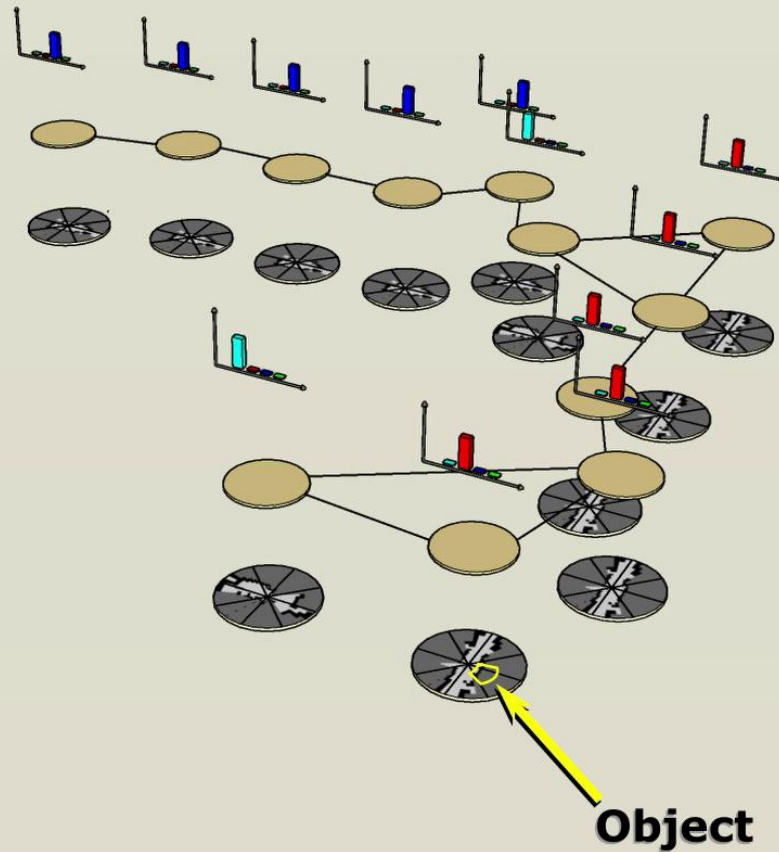


Ground truth: **Yes**

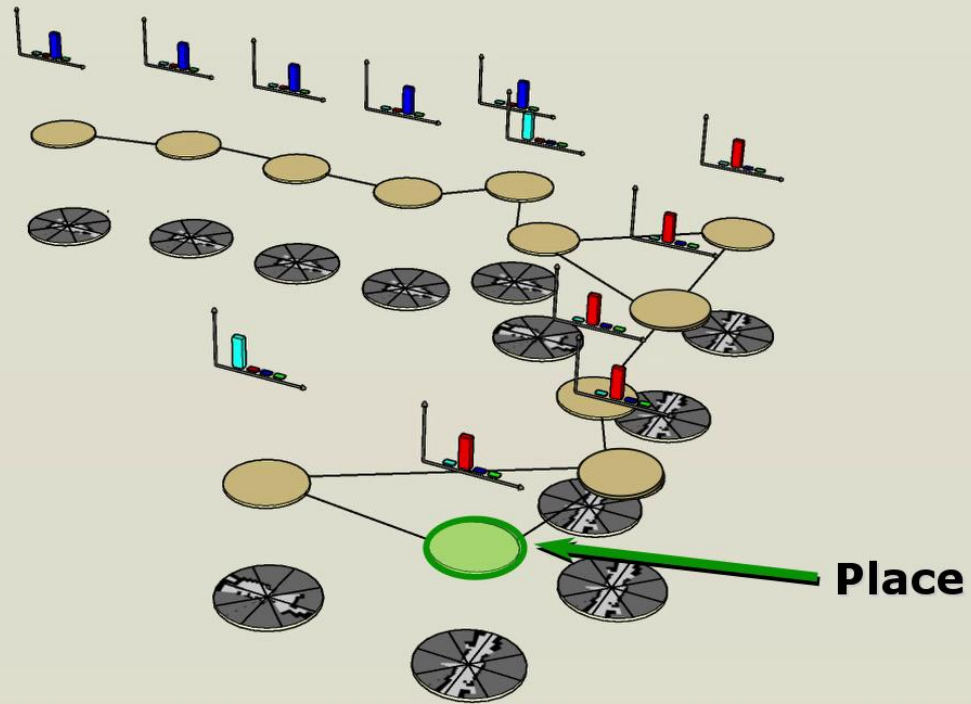
SPATIAL KNOWLEDGE REPRESENTATION



SPATIAL KNOWLEDGE REPRESENTATION

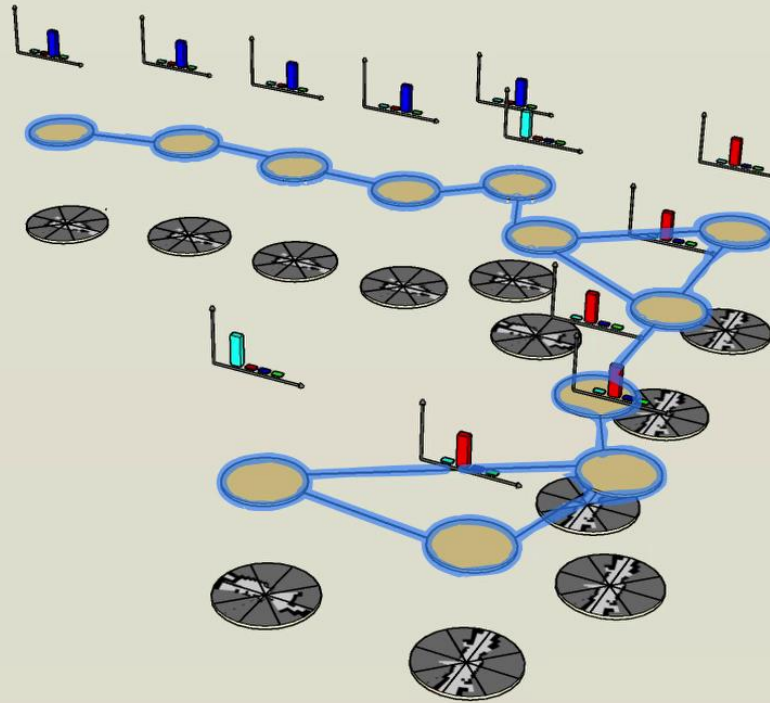


SPATIAL KNOWLEDGE REPRESENTATION

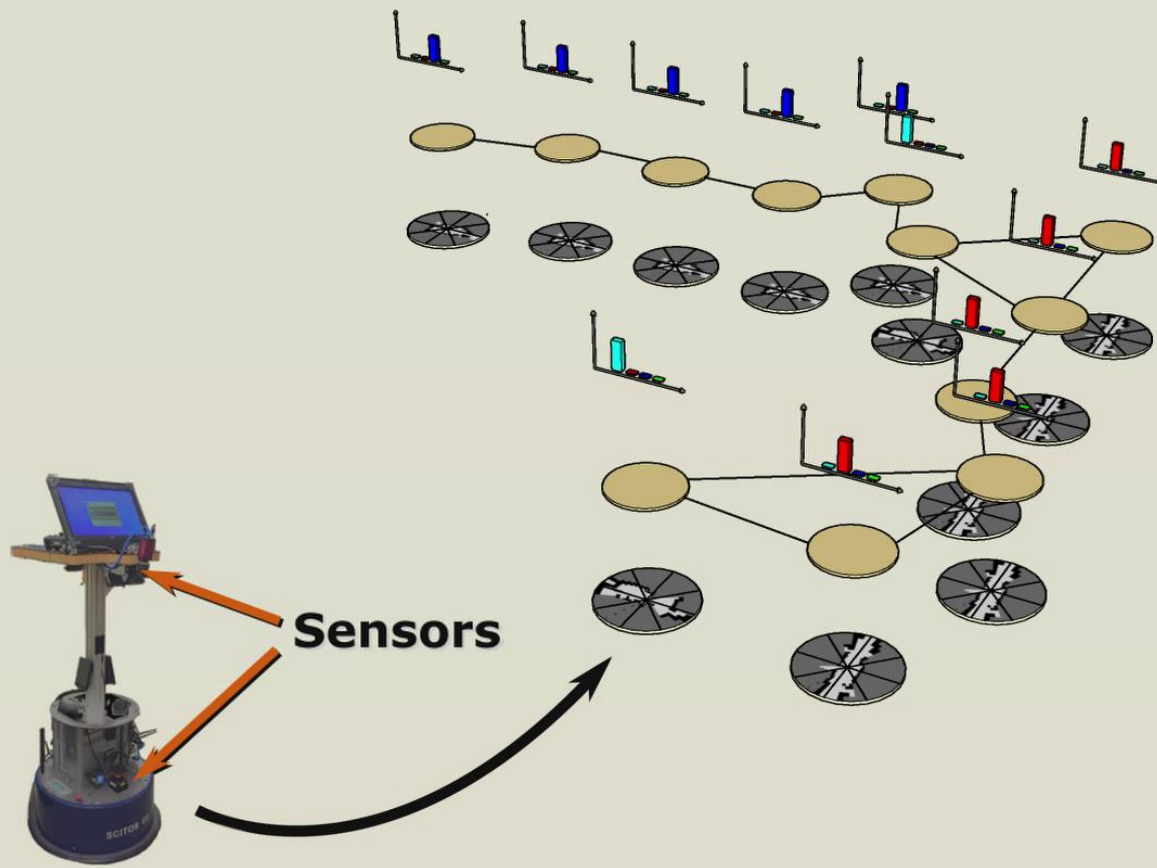


SPATIAL KNOWLEDGE REPRESENTATION

Topology

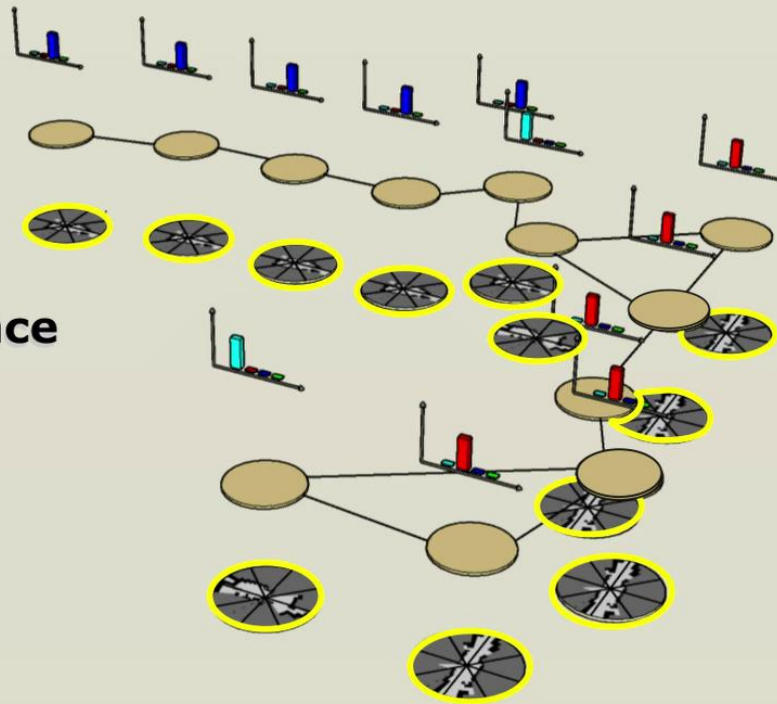


SPATIAL KNOWLEDGE REPRESENTATION



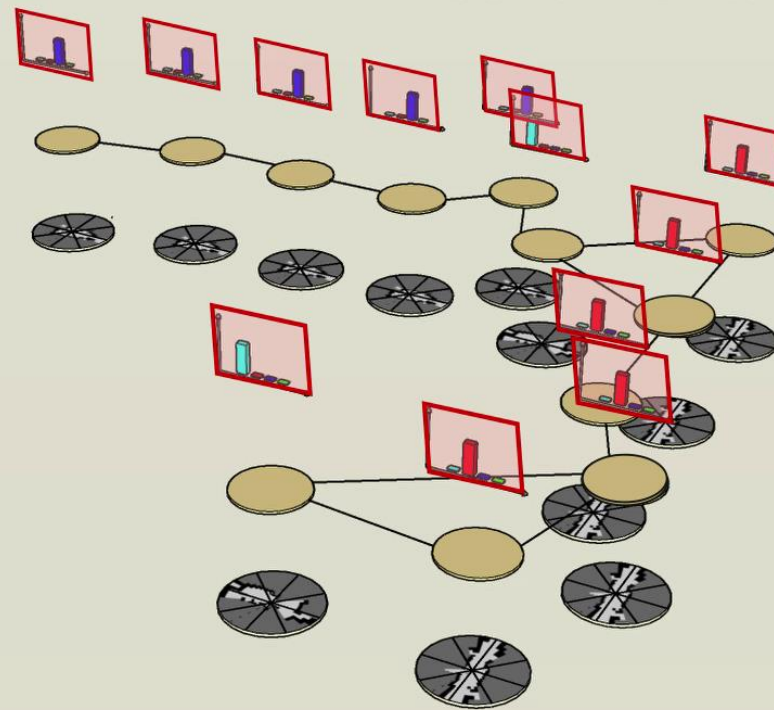
SPATIAL KNOWLEDGE REPRESENTATION

Place appearance

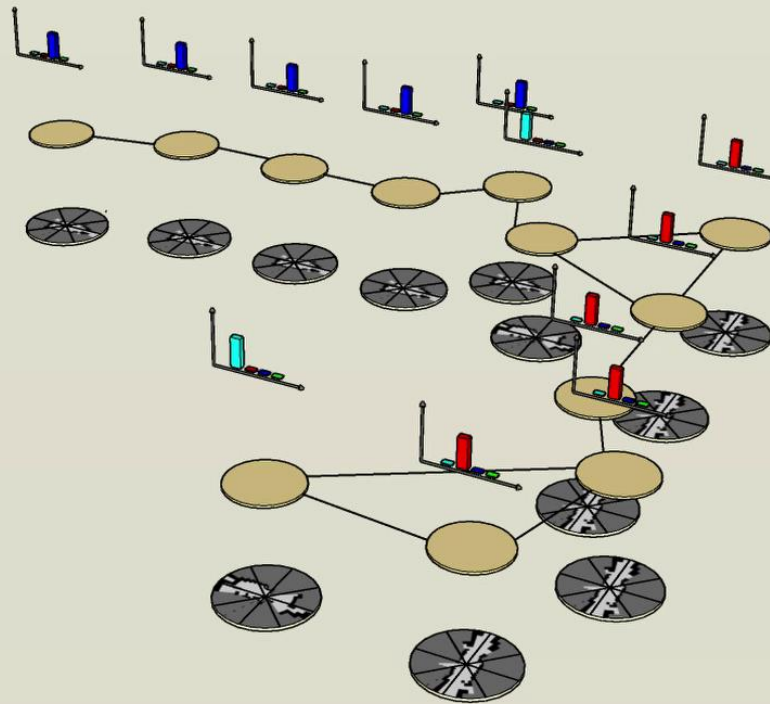


SPATIAL KNOWLEDGE REPRESENTATION

Semantic Attributes

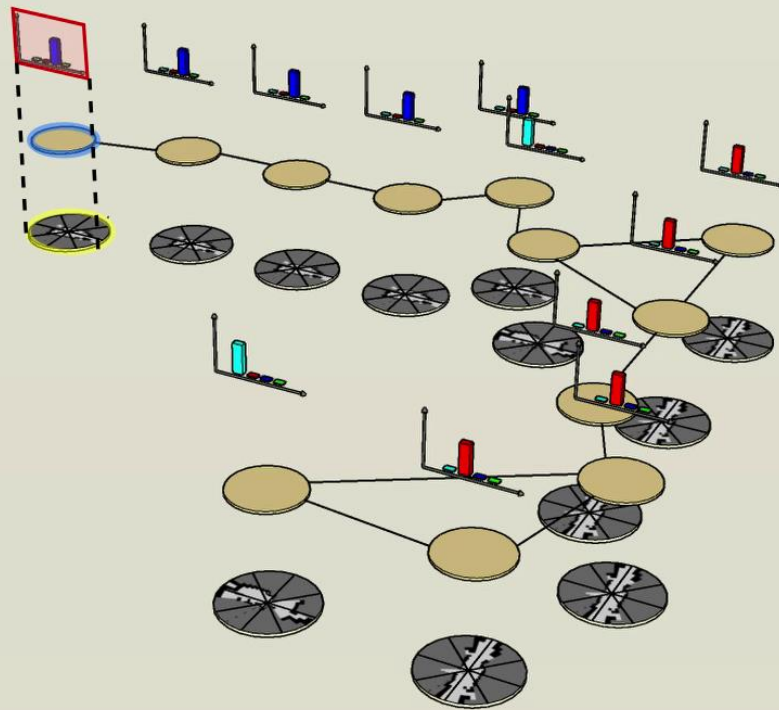


SPATIAL KNOWLEDGE REPRESENTATION

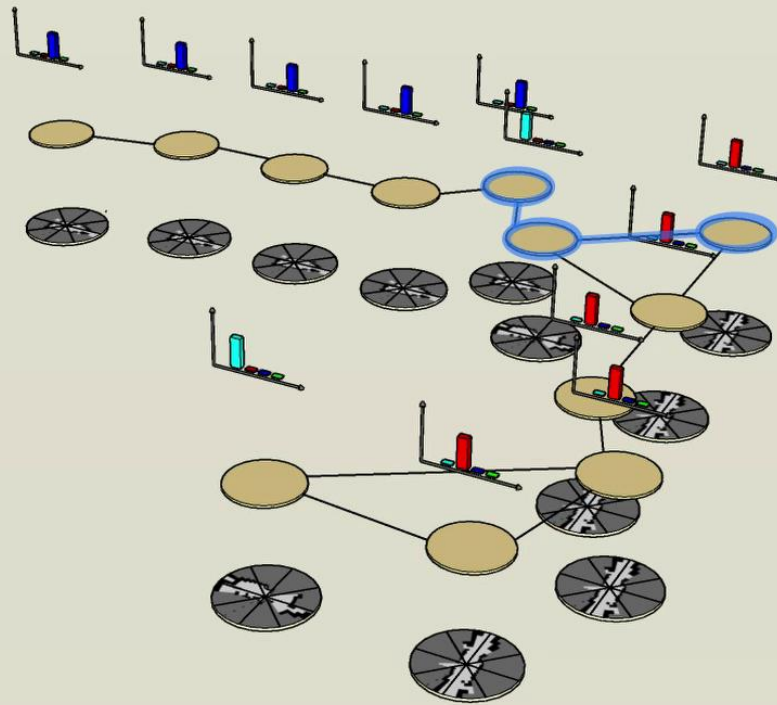


**Sensory observations:
local, partial, noisy**

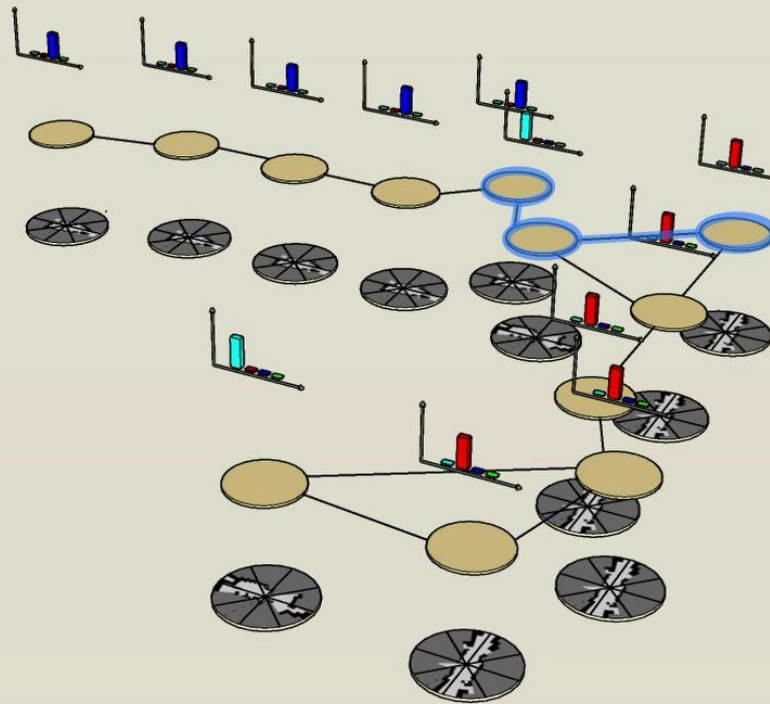
SPATIAL KNOWLEDGE REPRESENTATION



SPATIAL KNOWLEDGE REPRESENTATION



SPATIAL KNOWLEDGE REPRESENTATION

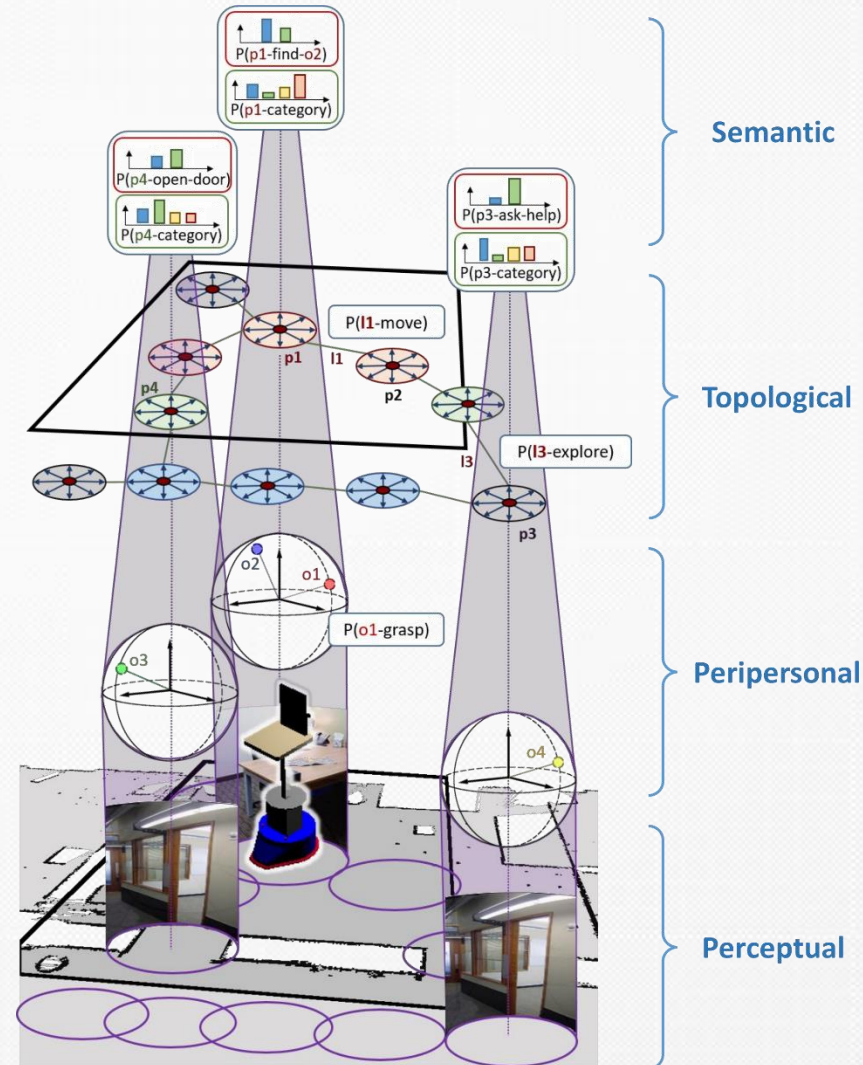


How to structure the body of spatial knowledge
for a deliberative agent?

