

Visual Question Answering and Visual Reasoning

Zhe Gan

6/15/2020



Overview

- Goal of this part of the tutorial:
 - Use VQA and visual reasoning as example tasks to understand Vision-and-Language representation learning
 - After the talk, everyone can confidently say: “yeah, I know VQA and visual reasoning pretty well now”
 - Focus on high-level intuitions, not technical details
 - Focus on static images, instead of videos
 - Focus on a selective set of papers, not a comprehensive literature review

Agenda

- Task Overview
 - *What are the main tasks that are driving progress in VQA and visual reasoning?*
- Method Overview
 - *What are the state-of-the-art approaches and the key model design principles underlying these methods?*
- Summary
 - *What are the core challenges and future directions?*

Agenda

- Task Overview

- *What are the main tasks that are driving progress in VQA and visual reasoning?*

- Method Overview

- *What are the state-of-the-art approaches and the key model design principles underlying these methods?*

- Summary

- *What are the core challenges and future directions?*


What is V+L about?

- V+L research is about how to train a smart AI system that can see and talk

AI Systems That Can See And Talk

Prof. Devi Parikh / Georgia Tech and Facebook AI Research

[Abstract & Bio](#)



AI Systems That Can See And Talk

Devi Parikh

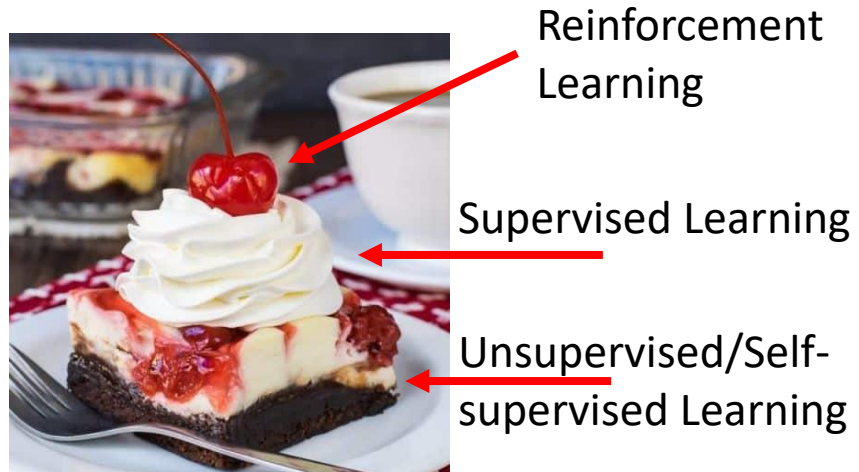
Georgia Tech

facebook Artificial Intelligence Research

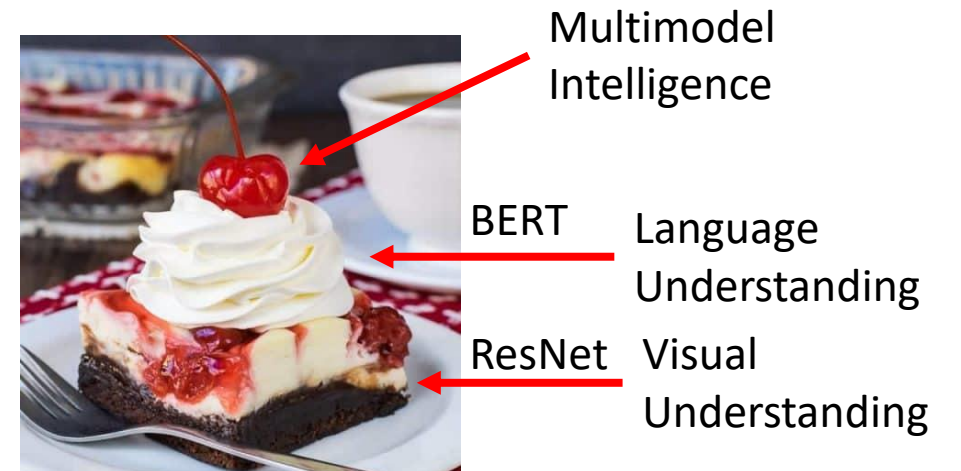
What is V+L about?

- V+L research is about how to train a smart AI system that can see and talk

Prof. Yann LeCun's cake theory

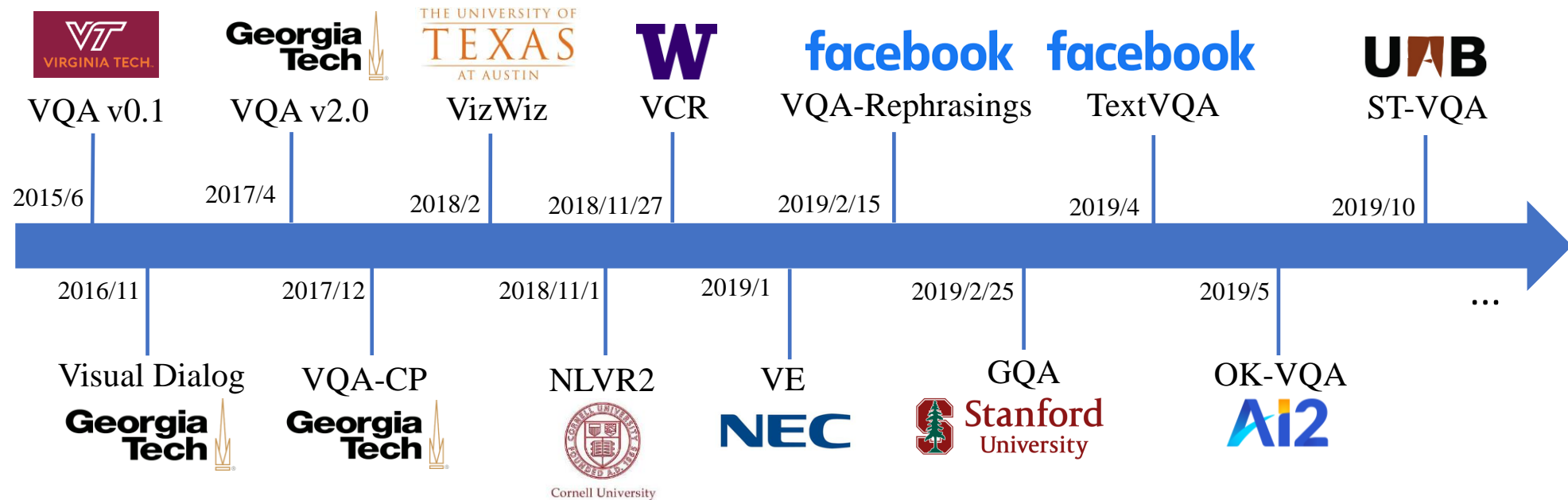


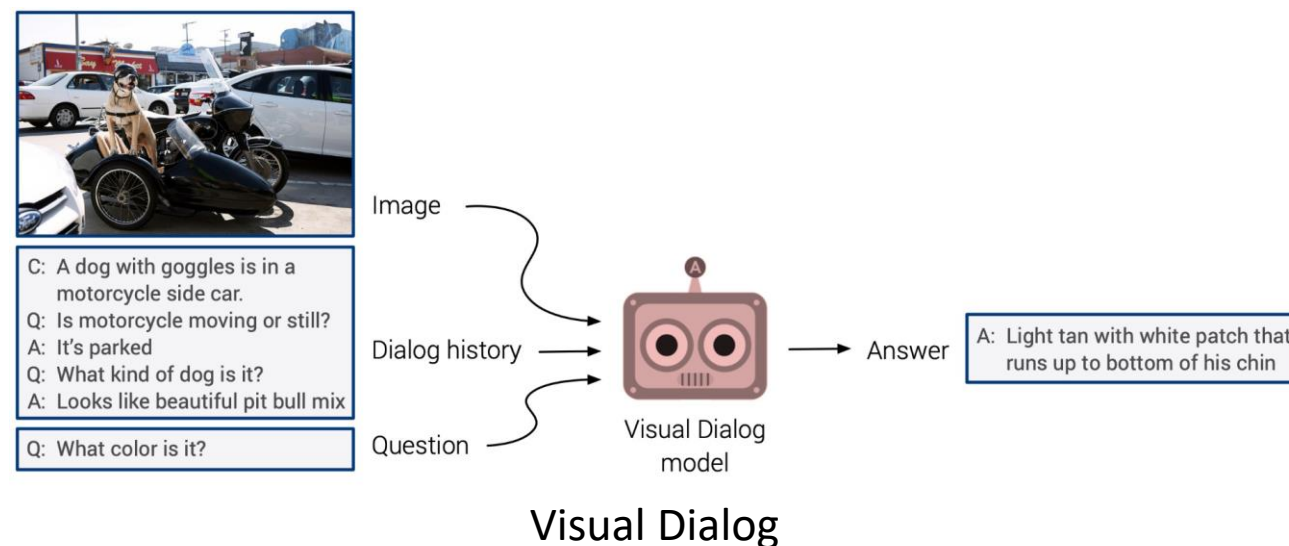
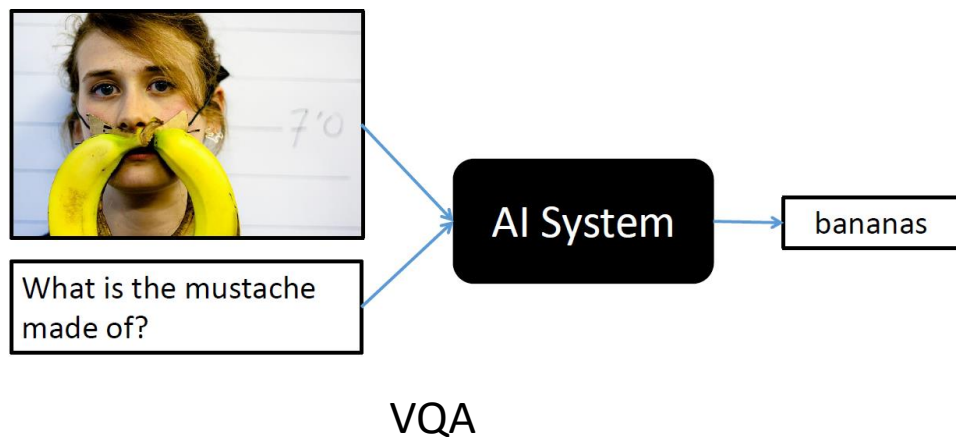
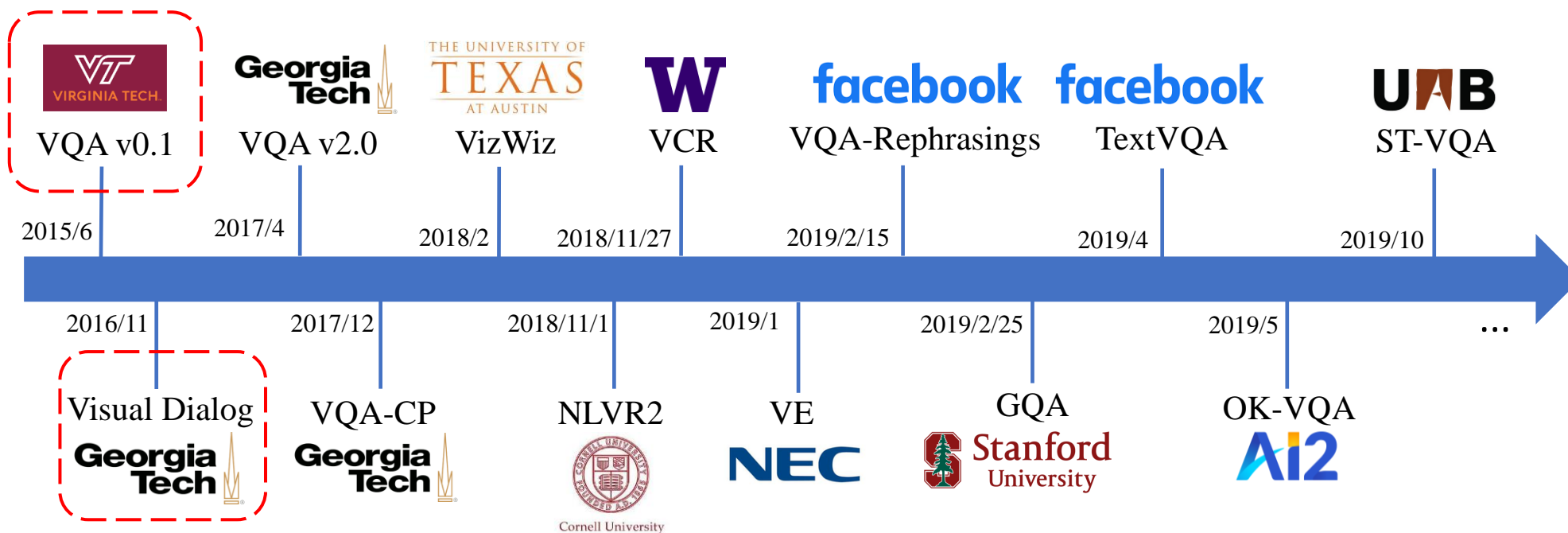
In our V+L context

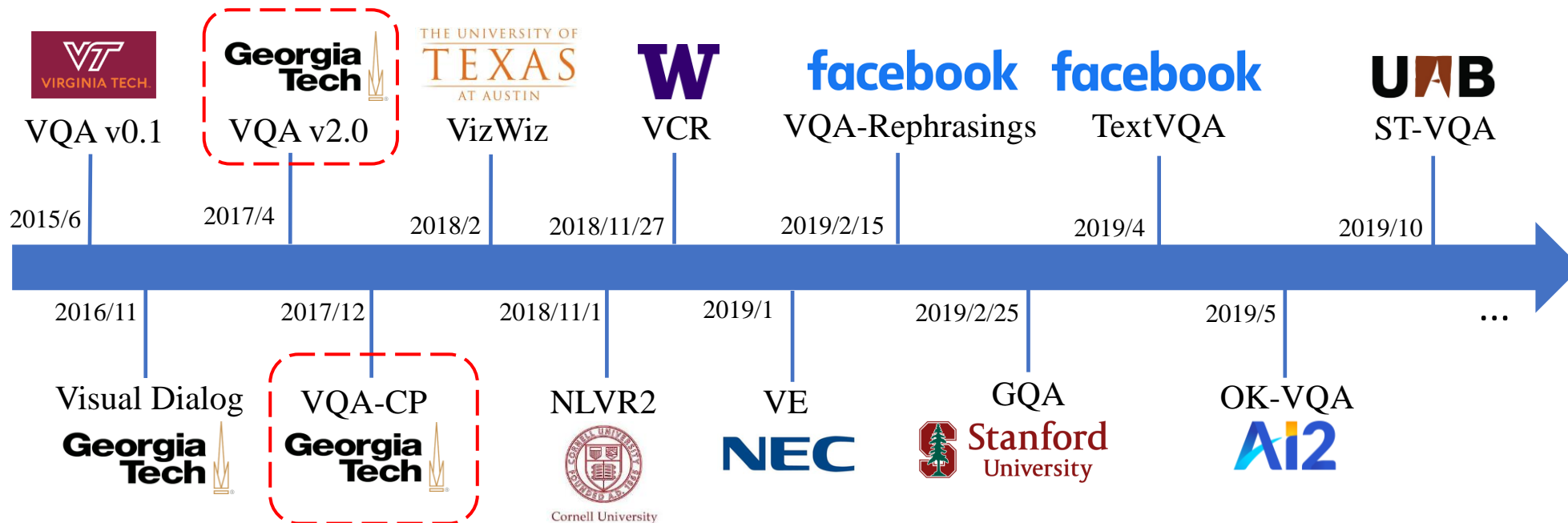


Task Overview: VQA and Visual Reasoning

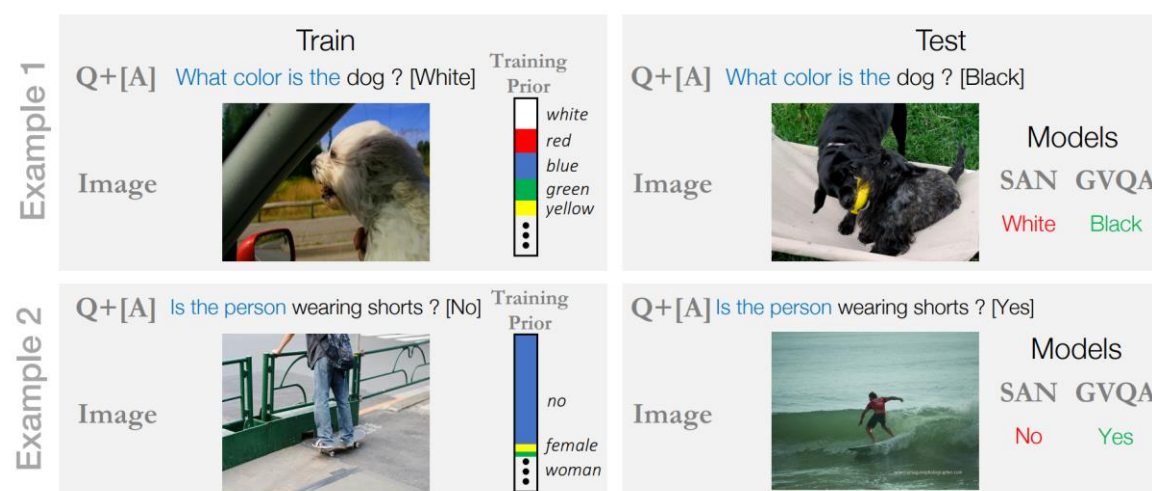
- Large-scale annotated datasets have driven tremendous progress in this field







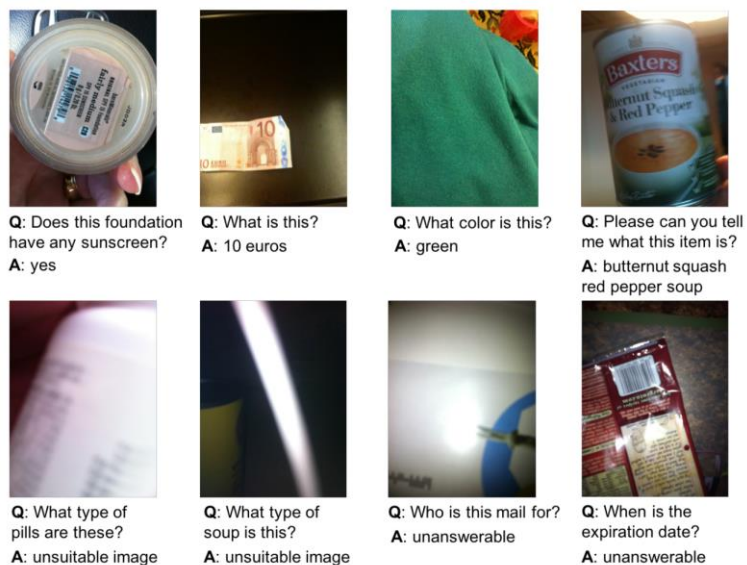
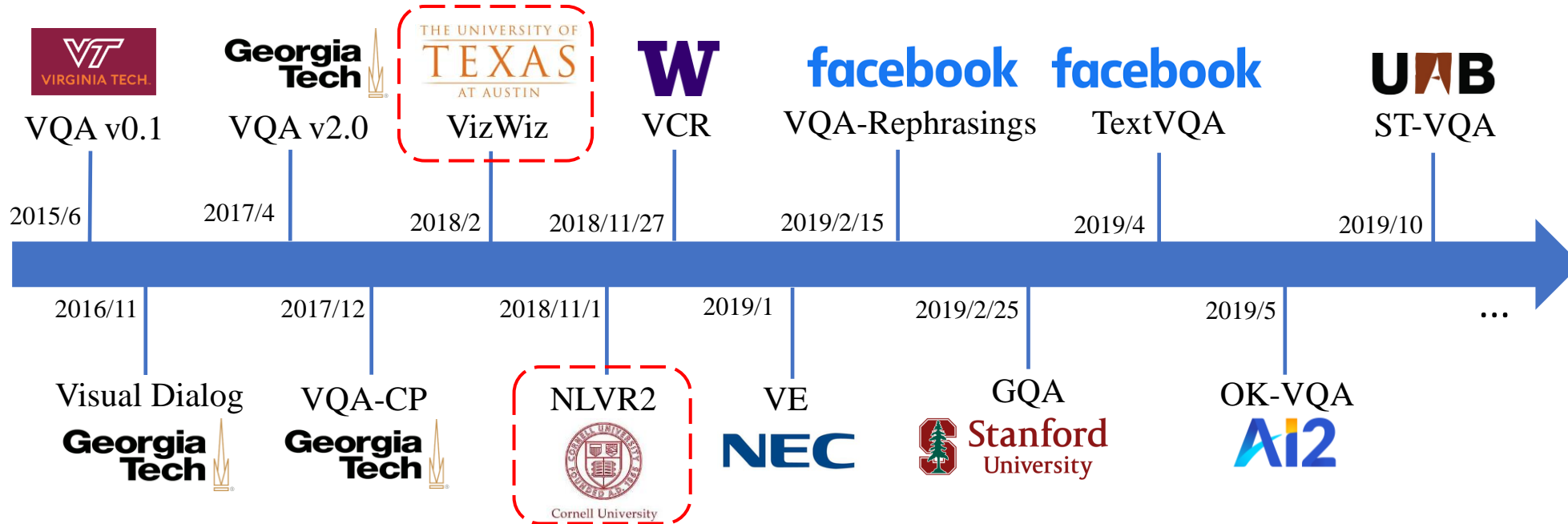
VQA v2.0



