

Spring 2021

ADVANCED TOPICS IN COMPUTER VISION

Atlas Wang

Assistant Professor, The University of Texas at Austin

Visual Informatics Group@UT Austin

<https://vita-group.github.io/>

What is Visual Recognition?

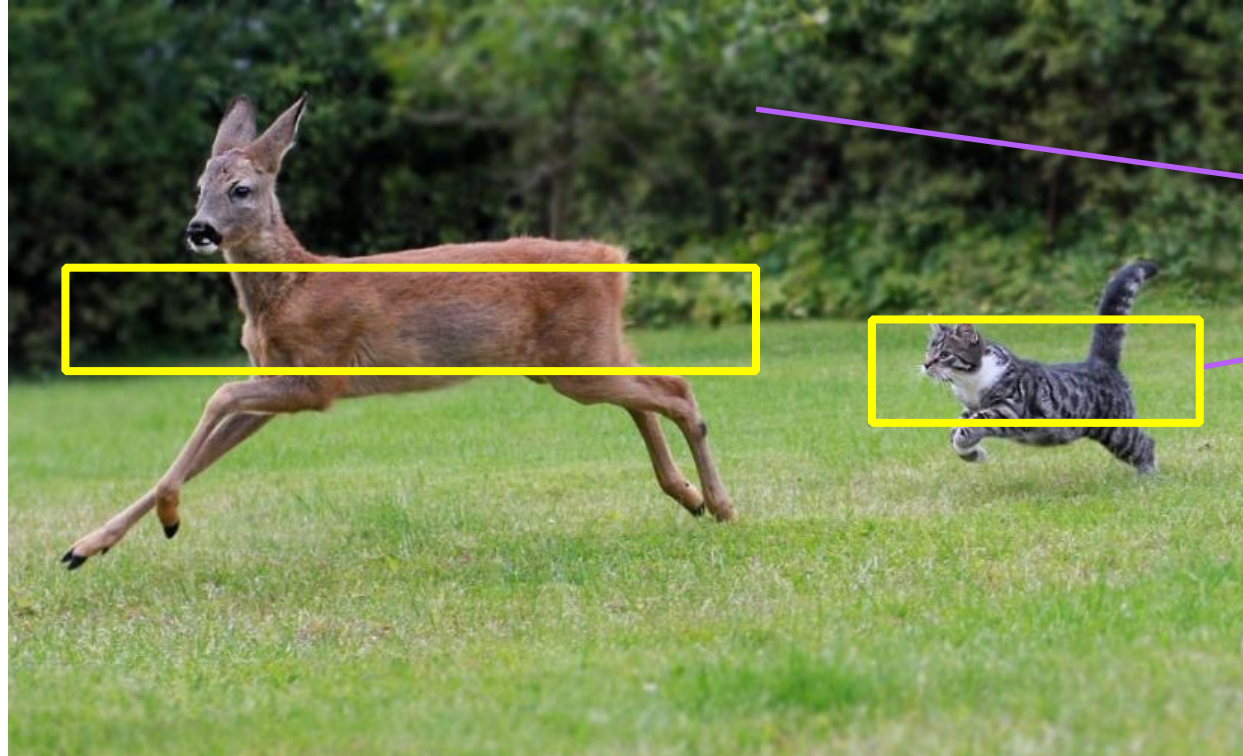
Image tagging



deer
cat
trees
grass

What is Visual Recognition?

Object detection



deer

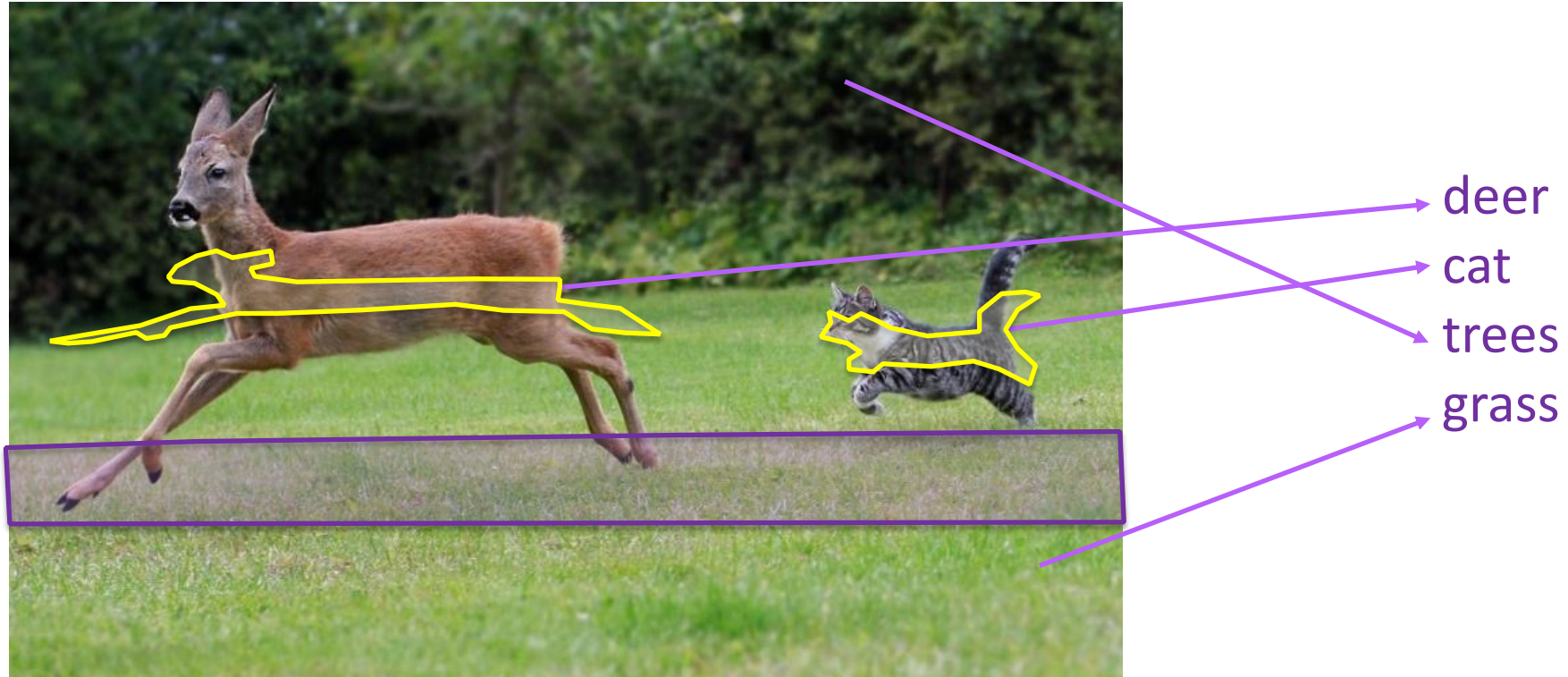
cat

trees

grass

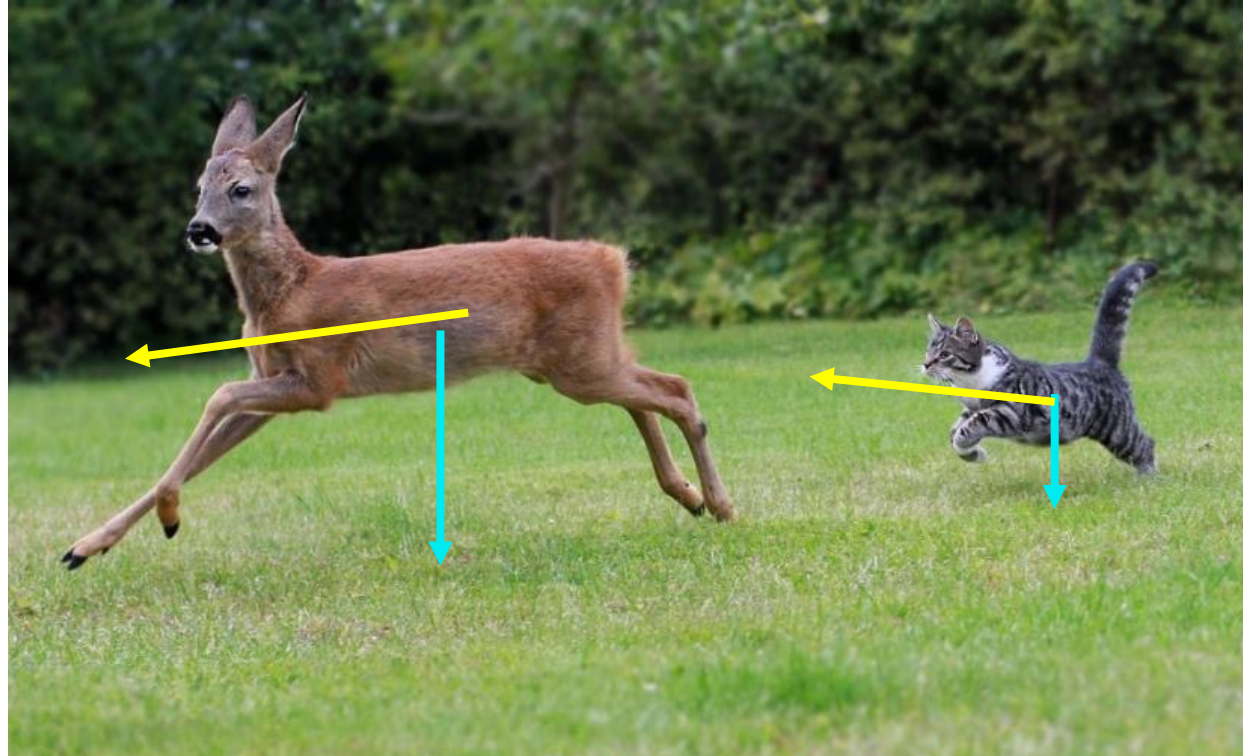
What is Visual Recognition?

Object segmentation



What is Visual Recognition?

Physics / Intuition



↓ Gravity
← Velocity

Pushing the Limits of Visual Recognition

Reasoning about Language!

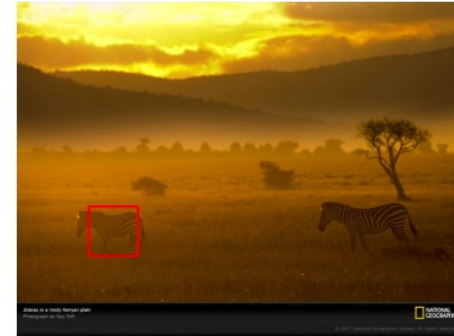


a cat is chasing a
young deer

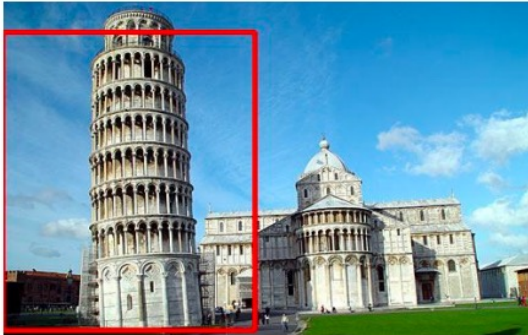
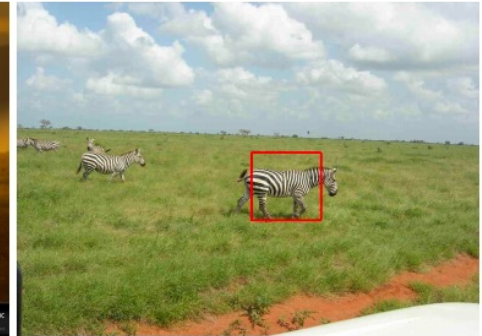
Vision + Language



Helicopter is found in **Airfield**



Zebra is found in **Savanna**



Leaning tower is found in **Pisa**

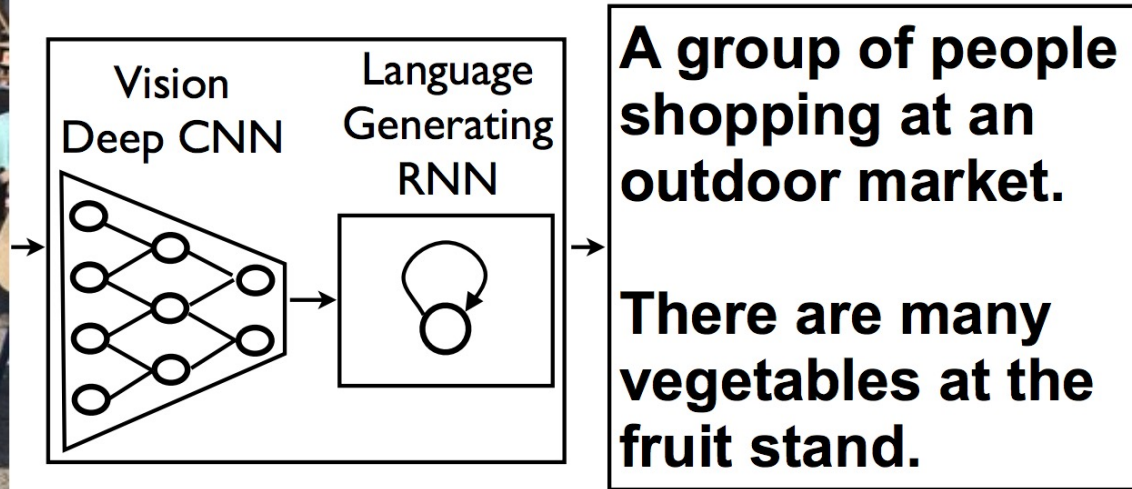


Opera house is found in **Sydney**



Knowledge from Images and Text: Chen et al. 2013

Vision + Language: Applications (1)



Visual Captioning: Vinyals et al. 2015

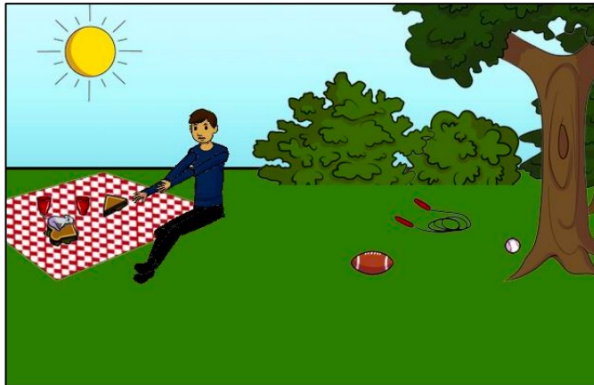
Vision + Language: : Applications (2)



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

Visual Question Answering: Agrawal et al. 2015

Vision + Language : Applications (3)

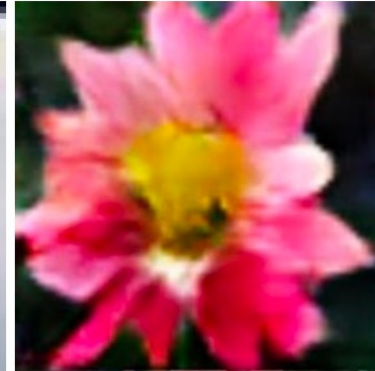
This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face



This bird is white with some black on its head and wings, and has a long orange beak



This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments



Text to Images: Zhang et al. 2016

Problem Overview (1): Visual Captioning

- Describe the content of an image or video with a natural language sentence.



A cat is sitting next to a pine tree, looking up.



A dog is playing piano with a girl.

Applications of Visual Captioning

- Alt-text generation (from PowerPoint)
- Content-based image retrieval (CBIR)
- Helping the visually impaired
- Or just for fun!