DASH: DEEP SPATIAL AFFORDANCE HIERARCHY

[Pronobis, Riccio, Rao, RSS-SSRR'17]

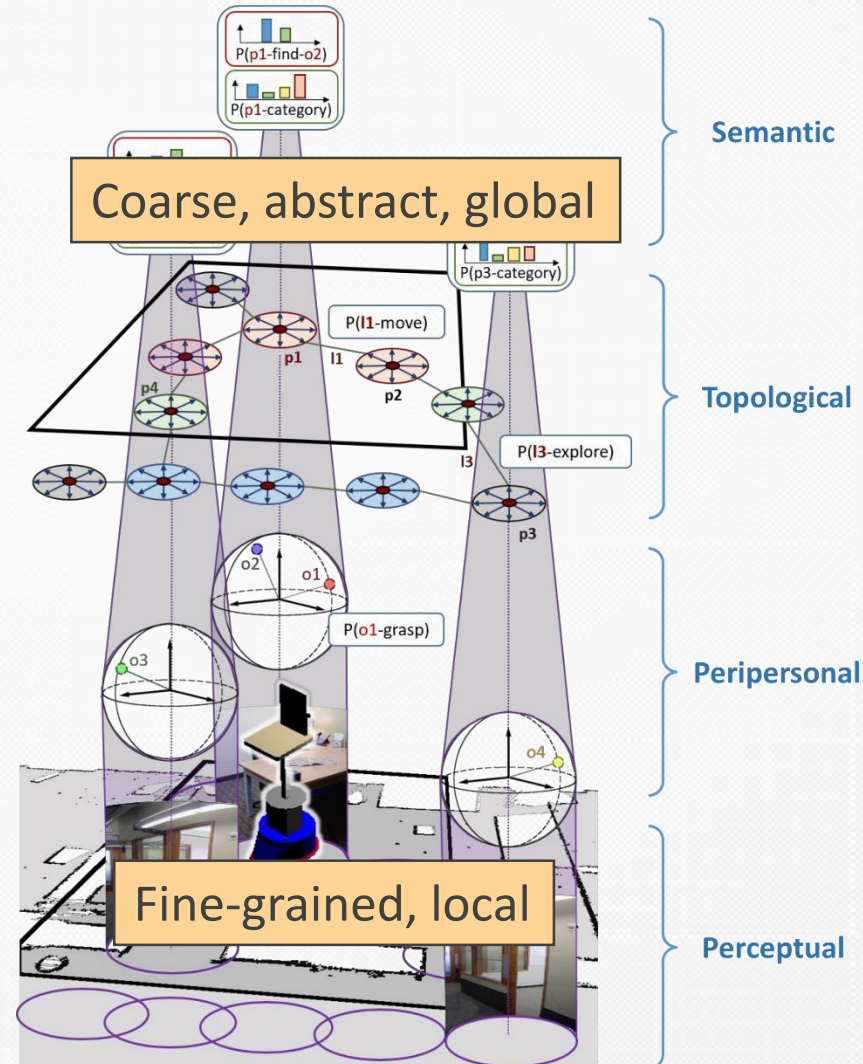
- Hierarchical, layered
- New representation and deep probabilistic model



DASH: DEEP SPATIAL AFFORDANCE HIERARCHY

[Pronobis, Riccio, Rao, RSS-SSRR'17]

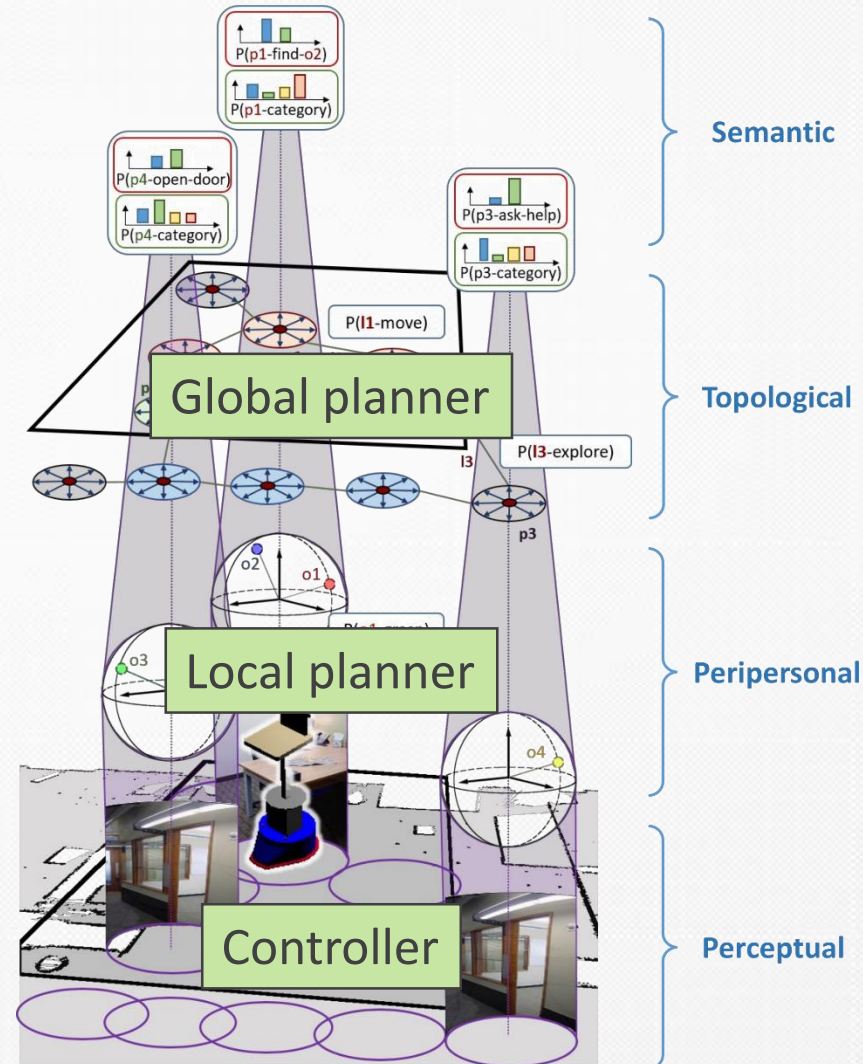
- Hierarchical, layered
- New representation and deep probabilistic model
- 4 layers, different:
 - Aspects of the world
 - Levels of abstraction
 - Spatial scales
 - Frames of reference



DASH: DEEP SPATIAL AFFORDANCE HIERARCHY

[Pronobis, Riccio, Rao, RSS-SSRR'17]

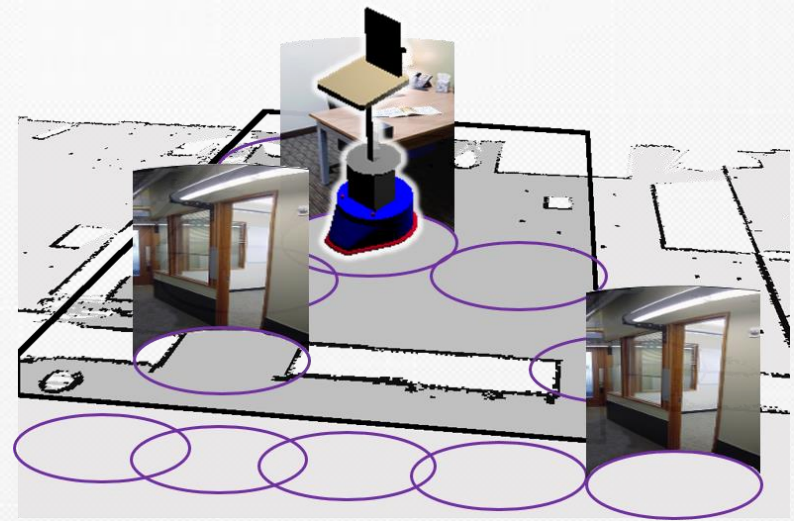
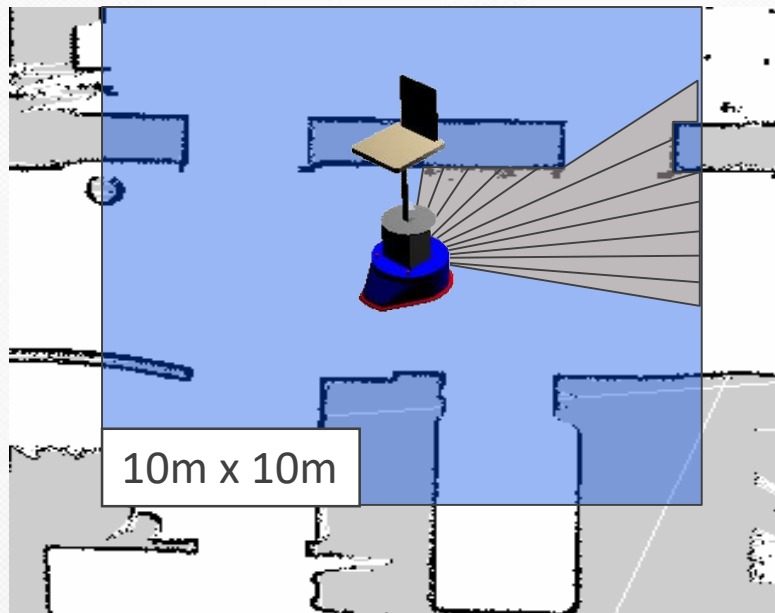
- Hierarchical, layered
- New representation and deep probabilistic model
- 4 layers, different:
 - Aspects of the world
 - Levels of abstraction
 - Spatial scales
 - Frames of reference



DASH: PERCEPTUAL LAYER

[Pronobis, Riccio, Rao, RSS-SSRR'17]

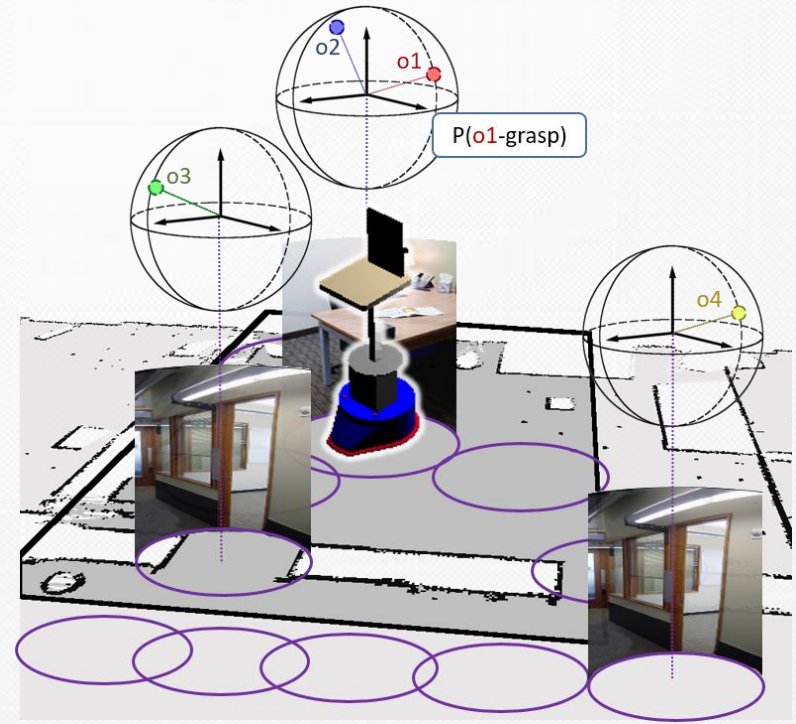
- Spatio-temporal integration of sensory data
- Accurate geometry and appearance
- Realized as sliding window following the robot



DASH: PERIPERSONAL LAYER

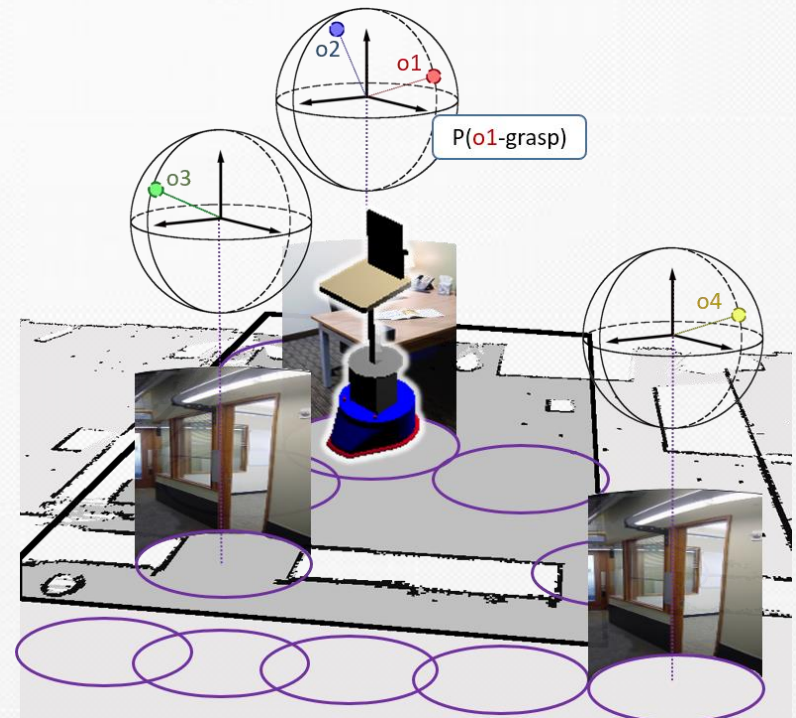
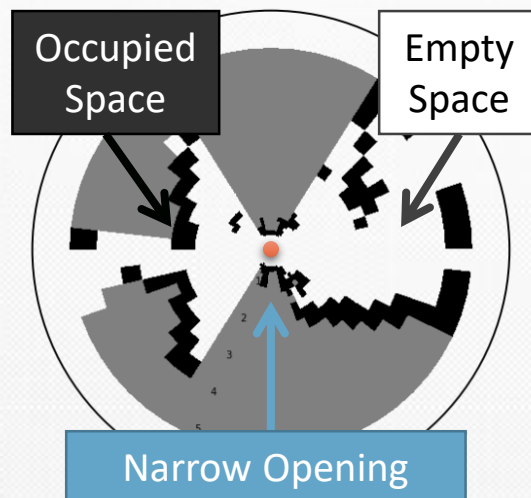
[Pronobis, Riccio, Rao, RSS-SSRR'17]

- Collection of ego-centric representations, each:
 - Models space immediately reachable or observable
 - From perspective of robot at a specific place
 - Updated when robot visits a place
- Realized using collection of polar occupancy grids



- Collection of ego-centric representations, each:
 - Models space immediately reachable or observable
 - From perspective of robot at a specific place
 - Updated when robot visits a place
- Realized using collection of polar occupancy grids

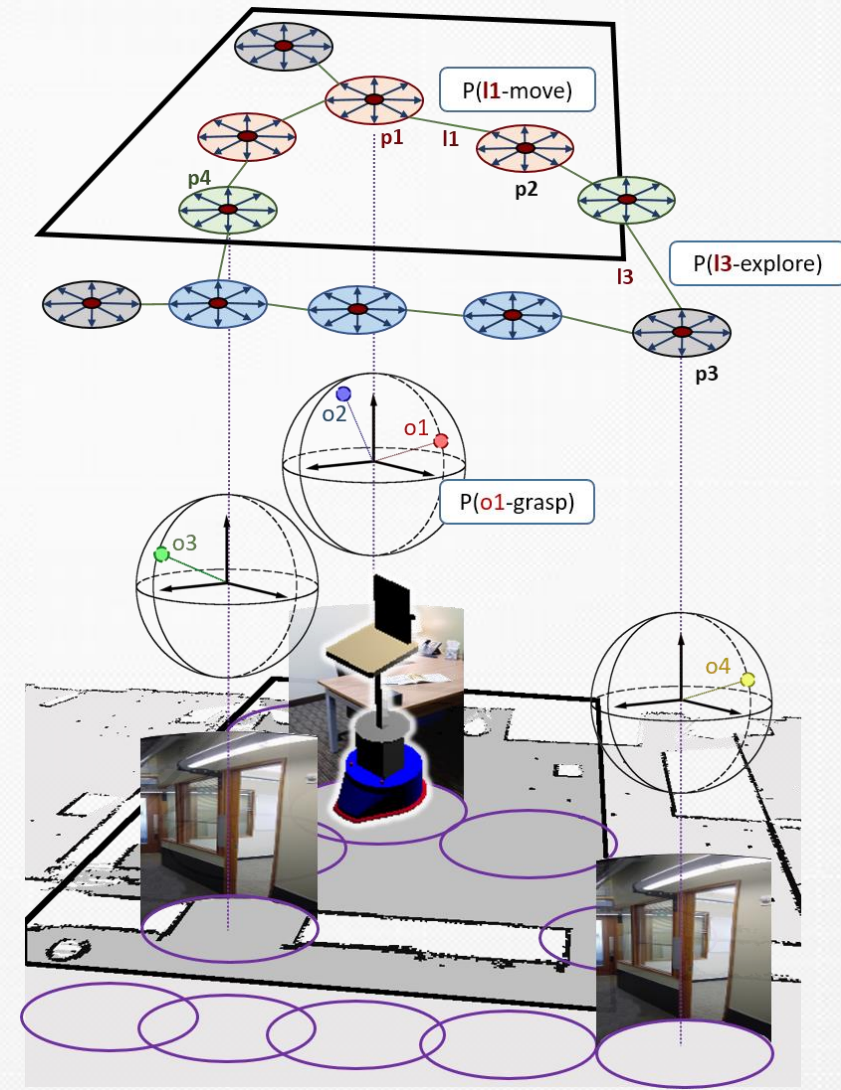
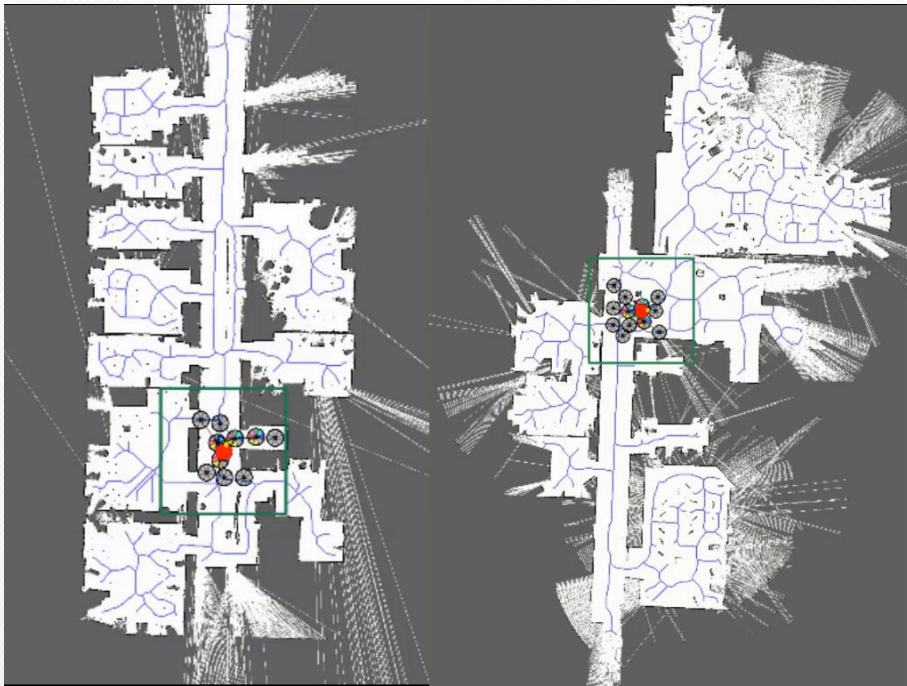
For a doorway:



DASH: TOPOLOGICAL LAYER

[Pronobis, Riccio, Rao, RSS-SSRR'17]

- Efficient representation of large-scale space
 - Coarse global geometry
 - Topology

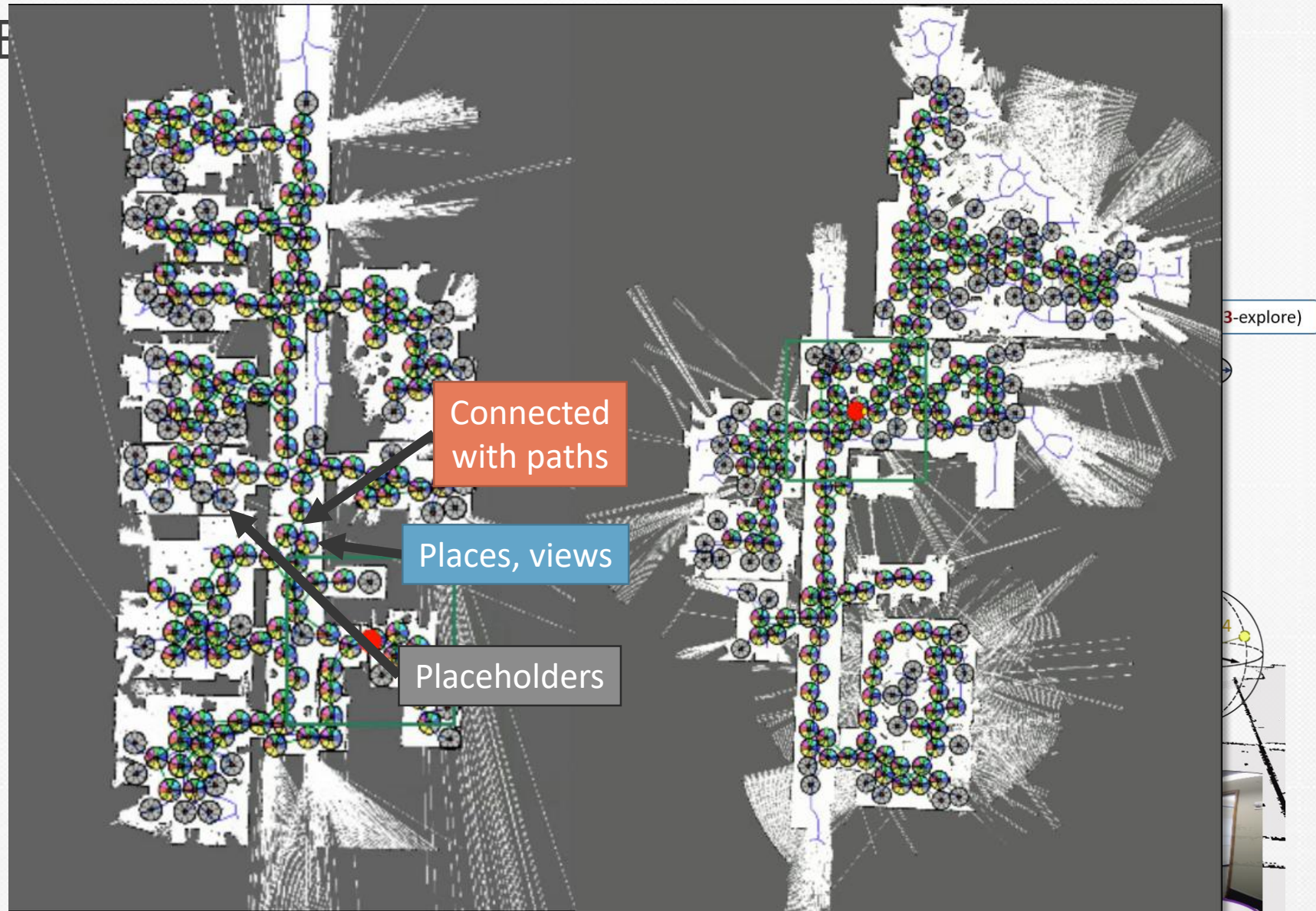


DASH: TOPOLOGICAL LAYER

[Pronobis, Riccio, Rao, RSS-SSRR'17]

