Language
Technologies
Institute

Carnegie
Mellon
University

# Multimodal Machine Learning

## Lecture 5.2: Structured Representations and Reasoning

Louis-Philippe Morency

# Administrative Stuff

Language Technologies Institute

Carnegie Mellon University

# Second Project Assignment (Due Sunday 10/8)

Main goals:

1.  Help clarify and expand your research ideas
    - Build qualitative intuitions by directly studying the original data
    - Perform analyses on your dataset, relevant to your research ideas
2.  Understand the structure in your data and modalities
    - Perform analyses and visualizations to understand each modality
    - Study representations from language and visual modalities

Two types of analyses:

- Idea-oriented analyses
- Modality-oriented analyses

# Lecture Schedule

| Classes | Tuesday Lectures | Thursday Lectures |
|---|---|---|
| **Week 1**<br>8/29 & 8/31 | **Course introduction**<br><ul><li>Multimodal core challenges</li><li>Course syllabus</li></ul> | **Multimodal applications and datasets**<br><ul><li>Research tasks and datasets</li><li>Team projects</li></ul> |
| **Week 2**<br>9/5 & 9/7<br>Read due: 9/9 | **Unimodal representations**<br><ul><li>Dimensions of heterogeneity</li><li>Visual representations</li></ul> | **Unimodal representations**<br><ul><li>Language representations</li><li>Signals, graphs and other modalities</li></ul> |
| **Week 3**<br>9/12 & 9/14<br>Read due: 9/16<br>Proj. Due: 9/13 | **Multimodal representations**<br><ul><li>Cross-modal interactions</li><li>Multimodal fusion</li></ul> | **Multimodal representations**<br><ul><li>Coordinated representations</li><li>Multimodal fission</li></ul> |
| **Week 4**<br>9/19 & 9/21<br>*Proj. due: 9/24* | **Multimodal alignment and grounding**<br><ul><li>Explicit alignment</li><li>Multimodal grounding</li></ul> | **Alignment and representations**<br><ul><li>Self-attention transformer models</li><li>Masking and self-supervised learning</li></ul> |
| **Week 5**<br>9/26 & 9/28<br>Read due: 9/30 | **Multimodal transformers – Part 1**<br><ul><li>Language pretraining</li><li>Multimodal transformers</li></ul> | **Multimodal Reasoning**<br><ul><li>Hierarchical and graph representations</li><li>Modular and neuro-symbolic models</li></ul> |
| **Week 6**<br>10/3 & 10/5<br>*Proj. due: 10/8* | **Multimodal transformers – Part 2**<br><ul><li>Image and video transformers</li><li>Vision-language transformers</li></ul> | ***Multimodal language grounding***<br><ul><li>Guest lecturer: Jack Hessel</li><li>Vision, language and grounding</li></ul> |

# Multimodal Machine Learning

## Lecture 5.2: Structured Representations and Reasoning
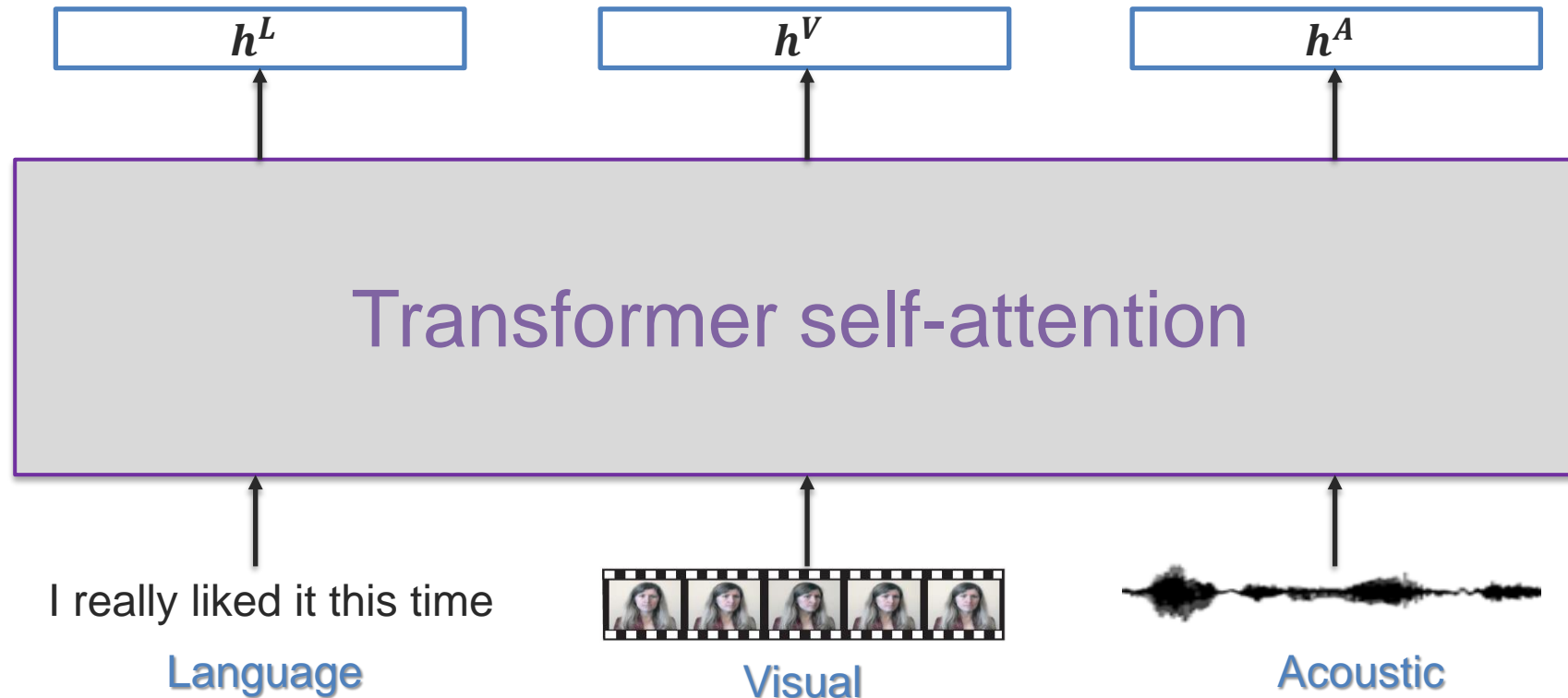
**Louis-Philippe Morency**

*\* Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yonatan Bisk.*

# Objectives of today's class

- Multimodal transformers
    - Modality-shift transformer (MAG-BERT)
- Sequence-to-sequence modeling with Transformers
- Going beyond sequences
    - Graph representations
        - Graph neural networks
    - Hierarchical representations
    - Modular representations
        - Neural module networks
        - Neuro-symbolic networks

# Language-Vision Transformers

# Simple Solution: Contextualized Multimodal Embeddings

$h^L$

$h^V$

$h^A$

Transformer self-attention

I really liked it this time

Language

Visual

Acoustic

# Multimodal Transformer – Pairwise Cross-Modal



**Visual**

**Vocal**

**Verbal**

"I like…"

Unimodal Representations

Cross-modally Contextualized Unimodal Representations

Multimodal representation

**Cross-Modal Attention Block    x N layers**

# Reminder: Modality-Shifting Fusion

Primary modality

Secondary modalities

$x_A$

$x_B$

$x_C$

gate

shift

$z$

**Example with language modality:**

Primary modality: language

Secondary modalities: acoustic and visual

Negative-shifted representation

word: "expectations"

Positive-shifted representation

Wang et al., Words Can Shift: Dynamically Adjusting Word Representations Using Nonverbal Behaviors, AAAI 2019

# Modality-Shifting with Transformers

Multimodal Adaptation Gate (MAG) + BERT



Rahman et al., Integrating Multimodal Information in Large Pretrained Transformers, ACL 2020

# Memory for Multimodal Sequences

**Memory + aligned contextualized representations**
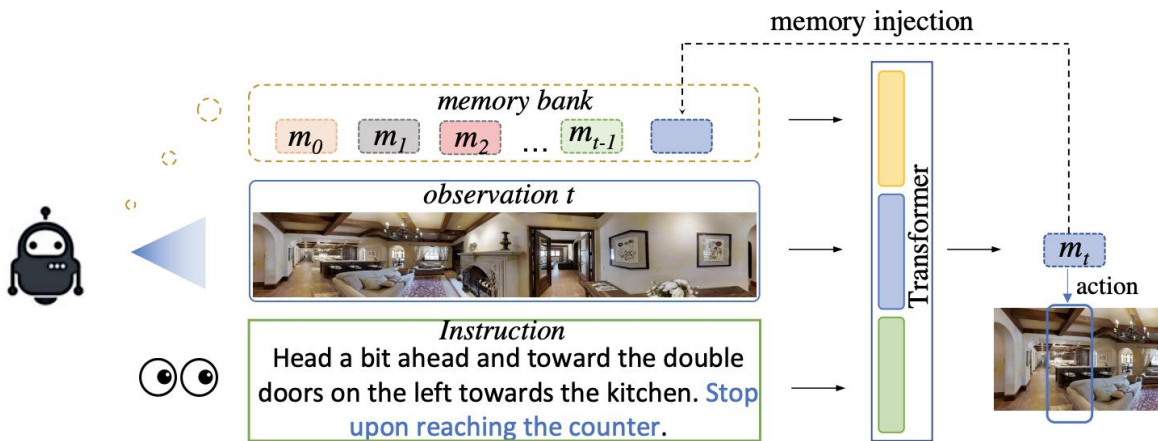
*Where have I visited previously?*



[Chen et al., History Aware Multimodal Transformer for Vision-and-Language Navigation. NeurIPS 2021]
[Lin et al., Multimodal Transformer with Variable-length Memory for Vision-and-Language Navigation. ECCV 2022]

# Memory for Multimodal Sequences

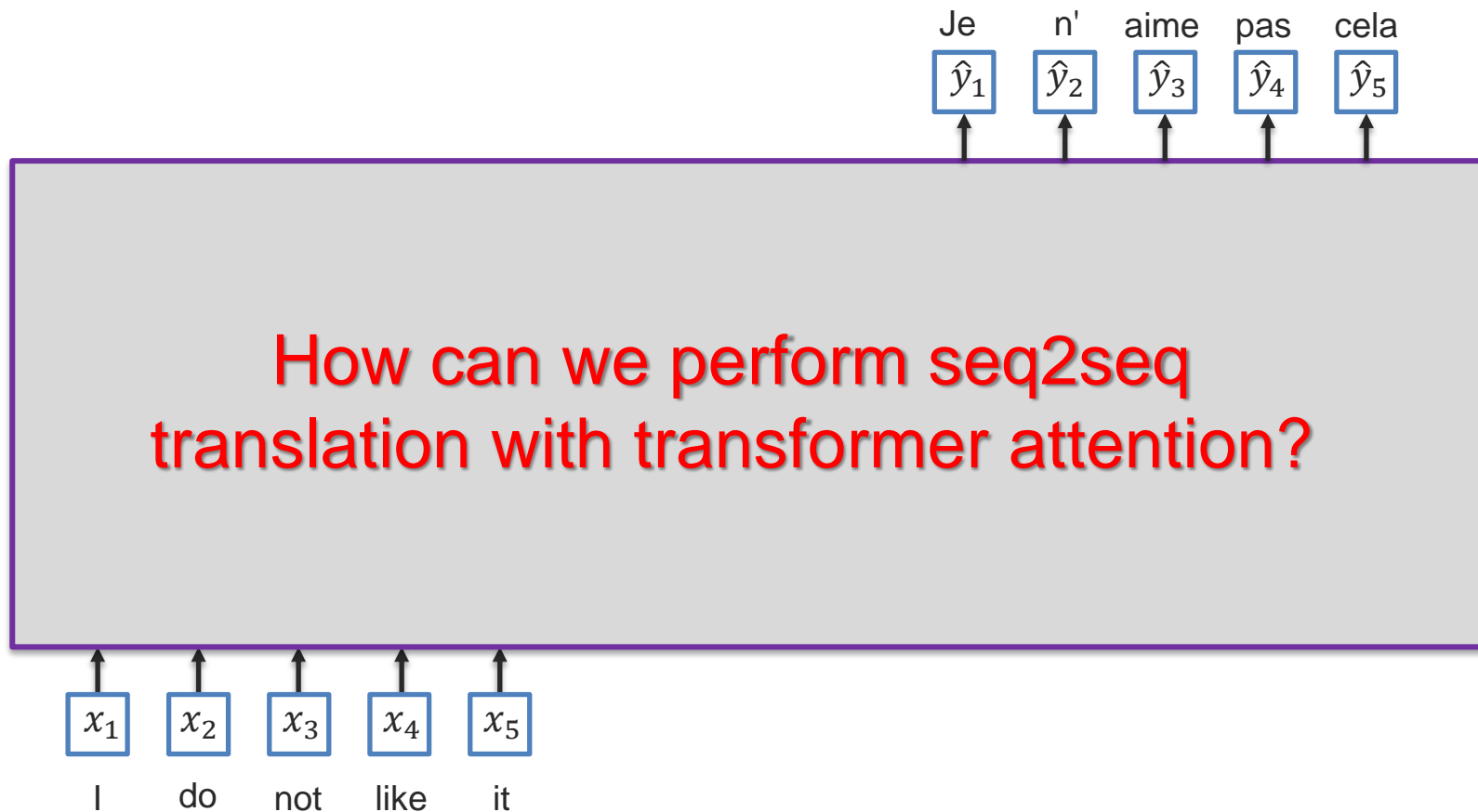**Memory + aligned contextualized representations**

