# Reasoning over graphs

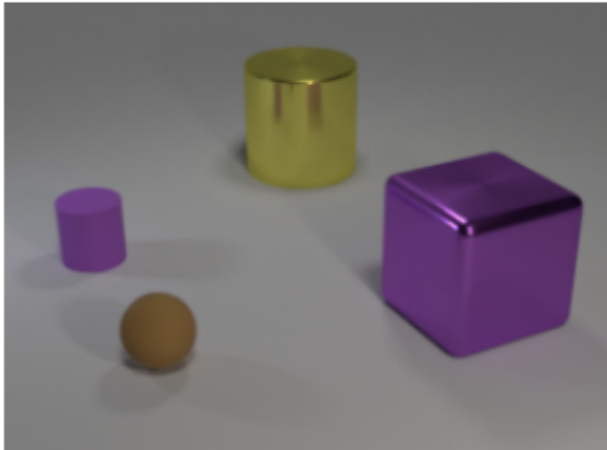https://neuralreasoning.github.io/

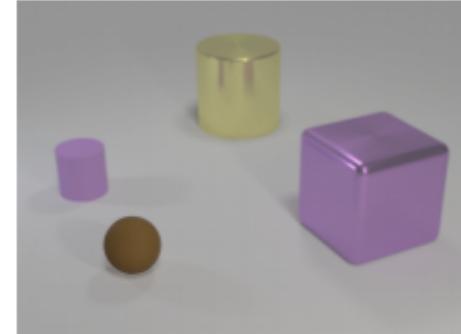Presented by Vuong Le

# Reasoning on Graphs

- Relational questions: requiring explicit reasoning about the relations between multiple objects
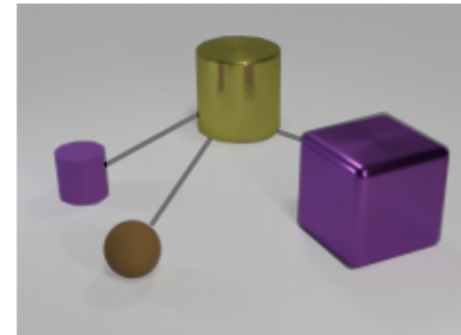