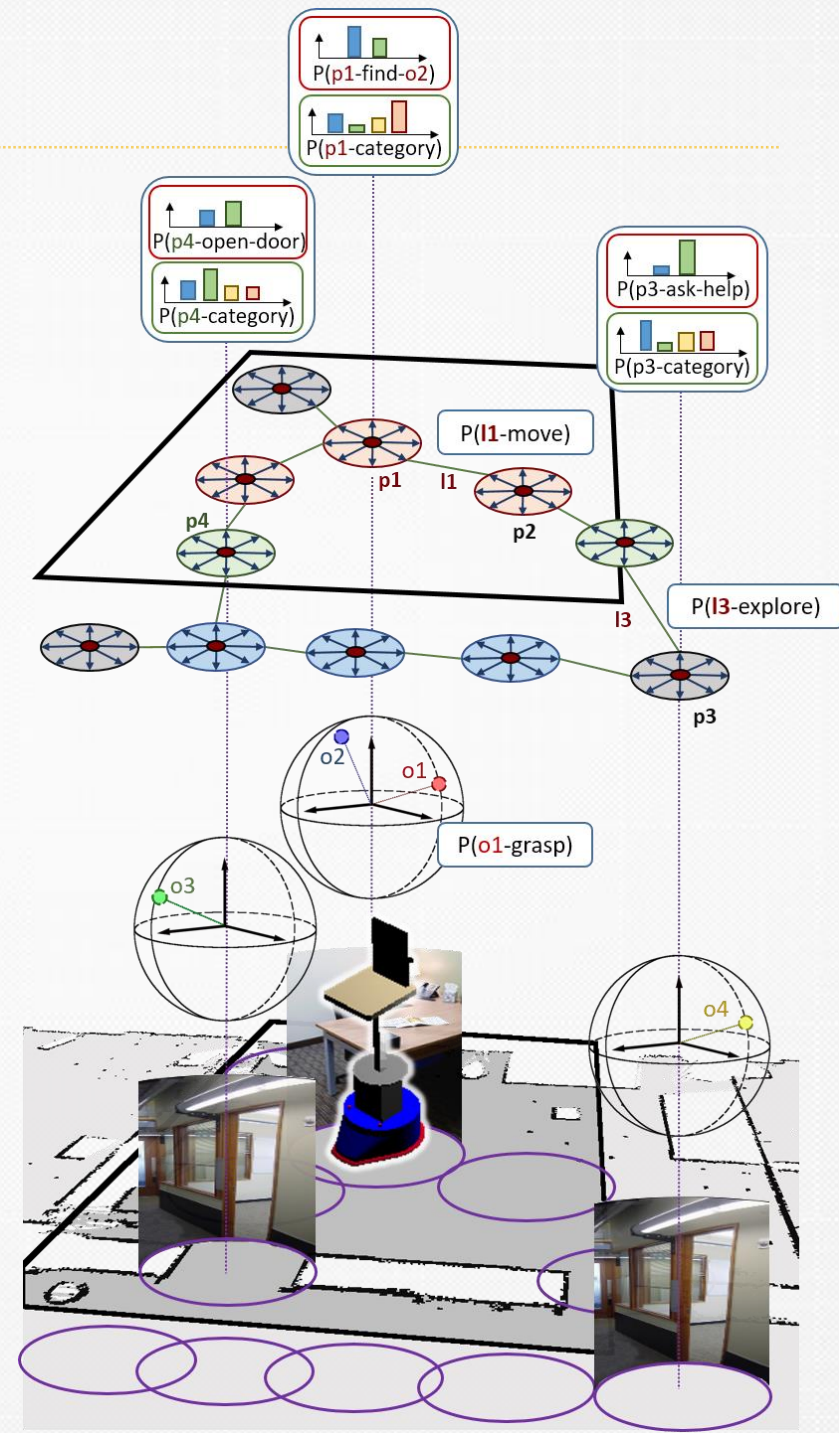
•

B



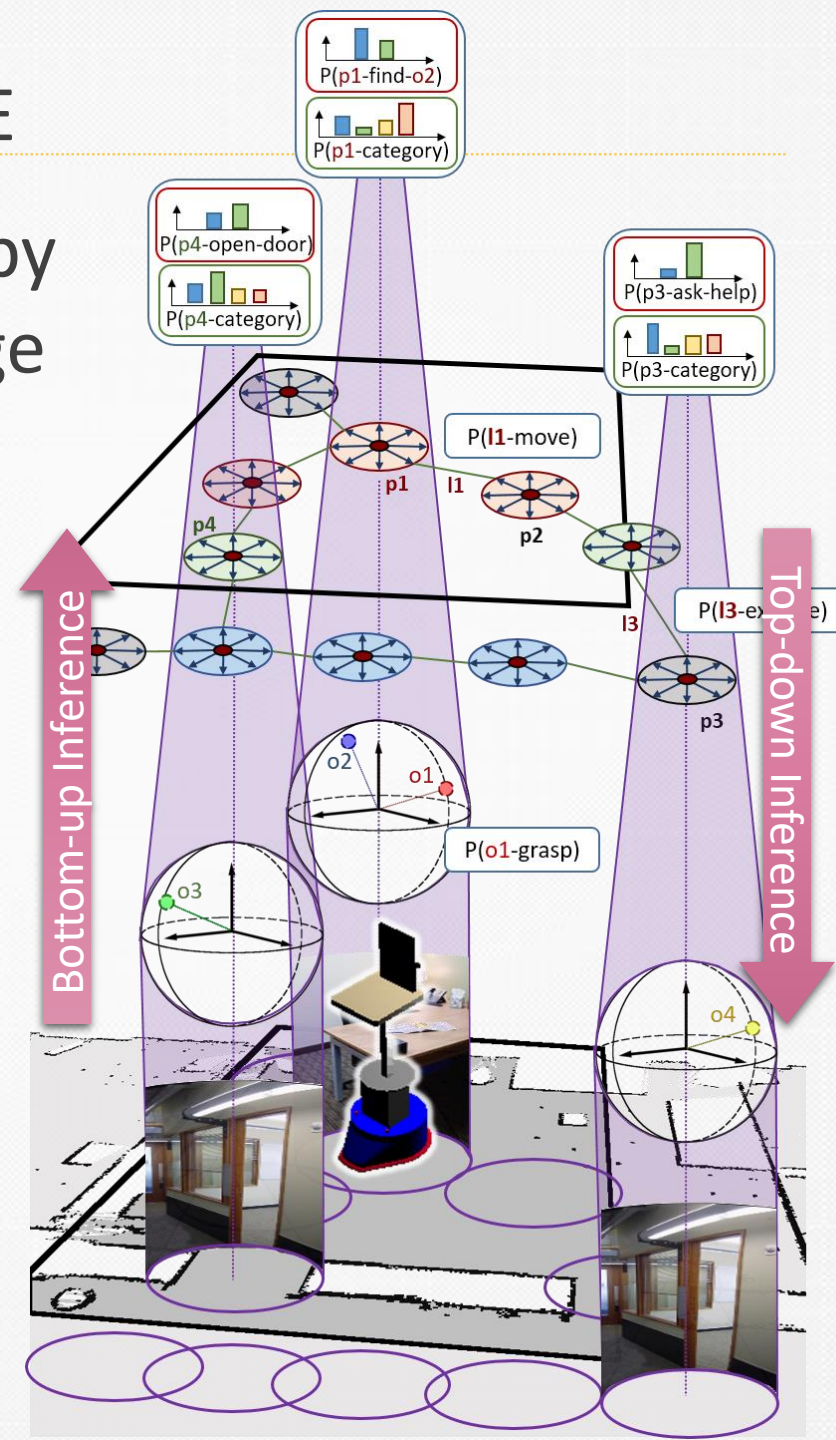
DASH: SEMANTIC LAYER

- Simple probabilistic relational representation
- Relates entities to semantic concepts
 - “place1 is-a kitchen”



DASH: GENERAL KNOWLEDGE

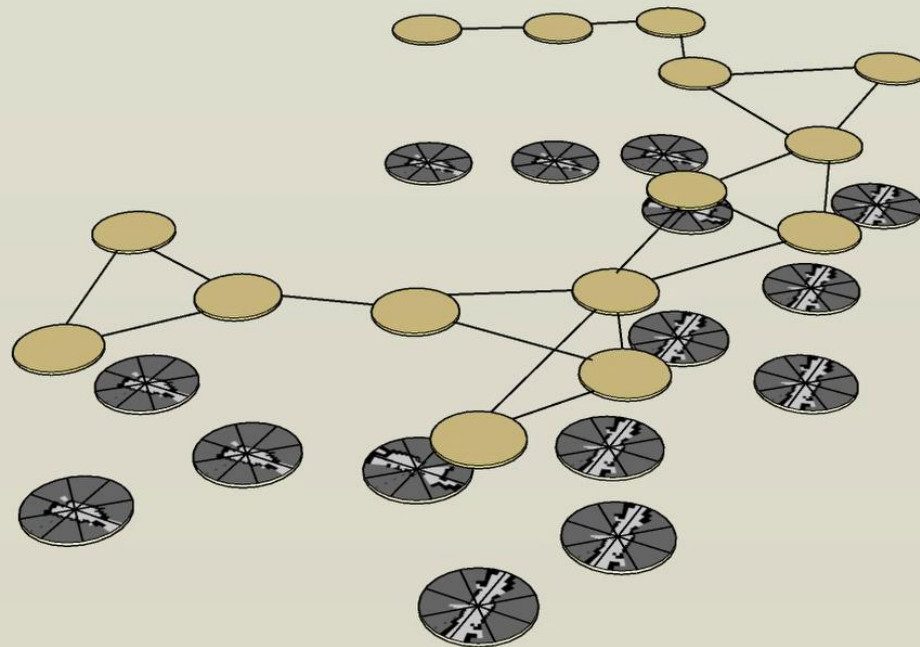
- **Instance** layers connected by model of **general** knowledge
- Enables top-down and bottom-up inferences
 - Filling missing data (**what's behind robot?**)
 - Inferring latent info (**what room is this?**)
 - Resolving ambiguities



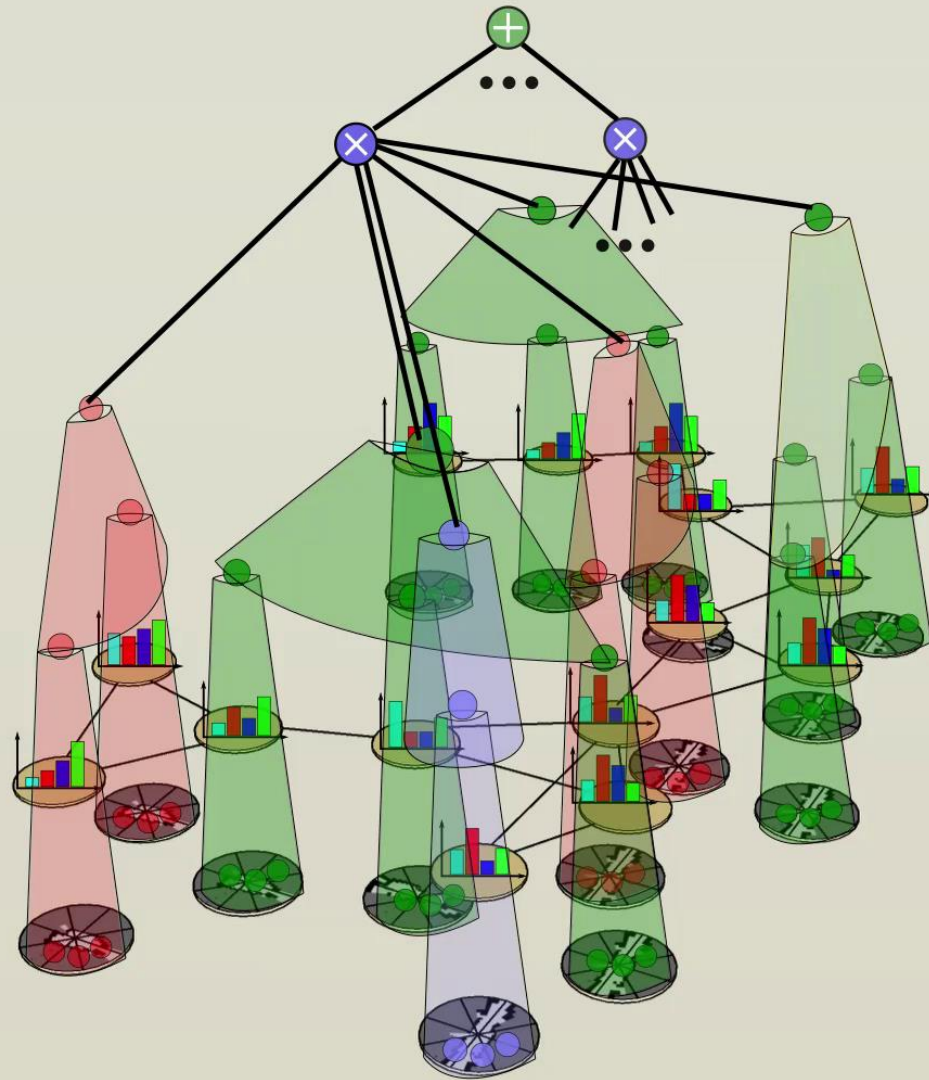
GENERAL KNOWLEDGE: END2END DEEP APPROACH



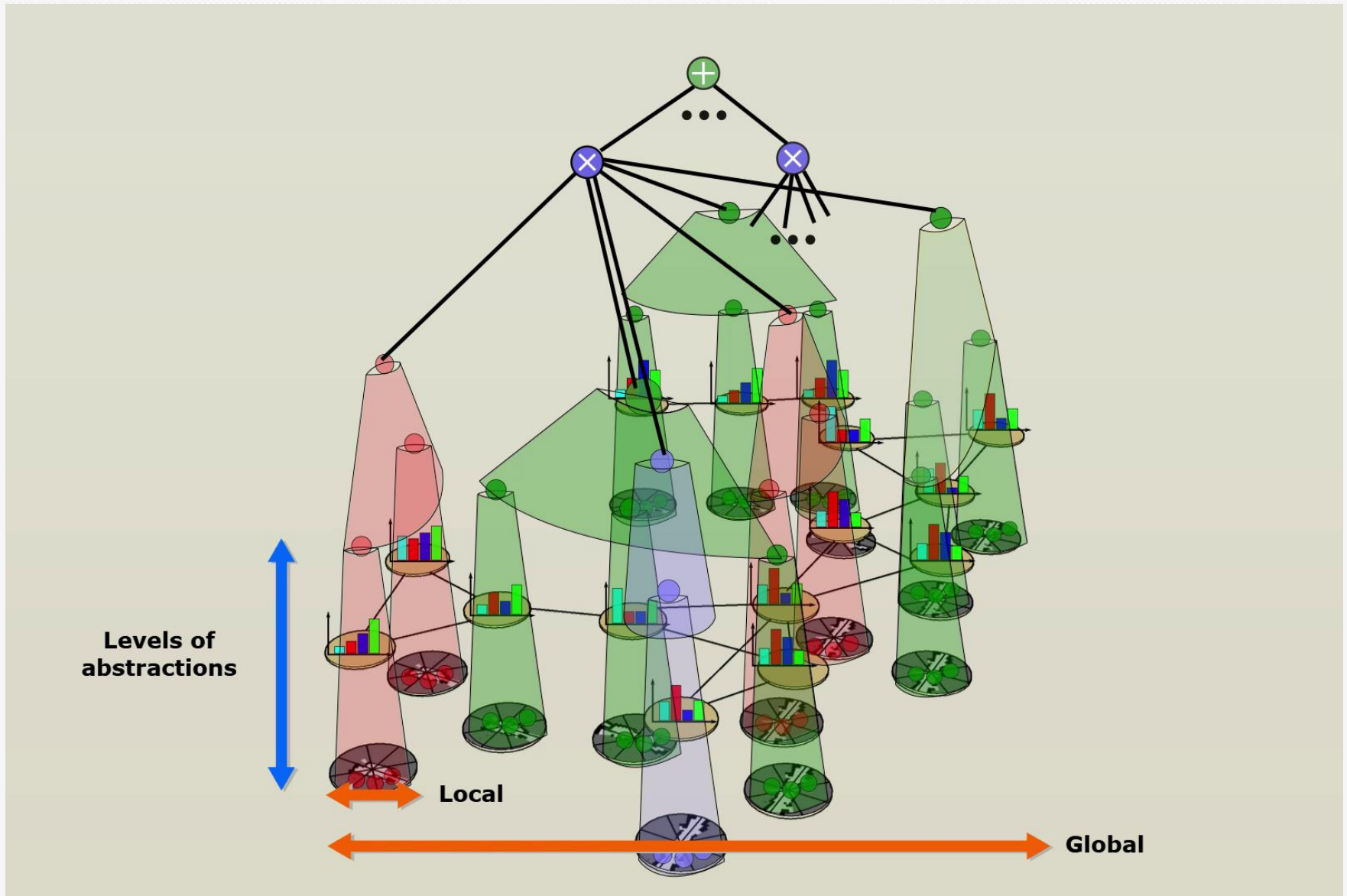
GENERAL KNOWLEDGE: END2END DEEP APPROACH



GENERAL KNOWLEDGE: END2END DEEP APPROACH

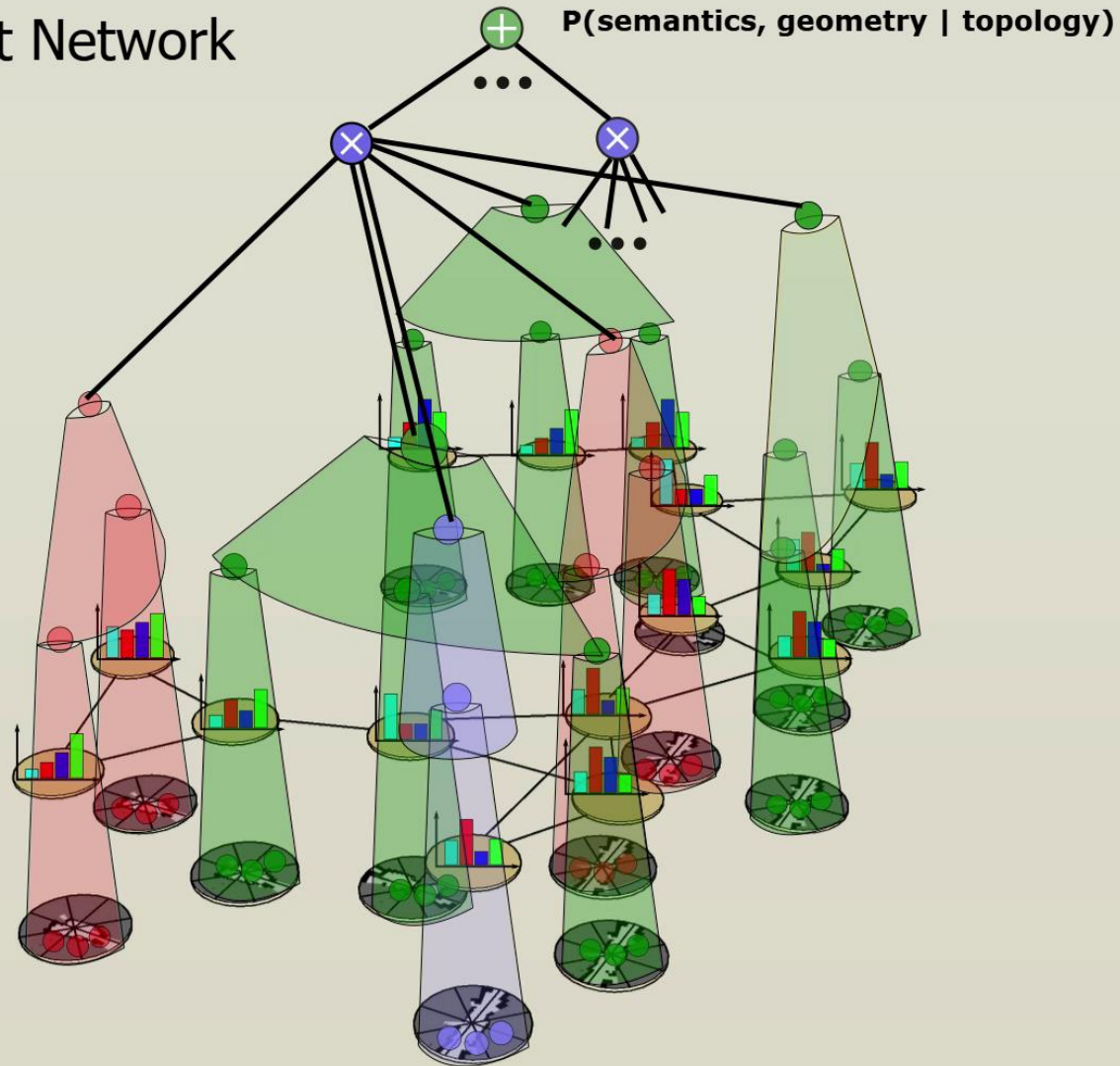


GENERAL KNOWLEDGE: END2END DEEP APPROACH



GENERAL KNOWLEDGE: END2END DEEP APPROACH

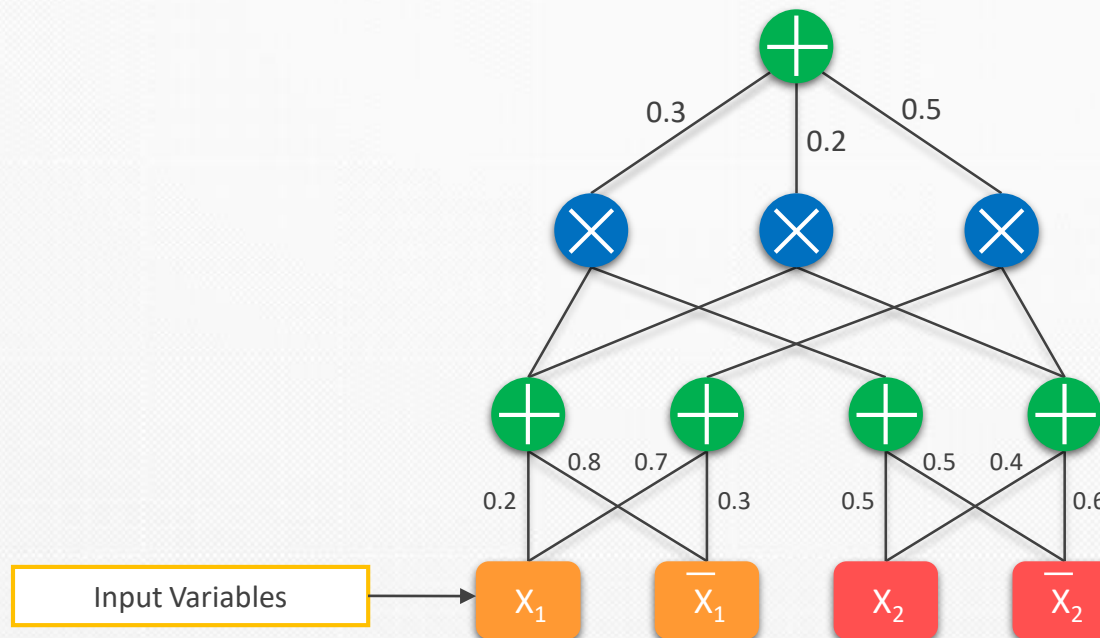
Sum-Product Network



SUM-PRODUCT NETWORKS

[Poon & Domingos, UAI'11,
Friesen & Domingos, ICML'16]

- 2 Views: Deep architecture and Graphical model
- Learn conditional or joint distributions
- Tractable partition function, exact inference
- Structure semantics: hierarchical mixture of parts

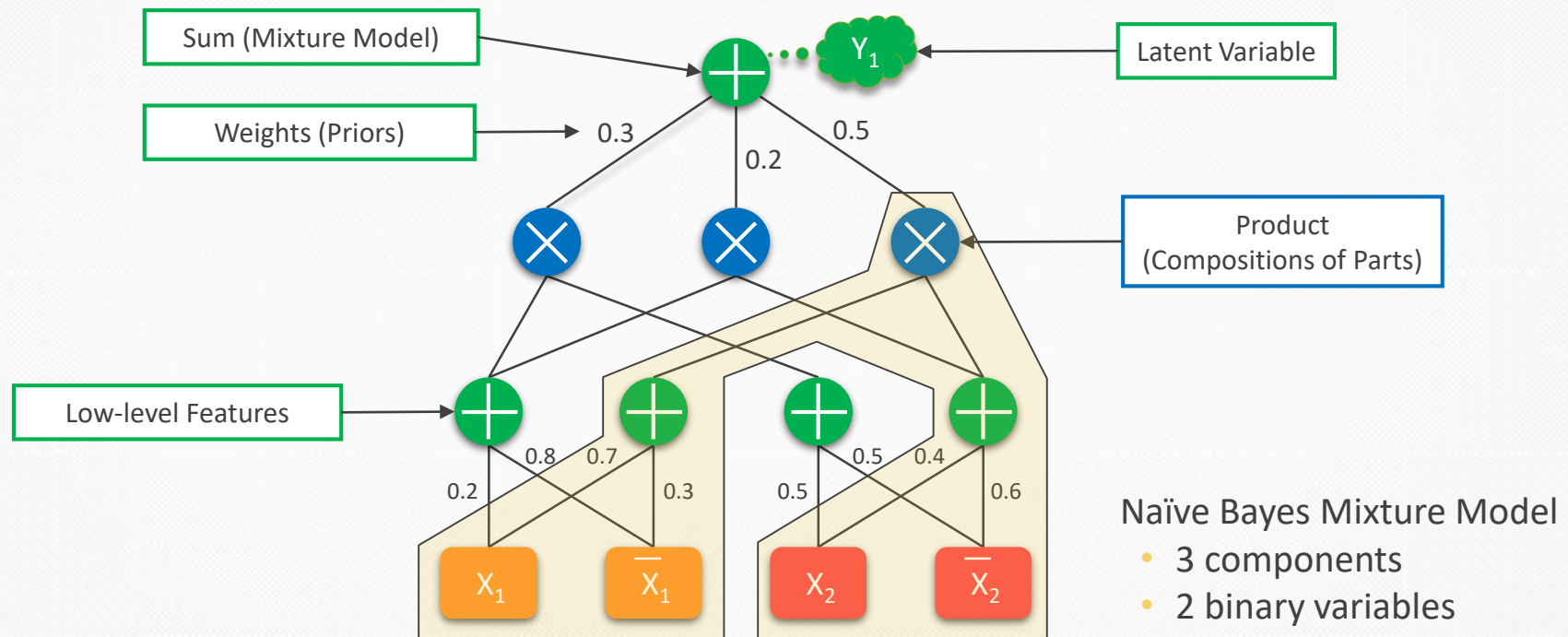


Naïve Bayes Mixture Model

- 3 components
- 2 binary variables

SUM-PRODUCT NETWORKS

[Poon & Domingos, UAI'11,
Friesen & Domingos, ICML'16]



- Large SPNs can be very deep

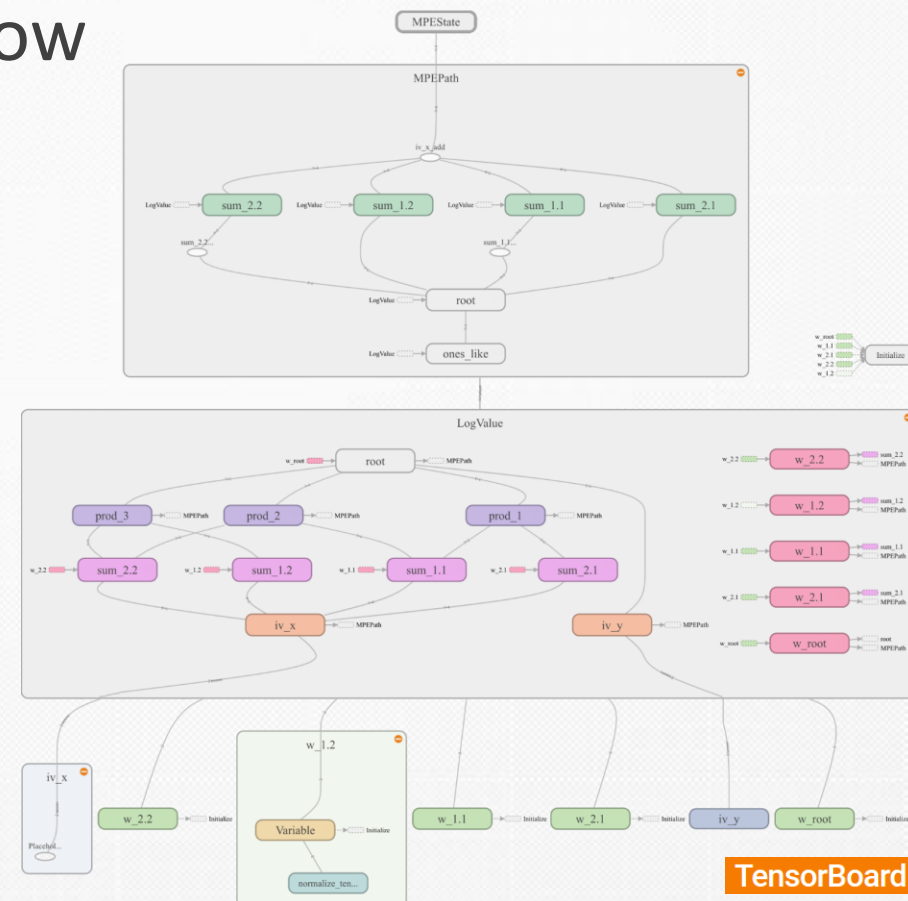
SPNS: LEARNING AND INFERENCE

- Learning
 - Generative (EM/GD) [Poon & Domingos, UAI'11]
 - Discriminative (GD) [Gens & Domingos, NIPS'12]
 - Simultaneous learning of parameters and structure [Gens & Domingos, ICML'13, Hsu et al., ICLR'17]
- Inference
 - Single up/down pass through the network
 - Upwards pass:
 - Probability of evidence
 - Downwards pass:
 - Gradients to obtain marginals
 - MPE state of variables