*Where have I visited previously?*



+ Contextualized representations

+ Memory mechanisms

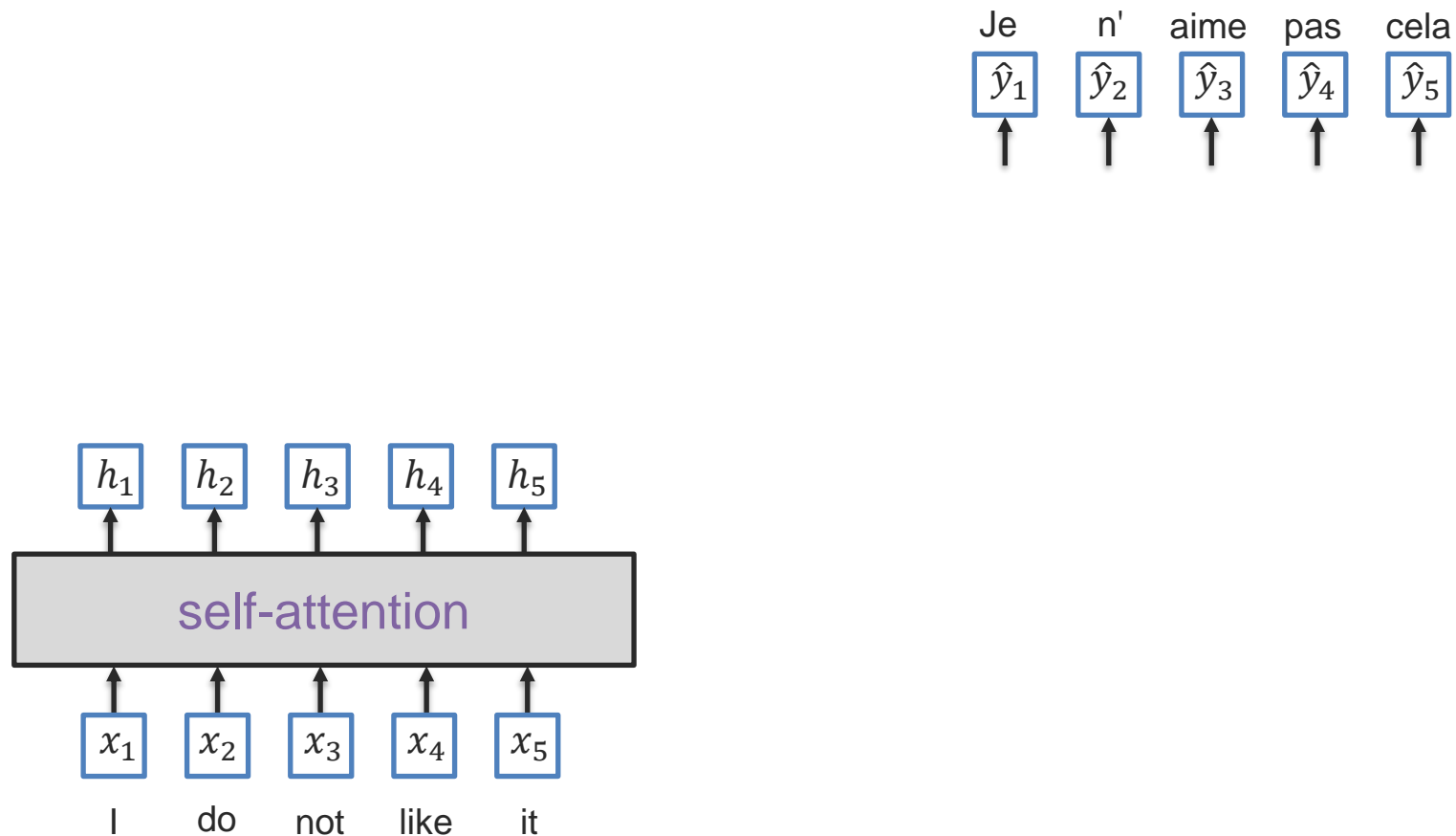[Chen et al., History Aware Multimodal Transformer for Vision-and-Language Navigation. NeurIPS 2021]
[Lin et al., Multimodal Transformer with Variable-length Memory for Vision-and-Language Navigation. ECCV 2022]
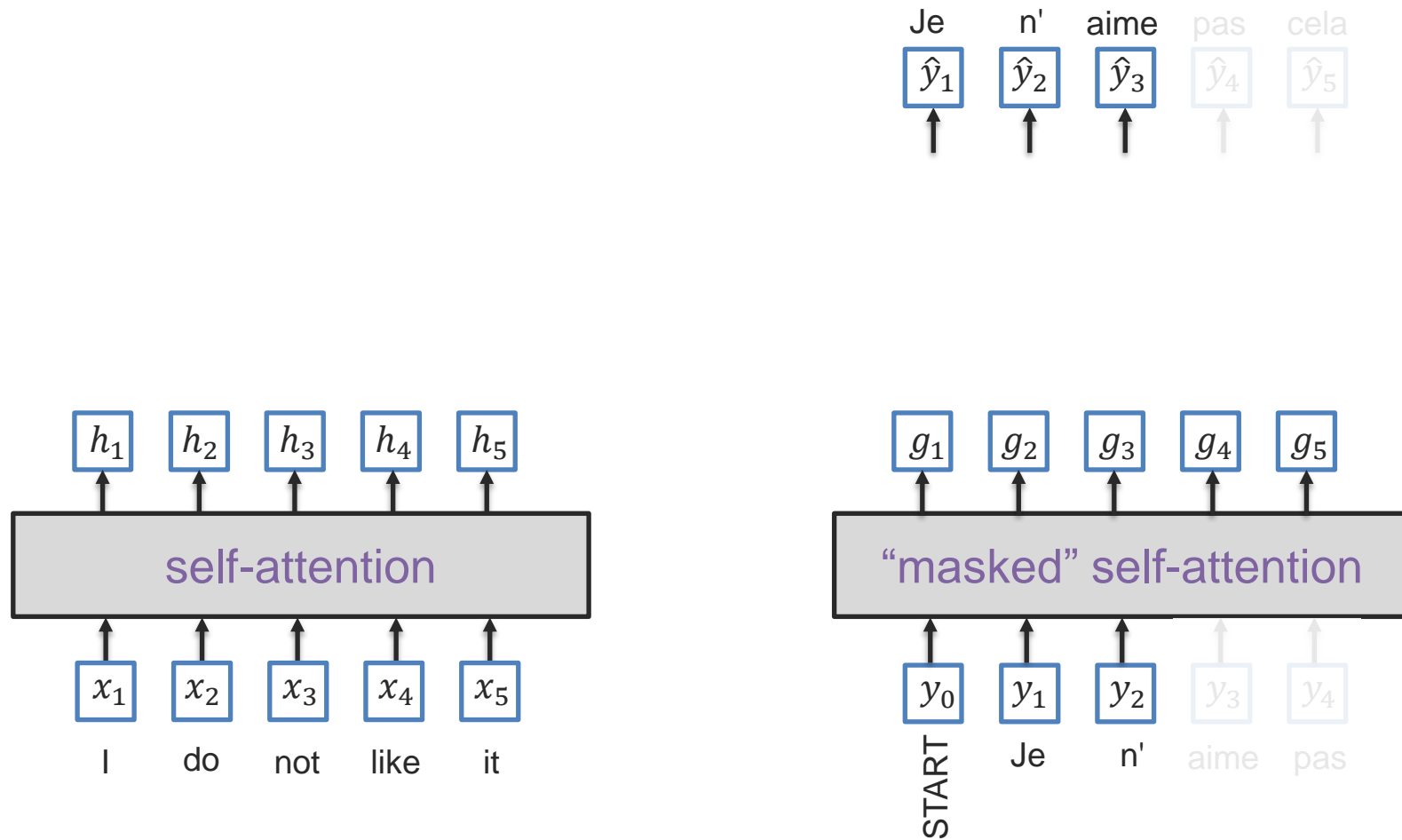
# Sequence-to-Sequence Using Transformer

# Sequence-to-Sequence Modeling

Je    n'    aime    pas    cela
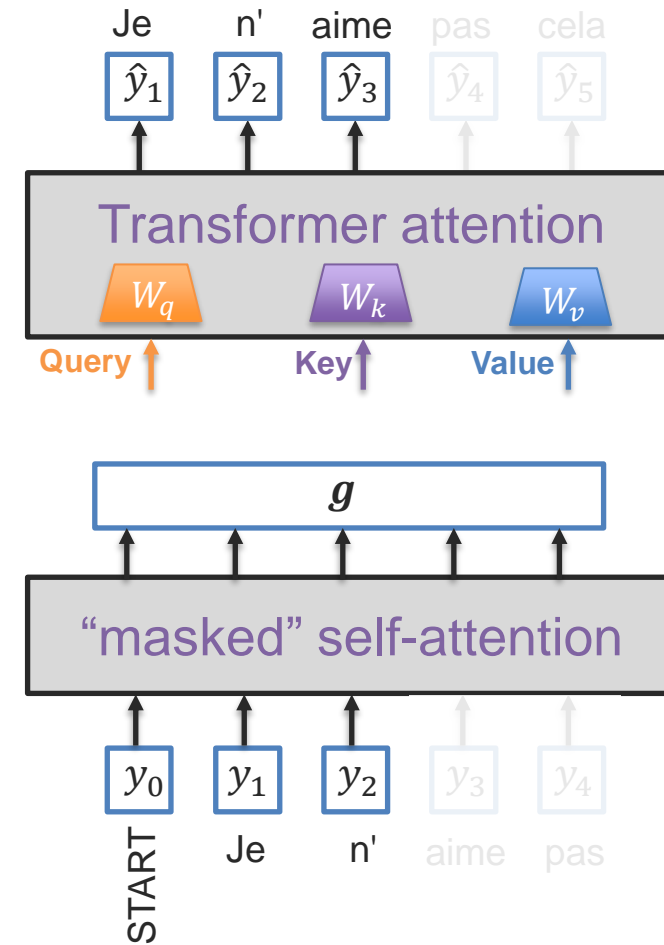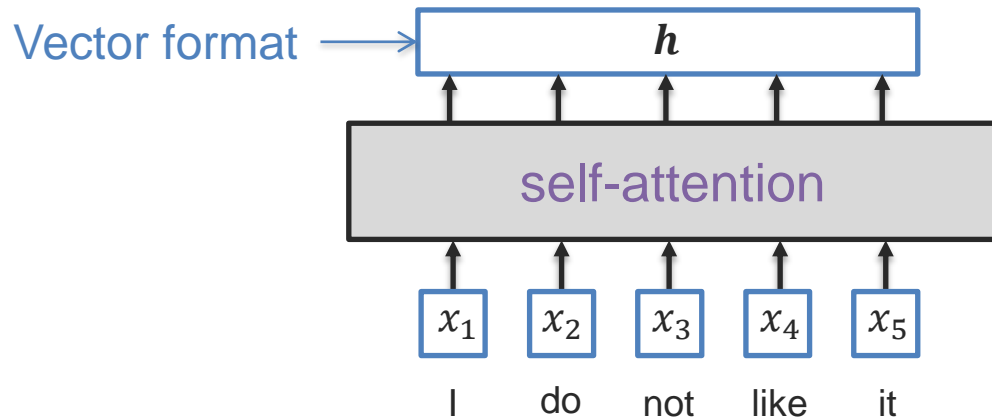
$\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$    $\hat{y}_4$    $\hat{y}_5$

## How can we perform seq2seq translation with transformer attention?

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$

I    do    not    like    it

# Seq2Seq with Transformer Attentions

Je   n'   aime   pas   cela

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$   $\hat{y}_4$   $\hat{y}_5$

$h_1$   $h_2$   $h_3$   $h_4$   $h_5$

self-attention

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$

I   do   not   like   it

# Seq2Seq with Transformer Attentions

Je    n'    aime    pas    cela

$\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$    $\hat{y}_4$    $\hat{y}_5$

$h_1$    $h_2$    $h_3$    $h_4$    $h_5$

self-attention

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$

I    do    not    like    it

$g_1$    $g_2$    $g_3$    $g_4$    $g_5$

"masked" self-attention

$y_0$    $y_1$    $y_2$    $y_3$    $y_4$

START    Je    n'    aime    pas

# Seq2Seq with Transformer Attentions

How should we connect the encoder and decoder self-attention to the transformer attention?

