



# What can go here?

**Evangelos Kalogerakis**

[based on “SceneGraphNet: Neural Message Passing for 3D Indoor Scene Augmentation”,  
Yang Zhou, Zachary While, Evangelos Kalogerakis, ICCV 2019]



# What can go here?

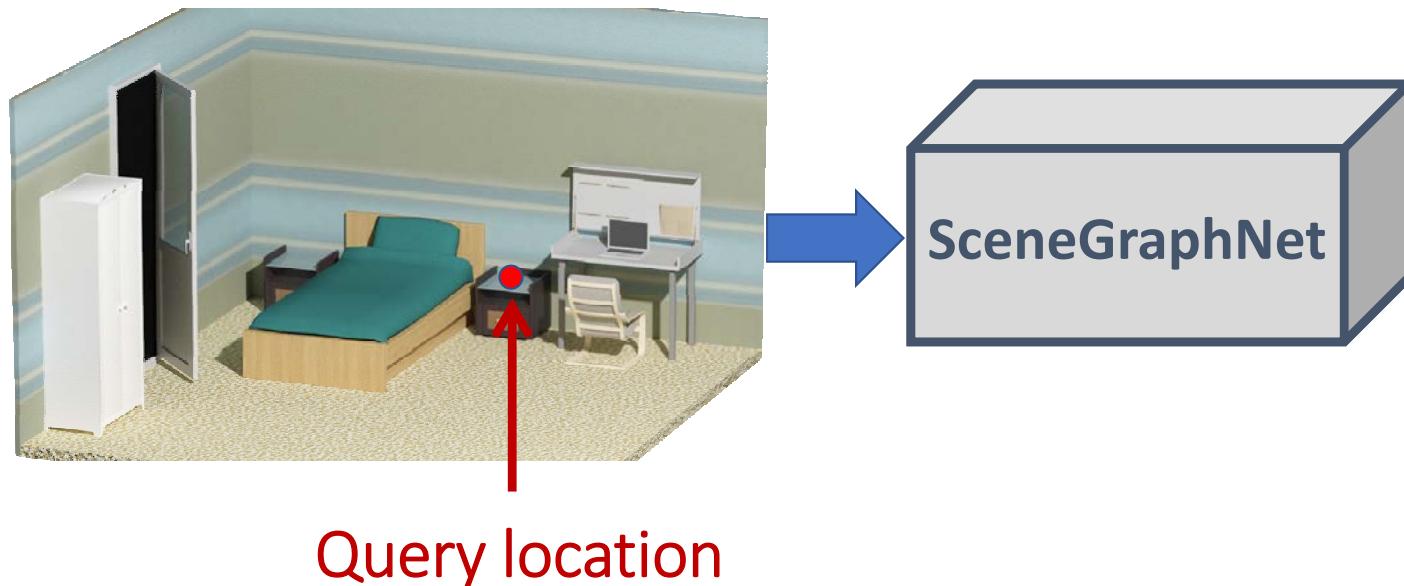
Evangelos Kalogerakis

[based on “SceneGraphNet: Neural Message Passing for 3D Indoor Scene Augmentation”,  
Yang Zhou, Zachary While, Evangelos Kalogerakis, ICCV 2019]



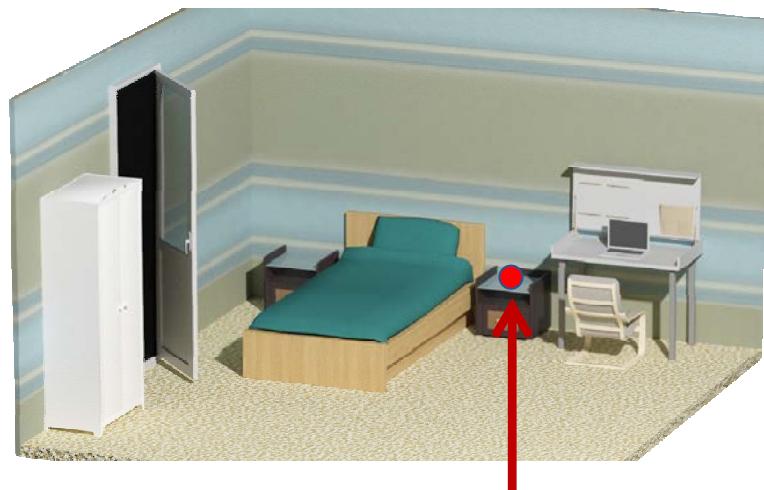
Goal: predict objects to populate indoor scenes

3D scene model



Goal: predict objects to populate indoor scenes

3D scene model



Query location

$P(\text{object type} \mid \text{query, scene})$

SceneGraphNet

lamp

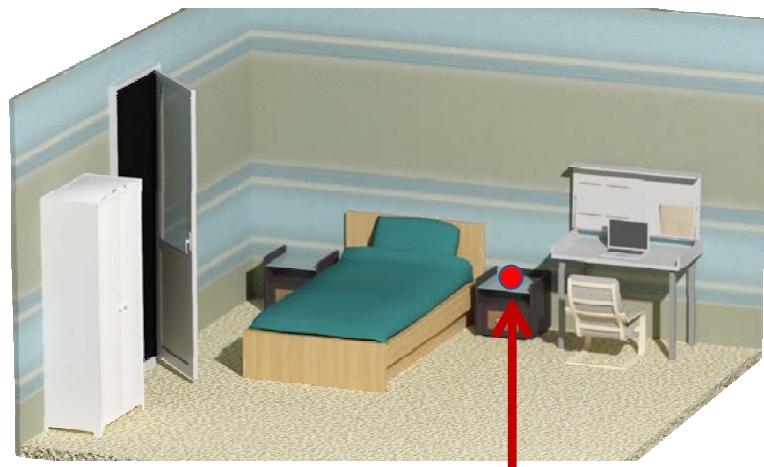
alarm clock

books

vase

Goal: predict objects to populate indoor scenes

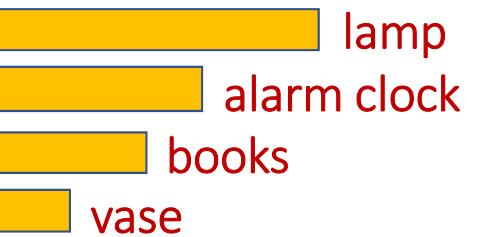
3D scene model



Query location

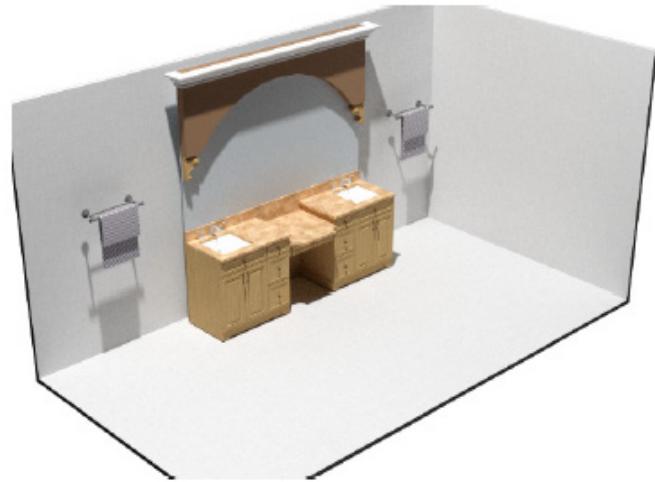
SceneGraphNet

$P(\text{object type} \mid \text{query, scene})$

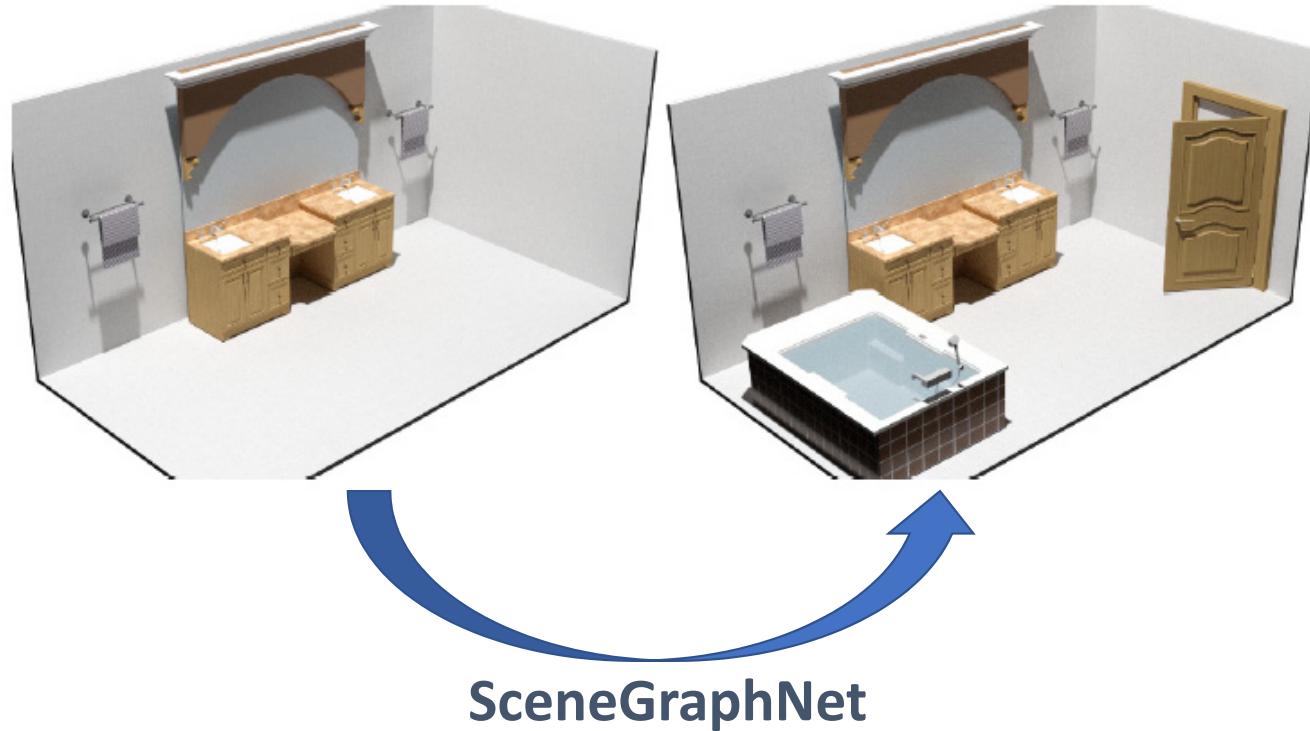


+ object  
dimensions

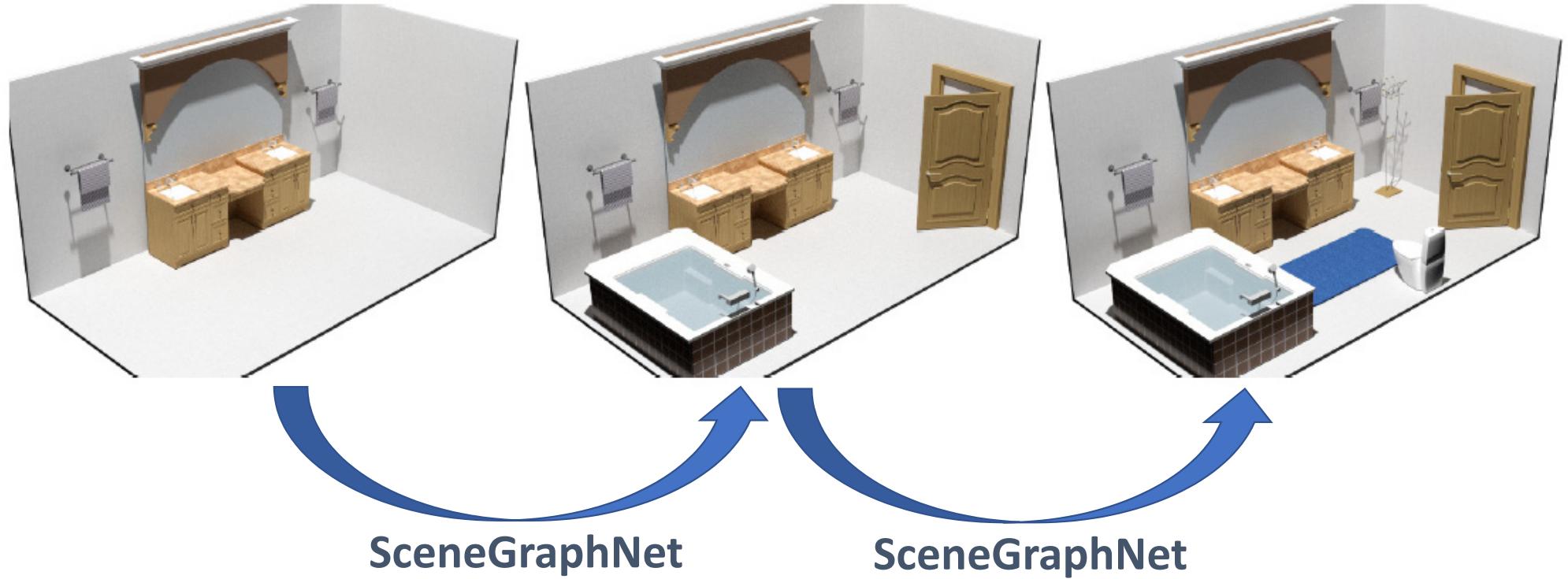
# Alternative goal: iterative scene synthesis



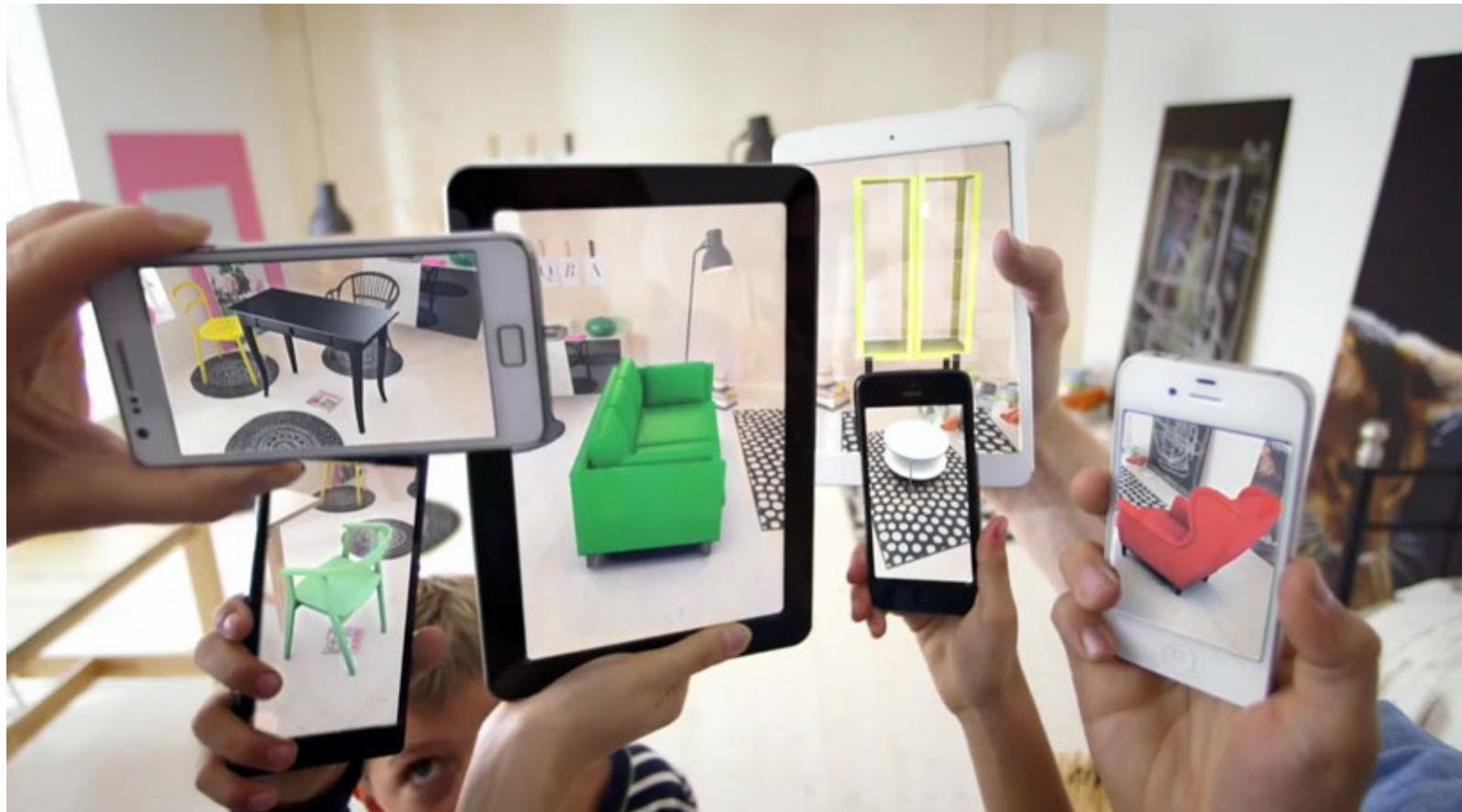
# Alternative goal: iterative scene synthesis



Alternative goal: iterative scene synthesis



# Motivation: AR/VR applications for indoor design



## Prior work: learning-based scene completion

- Probabilistic grammars e.g., [Fisher et al. 12, Qi et al. 18, Purkait et al. 19]

## Prior work: learning-based scene completion

- Probabilistic grammars e.g., [**Fisher et al. 12, Qi et al. 18, Purkait et al. 19**]
- View-based, image convolutional methods e.g., [**Wang et al. 18, Ritchie et al 19**]

## Prior work: learning-based scene completion

- Probabilistic grammars e.g., [Fisher et al. 12, Qi et al. 18, Purkait et al. 19]
- View-based, image convolutional methods e.g., [Wang et al. 18, Ritchie et al 19]
- Volumetric methods e.g., [Song et al. 16]

## Prior work: learning-based scene completion

- Probabilistic grammars e.g., [Fisher et al. 12, Qi et al. 18, Purkait et al. 19]
- View-based, image convolutional methods e.g., [Wang et al. 18, Ritchie et al 19]
- Volumetric methods e.g., [Song et al. 16]
- GANs on matrix representation of scene objects e.g., [Zhang et al. 20]

## Prior work: learning-based scene completion

- Probabilistic grammars e.g., [Fisher et al. 12, Qi et al. 18, Purkait et al. 19]
- View-based, image convolutional methods e.g., [Wang et al. 18, Ritchie et al 19]
- Volumetric methods e.g., [Song et al. 16]
- GANs on matrix representation of scene objects e.g., [Zhang et al. 20]
- Recursive autoencoders on trees graph representations e.g., [Li et al. 19]

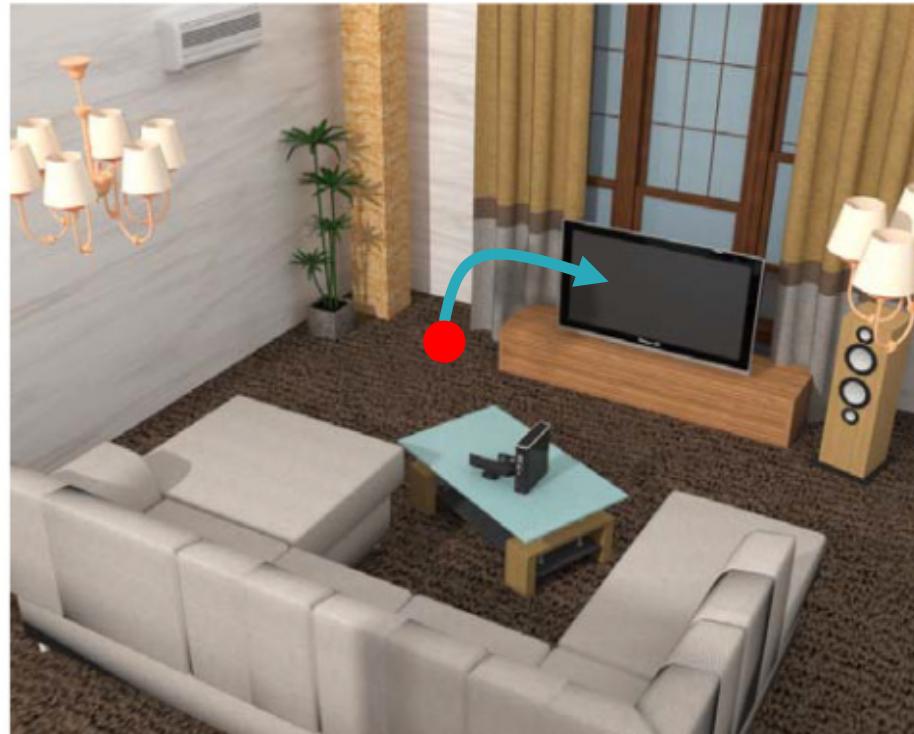
# Prior work: learning-based scene completion