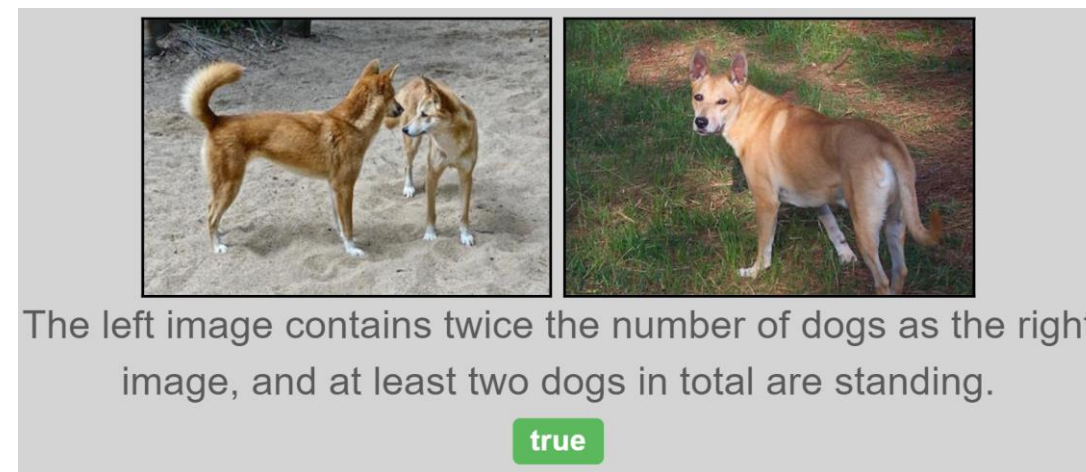
VQA-CP

[1] Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering, CVPR 2017

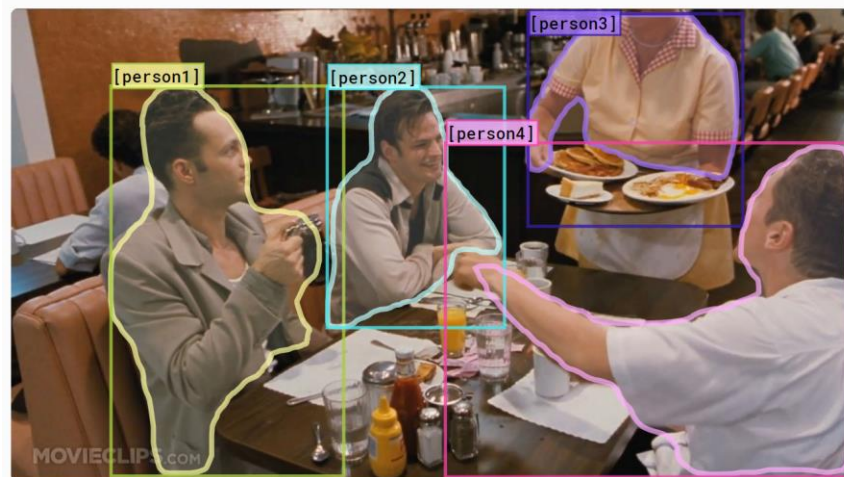
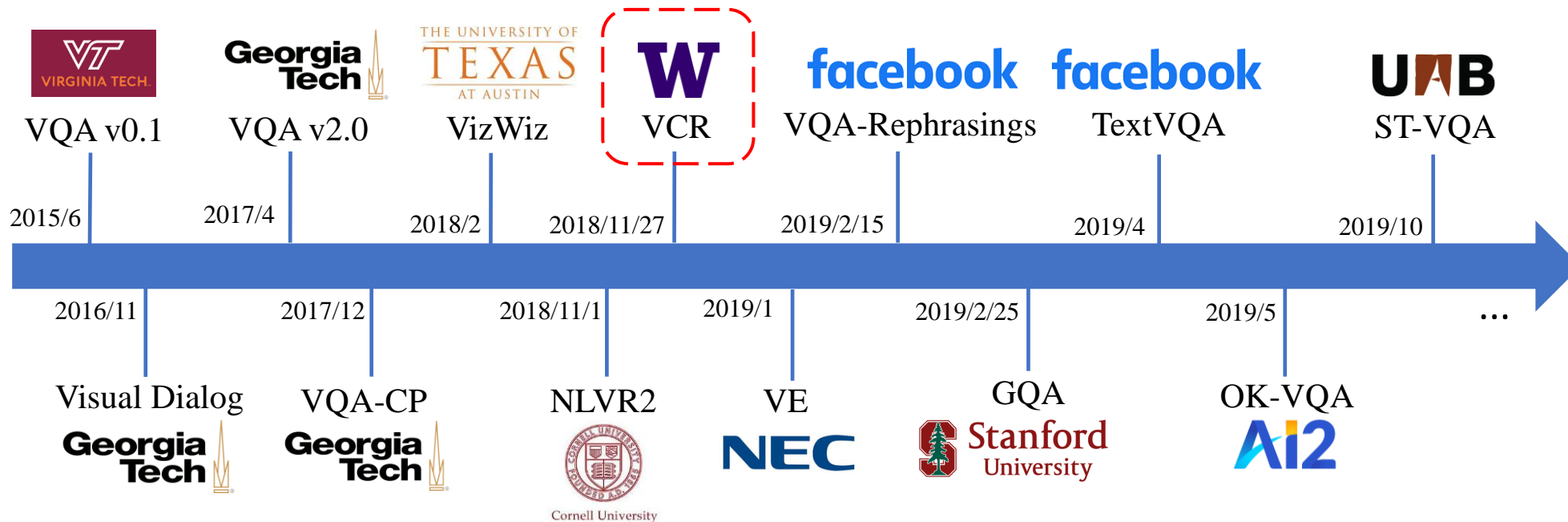
[2] Don't Just Assume; Look and Answer: Overcoming Priors for Visual Question Answering, CVPR 2018



VizWiz



NLVR2



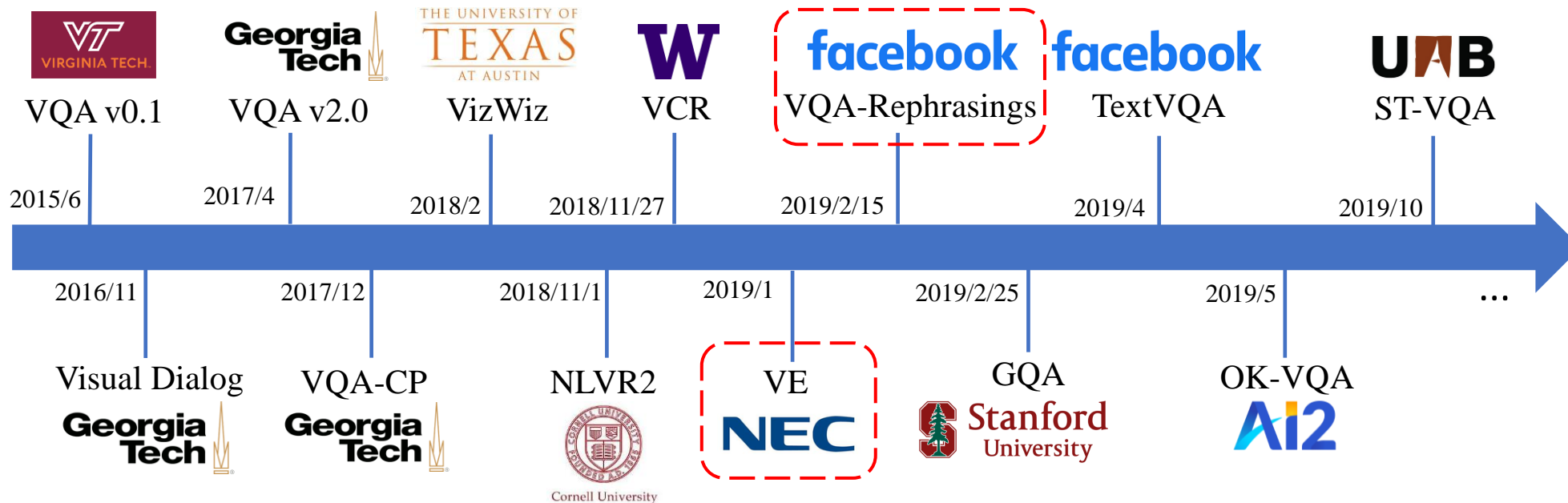
Why is [person4] pointing at [person1]?

- a) He is telling [person3] that [person1] ordered the pancakes.
- b) He just told a joke.
- c) He is feeling accusatory towards [person1].
- d) He is giving [person1] directions.

Rationale: I think so because...

- a) [person1] has the pancakes in front of him.
- b) [person4] is taking everyone's order and asked for clarification.
- c) [person3] is looking at the pancakes both she and [person2] are smiling slightly.
- d) [person3] is delivering food to the table, and she might not know whose order is whose.





Premise

+

- Two woman are holding packages.
- The sisters are hugging goodbye while holding to go packages after just eating lunch.
- The men are fighting outside a deli.

Hypothesis

=

- Entailment
- Neutral
- Contradiction

Answer

Visual Entailment

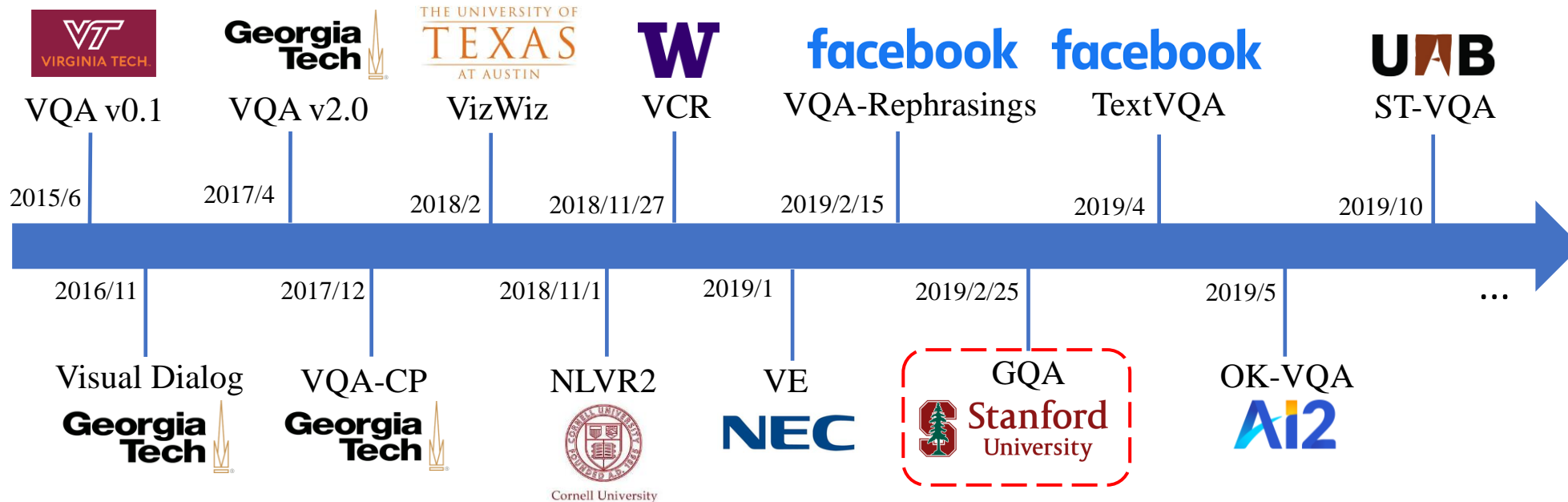


	Prediction
What is in the basket?	banana
What is contained in the basket?	pizza
What can be seen inside the basket?	remote
What does the basket mainly contain?	paper
Is it safe to turn left?	Yes
Can one safely turn left?	No
Would it be safe to turn left?	No
Would turning left considered safe in this picture?	Yes

VQA-Rephrasings

[1] Visual Entailment: A Novel Task for Fine-Grained Image Understanding, 2019

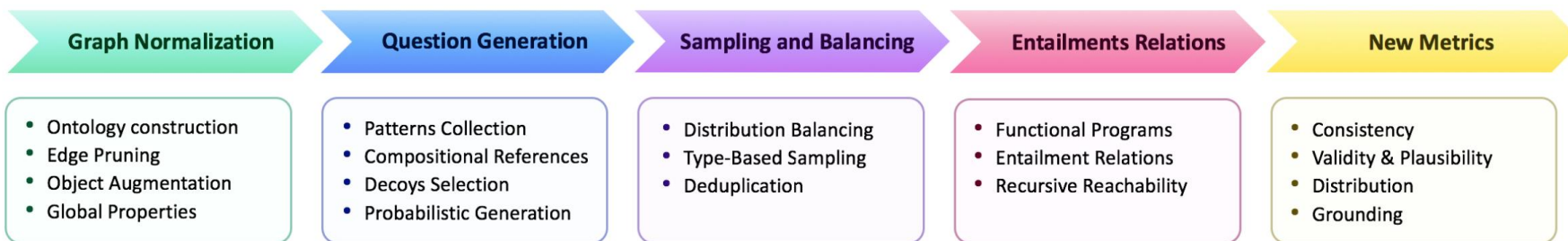
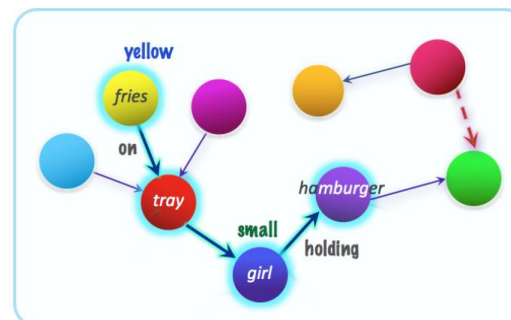
[2] Cycle-Consistency for Robust Visual Question Answering, CVPR 2019

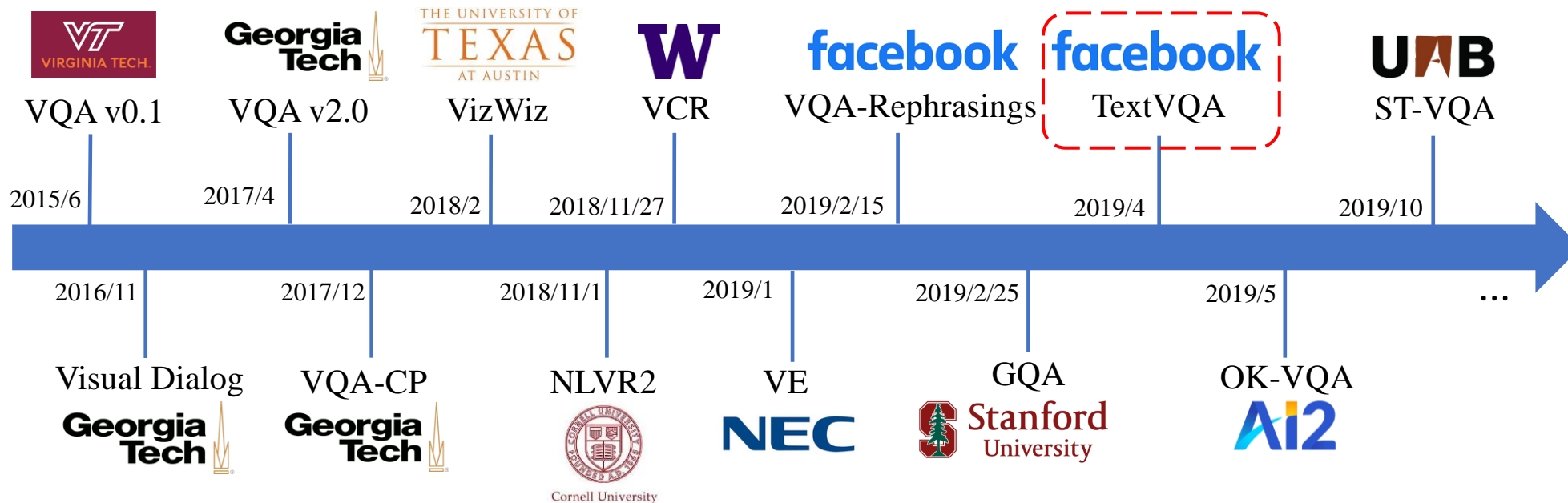


Pattern: What/Which <type> [do you think] <is> <dobject>, <attr> or <decoy>?
Program: Select: <dobject> → Choose <type>: <attr>|<decoy>
Reference: The food on the red object left of the small girl that is holding a hamburger
Decoy: brown

What color is the food on the red object left of the small girl that is holding a hamburger, yellow or brown?

Select: hamburger → Relate: girl1, holding → Filter size: small → Relate: object, left → Filter color: red → Relate: food, on → Choose color: yellow | brown





What is the top oz?

Ground Truth

16

Prediction

red



What is the largest denomination on table?

Ground Truth

500

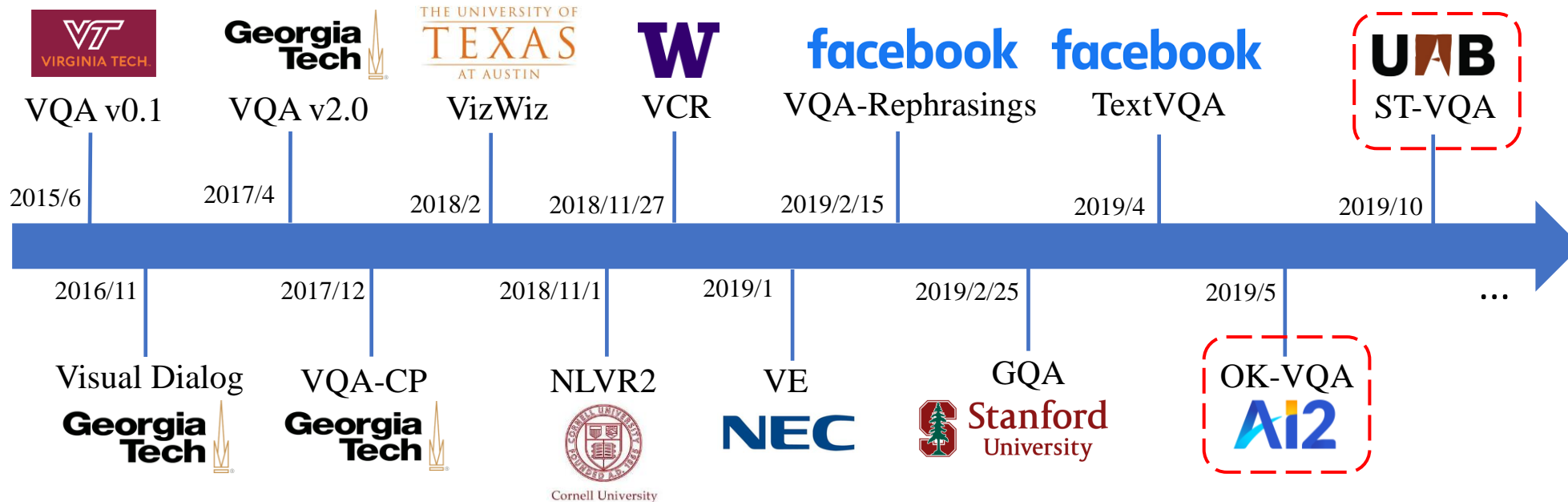
Prediction

unknown



TextVQA

A dataset to benchmark visual reasoning based on text in images.



Q: Which American president is associated with the stuffed animal seen here?

A: Teddy Roosevelt

Outside Knowledge

Another lasting, popular legacy of Roosevelt is the stuffed toy bears—teddy bears—named after him following an incident on a hunting trip in Mississippi in 1902.

Developed apparently simultaneously by toymakers ... and named after President Theodore "Teddy" Roosevelt, the teddy bear became an iconic children's toy, celebrated in story, song, and film.

At the same time in the USA, Morris Michtom created the first teddy bear, after being inspired by a drawing of Theodore "Teddy" Roosevelt with a bear cub.

OK-VQA



Q: What is the price of the bananas per kg?

A: \$11.98



Q: What does the red sign say?

A: Stop

Scene Text VQA

[1] OK-VQA: A Visual Question Answering Benchmark Requiring External Knowledge, CVPR 2019

[2] Scene Text Visual Question Answering, ICCV 2019

More datasets...

SQuINTing at VQA Models: Interrogating VQA Models with Sub-Questions

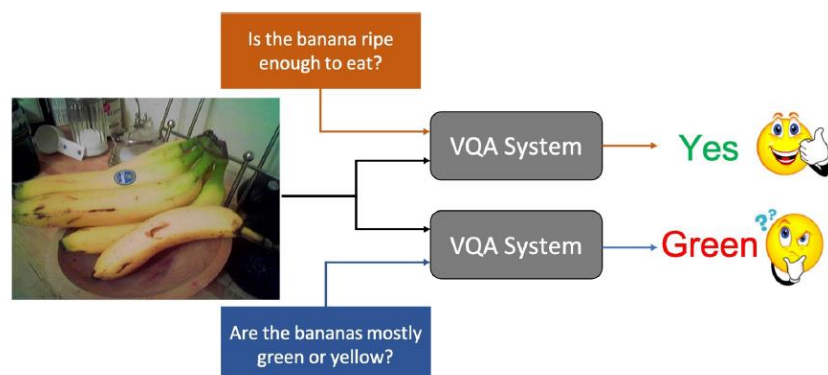


Figure 1: A potential reasoning failure: Current models answer “Yes” correctly to the Reasoning question “Is the banana ripe enough to eat?”. We might assume that correctly answering the Reasoning question stems from perceiving relevant concepts correctly – perceiving yellow bananas in this example. But when asked “Are the bananas mostly green or yellow?”, it answers “Green” incorrectly – indicating that the model possibly answered the original for the wrong reasons even if the answer was right. We quantify the extent to which this phenomenon occurs in VQA and introduce a new dataset aimed at stimulating research on well grounded reasoning.

VQA-LOL: Visual Question Answering under the Lens of Logic

Question	Pred. Answer	LXMERT accuracy
Q_1 : Is there beer?	YES (96.26 %) NO (3.74 %)	Original questions 86.65
Q_2 : Is the man wearing shoes?	NO (90.03 %) YES (9.97 %)	
$\neg Q_2$: Is the man <i>not</i> wearing shoes?	NO (80.23 %) YES (19.77 %)	Question Composition 50.79
$\neg Q_2 \wedge Q_1$: Is the man <i>not</i> wearing shoes <i>and</i> is there beer?	NO (62.00 %) YES (37.99 %)	
$Q_1 \wedge C$: Is there beer and does this seem like a man bending over to look inside of a fridge?	NO (100 %) YES (0.00 %)	Composition using COCO annotations 50.51
$\neg Q_2 \vee B$: Is the man <i>not</i> wearing shoes or is there a clock?	NO (100 %) YES (0.00 %)	
$Q_1 \wedge \text{antonym}(B)$: Is there beer and is there a wine glass?	YES (84.37 %) NO (15.60 %)	

Annotations from COCO

OBJECTS (B):

person, bottle, bowl, microwave, fridge, clock

CAPTIONS (C):

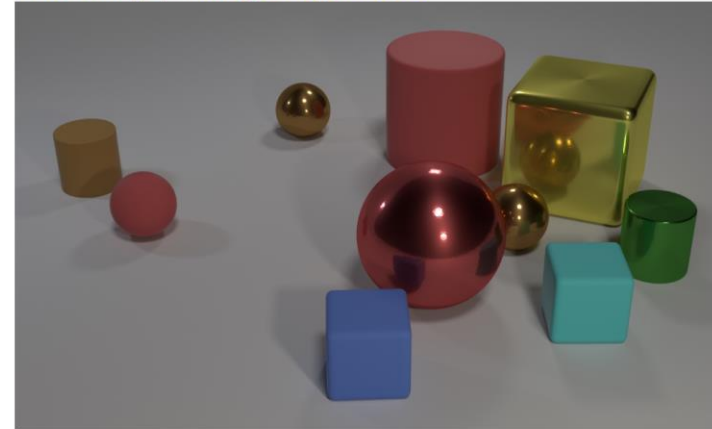
“A man bending over to look inside of a fridge.”