Je  n'  aime  pas  cela

$\hat{y}_1$  $\hat{y}_2$  $\hat{y}_3$  $\hat{y}_4$  $\hat{y}_5$

Transformer attention

$W_q$  $W_k$  $W_v$

**Query**  **Key**  **Value**

Vector format → $h$

self-attention

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

I  do  not  like  it

$g$

"masked" self-attention

$y_0$  $y_1$  $y_2$  $y_3$  $y_4$

START  Je  n'  aime  pas

# Seq2Seq with Transformer Attentions

# Going Beyond Sequences: Graph Representations

# Transformers – Fully-Connected Sequences



Should everything be connected to everything?

What if we have domain knowledge about connections?
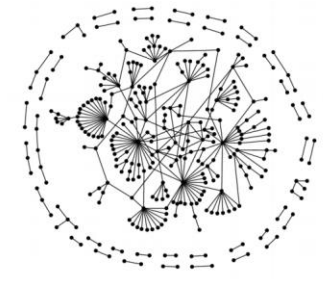
# Tree and Graph Networks

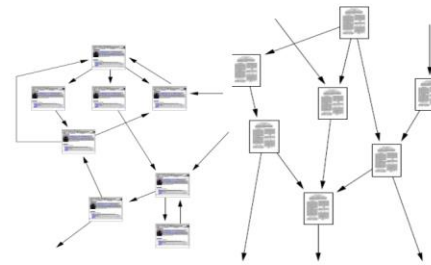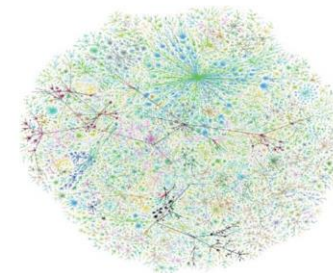**From linear chain models to tree and graph-structured models**
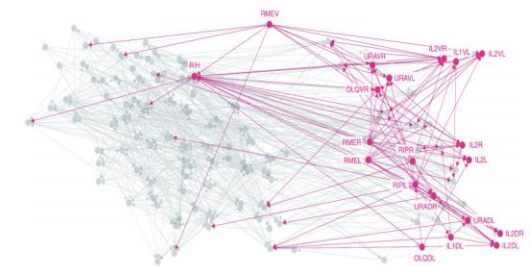


Social networks

Economic networks

Biomedical networks

Information networks: Web & citations
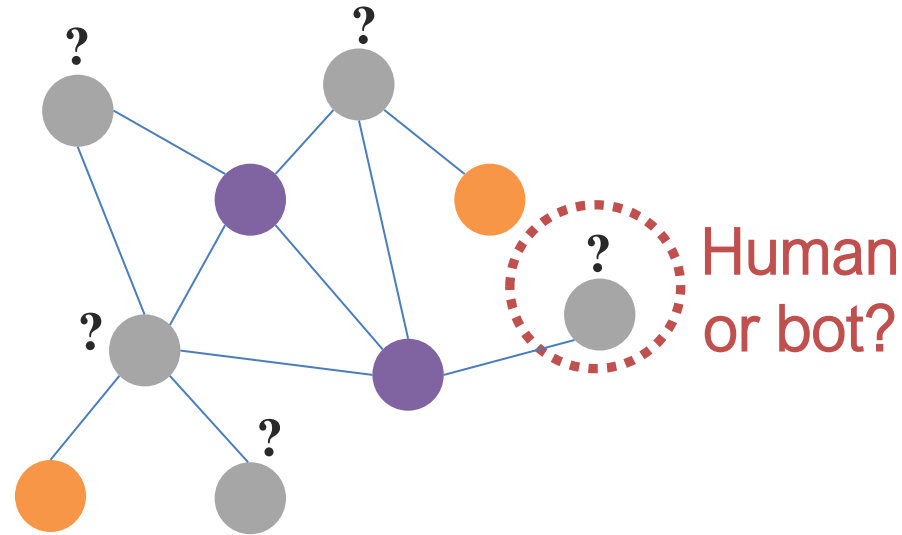
Internet

Networks of neurons

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]

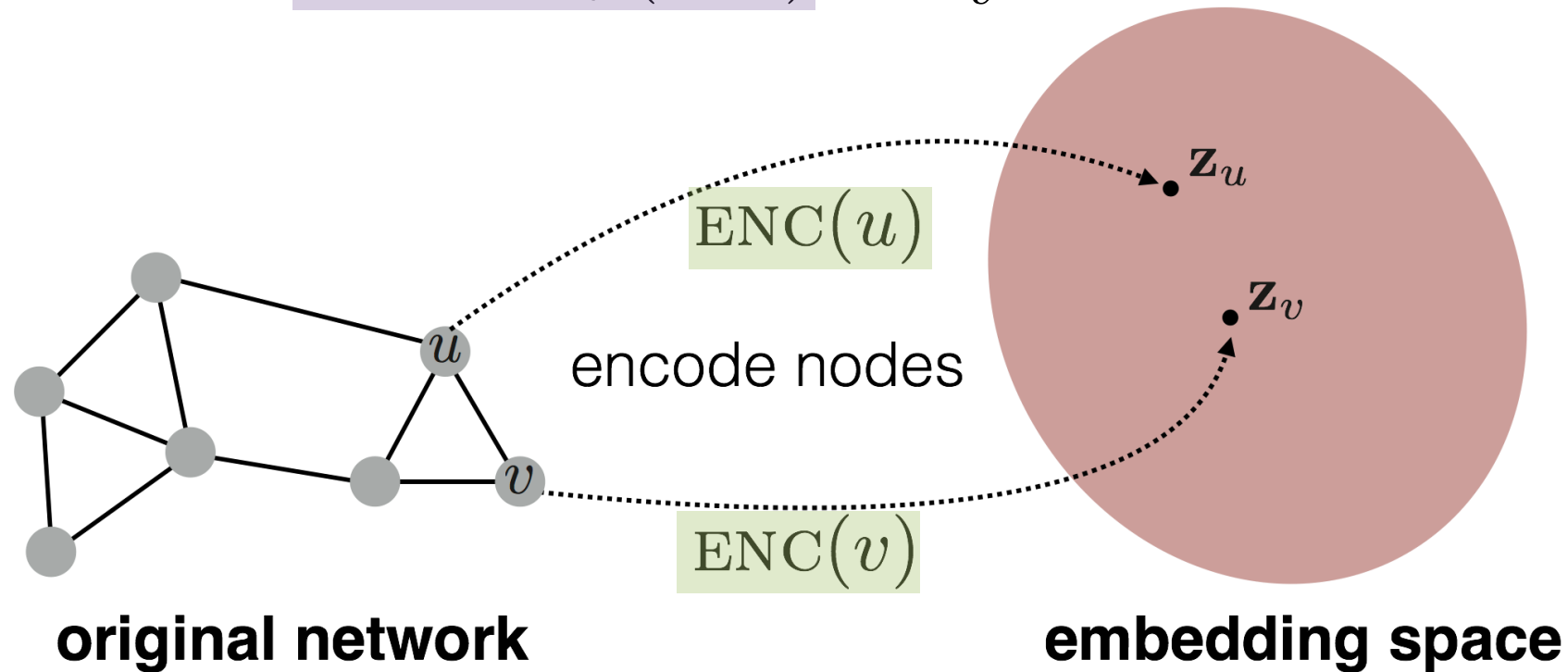Goal: Learn from labels associated with a subset of nodes (or with all nodes)



e.g., an online social network

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]

# Graphs – Unsupervised Task

Goal: Learn an embedding space where

$$\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$$



$\text{ENC}(u)$

$\text{ENC}(v)$

encode nodes

$\mathbf{z}_u$

$\mathbf{z}_v$

**original network**

**embedding space**

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]

# Graph Neural Nets

Assume we have a graph **G:**

    **V** is the set of vertices

    **A** is the binary adjacency matrix

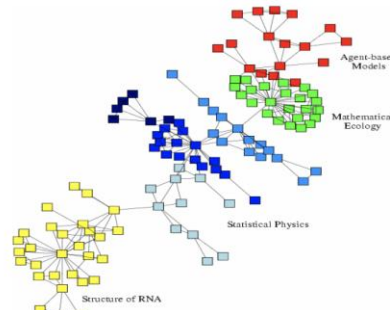    **X** is a matrix of node features:

- Categorical attributes, text, image data
  e.g. profile information in a social network
- …
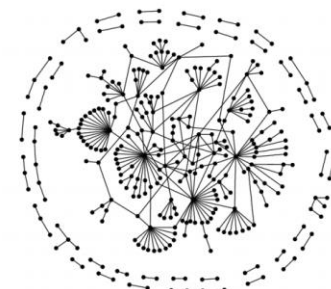
    **Y** is a vector of node labels (optional)



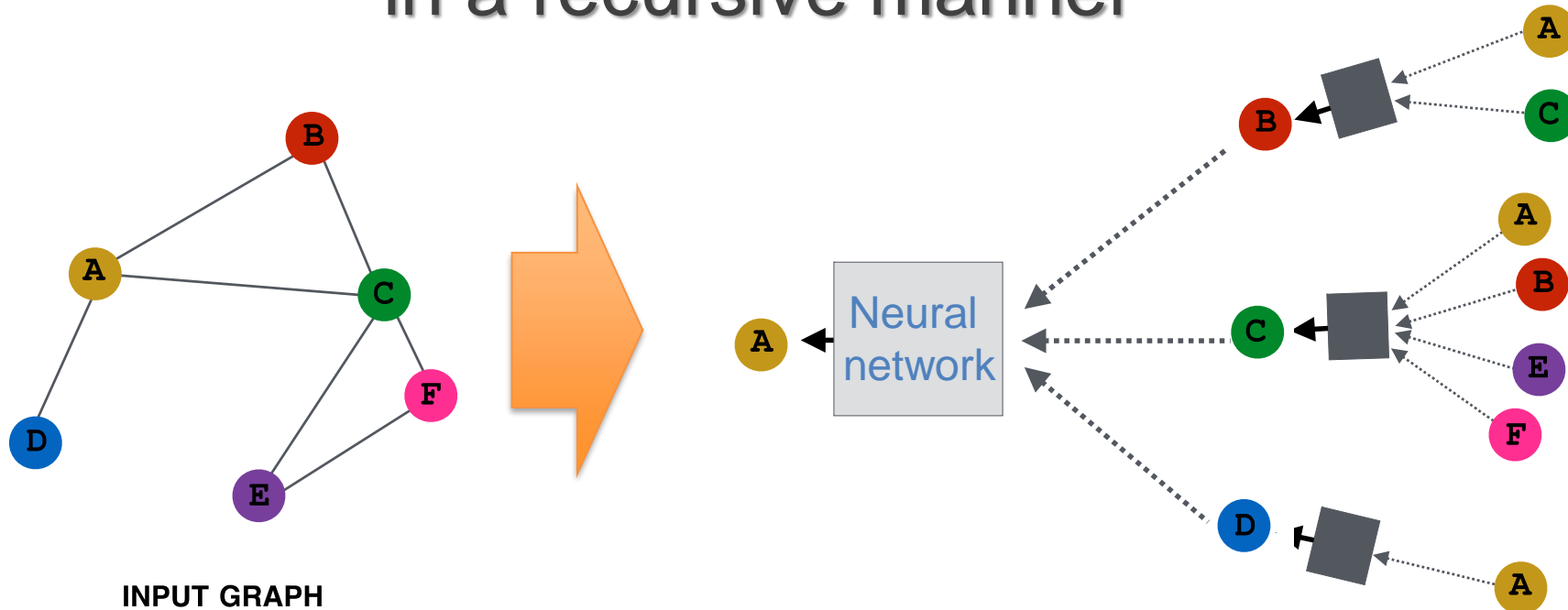Social networks    Economic networks    Biomedical networks