**Original Image:**

**Non-relational question:**
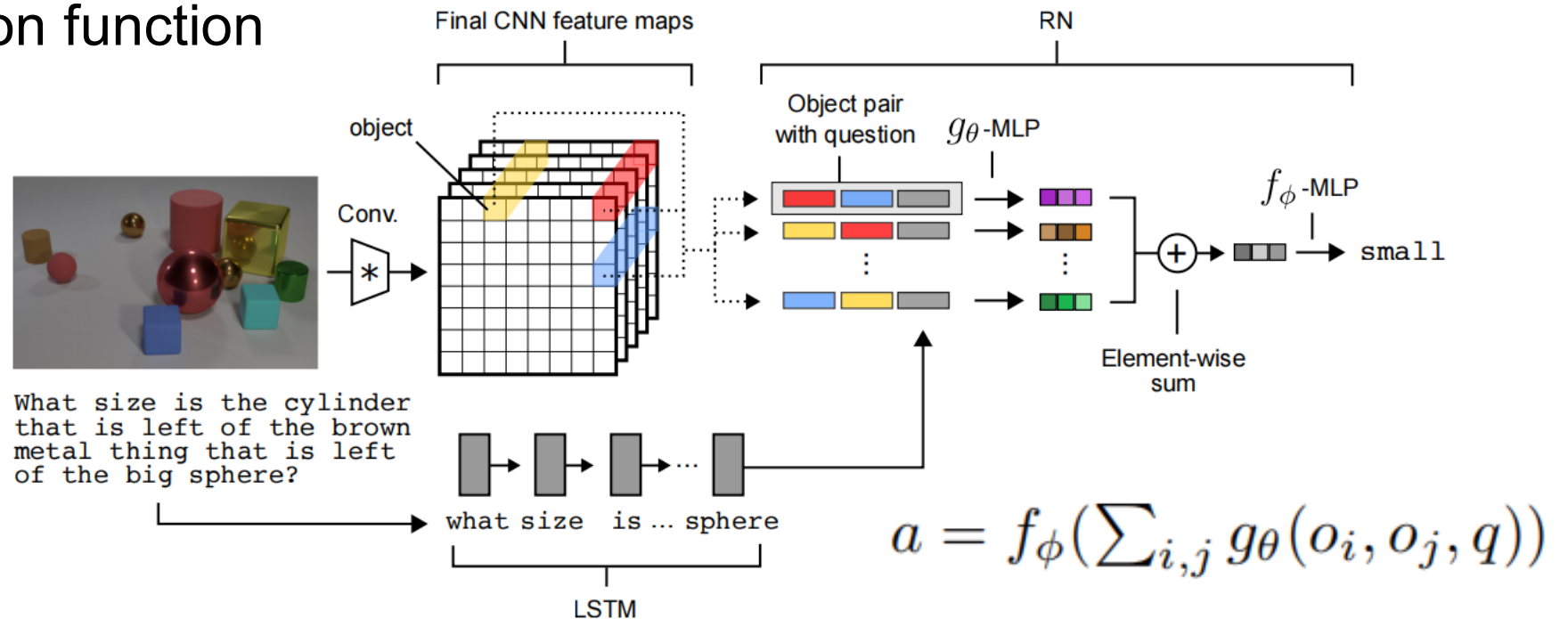
What is the size of the brown sphere?

**Relational question:**

Are there any rubber things that have the same size as the yellow metallic cylinder?

Figure credit: Santoro et al 2017

# Relation networks (Santoro et al 2017)

- Relation networks $\quad \mathrm{RN}(O) = f_\phi \left( \sum_{i,j} g_\theta(o_i, o_j) \right)$