- Probabilistic grammars e.g., [Fisher et al. 12, Qi et al. 18, Purkait et al. 19]
- View-based, image convolutional methods e.g., [Wang et al. 18, Ritchie et al 19]
- Volumetric methods e.g., [Song et al. 16]
- GANs on matrix representation of scene objects e.g., [Zhang et al. 20]
- Recursive autoencoders on trees graph representations e.g., [Li et al. 19]
- GNNs on object relation graphs [Wang et al. 19]

Key idea: capture surrounding object relations

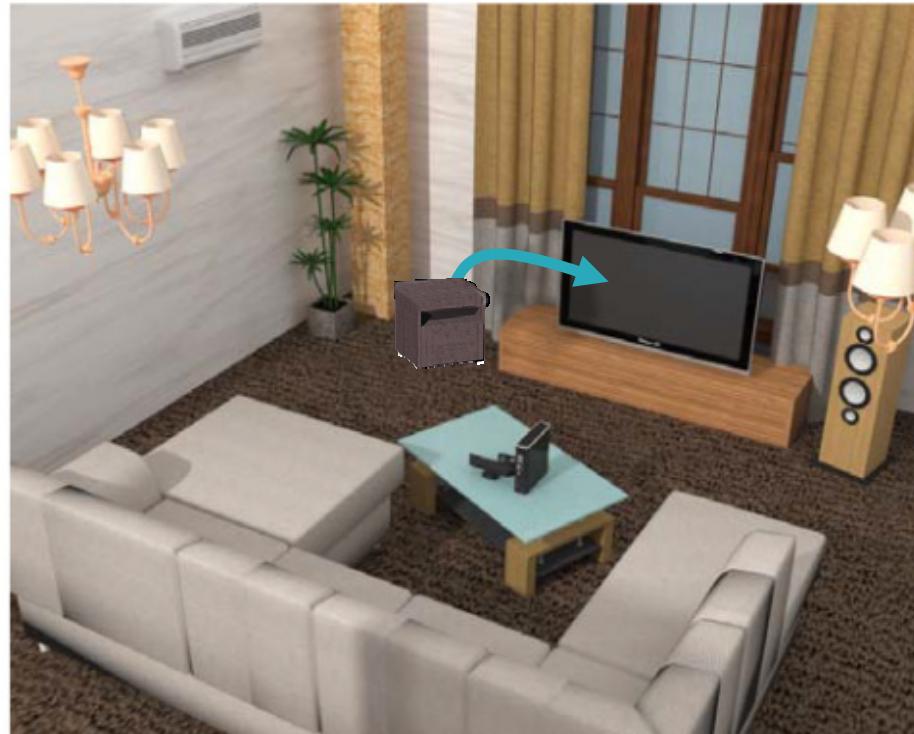


Key idea: capture surrounding object relations



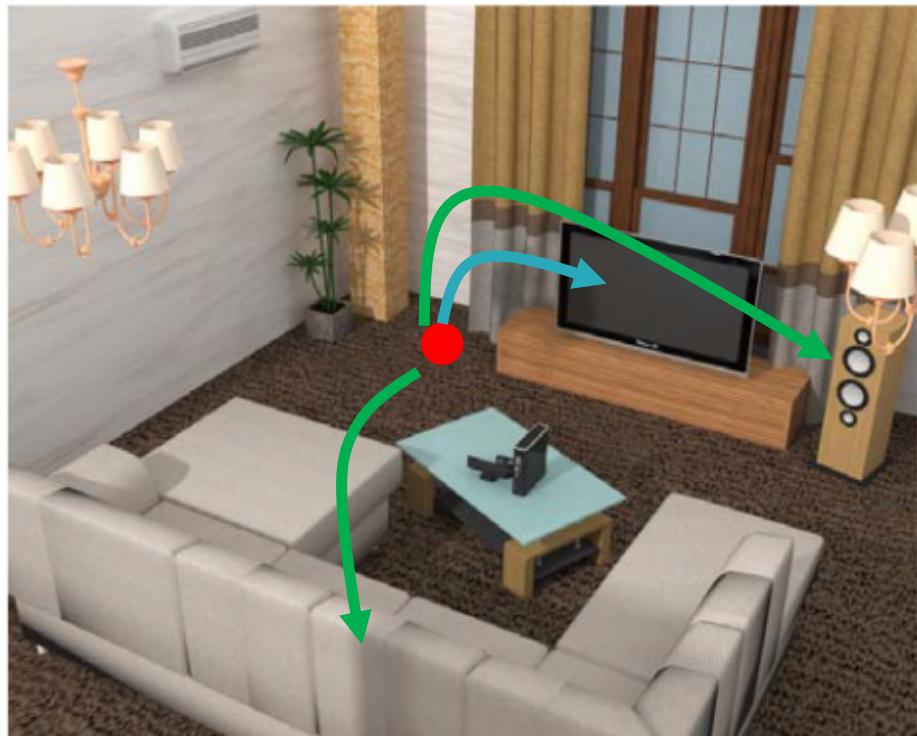
**Short-range relations?**

Key idea: capture surrounding object relations



**Short-range relations?**

Key idea: capture surrounding object relations



**Short-range relations**

+

**Longer-range relations**

Key idea: capture surrounding object relations

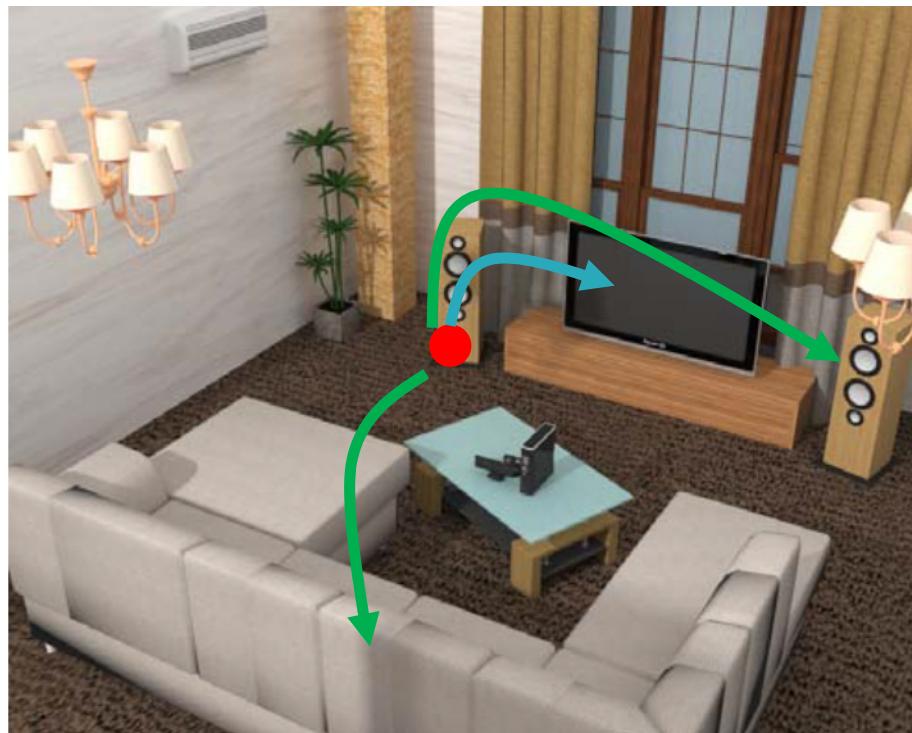


**Short-range relations**

+

**Longer-range relations**

# Key Challenges



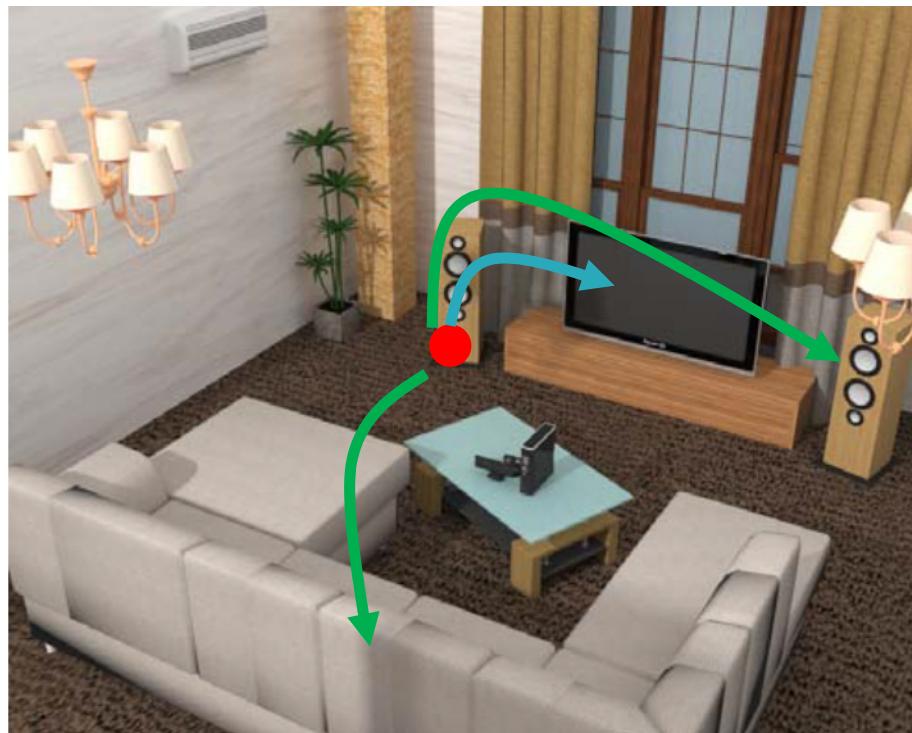
**Short-range relations**

+

**Longer-range relations**

**How to learn these relations?**

# Key Challenges



**Short-range relations**

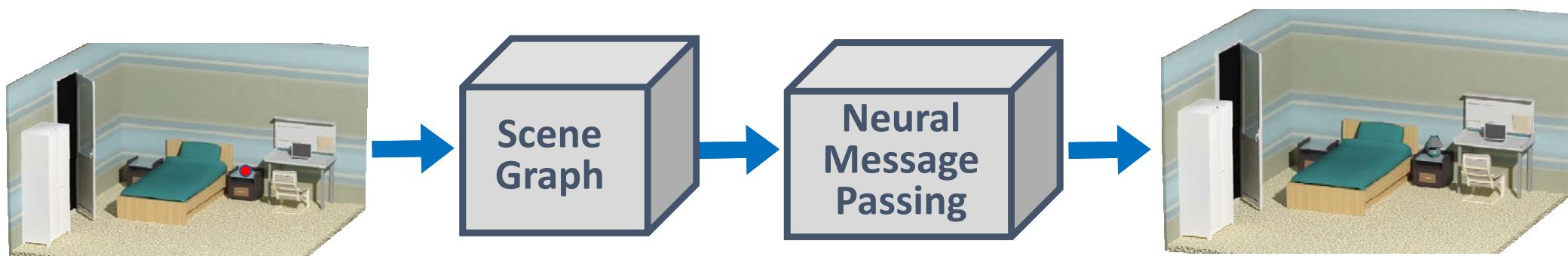
+

**Longer-range relations**

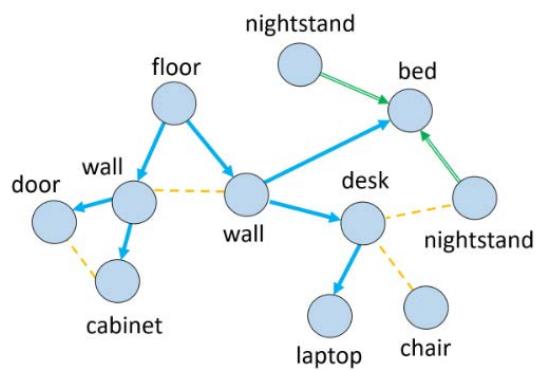
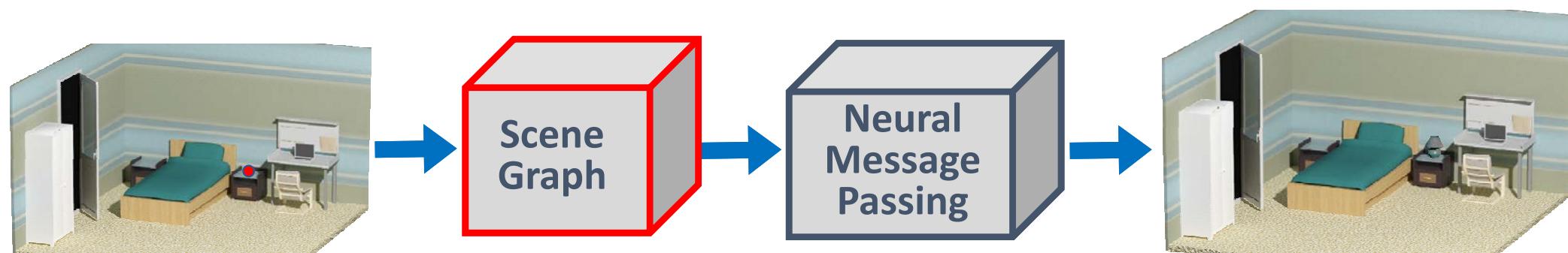
**How to learn these relations?**

**How to weigh & aggregate them?**

# Pipeline



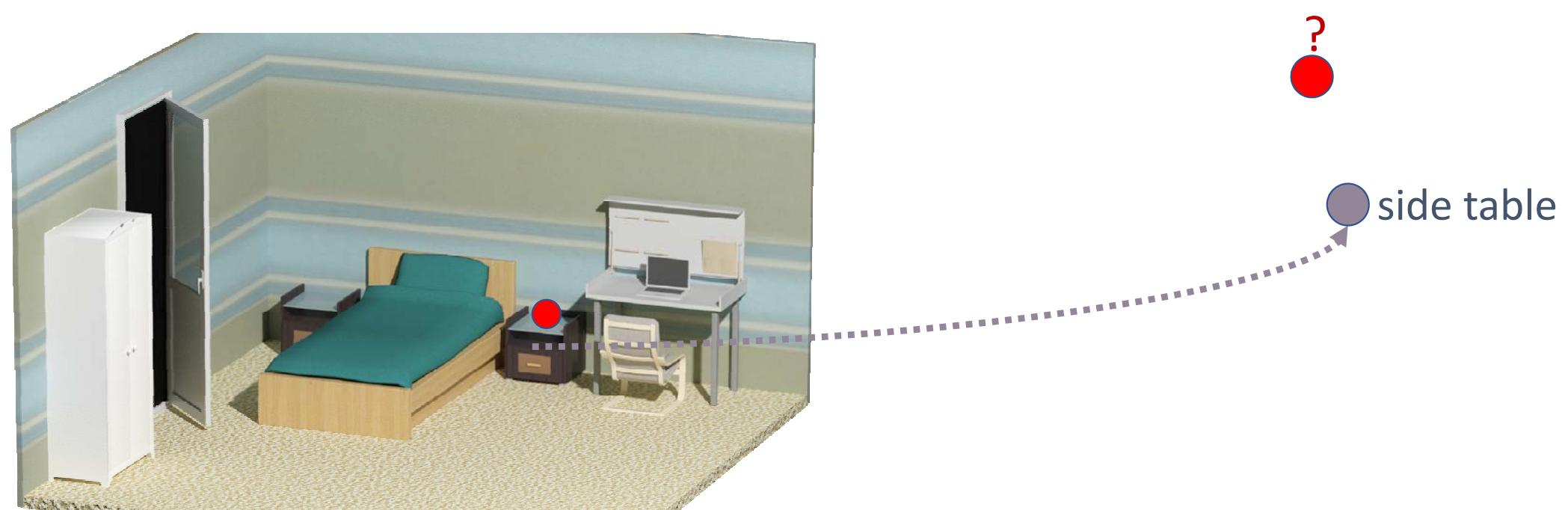
# Pipeline



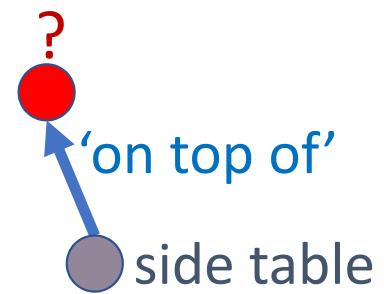
# Scene graph construction



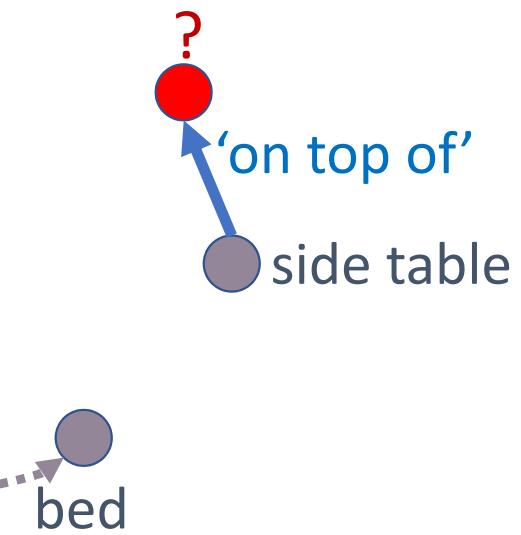
Nodes represent existing objects + queries



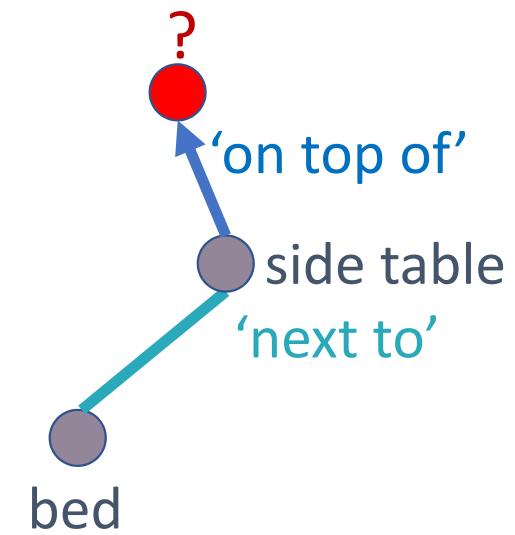
# Edges represent relationships



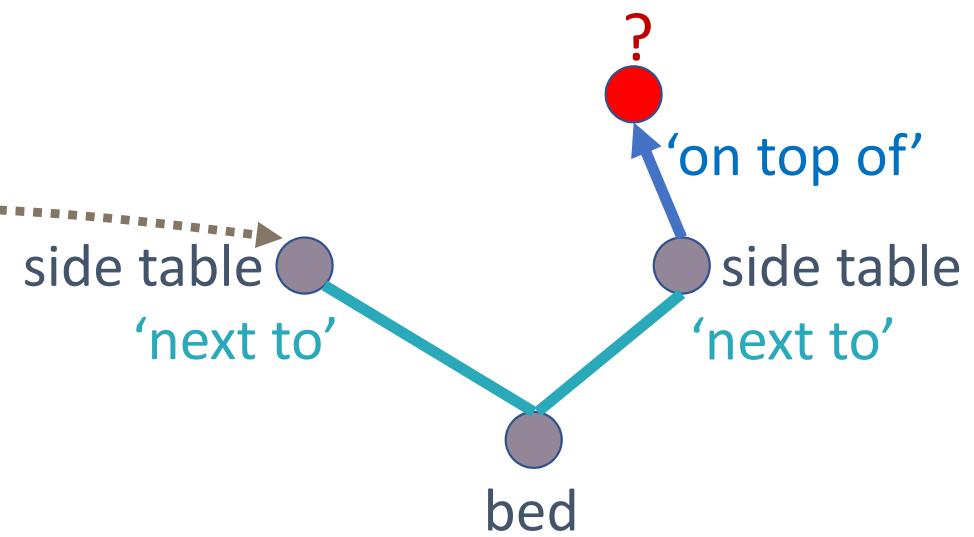
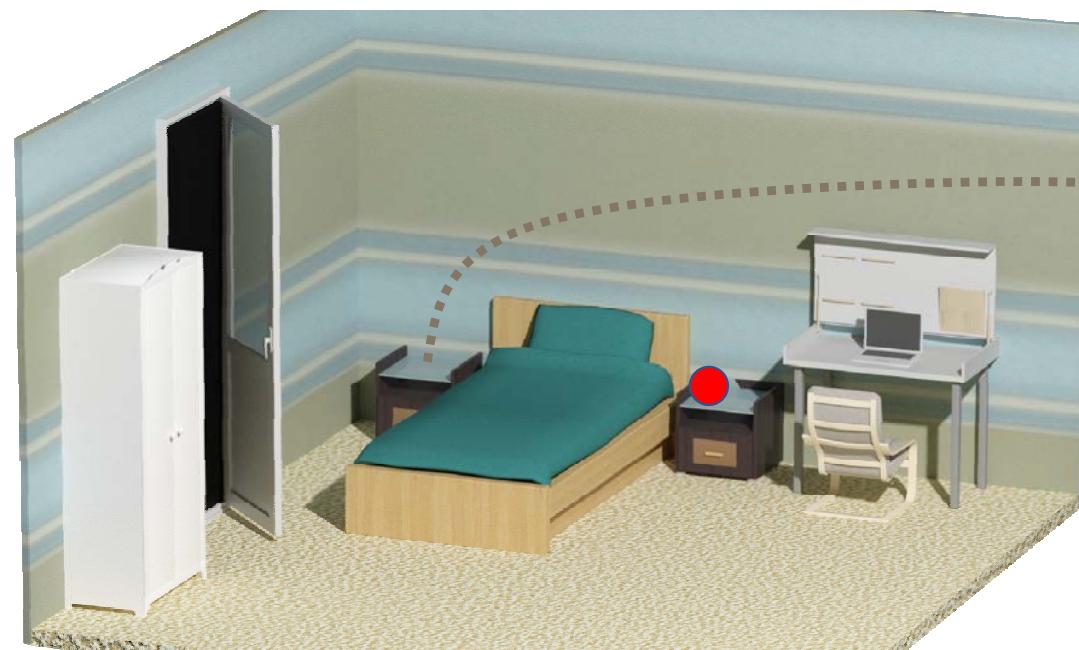
Bed: “I am also here – I may affect decisions, too”