Alt Text: A cat sitting on top of a grass covered field



Image Captioning with CNN-LSTM

- Problem Formulation

$$\theta^* = \arg \max_{\theta} \sum_{(I,S)} \log p(S|I; \theta)$$
$$\log p(S|I) = \sum_{t=0}^N \log p(S_t|I, S_0, \dots, S_{t-1})$$

- The Encoder-Decoder framework



"Show and Tell"

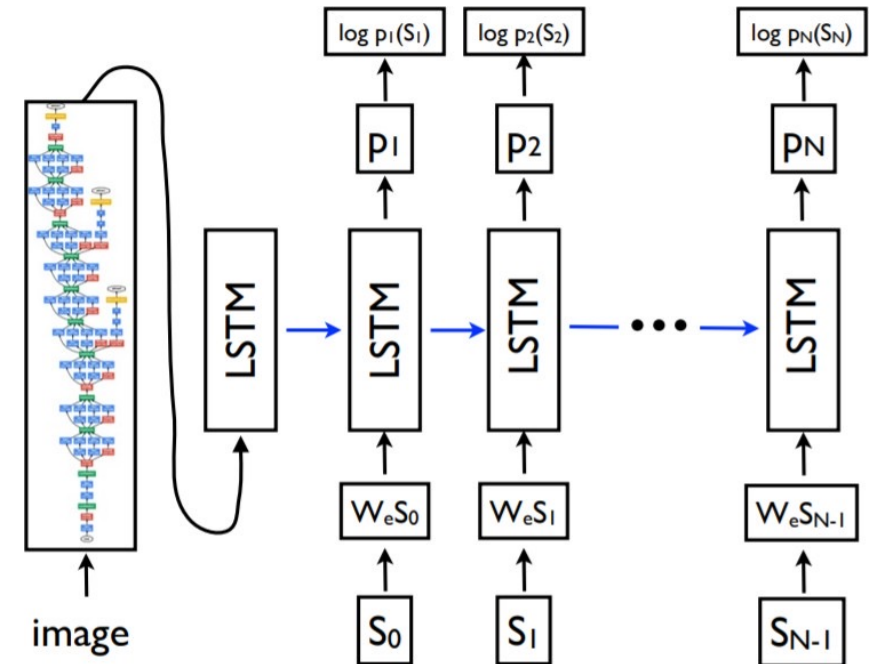
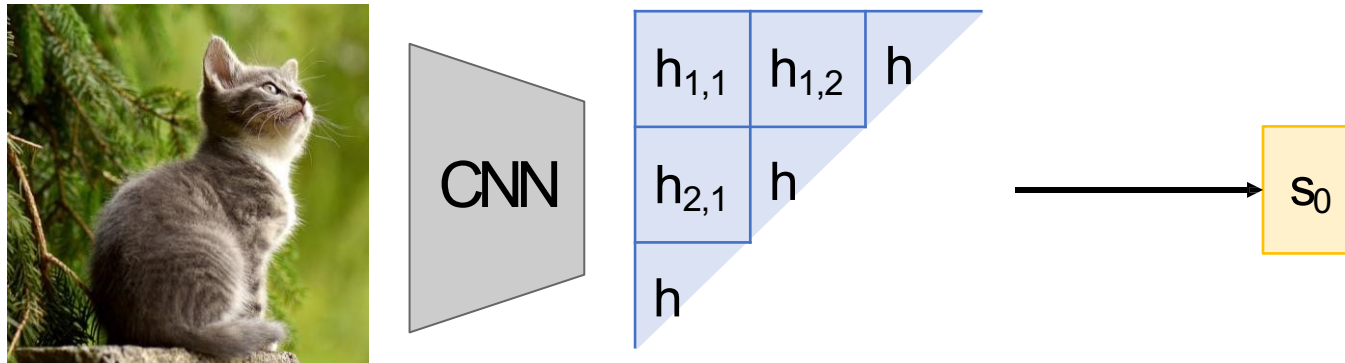


Image Captioning with Soft Attention

- Soft Attention – Dynamically attend to input content based on query.
- Basic elements: query – q , keys – K , and values – V
- In our case, keys and values are usually identical. They come from the CNN activation map.
- Query q is determined by the global image feature or LSTM's hidden states.

Image Captioning with Soft Attention



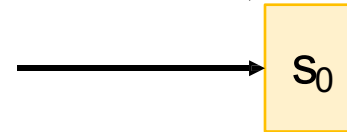
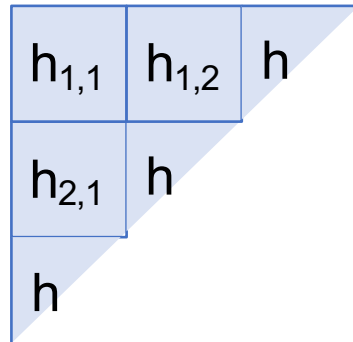
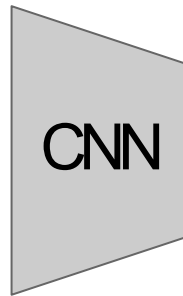
Use a CNN to compute a
grid of features for an image

Image Captioning with Soft Attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

Alignment scores

$e_{1,1,1}$	$e_{1,1,2}$	$e_{1,1,3}$
$e_{1,2,1}$	$e_{1,2,2}$	$e_{1,2,3}$
$e_{1,3,1}$	$e_{1,3,2}$	$e_{1,3,3}$

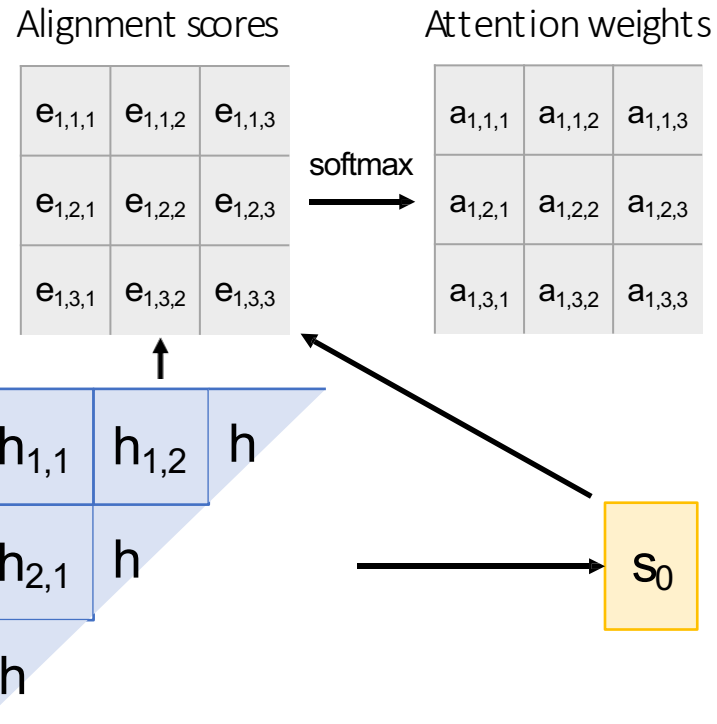


Use a CNN to compute a
grid of features for an image

Image Captioning with Soft Attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

$$a_{t,:} = \text{softmax}(e_{t,:,:})$$



Use a CNN to compute a grid of features for an image

Image Captioning with Soft Attention

$$\begin{aligned}e_{t,i,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\a_{t,:} &= \text{softmax}(e_{t,:,:}) \\c_t &= \sum_{i,j} a_{t,i,j} h_{i,j}\end{aligned}$$



CNN

Use a CNN to compute a grid of features for an image

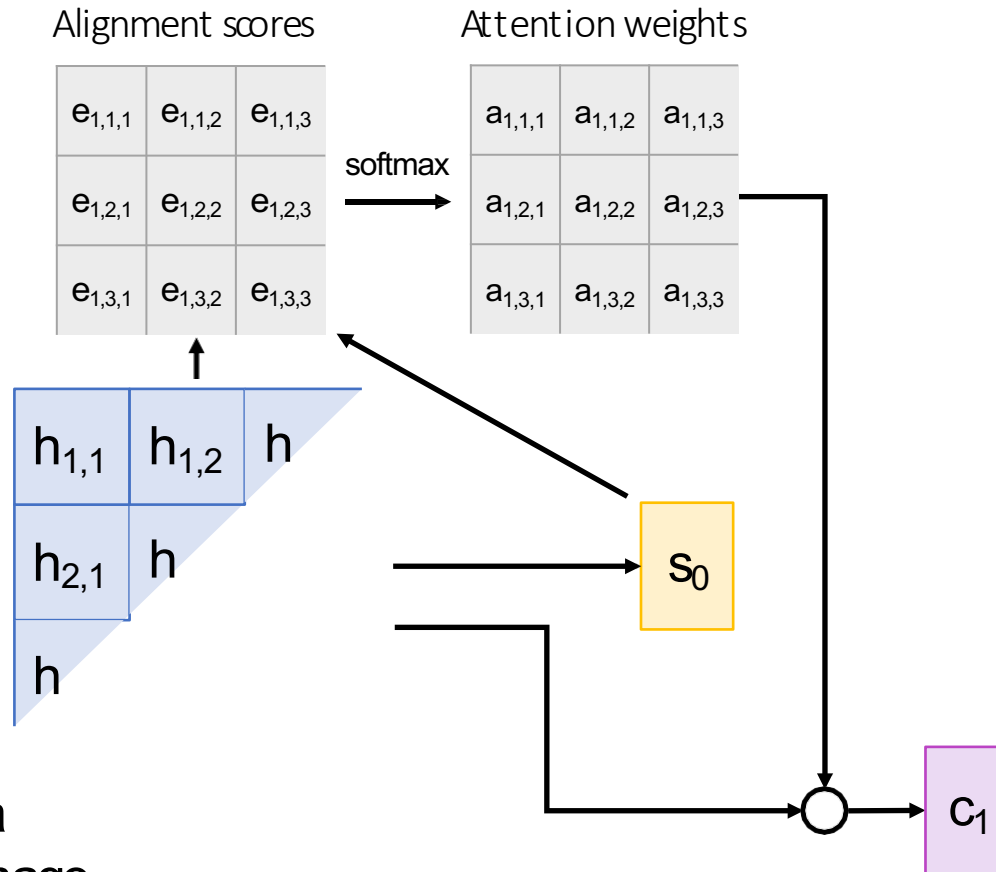


Image Captioning with Soft Attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

$$a_{t,:} = \text{softmax}(e_{t,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



Use a CNN to compute a grid of features for an image

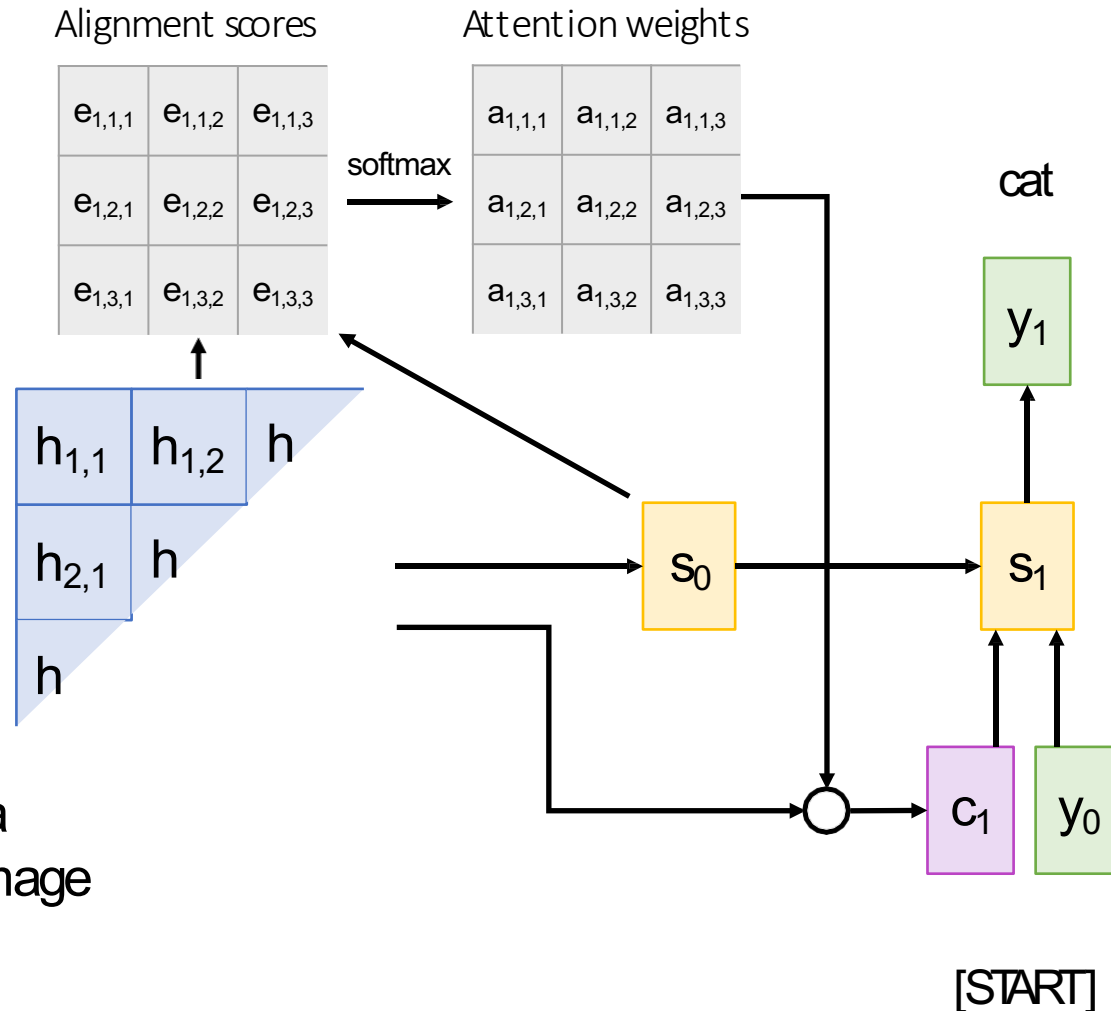


Image Captioning with Soft Attention

$$\begin{aligned} e_{t,i,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\ a_{t,:} &= \text{softmax}(e_{t,:}) \\ c_t &= \sum_{i,j} a_{t,i,j} h_{i,j} \end{aligned}$$