“A person standing in front of an opened refrigerator?”

Diagnostic Datasets

- CLEVR (Compositional Language and Elementary Visual Reasoning)
 - Has been extended to visual dialog (CLEVR-Dialog), referring expressions (CLEVR-Ref+), and video reasoning (CLEVRER)

Questions in CLEVR test various aspects of visual reasoning including **attribute identification**, **counting**, **comparison**, **spatial relationships**, and **logical operations**.



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 **red** things?

- [1] CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning, CVPR 2017
- [2] CLEVR-Dialog: A Diagnostic Dataset for Multi-Round Reasoning in Visual Dialog, NAACL 2019
- [3] CLEVR-Ref+: Diagnosing Visual Reasoning with Referring Expressions, CVPR 2019
- [4] CLEVRER: CoLLision Events for Video REpresentation and Reasoning, ICLR 2020

Beyond VQA: Visual Grounding

- Referring Expression Comprehension: RefCOCO(+/g)
 - ReferIt Game: Referring to Objects in Photographs of Natural Scenes
- Flickr30k Entities



A man with pierced ears is wearing glasses and an orange hat.
A man with glasses is wearing a beer can croched hat.
A man with gauges and glasses is wearing a Blitz hat.
A man in an orange hat starring at something.
A man wears an orange hat and glasses.

[1] OK-VQA: A Visual Question Answering Benchmark Requiring External Knowledge, EMNLP 2014

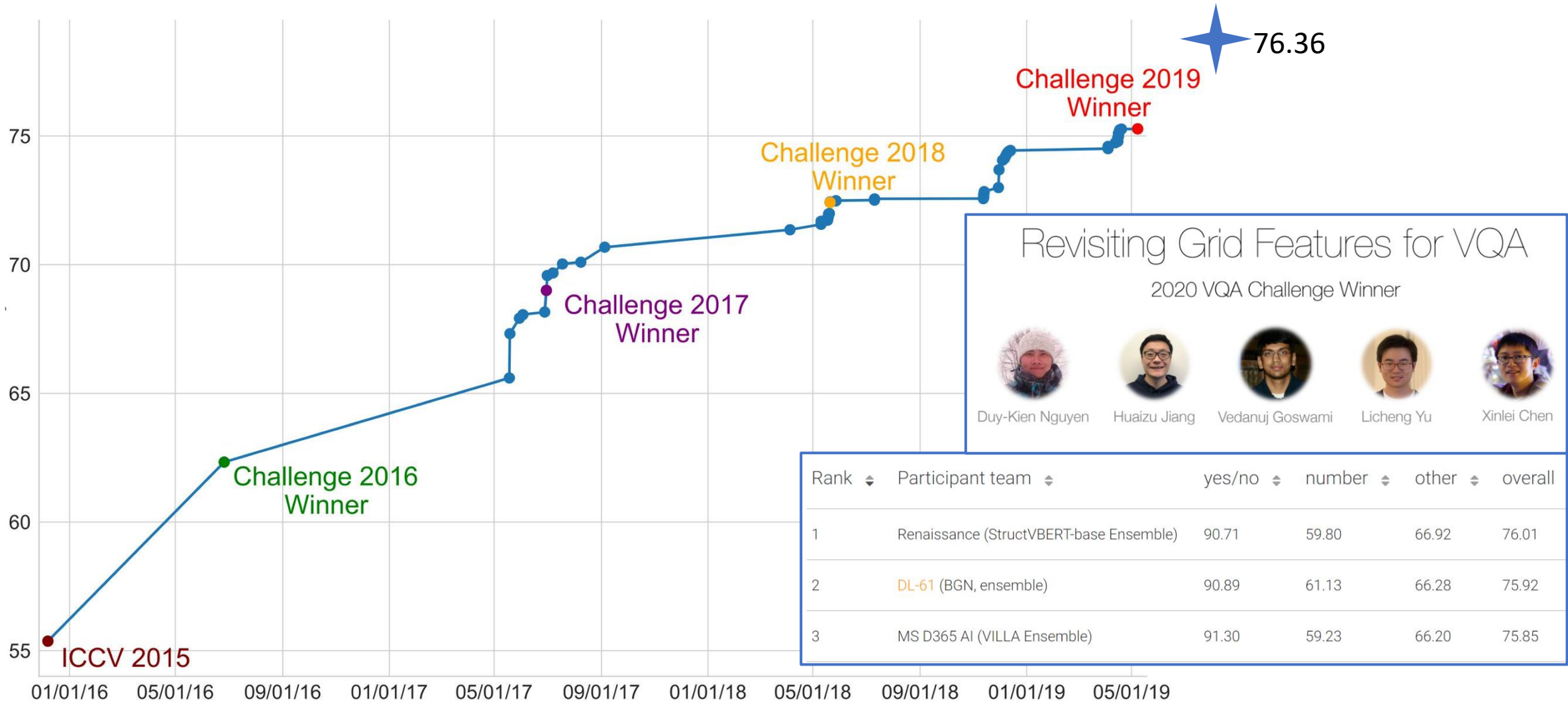
[2] Flickr30K Entities: Collecting Region-to-Phrase Correspondences for Richer Image-to-Sentence Models, IJCV 2017

Beyond VQA: Visual Grounding

- PhraseCut: Language-based image segmentation



Visual Question Answering



Agenda

- Task Overview

- *What are the main tasks that are driving progress in VQA and visual reasoning?*

- Method Overview

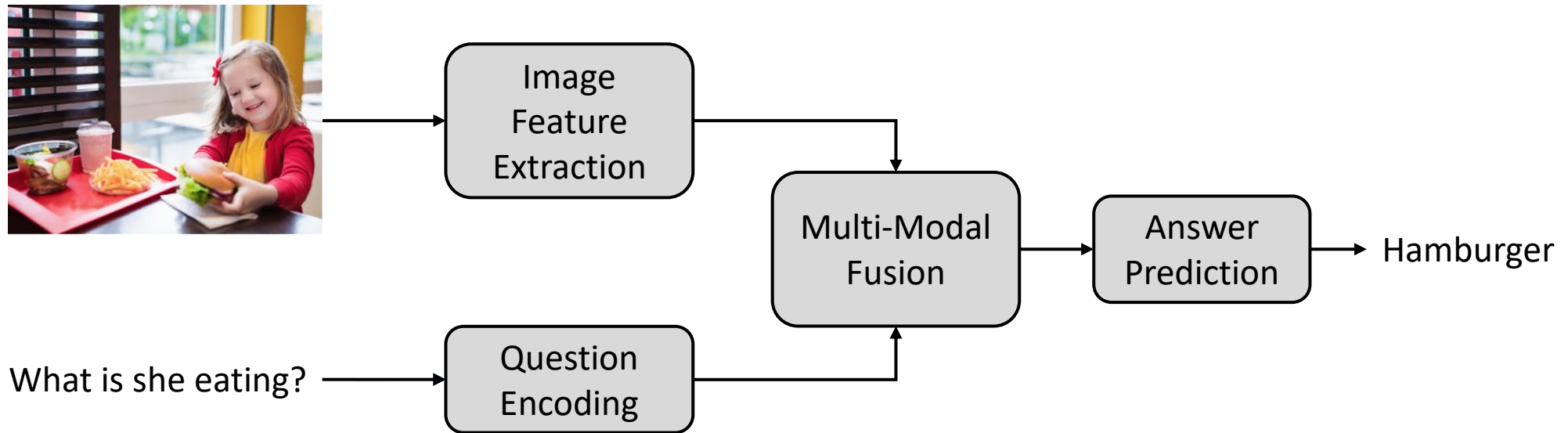
- *What are the state-of-the-art approaches and the key model design principles underlying these methods?*

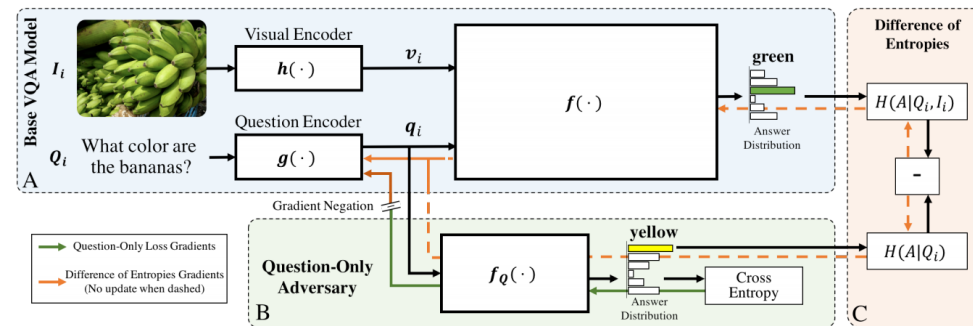
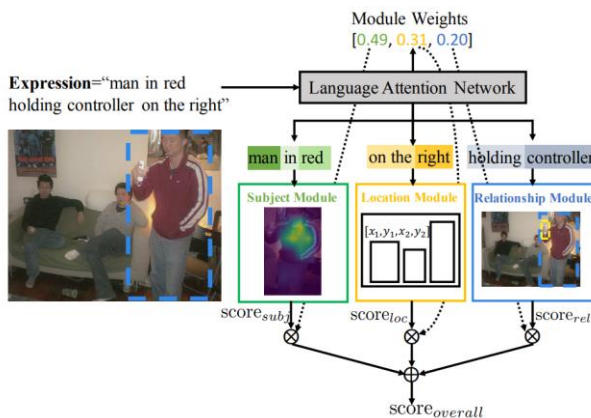
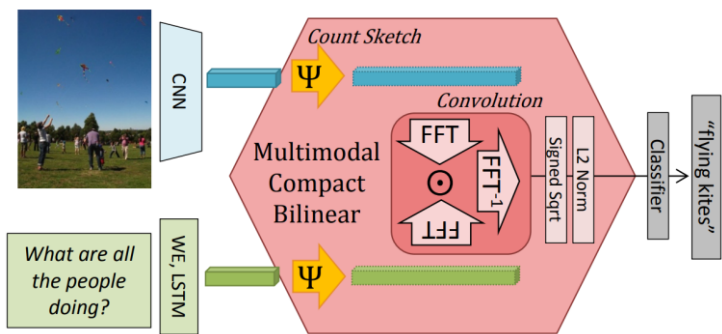
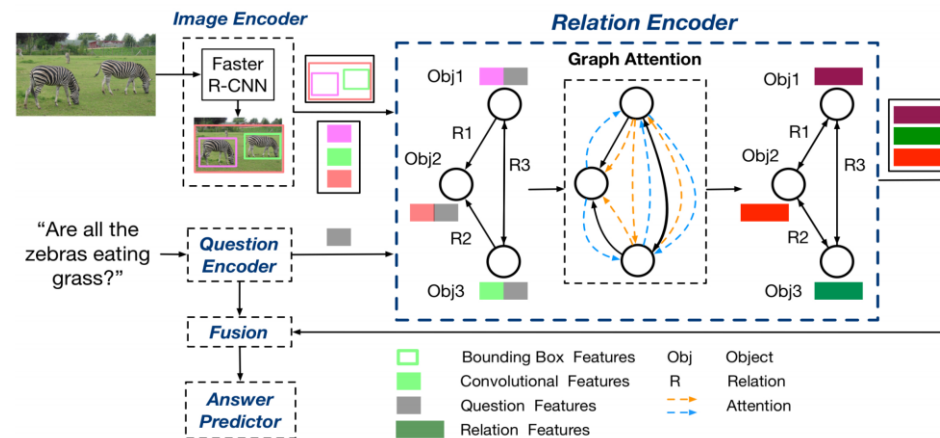
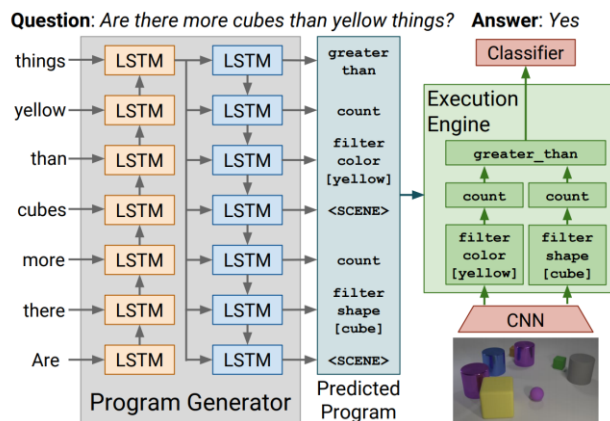
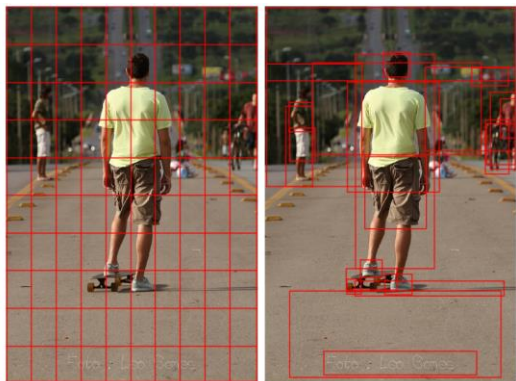
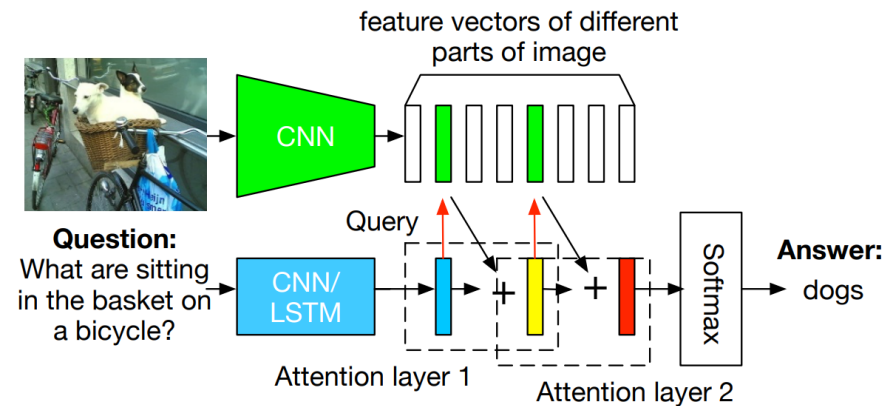
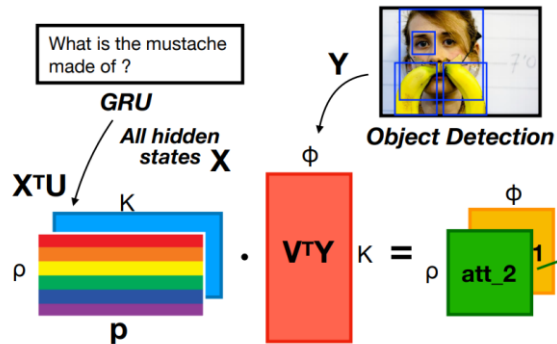
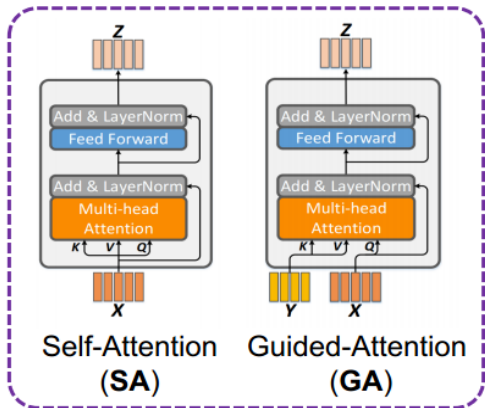
- Summary

- *What are the core challenges and future directions?*

Overview

- How a typical system looks like



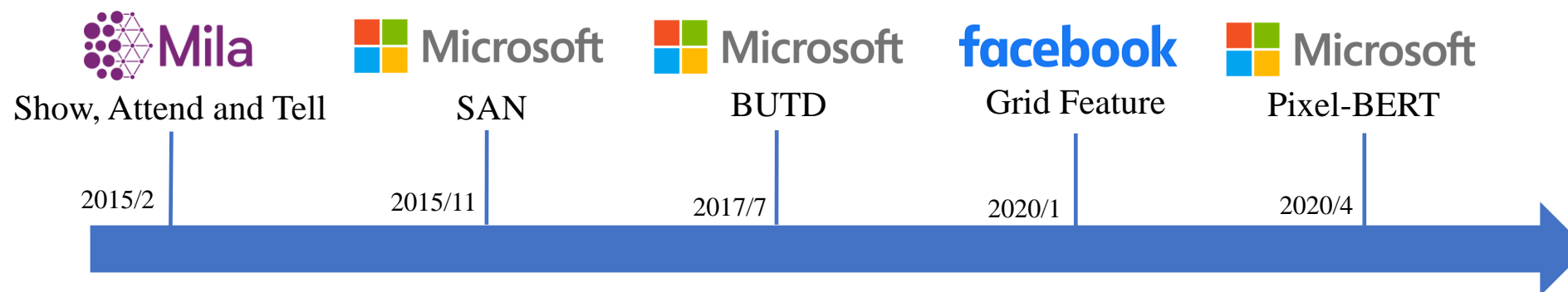


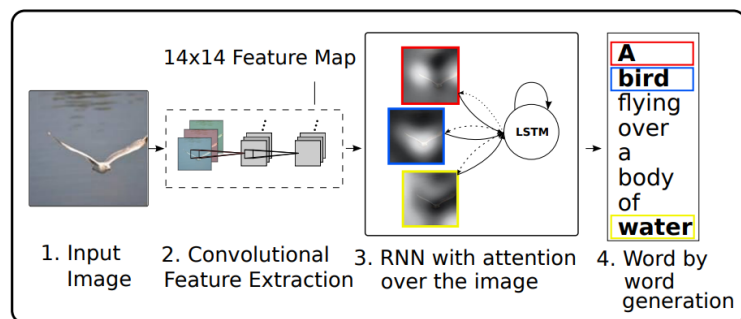
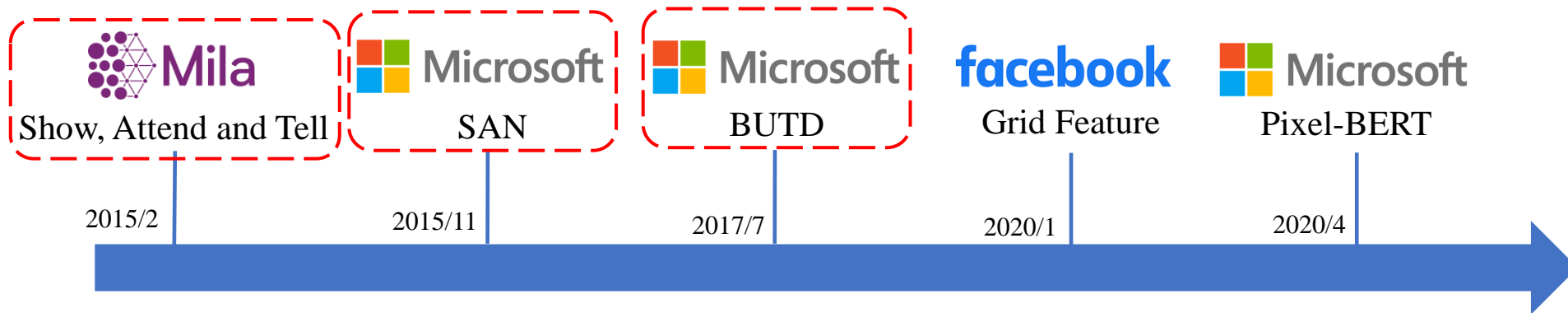
Overview

- Better image feature preparation
- Enhanced multimodal fusion
 - Bilinear pooling: how to fuse two vectors into one
 - Multimodal alignment: *cross-modal* attention
 - Incorporation of object relations: *intra-modal* self-attention, graph attention
 - Multi-step reasoning
- Neural module networks for compositional reasoning
- Robust VQA (briefly mention)
- Multimodal pre-training (briefly mention)

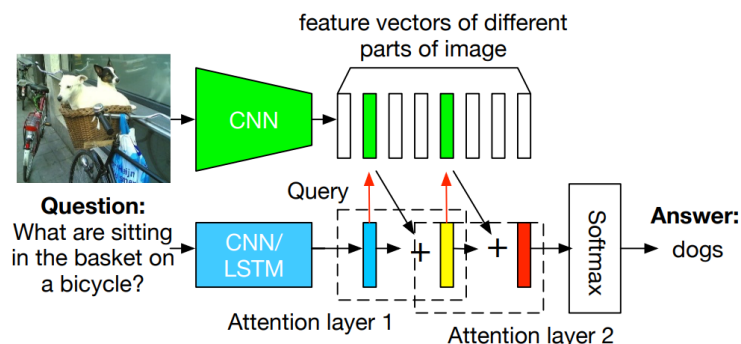
Better Image Feature Preparation

- From *grid* features to *region* features, and to *grid* features again

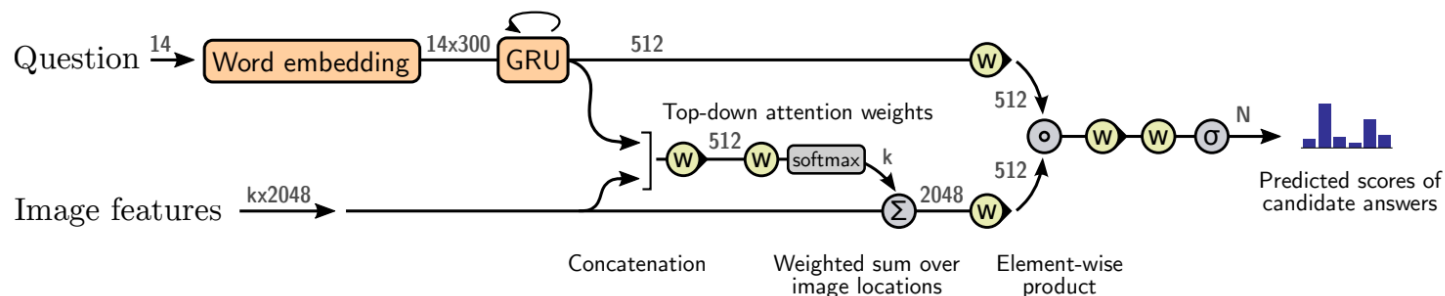




Show, Attend and Tell



Stacked Attention Network



2017 VQA Challenge Winner

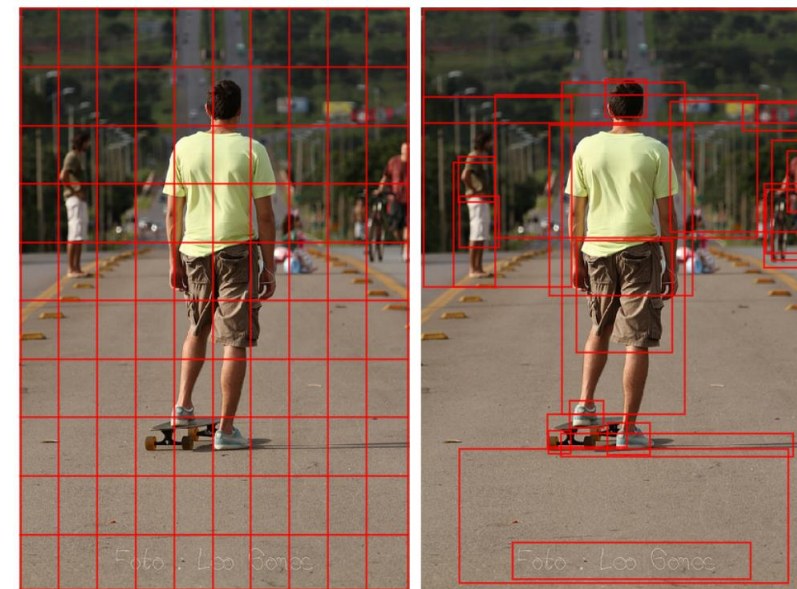


Figure 1. Typically, attention models operate on CNN features corresponding to a uniform grid of equally-sized image regions (left). Our approach enables attention to be calculated at the level of objects and other salient image regions (right).