LIBSPN

[Pronobis, Ranganath, Rao, ICML-PADL'17]

- New general-purpose Python library for SPNs
- Learning and inference in large networks
- Integrated with TensorFlow
 - Multi-GPU computations
 - Integration of SPNs with other models
- Open-source soon at:
<http://libspn.org>

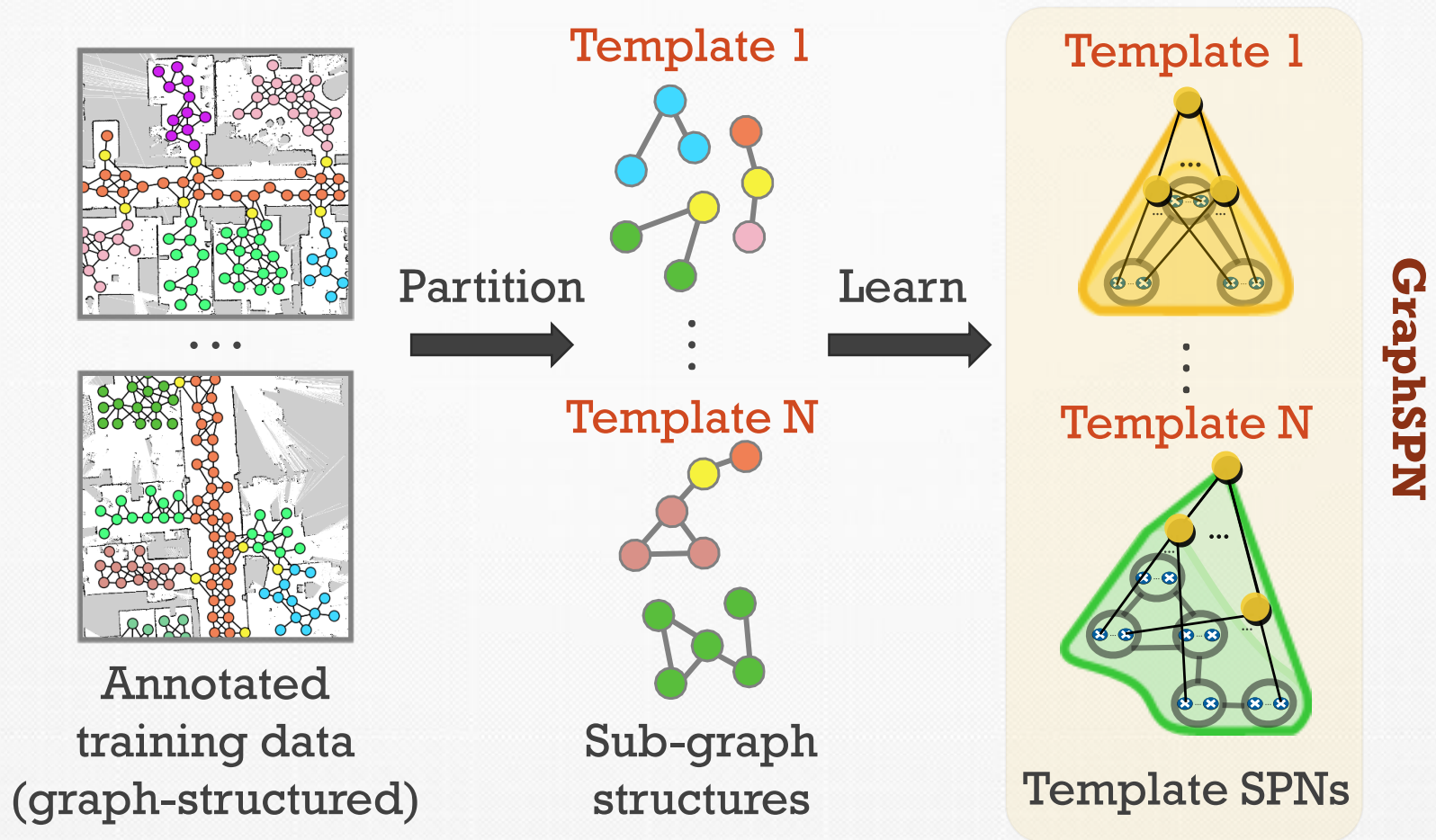


GRAPHSPNS: STRUCTURED PREDICTION WITH SPNS

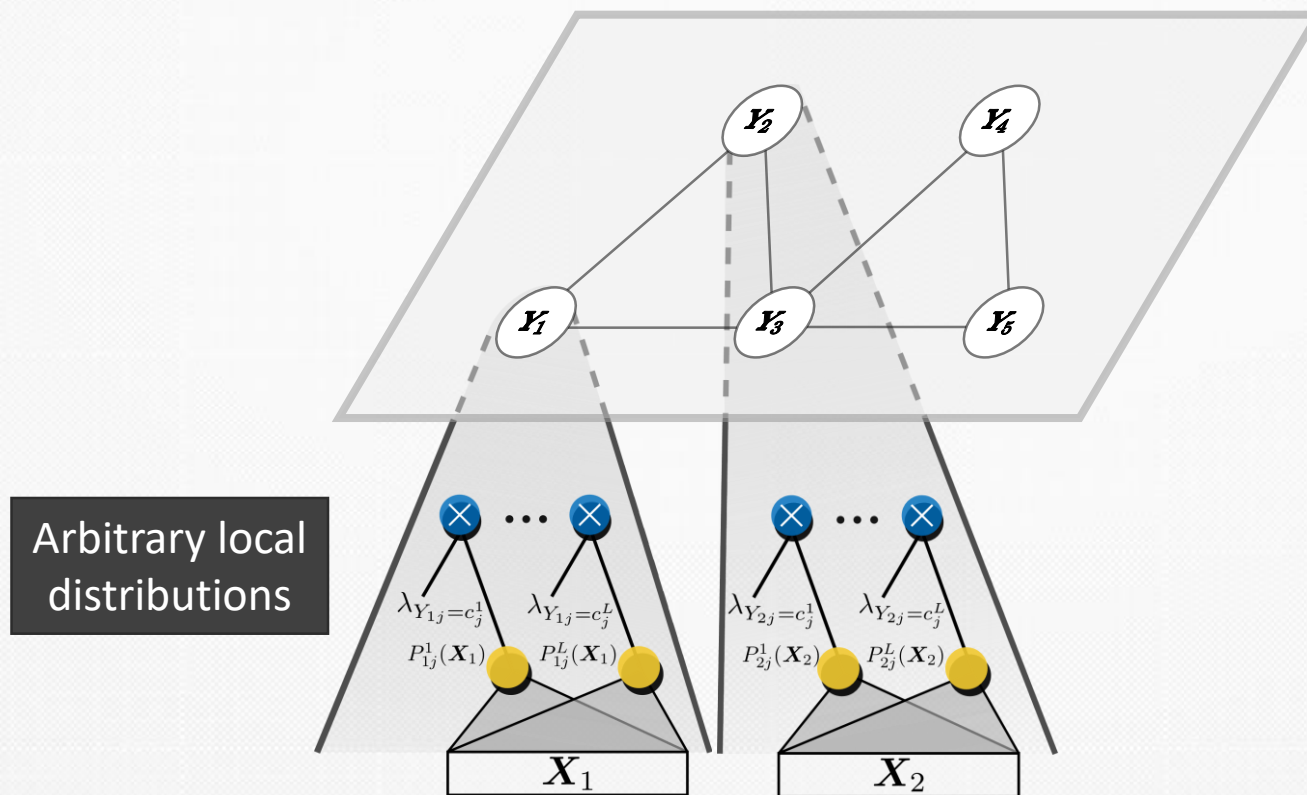
- Large environments pose complex structured prediction problems
 - Structured by graphs of varying size
 - Contaminated with noise
- Many approaches (deep), but:
 - Strict constraints on variable interactions
 - Fixed number of variables
 - Static global structure
- GraphSPNs [Zheng, Pronobis, Rao, AAAI'18]
 - Probabilistic SP approach
 - Dependencies between latent variables expressed in terms of **arbitrary, dynamic graphs**
 - Applicable to many problems (3D scenes, segmentation)



- Template model defined as a set of template SPNs representing learned, higher-order relations

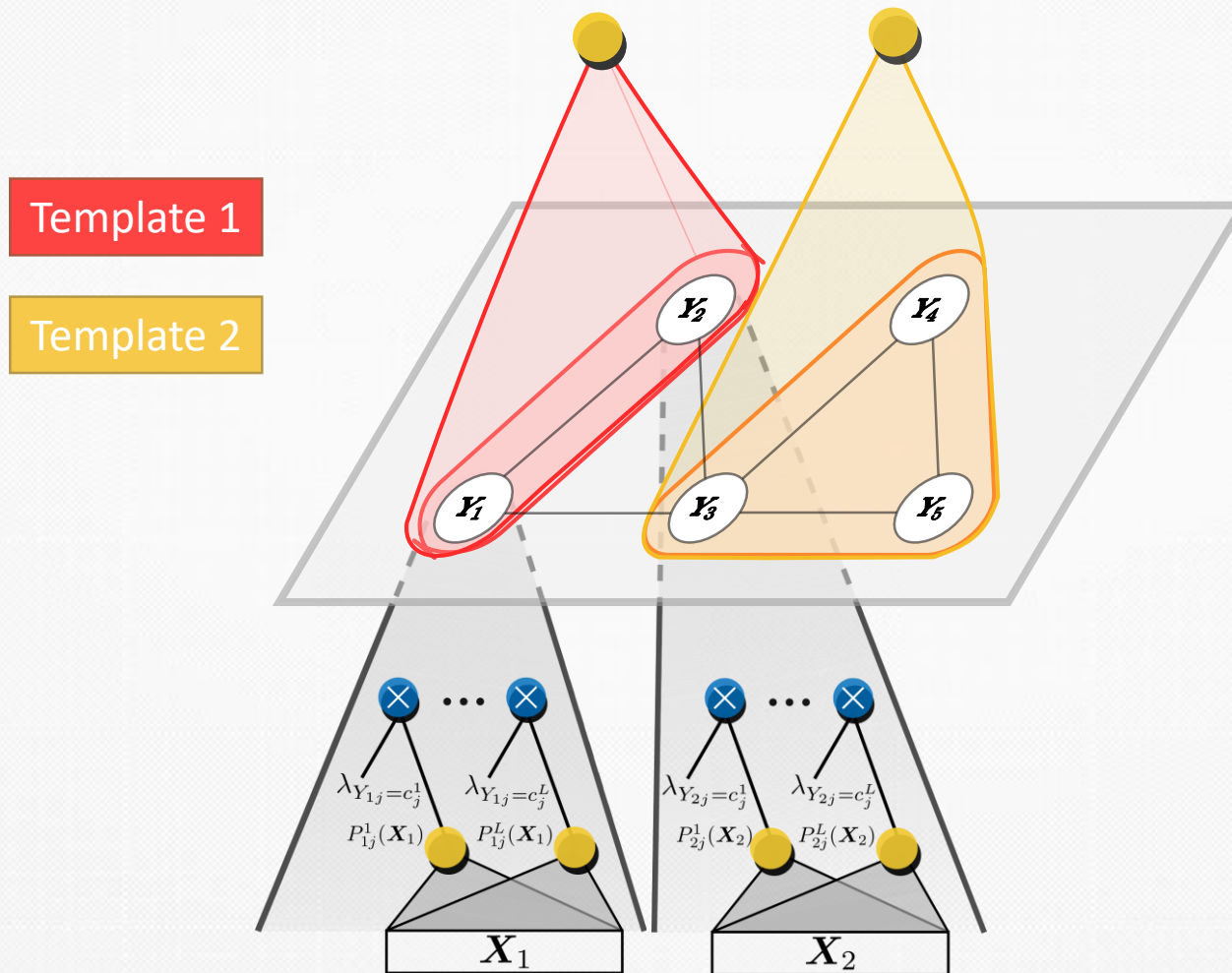


- Learned template models form a single distribution $P_D^{\mathcal{G}}(\mathbf{X}, \mathbf{Y})$ for a particular graph-structured problem



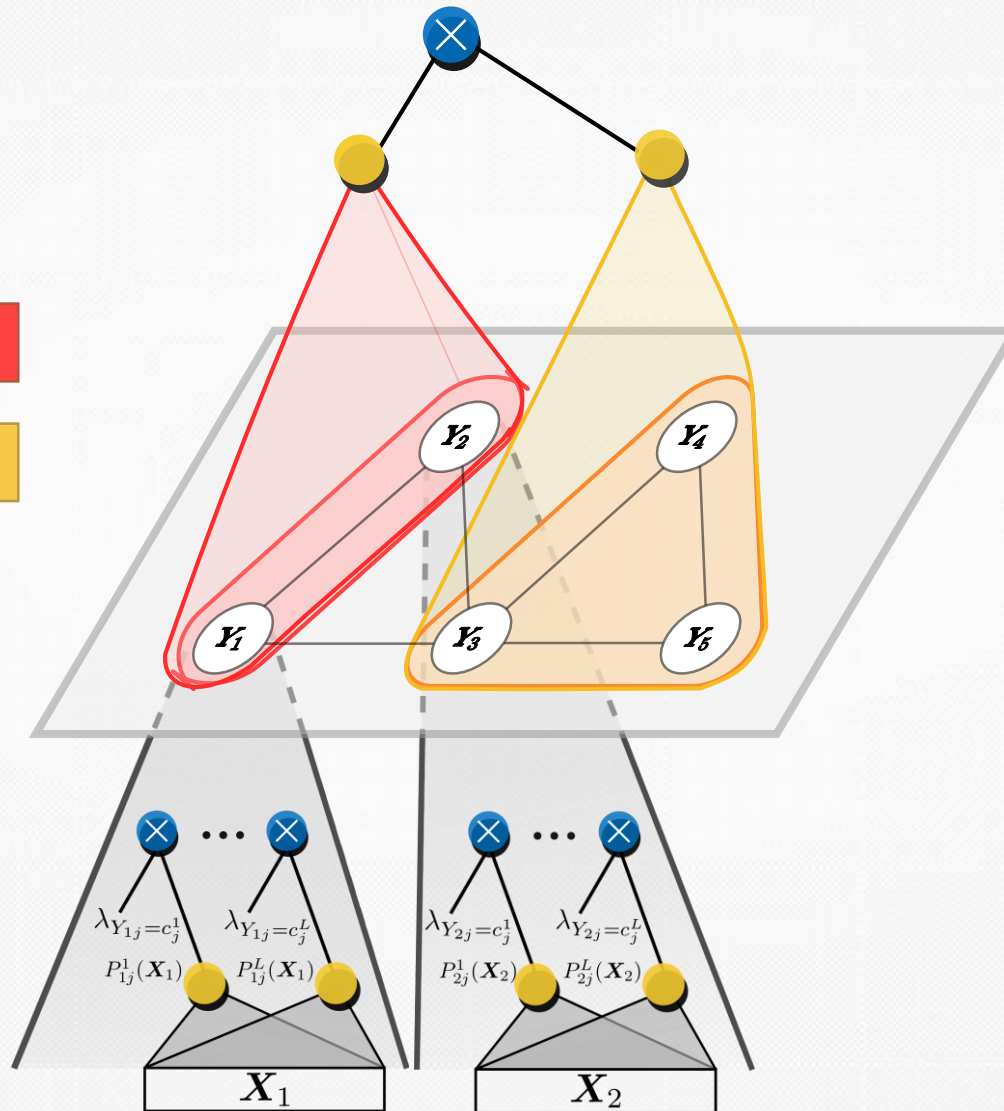
GRAPHSPNS: INFERENCE

[Zheng, Pronobis, Rao, AAAI'18]



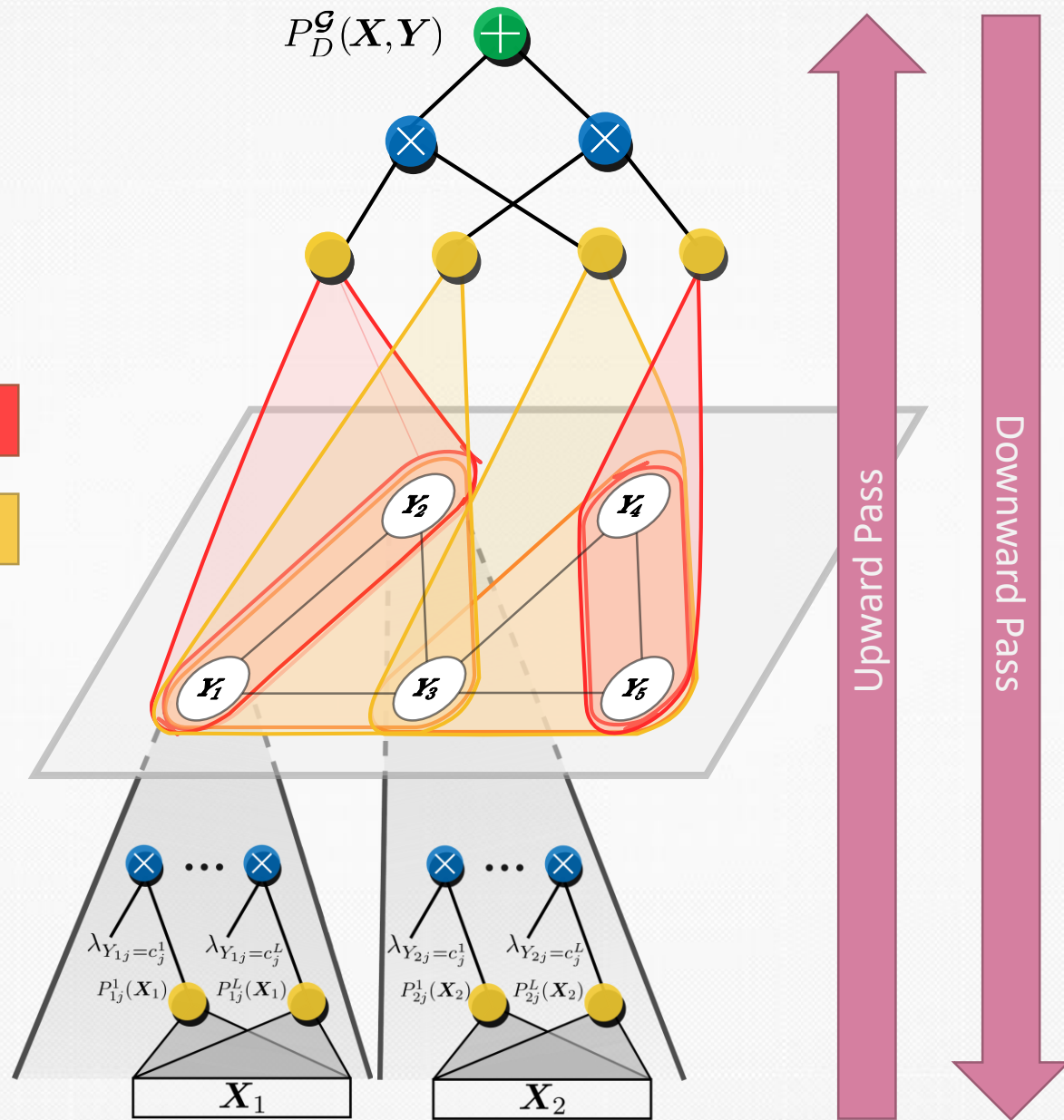
Template 1

Template 2



Template 1

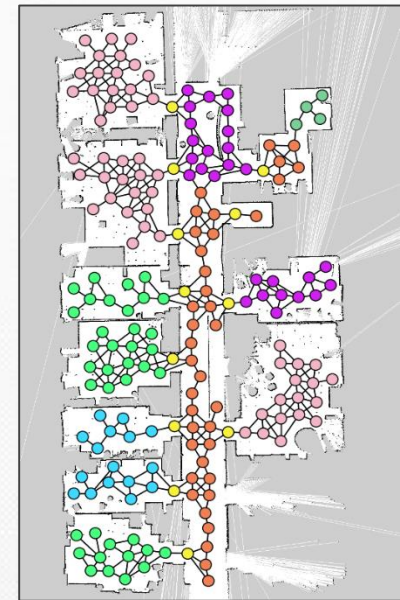
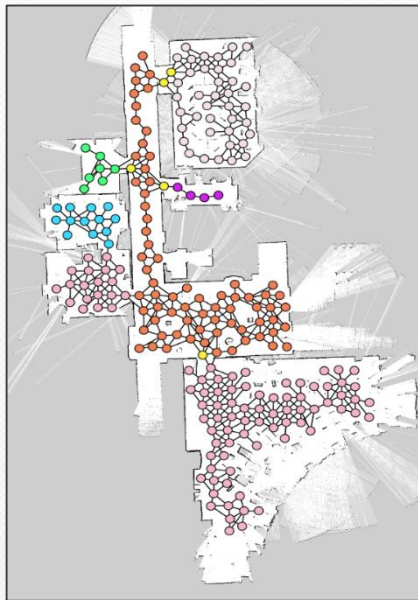
Template 2



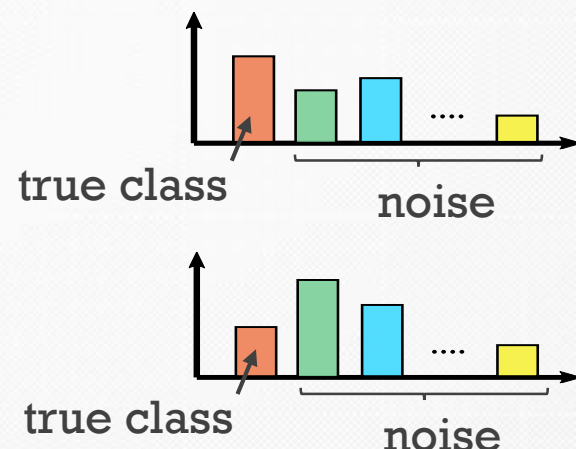
GRAPHSPNS: EVALUATION

[Zheng, Pronobis, Rao, AAAI'18]

- GraphSPNs vs graphical models (MRFs)
- Data structured by topological graphs
 - 99 graphs: 11 floors in 3 buildings
 - Places labeled with 10 semantic categories
 - Leave-one-building-out procedure

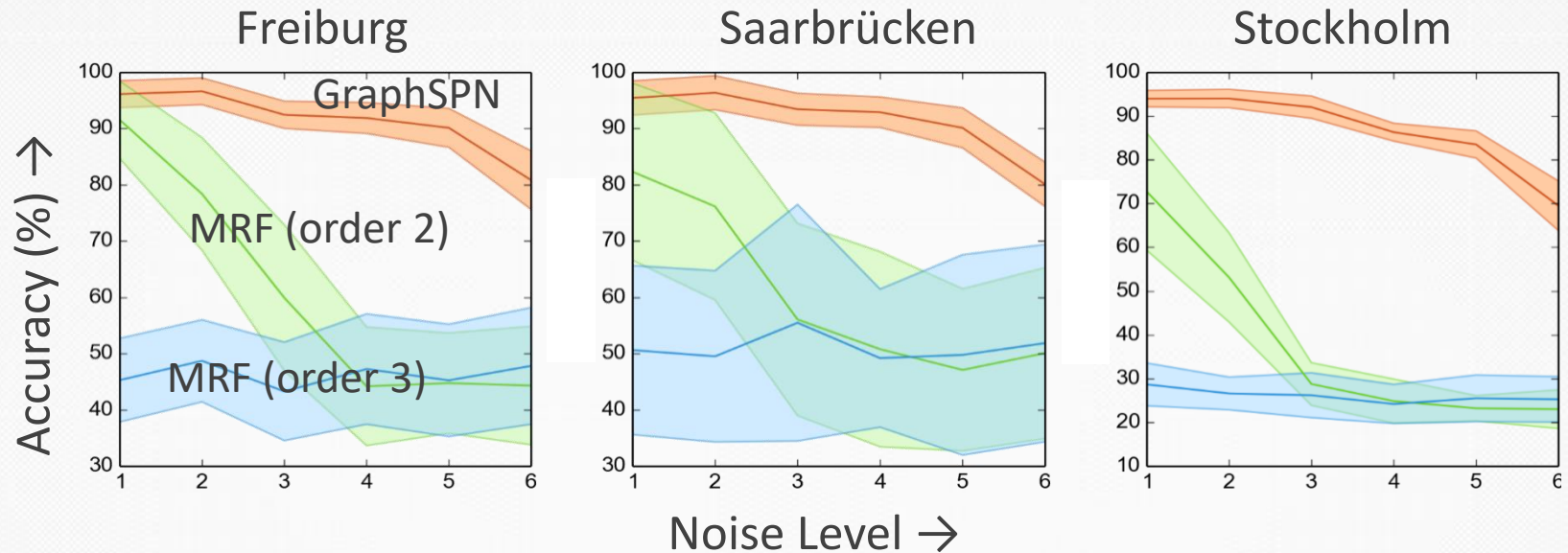


- GraphSPNs vs graphical models (MRFs)
- Data structured by topological graphs
 - 99 graphs: 11 floors in 3 buildings
 - Places labeled with 10 semantic categories
 - Leave-one-building-out procedure
- Local evidence:
 - Distributions over semantic categories
 - Noisified ground truth
- Local distributions:
 - Unary potentials (MRFs)
 - Corresponding basic distributions (GraphSPN)
- Task: Recover true semantic labels of places