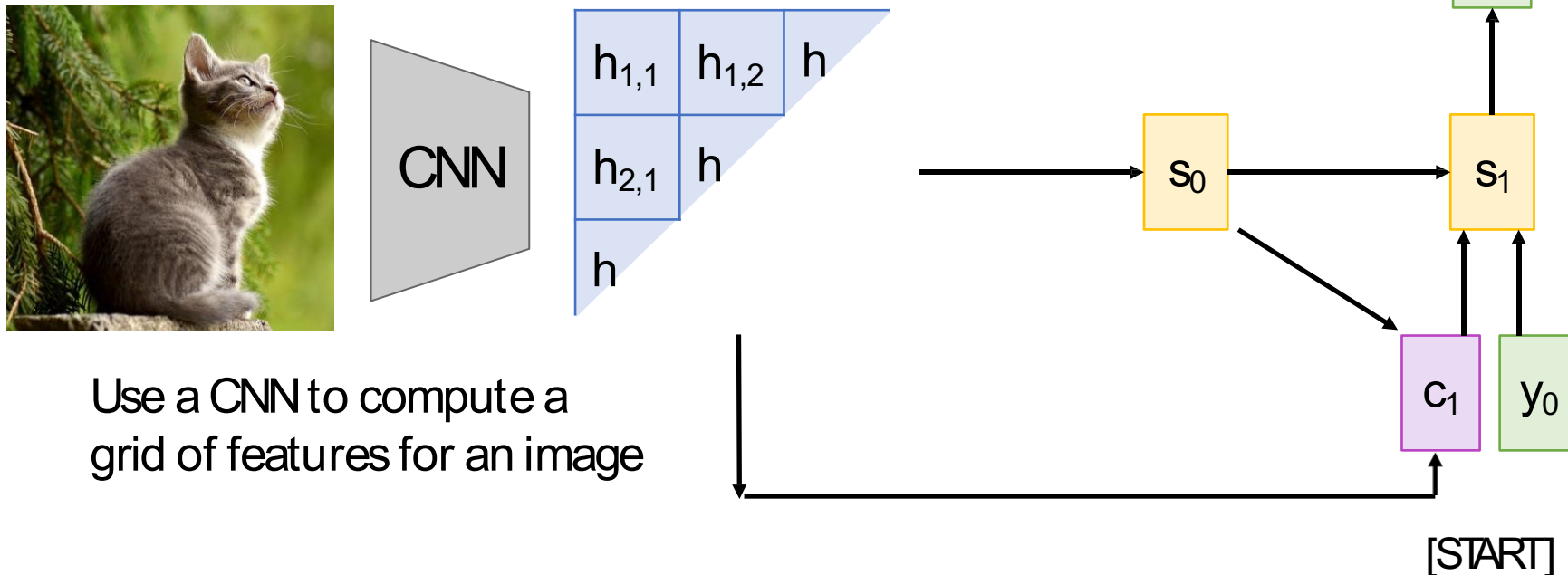


Image Captioning with Soft Attention

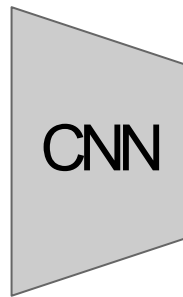
$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

$$a_{t,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$

Alignment scores

$e_{2,1,1}$	$e_{2,1,2}$	$e_{2,1,3}$
$e_{2,2,1}$	$e_{2,2,2}$	$e_{2,2,3}$
$e_{2,3,1}$	$e_{2,3,2}$	$e_{2,3,3}$



Use a CNN to compute a grid of features for an image

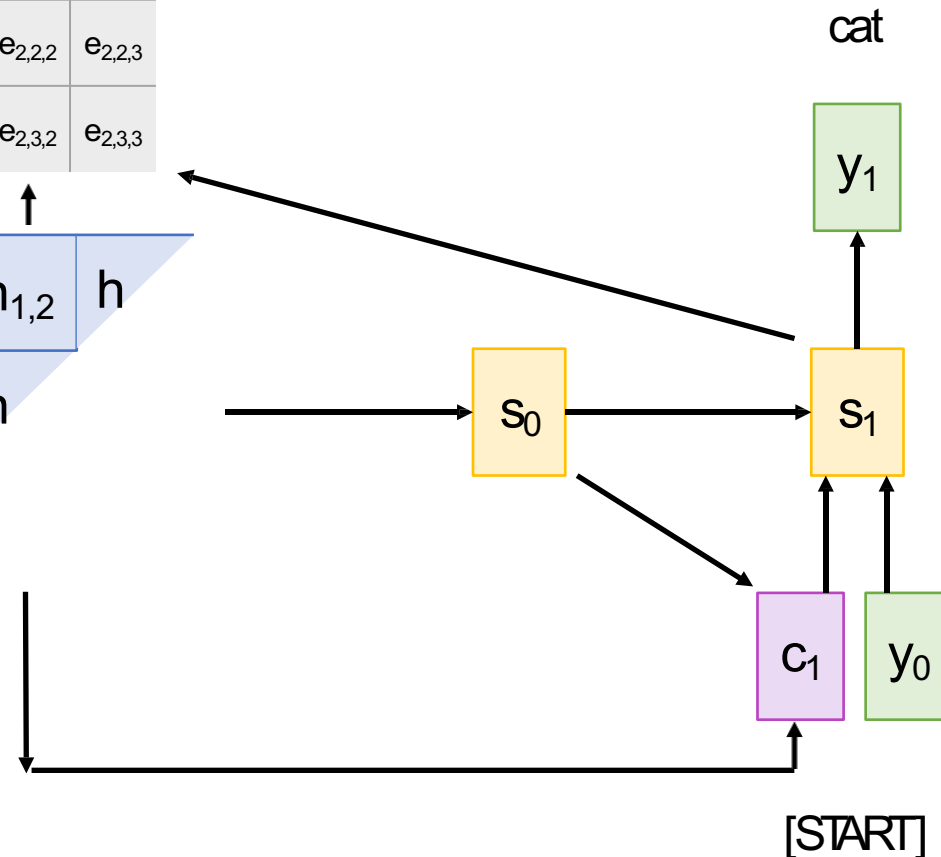
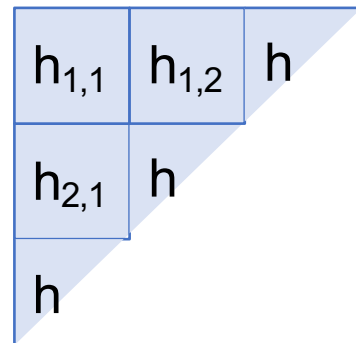
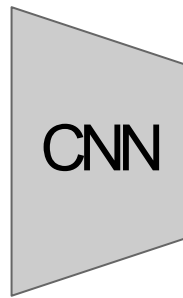


Image Captioning with Soft Attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

$$a_{t,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



Use a CNN to compute a grid of features for an image

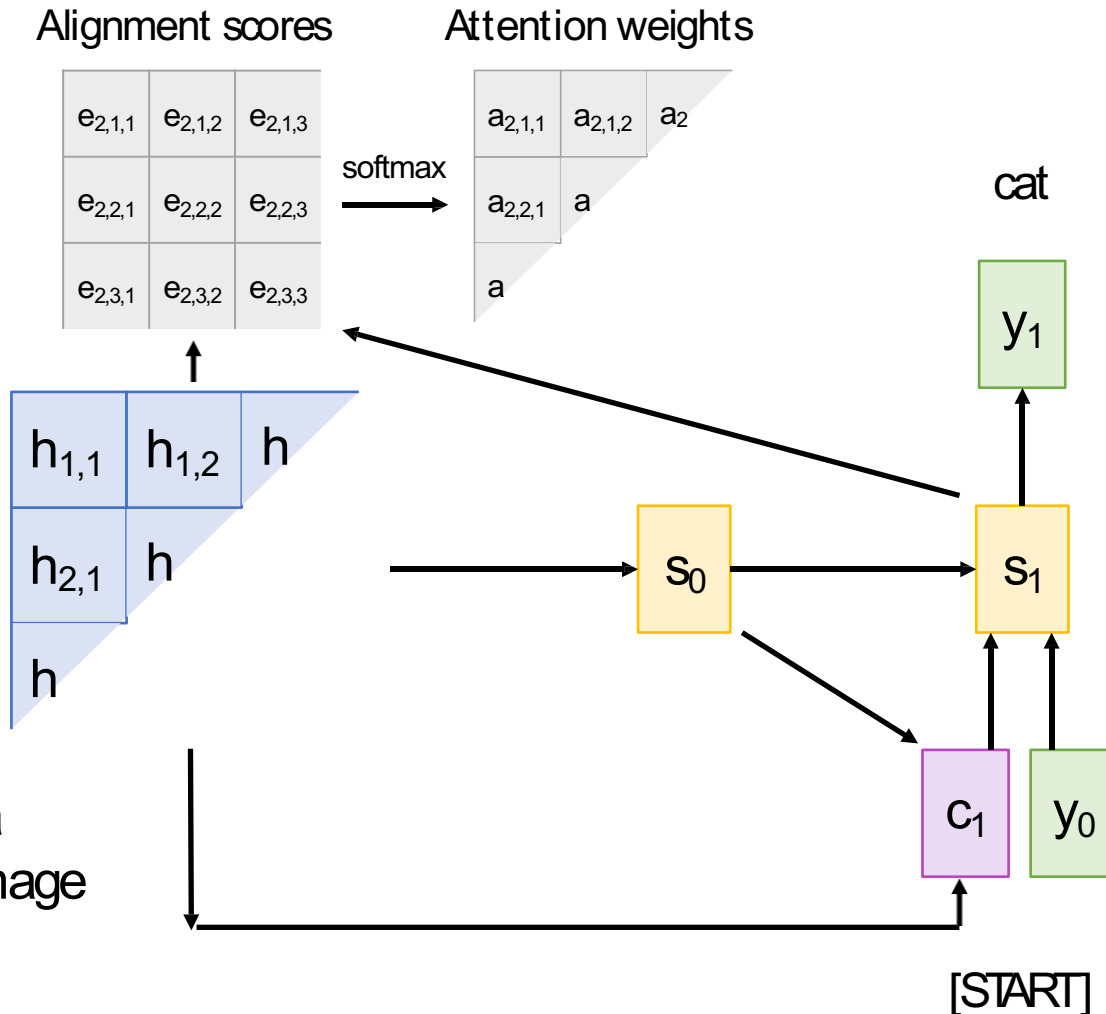


Image Captioning with Soft Attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

$$a_{t,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



Use a CNN to compute a grid of features for an image

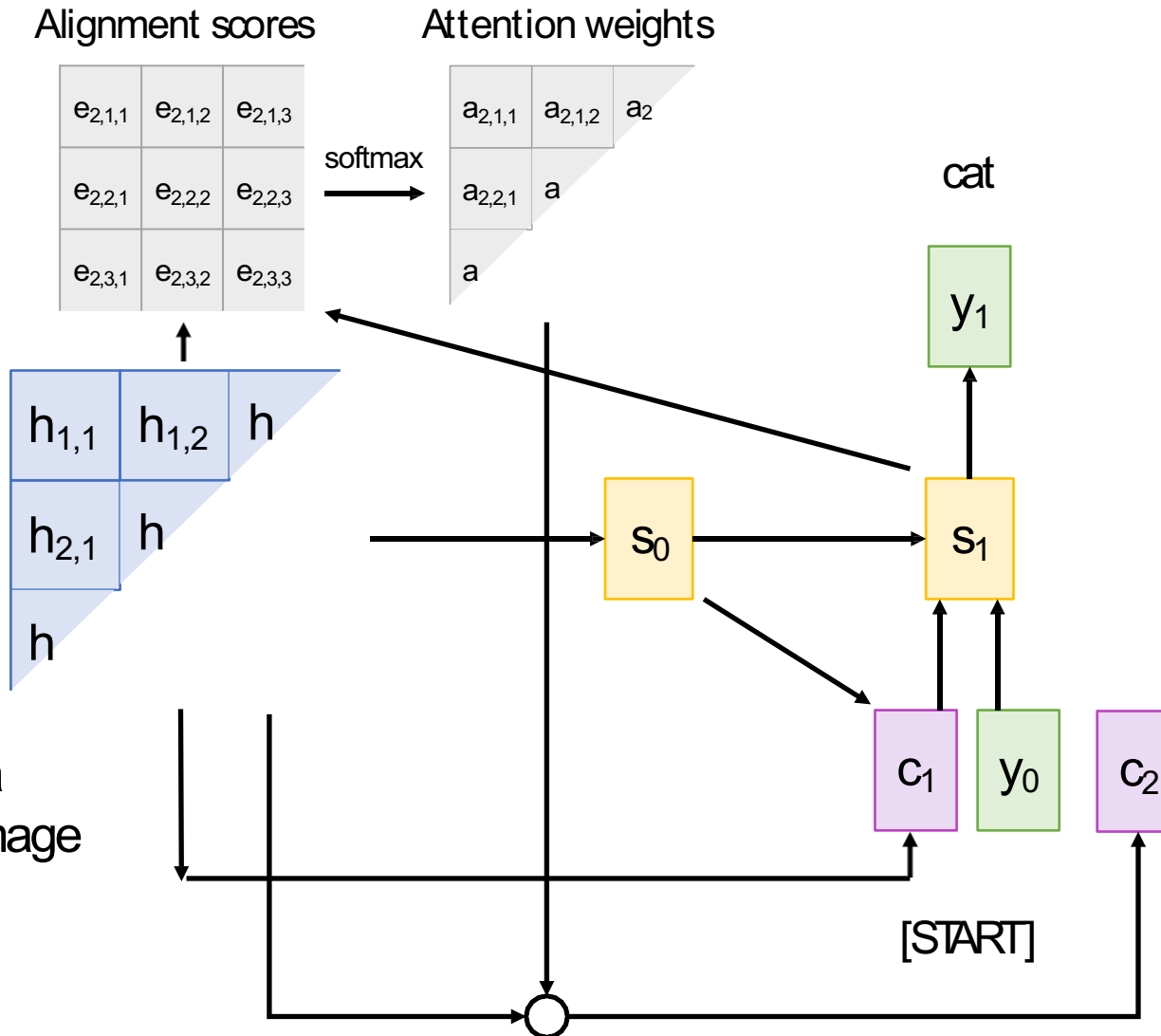


Image Captioning with Soft Attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

$$a_{t,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



Use a CNN to compute a grid of features for an image

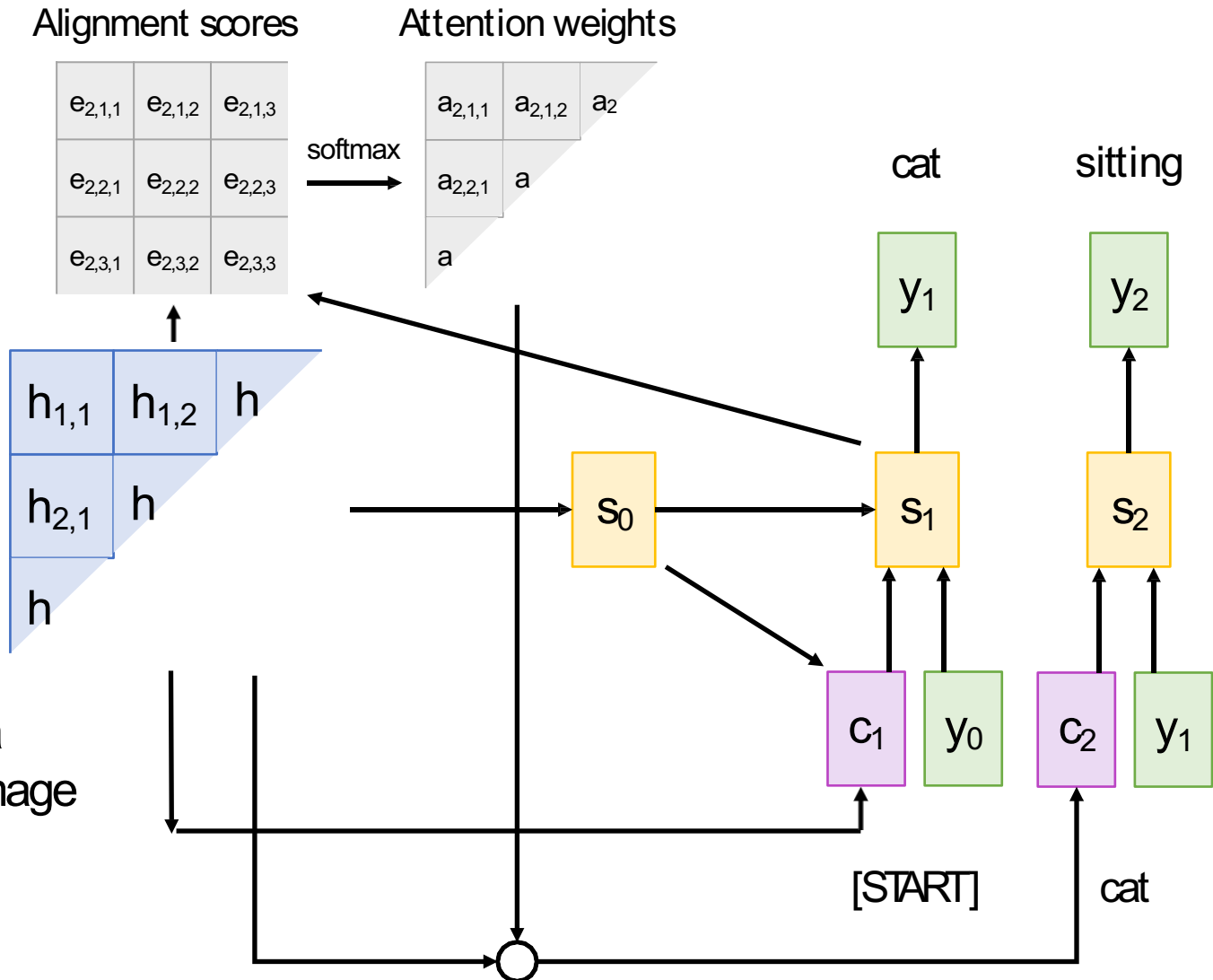
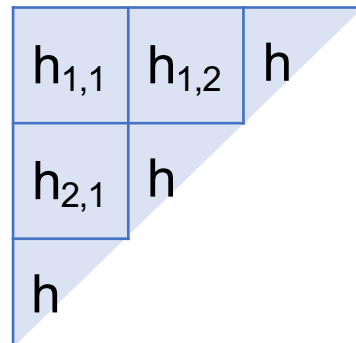
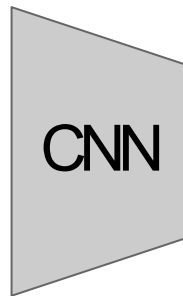