[1] Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

[2] Stacked Attention Networks for Image Question Answering, CVPR 2016

[3] Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering, CVPR 2018

 **Mila**
Show, Attend and Tell

2015/2

 **Microsoft**

SAN

2015/11

 **Microsoft**

BUTD

2017/7

 **facebook**

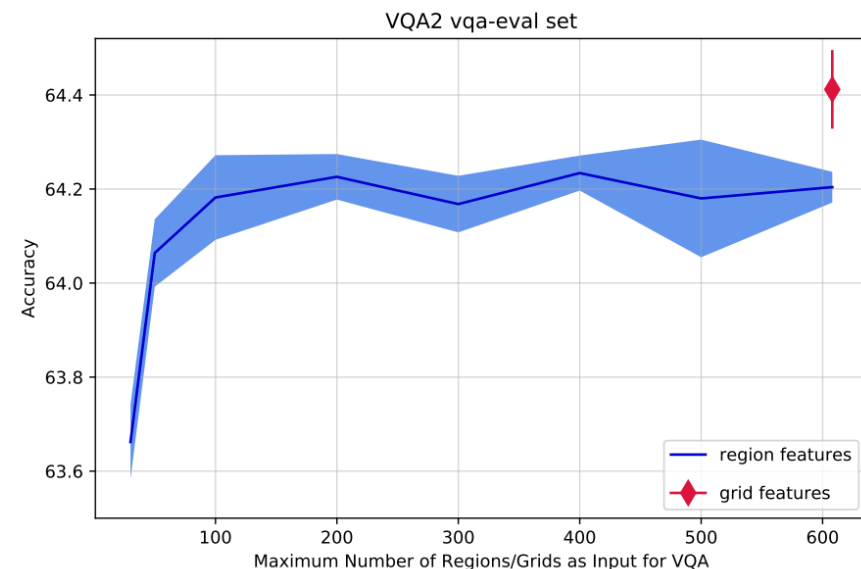
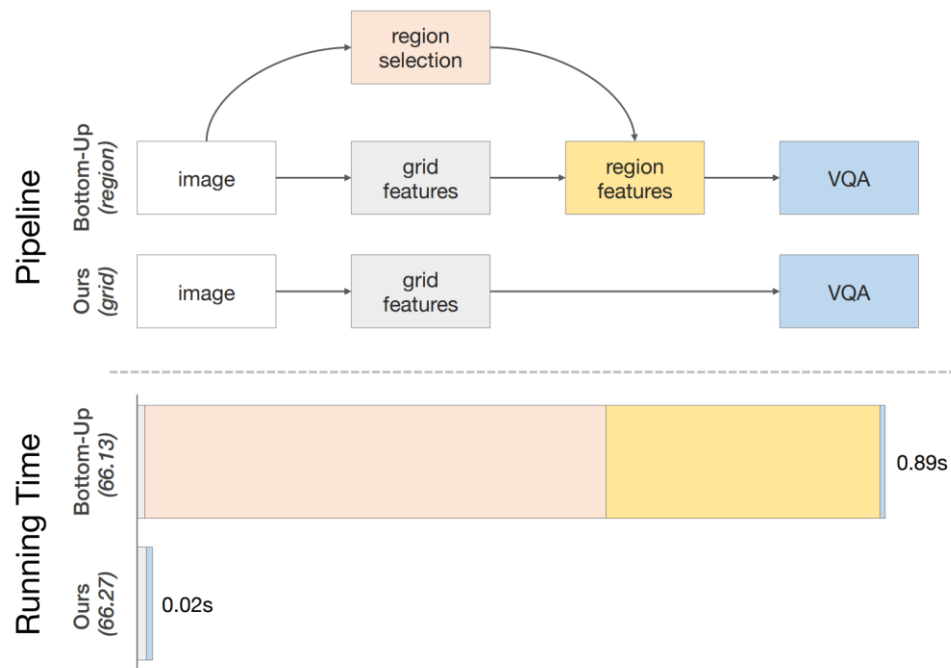
Grid Feature

2020/1

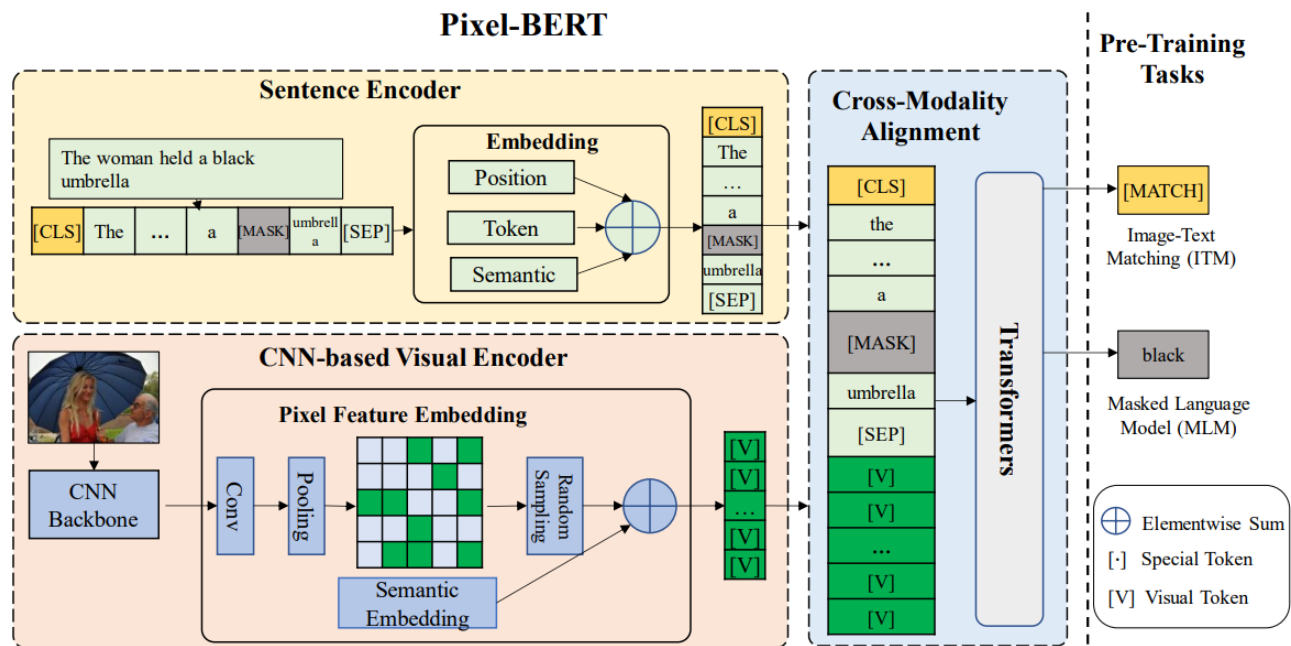
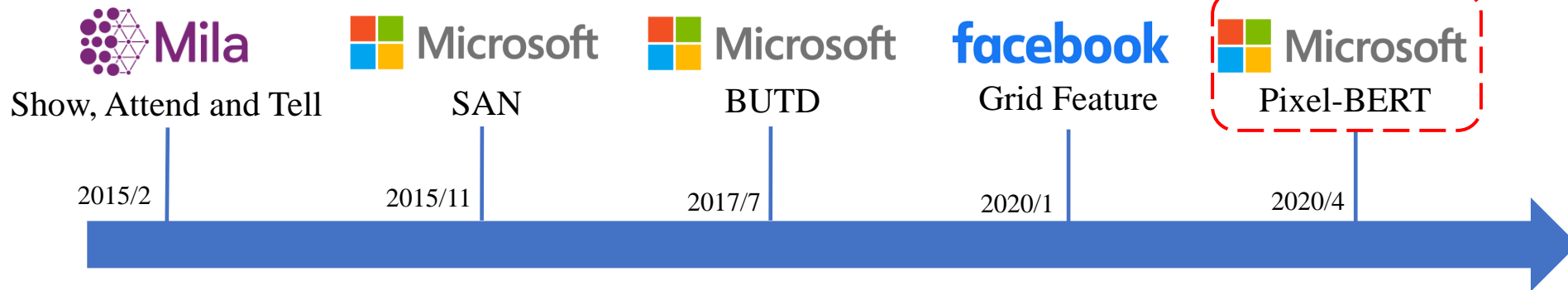
 **Microsoft**

Pixel-BERT

2020/4



In Defense of Grid Features for VQA

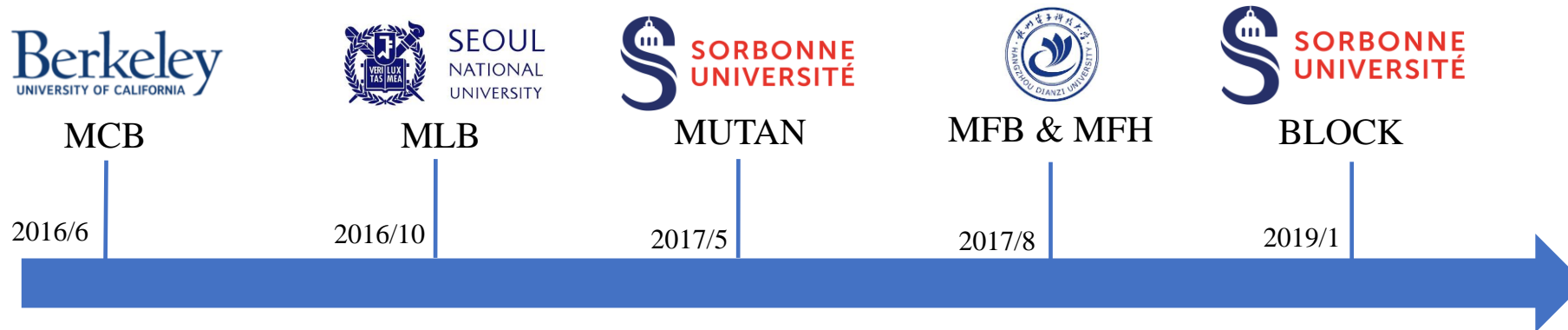


Model	test-dev	test-std
MUTAN[5]	60.17	-
BUTD[2]	65.32	65.67
ViLBERT[21]	70.55	70.92
VisualBERT[19]	70.80	71.00
VLBERT[29]	71.79	72.22
LXMERT[33]	72.42	72.54
UNITER[6]	72.27	72.46
Pixel-BERT (r50)	71.35	71.42
Pixel-BERT (x152)	74.45	74.55

Table 2. Evaluation of Pixel-BERT with other methods on VQA.

Bilinear Pooling

- Instead of simple concatenation and element-wise product for fusion, bilinear pooling methods have been studied
- Bilinear pooling and attention mechanism can be enhanced with each other





2016/6



2016/10



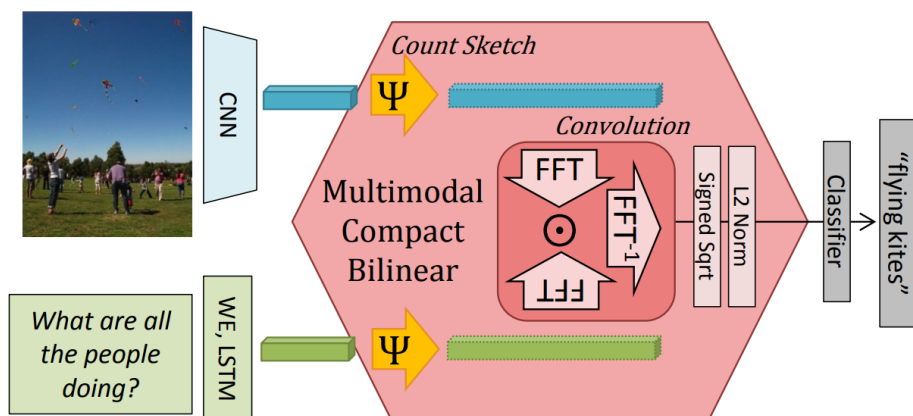
2017/5



2017/8



2019/1



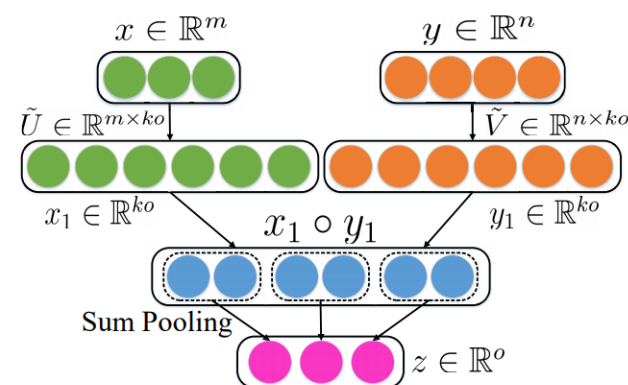
Multimodal Compact Bilinear Pooling

2016 VQA Challenge Winner

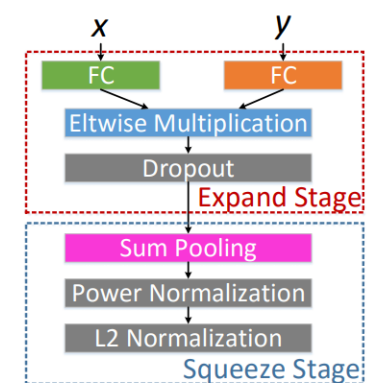
However, the feature after FFT is very high dimensional.

$$\mathbf{f} = \mathbf{P}^T (\mathbf{U}^T \mathbf{x} \circ \mathbf{V}^T \mathbf{y}) + \mathbf{b}$$

Multimodal Low-rank Bilinear Pooling



(a) Multi-modal Factorized Bilinear Pooling

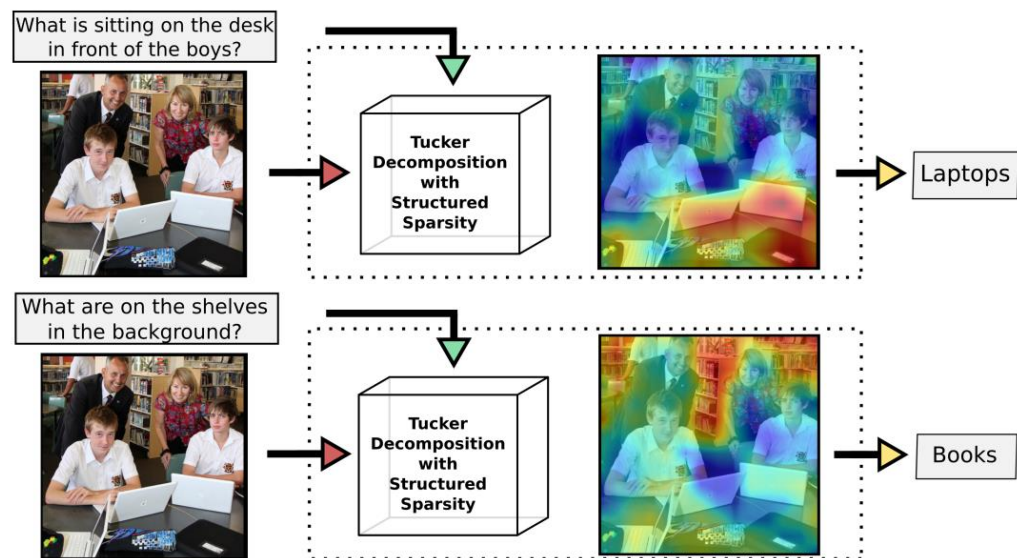


(b) MFB module

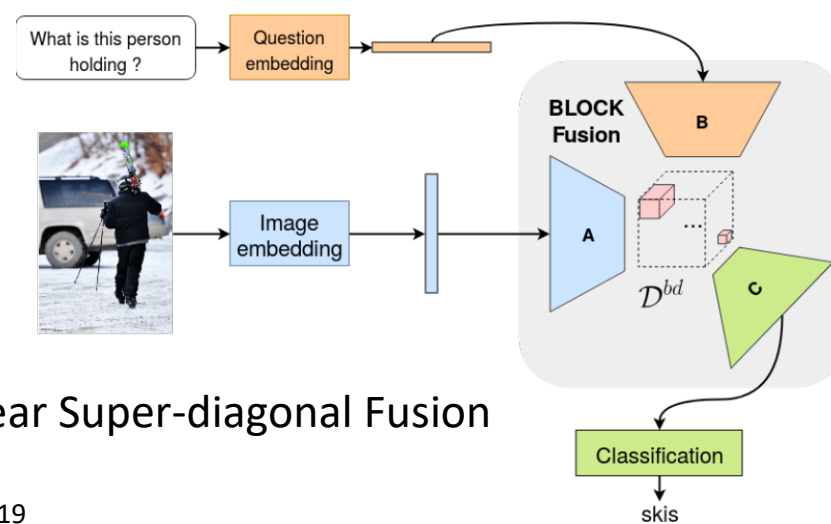
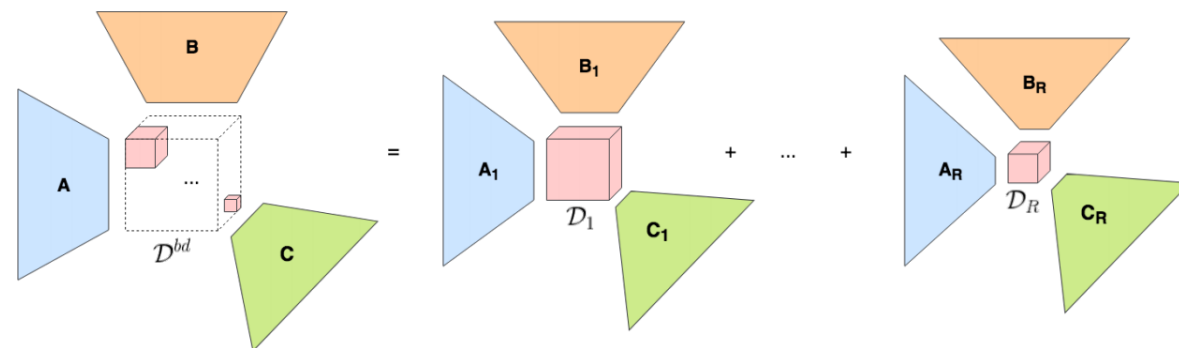
[1] Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding, EMNLP 2016

[2] Hadamard Product for Low-rank Bilinear Pooling, ICLR 2017

[3] Multi-modal Factorized Bilinear Pooling with Co-Attention Learning for Visual Question Answering, ICCV 2017

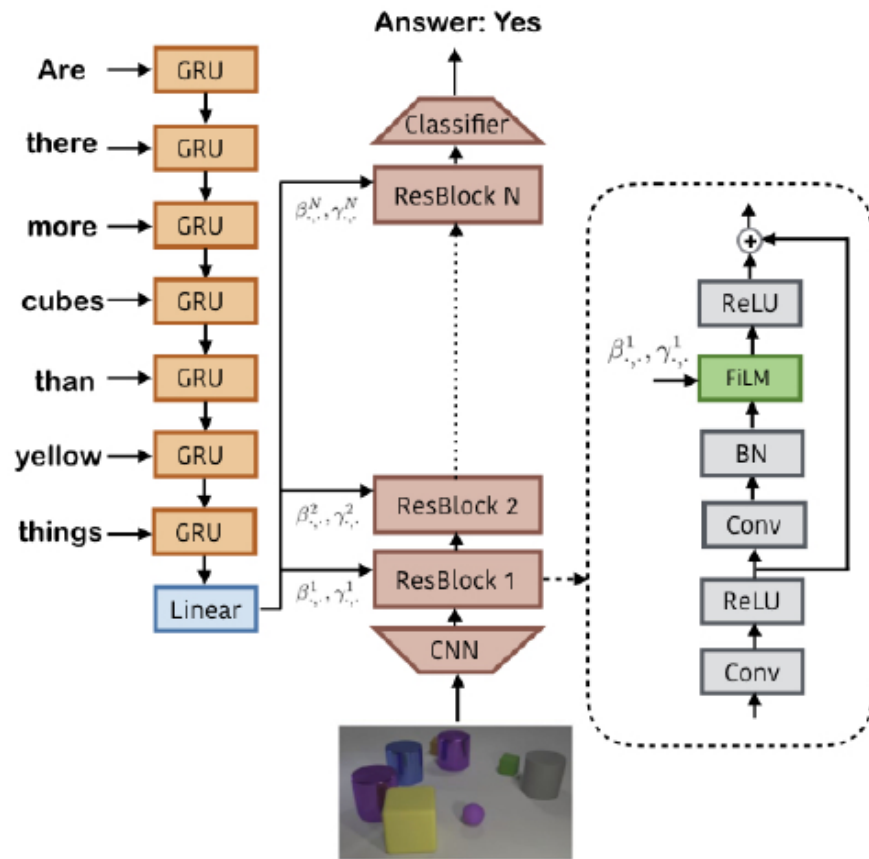


Multimodal Tucker Fusion



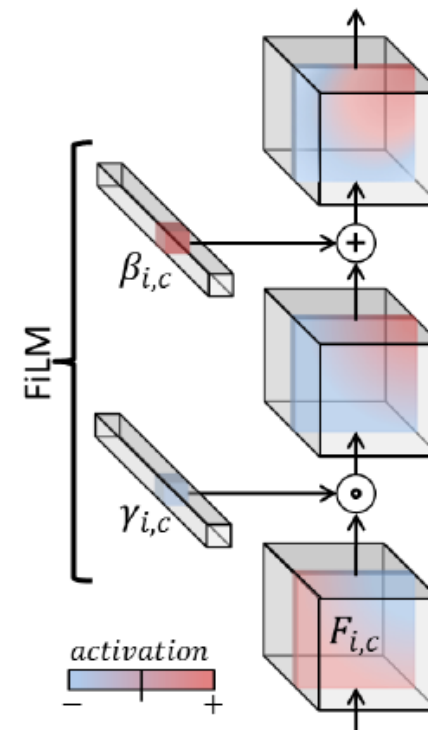
Bilinear Super-diagonal Fusion

FiLM: Feature-wise Linear Modulation



$$\gamma_{i,c} = f_c(\mathbf{x}_i) \quad \beta_{i,c} = h_c(\mathbf{x}_i),$$
$$FiLM(\mathbf{F}_{i,c} | \gamma_{i,c}, \beta_{i,c}) = \gamma_{i,c} \mathbf{F}_{i,c} + \beta_{i,c}.$$

Something similar to conditional batch normalization



Multimodal Alignment

- Cross-modal attention:
 - Tons of work in this area
 - Early work: questions attend to image grids/regions
 - Current focus: image-text co-attention





2015/11



2016/5

NAVER

DAN

2016/11



東北大学
TOHOKU UNIVERSITY

DCN

2018/4

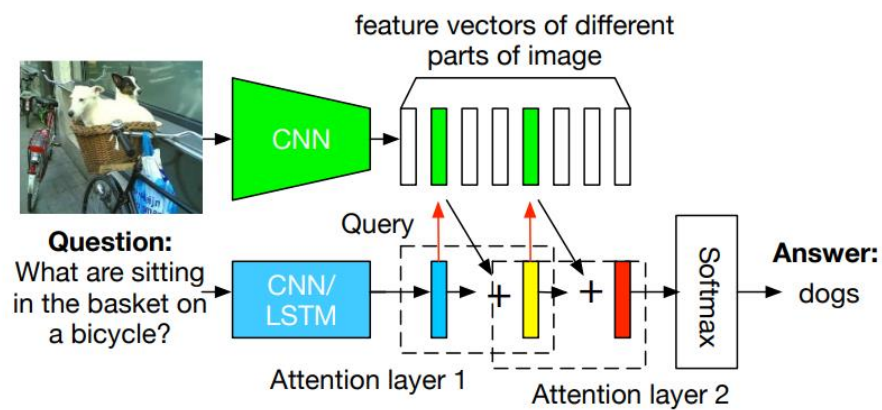


SEOUL
NATIONAL
UNIVERSITY

BAN

2018/5

...