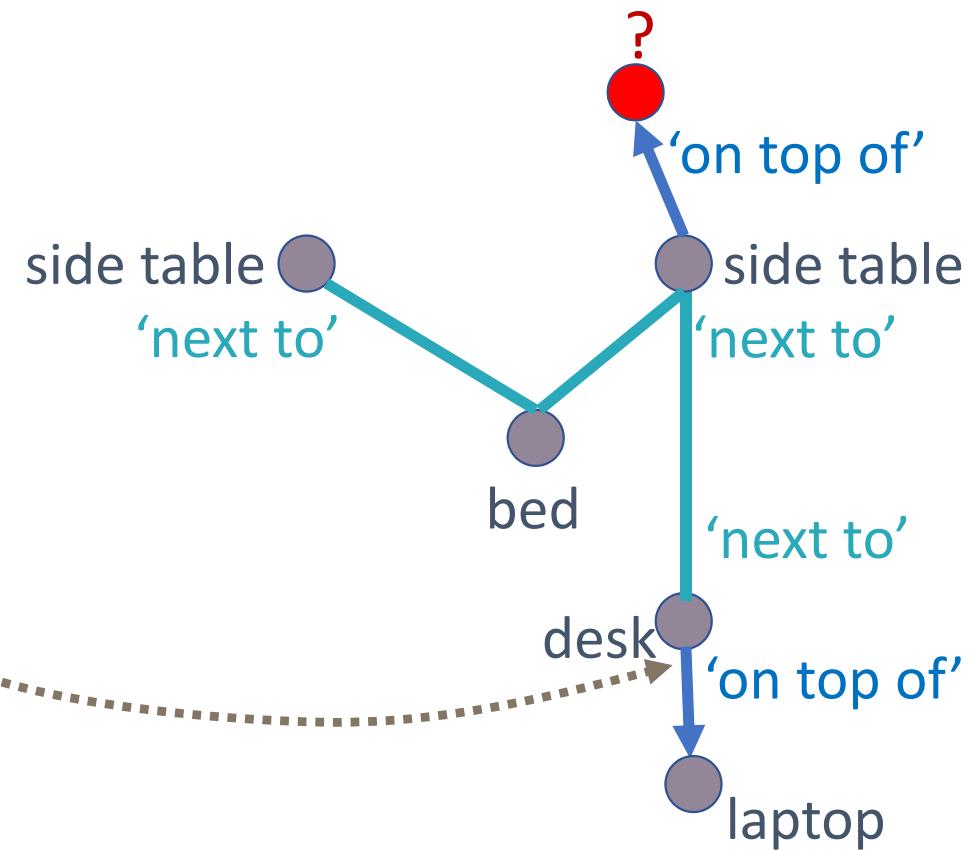
# More nodes / edges



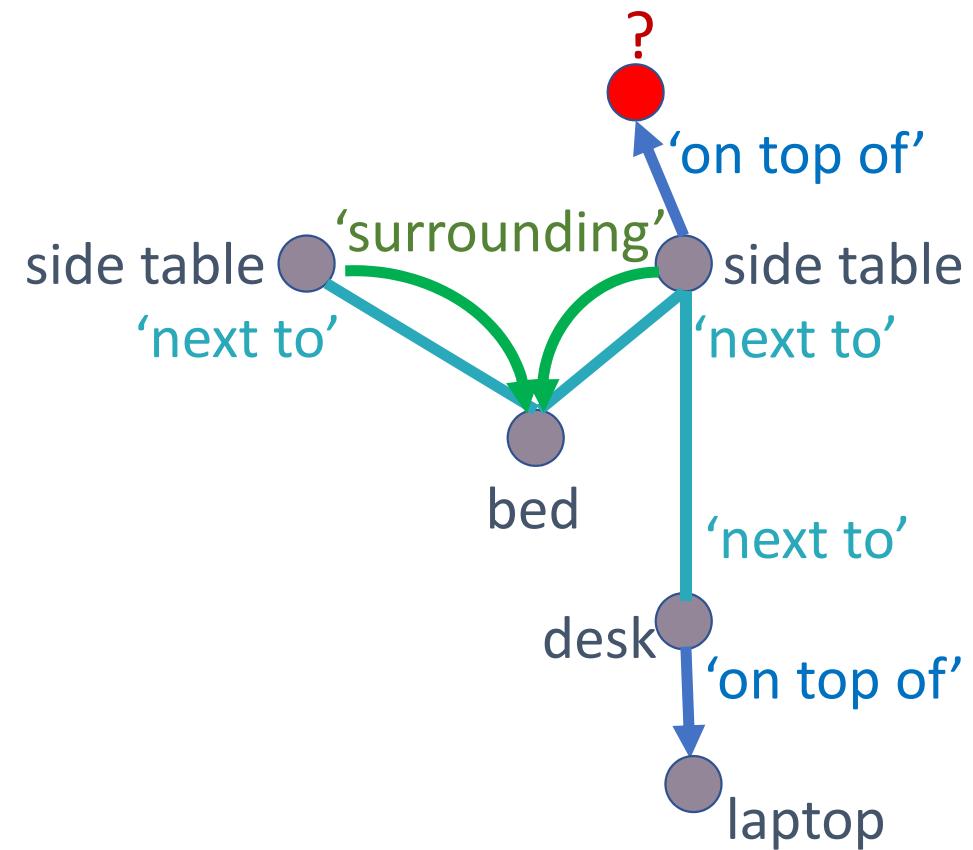
## More nodes / edges



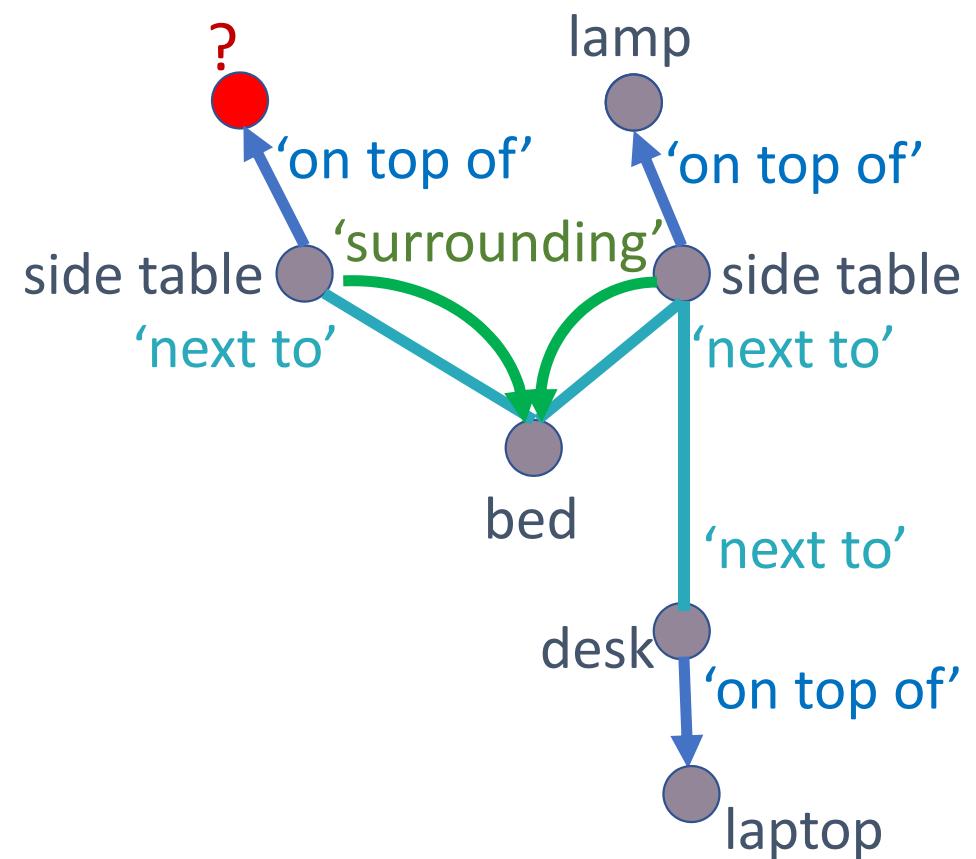
# Many more nodes / edges



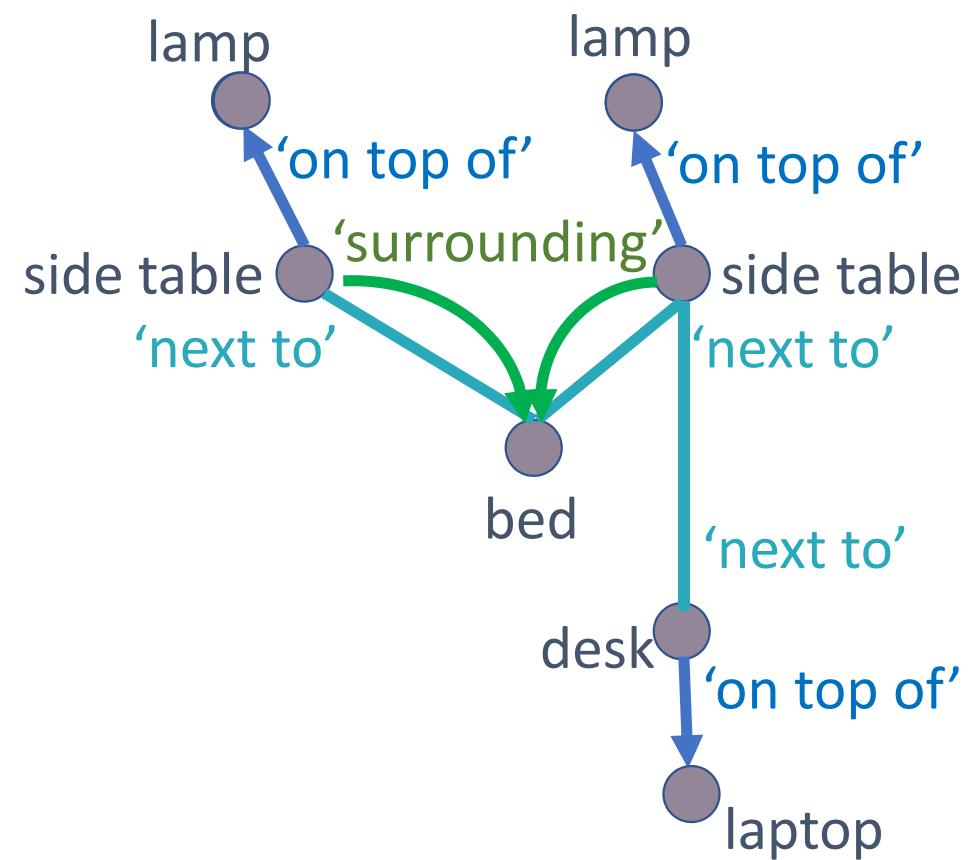
# Multiple edges for each pair of objects



# More queries



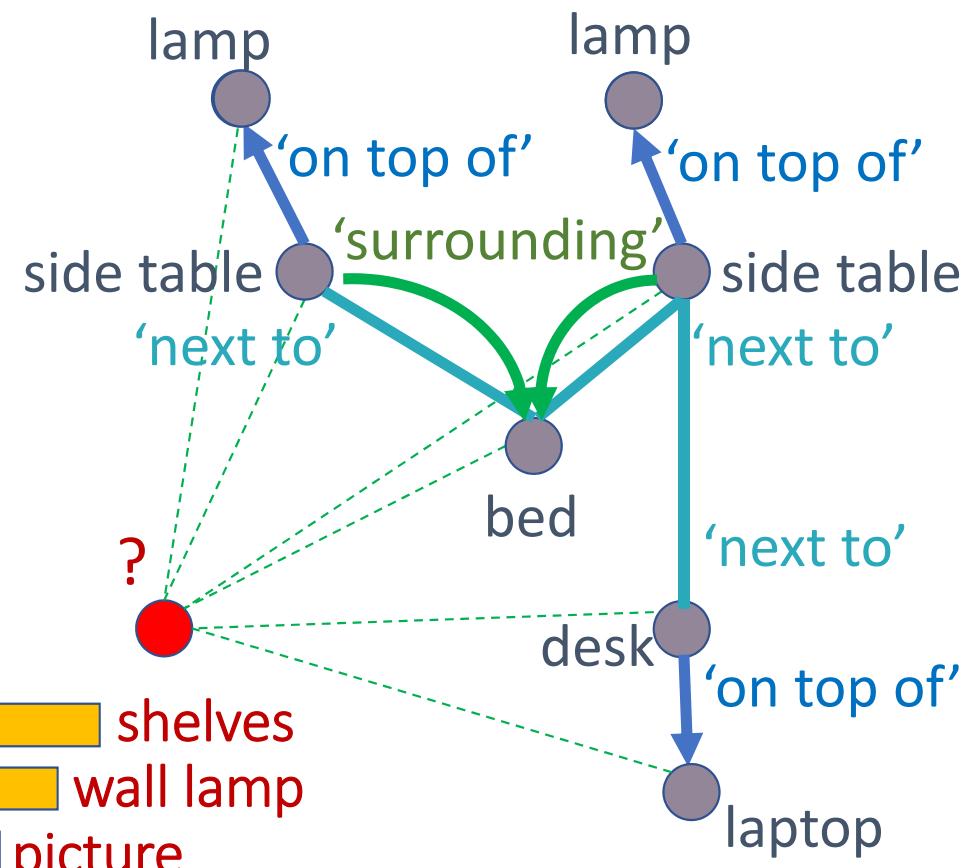
# More queries



# Denser graphs yielded better predictions

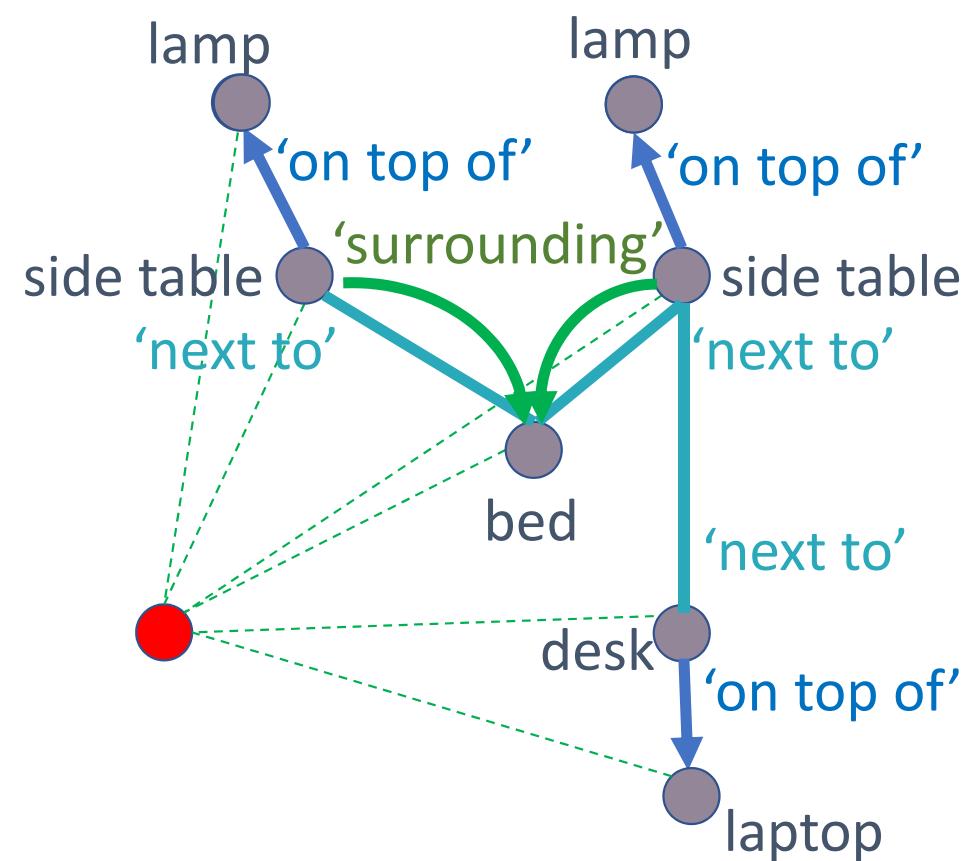


shelves  
wall lamp  
picture



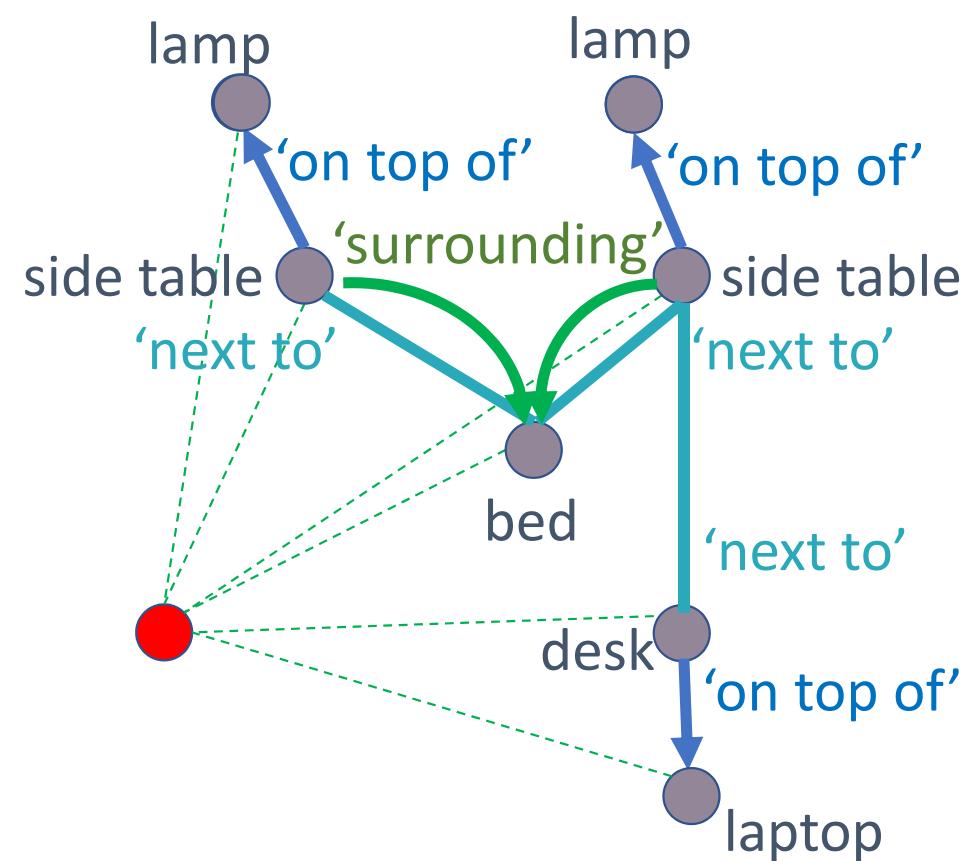
## Edge types

- “On top of”
- “Under”