- $f_\phi$ and $g_\theta$ are neural functions

- $g_\theta$ generate "relation" between the two objects

- $f_\phi$ is the aggregation function
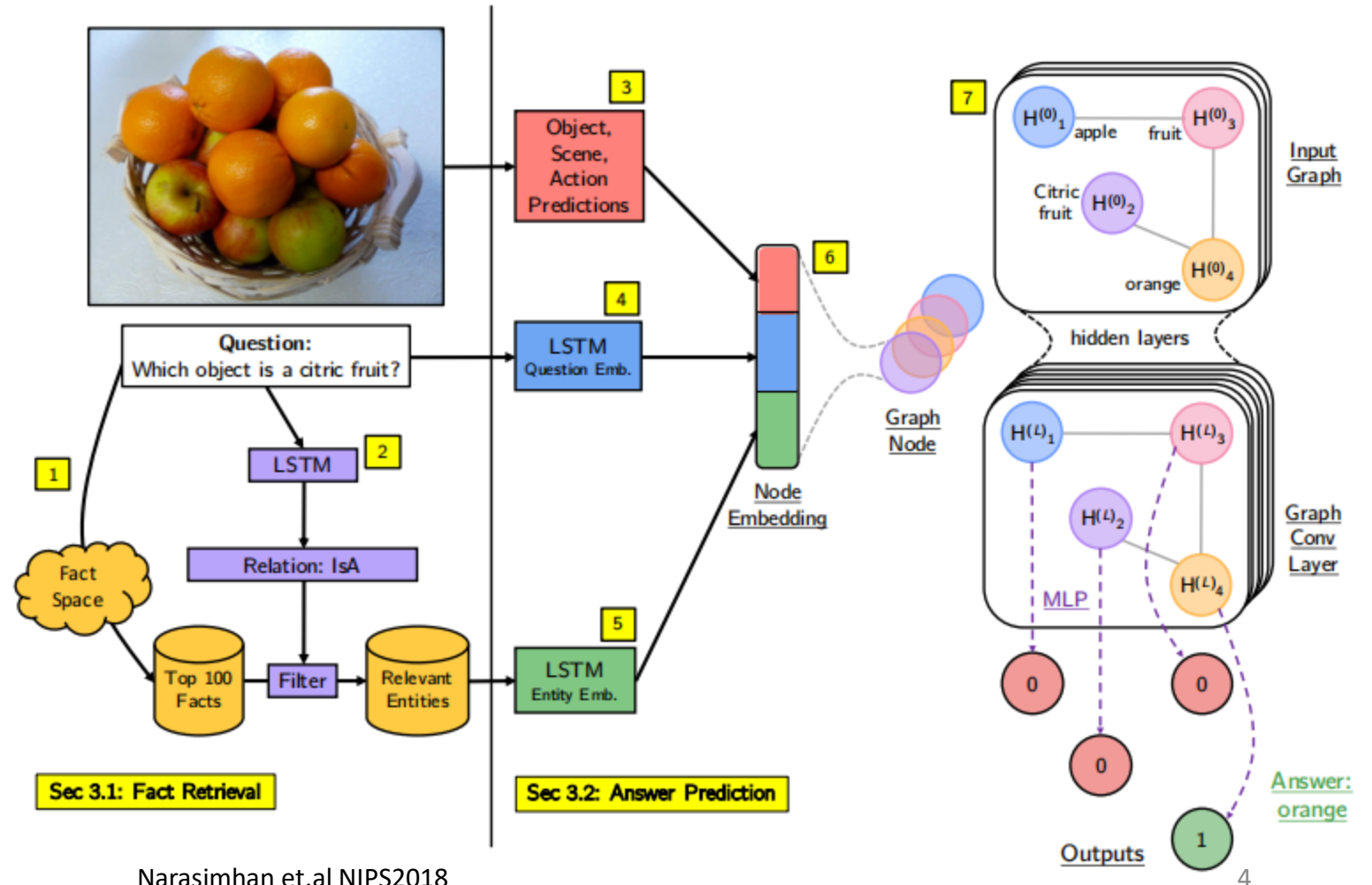


$$a = f_\phi(\sum_{i,j} g_\theta(o_i, o_j, q))$$

→ The relations here are implicit, over-complete, pair-wise
→ inefficient, and lack expressiveness

# Reasoning with Graph convolution networks

- Input graph is built from image entities and question
- GCN is used to gather facts and produce answer

→ The relations are now explicit and pruned

→ But the graph building is very stiff:
- Unrecoverable from mistakes
- Information during reasoning are not used to build graphs
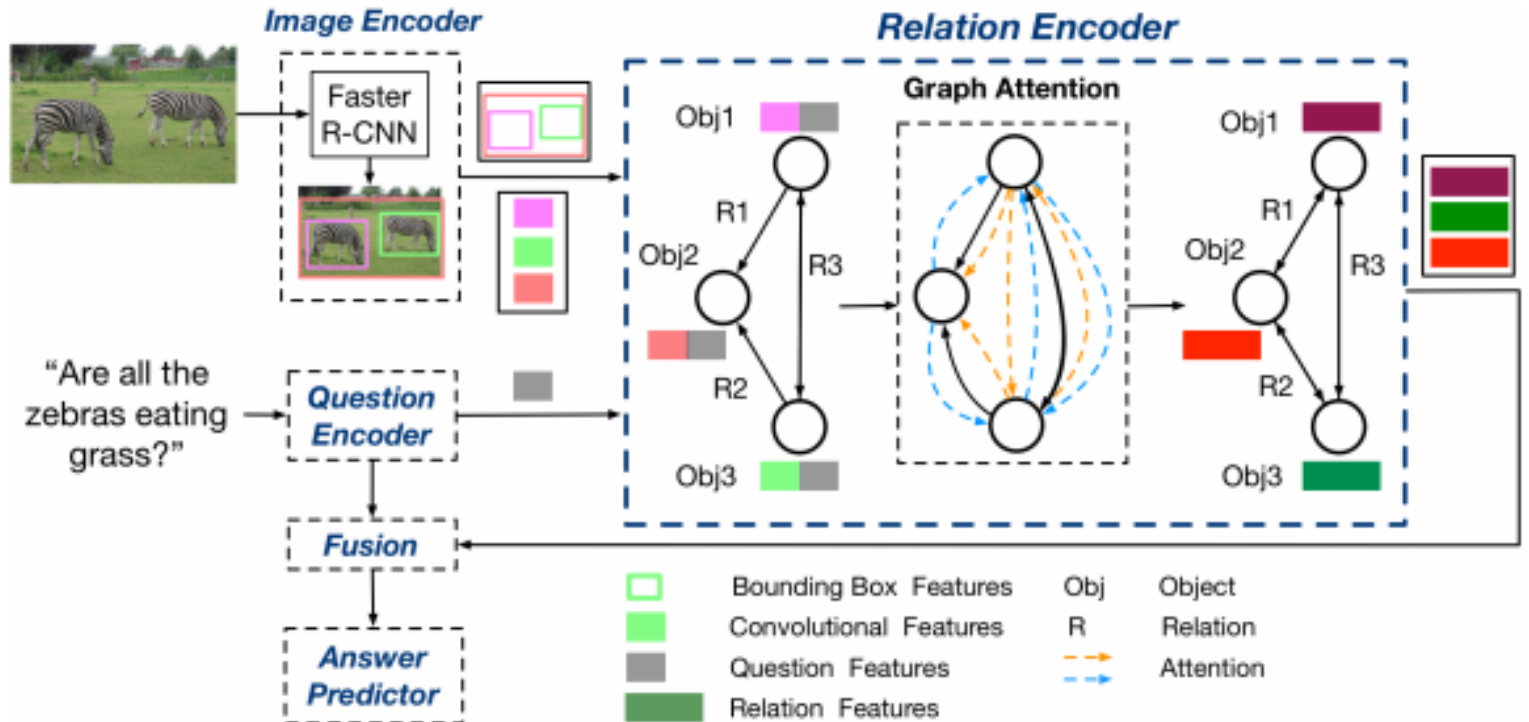
Narasimhan et.al NIPS2018