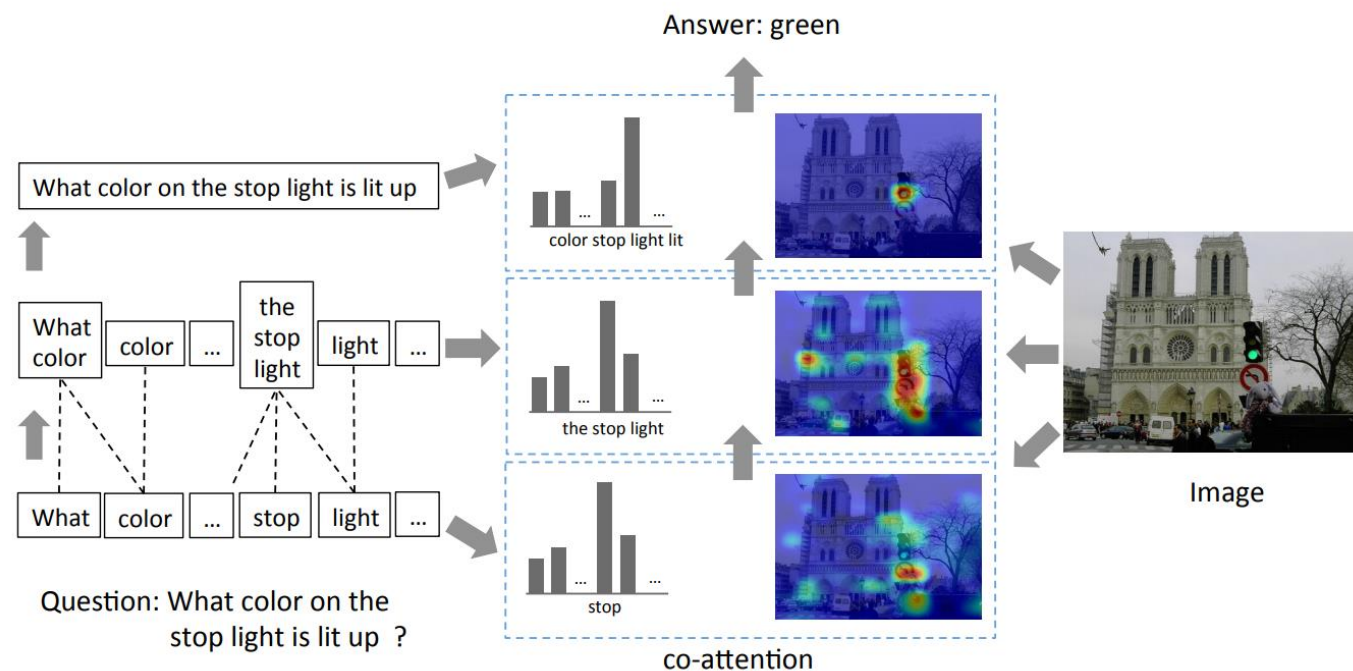


(a) Stacked Attention Network for Image QA



Original Image First Attention Layer Second Attention Layer

(b) Visualization of the learned multiple attention layers.



Parallel Co-attention and Alternative Co-attention

- [1] Stacked Attention Networks for Image Question Answering, CVPR 2016
- [2] Hierarchical Question-Image Co-Attention for Visual Question Answering, NeurIPS 2016



SAN

2015/11



HierCoAttn

2016/5

NAVER

DAN

2016/11



東北大学
TOHOKU UNIVERSITY

DCN

2018/4

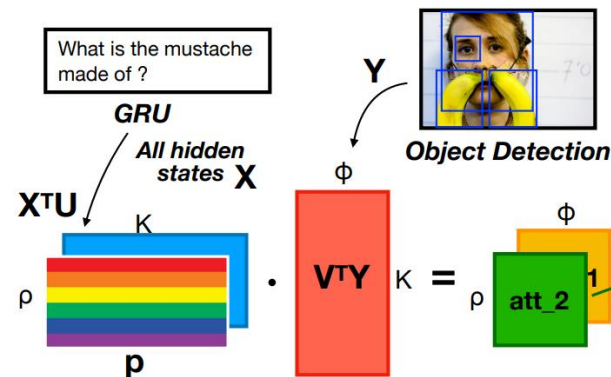
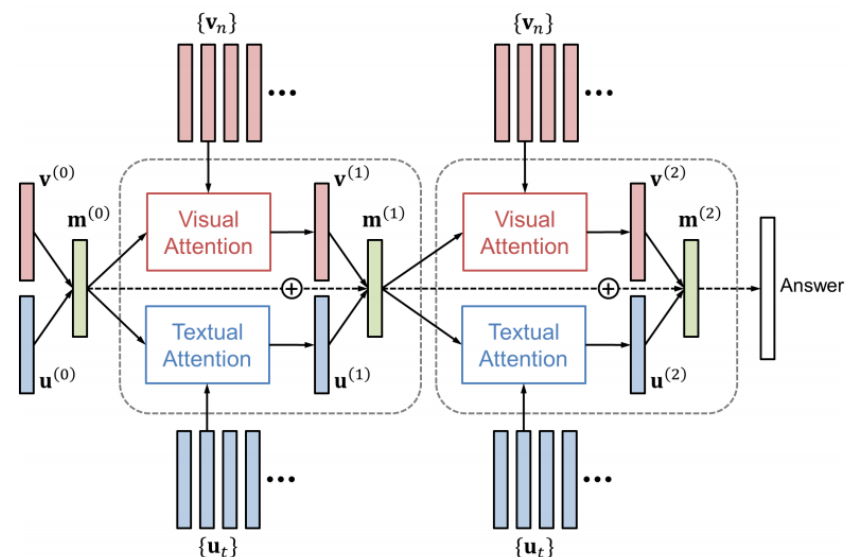


SEOUL
NATIONAL
UNIVERSITY

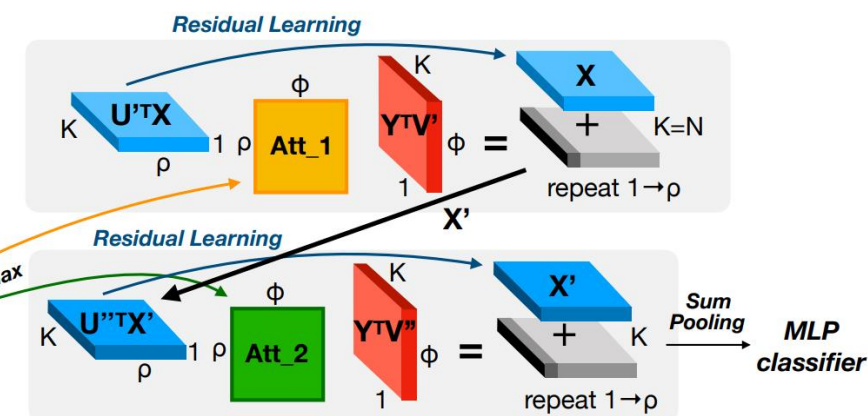
BAN

2018/5

...



Step 1. Bilinear Attention Maps



Step 2. Bilinear Attention Networks

2018 VQA Challenge Runner-Up

- Multiple Glimpses
- Counter Module
- Residual Learning
- Glove Embeddings

DAN: Dual Attention Network

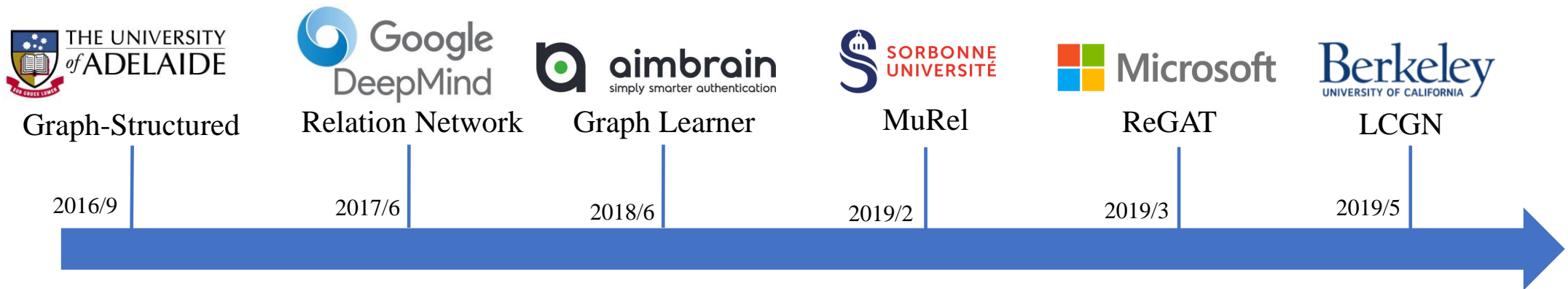
DCN: Dense Co-attention Network

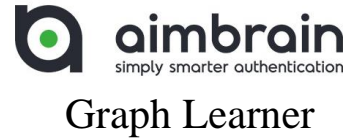
[1] Stacked Attention Networks for Image Question Answering, CVPR 2016

[2] Improved Fusion of Visual and Language Representations by Dense Symmetric Co-Attention for Visual Question Answering, CVPR 2018

Relational Reasoning

- Intra-modal attention
 - Recently becoming popular
 - Representing image as a graph
 - Graph Convolutional Network & Graph Attention Network
 - Self-attention used in Transformer





2016/9

2017/6

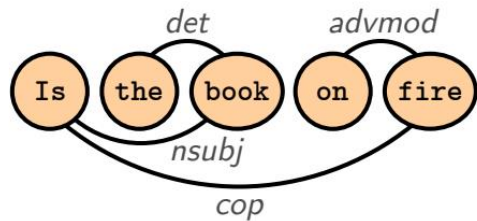
2018/6

2019/2

2019/3

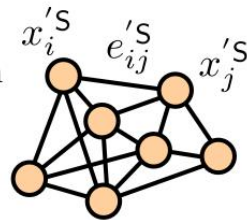
2019/5

Input scene description
and parsed question

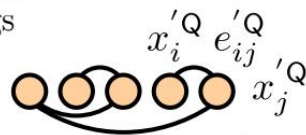


Initial
embedding

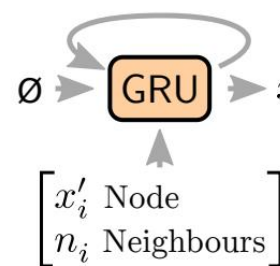
Affine
projection



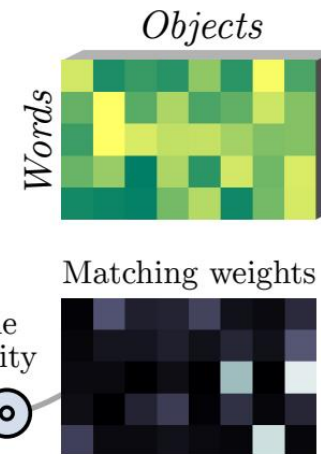
Word/vector
embeddings



Graph
processing



Combined
features

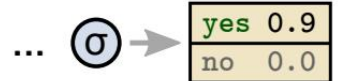


Weighted
sum



Prediction over
candidate answers

Sigmoid or
softmax



Graph-Structured Representations for Visual Question Answering



Graph-Structured

2016/9



Relation Network

2017/6



Graph Learner

2018/6



MuRel

2019/2



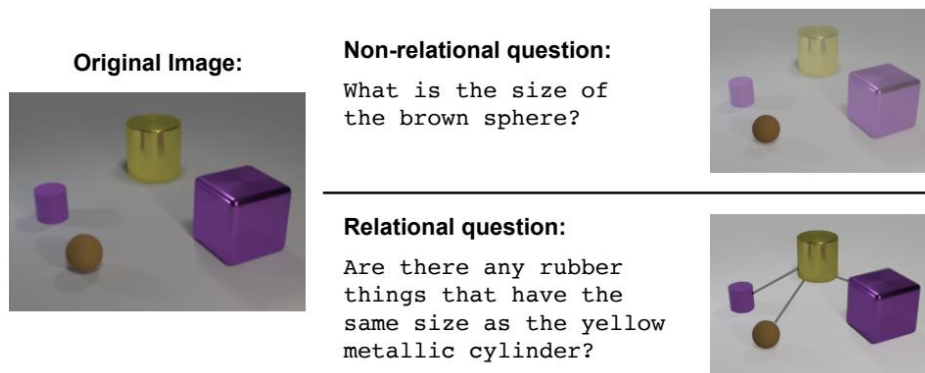
ReGAT

2019/3

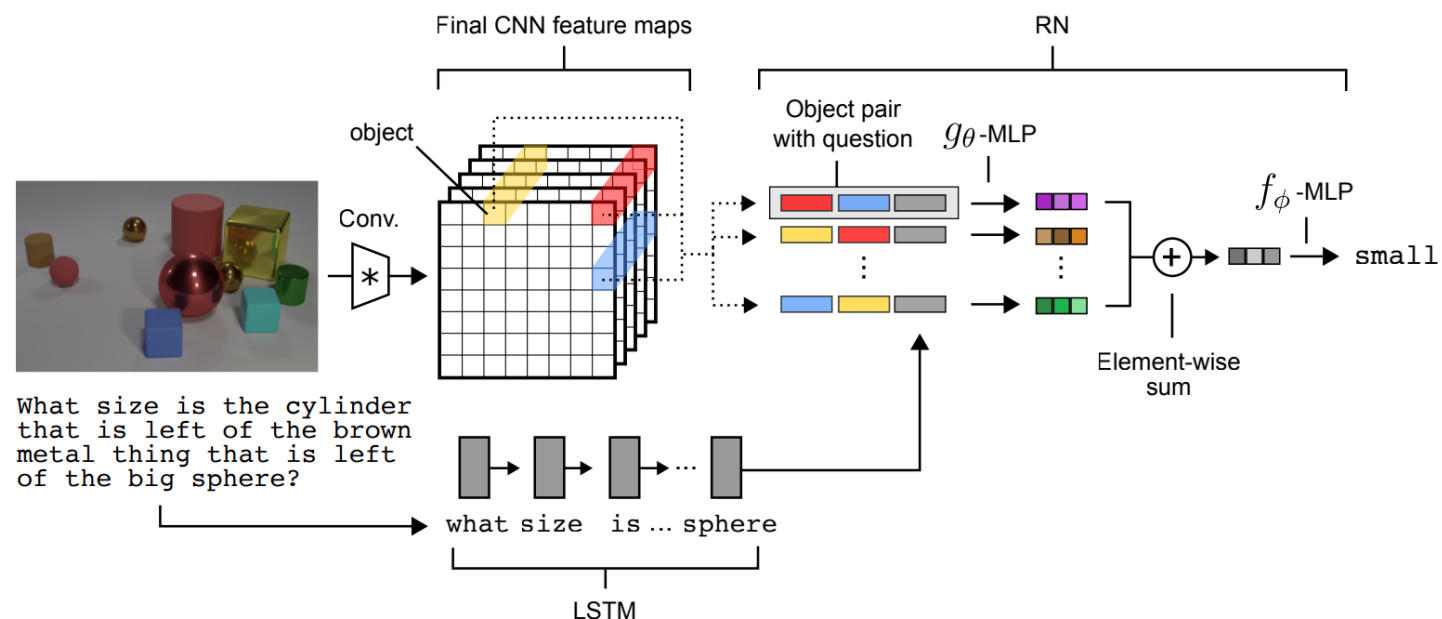


LCGN

2019/5



$$\text{RN}(O) = f_{\phi} \left(\sum_{i,j} g_{\theta}(o_i, o_j) \right)$$



Relational Network: A fully-connected graph is constructed



Graph-Structured

2016/9



Relation Network

2017/6



Graph Learner

2018/6



MuRel

2019/2



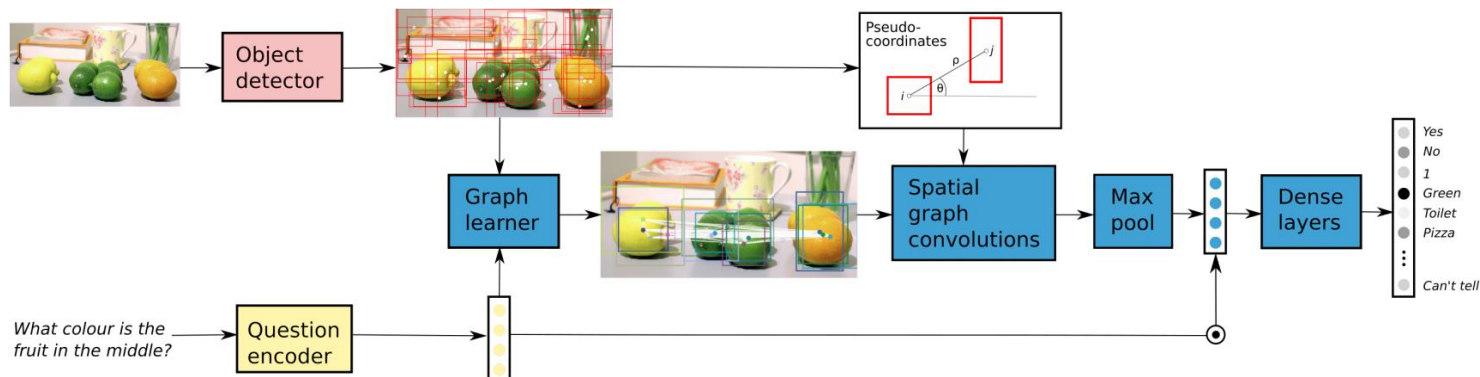
ReGAT

2019/3



LCGN

2019/5

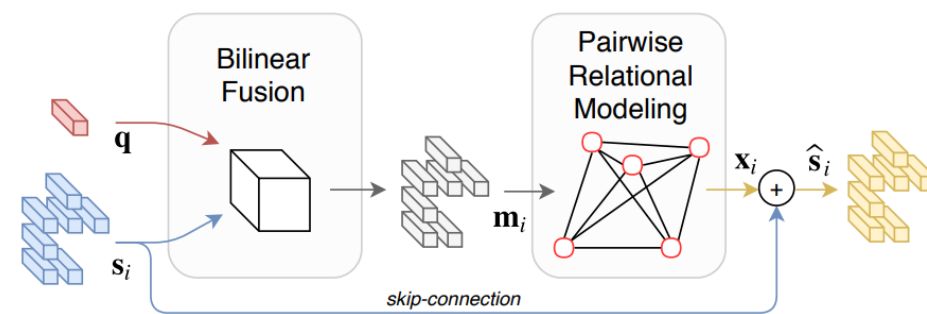
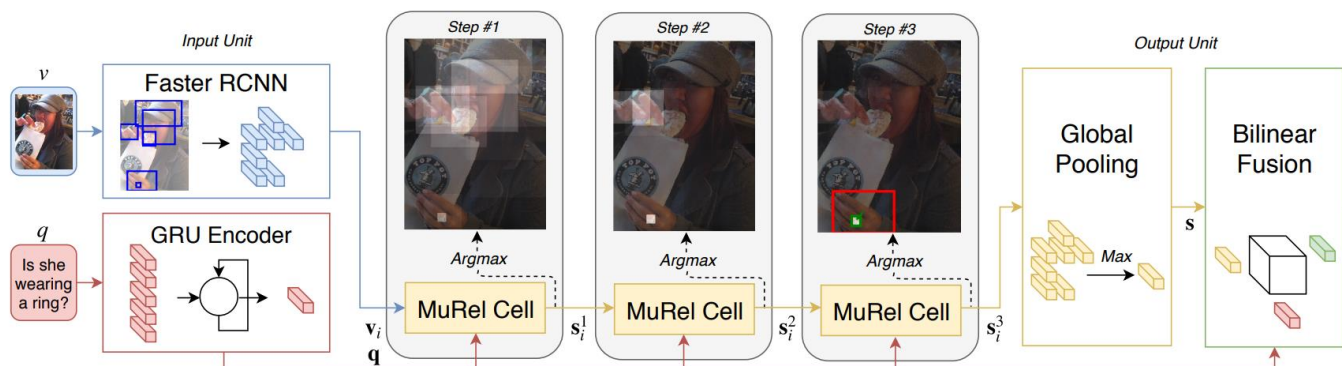


$$\mathbf{e}_n = F([\mathbf{v}_n \| \mathbf{q}]), \quad n = 1, 2, \dots, N$$

$$\mathbf{E} \in \mathbb{R}^{N \times d_e}$$

$$\mathbf{A} = \mathbf{E}\mathbf{E}^T \text{ so that } \hat{A}_{i,j} = \mathbf{e}_i^T \mathbf{e}_j.$$

$$\mathcal{N}(i) = \text{topm}(\mathbf{a}_i)$$



[1] Learning Conditioned Graph Structures for Interpretable Visual Question Answering, NeurIPS 2018