Image Captioning with Soft Attention

$$\begin{aligned}e_{t,i,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\ a_{t,:} &= \text{softmax}(e_{t,:}) \\ c_t &= \sum_{i,j} a_{t,i,j} h_{i,j}\end{aligned}$$

Each timestep of decoder uses a different context vector that looks at different parts of the input image



Use a CNN to compute a grid of features for an image

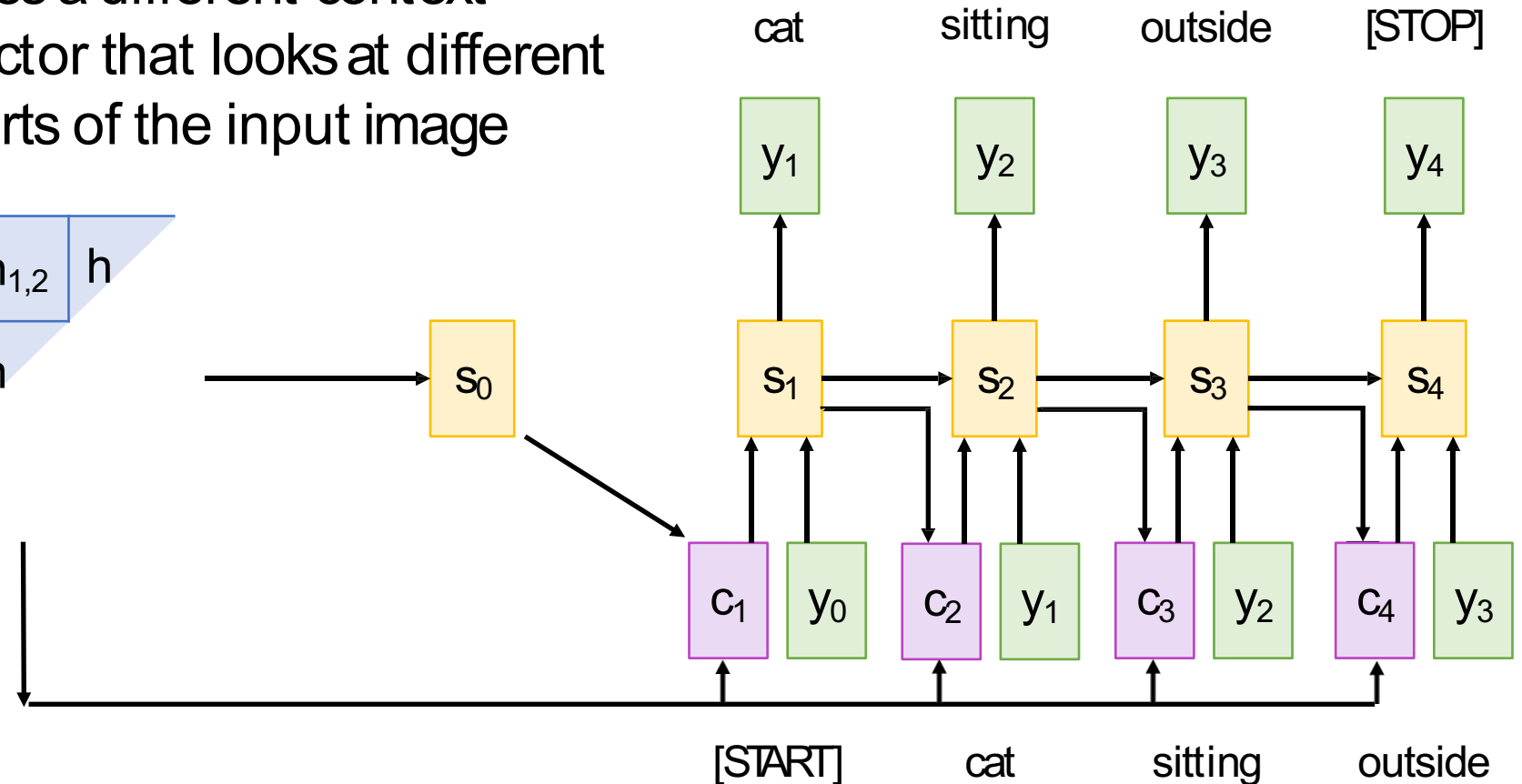


Image Captioning with Soft Attention

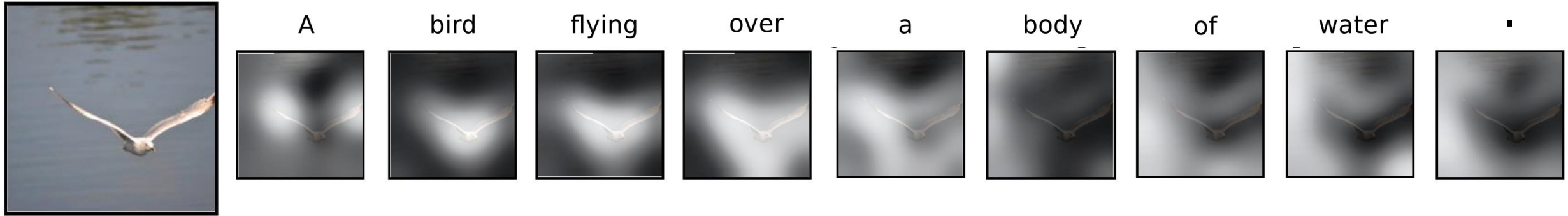


Image Captioning with Soft Attention



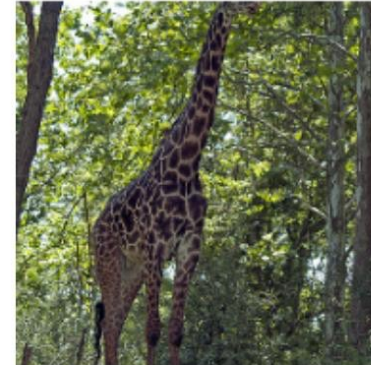
A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Image Captioning with Region Attention

- Variants of Soft Attention based on the feature input
 - Grid activation features (covered)
 - Region proposal features

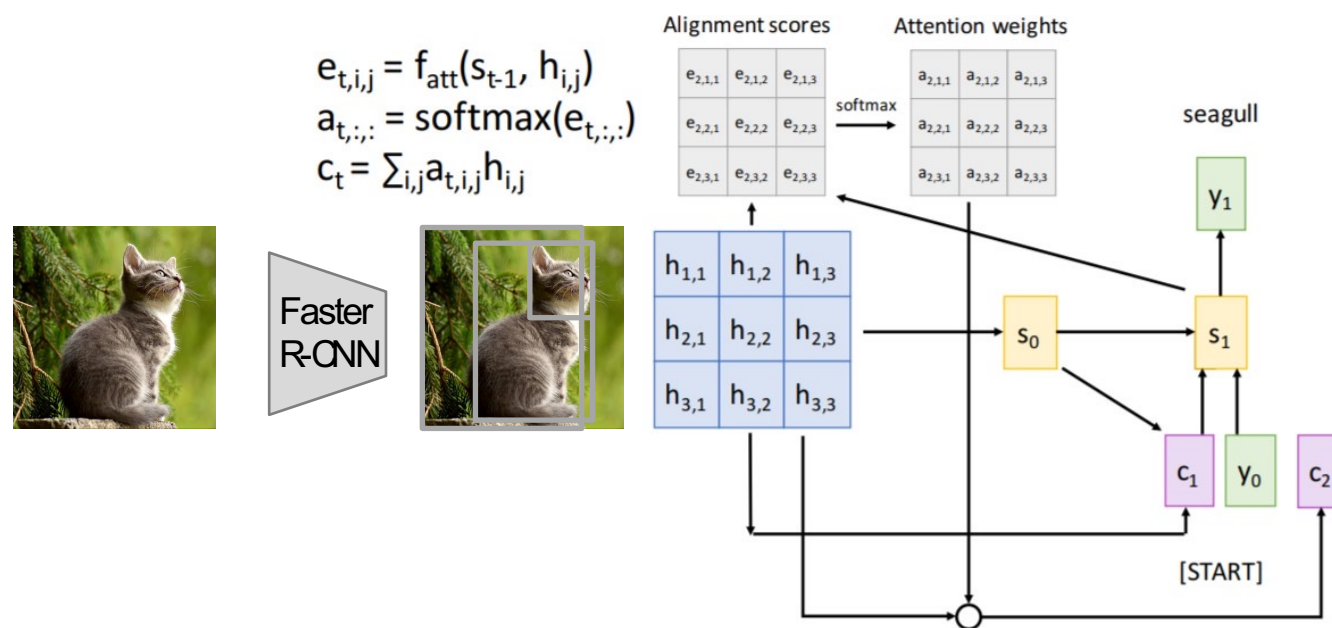


Image Captioning with Transformer

- Transformer performs sequence-to-sequence generation.
- Self-Attention – A type of soft attention that “attends to itself”.
- Self-Attention is a special case of Graph Neural Networks (GNNs) that has a fully-connected graph.
- Self-attention is sometimes used to model relationship between object regions, similar to GCNs.

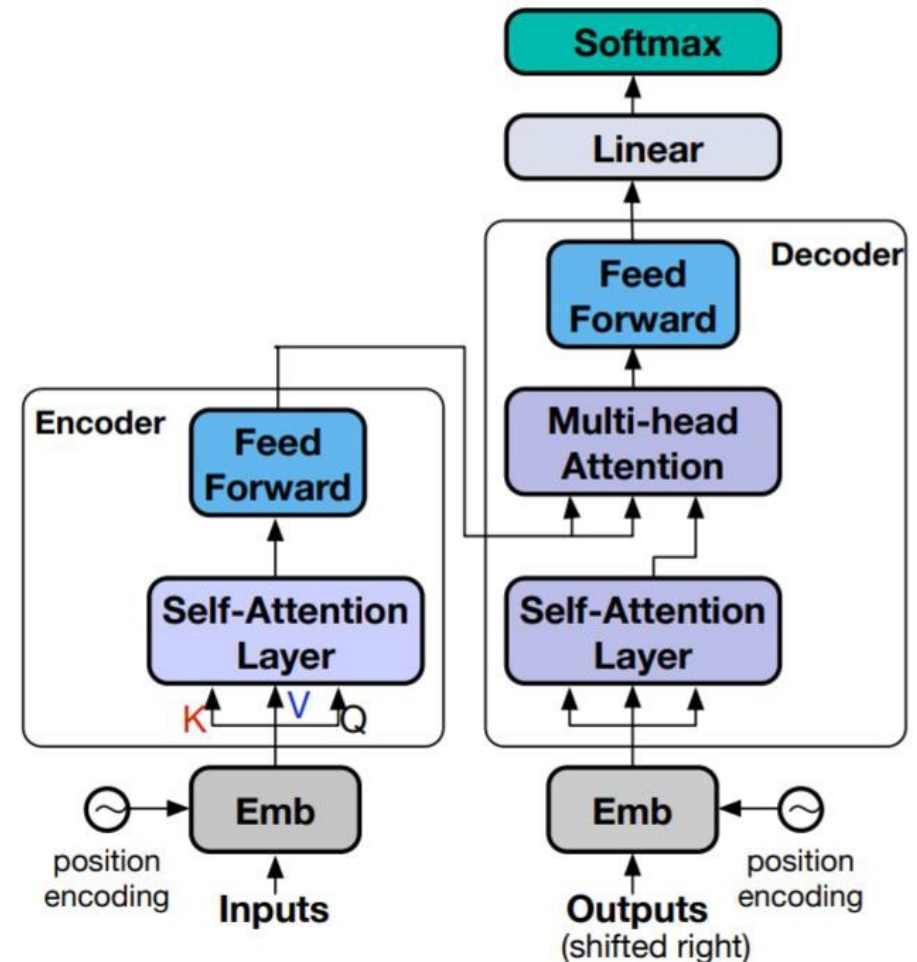
Vaswani et al. “Attention is all you need”, NIPS2017.

Yao et al. “Exploring visual relationship for image captioning”, ECCV 2018.

Further readings: <https://graphdeeplearning.github.io/post/transformers-are-gnns/>

Image Captioning with Transformer

- Transformer is first adapted for captioning in Zhou et al.
- Others: Object Relation Transformer, Meshed-Memory Transformer



Zhou et al. "End-to-end dense video captioning with masked transformer", CVPR2018.

Herdade et al. "Image Captioning: Transforming Objects into Words", NeurIPS 2019.

Comia et al. "Meshed-Memory Transformer for Image Captioning", CVPR2020.

Vision-Language Pre-training (VLP)

- Two-stage training strategy: **pre-training** and **fine-tuning**.
- **Pre-training** is performed on a large dataset. Usually with auto-generated captions. The training objective is *unsupervised*.
- **Fine-tuning** is task-specific *supervised* training on downstream tasks.
- All methods are based on BERT (a variant of Transformer).