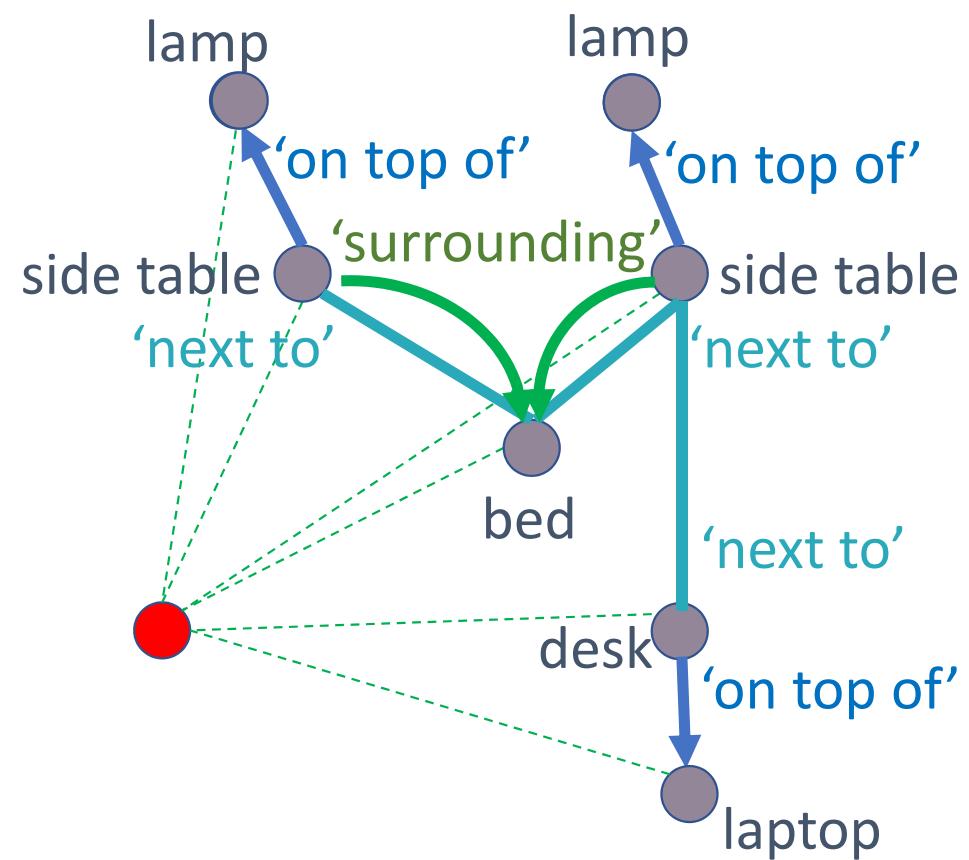
## Edge types

- “On top of”
- “Under”
- “Surrounded-by”
- “Surrounding”



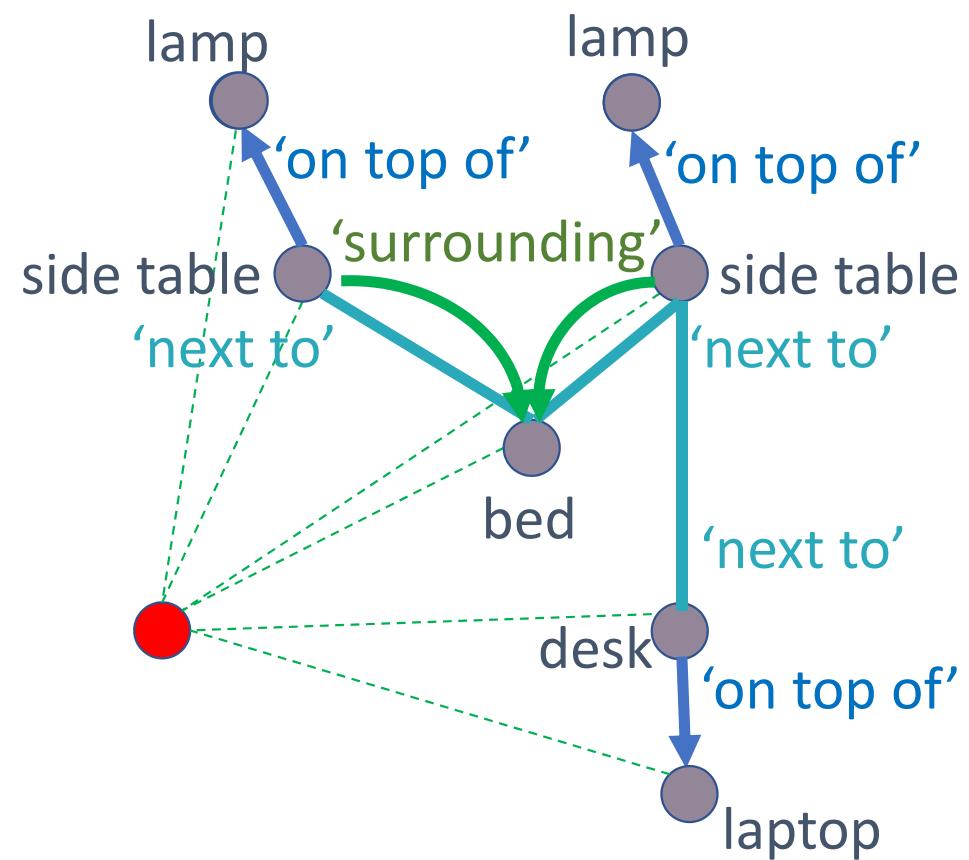
## Edge types

- “On top of”
- “Under”
- “Surrounded-by”
- “Surrounding”
- “Next-to”

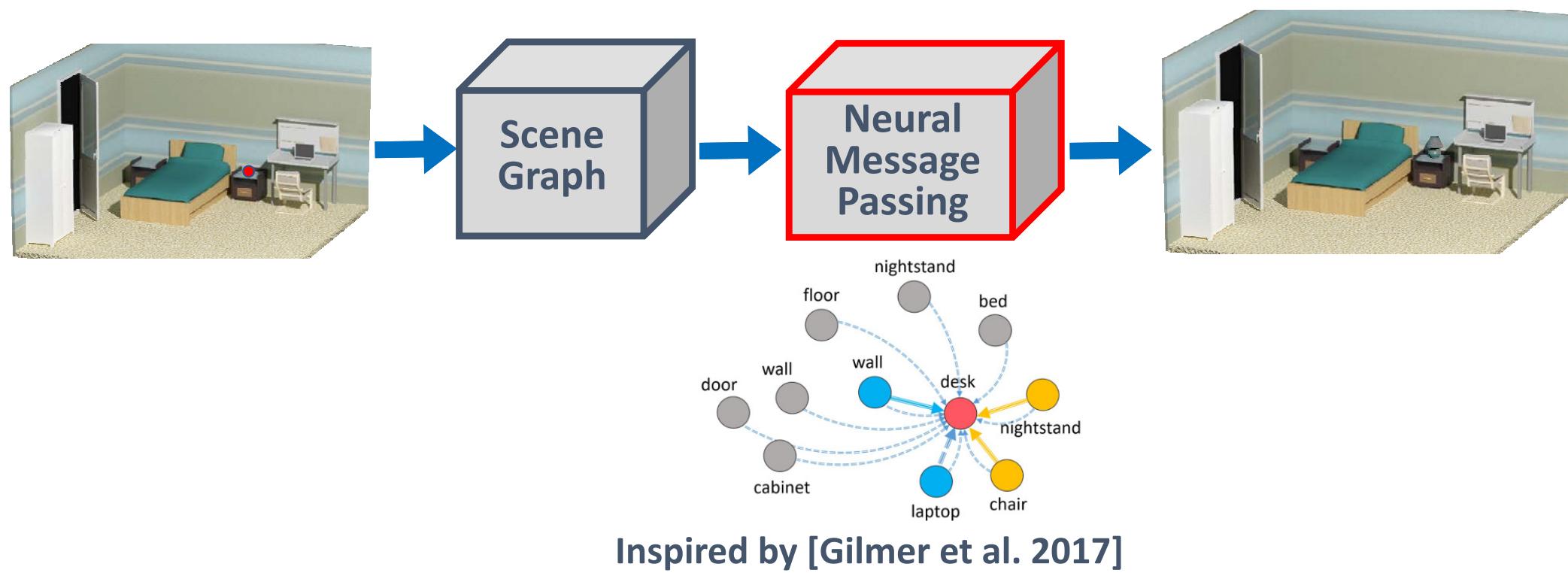


## Edge types

- “On top of”
- “Under”
- “Surrounded-by”
- “Surrounding”
- “Next-to”
- “Co-occurring”



# Pipeline



# Node representation



lamp

# Node representation



lamp

$h_{lamp}$



## Raw node features



lamp



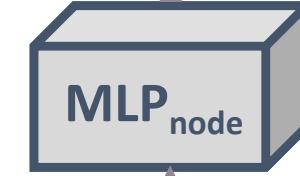
$\{(x,y,z), \text{size}, \text{object category}, \text{shape features}\}$

# Initial node representation



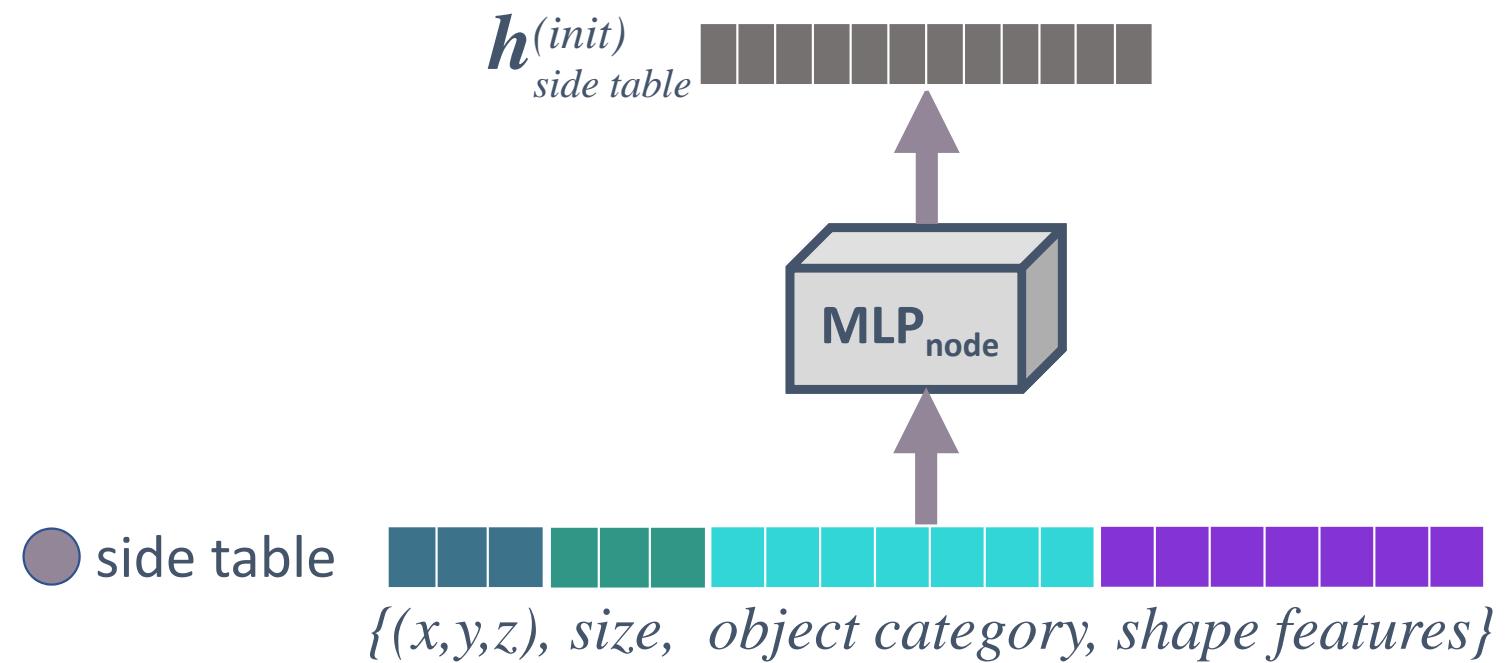
lamp

$h_{lamp}^{(init)}$

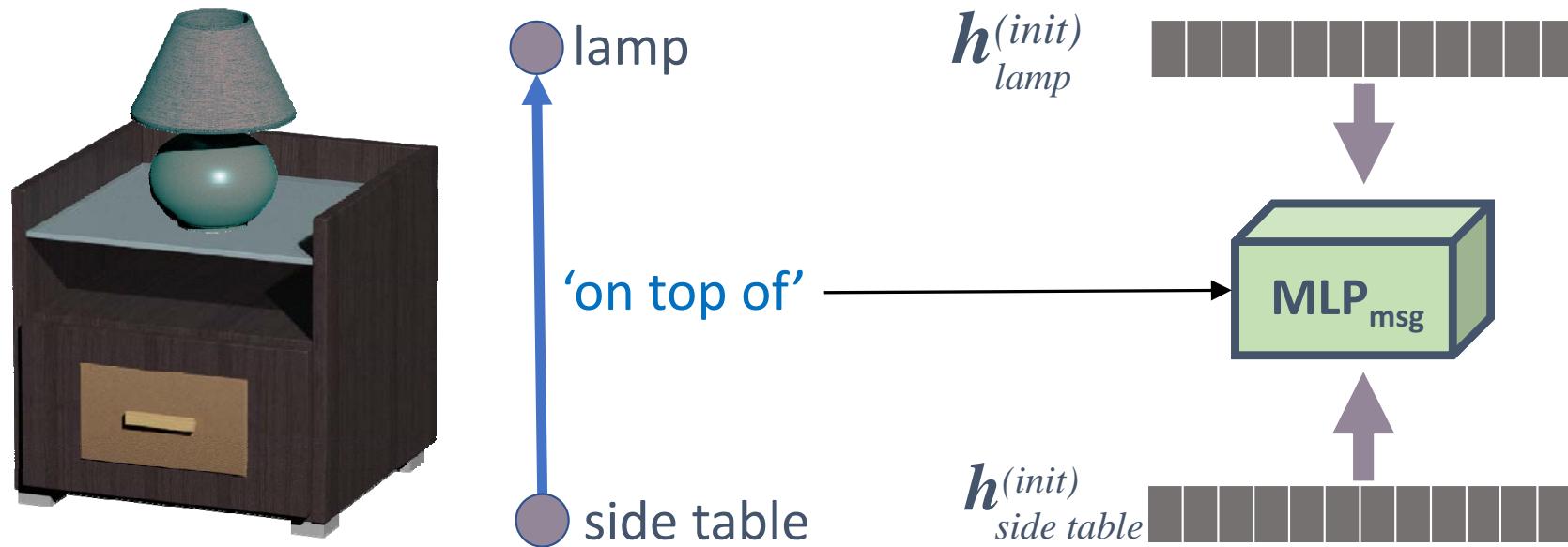


$\{(x,y,z), \text{size}, \text{object category}, \text{shape features}\}$

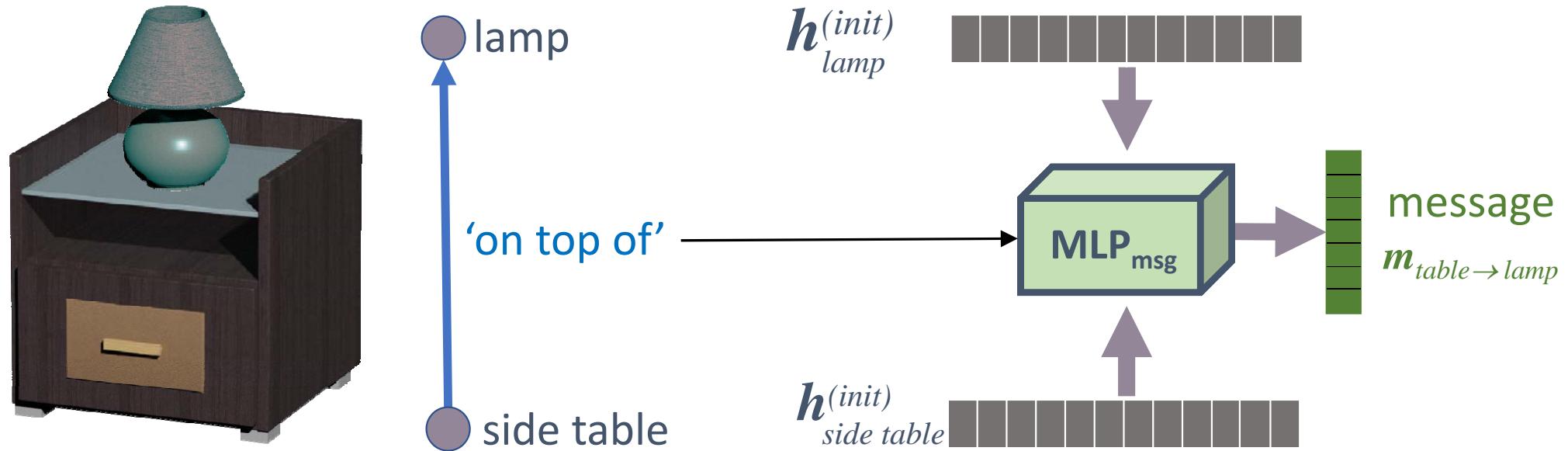
# Initial node representation



# Messages between objects



# Messages between objects



# Edge weights



lamp

'on top of'