[2] MUREL: Multimodal Relational Reasoning for Visual Question Answering, CVPR 2019



Graph-Structured

2016/9



Relation Network

2017/6



Graph Learner

2018/6



MuRel

2019/2



ReGAT

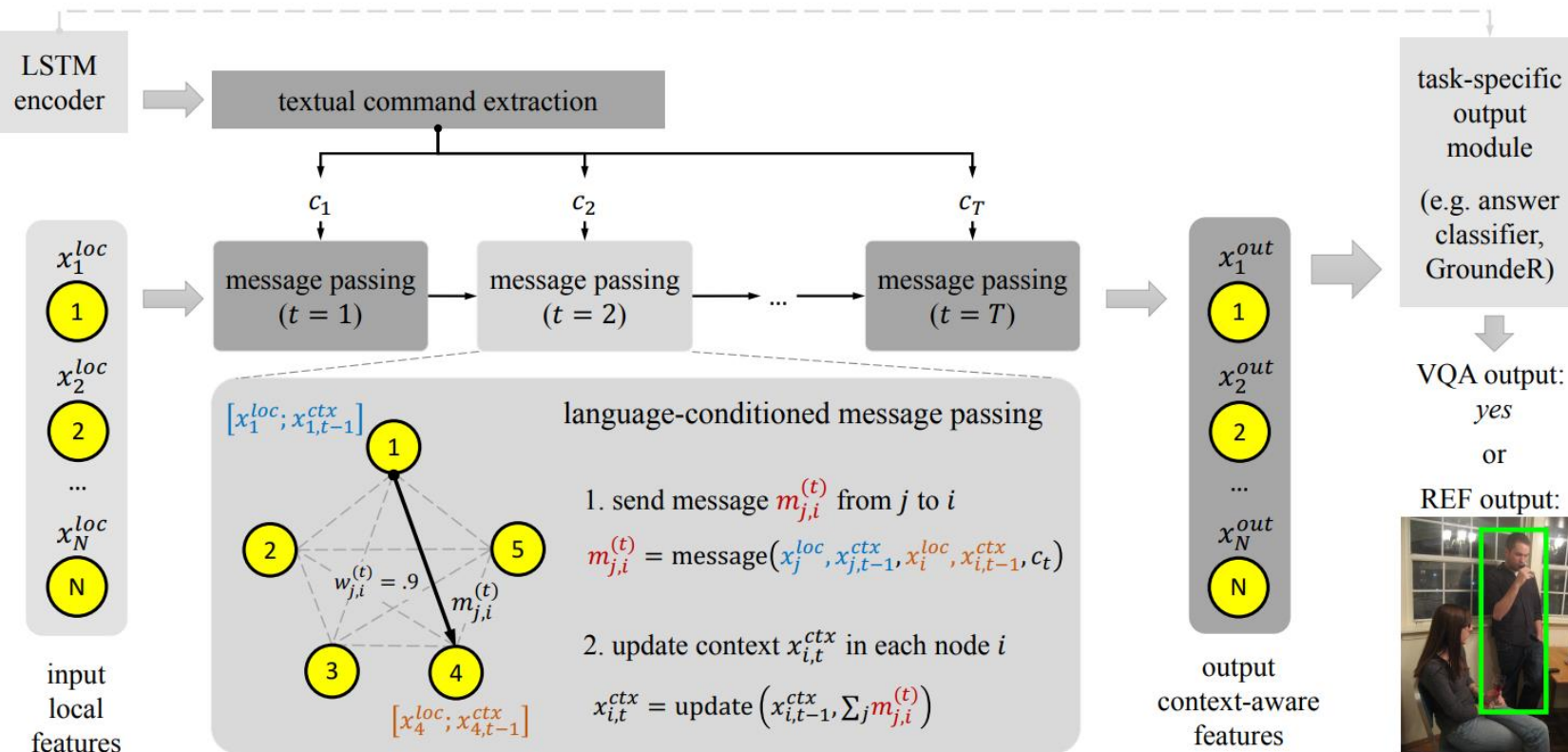
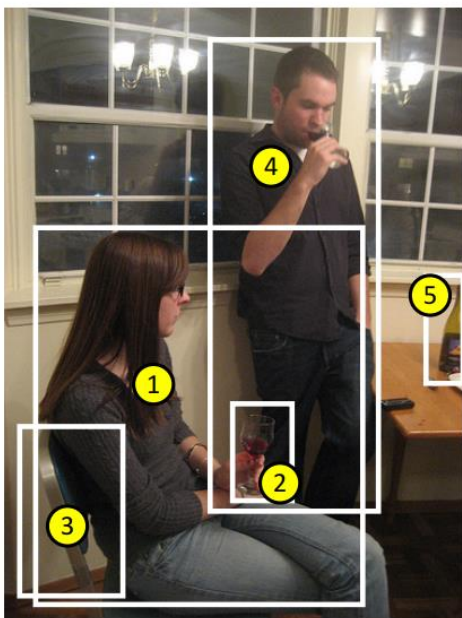
2019/3



LCGN

2019/5

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





Graph-Structured

2016/9



Relation Network

2017/6



Graph Learner

2018/6



MuRel

2019/2

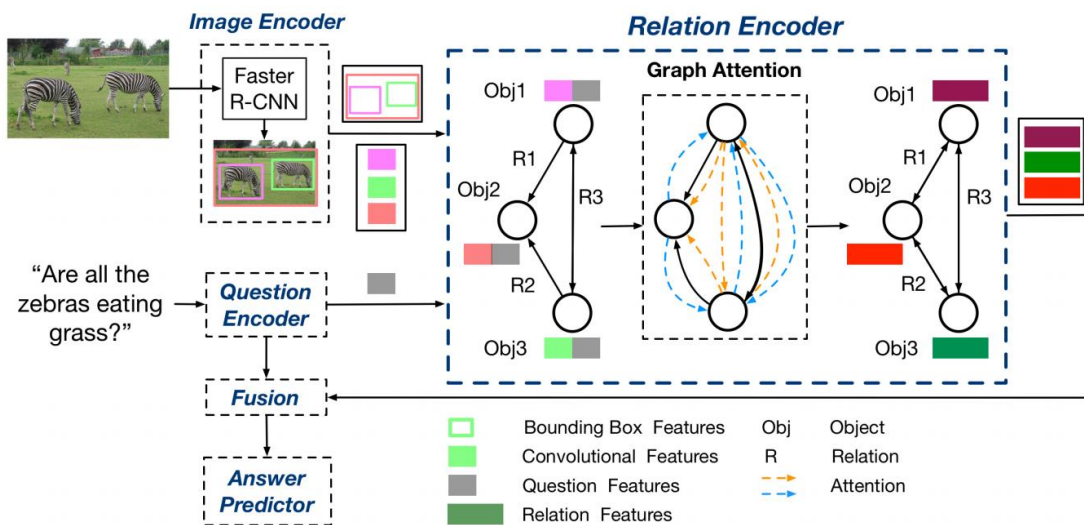


2019/3

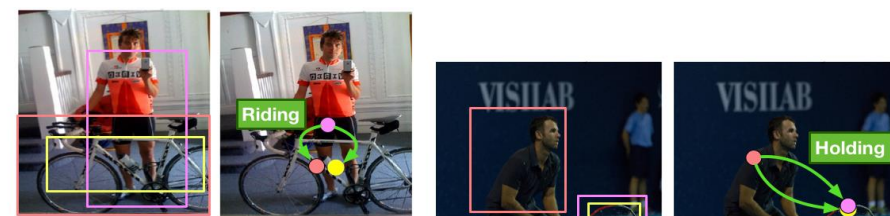


LCGN

2019/5



- Explicit Relation: Semantic & Spatial relation
- Implicit Relation: Learned dynamically during training



(a) Semantic Relation



(b) Spatial Relation



(c) Implicit Relation



Graph-Structured

2016/9



Relation Network

2017/6



Graph Learner

2018/6



MuRel

2019/2

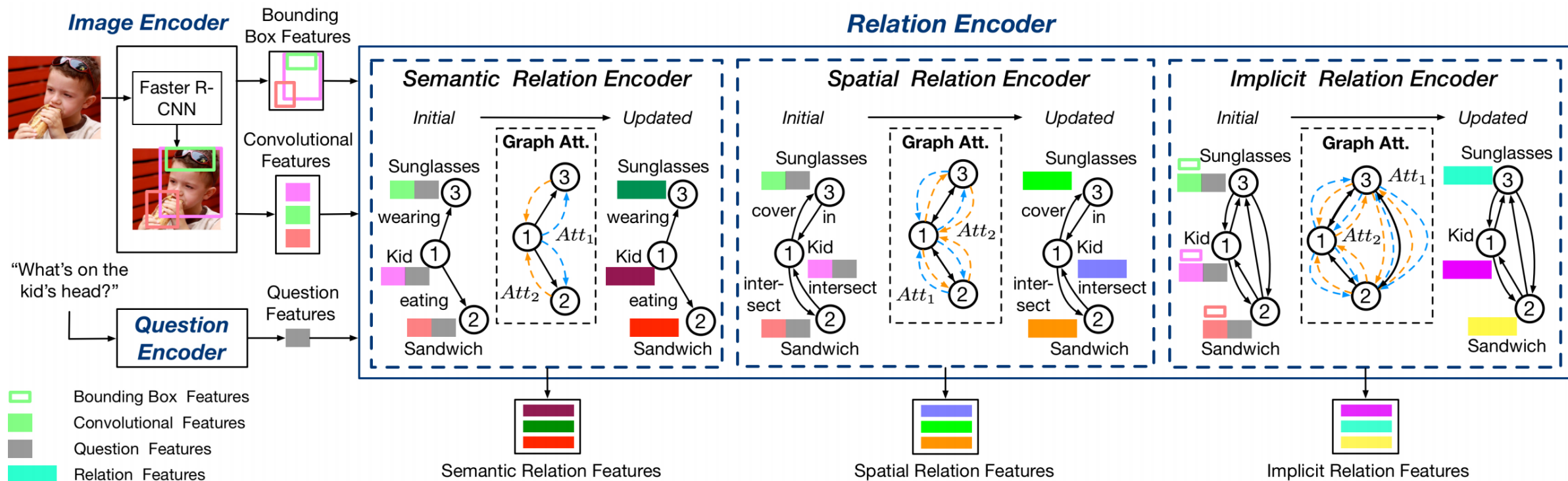


2019/3



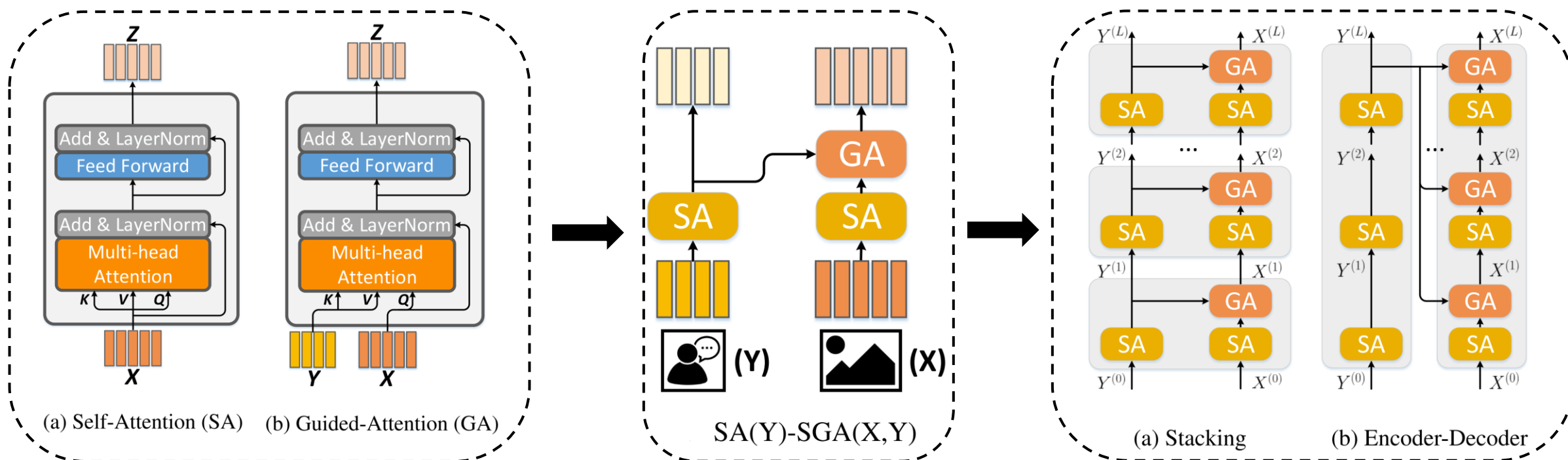
LCGN

2019/5



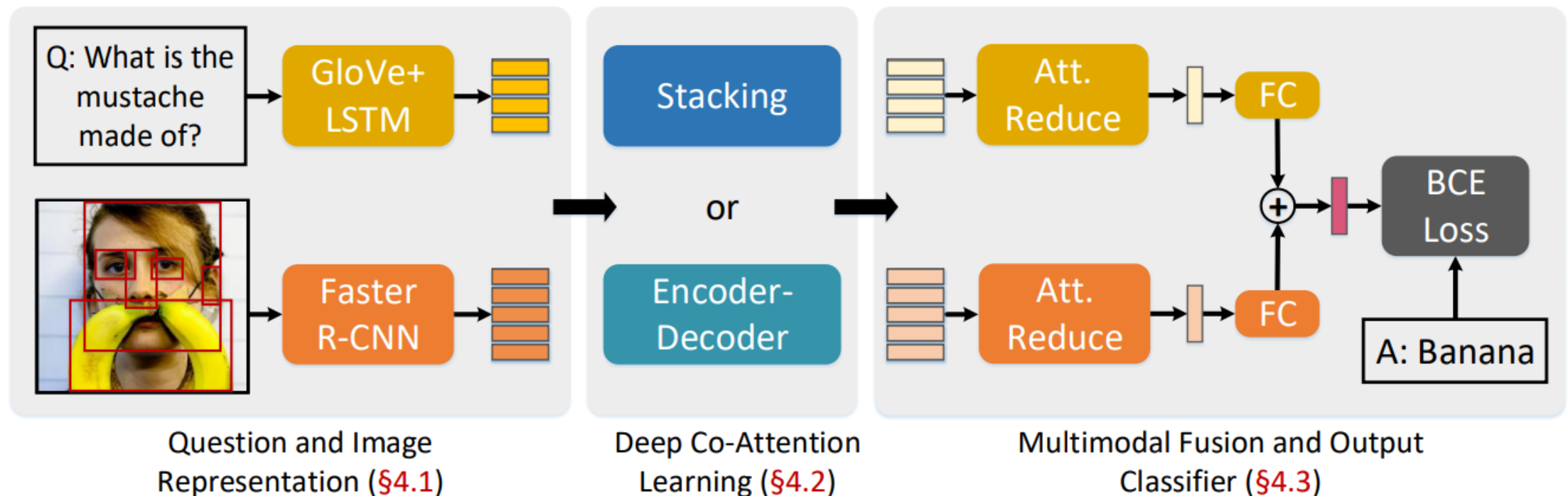
MCAN: Deep Modular Co-Attention Network

- Winning entry to VQA Challenge 2019
- Similar idea also explored in DFAF, close to [V+L pre-training](#) models



MCAN: Deep Modular Co-Attention Network

- Winning entry to VQA Challenge 2019
- Similar idea also explored in DFAF, close to [V+L pre-training](#) models

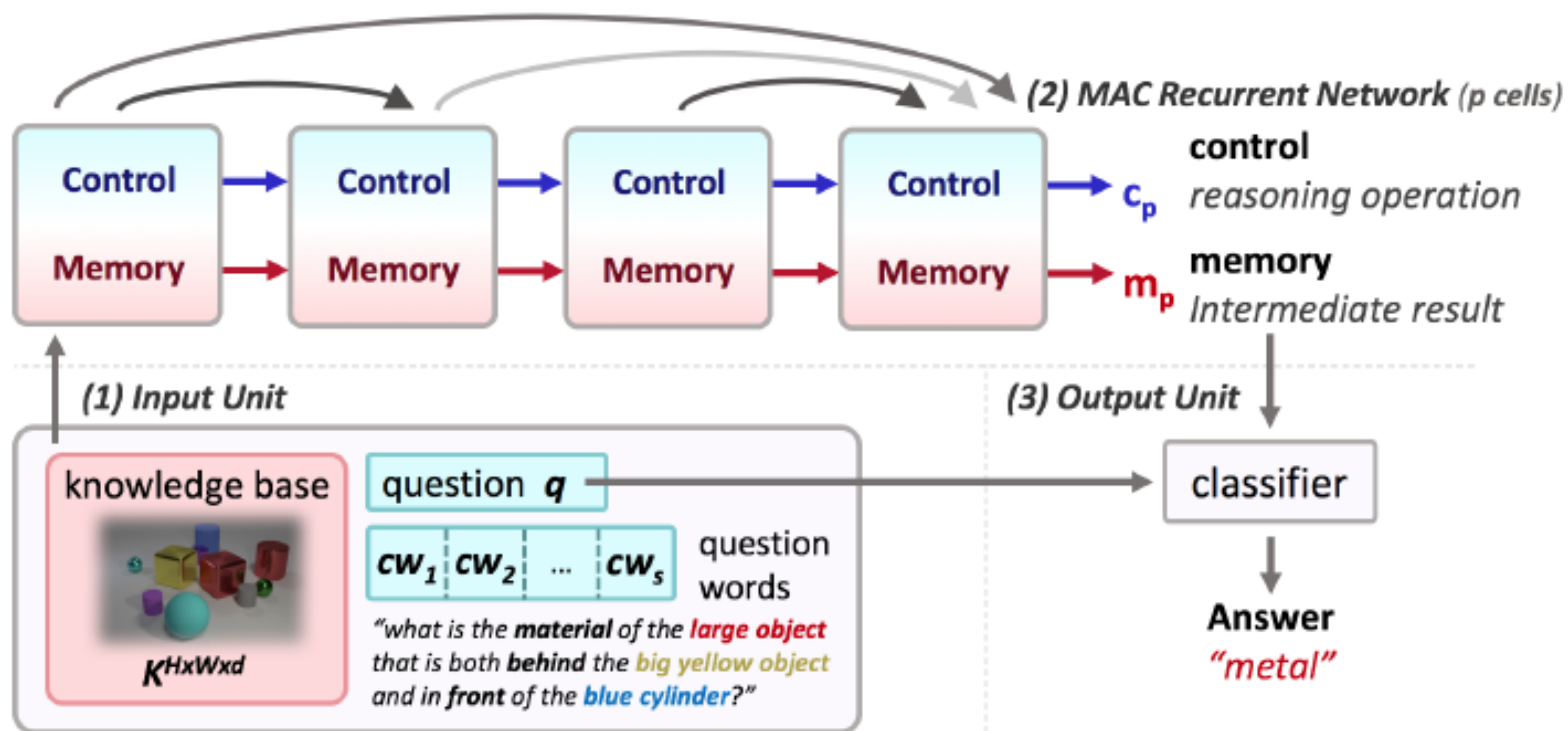


[1] Deep Modular Co-Attention Networks for Visual Question Answering, CVPR 2019

[2] Dynamic Fusion with Intra- and Inter- Modality Attention Flow for Visual Question Answering, CVPR 2019

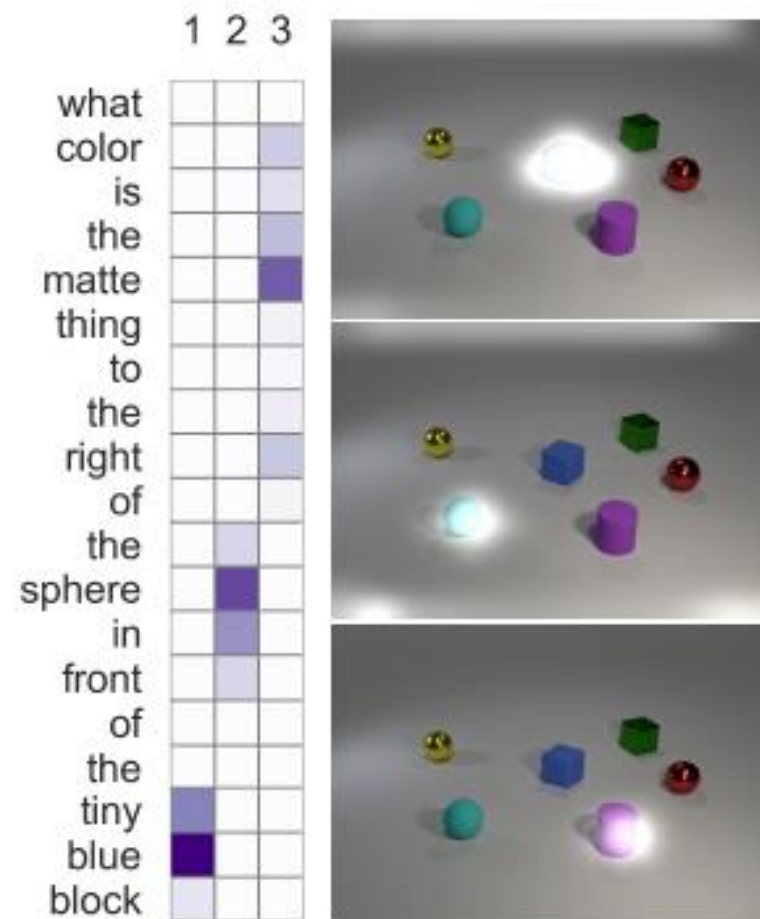
MAC: Memory, Attention and Composition

- Multi-step reasoning via recurrent MAC cells, while retaining end-to-end differentiability