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]
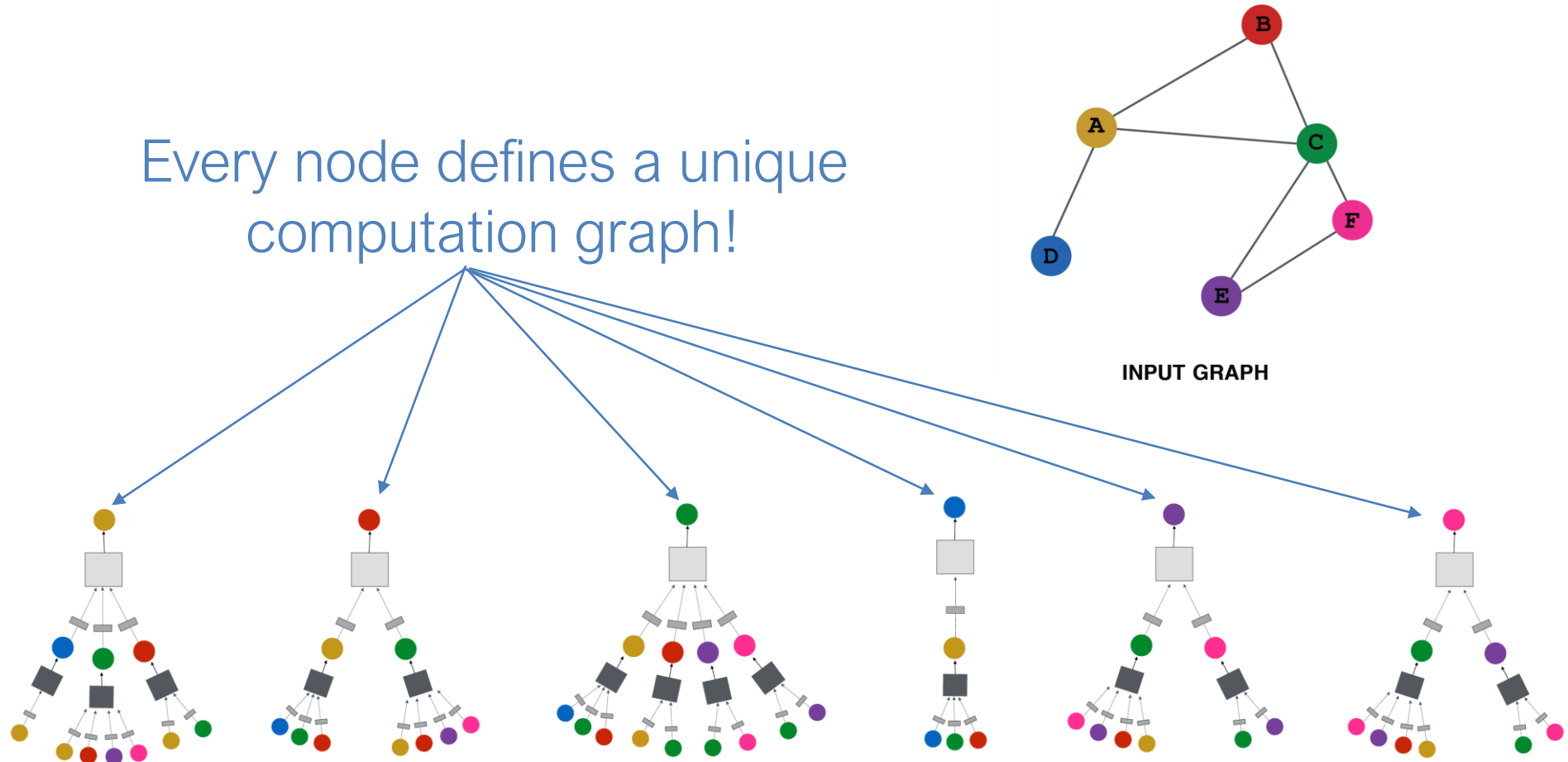
# Graph Neural Nets

**Key idea:** Generate node embeddings based on local neighborhoods in a recursive manner



**INPUT GRAPH**

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]

# Graph Neural Nets

Every node defines a unique computation graph!
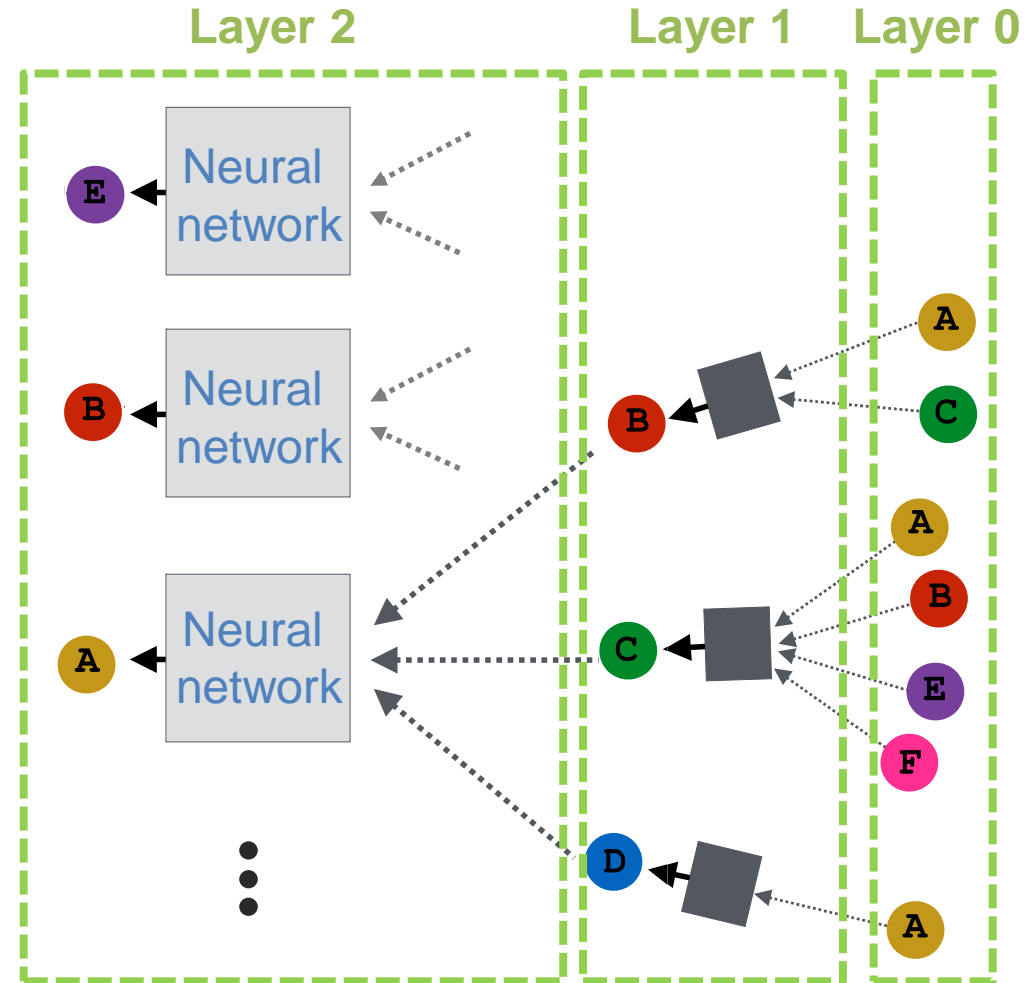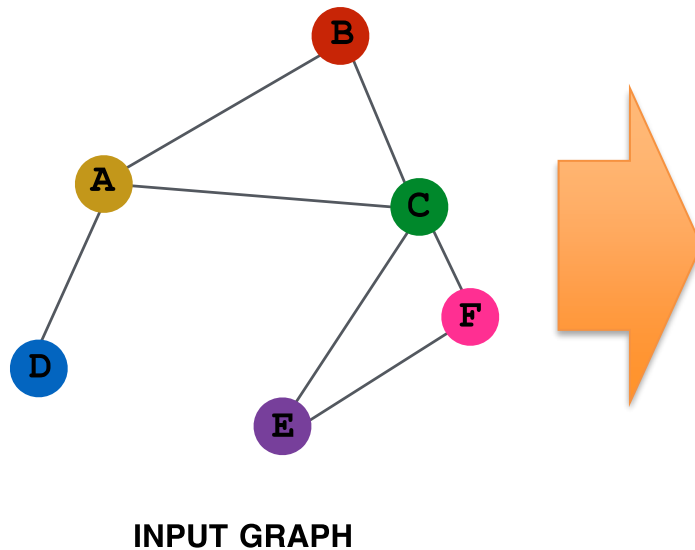
**INPUT GRAPH**

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]

# Graph Neural Nets

## And multiple layers!

➡ Shared parameters within a specific layer

➡ "layer-0" is the input feature $x_u$



INPUT GRAPH

Layer 2        Layer 1    Layer 0

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]
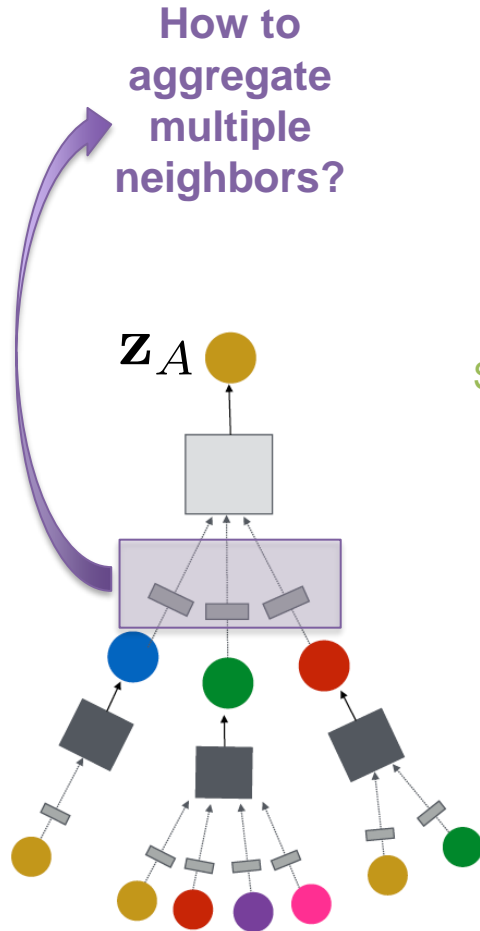
# Graph Neural Nets – Neighborhood Aggregation

**How to aggregate multiple neighbors?**

$\mathbf{z}_A$

### Average pooling (Scarselli et al., 2005)

Different weights for neighbors and self

$$\mathbf{h}_v^k = \sigma\left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1}\right)$$