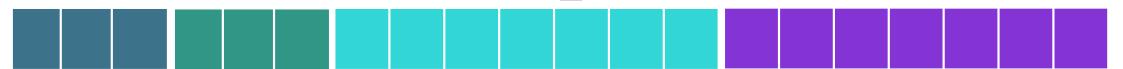
side table

$\{(x,y,z), \text{size}, \text{object category}, \text{shape features}\}$



weight  
 $a_{\text{table}, \text{lamp}}$

$\{(x,y,z), \text{size}, \text{object category}, \text{shape features}\}$



# Node update

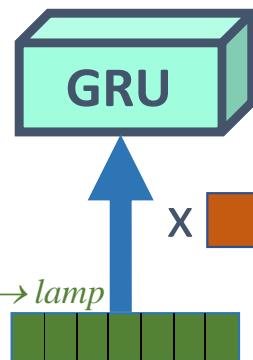


lamp

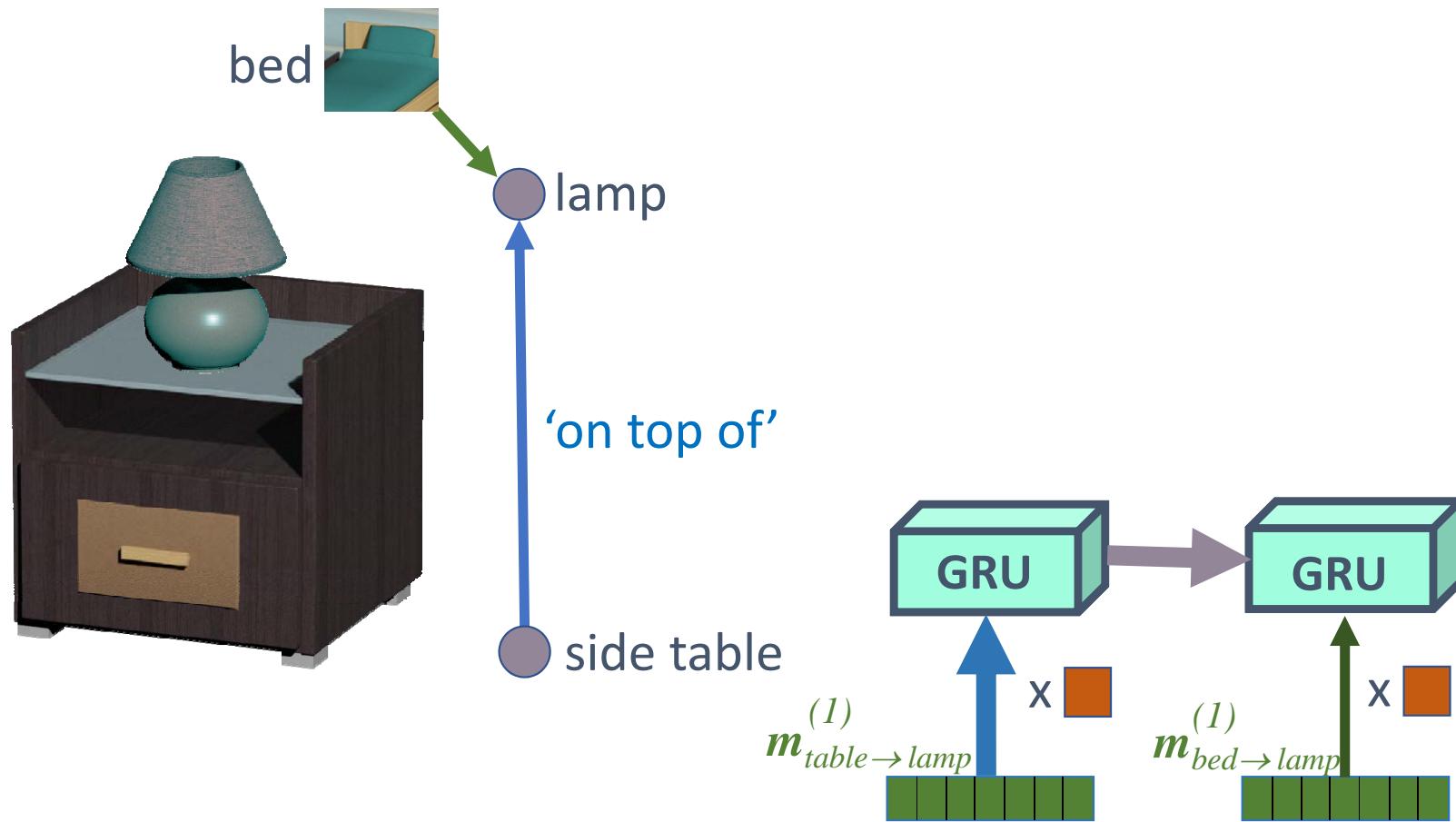
'on top of'

side table

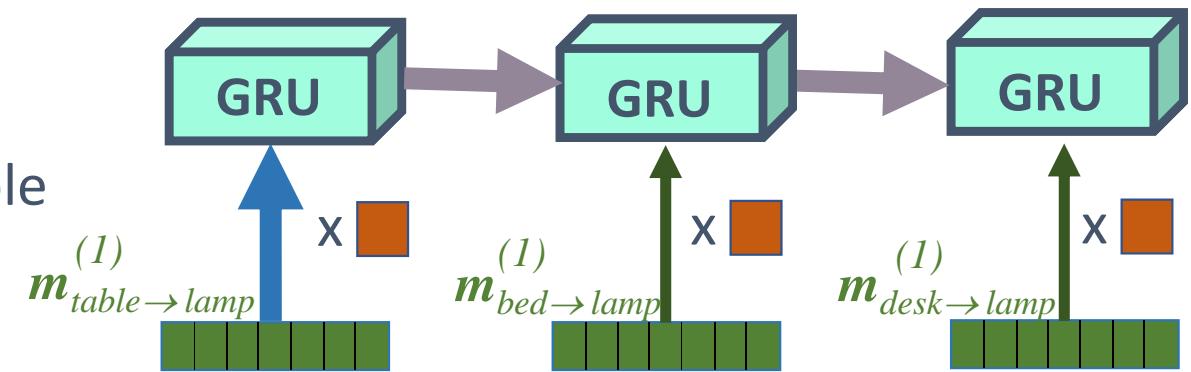
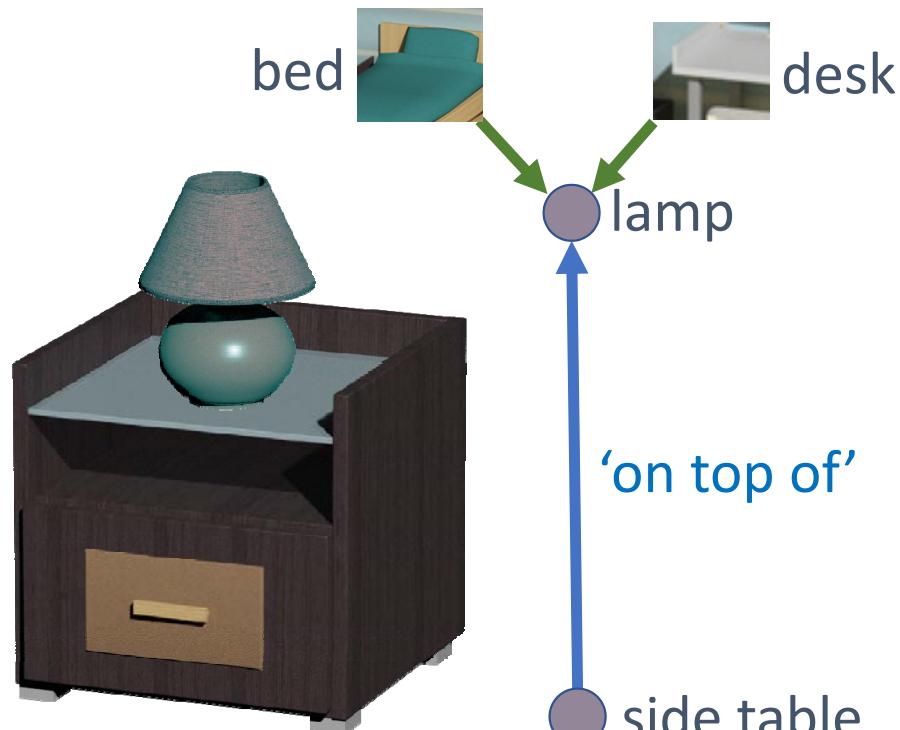
$m_{table \rightarrow lamp}^{(1)}$



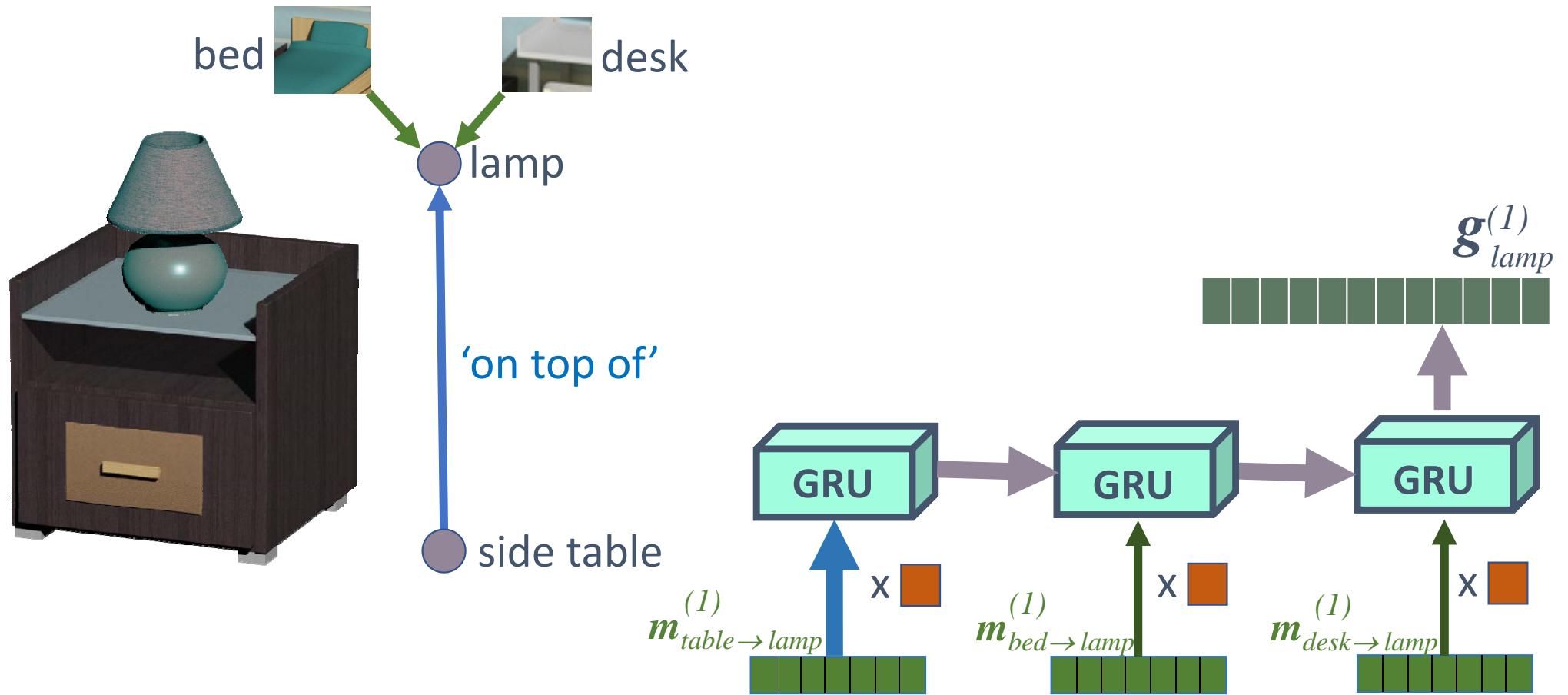
# Node update



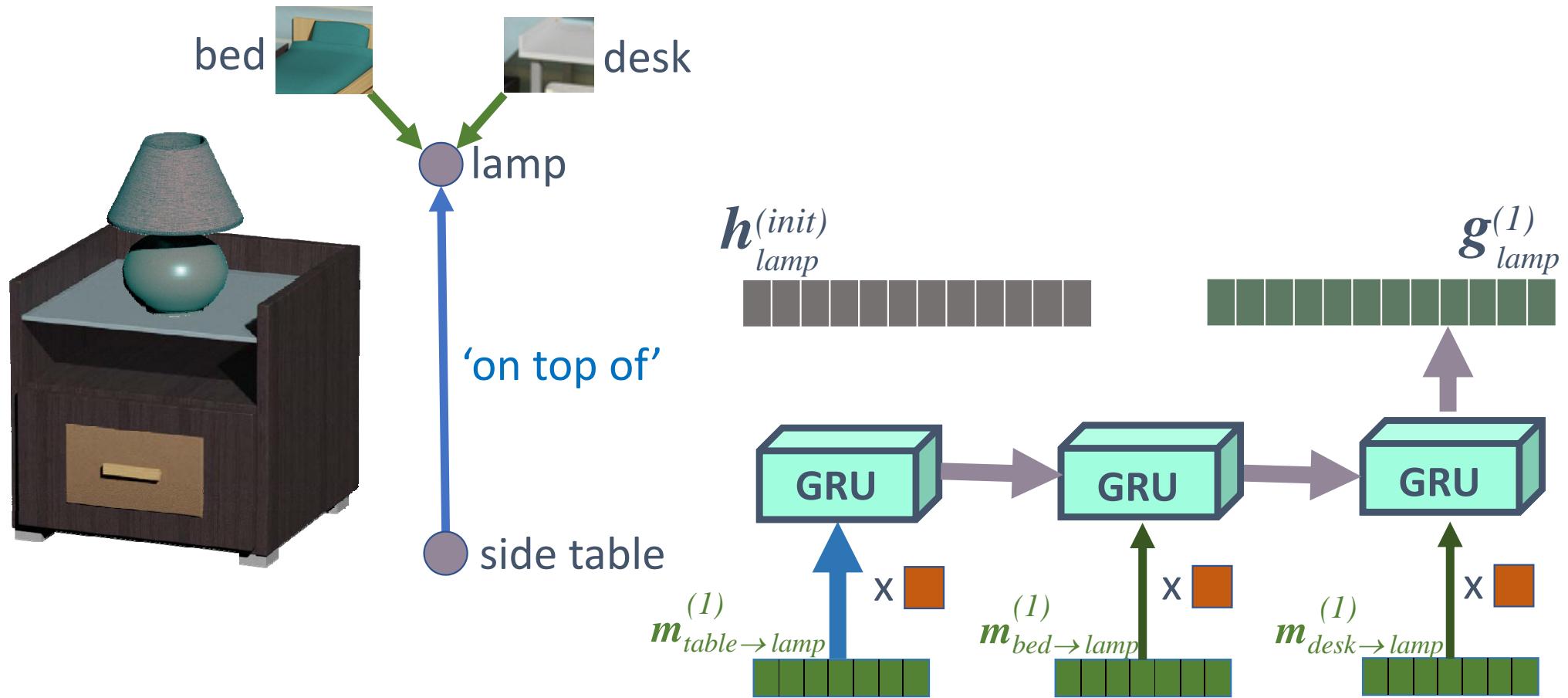
# Node update



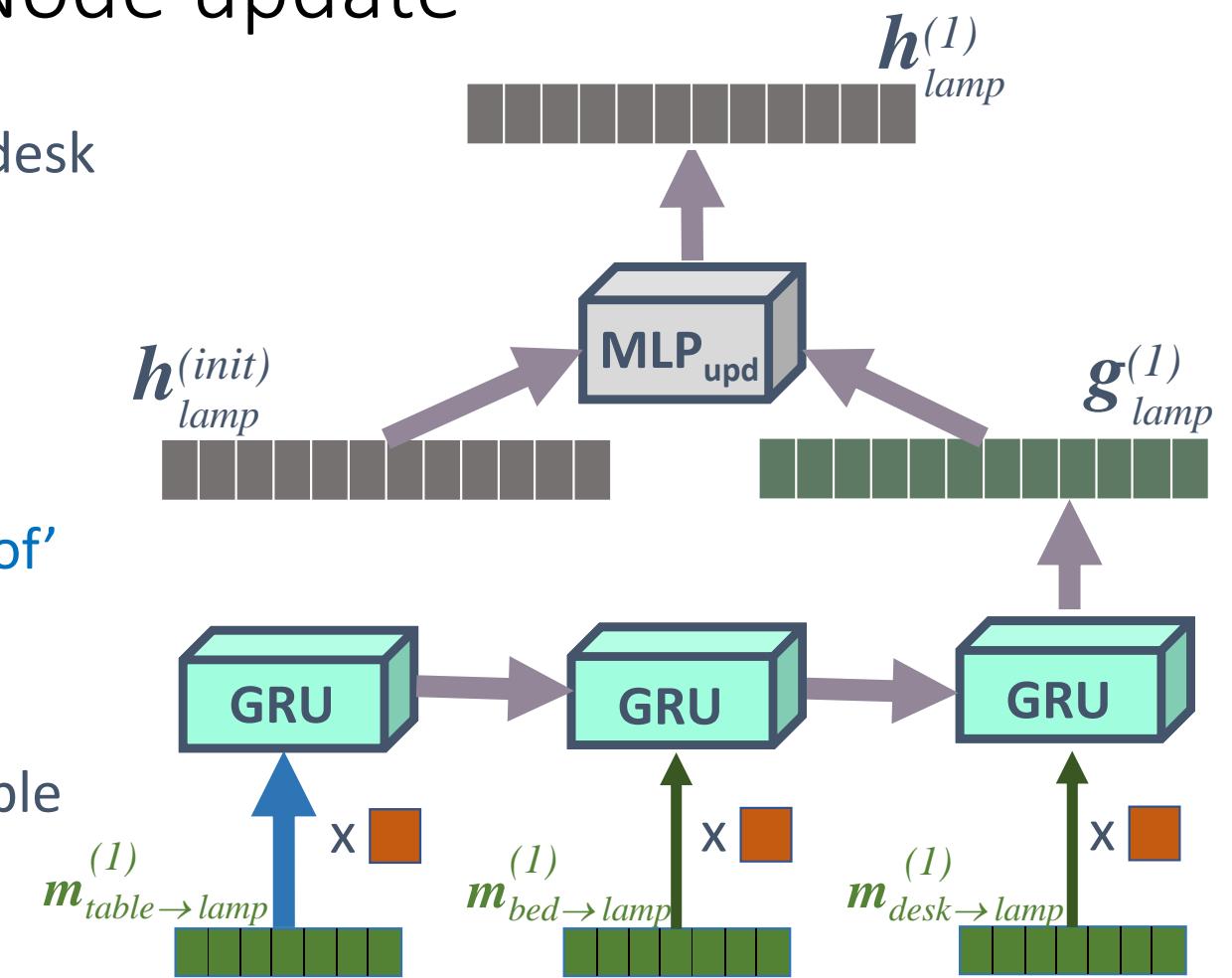
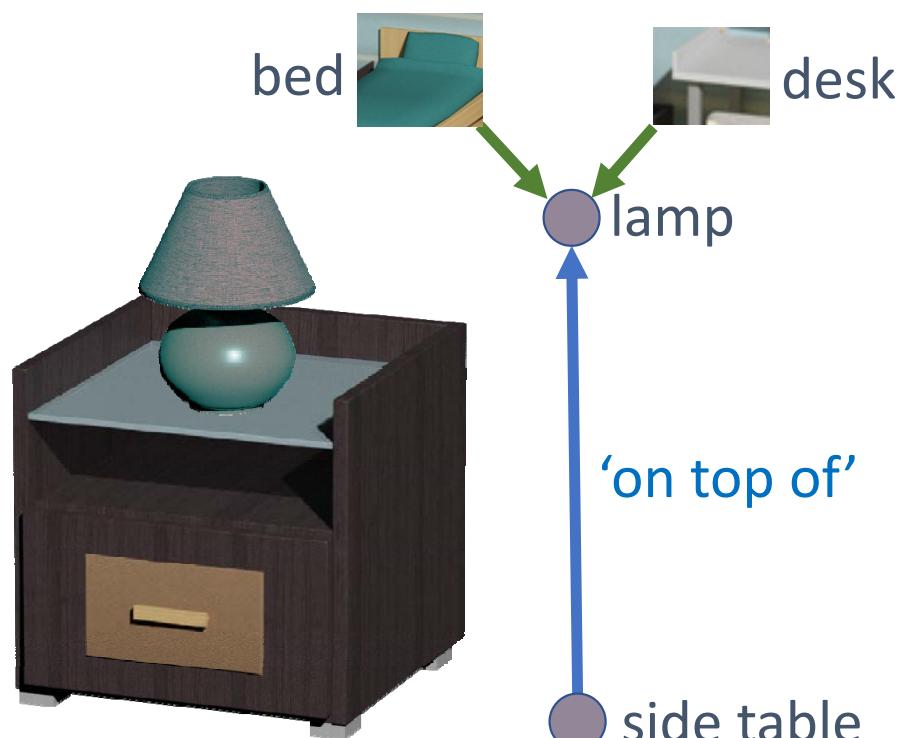
# Node update