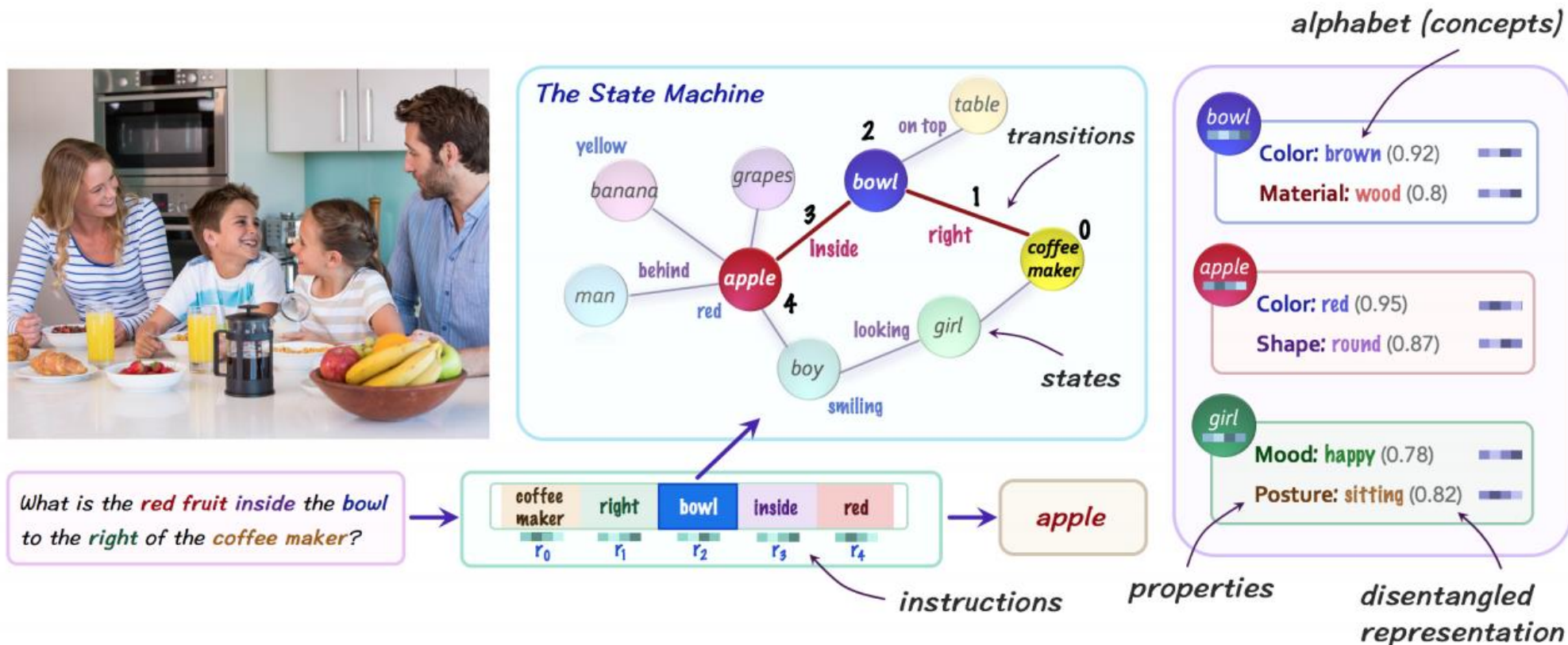
MAC: Memory, Attention and Composition

- Each cell maintains recurrent dual states:
 - Control* c_i : the reasoning operation that should be accomplished at this step.
 - Memory* m_i : the retrieved information relevant to the query, accumulated over previous iterations.
 - Implementation-wise*:
 - Attention-based average** of a given query (question)
 - Attention-based average** of a given Knowledge Base (image)



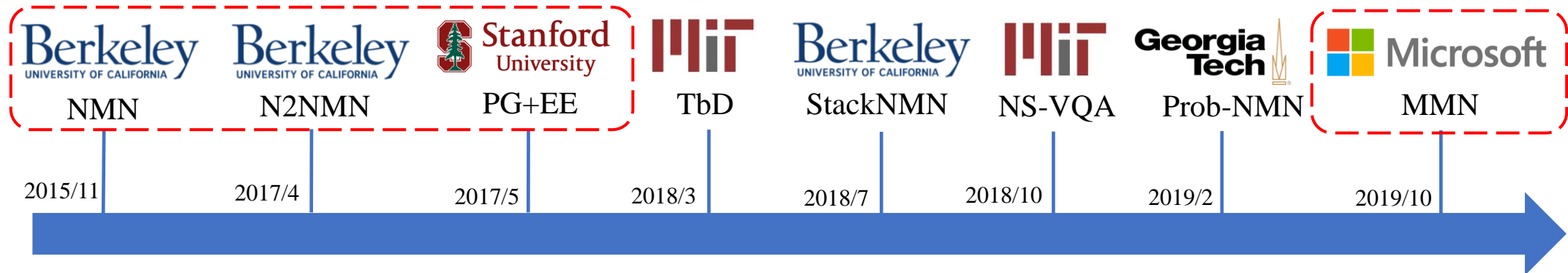
Neural State Machine

- We see and reason with **concepts**, not visual details, 99% of the time
- We build semantic **world models** to represent our environment



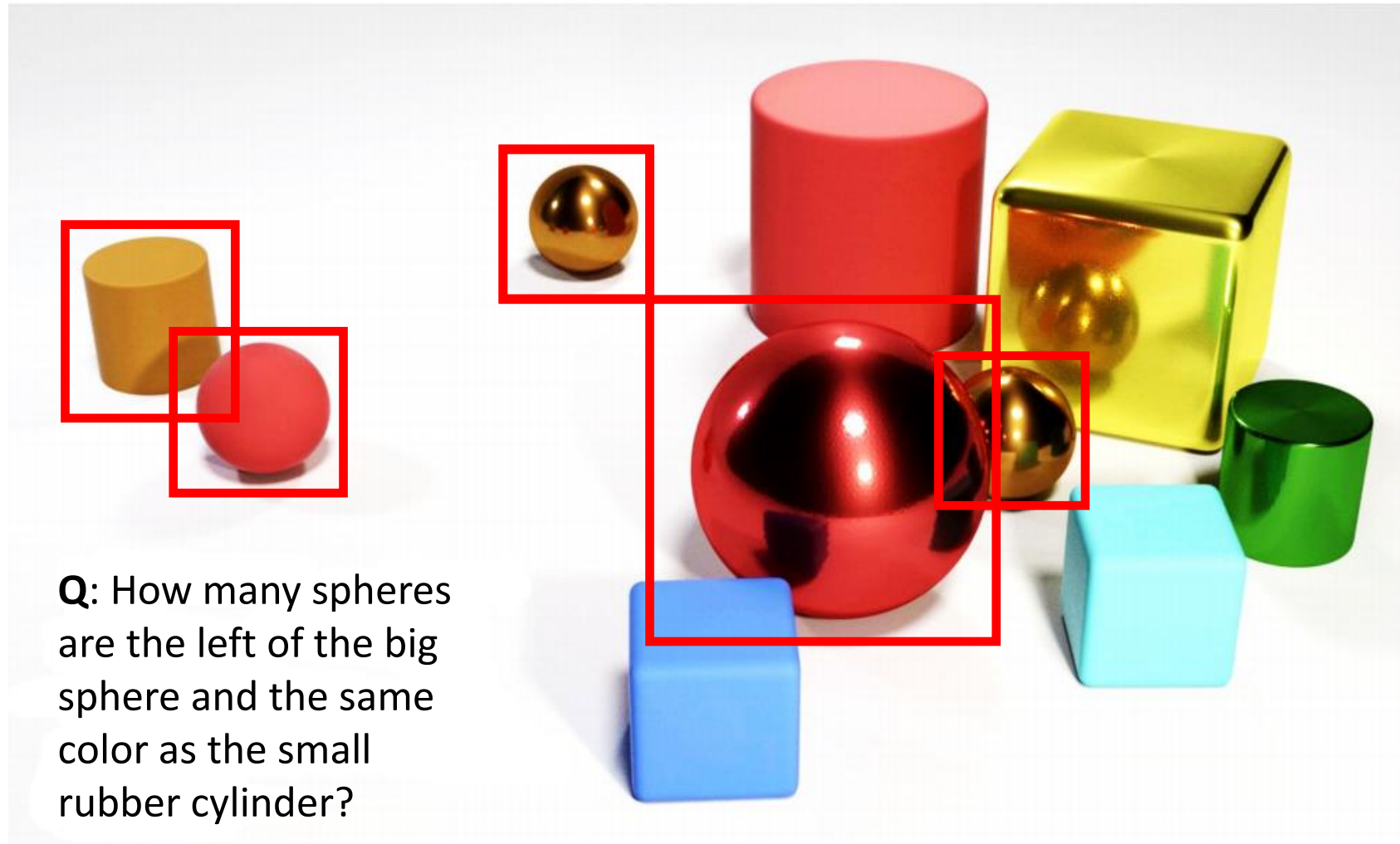
Neural Module Network

- All the previously mentioned work can be considered as [*Monolithic Network*](#)
- Design [*Neural Modules*](#) for compositional visual reasoning



- [1] Deep Compositional Question Answering with Neural Module Networks, CVPR, 2016
- [2] Learning to Reason: End-to-End Module Networks for Visual Question Answering, ICCV 2017
- [3] Inferring and Executing Programs for Visual Reasoning, ICCV 2017
- [4] Transparency by Design: Closing the Gap Between Performance and Interpretability in Visual Reasoning, CVPR 2018
- [5] Explainable Neural Computation via Stack Neural Module Networks, ECCV 2018
- [6] Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding, NeurIPS 2018
- [7] Probabilistic Neural-symbolic Models for Interpretable Visual Question Answering, ICML 2019
- [8] Meta Module Network for Compositional Visual Reasoning, 2019

Compositional Visual Reasoning



Identify big sphere
↓
Spheres on left
↓
Rubber cylinder
↓
Sphere of same color
↓
Count
A: 1

Consider a compositional model

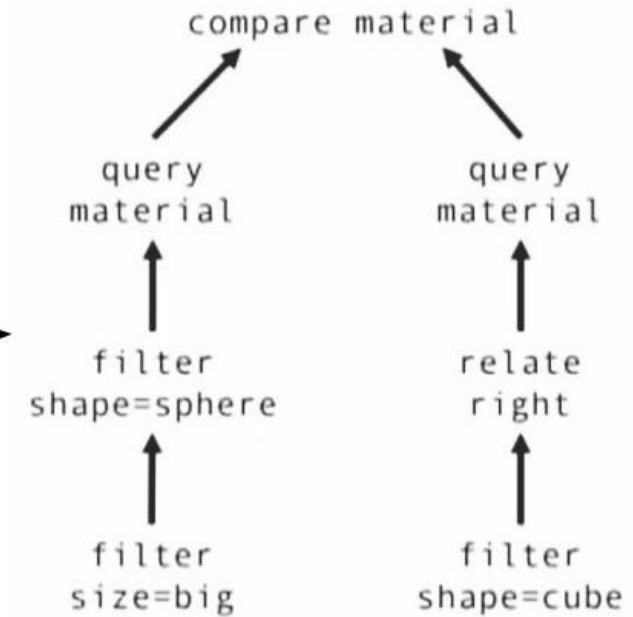
Q: How many spheres are the left of the big sphere and the same color as the small rubber cylinder?

Q: How many spheres are the right of the big sphere and the same color as the small rubber cylinder?

Q: Is the big sphere the same material as the thing on the right of the cube?

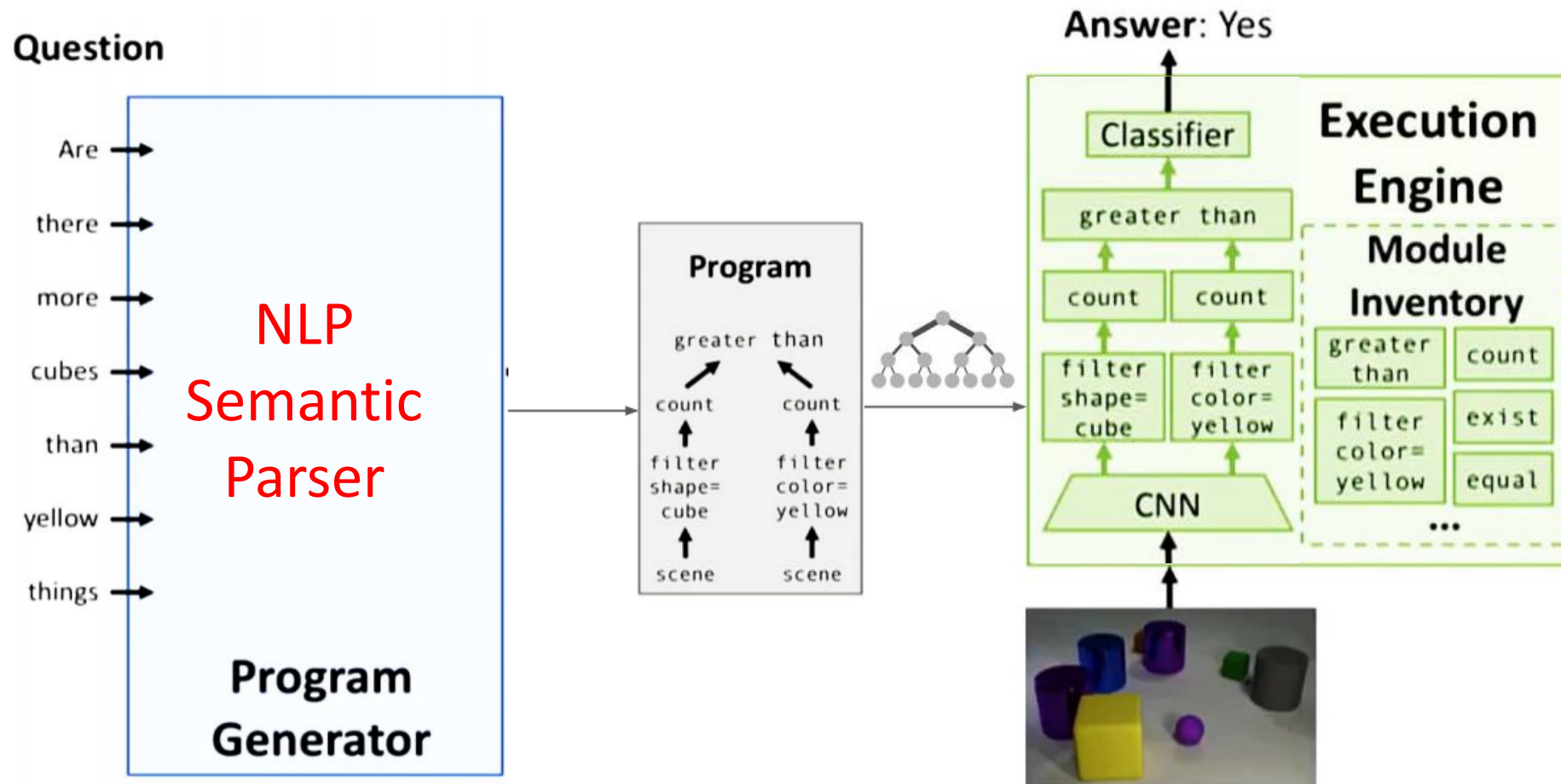
Common operations

Attributes identification
Counting objects
Comparisons
Spatial relationships
Logical operations

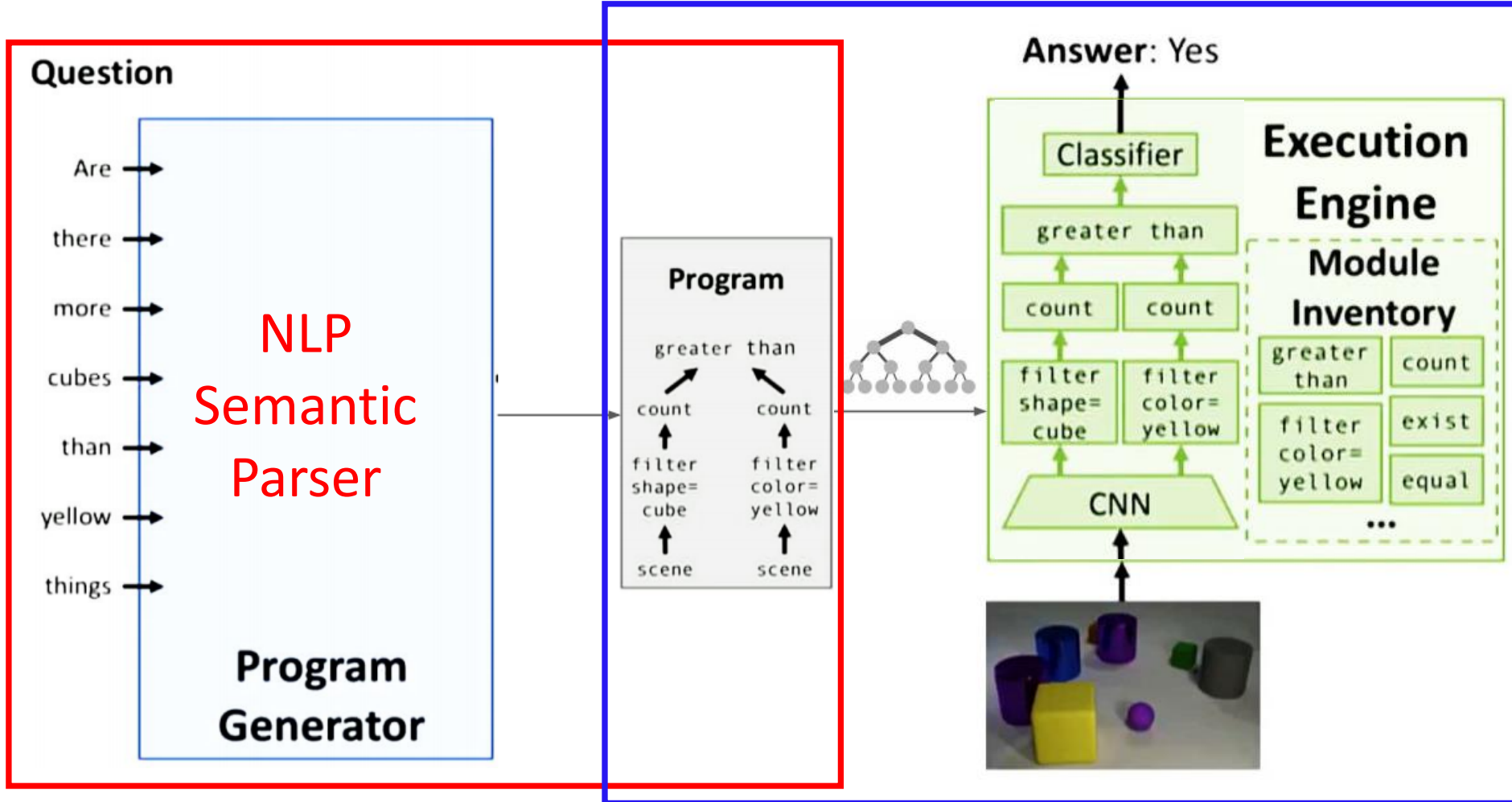


**Network architecture
corresponding to the
third question**

Overview of the NMN approach



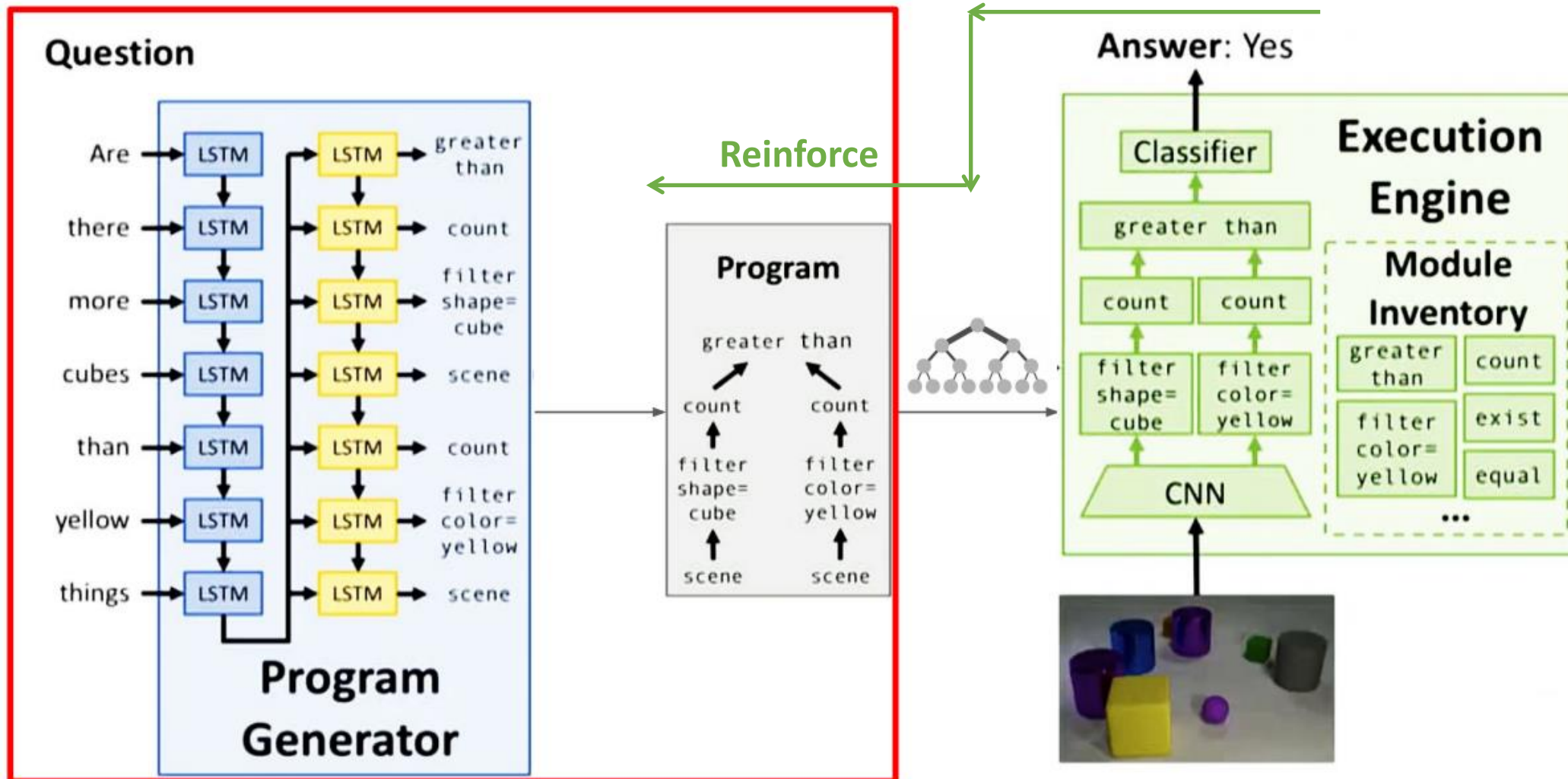
Overview of the NMN approach



Uses some pre-trained parser

Trained separately

Inferring and Executing Programs

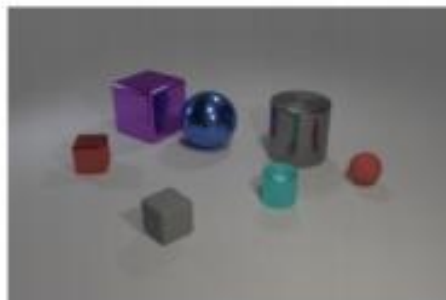


What do the modules learn?

Q: What shape is the...