# Reasoning with Graph attention networks

- The graph is determined during reasoning process with attention mechanism

→The relations are now adaptive and integrated with reasoning

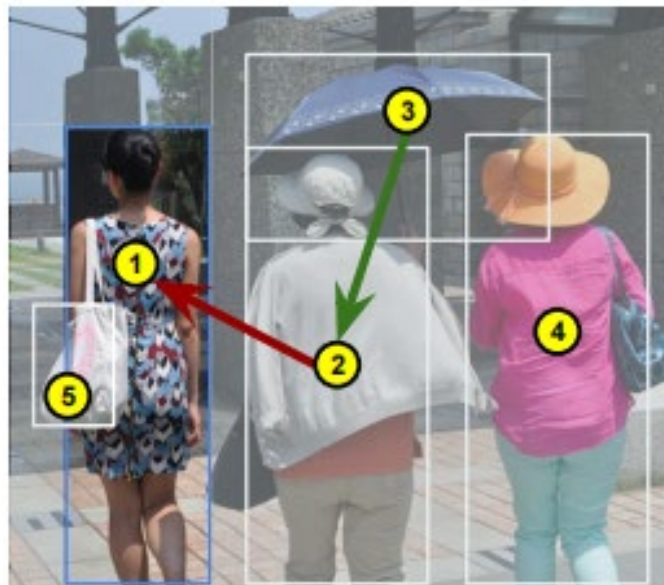→ Are the relations singular and static?

# Dynamic reasoning graphs

- On complex questions, multiple sets of relations are needed
- We need not only multi-step but also multi-form structures
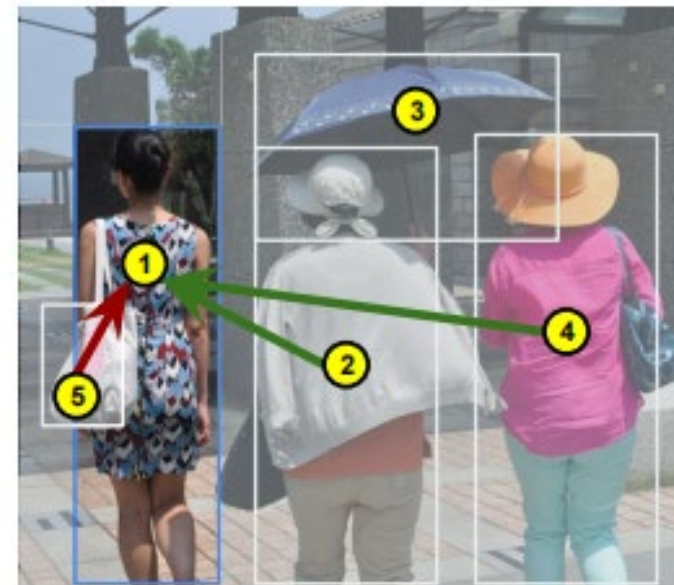- Let's do multiple dynamically–built graphs!