VideoBERT: A Joint Model for Video and Language Representation Learning

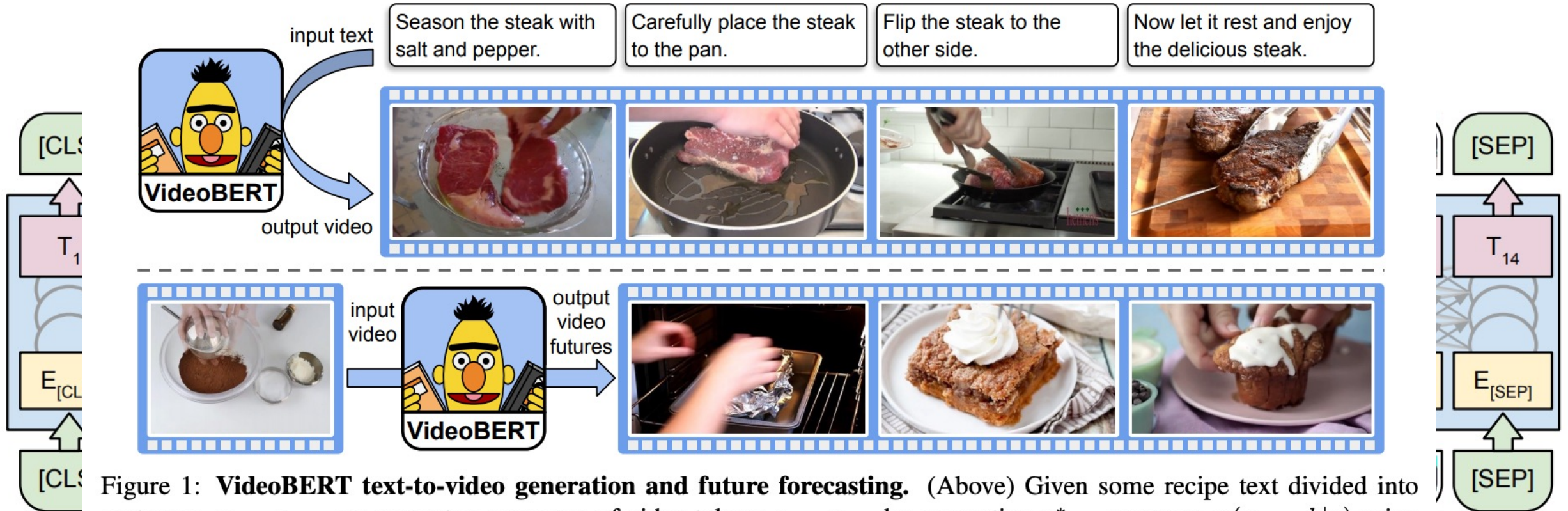


Figure 1: **VideoBERT text-to-video generation and future forecasting.** (Above) Given some recipe text divided into sentences, $y = y_{1:T}$, we generate a sequence of video tokens $x = x_{1:T}$ by computing $x_t^* = \arg \max_k p(x_t = k|y)$ using VideoBERT. (Below) Given a video token, we show the top three future tokens forecasted by VideoBERT at different time scales. In this case, VideoBERT predicts that a bowl of flour and cocoa powder may be baked in an oven, and may become a brownie or cupcake. We visualize video tokens using the images from the training set closest to centroids in feature space.

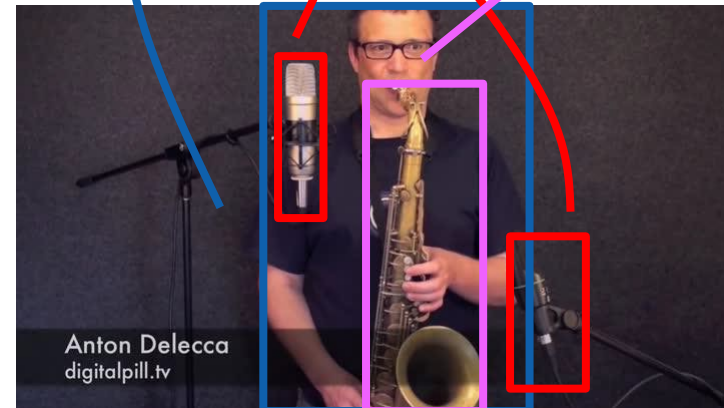
Grounded Visual Description

- Essentially, visual description + object grounding or detection
- To achieve better result interpretability, we need grounding!
 - Image domain: Neural Baby Talk, etc.
 - Video domain: Grounded Video Description, etc.
- Requires special dataset that has both description and bounding box

Single-Frame Annotation



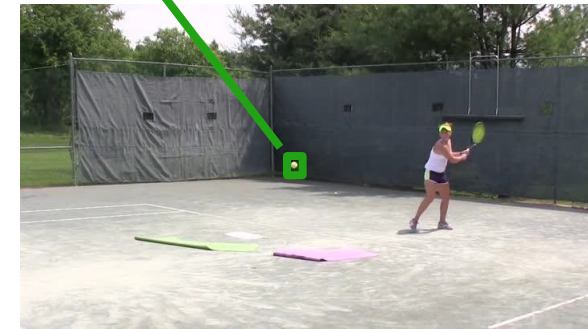
We see a man playing a saxophone
in front of microphones.



Multi-Frame Annotation



Two women are on a tennis court, showing the technique to posing and hitting the ball.



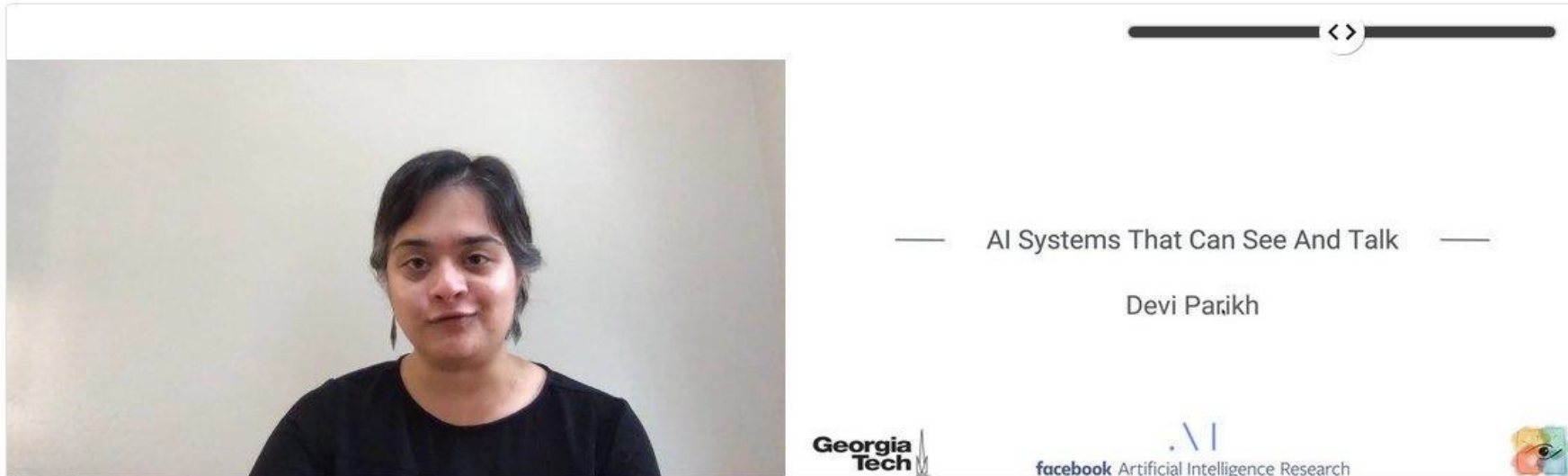
Problem Overview (2): VQA and Visual Reasoning

- How to train a smart multi-modal AI system that can both see and talk?

AI Systems That Can See And Talk

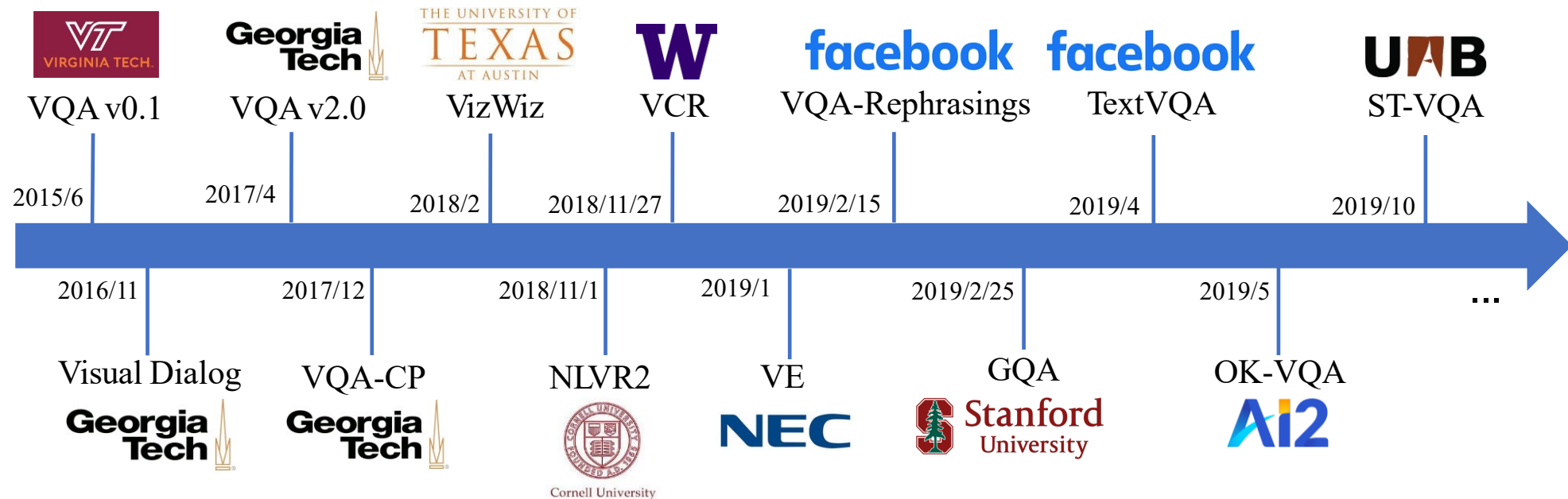
Prof. Devi Parikh / Georgia Tech and Facebook AI Research

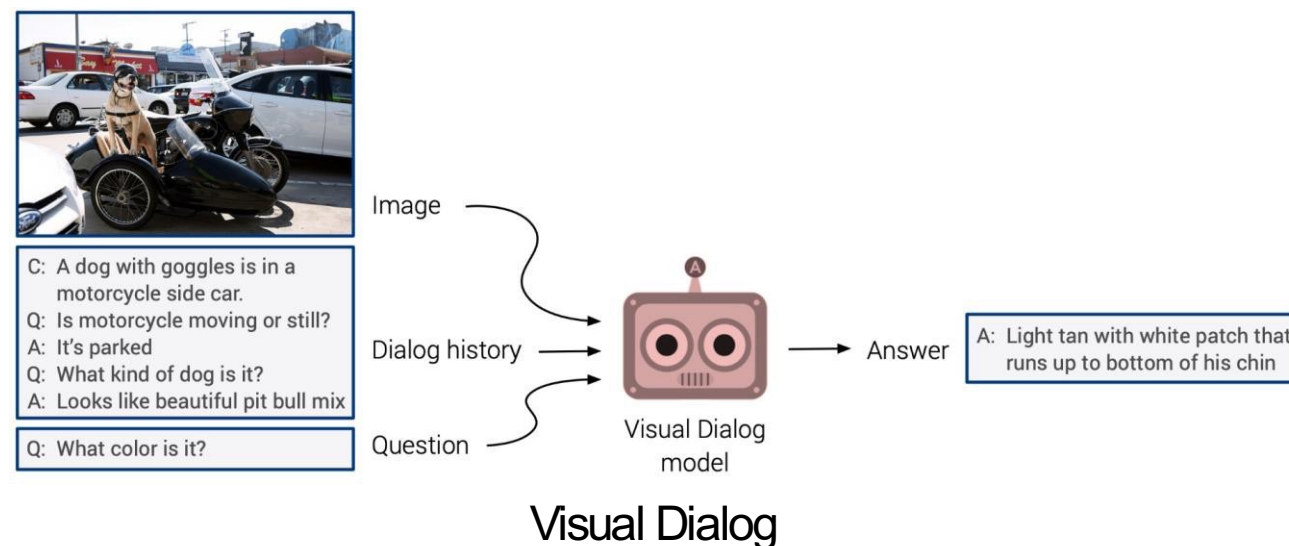
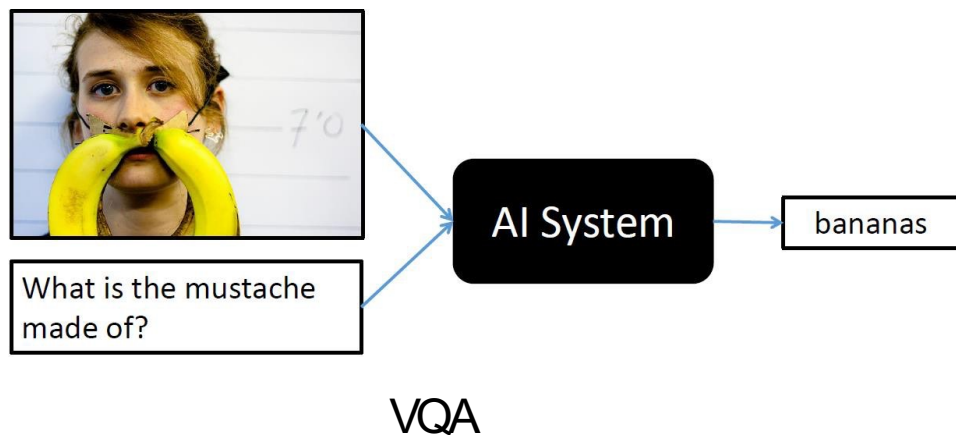
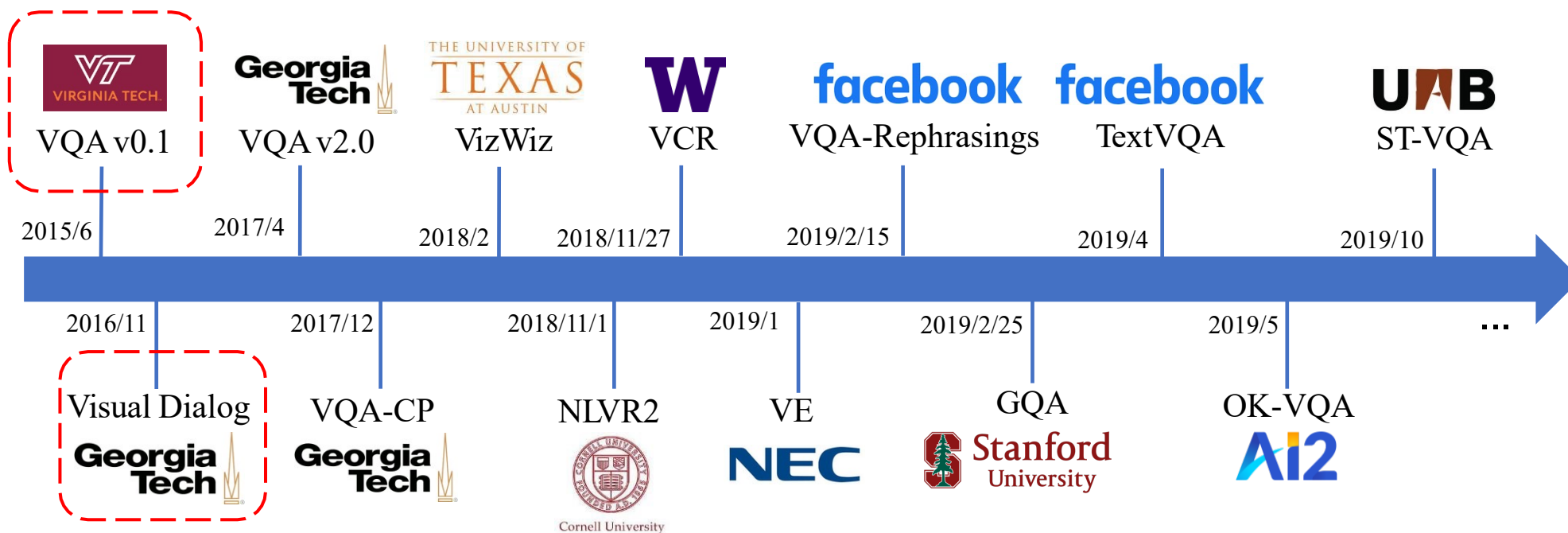
[Abstract & Bio](#)

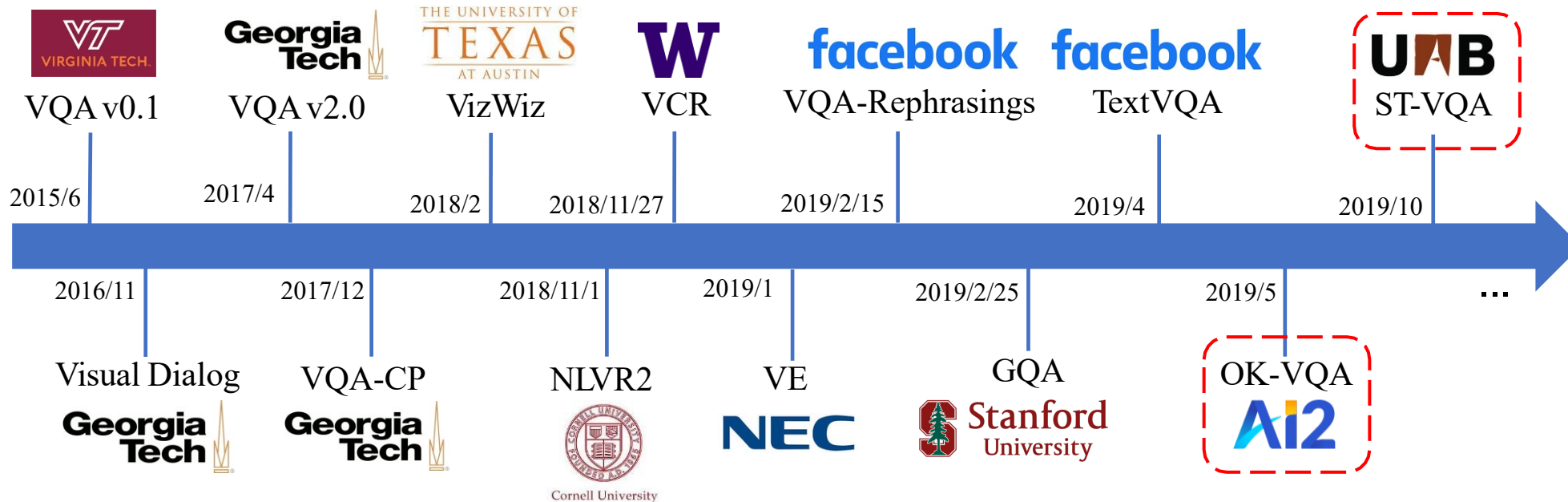


Problem Overview (2): VQA and Visual Reasoning

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







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

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.