**K is num layers**

### Graph Convolution Network (Kipf et al., 2017)

Same weights

$$\mathbf{h}_v^k = \sigma\left(\mathbf{W}_k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)||N(v)|}}\right)$$
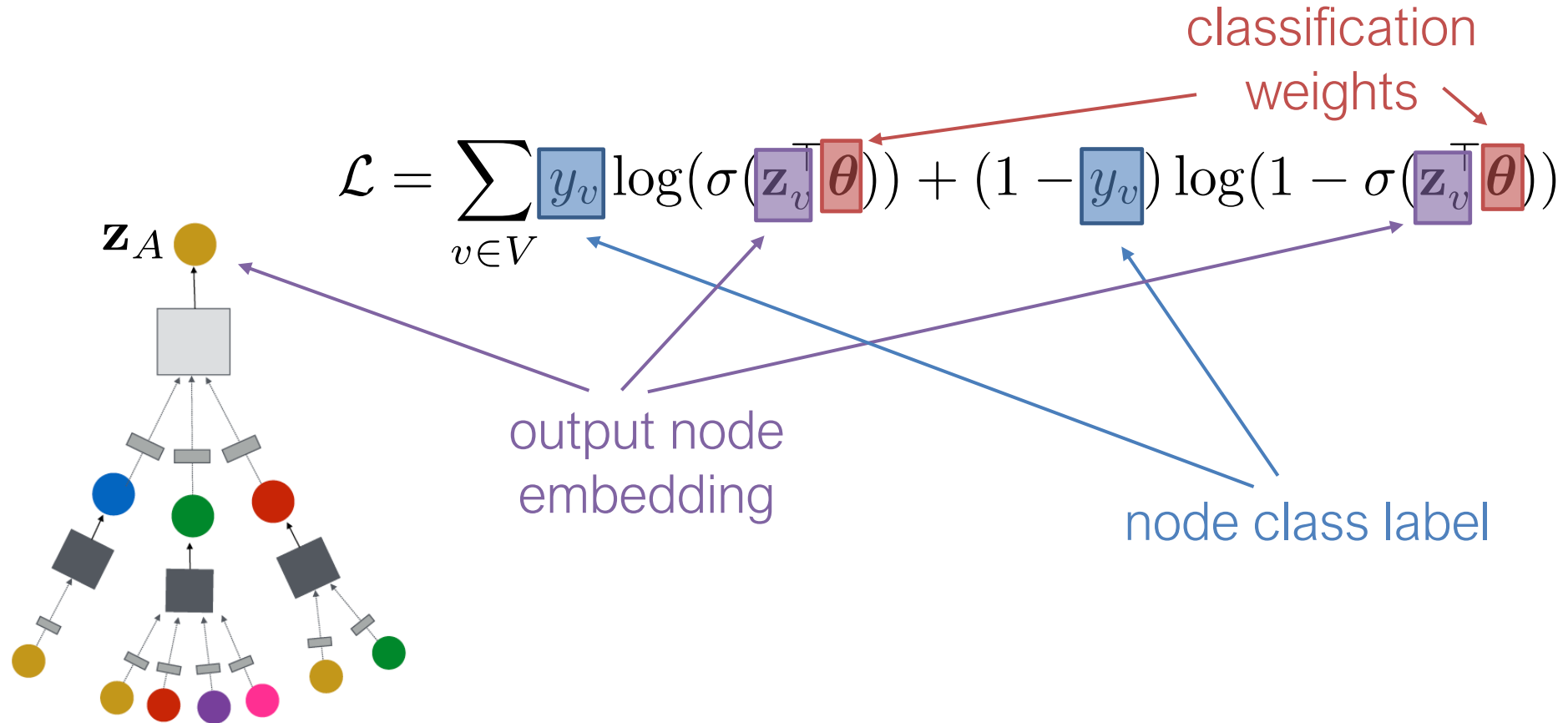
Different normalization

It can be efficiently implemented

### Graph Attention Network (Velickovic et al., 2018)

Attention weights

$$\mathbf{h}_v^k = \sigma\left(\mathbf{W}_k \sum_{u \in N(v) \cup v} \frac{\alpha_{uv} \mathbf{h}_u^{k-1}}{\sqrt{|N(u)||N(v)|}}\right)$$

Very similar to a self-attention transformer

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]

# Graph Neural Nets – Supervised Training

classification weights

$$\mathcal{L} = \sum_{v \in V} y_v \log(\sigma(\mathbf{z}_v^\top \boldsymbol{\theta})) + (1 - y_v) \log(1 - \sigma(\mathbf{z}_v^\top \boldsymbol{\theta}))$$

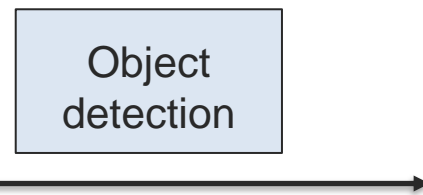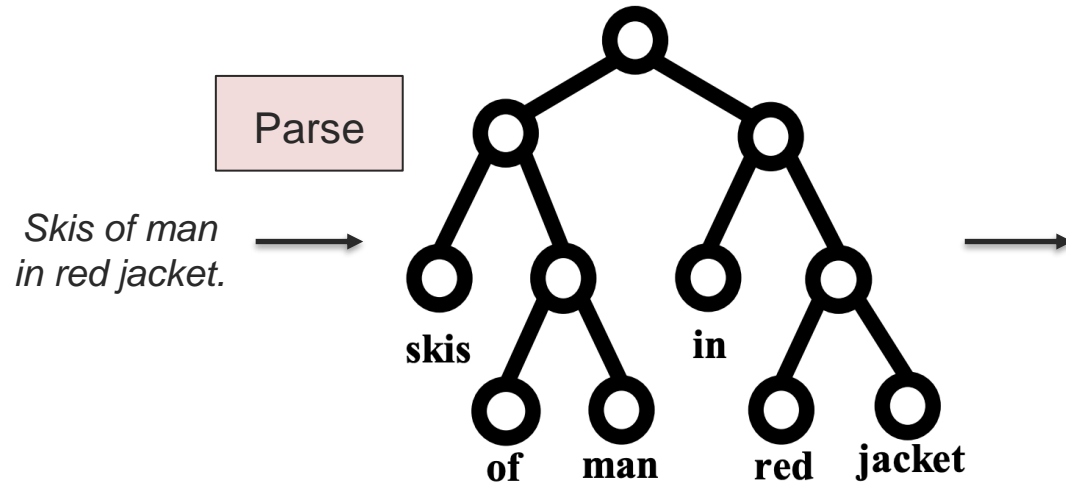$\mathbf{z}_A$

output node embedding

node class label

[Leskovec. Representation Learning on Networks. WWW 2018; Hamilton and Tang, Tutorial on Graph Representation Learning. AAAI 2019]

# Going Beyond Sequences: Hierarchical Structure

# Hierarchical Structure

**Leverage syntactic structure of language**
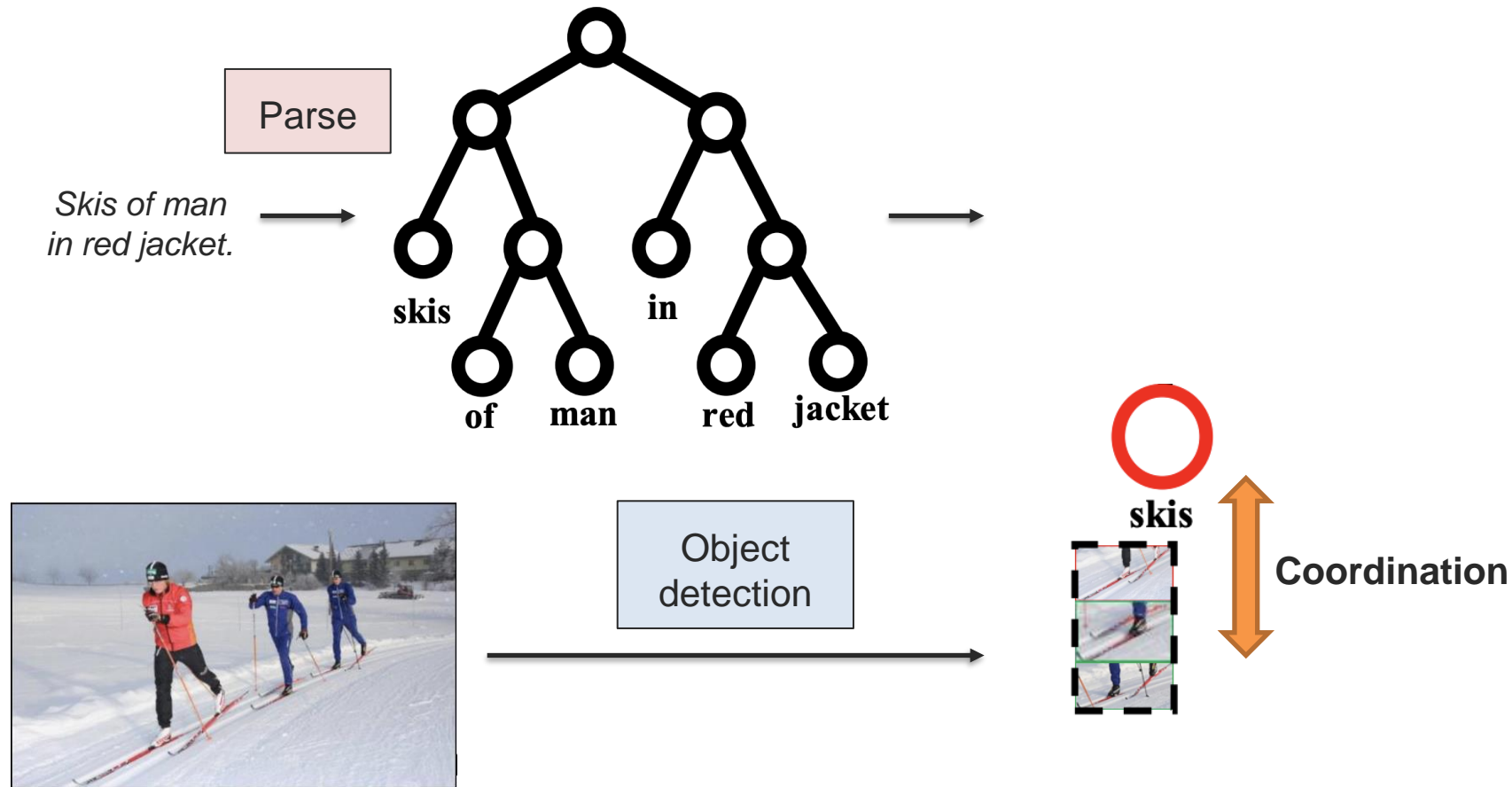


*Skis of man in red jacket.*

Parse

skis    in

of    man    red    jacket

Object detection

[Hong et al., Learning to Compose and Reason with Language Tree Structures for Visual Grounding. IEEE TPAMI 2019]

# Hierarchical Structure

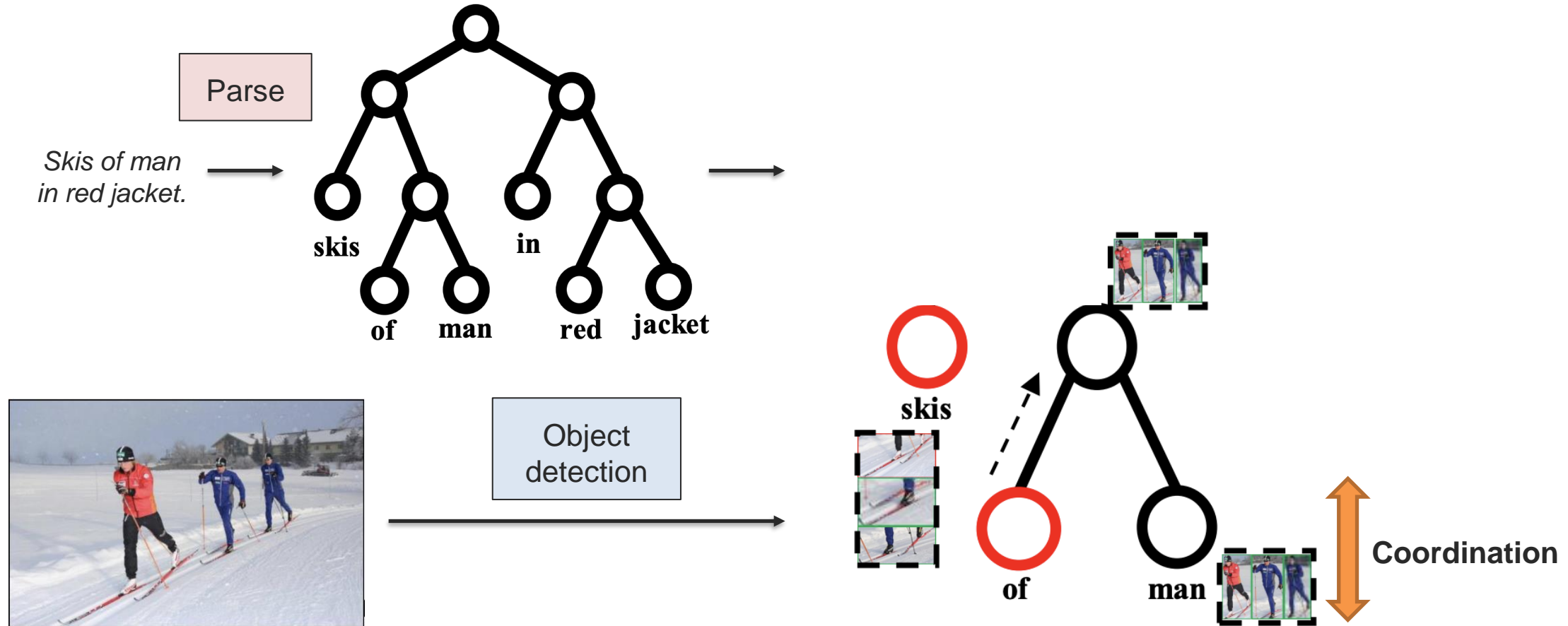**Leverage syntactic structure of language**