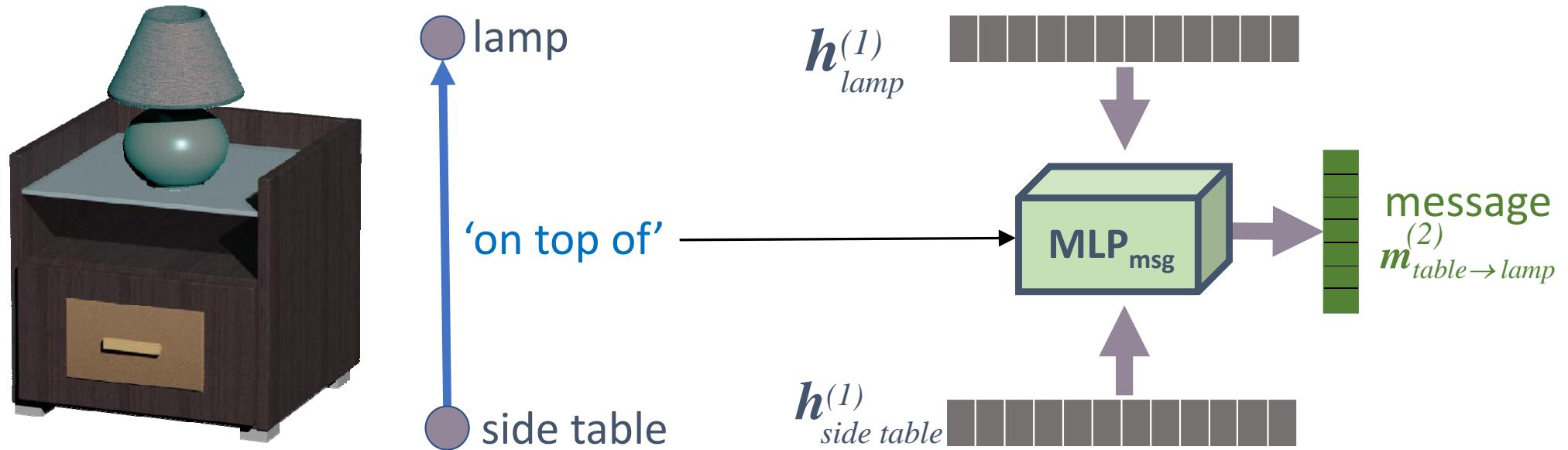
# Node update



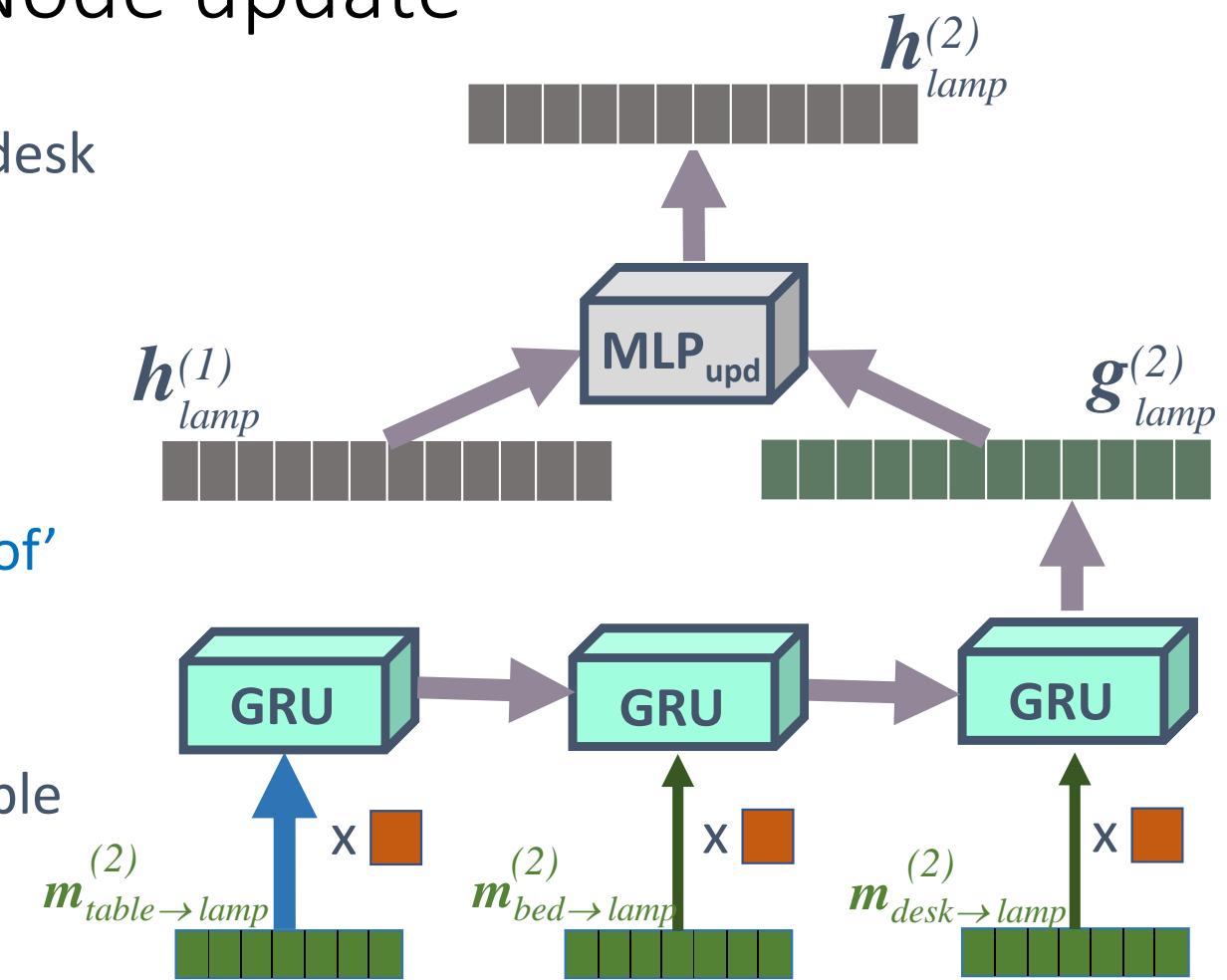
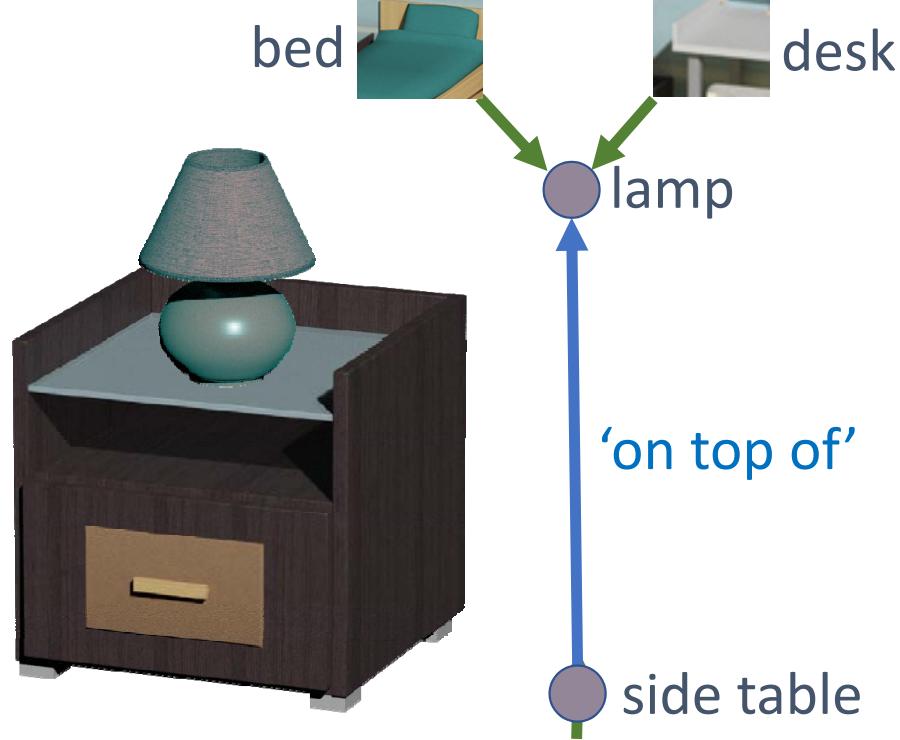
# Node update



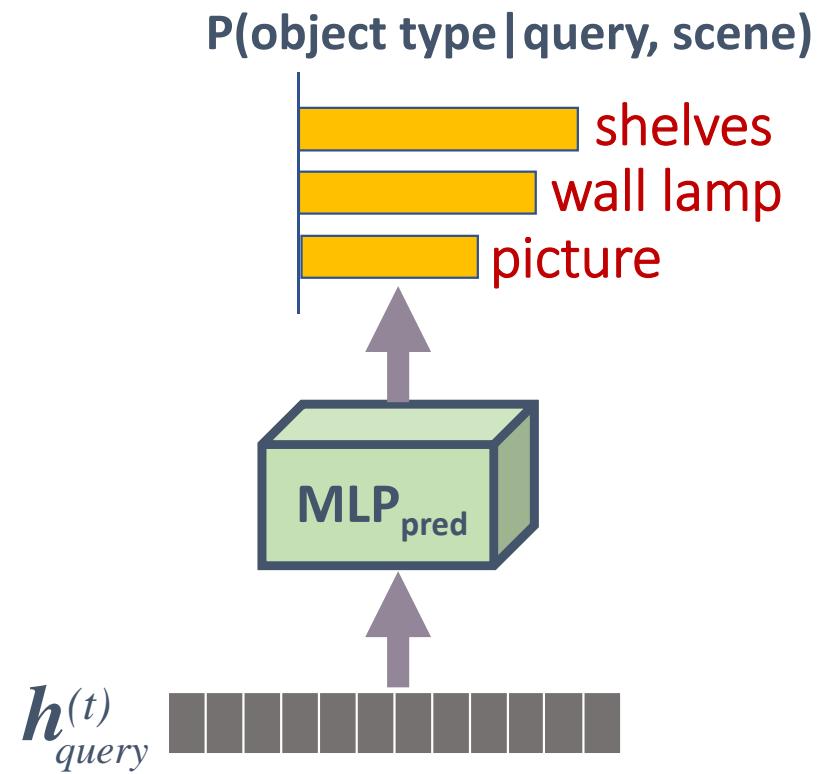
# Messages between objects



# Node update



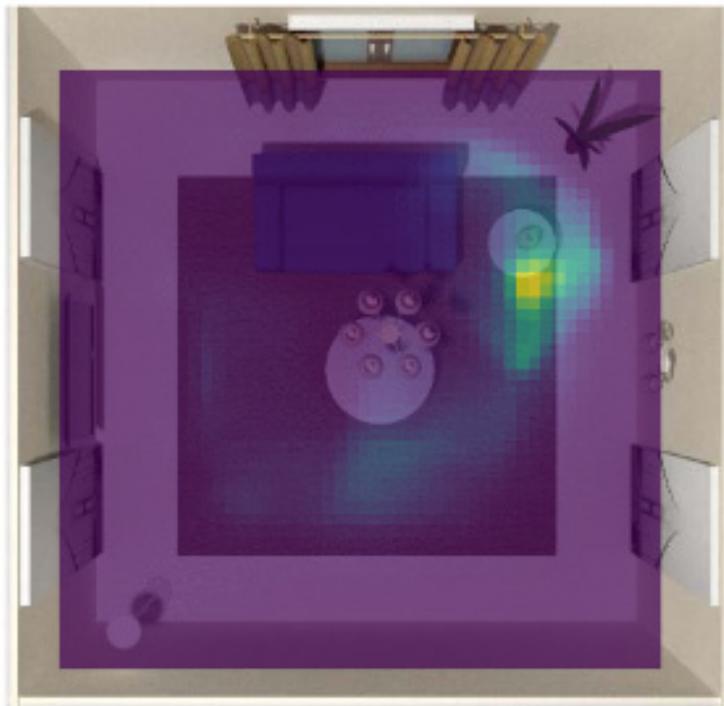
# Decoding query node representations



# Iterative synthesis



Incomplete Scene



Evaluate output probability in a dense grid of locations e.g., chair



# Iterative synthesis



Incomplete Scene

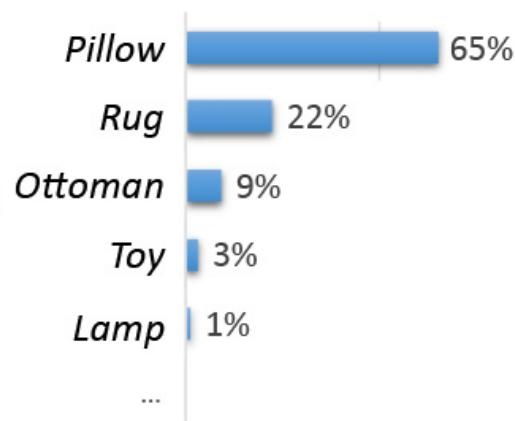


Augmented Scene

# Application: context-driven object recognition



Ottoman



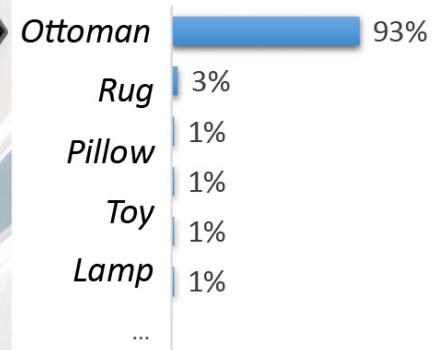
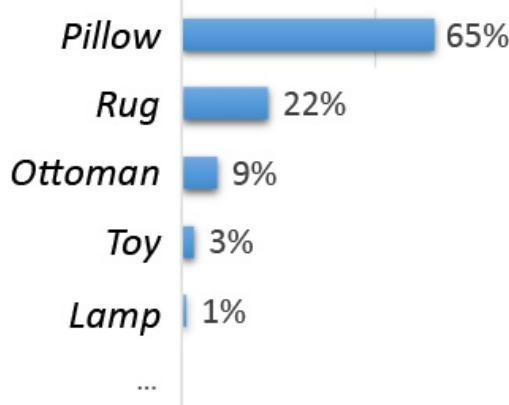
$P(\text{object type} \mid \text{shape})$

(multi-view CNN, Su et al. 2015])

# Application: context-driven object recognition



Ottoman



$P(\text{object type} \mid \text{shape})$   
(multi-view CNN, Su et al. 2015])

$P(\text{object type} \mid \text{query, scene}) \times$   
 $P(\text{object type} \mid \text{shape})$



# Training

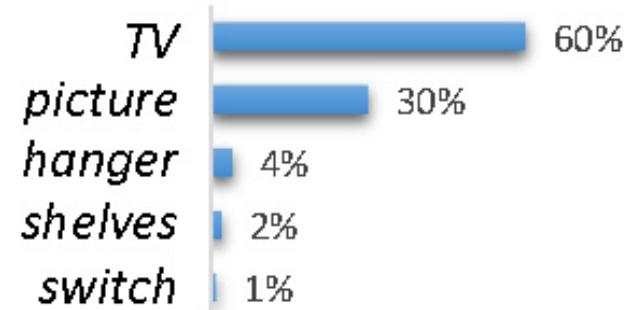
All MLP/GRU modules are learned from a **large dataset of indoor rooms** (SUNCG)

We remove random objects, then train our network to predict them

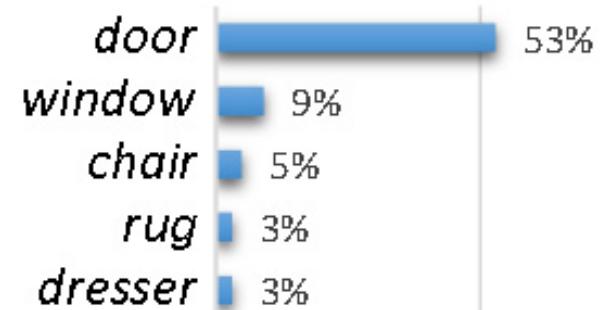


[Song et al. 16]

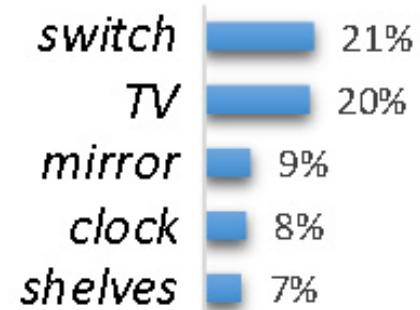
# SceneGraphNet Results



# GRAINS Results [Li et al. 19]



# “Deep Priors” Results [Wang et al. 19]



# Numerical Evaluation

How well do we predict objects removed from a test scene?

Method	Correct object type
“GRAINS” *	44.2%

\* GRAINS: Generative Recursive Autoencoders for INdoor Scenes, TOG 2019

# Numerical Evaluation

How well do we predict objects removed from a test scene?

Method	Correct object type
“GRAINS” *	44.2%
Wang et al.**	50.9%

\* GRAINS: Generative Recursive Autoencoders for INdoor Scenes, TOG 2019

\*\* Deep convolutional priors for indoor scene synthesis, SIGGRAPH 2018

# Numerical Evaluation

How well do we predict objects removed from a test scene?

Method	Correct object type
“GRAINS” *	44.2%
Wang et al.**	50.9%
<b>SceneGraphNet</b>	<b>67.3%</b>

\* GRAINS: Generative Recursive Autoencoders for INdoor Scenes, TOG 2019

\*\* Deep convolutional priors for indoor scene synthesis, SIGGRAPH 2018

# Numerical Evaluation

How well do we predict objects removed from a test scene?

Method	Correct object within top-5 predictions
“GRAINS” *	73.6%
Wang et al.**	79.3%
<b>SceneGraphNet</b>	<b>89.5%</b>

\* GRAINS: Generative Recursive Autoencoders for INdoor Scenes, TOG 2019

\*\* Deep convolutional priors for indoor scene synthesis, SIGGRAPH 2018

# Numerical Evaluation

Object classification accuracy averaged over all objects and test scenes

Method	Classification accuracy
MVCNN *	59.2%
<b>MVCNN + SceneGraphNet</b>	<b>72.2%</b>

\* Multi-view Convolutional Neural Networks for 3D Shape Recognition, ICCV 2015

# Summary

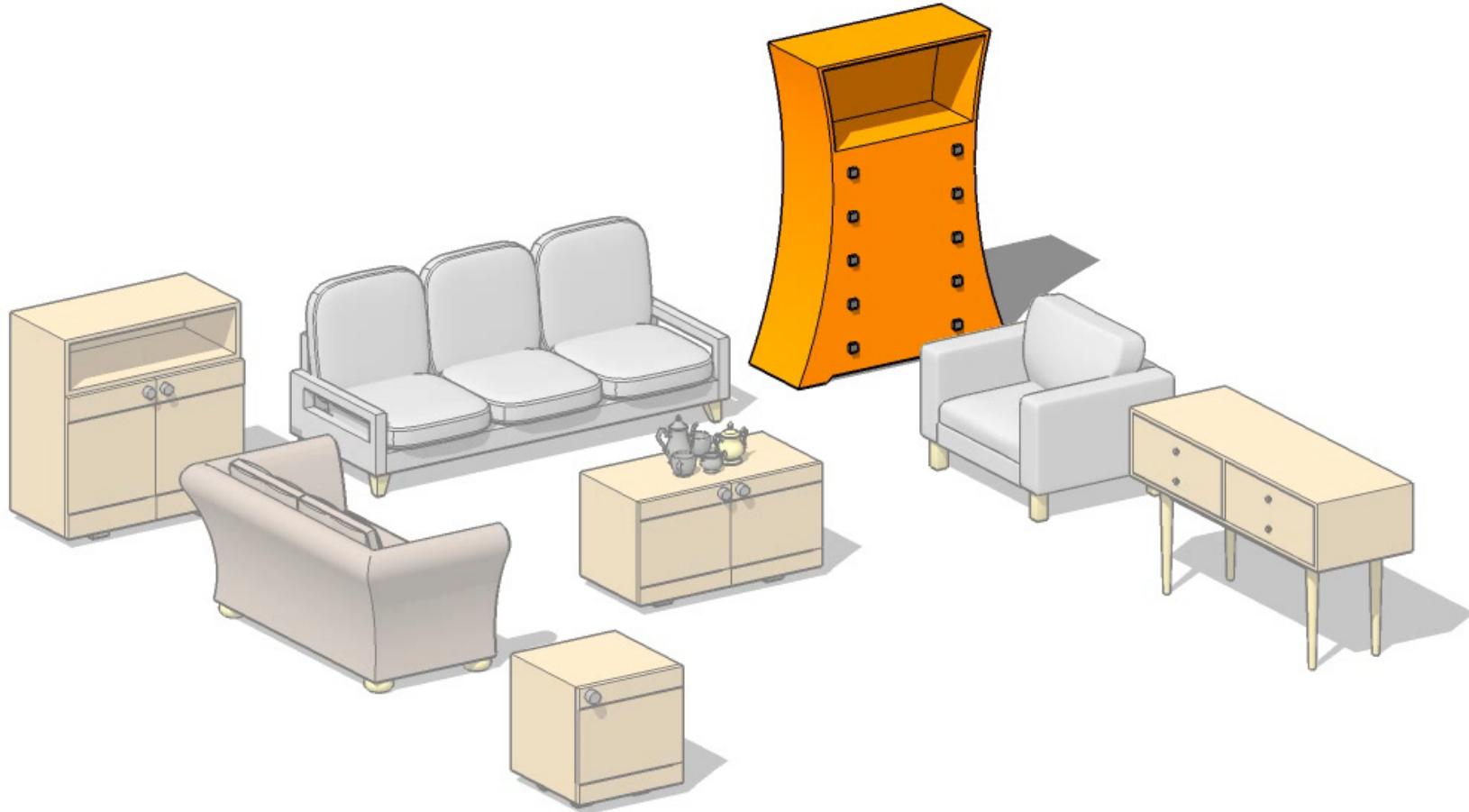
- Network for context-based scene completion & iterative scene synthesis
- Key component: neural message passing on a graph whose nodes represent objects and edges represent spatial and structural relationships
- Demonstrated also for context-based object recognition

# Limitations

- Graphics-oriented: input scenes are synthetic & 3D models are labeled
- Our graph edges are derived from simple heuristics
- Limited to predicting only object categories and sizes – it does not generate geometry or style!