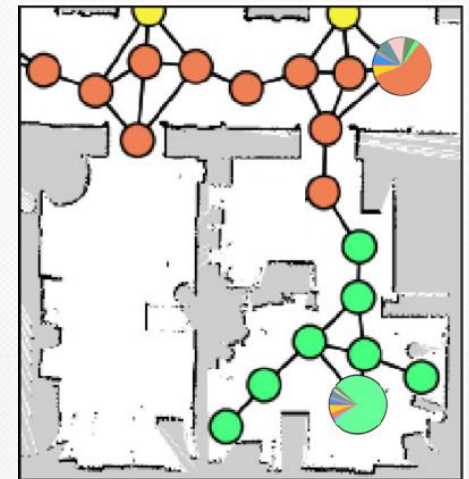
GRAPHSPNS: RESULTS

[Zheng, Pronobis, Rao, AAAI'18]

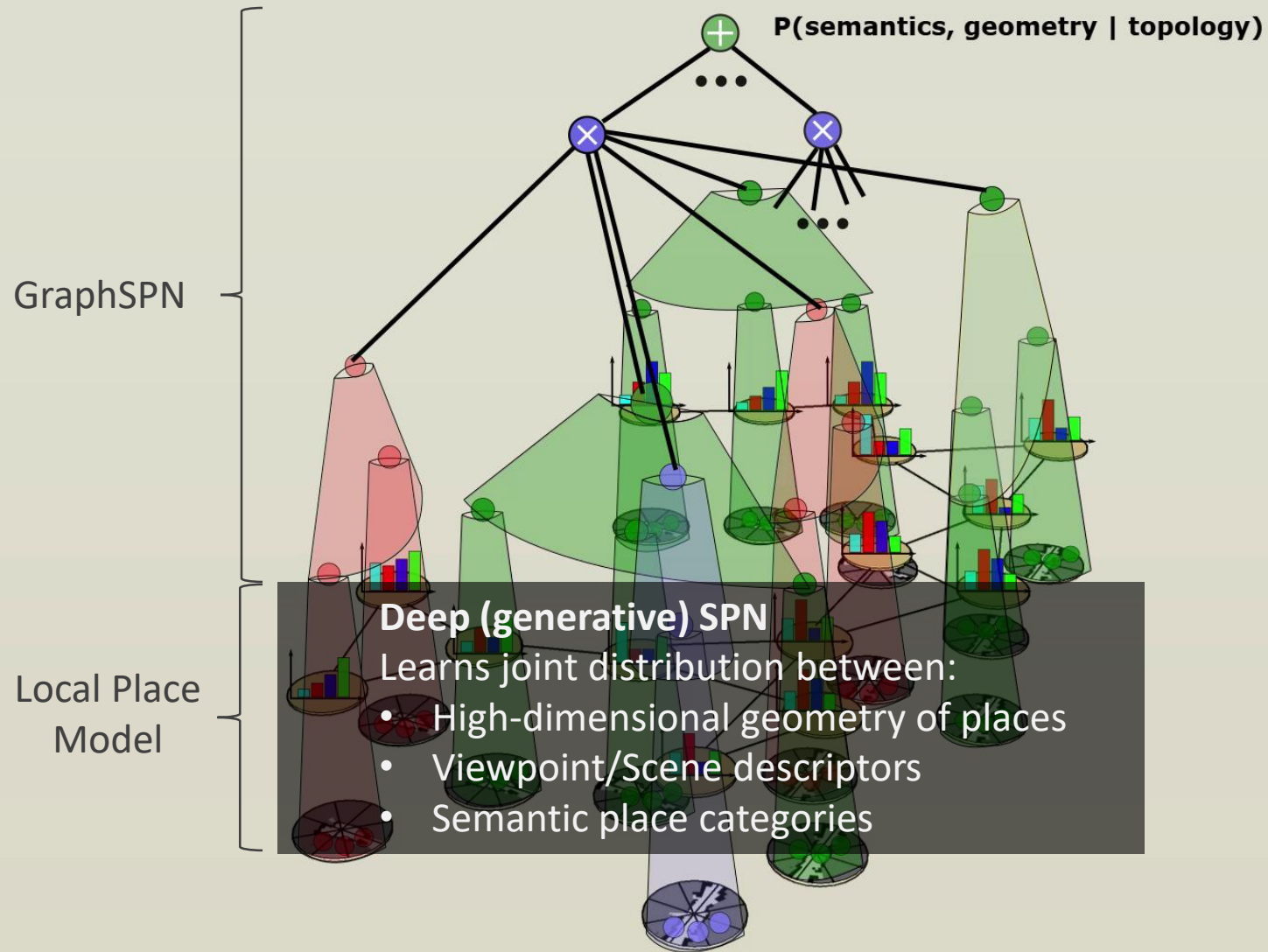


Predictions for nodes without evidence (unexplored placeholders):

GraphSPN		
Freiburg	Saarbrücken	Stockholm
67.58%(+/-10.42)	78.15%(+/-9.95)	67.57%(+/-11.11)
MRF-2		
Freiburg	Saarbrücken	Stockholm
28.32%(+/-7.53)	39.85%(+/-19.42)	12.44%(+/-3.46)

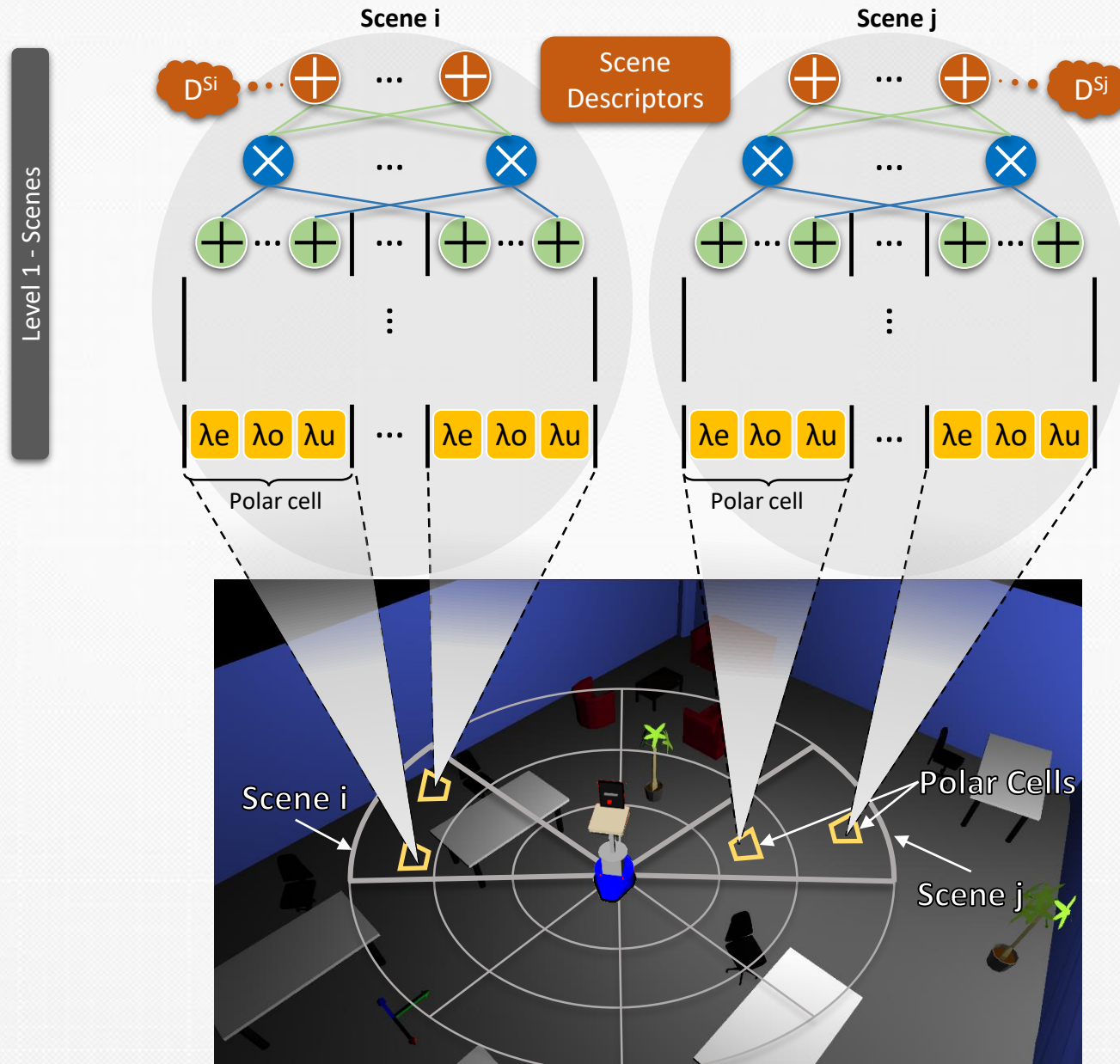


GENERAL KNOWLEDGE: END2END DEEP APPROACH



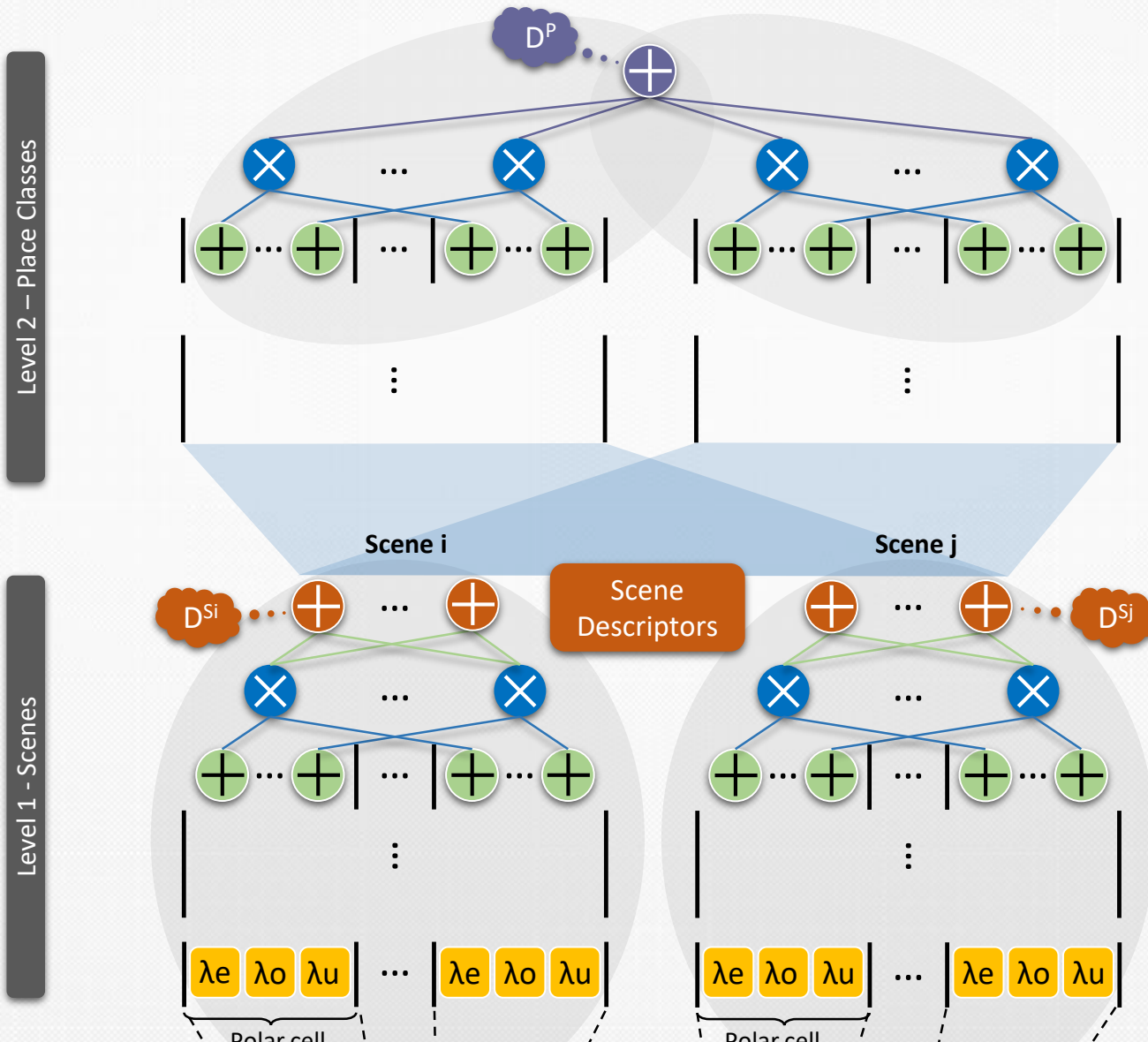
LOCAL PLACE MODEL: ARCHITECTURE

[Pronobis, Rao, IROS'17]

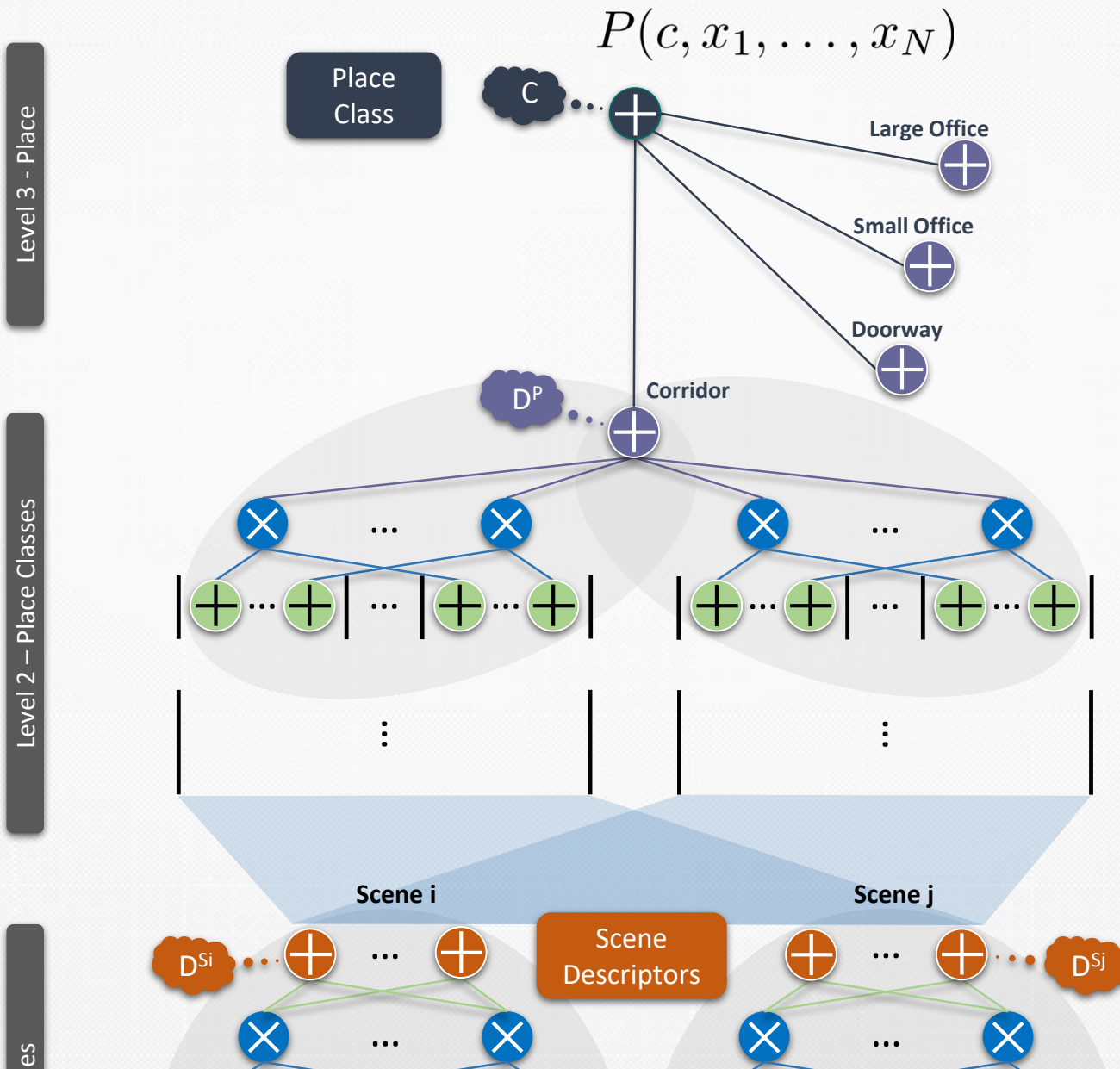


LOCAL PLACE MODEL: ARCHITECTURE

[Pronobis, Rao, IROS'17]



LOCAL PLACE MODEL: ARCHITECTURE [Pronobis, Rao, IROS'17]



- Semantically annotated sequences of sensory data
 - Robot navigating
4 floors of office building
 - Each floor contains
multiple instances of:
small office, large office,
corridor, doorway
(considered known)
 - Additional few instances of:
kitchen, elevator, lab,
living room, meeting room
(considered novel)
 - Leave-one-floor-out procedure
- Sensor: Laser-range scanner
- Learning: Generative hard EM

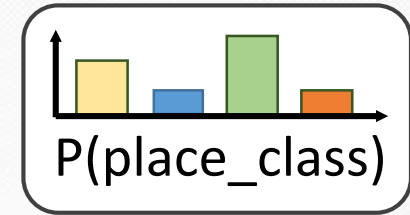


BOTTOM-UP INFERENCE: PLACE CATEGORIZATION

[Pronobis, Rao, IROS'17]



$$c^* = \underset{c}{\operatorname{argmax}} P(c, x_1, \dots, x_N)$$



- Baseline:
 - SVM + RBF Kernel
 - Geometric features [Mozos et al. '05] from high-res 360° virtual scans
- Average CR:
 - SVM (discriminative): 85.9%
 - DGSM (generative): 92.7%

Confusion Matrix (DGSM)

True Class	Predicted Class			
	Corridor	Doorway	Small Office	Large Office
Corridor	98.3	1.5	0.2	0.0
Doorway	0.0	95.7	2.2	2.2
Small Office	0.0	0.6	88.3	11.1
Large Office	0.0	0.6	5.0	94.4

BOTTOM-UP INFERENCE: NOVELTY DETECTION

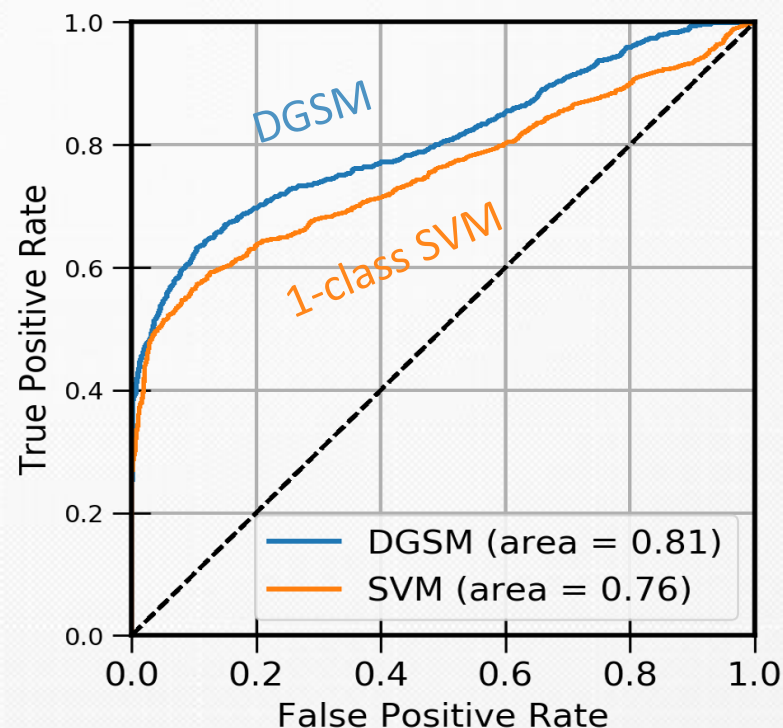
[Pronobis, Rao, IROS'17]



$$P(x_1, \dots, x_N) > threshold$$

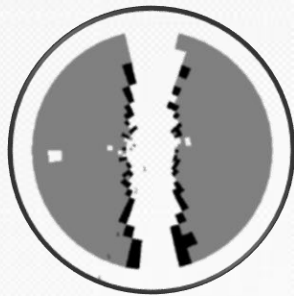
Known / Novel

- Baseline:
 - One-class SVM + RBF Kernel
 - Geometric features



TOP-DOWN INFERENCE: GENERATING PROTOTYPES

[Pronobis, Rao, IROS'17]



Corridor

$$x_1^*, \dots, x_N^* = \operatorname{argmax}_{x_1, \dots, x_N} P(x_1, \dots, x_N | c)$$

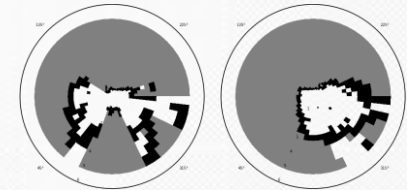
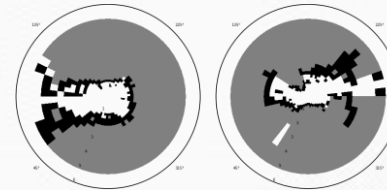
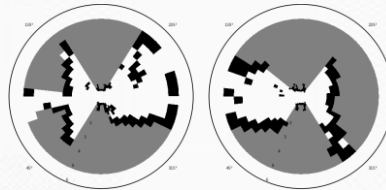
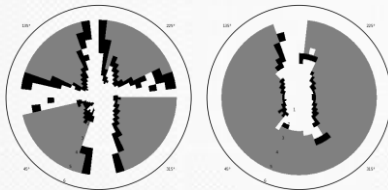
Corridor

Doorway

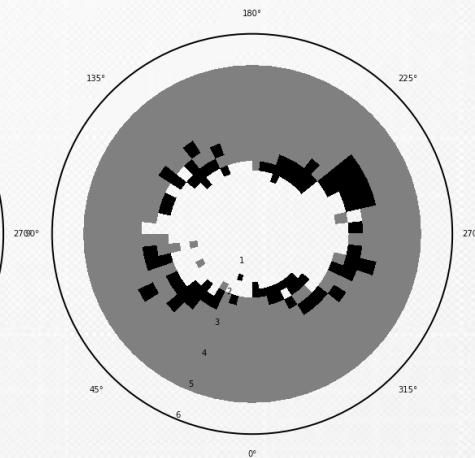
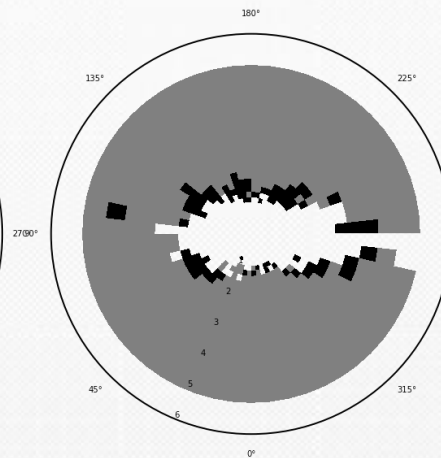
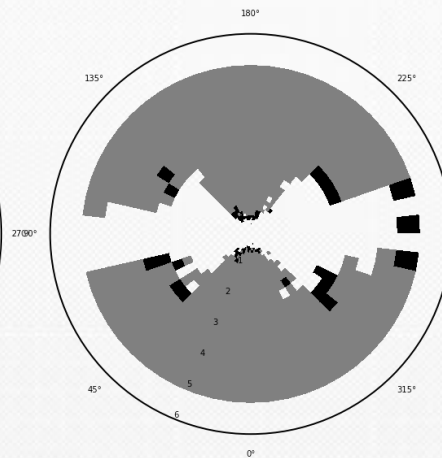
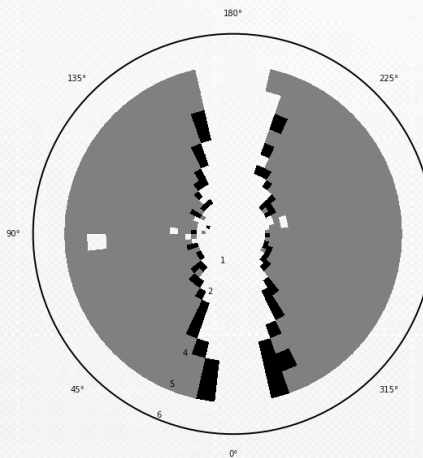
Small Office