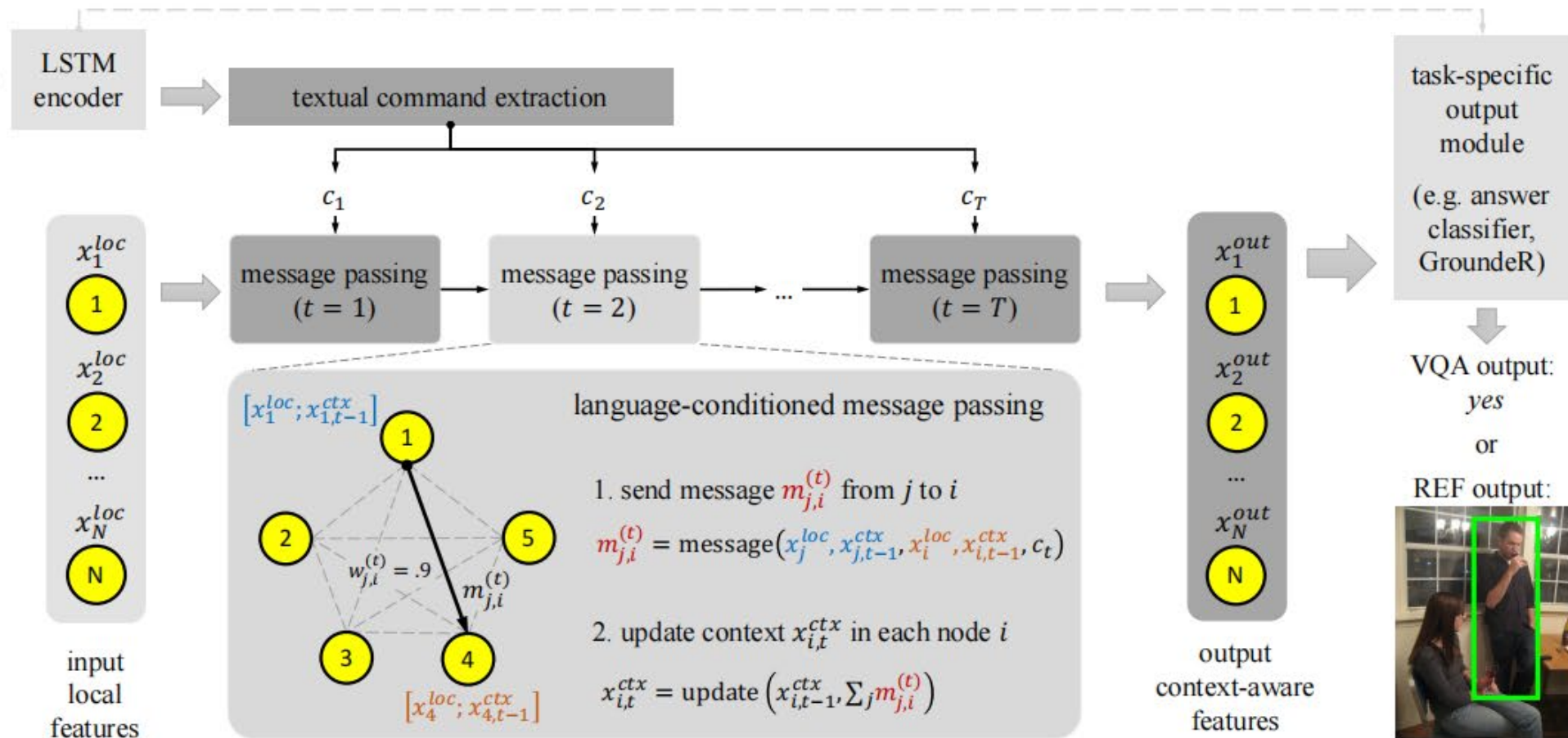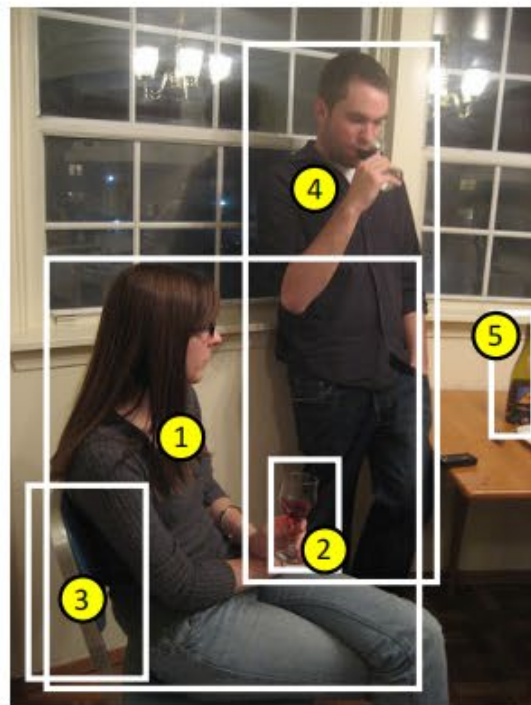**Question**: Is there a person to the left of the woman holding a blue umbrella?