# Shape style transfer





**Functionality Preserving Shape Style Transfer, TOG 2016**

# Challenge: separate style from function



exemplar (style)



target (function)



output

Observation: similar style elements



# Observation: similar style elements

similarly shaped elements



# Observation: similar style elements

similar dominant curves

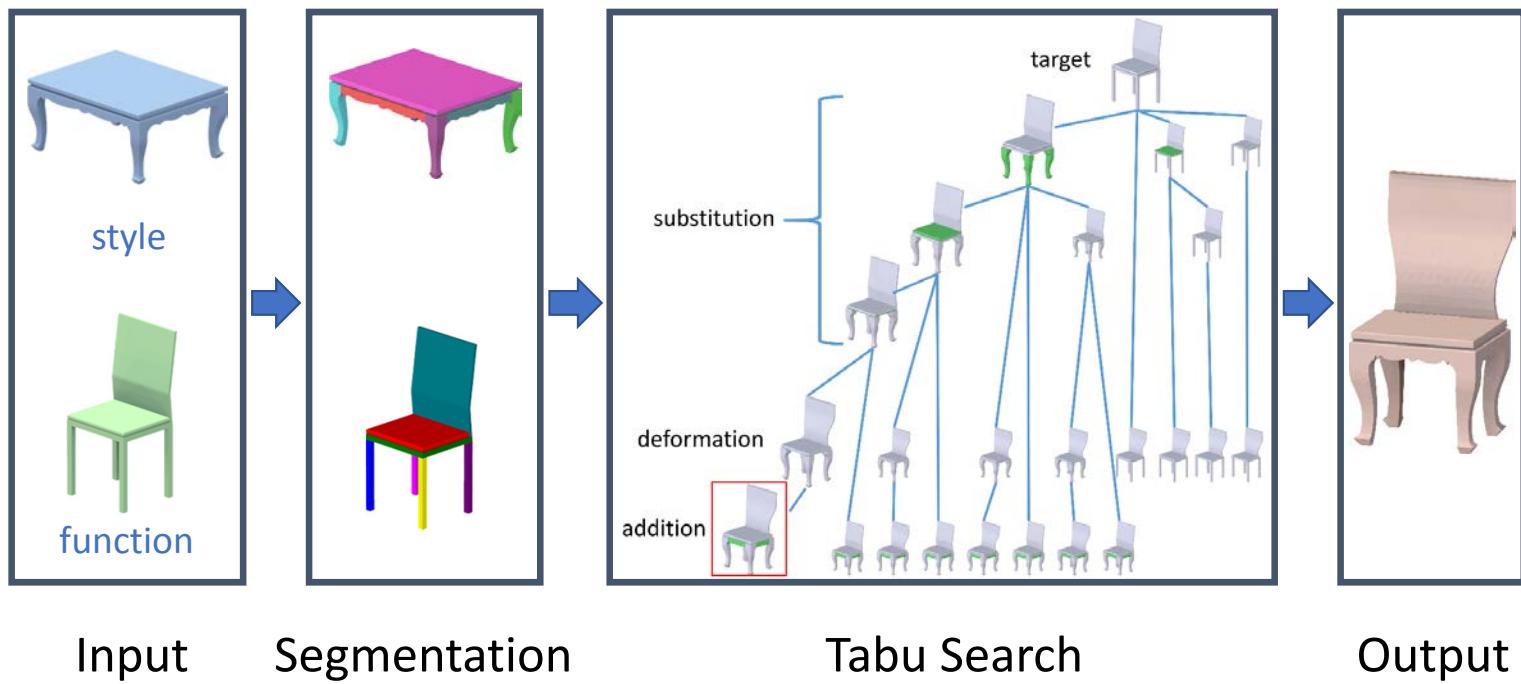


# Key idea

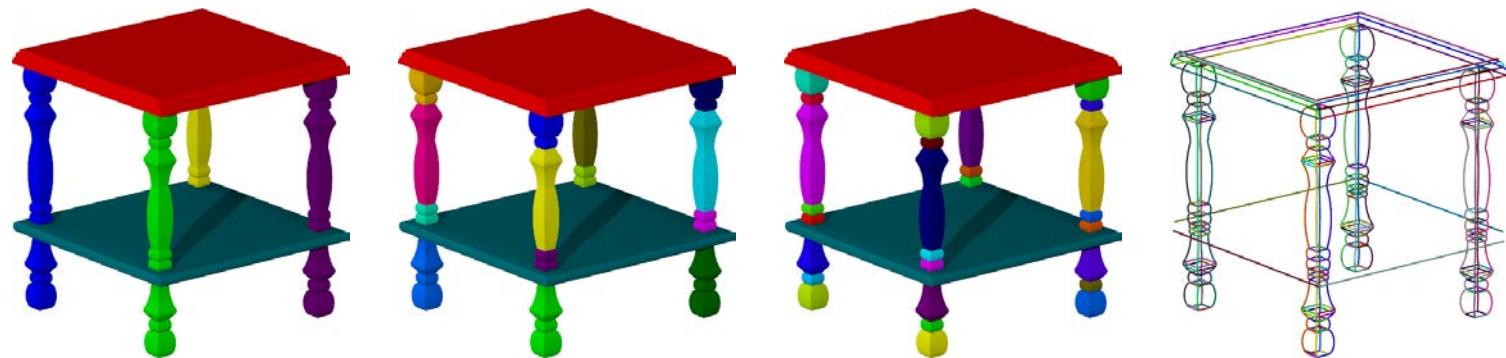
Search for **compatible** element operations that:

- **Maintain gross form & part structure of target shape**  
=> **functionality preservation**
- **Increase style similarity to exemplar**  
=> **style transfer**

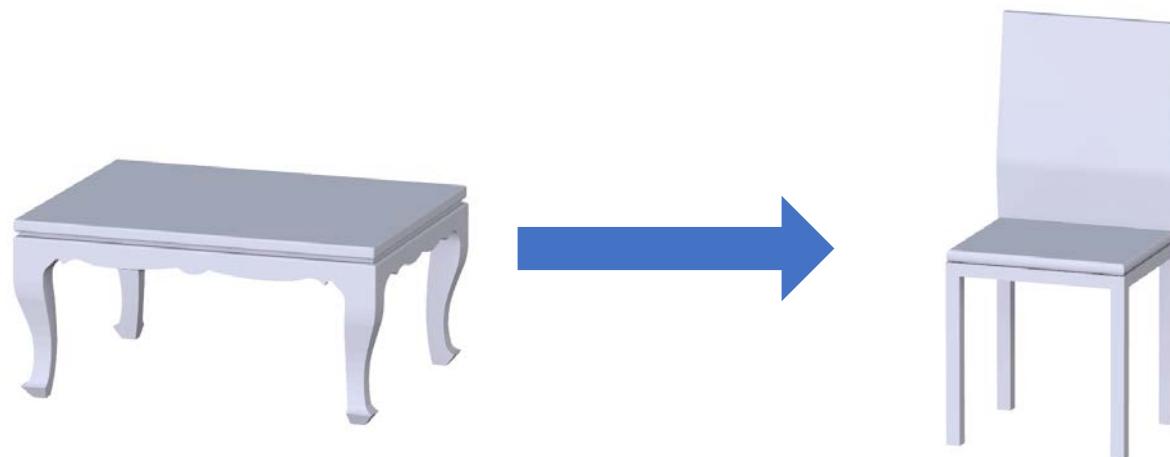
# Framework overview



# Hierarchical segmentation and extracted curves

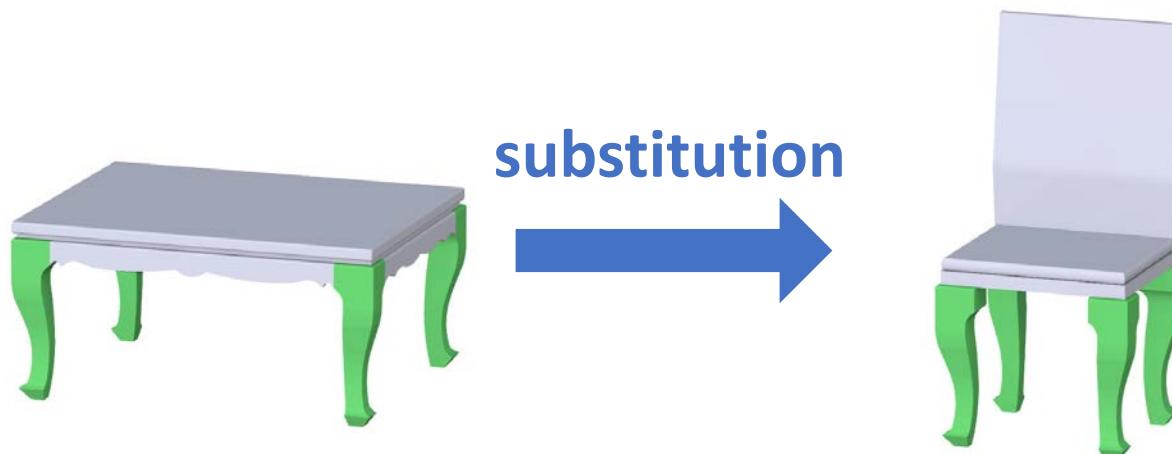


# Element-level operations



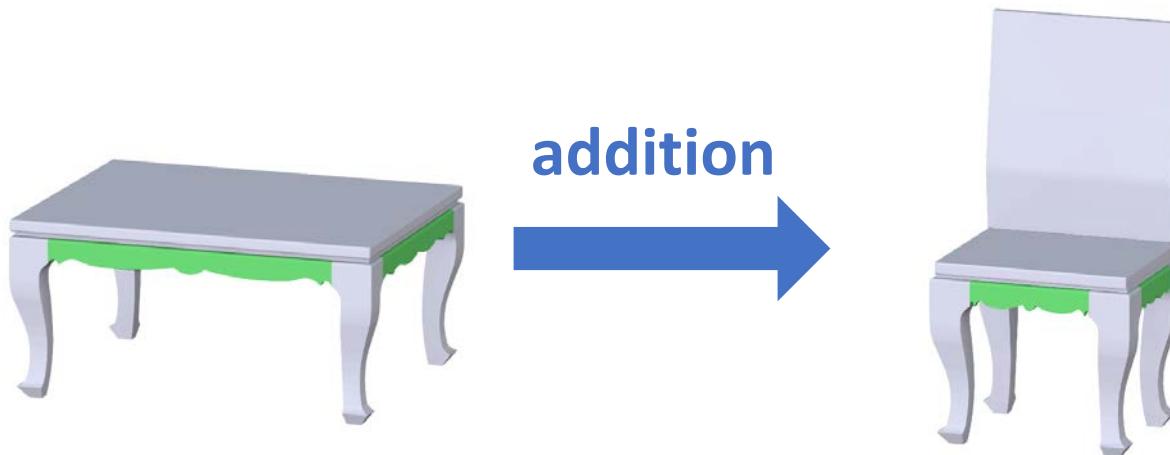
# Element-level operations

- part substitution



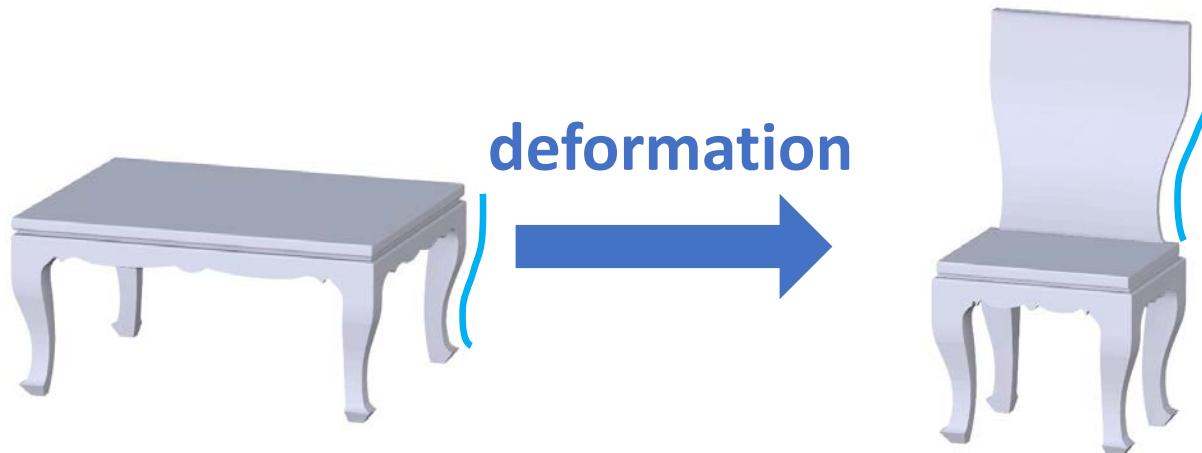
# Element-level operations

- part substitution
- part addition / removal



# Element-level operations

- part substitution
- part addition / removal
- curve-based deformation



# Unconstrained element edits



style



function



output

## Compatible element edits



style



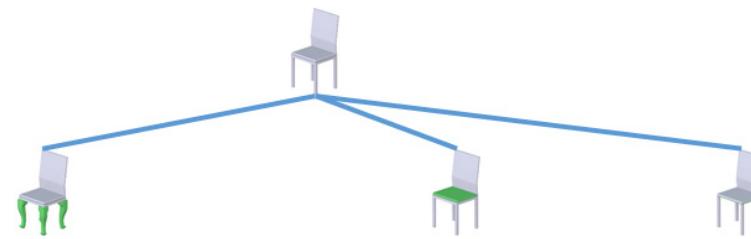
function



output

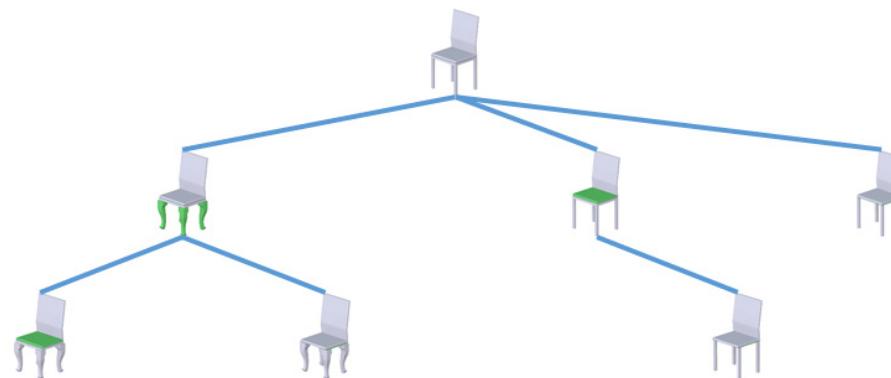
# Search for compatible element-level operations