Large Office

True



Generated



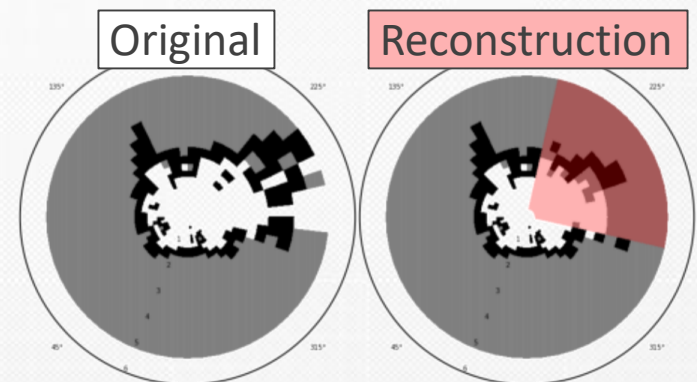
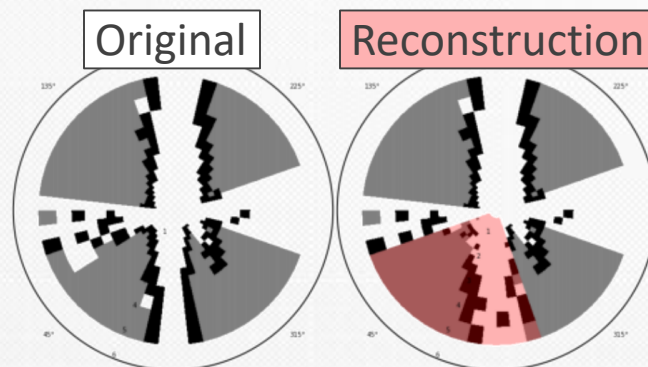
TOP-DOWN INFERENCE: MISSING OBSERVATIONS

[Pronobis, Rao, IROS'17]

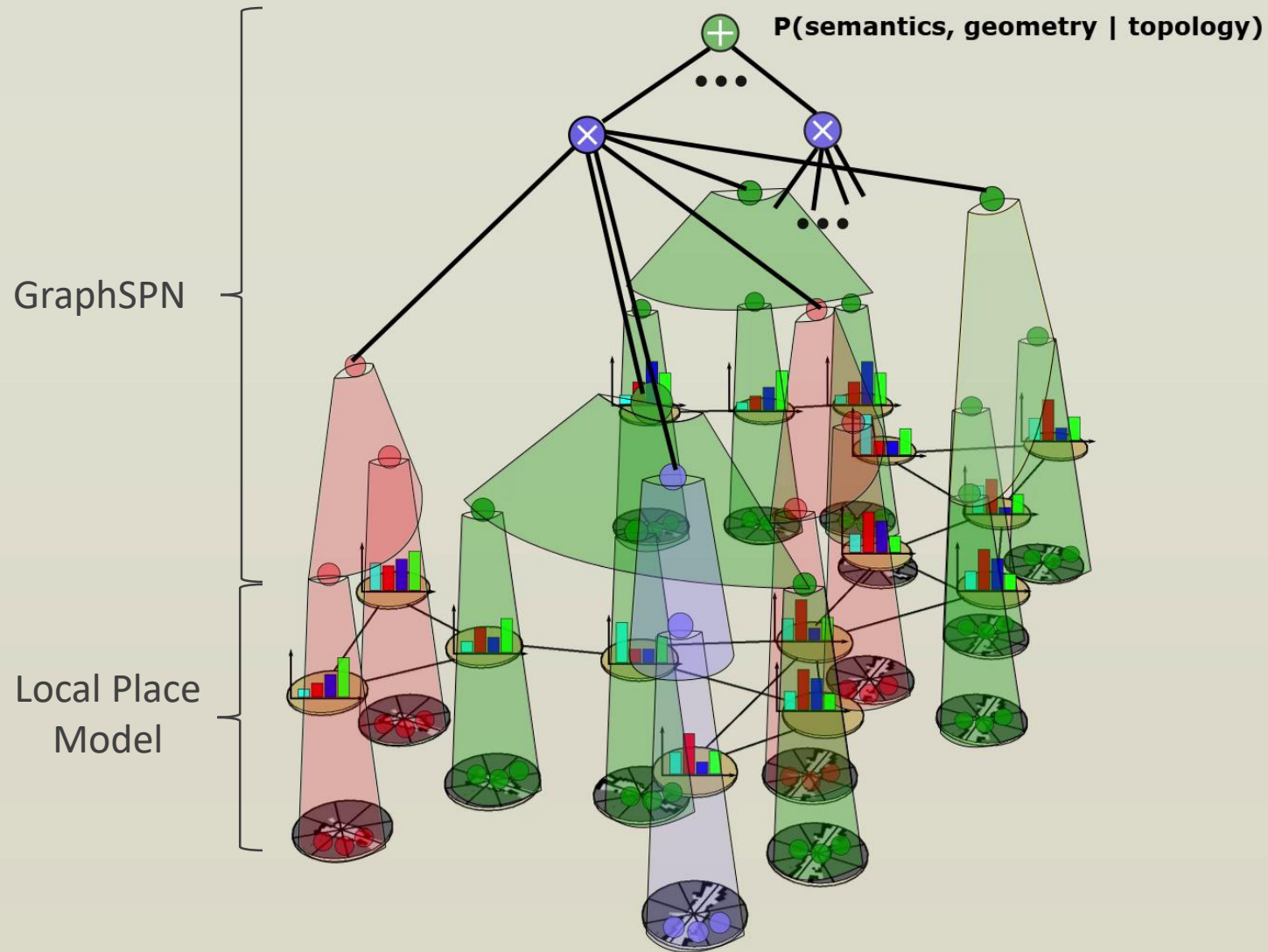


$$x_J^*, \dots, x_N^* = \underset{x_J, \dots, x_N}{\operatorname{argmax}} P(x_1, \dots, x_J, \dots, x_N)$$

- Baseline: DC-GANs + GD-based inpainting [Yeh et al. '16]
- Correctly predicted: **DGSM: 77.1%** **DC-GAN: 75.8%**

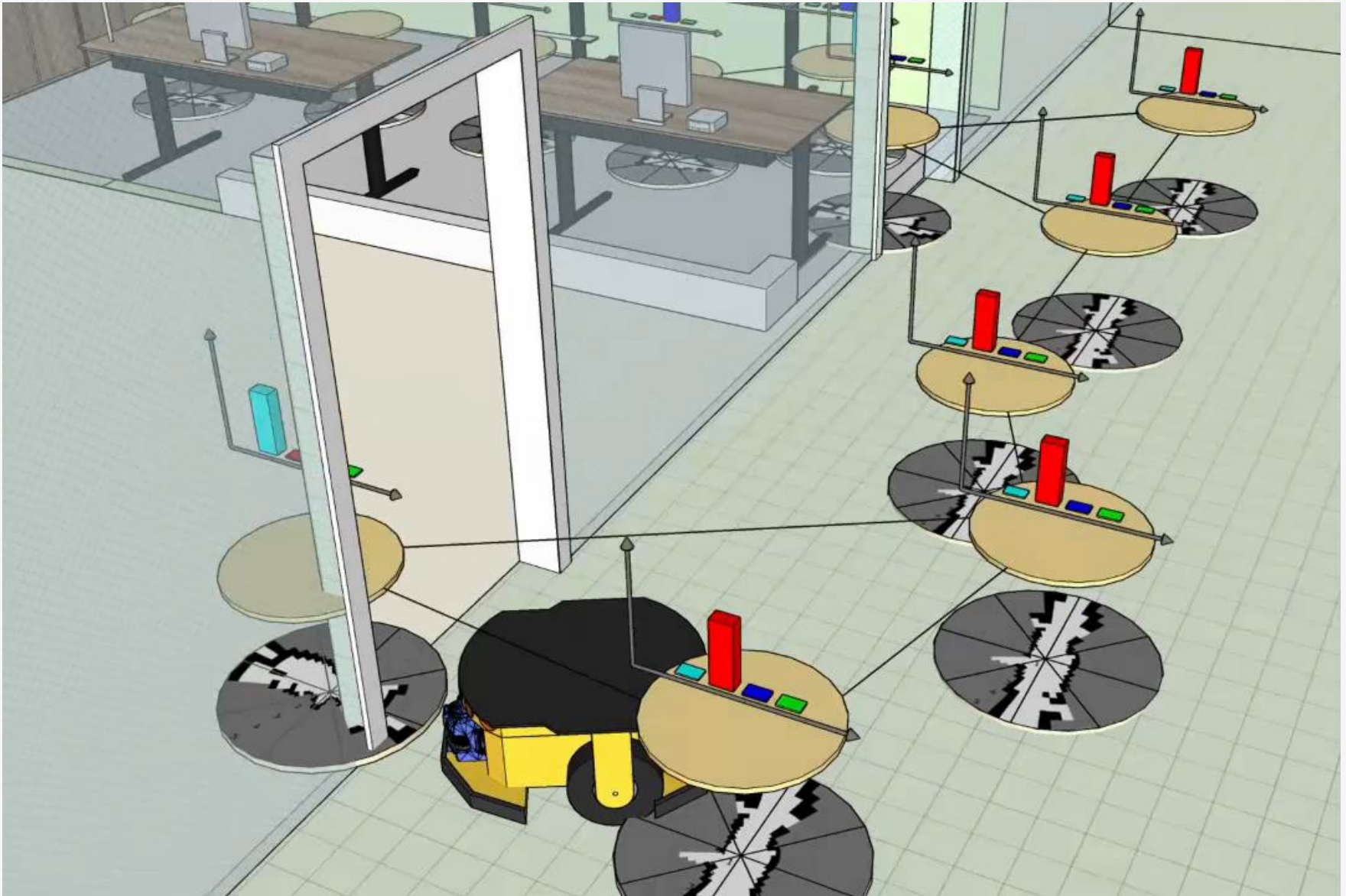


GENERAL KNOWLEDGE: END2END DEEP APPROACH



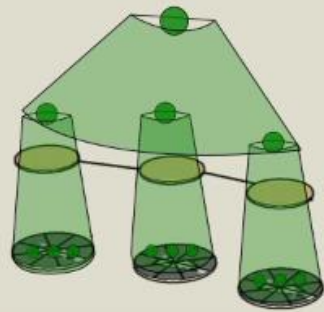
END2END: LEARNING

[Zheng, Pronobis '19]

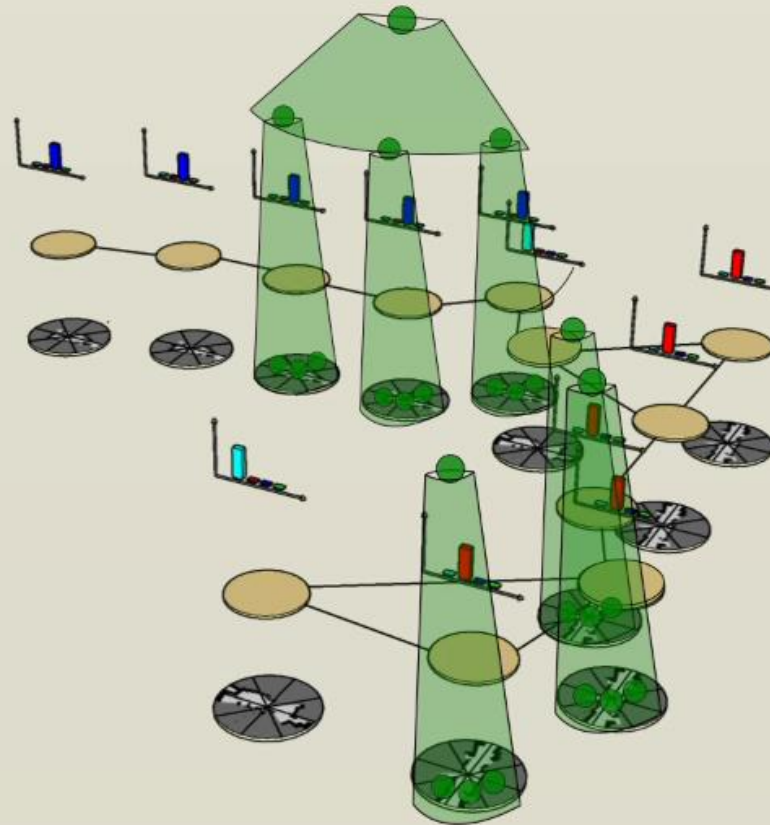


END2END: LEARNING

[Zheng, Pronobis '19]

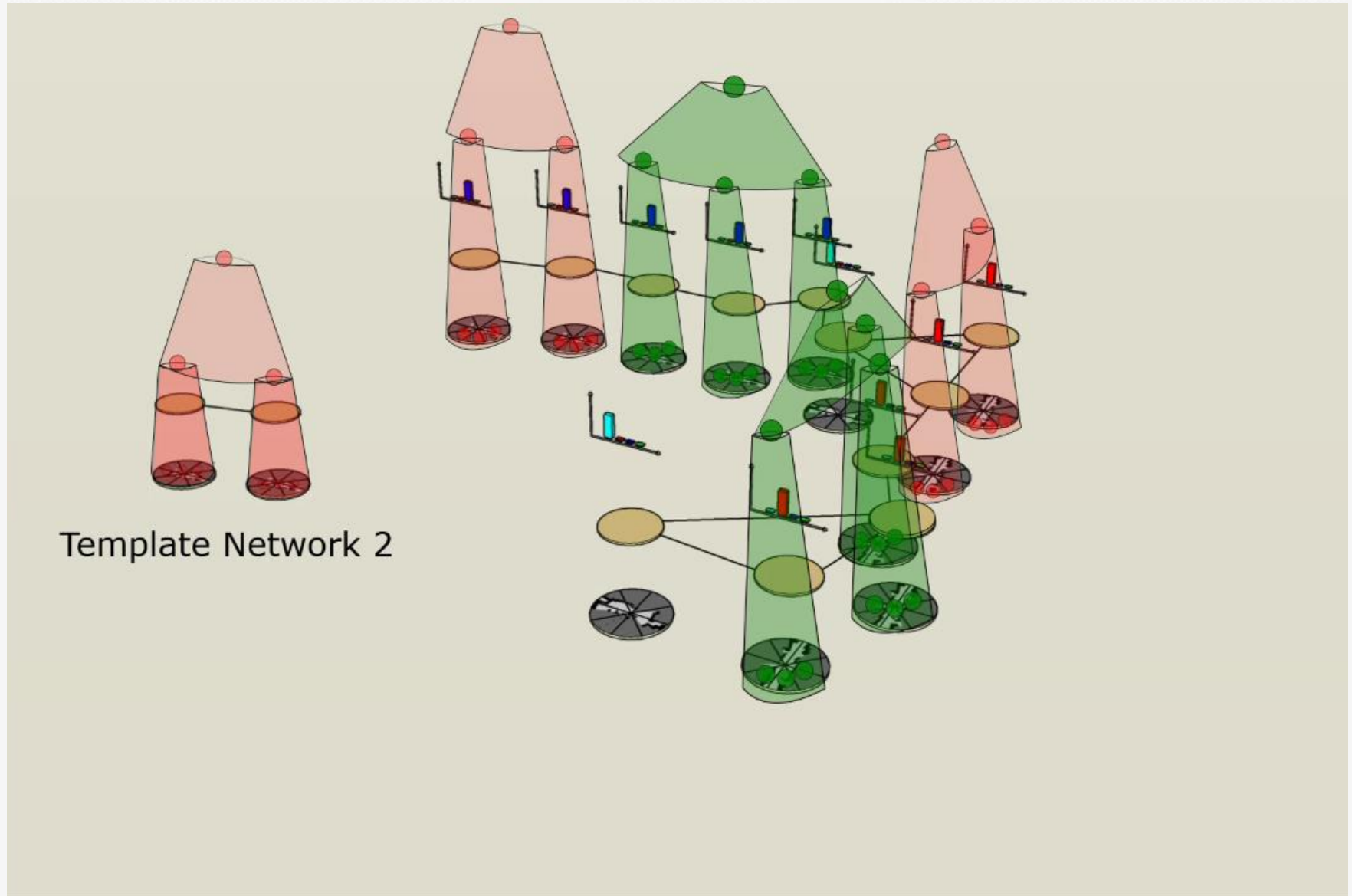


Template Network 1



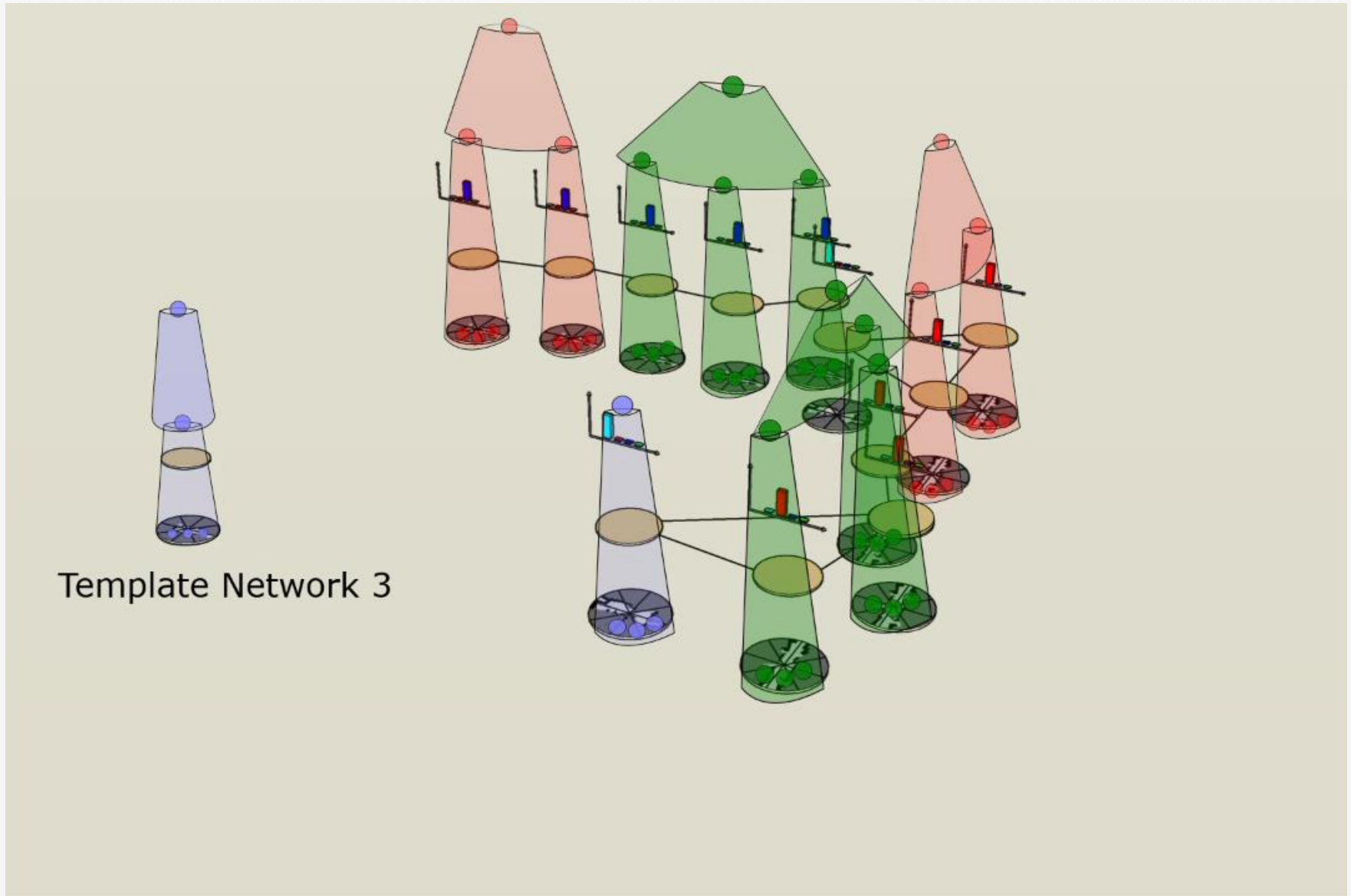
END2END: LEARNING

[Zheng, Pronobis '19]



END2END: LEARNING

[Zheng, Pronobis '19]



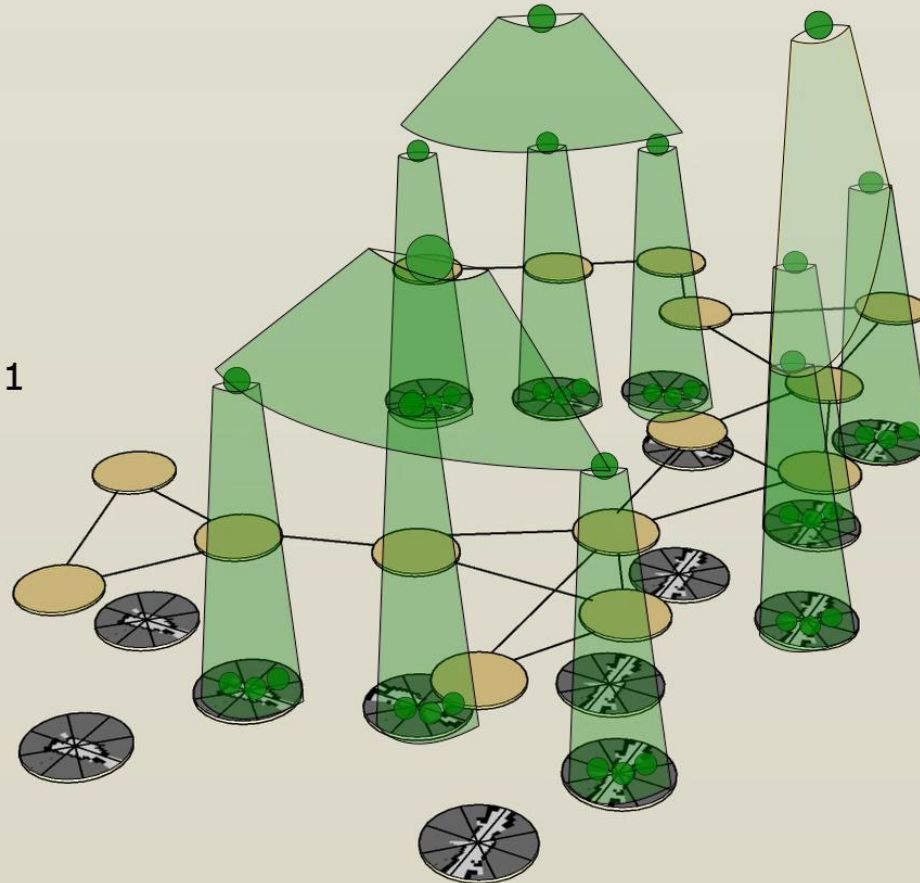
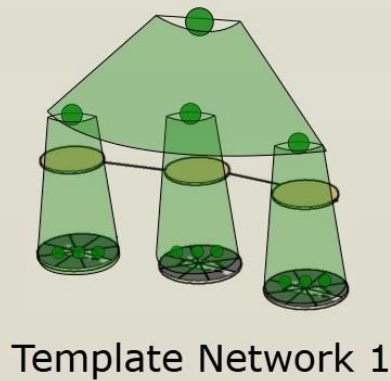
END2END: INFERENCE

[Zheng, Pronobis '19]