**Answer**: Yes

**Question**: Is the left-most person holding a red bag?

**Answer**: No

LCGN, Hu et.al. ICCV19

# Dynamic reasoning graphs



Is there *a man on the right of a person sitting on a chair holding a wine glass?*

LSTM encoder → textual command extraction

$c_1$    $c_2$    $c_T$

message passing $(t = 1)$ → message passing $(t = 2)$ → ... → message passing $(t = T)$

$x_1^{loc}$   1
$x_2^{loc}$   2
...
$x_N^{loc}$   N

input local features

$[x_1^{loc}; x_{1,t-1}^{ctx}]$

language-conditioned message passing

$w_{j,i}^{(t)} = .9$   $m_{j,i}^{(t)}$

$[x_4^{loc}; x_{4,t-1}^{ctx}]$

1. send message $m_{j,i}^{(t)}$ from $j$ to $i$

$$m_{j,i}^{(t)} = \text{message}\left(x_j^{loc}, x_{j,t-1}^{ctx}, x_i^{loc}, x_{i,t-1}^{ctx}, c_t\right)$$

2. update context $x_{i,t}^{ctx}$ in each node $i$

$$x_{i,t}^{ctx} = \text{update}\left(x_{i,t-1}^{ctx}, \sum_j m_{j,i}^{(t)}\right)$$

$x_1^{out}$   1
$x_2^{out}$   2
...
$x_N^{out}$   N

output context-aware features

task-specific output module

(e.g. answer classifier, GroundeR)
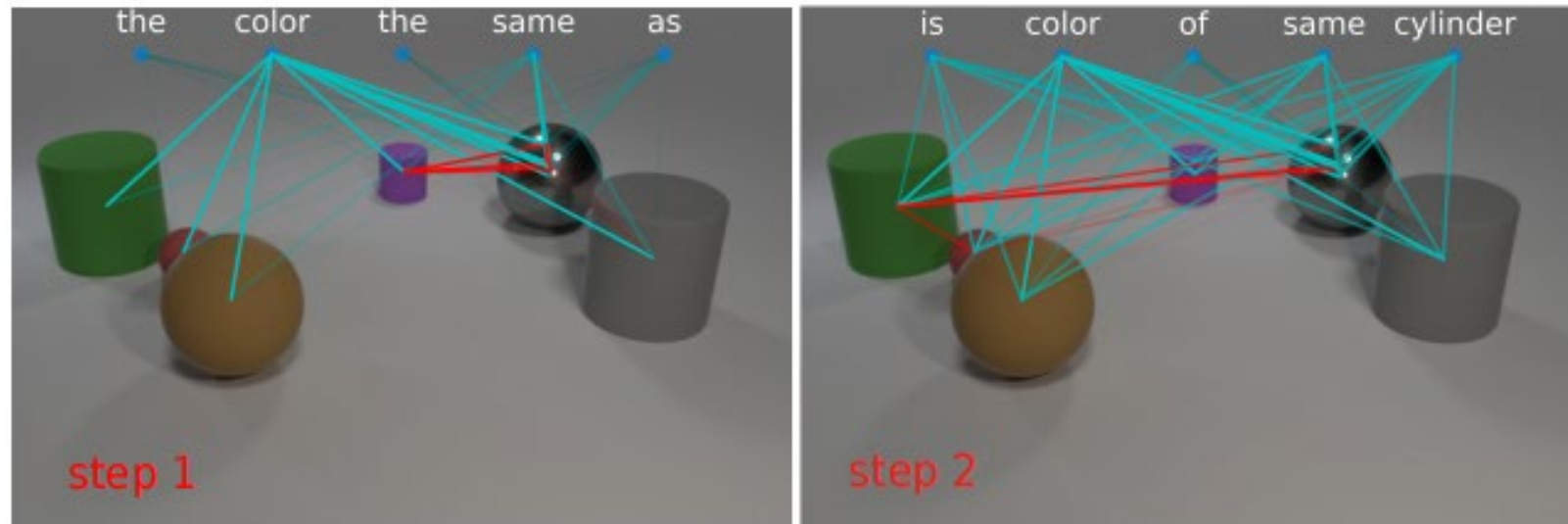
VQA output: *yes*

or

REF output:

→ The questions so far act as an unstructured command in the process
→ *Aren't their structures and relations important too?*

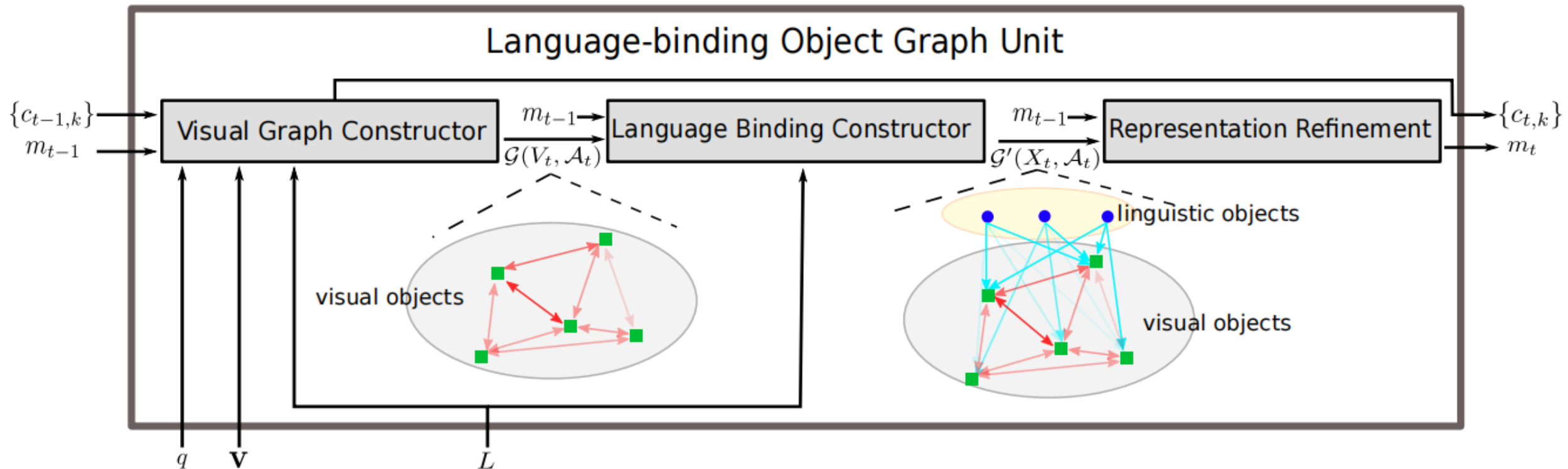# Reasoning on cross-modality graphs

- Two types of nodes: Linguistic entities and visual objects
- Two types of edges:
  - Visual relations
  - Linguistic-visual binding *(as a fuzzy grounding)*
- Adaptively updated during reasoning
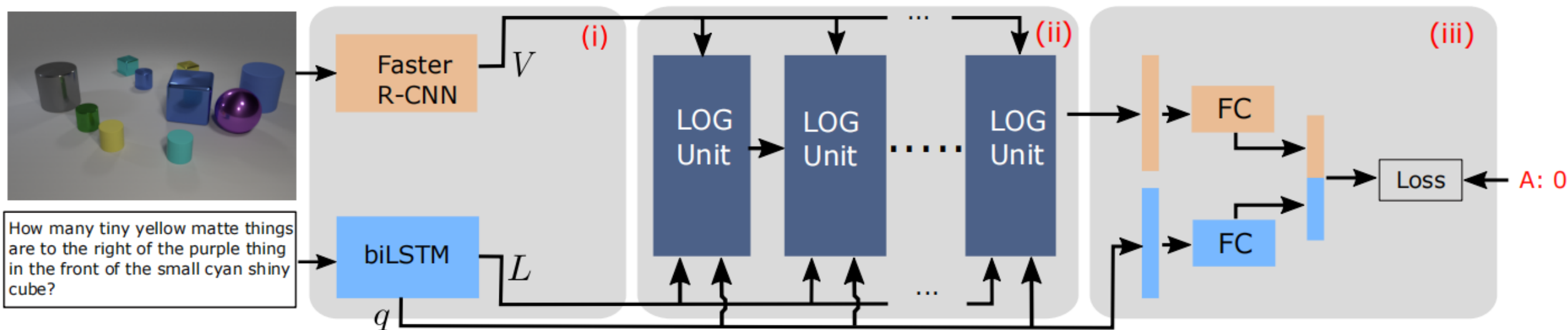