Parse

*Skis of man in red jacket.*

skis     in

of   man    red   jacket

Object detection

skis

**Coordination**

[Hong et al., Learning to Compose and Reason with Language Tree Structures for Visual Grounding. IEEE TPAMI 2019]
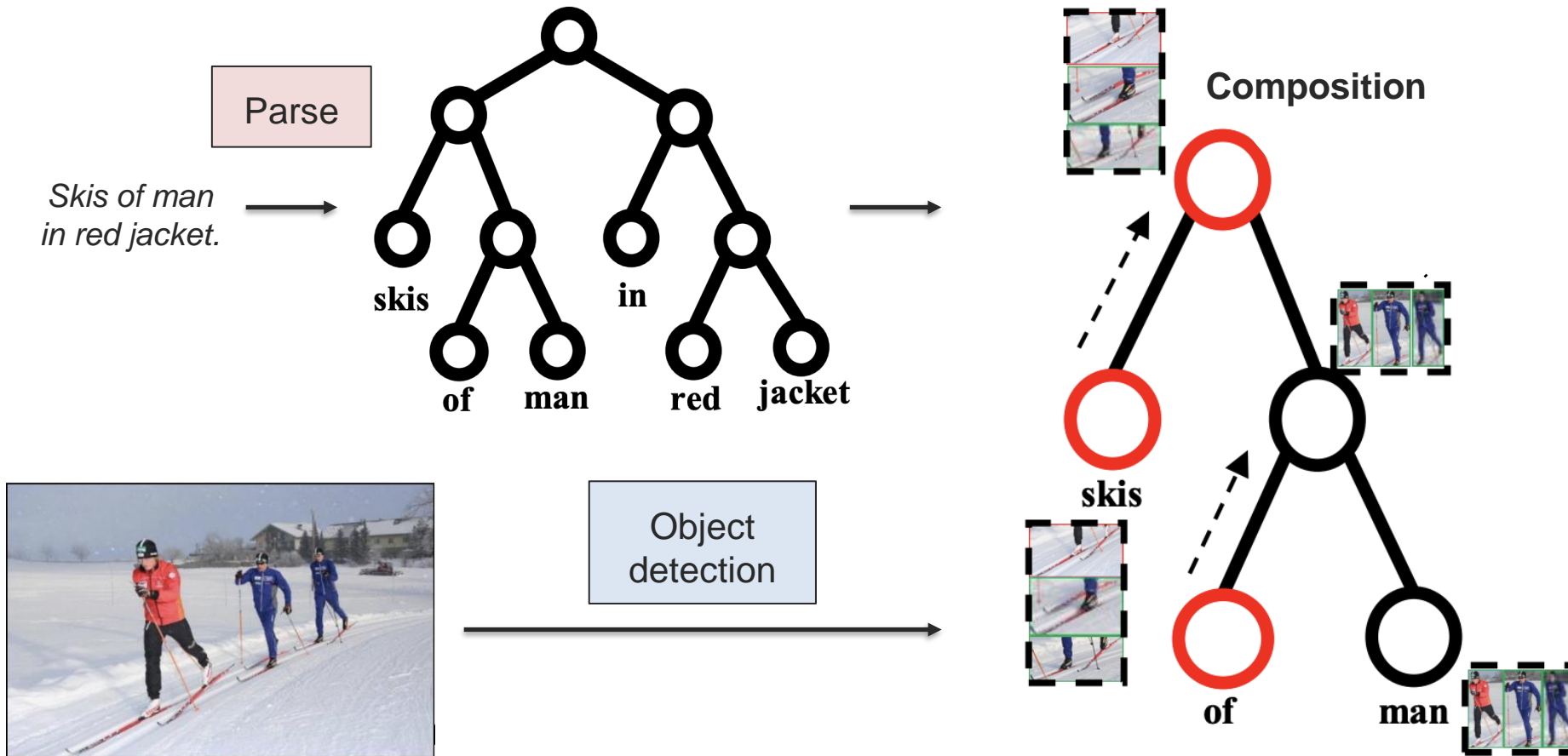
# Hierarchical Structure

**Leverage syntactic structure of language**



[Hong et al., Learning to Compose and Reason with Language Tree Structures for Visual Grounding. IEEE TPAMI 2019]

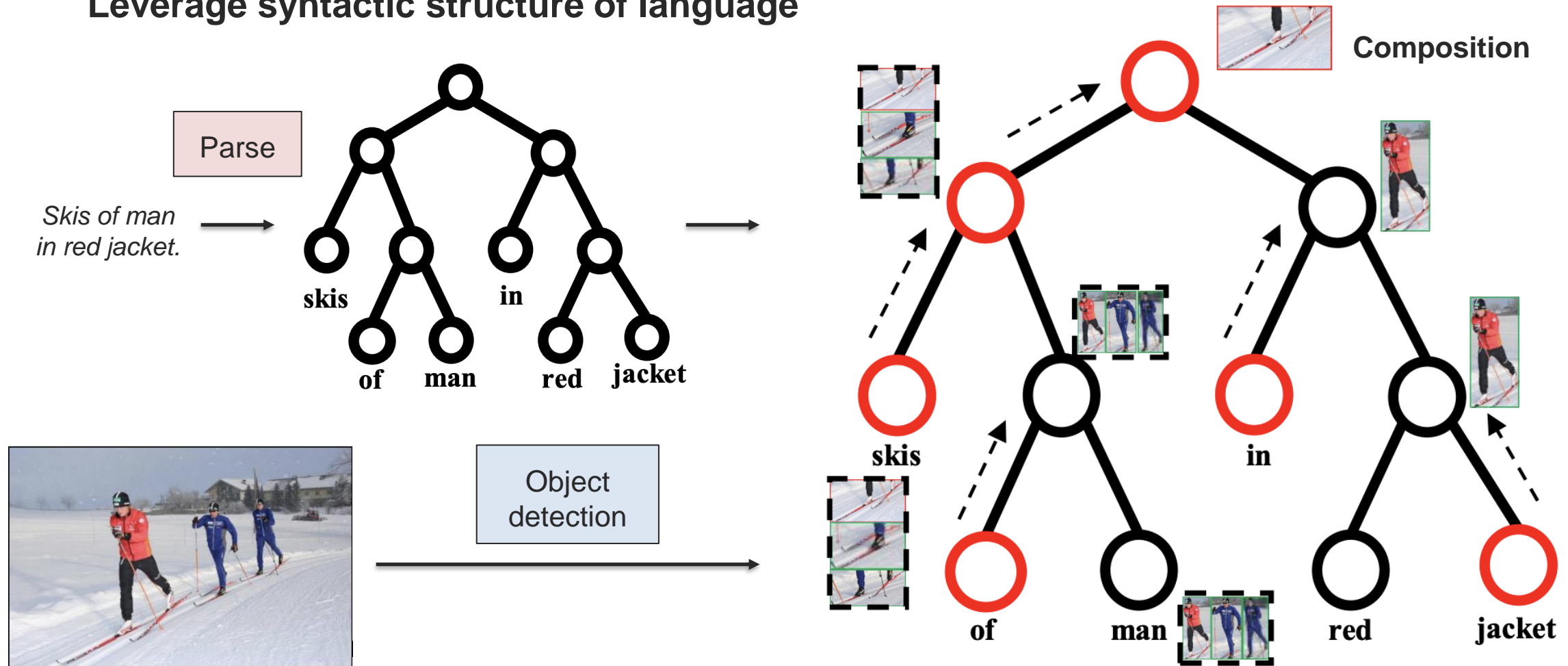# Hierarchical Structure

**Leverage syntactic structure of language**



[Hong et al., Learning to Compose and Reason with Language Tree Structures for Visual Grounding. IEEE TPAMI 2019]

# Hierarchical Structure
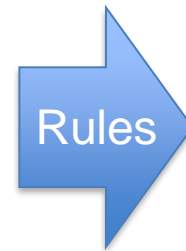
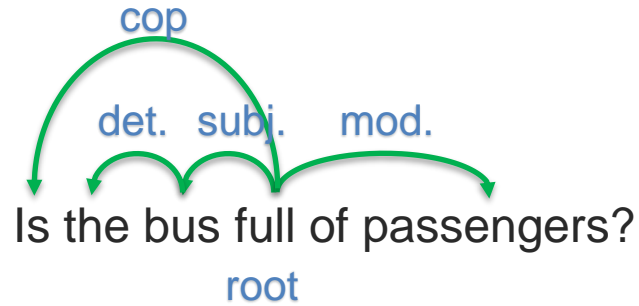**Leverage syntactic structure of language**



[Hong et al., Learning to Compose and Reason with Language Tree Structures for Visual Grounding. IEEE TPAMI 2019]
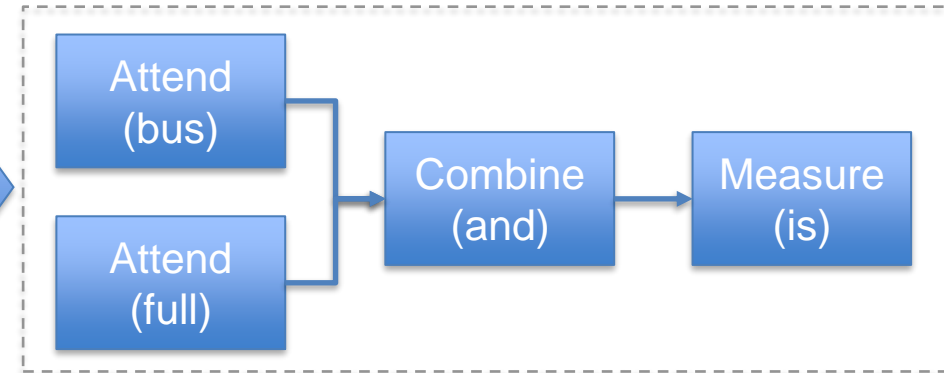
# Going Beyond Sequences: Modular Structure

# Neural Module Network

cop

det. subj. mod.

Is the bus full of passengers?

root

Rules

**Computation layout**

Attend (bus)

Attend (full)

Combine (and)

Measure (is)

Each module work on the attention map(s):
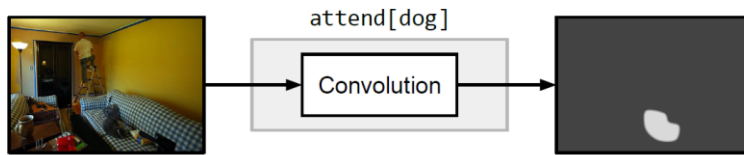
"tie"

Attend (tie)

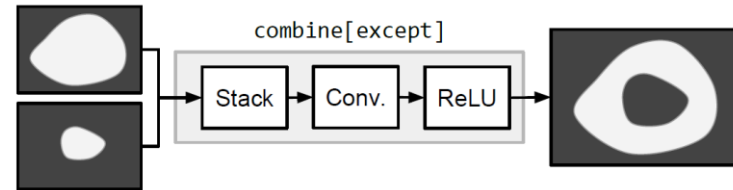Andreas et al., Deep Compositional Question Answering with Neural Module Networks, 2016

# Predefined Set of Modules

**1) Analyze the image:**



attend : $Image \rightarrow Attention$

attend[dog]

Convolution

combine : $Attention \times Attention \rightarrow Attention$

combine[except]

Stack → Conv. → ReLU

**2) Make a prediction**



classify[where]

Attend → FC → Softmax → couch

measure : $Attention \rightarrow Label$

measure[exists]

FC → ReLU → FC → Softmax → yes

Andreas et al., Deep Compositional Question Answering with Neural Module Networks, 2016

# CLEVR: Dataset for Visual Reasoning

**Perfect for a neural module network!**



**Q:** Are there an equal number of large things and metal spheres?
**Q:** What size is the cylinder that is left of the brown metal thing that is left of the big sphere? **Q:** There is a sphere with the same size as the metal cube; is it made of the same material as the small red sphere?
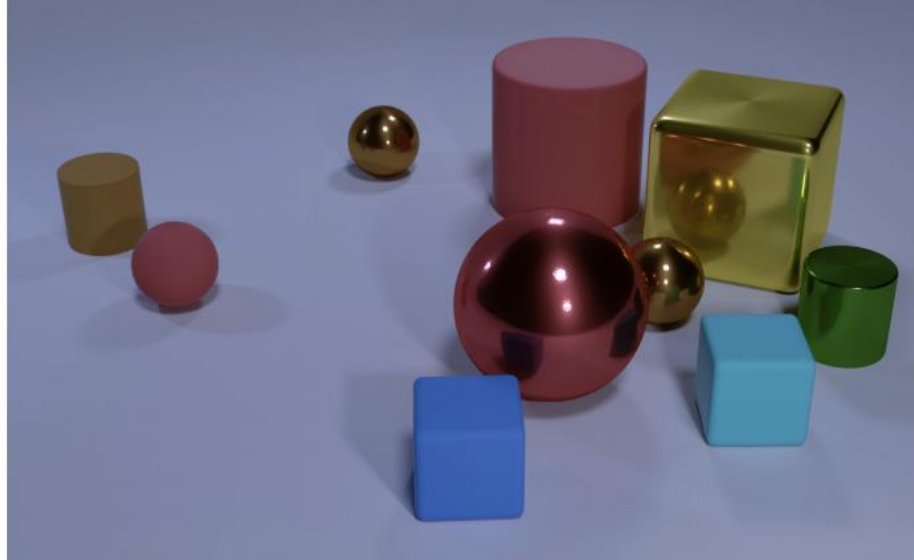**Q:** How many objects are either small cylinders or metal things?