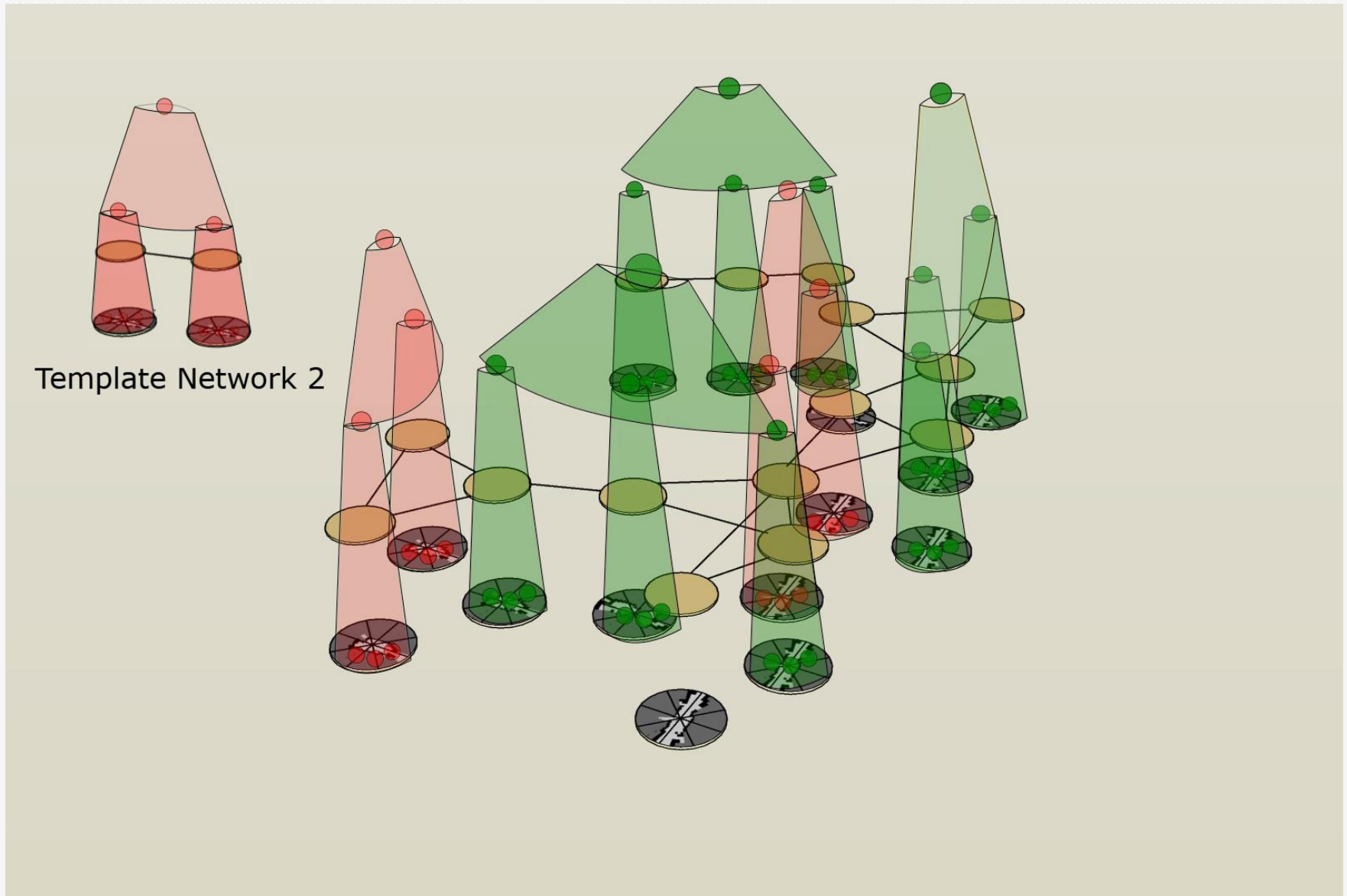
END2END: INFERENCE

[Zheng, Pronobis '19]



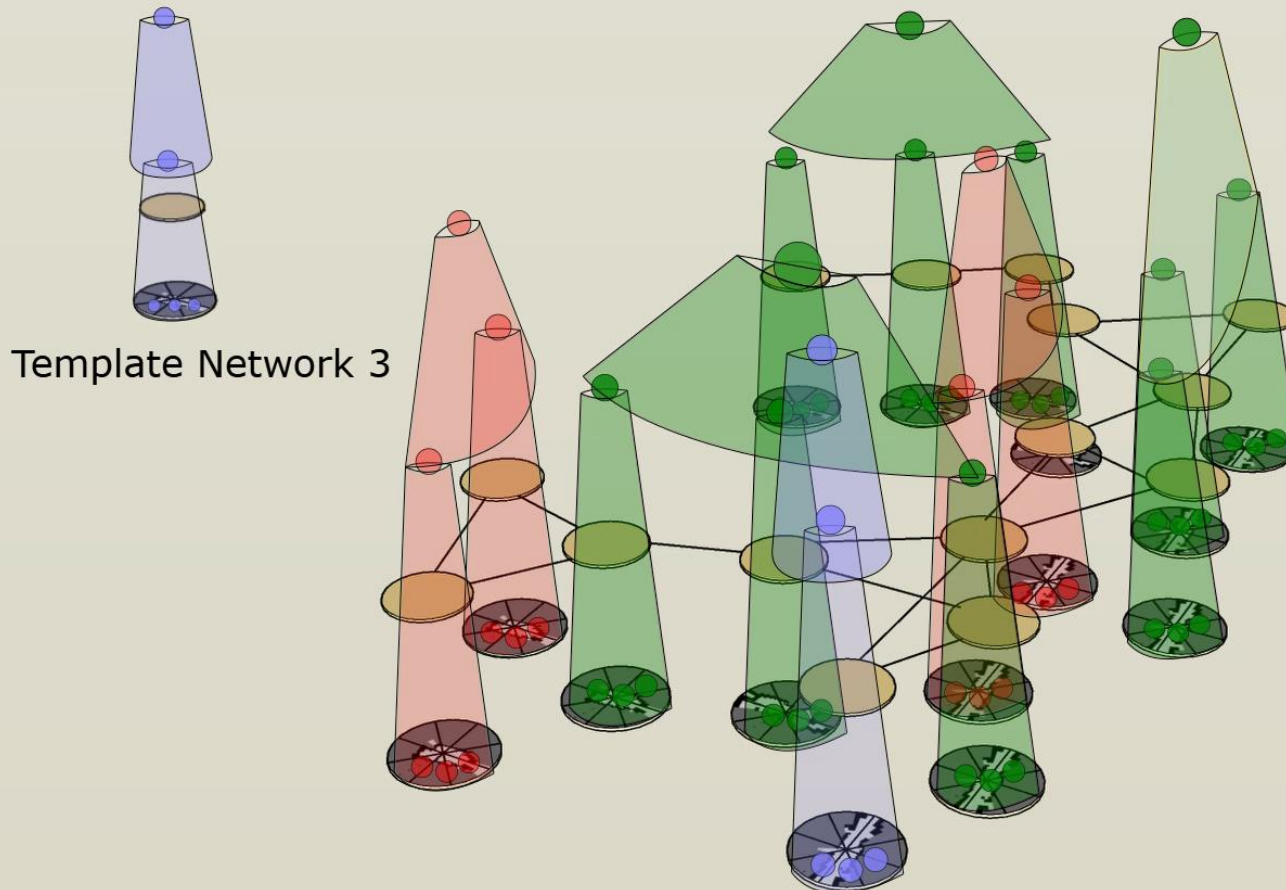
END2END: INFERENCE

[Zheng, Pronobis '19]



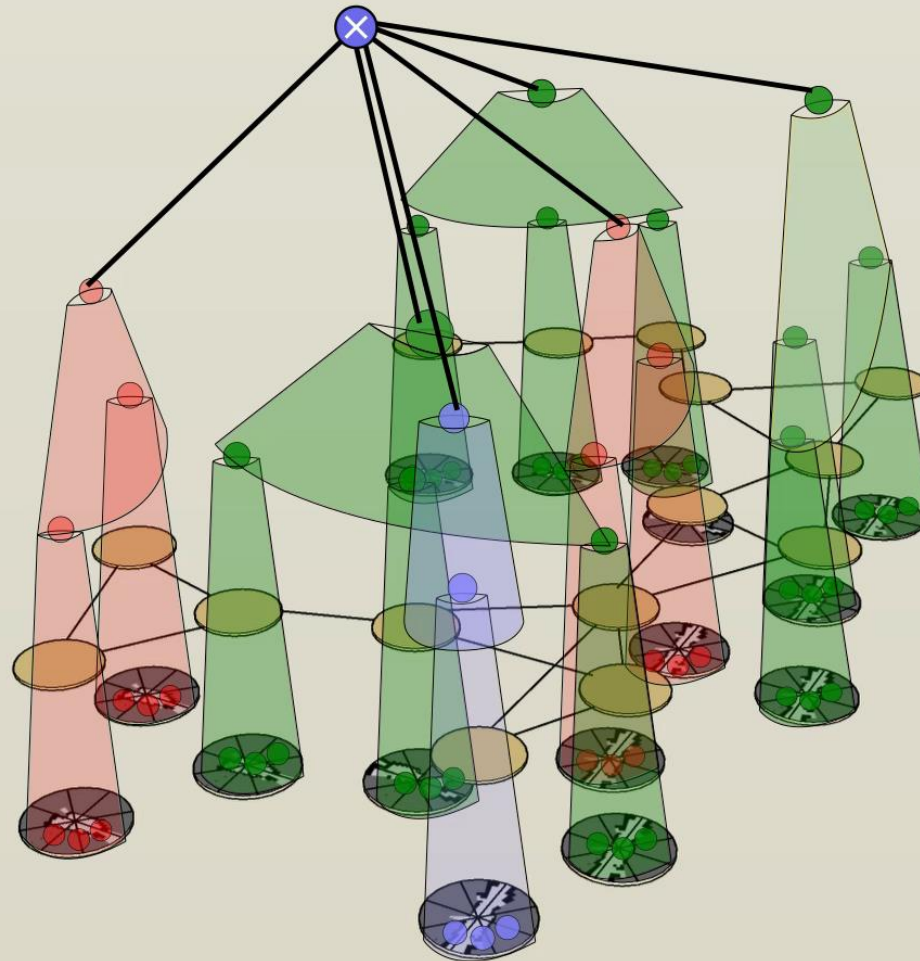
END2END: INFERENCE

[Zheng, Pronobis '19]



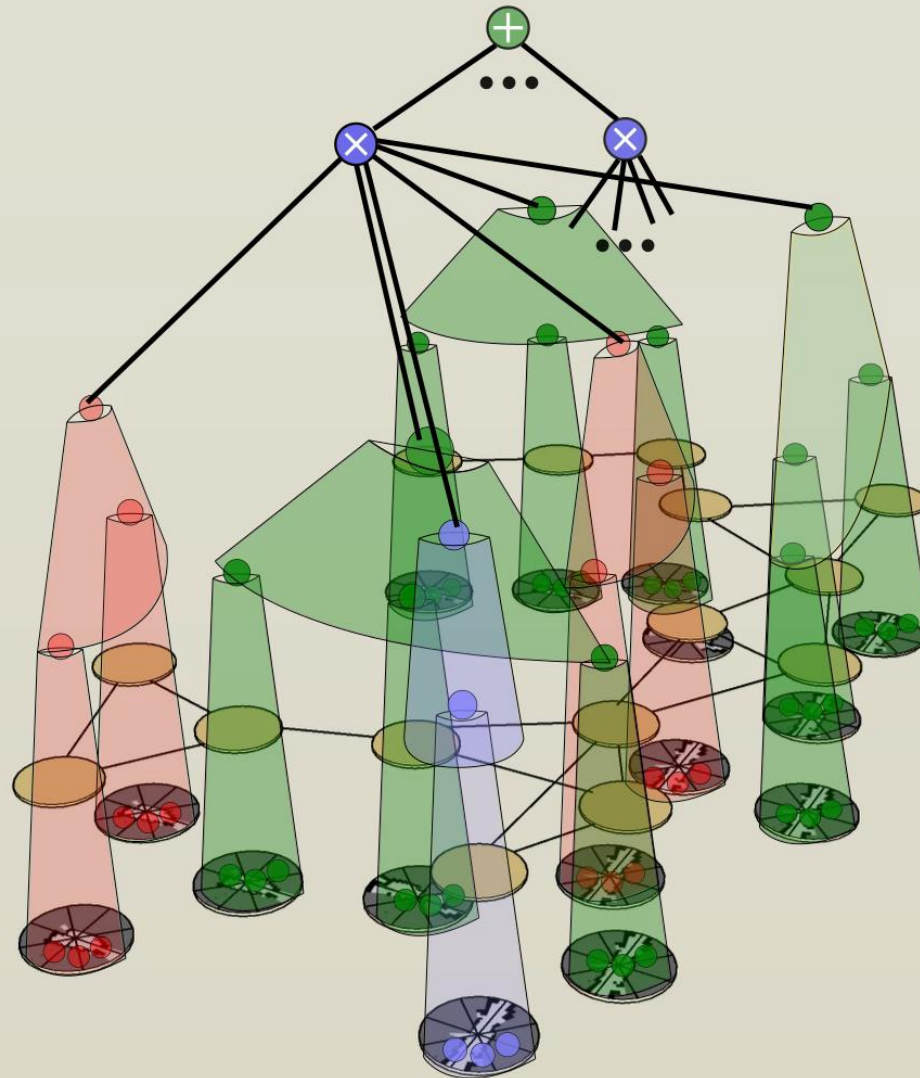
END2END: INFERENCE

[Zheng, Pronobis '19]



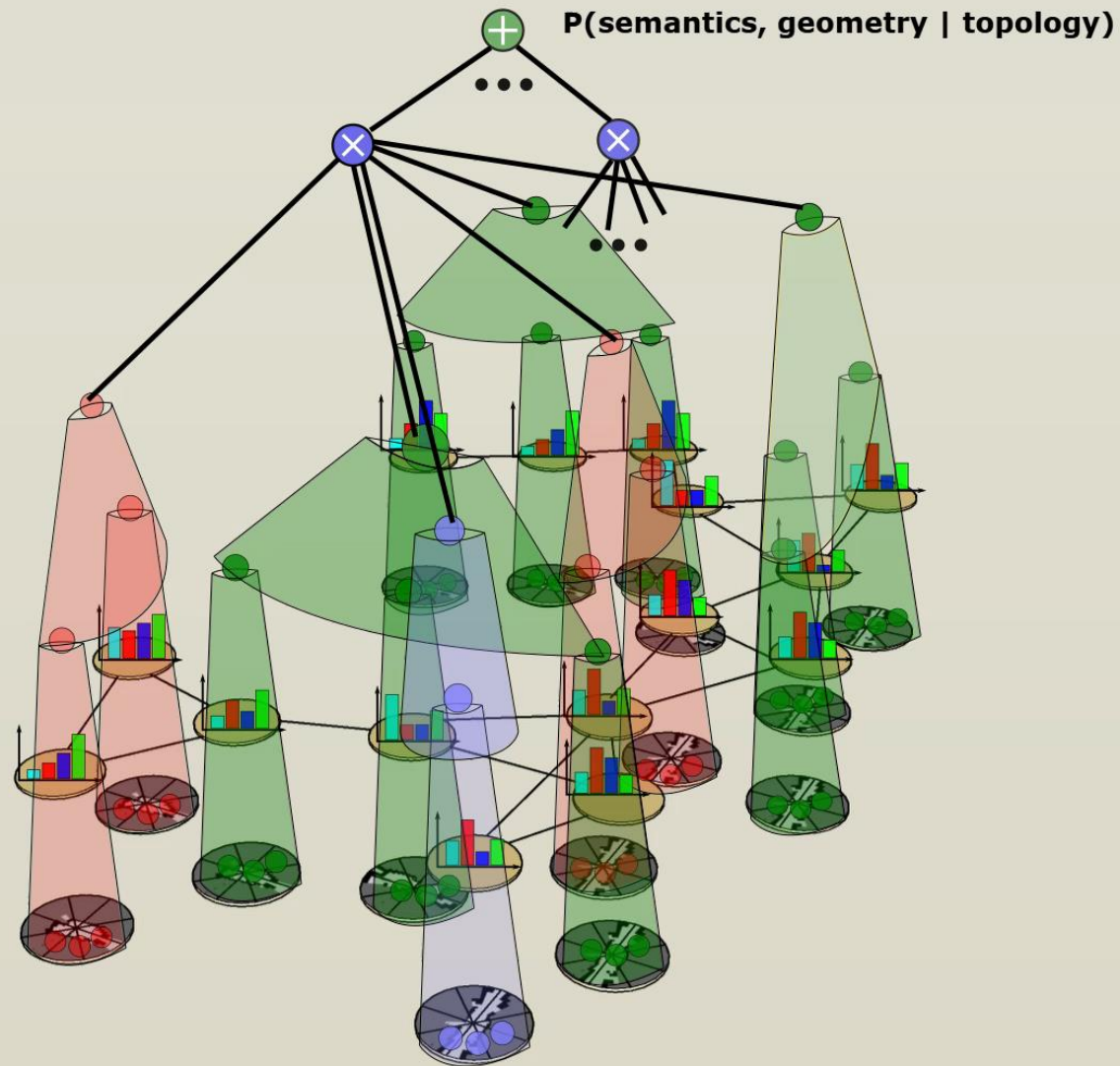
END2END: INFERENCE

[Zheng, Pronobis '19]



END2END: INFERENCE

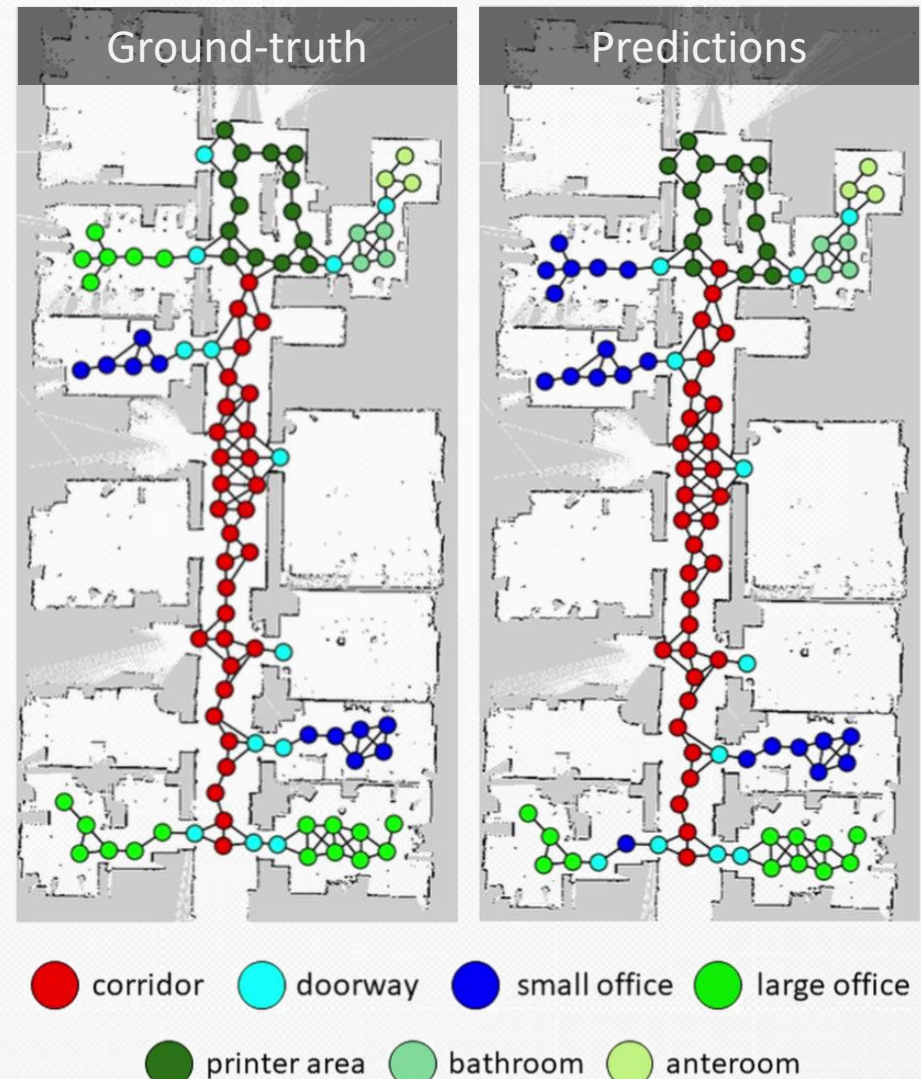
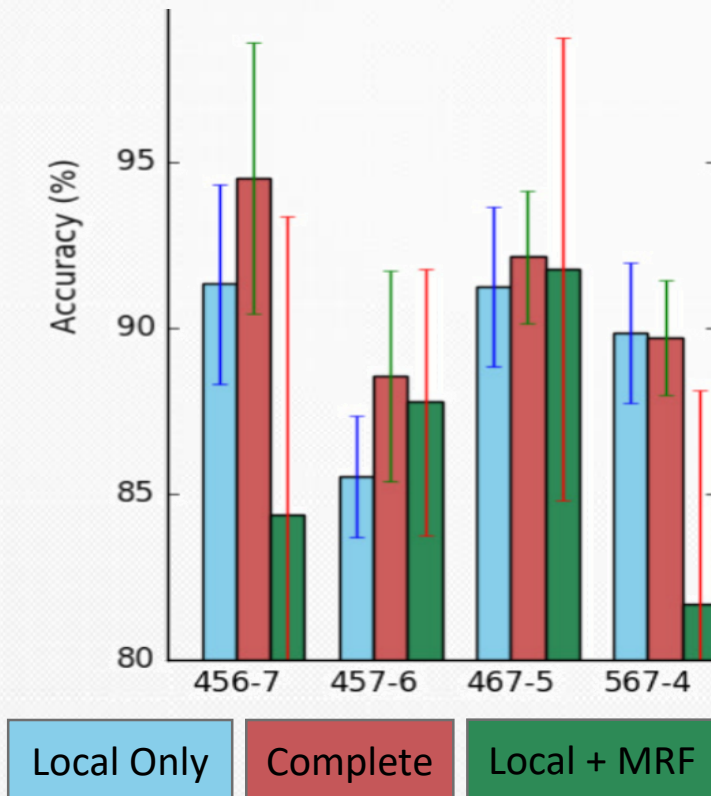
[Zheng, Pronobis '19]



END2END: SEMANTIC MAPPING

[Zheng, Pronobis '19]

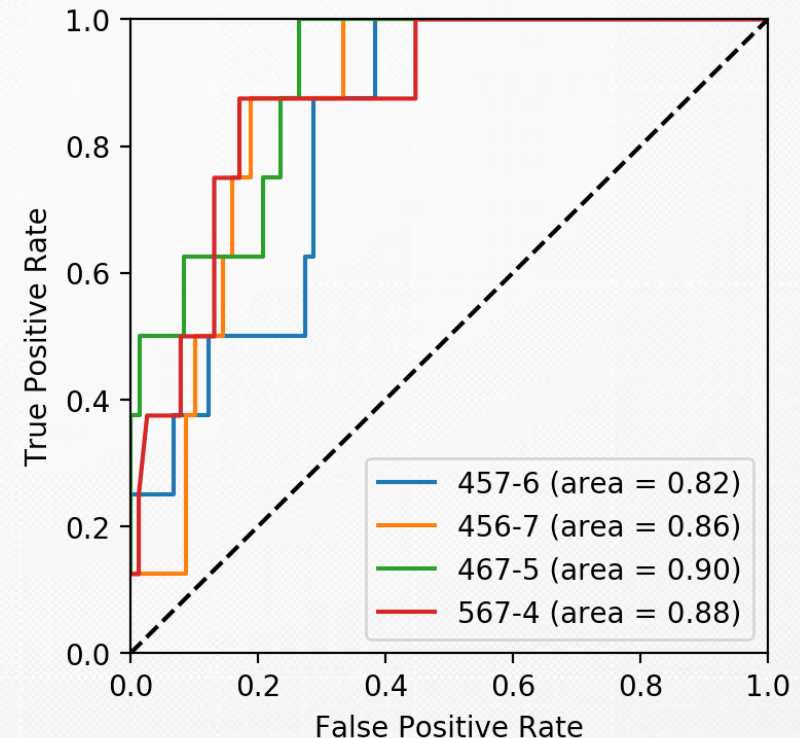
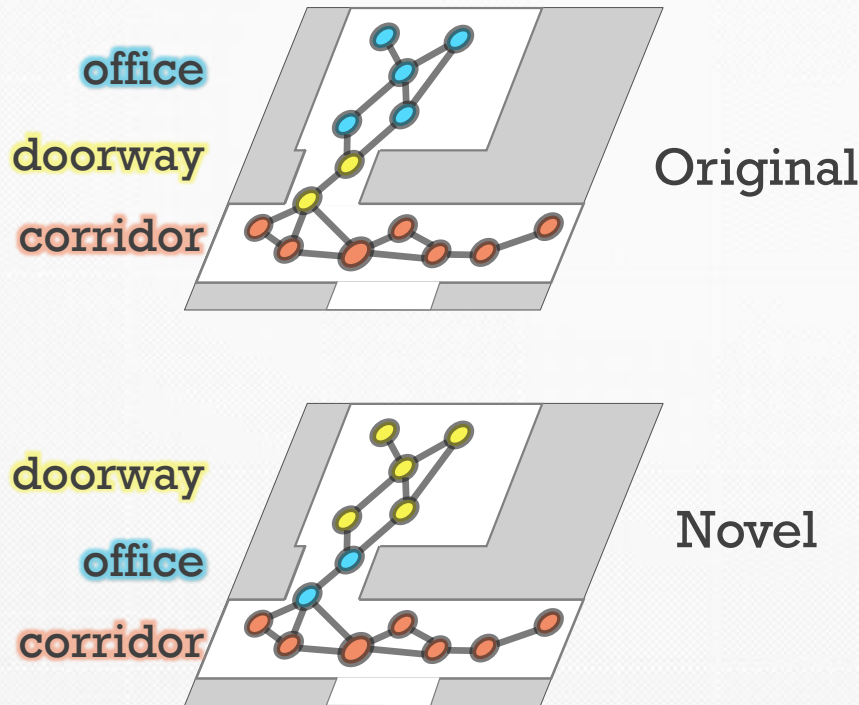
- Office building, 4 floors
 - 7 known room classes
- Overall: 93%



END2END: NOVEL GLOBAL STRUCTURE DETECTION

[Zheng, Pronobis '19]

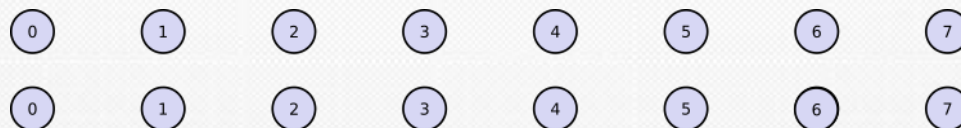
- $P(\text{semantics, geometry}) > \text{threshold} ?$
- Novel floor structures by swapping labels of predictions for two random rooms in test map