substitution

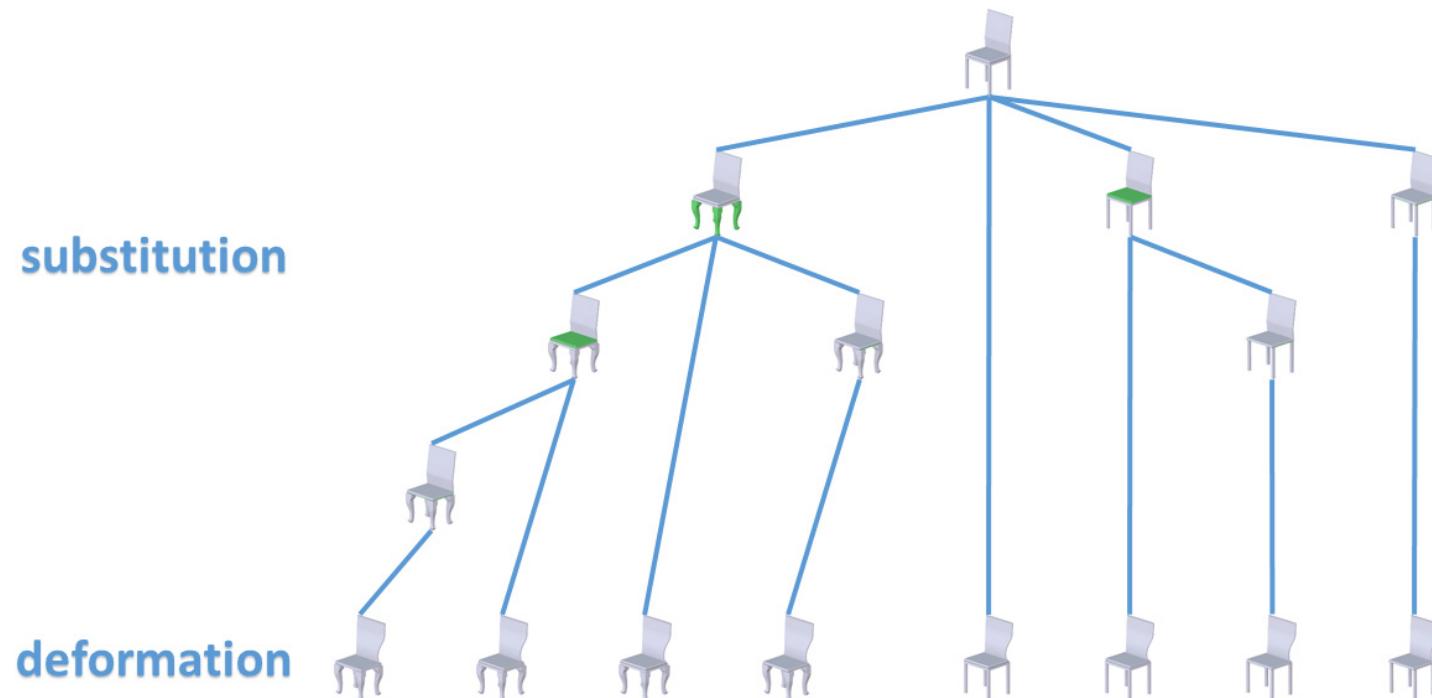


# Search for compatible element-level operations

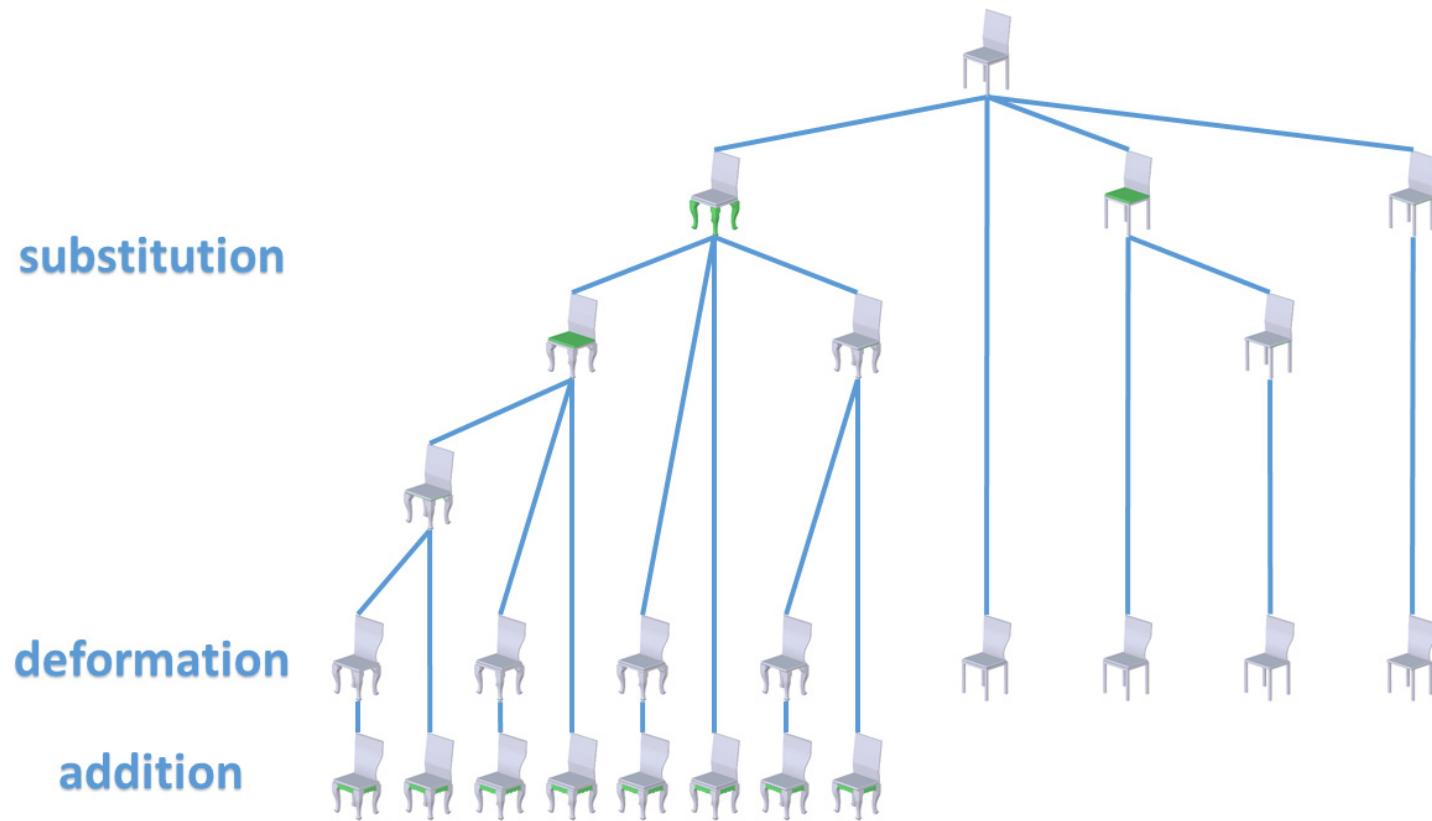
substitution



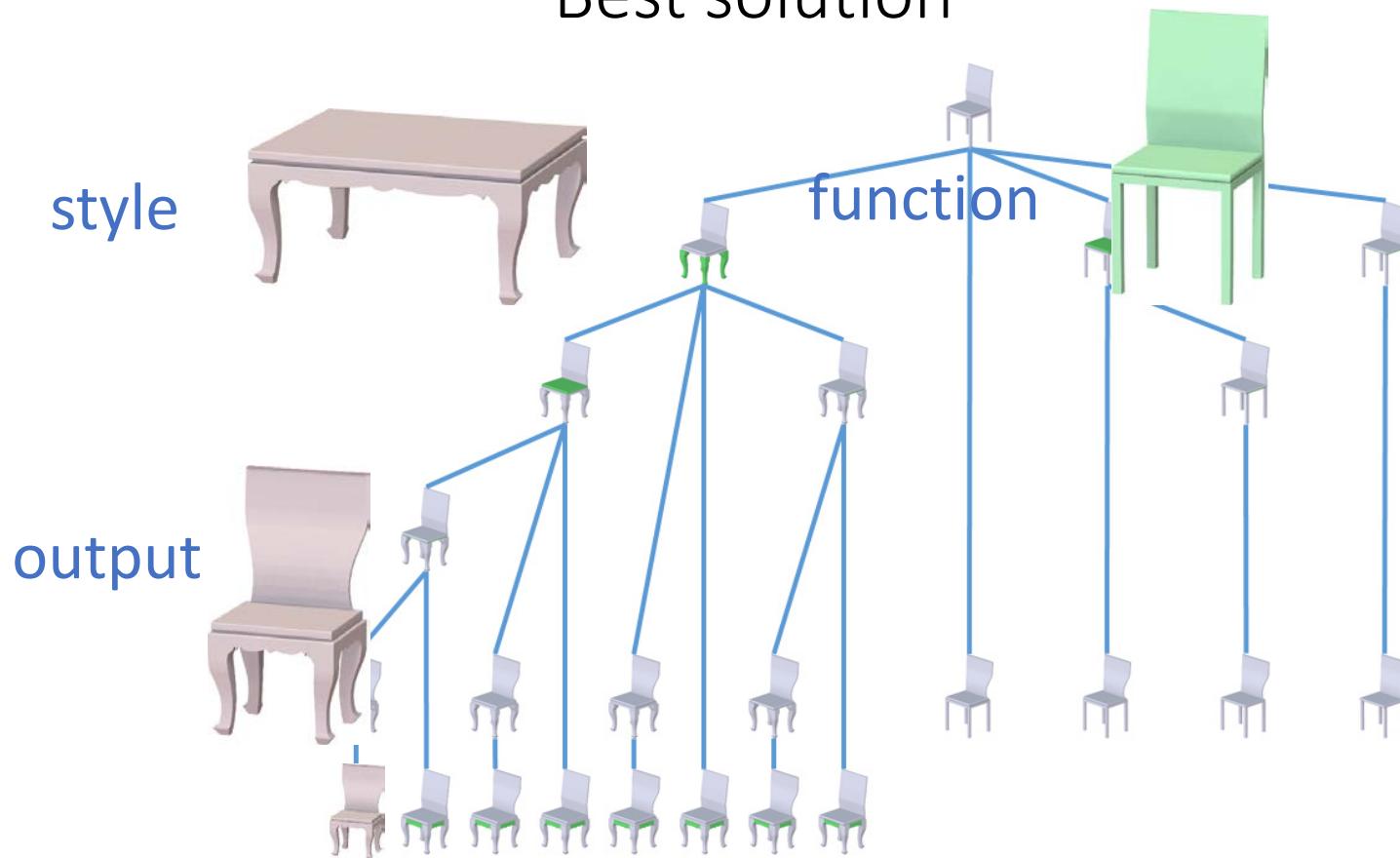
# Search for compatible element-level operations



# Search for compatible element-level operations



Best solution



Multiple solutions

style



function



outputs



# Measures

- Style similarity measure from:  
“Elements of Style: Learning Perceptual Shape Style Similarity”,  
Lun, Kalogerakis, Sheffer, TOG 2015
- Functional compatibility measure



$$D_{func}(\text{green chair}, \text{purple chair}) = 1.32$$



$$D_{func}(\text{green chair}, \text{purple chair}) = 0.57$$

# Measures

- Style similarity measure from:  
“Elements of Style: Learning Perceptual Shape Style Similarity”,  
Lun, Kalogerakis, Sheffer, TOG 2015

- Functional compatibility measure

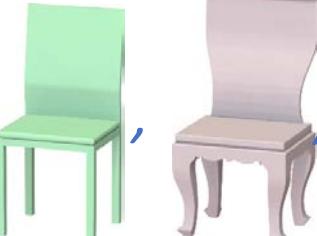
$$D_{func}(\text{green chair}, \text{purple stool}) = 1.32 \quad D_{func}(\text{green chair}, \text{purple chair}) = 0.57$$


# Measures

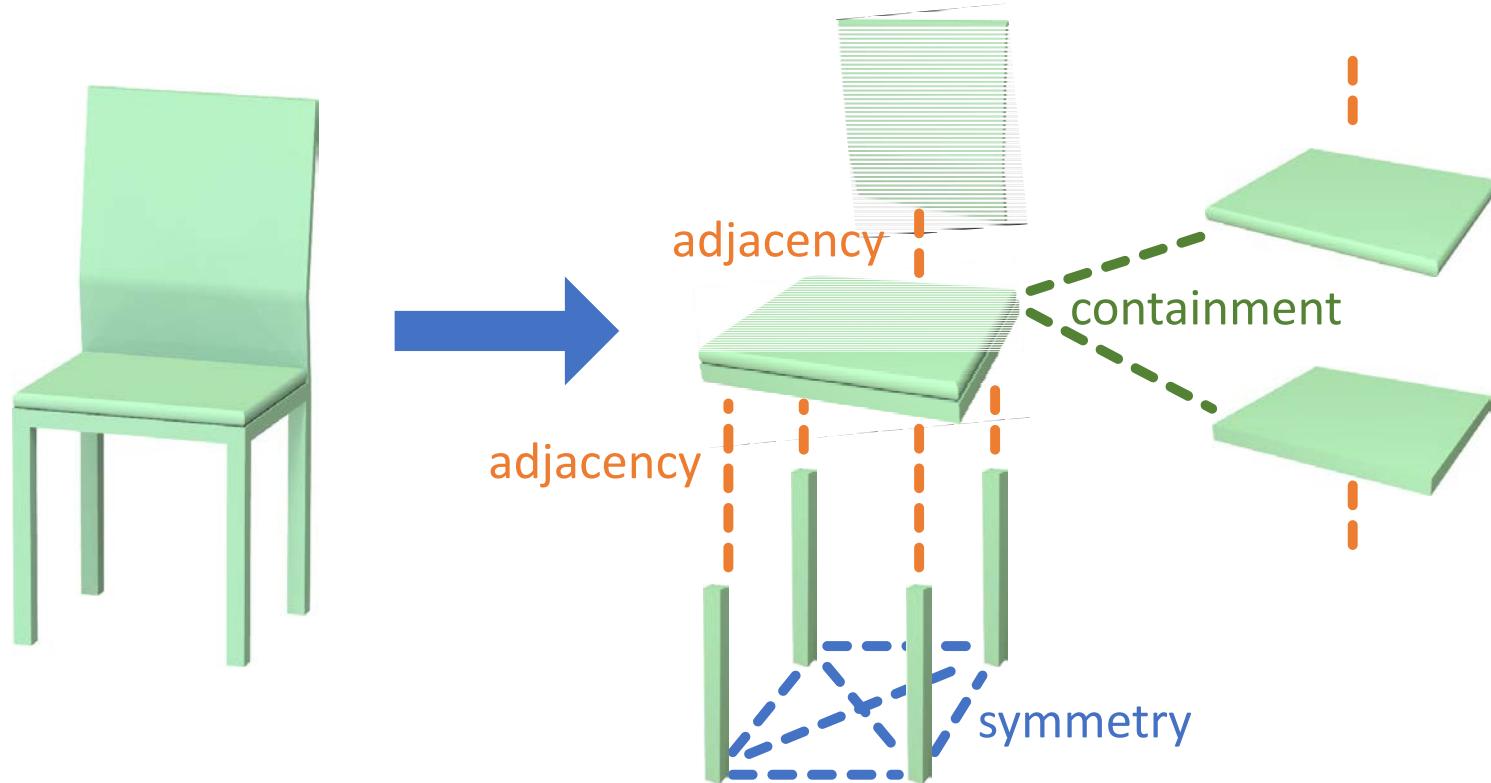
- Style similarity measure from:  
“Elements of Style: Learning Perceptual Shape Style Similarity”,  
Lun, Kalogerakis, Sheffer, TOG 2015

- Functional compatibility measure

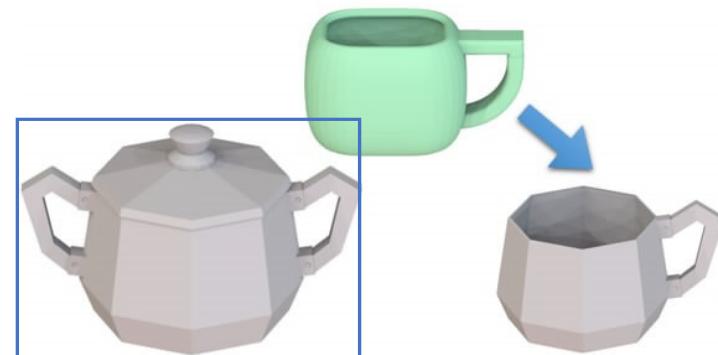
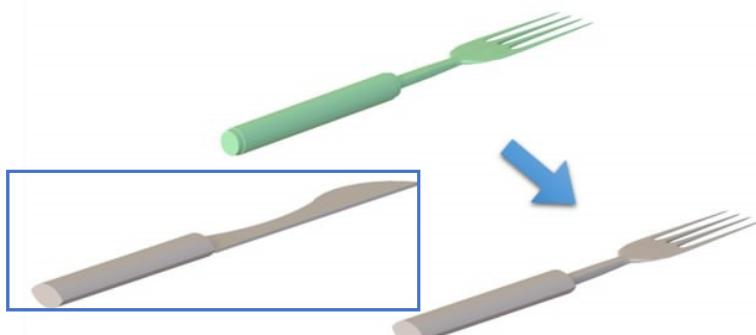
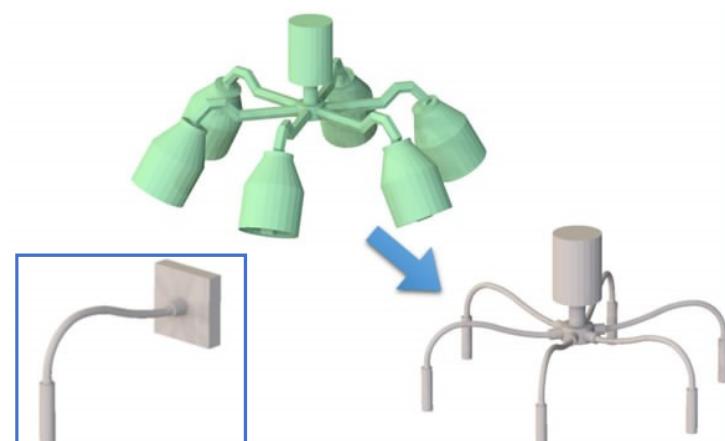
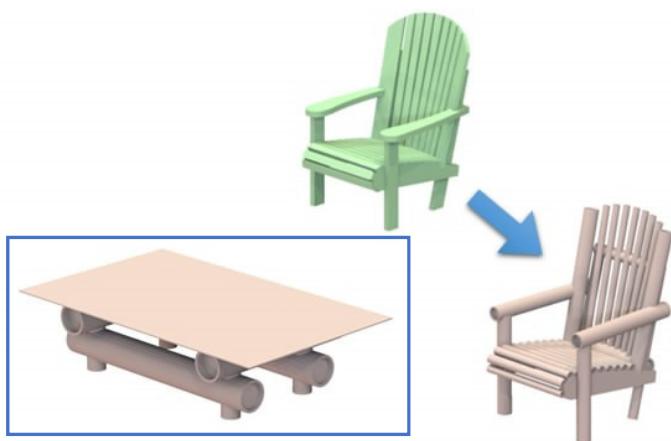
$$D_{func}(\text{green chair}, \text{purple stool}) = 1.32$$
A 3D rendering of a simple green chair and a small purple stool with curved legs.

$$D_{func}(\text{green chair}, \text{purple chair}) = 0.57$$
A 3D rendering of a green chair and a purple chair with a similar curved backrest and four legs.

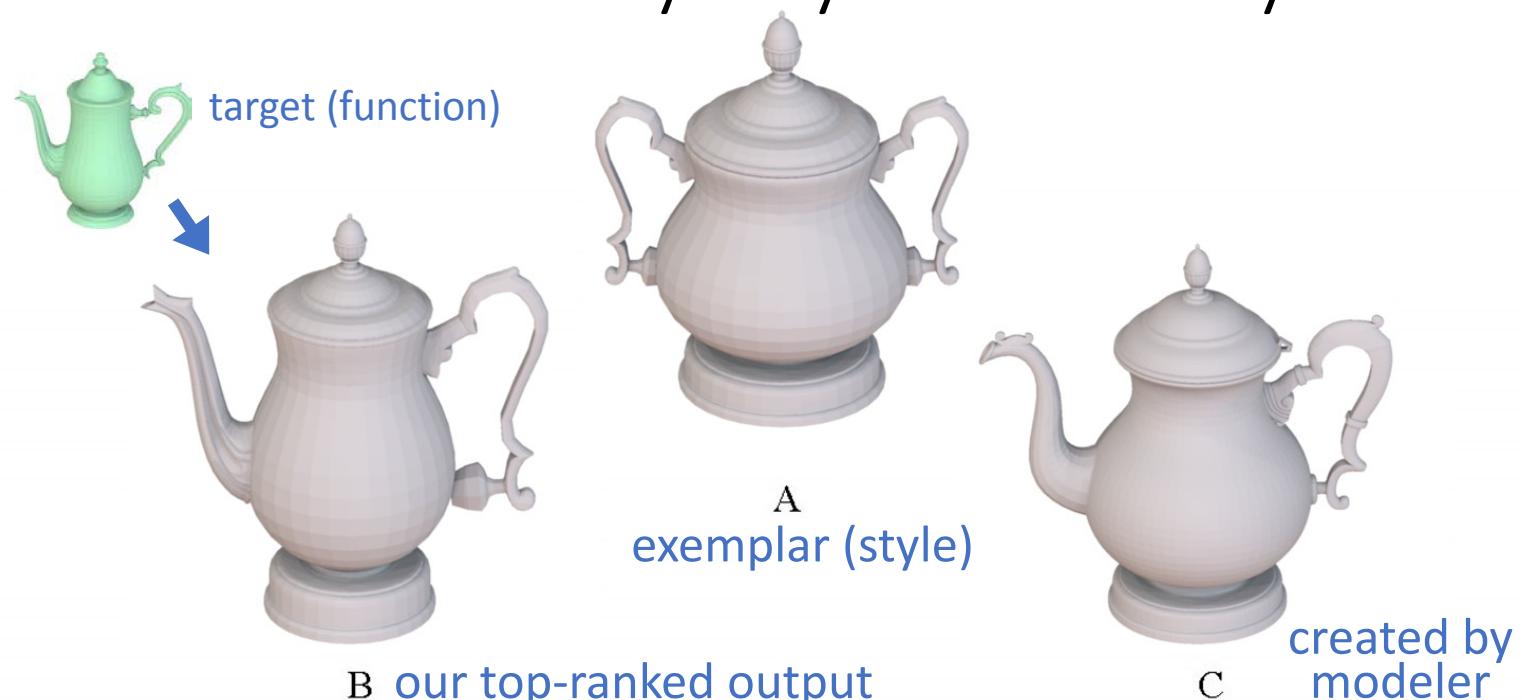
# Element relation graph



# Results



# User study: style similarity

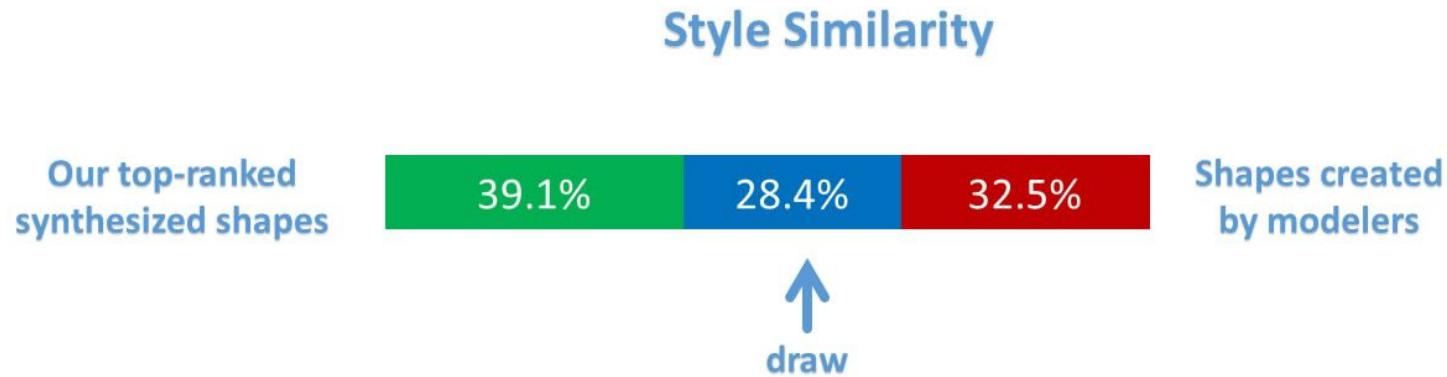


Which of the two shapes on the bottom (B or C) is more similar style-wise to the shape on the top (A)?

- B
- C
- can't tell - Both B and C are very similar to A in style
- can't tell - Neither B nor C are similar to A in style

[NEXT](#)

# User study: style similarity



## Summary & Future work

- Algorithm for synthesizing shapes by **transferring style** between objects with **different structure and functionality**
- Develop **style-aware generative models** that are capable of generating plausible shape structure & accurate surface geometric details

# Thank you!

Our project web pages:

<https://people.umass.edu/~yangzhou/scenegraphnet/>

&

<http://people.cs.umass.edu/~zlun/papers/StyleTransfer>