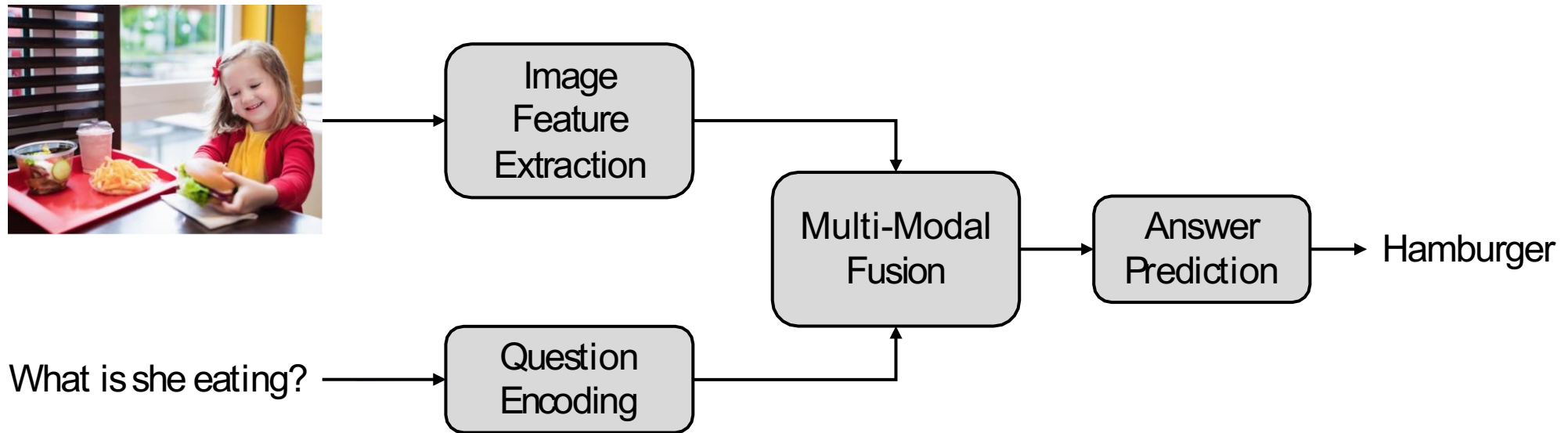
Beyond VQA: Visual Grounding

- PhraseCut: Language-based image segmentation



Approach Overview

- How a typical system looks like

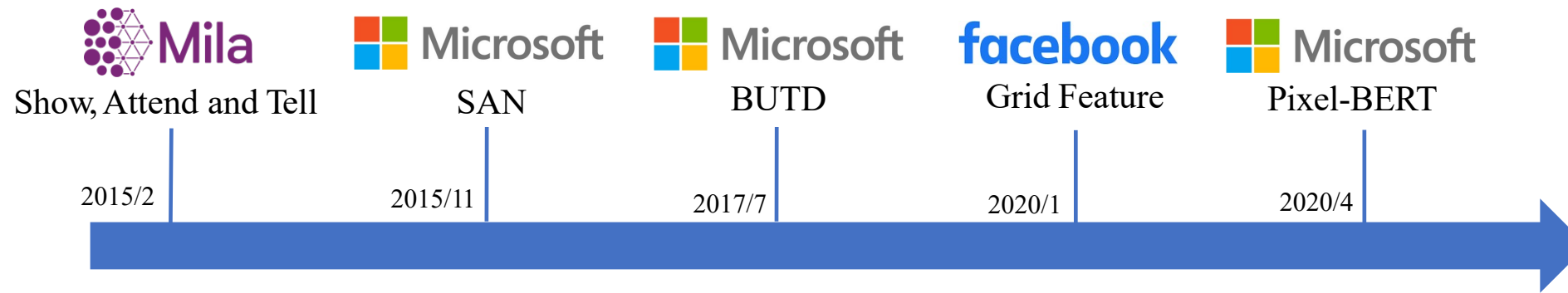


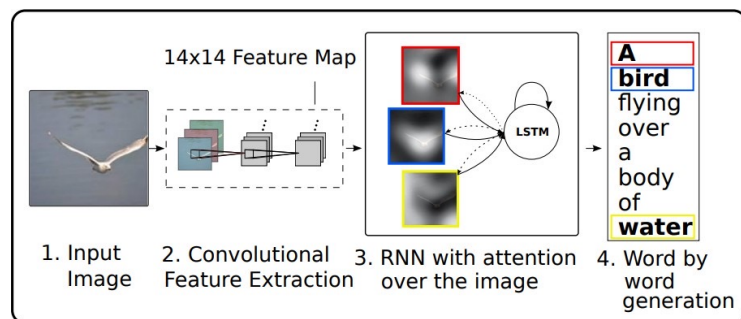
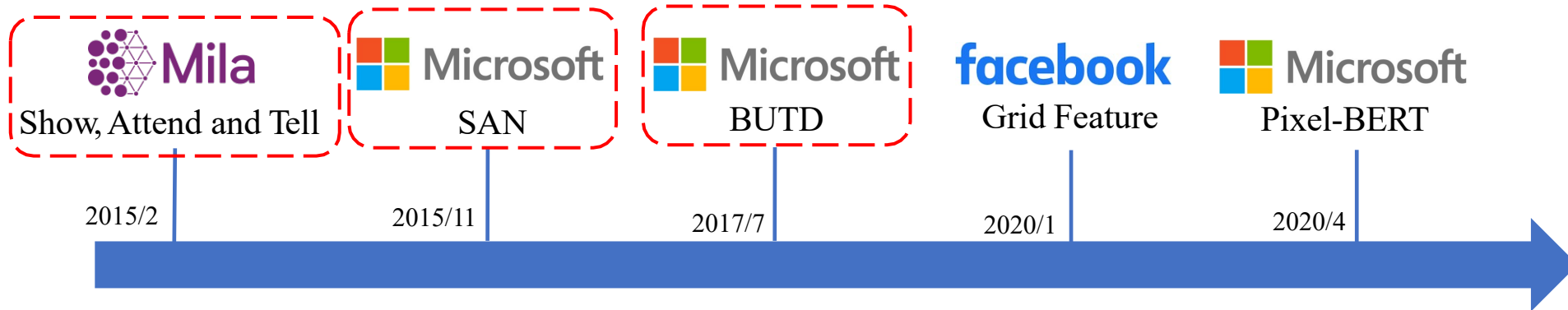
Research Challenges & Opportunities

- 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
- Multimodal pre-training

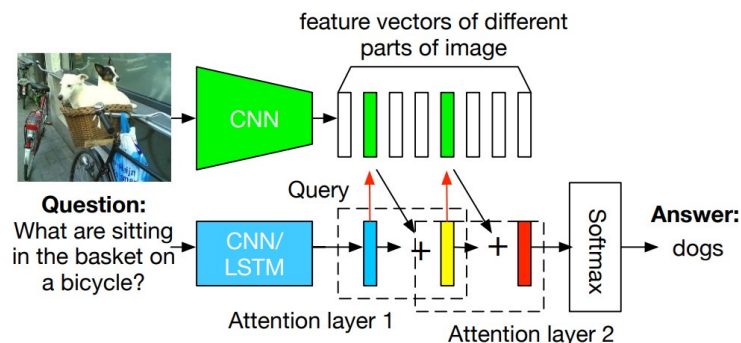
Better Image Feature Preparation

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

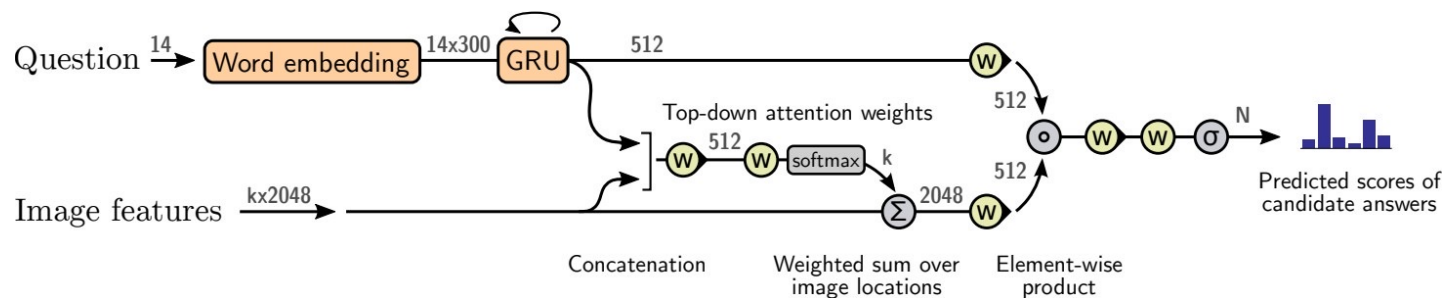




Show, Attend and Tell



Stacked Attention Network



2017 VQA Challenge Winner

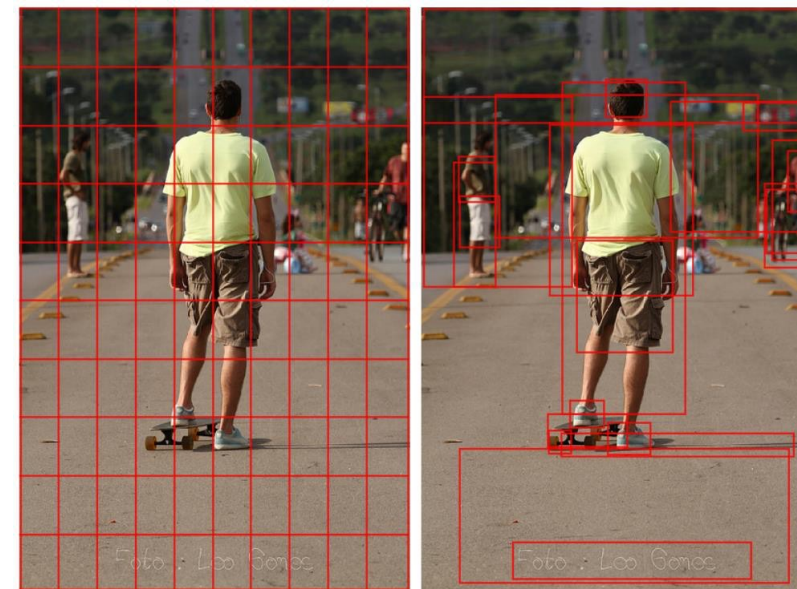


Figure 1. Typically, attention models operate on CNN features corresponding to a uniform grid of equally-sized image regions (left). Our approach enables attention to be calculated at the level of objects and other salient image regions (right).

- 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

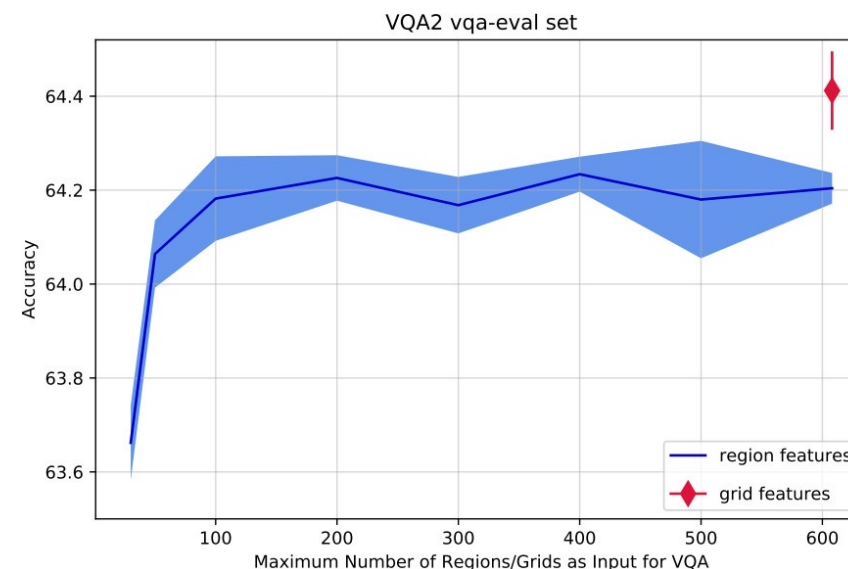
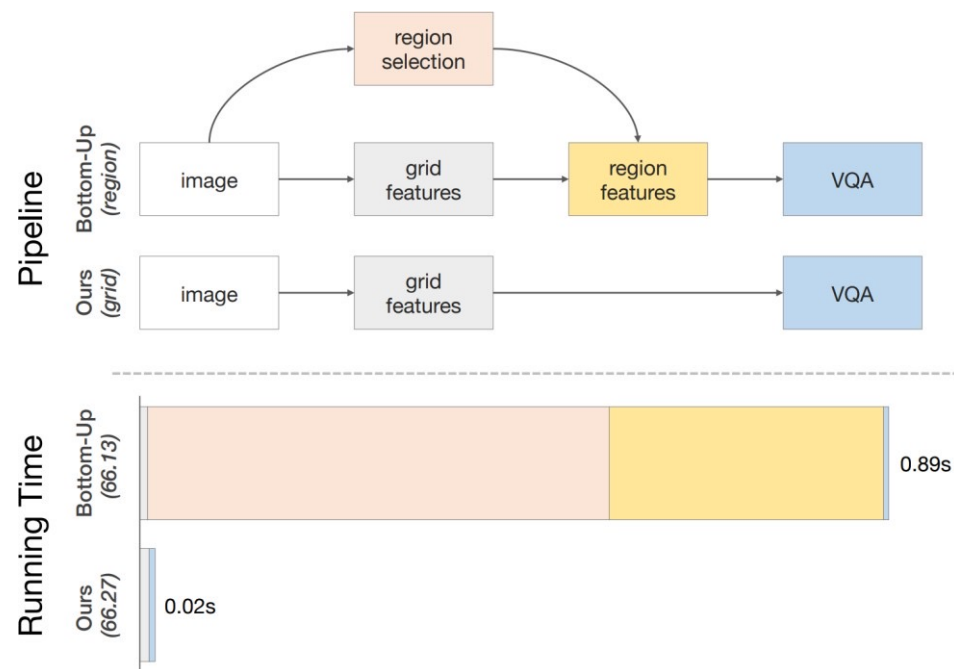
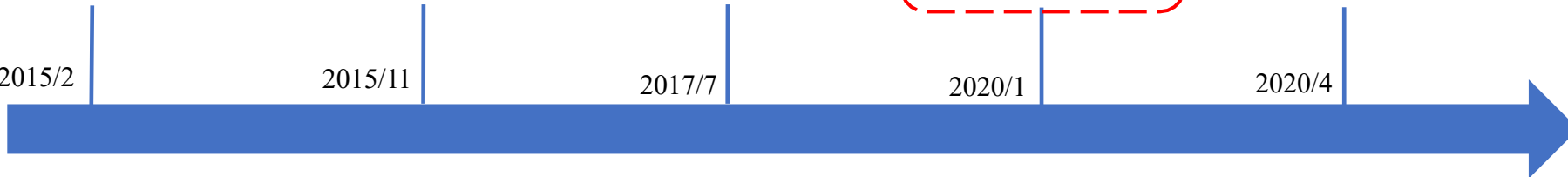
2017/7