SPATIAL SPNS FOR VISUAL SIGNALS

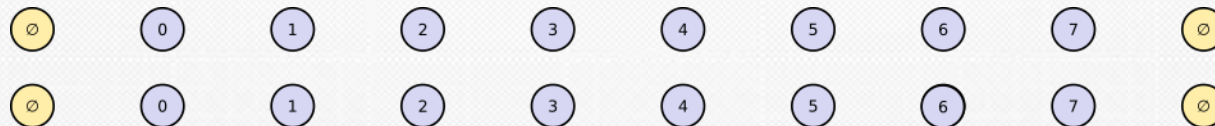
[van de Wolfshaar, Pronobis '19]

- New SPNs architecture
 - Resembles convolutional neural networks
 - For spatial signals, such as images
- 1D example:
 - Input layer of size 8 with 2 channels/pixel



SPATIAL SPNS FOR VISUAL SIGNALS

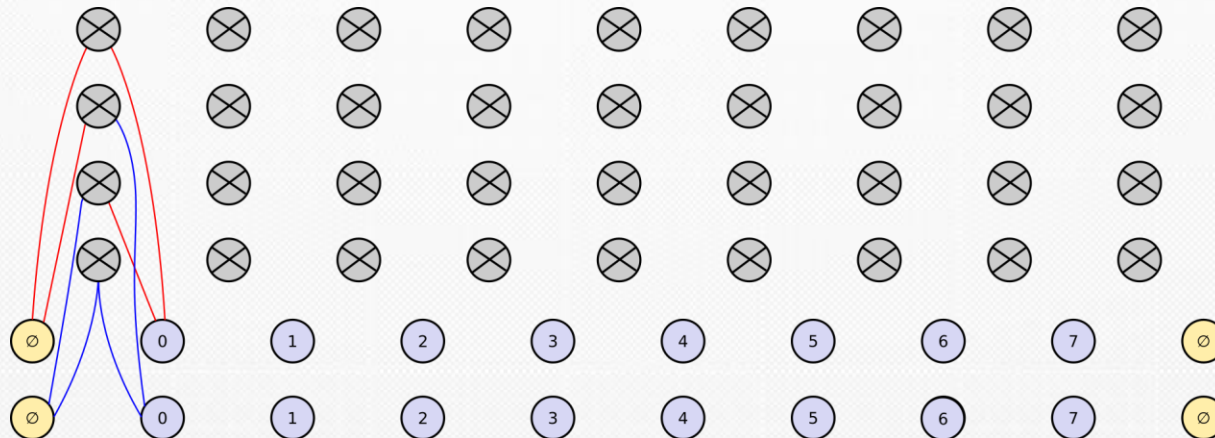
[van de Wolfshaar, Pronobis '19]



SPATIAL SPNS FOR VISUAL SIGNALS

[van de Wolfshaar, Pronobis '19]

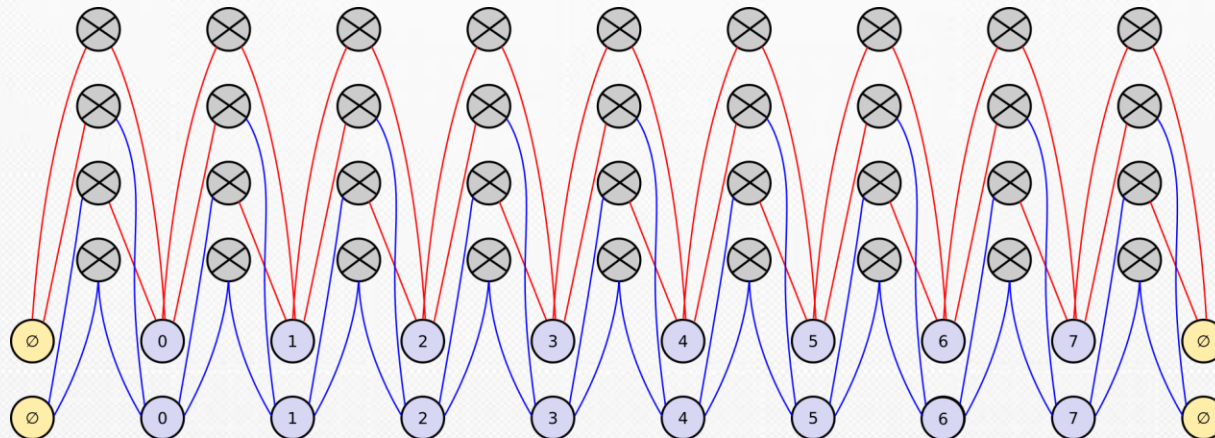
- Spatial products: stride 1, kernel size 2



SPATIAL SPNS FOR VISUAL SIGNALS

[van de Wolfshaar, Pronobis '19]

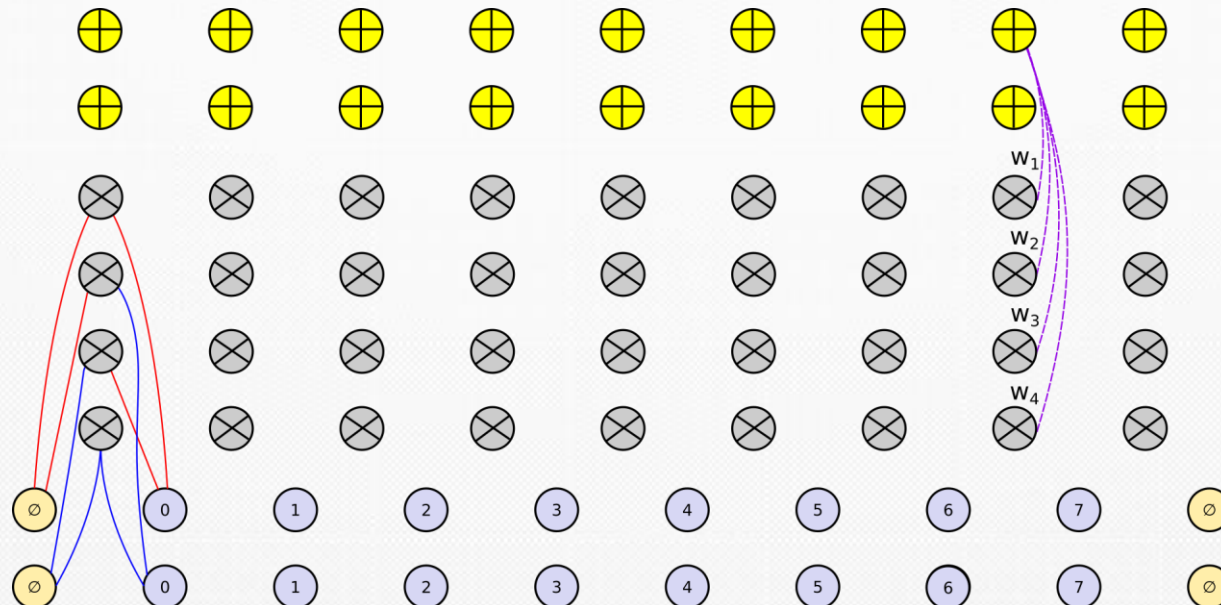
- Spatial products: stride 1, kernel size 2



SPATIAL SPNS FOR VISUAL SIGNALS

[van de Wolfshaar, Pronobis '19]

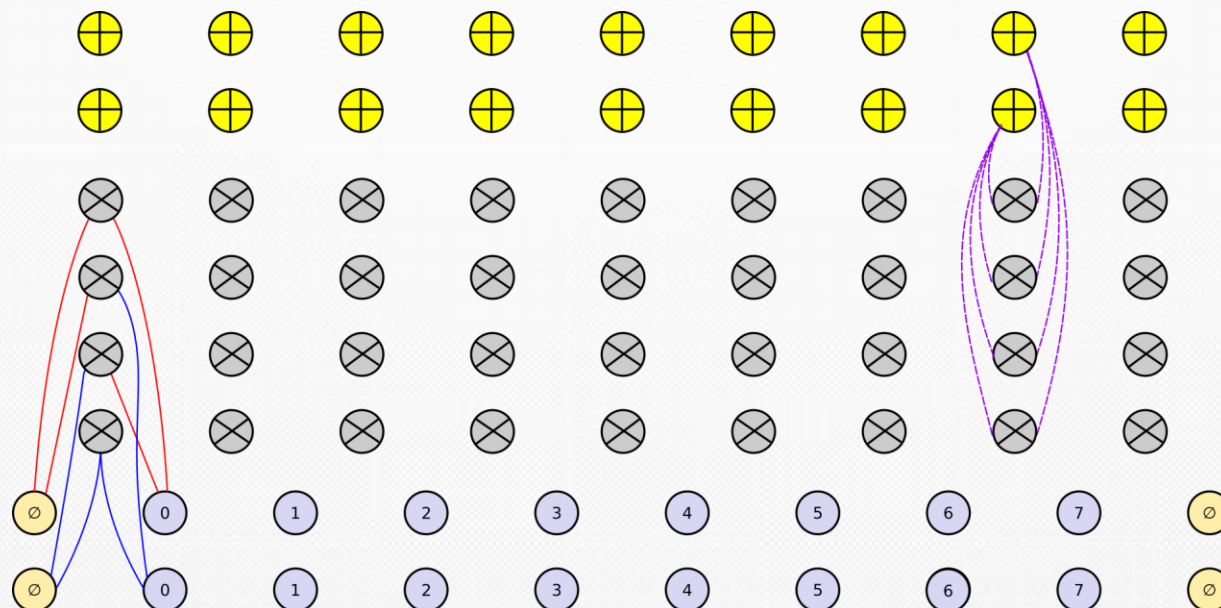
- Sums: 1x1 convolutions



SPATIAL SPNS FOR VISUAL SIGNALS

[van de Wolfshaar, Pronobis '19]

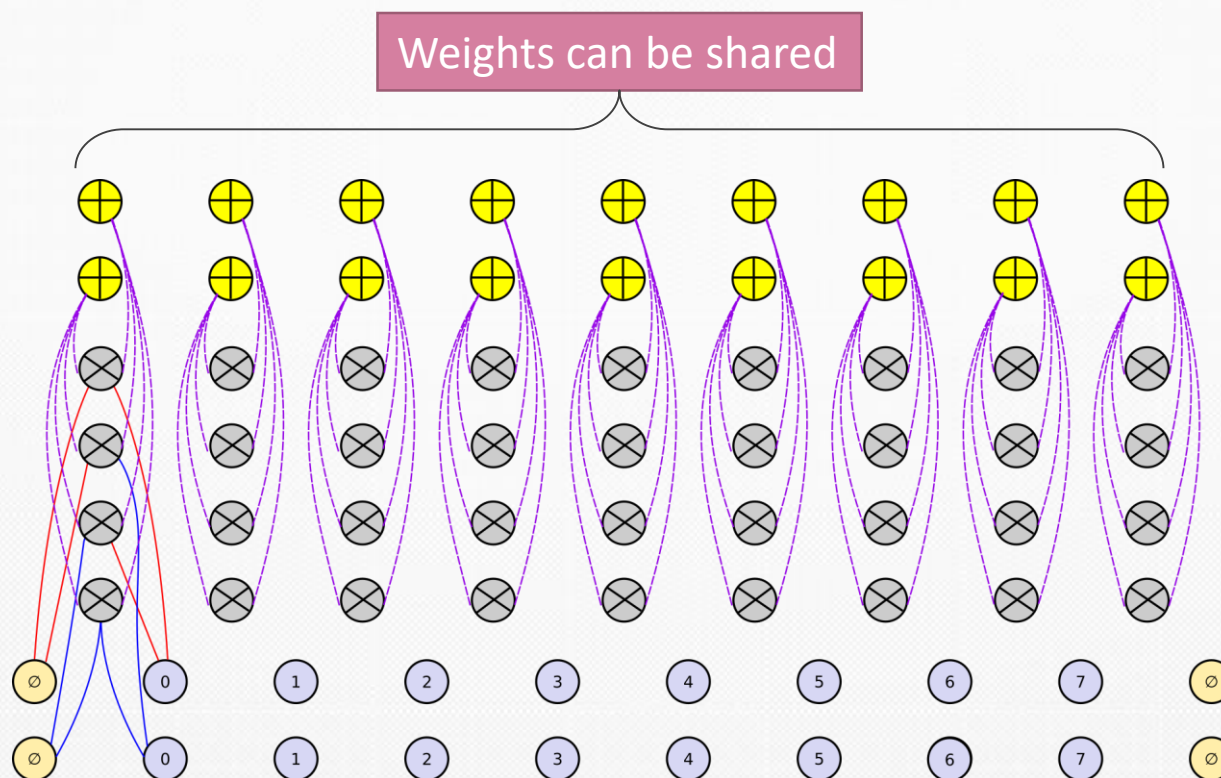
- Sums: 1x1 convolutions



SPATIAL SPNS FOR VISUAL SIGNALS

[van de Wolfshaar, Pronobis '19]

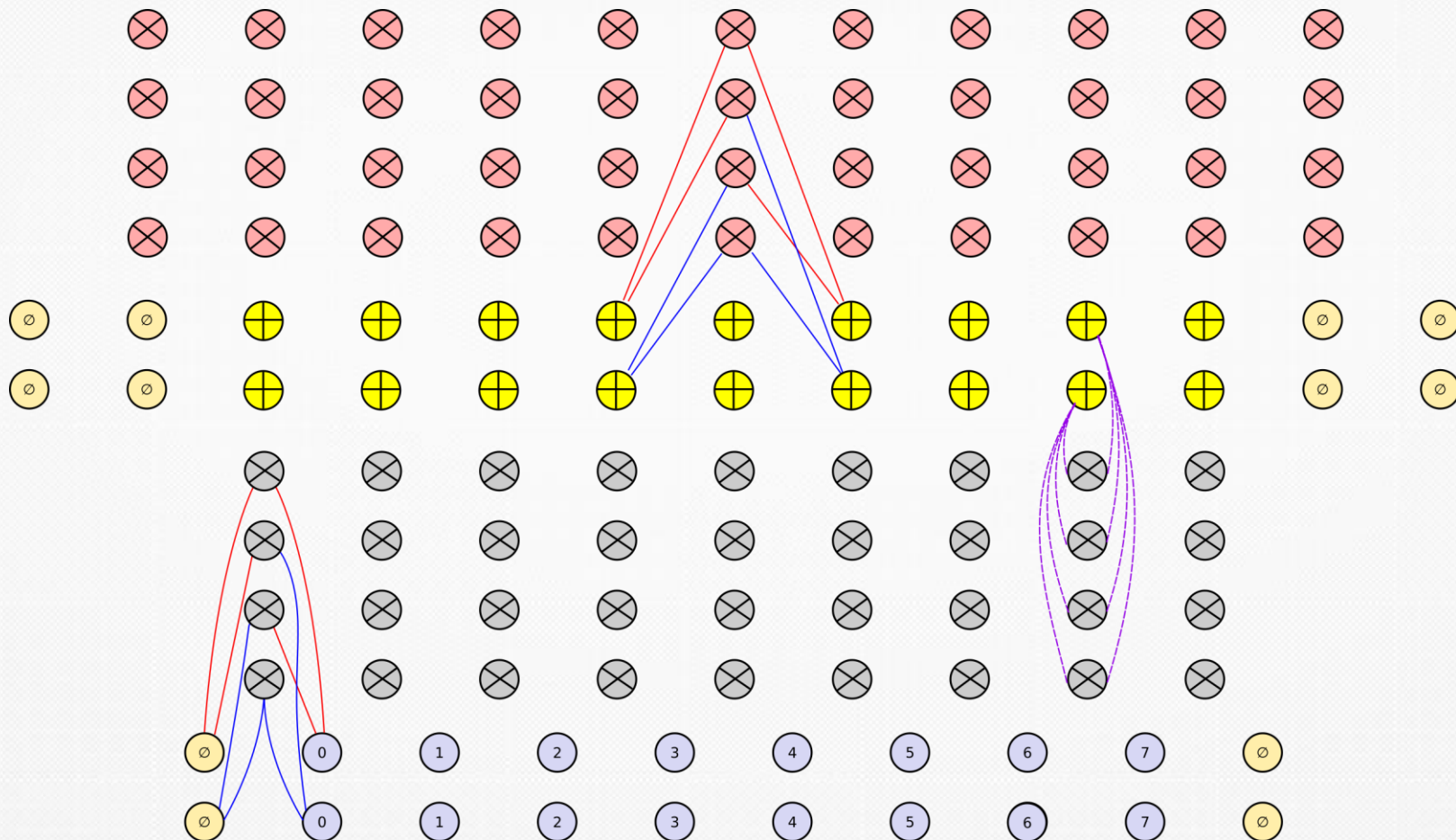
- Sums: 1x1 convolutions



SPATIAL SPNS FOR VISUAL SIGNALS

[van de Wolfshaar, Pronobis '19]

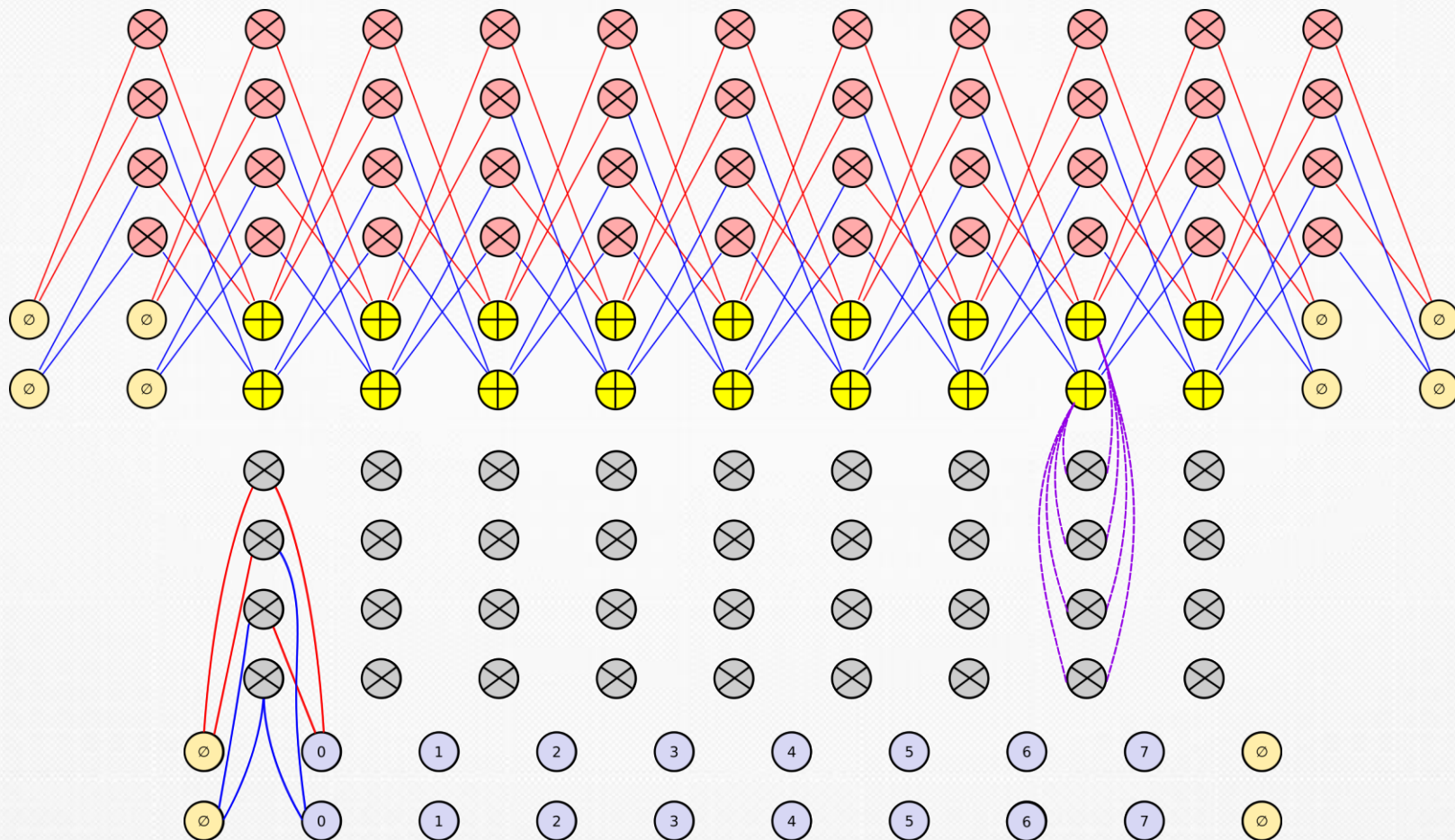
- Spatial products: stride 1, dilation rate **2**



SPATIAL SPNS FOR VISUAL SIGNALS

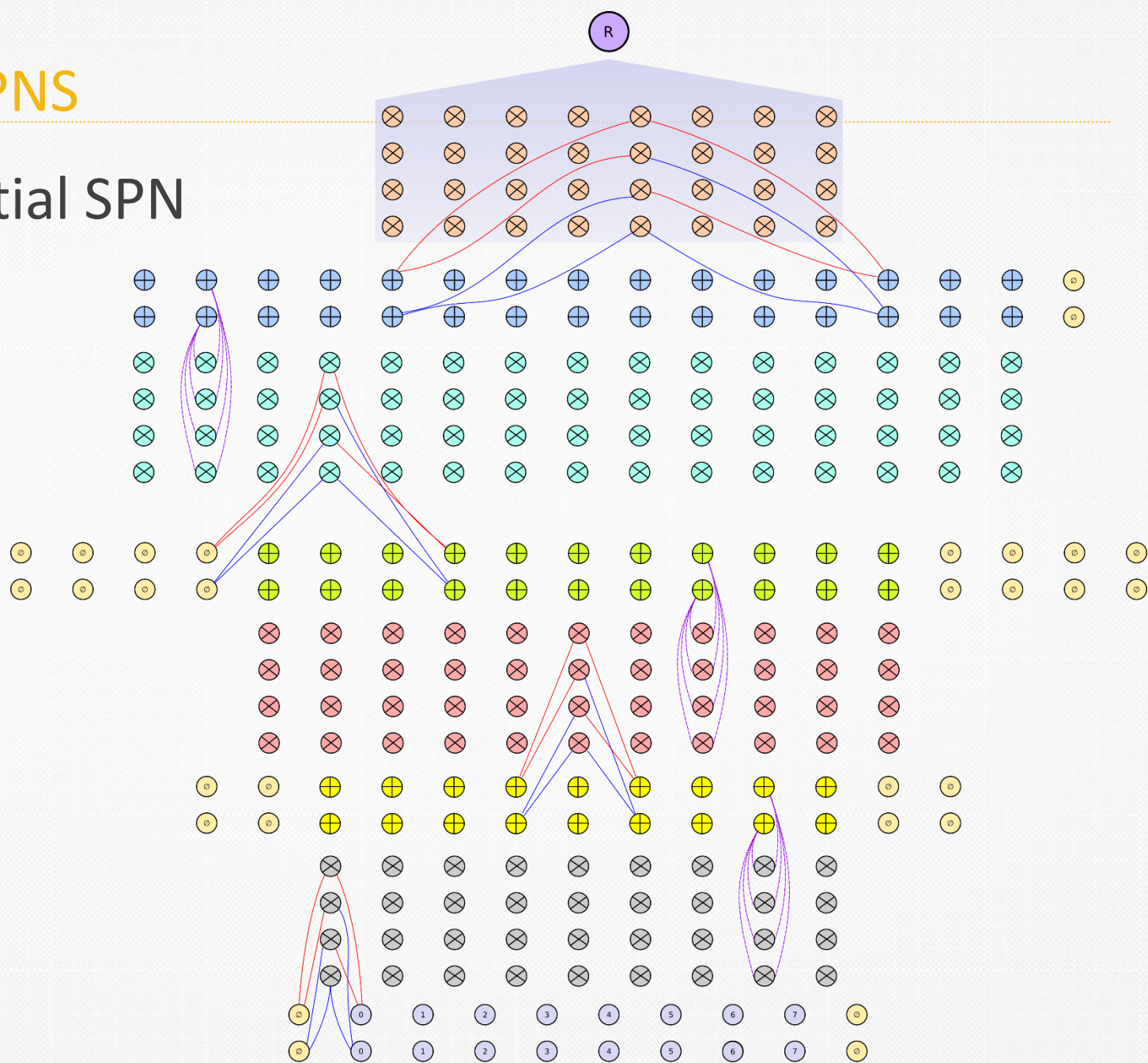
[van de Wolfshaar, Pronobis '19]

- Spatial products: stride 1, dilation rate **2**



SPATIAL SPNS

- Full Spatial SPN



SPATIAL SPNS: INITIAL RESULTS

[van de Wolfshaar, Pronobis '19]

- Dataset: MNIST

Model	Accuracy
RAT-SPN + discriminative GD [Peharz et al. 2018]	98.1 %
Vertical/horizontal splits + discriminative EBW [Rashwan et al. 2018]	95.0 %
Spatial SPN + discriminative GD	98.9 %
Spatial SPN + generative hard EM	96.1 %

CONCLUSIONS

- Comprehensive spatial representation and its realization using deep probabilistic networks
 - Learns general relationships
 - Between place geometry and semantic descriptions
 - Between pixels of places to topology of buildings
 - Up & down inferences across levels of abstraction
- Pioneers SPNs in the domain of robotics
 - SOTA model for complex structured prediction
 - Even in high-dimensional visual data
- Designed to support planning and communication
- Ongoing work
 - Probabilistic planning in the same architecture
 - Multi-modal representation: visual and depth sensors



THANK YOU

www.pronobis.pro