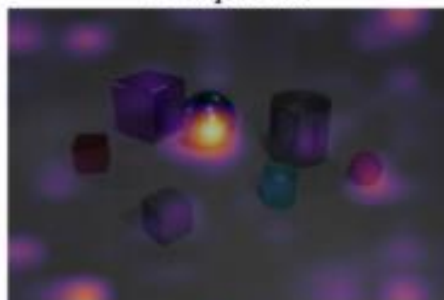
...purple thing?

A: cube



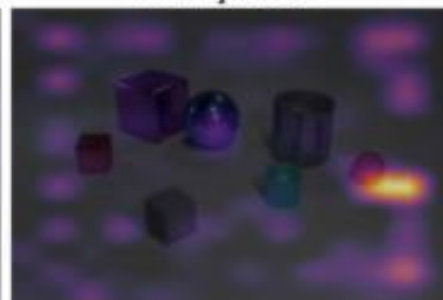
...blue thing?

A: sphere



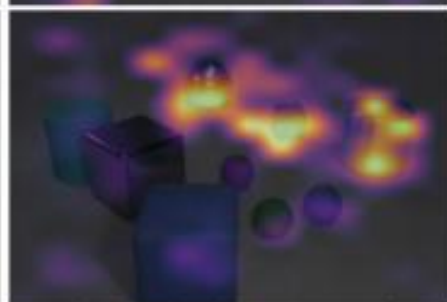
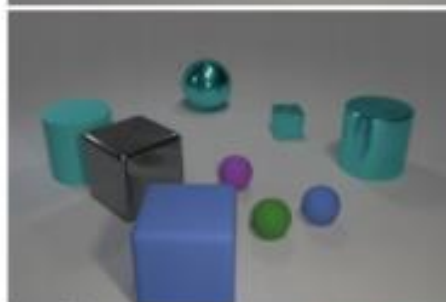
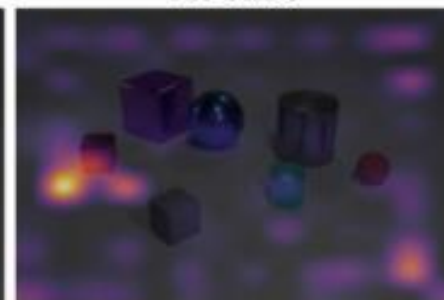
...red thing right of
the blue thing?

A: sphere



...red thing left of
the blue thing?

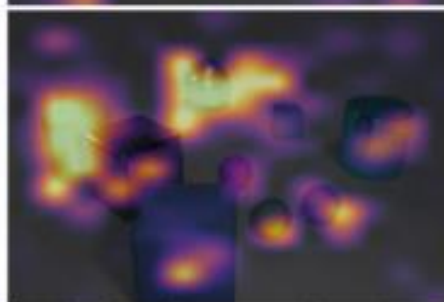
A: cube



Q: How many cyan
things are...

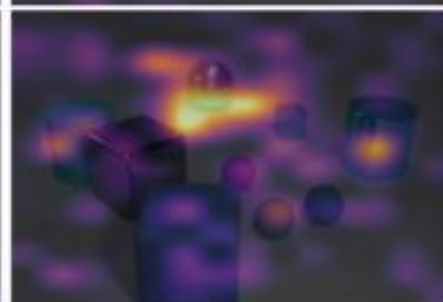
...right of the gray cube?

A: 3



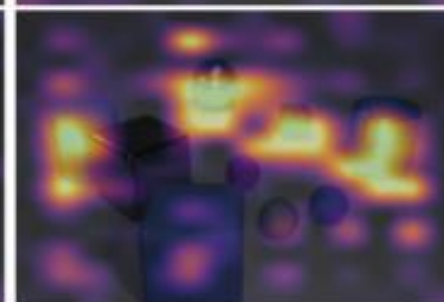
...left of the small cube?

A: 2



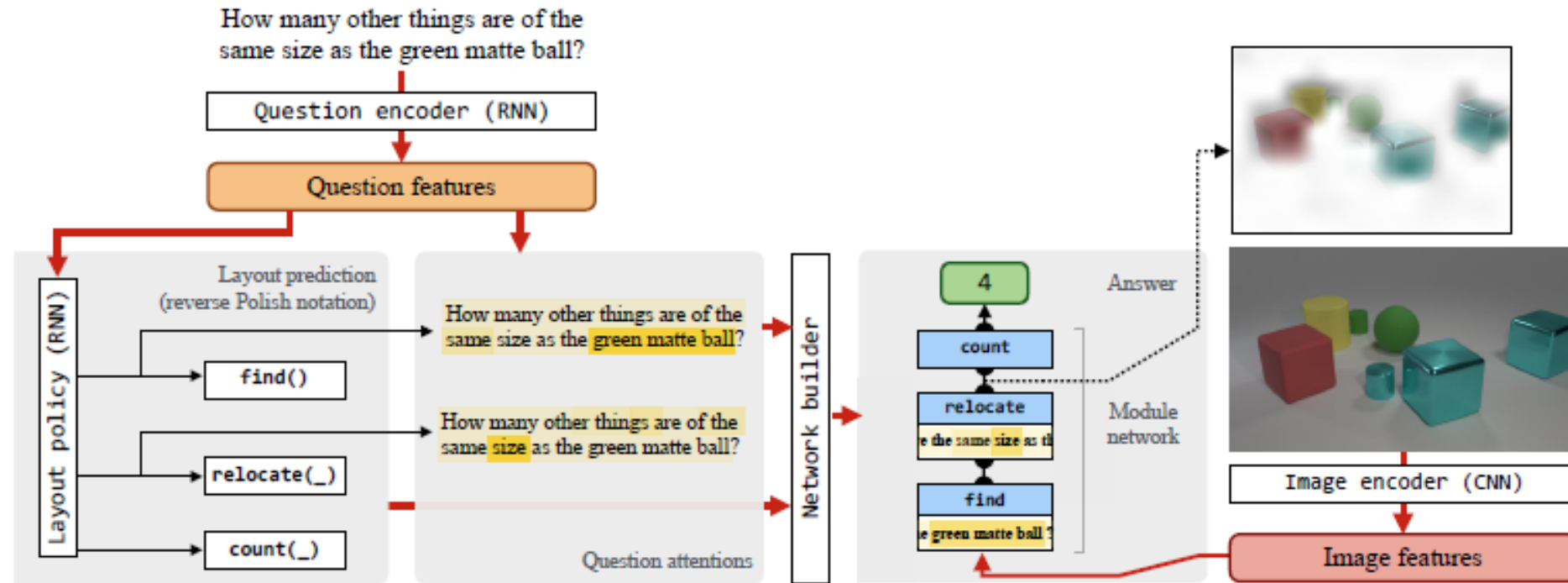
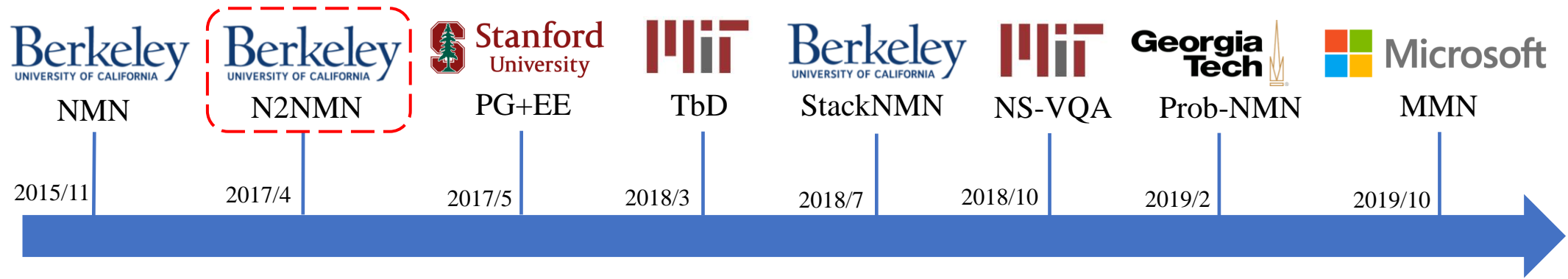
...right of the gray cube
and left of the small cube?

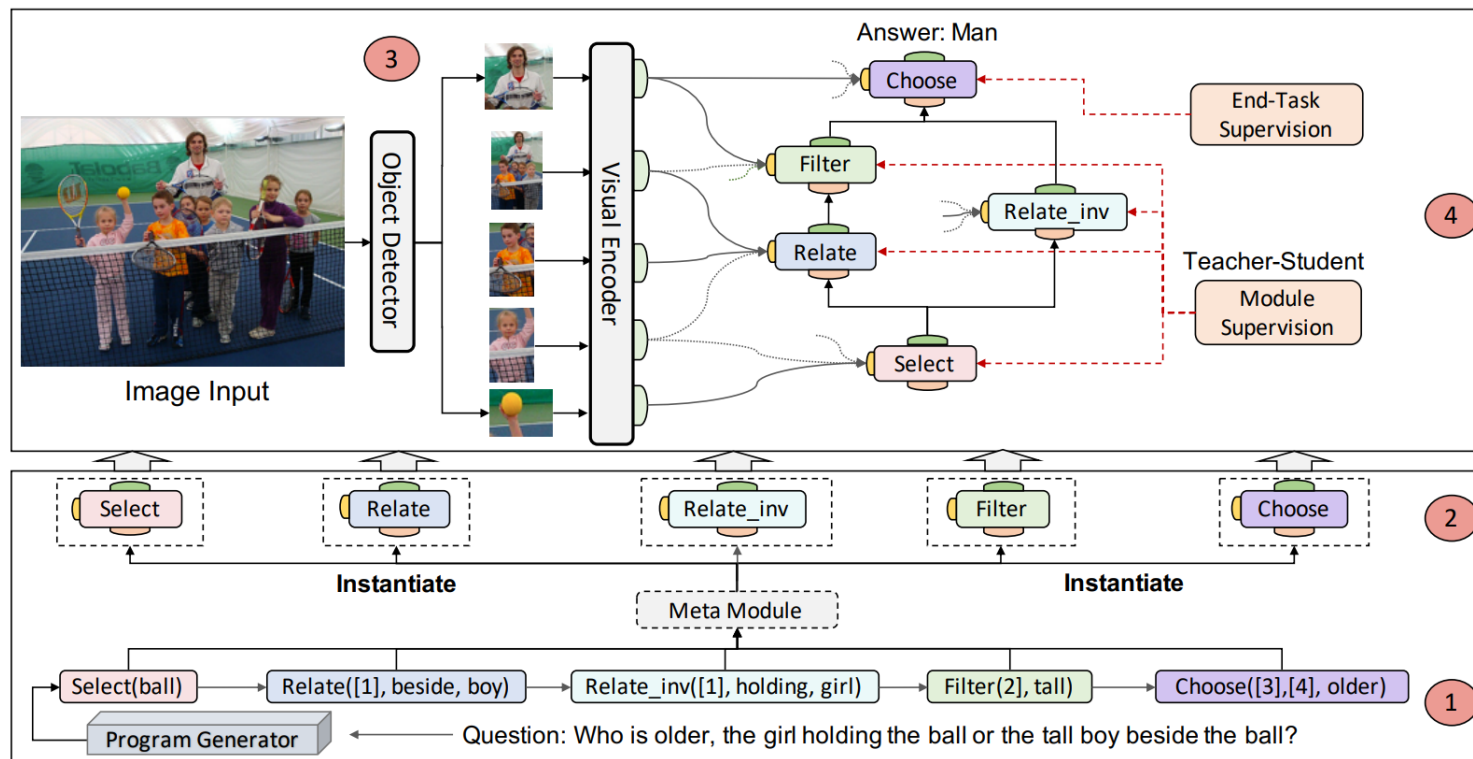
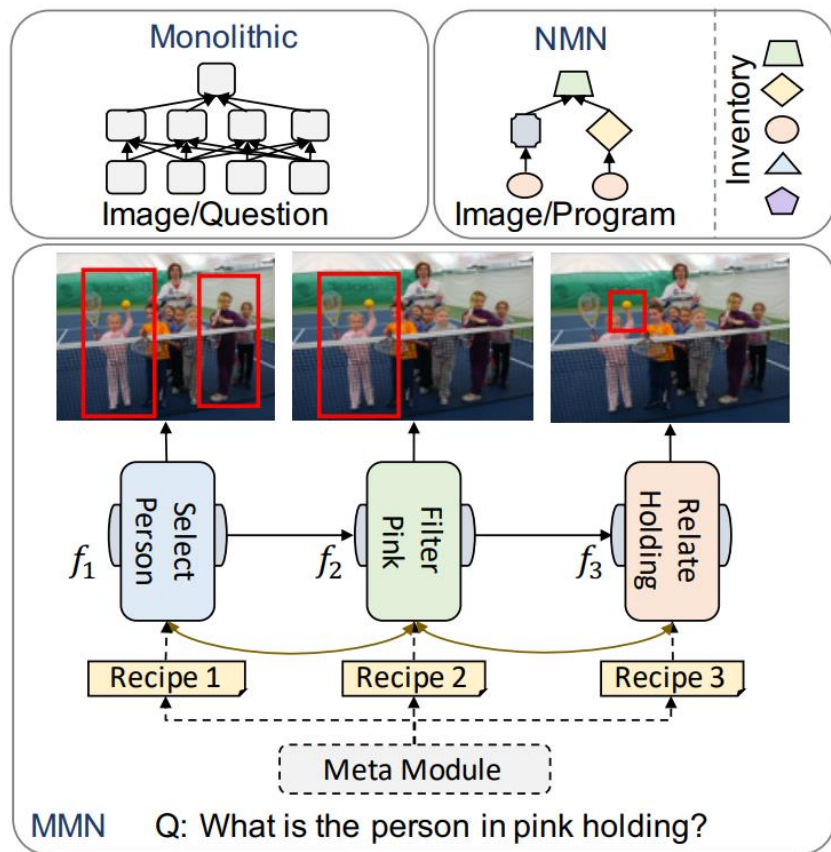
A: 1



...right of the gray cube
or left of the small cube?

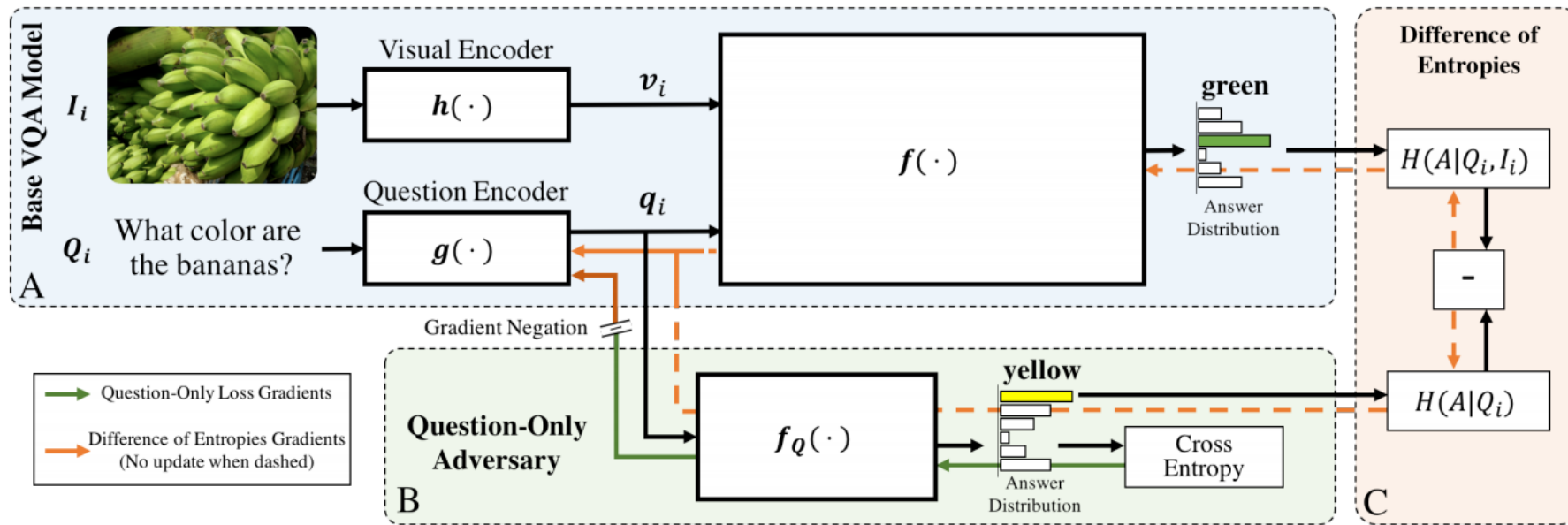
A: 4





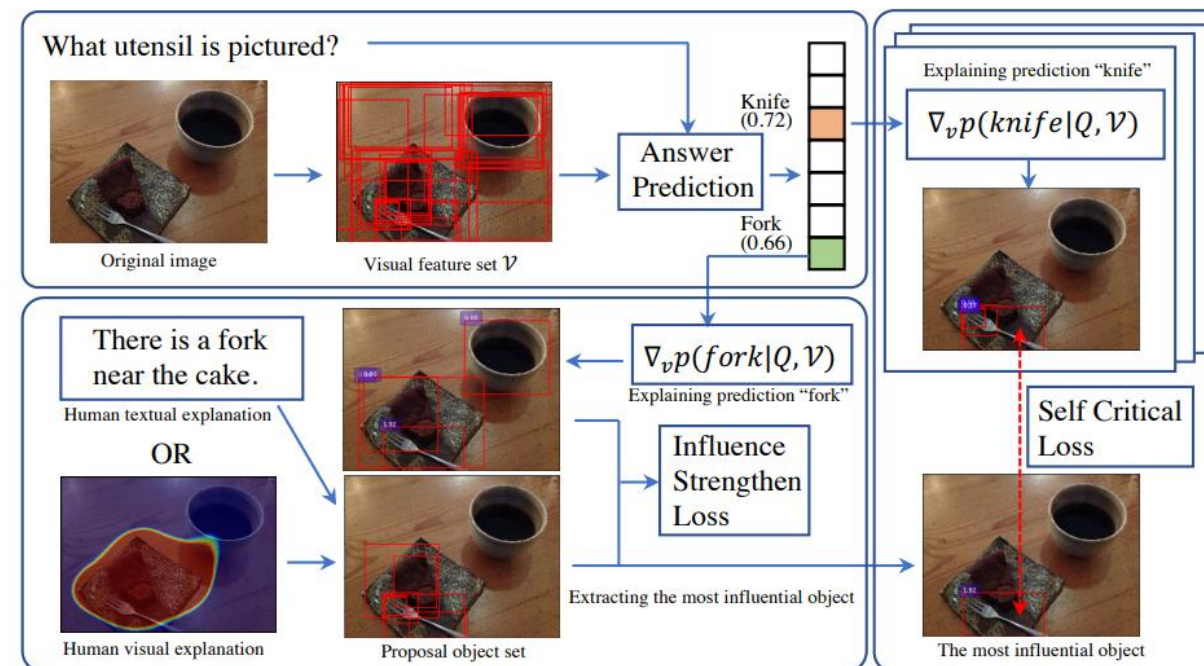
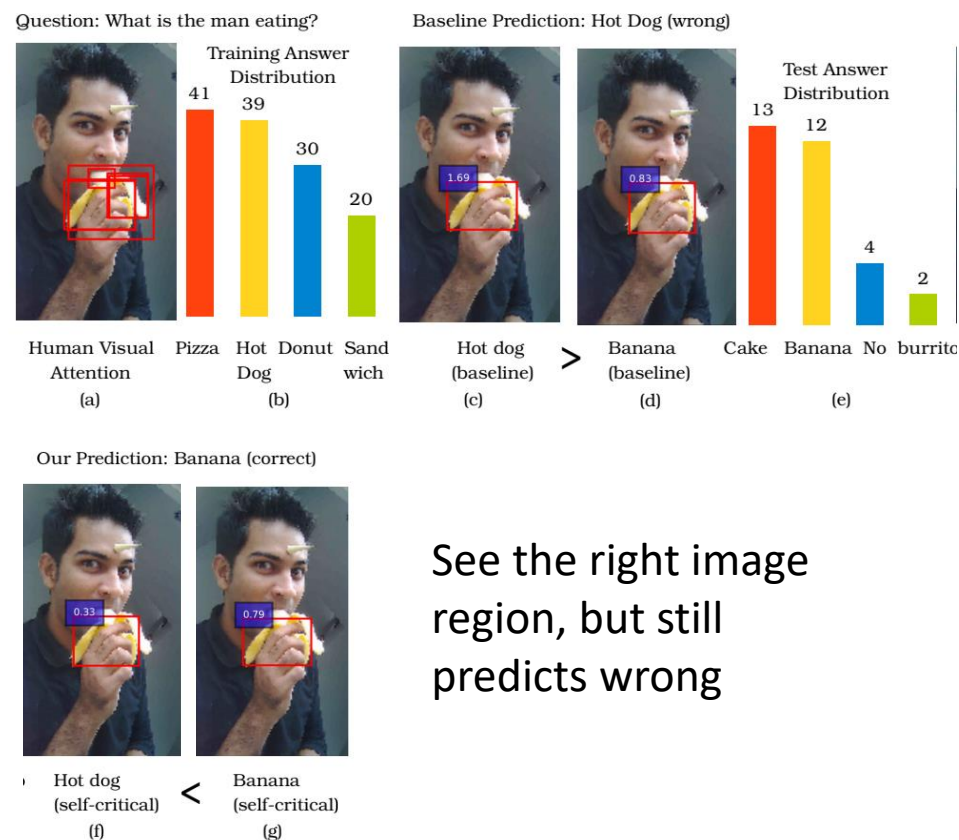
Robust VQA: two examples

- Overcoming language prior with adversarial regularization



Robust VQA: two examples

- Self-critical reasoning



Agenda

- Task Overview

- *What are the main tasks that are driving progress in V+L representation learning?*

- Method Overview

- *What are the state-of-the-art approaches and the key model design principles underlying these methods?*

- Summary

- *What are the core challenges and future directions?*

Take-away Messages

- Popular tasks:
 - VQA, GQA, VCR, RefCOCO, NLVR2, etc.
- Methods:
 - Grid vs. region features
 - Bilinear pooling and FiLM
 - Multimodal alignment with cross-modal attention
 - Relational reasoning with intra-modal attention (self-attention, graph attention)
 - Transformer model becomes popular in the field
 - Multi-step reasoning
 - Neural state machine
 - Neural module network

Challenges & Future Directions

- Can we have something like GLUE and SuperGLUE?
- Can we use a Visual Transformer to encode images to train a large V+L Transformer model end-to-end?
- Instead of Transformer, can we perform FiLM-like fusion for multi-modal pre-training?
- Since all the reasoning is performed in the embedding/neural space, it is not clear whether the model “truly” learns how to reason
- Adversarial robustness of V+L models is less explored in the current literature

Thank you!
Any Questions?