LOGNet, T.M Le et.al. IJCAI 2020

# Language-binding Object Graph (LOG) Unit

- Graph constructor: build the dynamic vision graph
- Language binding constructor: find the dynamic L-V relations
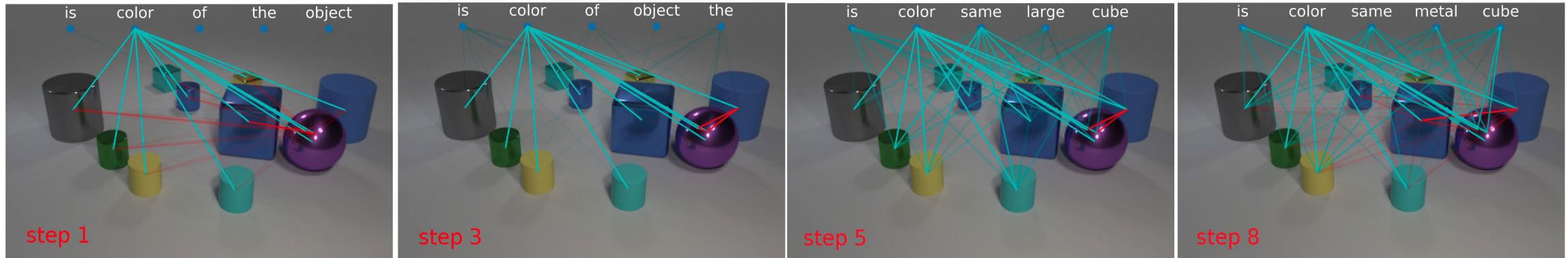


LOGNet, T.M Le et.al. IJCAI2020

# LOGNet: multi-step visual-linguistic binding

- Object-centric representation ✓

- Multi-step/multi-structure compositional reasoning ✓

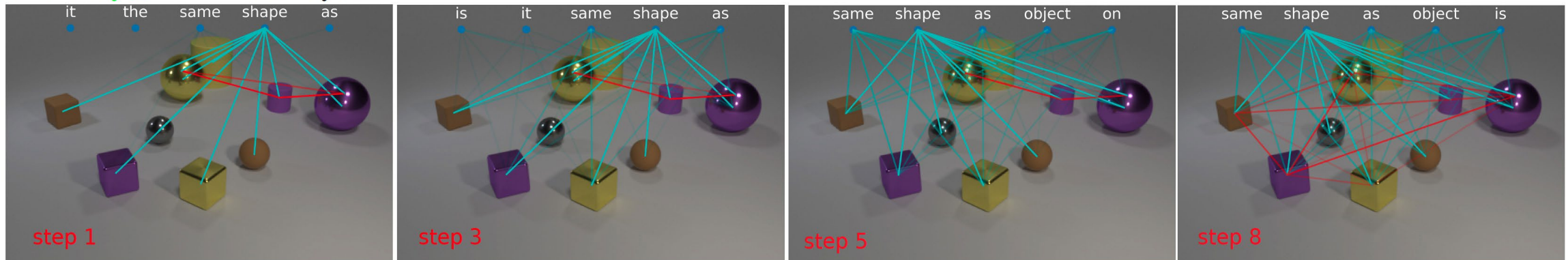- Linguistic-vision detail interaction ✓



LOGNet, T.M Le et.al. IJCAI2020

# Dynamic language-vision graphs in actions



**Question**: Is the color of the big matte object the same as the large metal cube?
**Prediction**: yes     **Answer**: yes



**Question**: There is a tiny purple rubber thing; does it have the same shape as the brown object that is on the left side of the rubber sphere?
**Prediction**: no     **Answer**: no

# We got sets and graphs, how about sequences?

- Videos pose another challenge for visual reasoning: the dynamics through time.
- Sets and graphs now becomes sequences of such.
- Temporal relations are the key factors
- The size of context is a core issue

→Lecture 8 will address these



(a) Question: What does the girl do 9 times?
Ground truth: blocks a person's punch

(b) Question: What does the man do before turning body to left?
Ground truth: breath