Johnson et al., CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning, CVPR 2017

# Module Network V2: End-to-End Learning

**Computation layout**

Is the bus full of passengers? → RNN →

- Attend (bus)
- Attend (full)
- Combine (and)
- Measure (is)

**No need to parse the question!**

**No rule-based creation of the layout!**
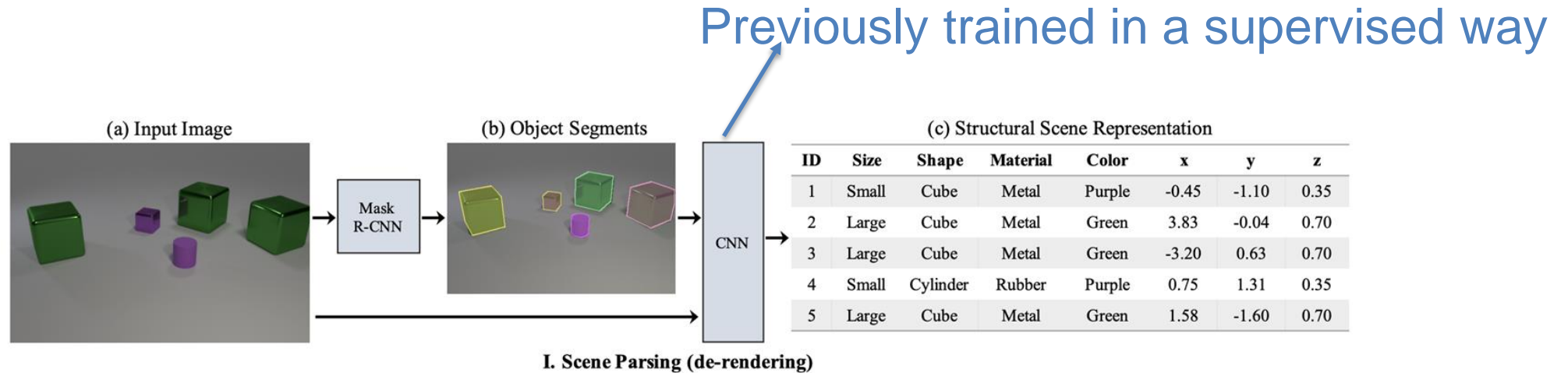
Hu et al., Learning to Reason: End-to-End Module Networks for Visual Question Answering, 2017

## 1) Image Attributes

Previously trained in a supervised way



(a) Input Image → Mask R-CNN → (b) Object Segments → CNN → (c) Structural Scene Representation

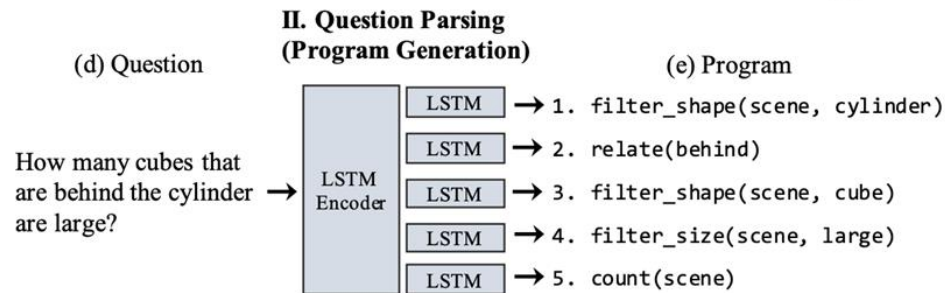| ID | Size | Shape | Material | Color | x | y | z |
|----|------|-------|----------|-------|------|-------|------|
| 1 | Small | Cube | Metal | Purple | -0.45 | -1.10 | 0.35 |
| 2 | Large | Cube | Metal | Green | 3.83 | -0.04 | 0.70 |
| 3 | Large | Cube | Metal | Green | -3.20 | 0.63 | 0.70 |
| 4 | Small | Cylinder | Rubber | Purple | 0.75 | 1.31 | 0.35 |
| 5 | Large | Cube | Metal | Green | 1.58 | -1.60 | 0.70 |

I. Scene Parsing (de-rendering)

Kexin Yi, et al. "Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding." Neurips 2018

## 2) Parsing questions into programs

**Similar to neural module networsk**



(a) Input Image

(b) Object Segments

Mask R-CNN

CNN

(c) Structural Scene Representation

| ID | Size | Shape | Material | Color | x | y | z |
|---|---|---|---|---|---|---|---|
| 1 | Small | Cube | Metal | Purple | -0.45 | -1.10 | 0.35 |
| 2 | Large | Cube | Metal | Green | 3.83 | -0.04 | 0.70 |
| 3 | Large | Cube | Metal | Green | -3.20 | 0.63 | 0.70 |
| 4 | Small | Cylinder | Rubber | Purple | 0.75 | 1.31 | 0.35 |
| 5 | Large | Cube | Metal | Green | 1.58 | -1.60 | 0.70 |

**I. Scene Parsing (de-rendering)**

**II. Question Parsing (Program Generation)**

(d) Question

How many cubes that are behind the cylinder are large?

LSTM Encoder

(e) Program

LSTM → 1. `filter_shape(scene, cylinder)`
LSTM → 2. `relate(behind)`
LSTM → 3. `filter_shape(scene, cube)`
LSTM → 4. `filter_size(scene, large)`
LSTM → 5. `count(scene)`

Kexin Yi, et al. "Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding." Neurips 2018

# 3) Program execution

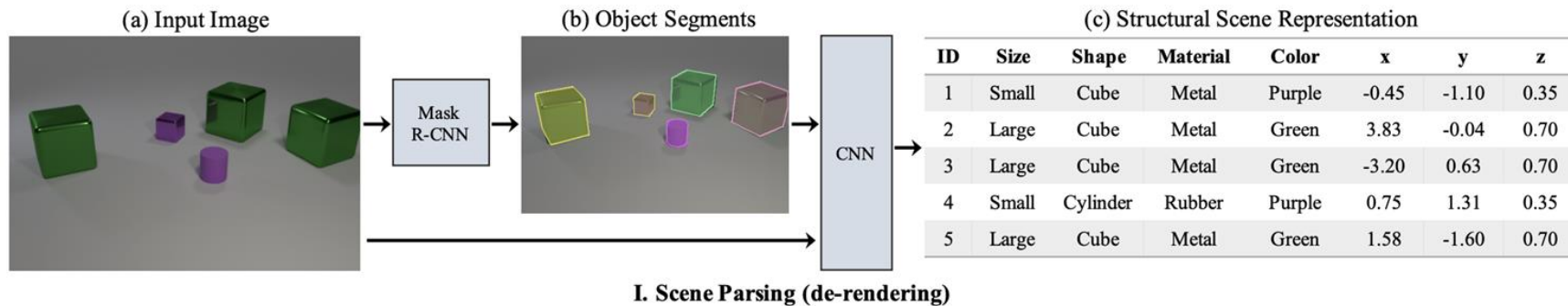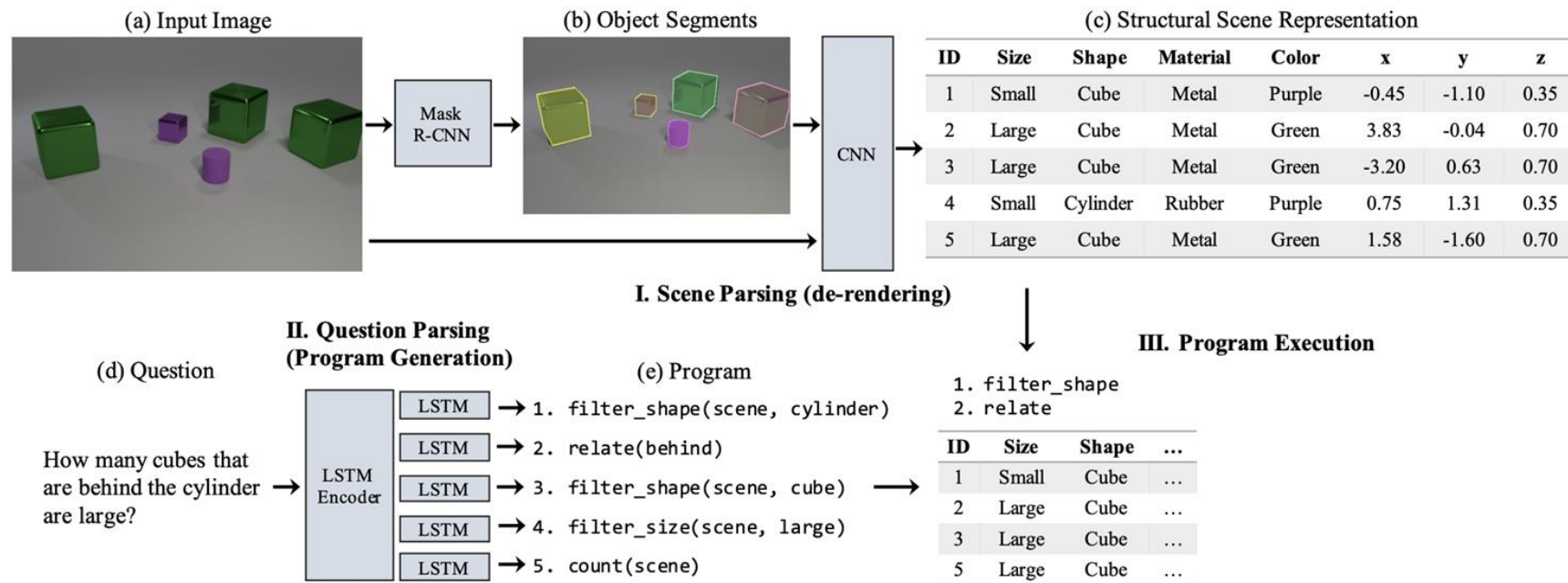**Execution of the program is somewhat easier given the "symbolic" representation of the image**



Kexin Yi, et al. "Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding." Neurips 2018

# Module Networks V4: The Neural State Machine



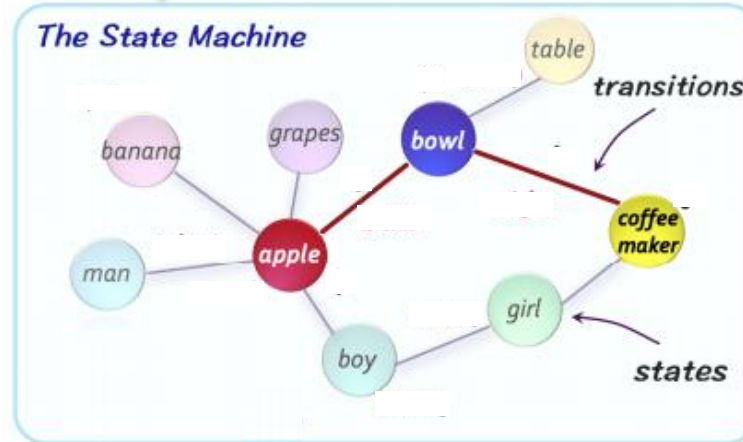How to solve this question using visual reasoning?

What is the red fruit inside the bowl to the right of the coffee maker?

1. Given an **image**, generate a probabilistic **scene graph** that captures the semantic concepts.

2. Treat the graph as a **state machine** and simulate iterative computation over it to *answer questions* or *draw inferences*.

3. Natural language questions are translated into *soft instructions* and used to perform sequential reasoning over the scene graph/state machine.

Hudson, Drew, and Christopher D. Manning. "Learning by abstraction: The neural state machine." NeurIPS 2019

# Module Networks V4: The Neural State Machine

Detect objects and create proximity graph



What is the red fruit inside the bowl to the right of the coffee maker?