 **facebook**
Grid Feature

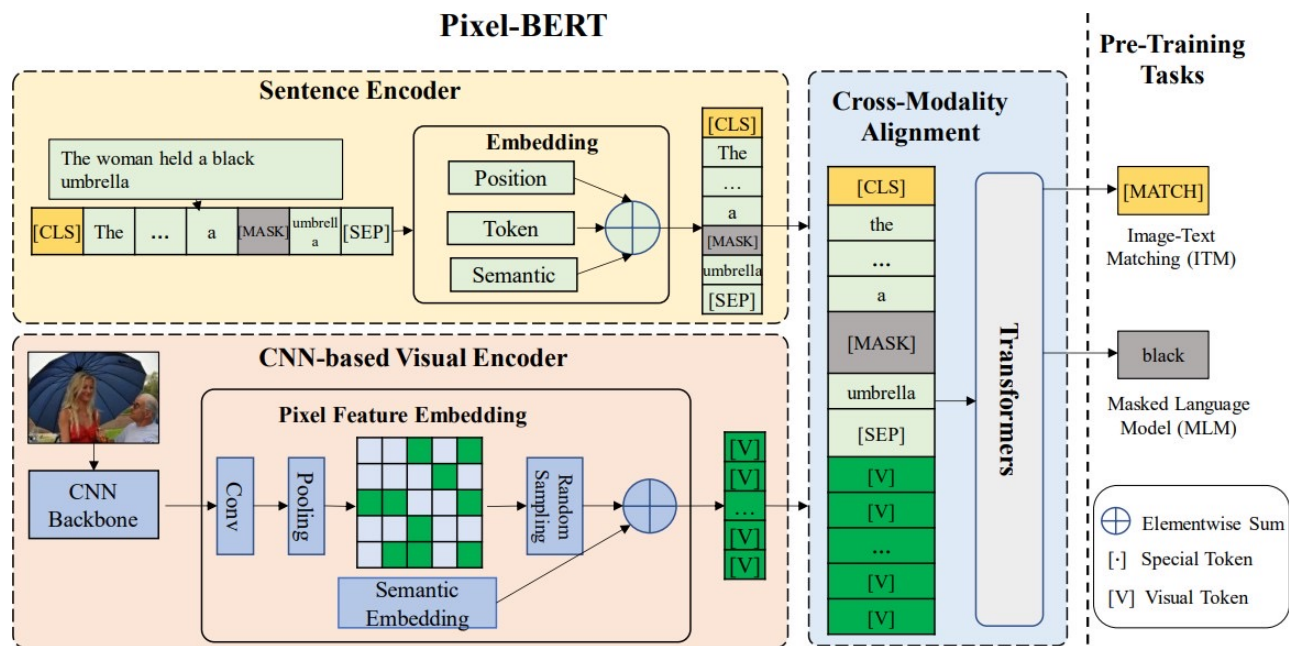
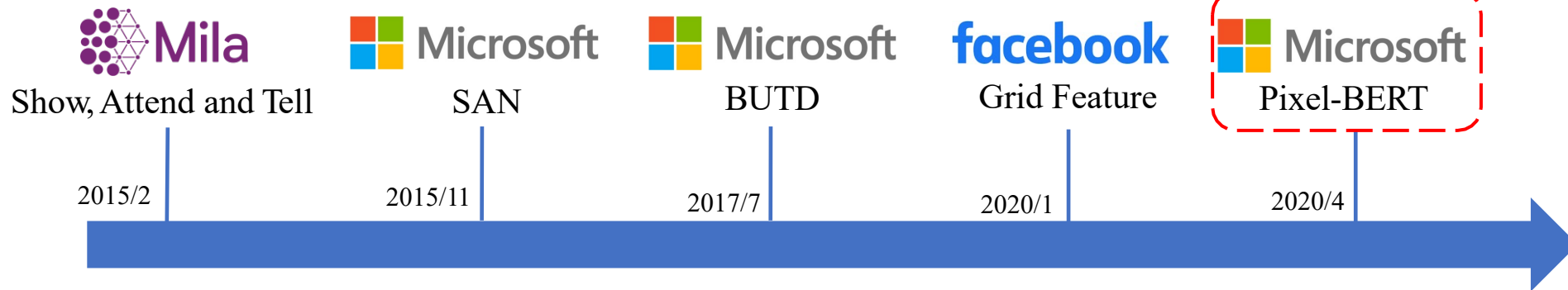
2020/1

 **Microsoft**
Pixel-BERT

2020/4



In Defense of Grid Features for VQA

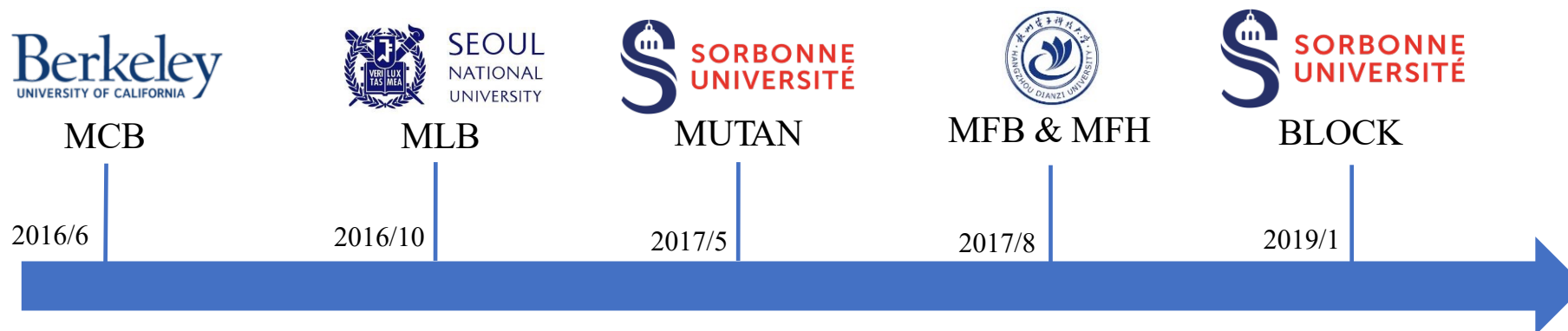


Model	test-dev	test-std
MUTAN[5]	60.17	-
BUTD[2]	65.32	65.67
ViLBERT[21]	70.55	70.92
VisualBERT[19]	70.80	71.00
VLBERT[29]	71.79	72.22
LXMERT[33]	72.42	72.54
UNITER[6]	72.27	72.46
Pixel-BERT (r50)	71.35	71.42
Pixel-BERT (x152)	74.45	74.55

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



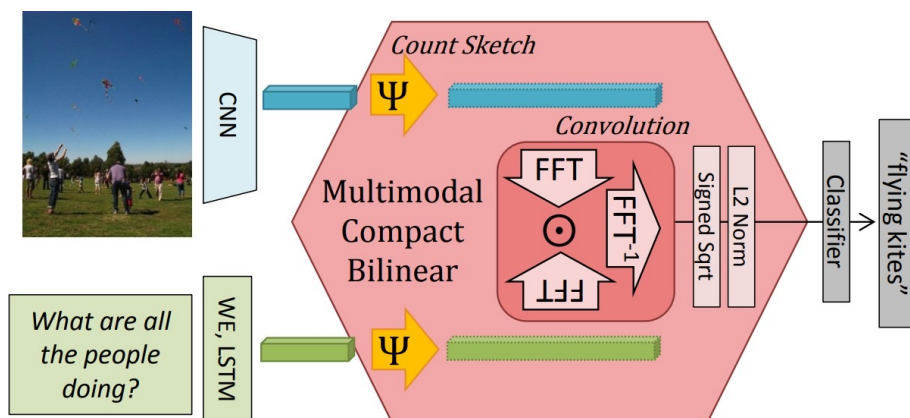
2017/5



2017/8



2019/1



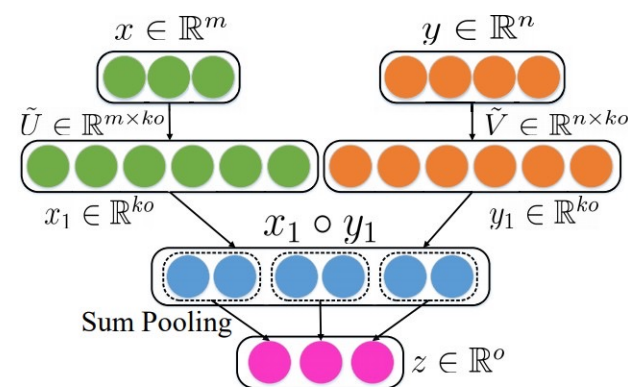
Multimodal Compact Bilinear Pooling

2016 VQA Challenge Winner

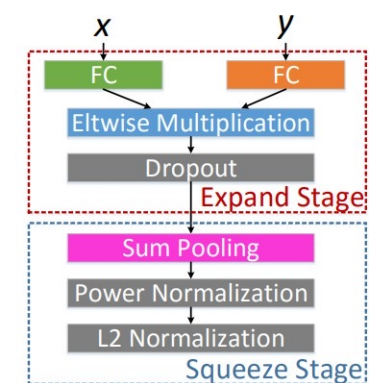
However, the feature after FFT is very high dimensional.

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

Multimodal Low-rank Bilinear Pooling

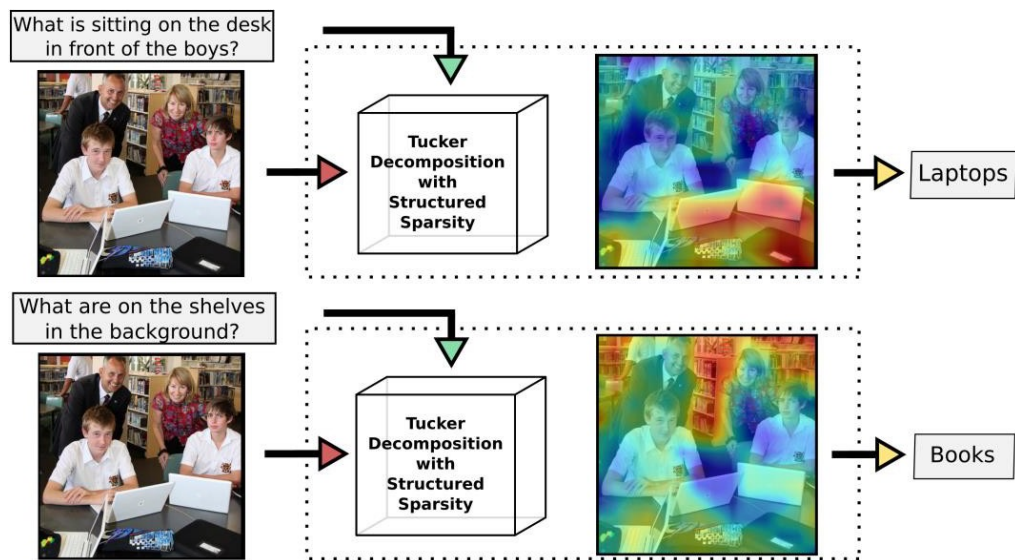


(a) Multi-modal Factorized Bilinear Pooling

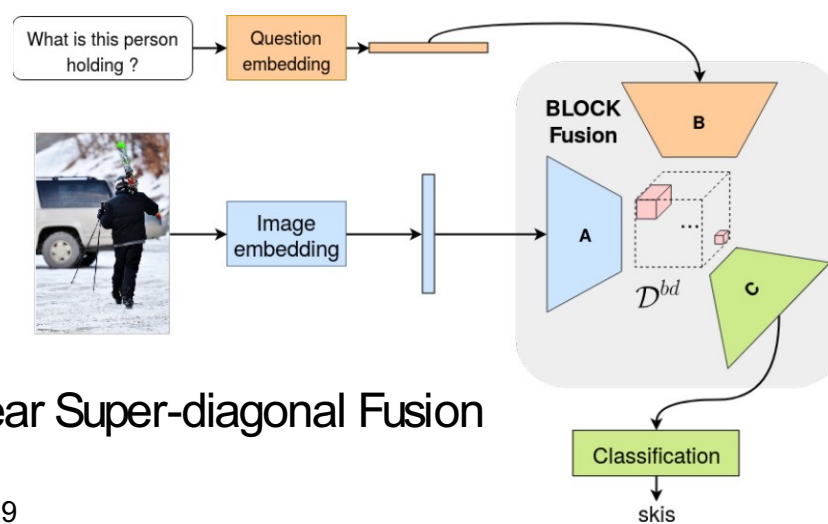
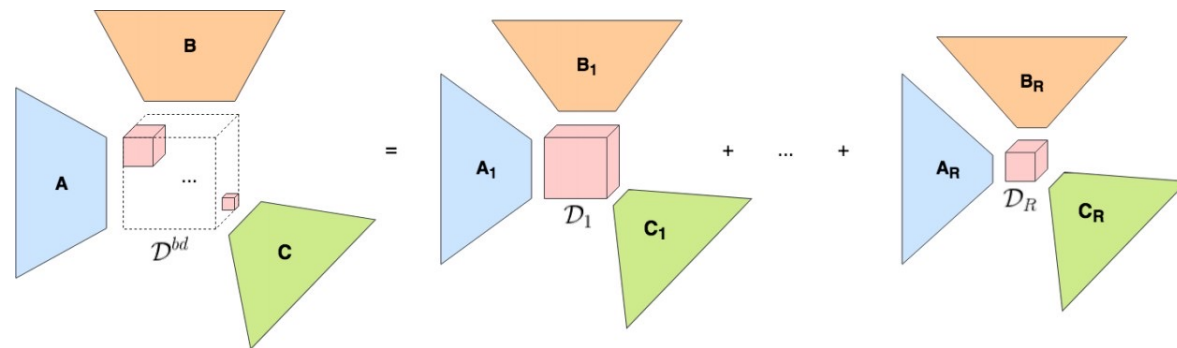