Hudson, Drew, and Christopher D. Manning. "Learning by abstraction: The neural state machine." NeurIPS 2019

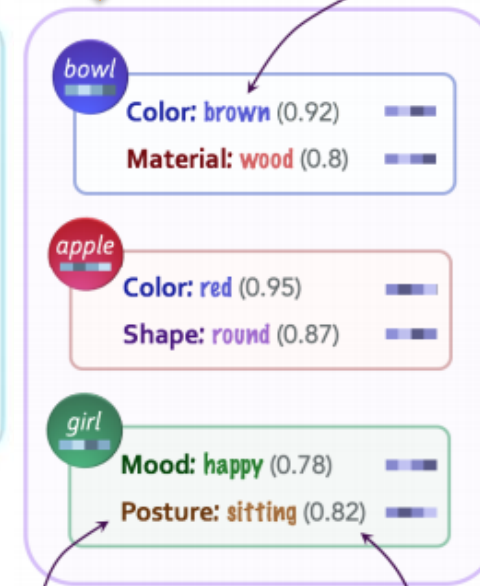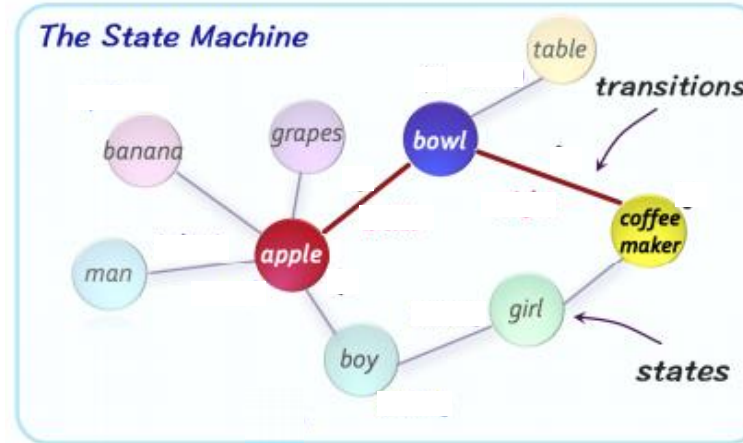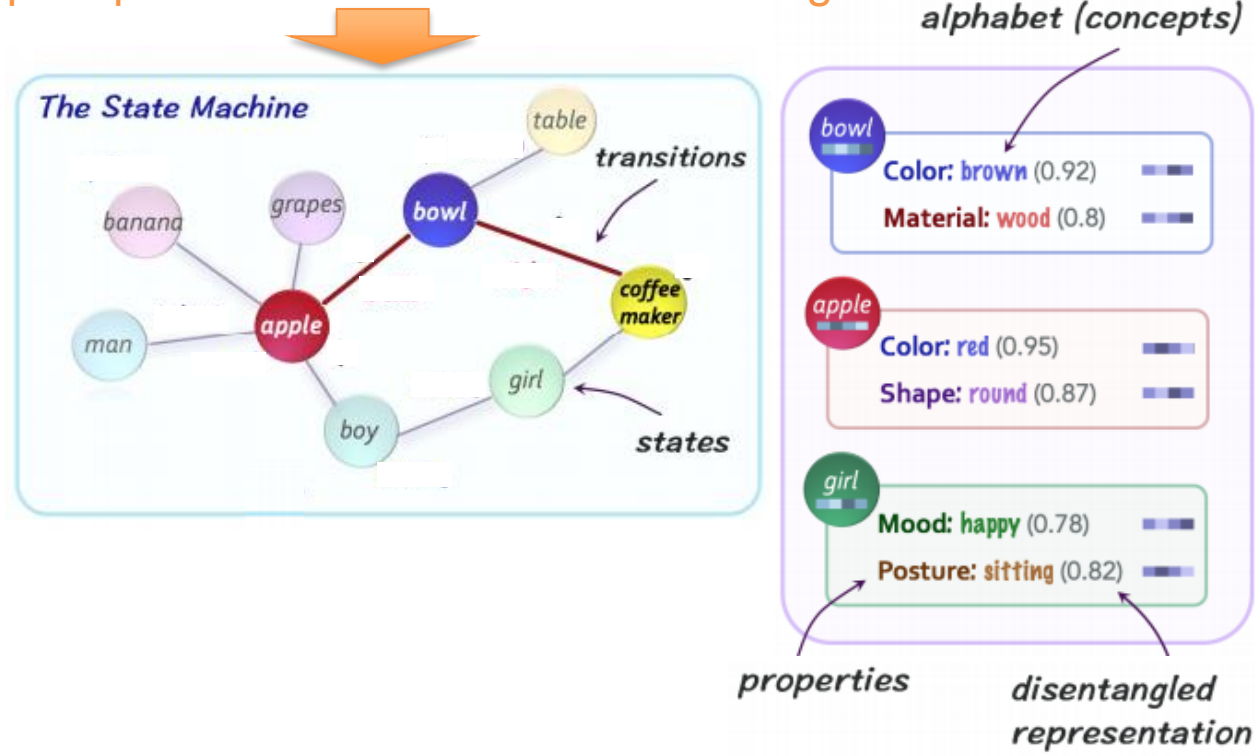# Module Networks V4: The Neural State Machine



Pre-trained an alphabet of concepts
(Visual Genome)

Manually grouped
by "properties"

Probabilities
computed at
runtime for each
object instance

Hudson, Drew, and Christopher D. Manning. "Learning by abstraction: The neural state machine." NeurIPS 2019

# Module Networks V4: The Neural State Machine

Predefined an alphabet of relations
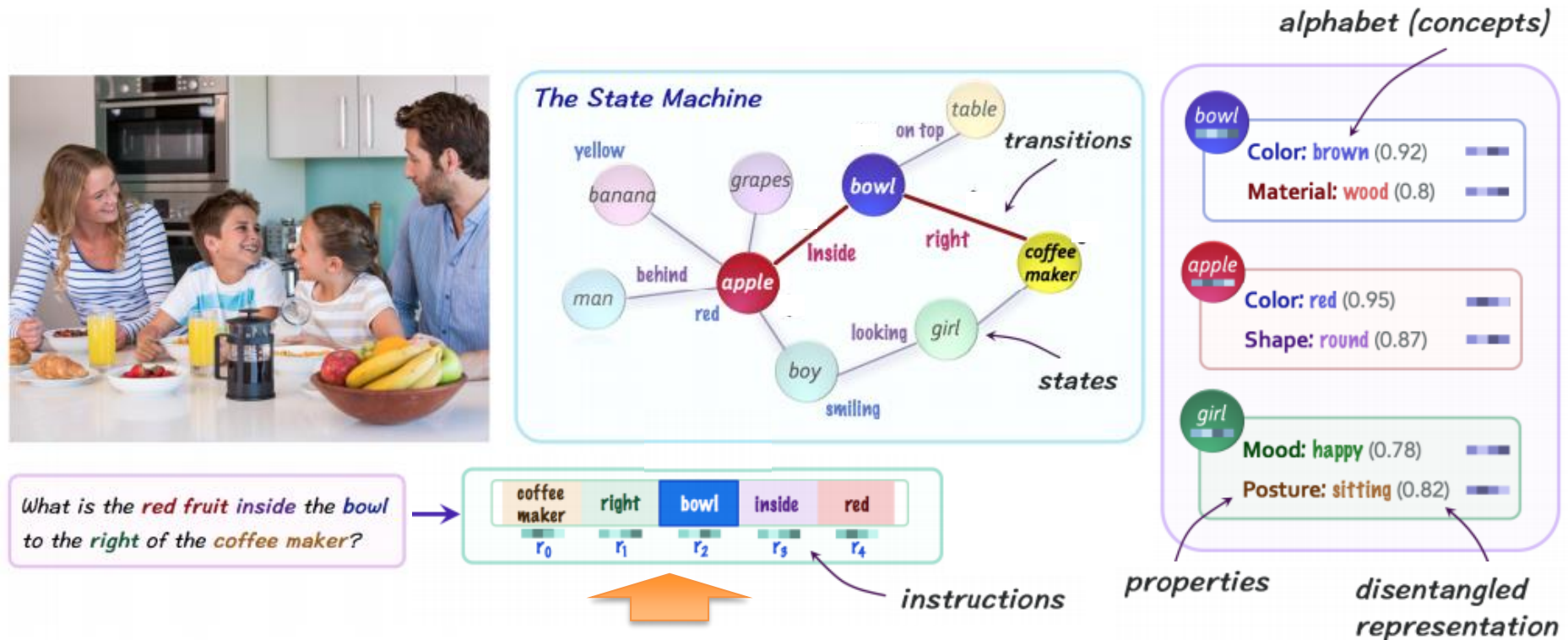and compute probabilities for each directed edges



Hudson, Drew, and Christopher D. Manning. "Learning by abstraction: The neural state machine." NeurIPS 2019

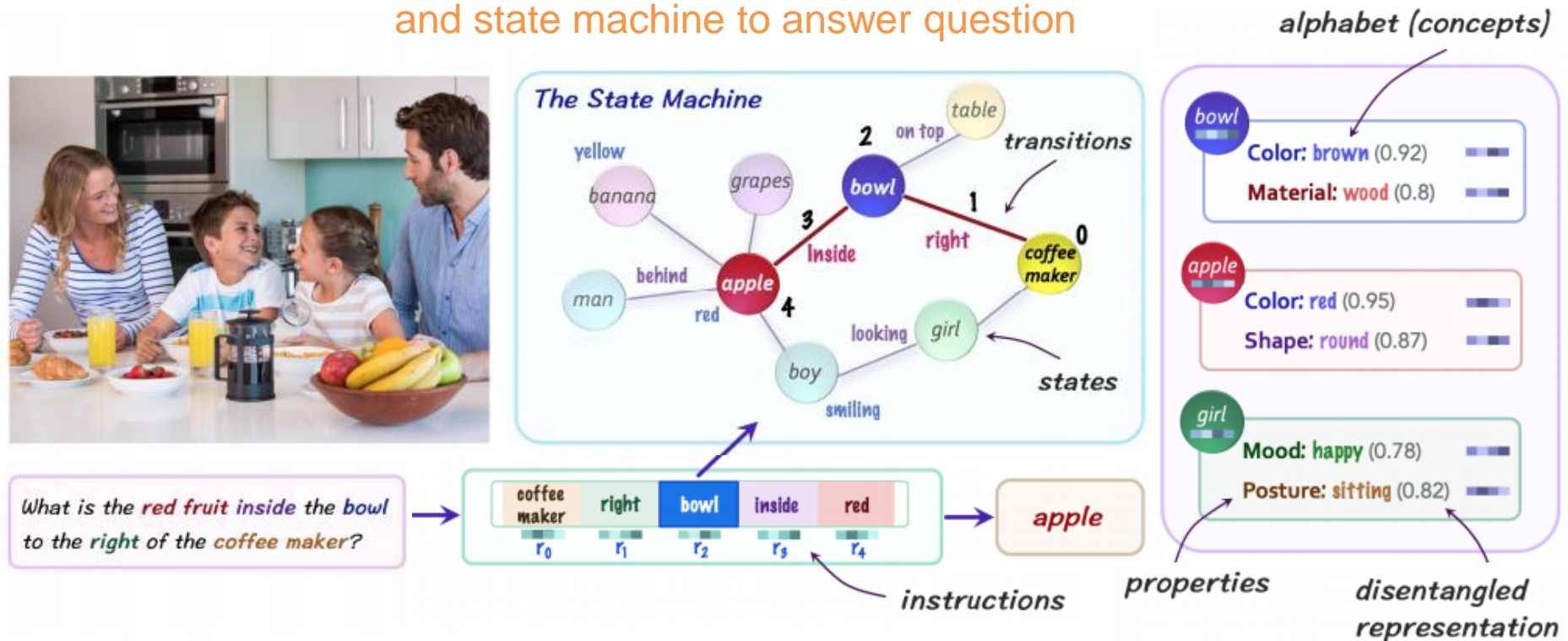# Module Networks V4: The Neural State Machine



Translate each word in a concept-based representation
and group in a fixed number of instruction steps

Hudson, Drew, and Christopher D. Manning. "Learning by abstraction: The neural state machine." NeurIPS 2019

# Module Networks V4: The Neural State Machine

Finally, perform reasoning using instructions
and state machine to answer question



Hudson, Drew, and Christopher D. Manning. "Learning by abstraction: The neural state machine." NeurIPS 2019