(b) MFB module

- 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



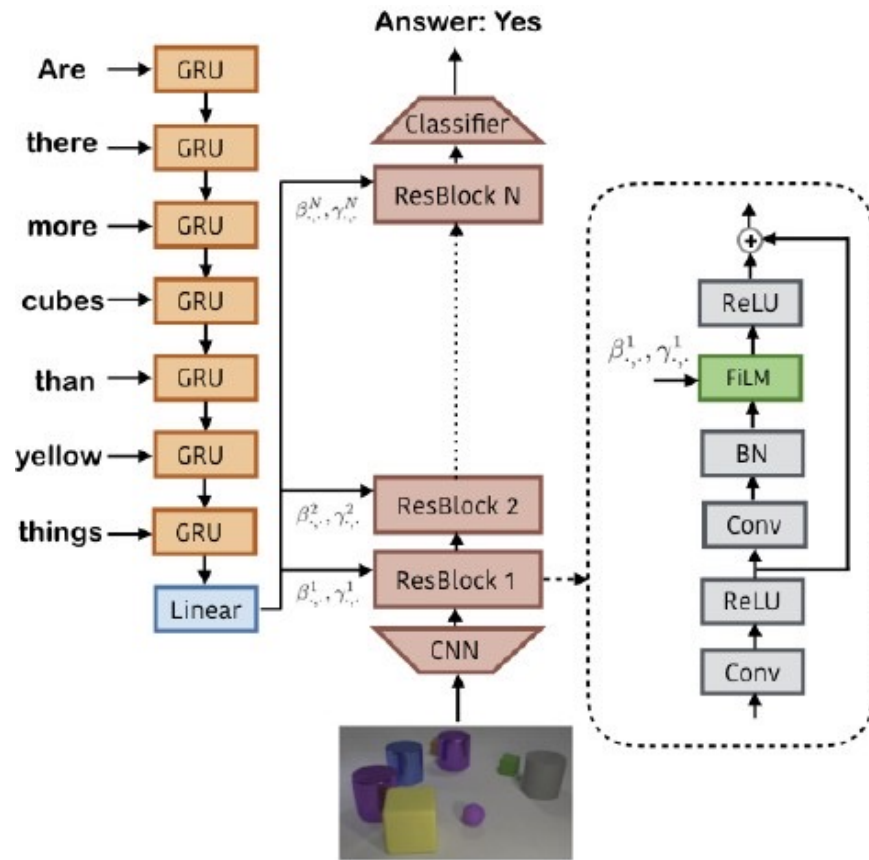
Multimodal Tucker Fusion



Bilinear Super-diagonal Fusion

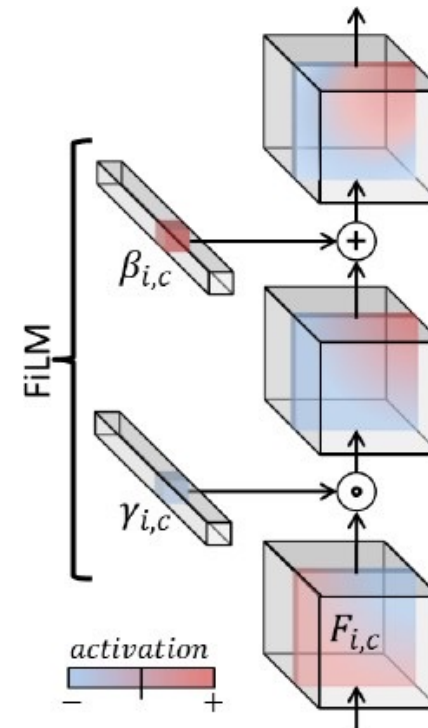
- 1 MUTAN: Multimodal Tucker Fusion for Visual Question Answering, ICCV 2017
- 2 BLOCK: Bilinear Superdiagonal Fusion for Visual Question Answering and Visual Relationship Detection, AAAI 2019

FiLM: Feature-wise Linear Modulation



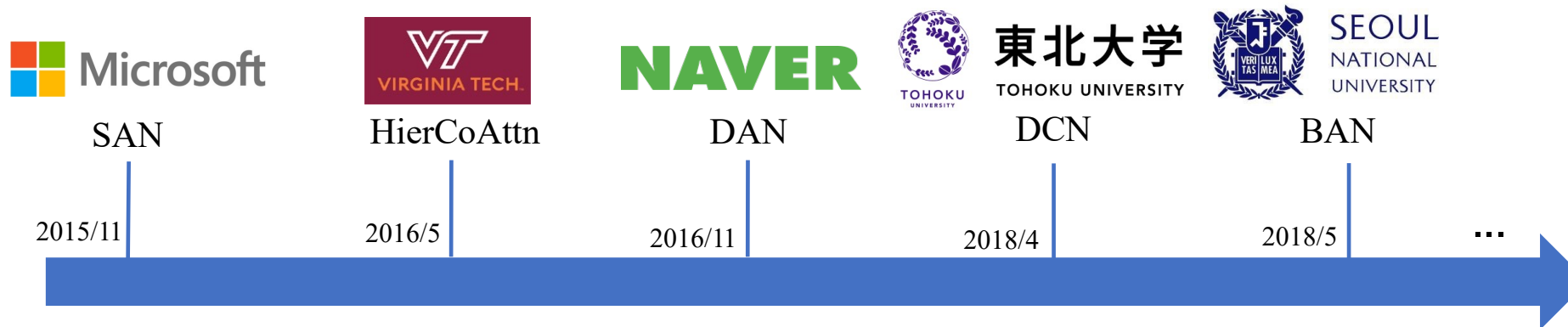
$$\gamma_{i,c} = f_c(x_i) \quad \beta_{i,c} = h_c(x_i),$$
$$FiLM(F_{i,c} | \gamma_{i,c}, \beta_{i,c}) = \gamma_{i,c} F_{i,c} + \beta_{i,c}.$$

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





SAN

2015/11



HierCoAttn

2016/5

NAVER

DAN

2016/11



東北大学

TOHOKU UNIVERSITY

DCN

2018/4

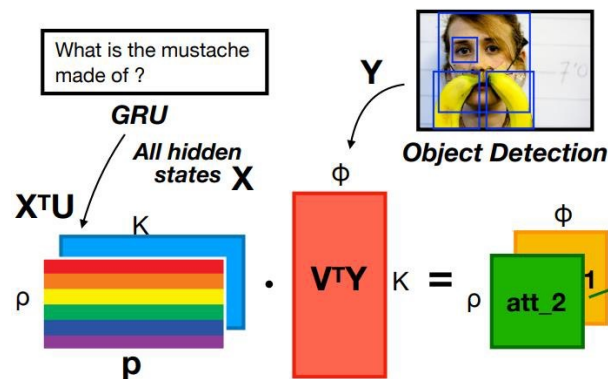
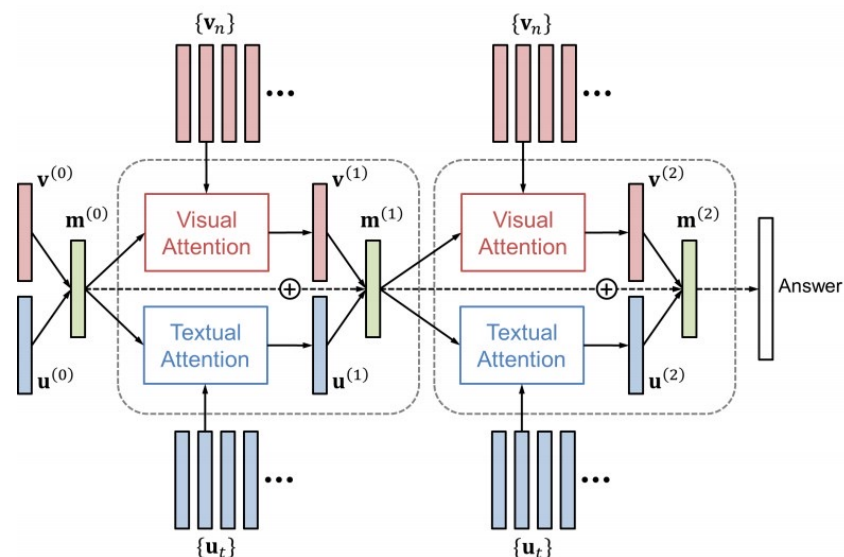


SEOUL
NATIONAL
UNIVERSITY

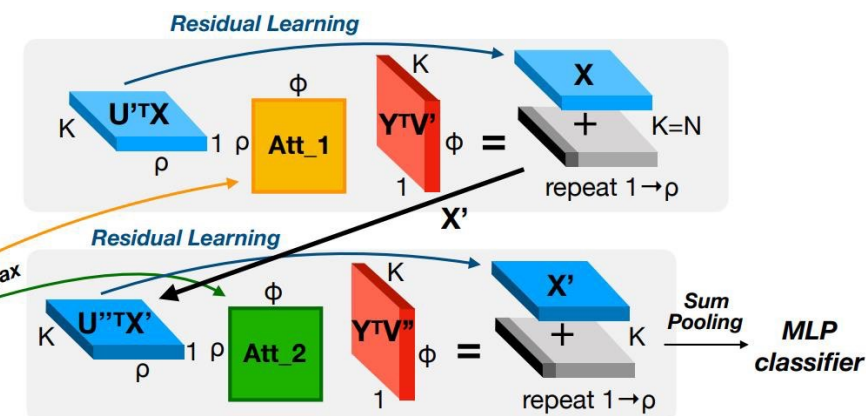
BAN

2018/5

...



Step 1. Bilinear Attention Maps



Step 2. Bilinear Attention Networks

2018 VQA Challenge Runner-Up

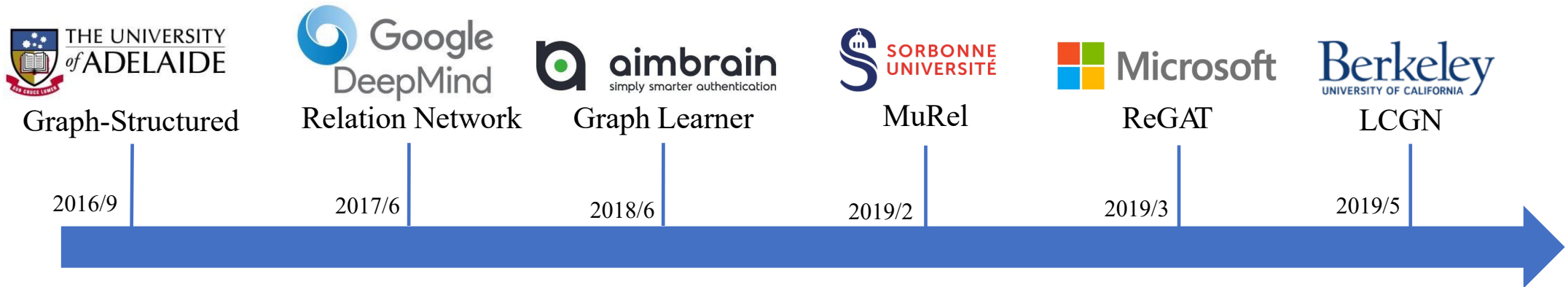
- Multiple Glimpses
- Residual Learning
- Counter Module
- Glove Embeddings

DAN: Dual Attention Network

DCN: Dense Co-attention Network

Relational Reasoning

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





Graph-Structured

2016/9



Relation Network

2017/6



Graph Learner

2018/6



MuRel

2019/2



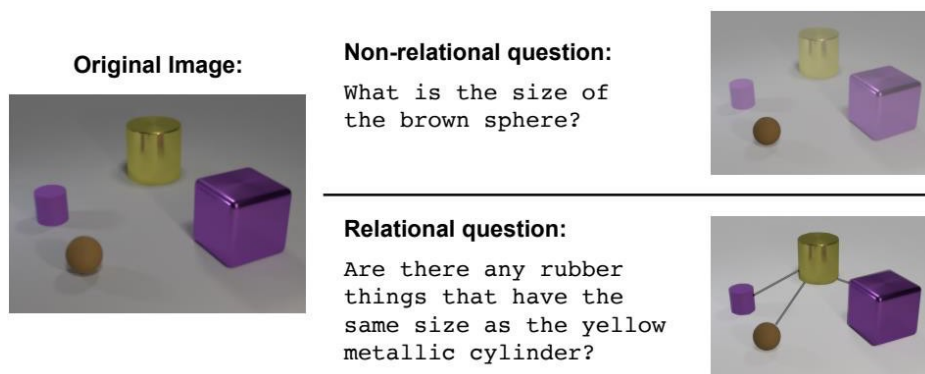
ReGAT

2019/3

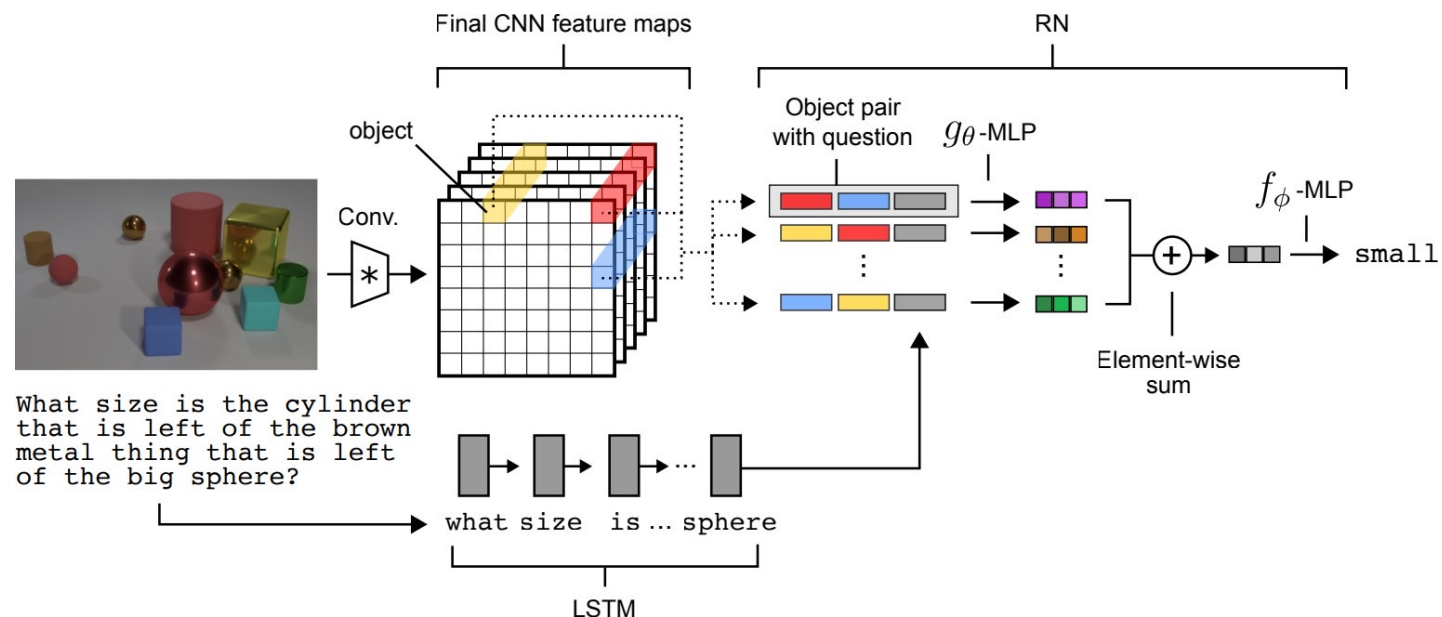


LCGN

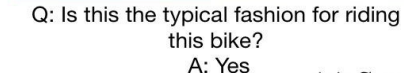
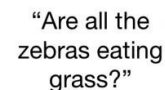
2019/5



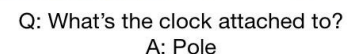
$$\text{RN}(O) = f_{\phi} \left(\sum_{i,j} g_{\theta}(o_i, o_j) \right)$$



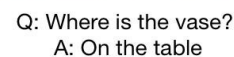
Relational Network: A fully-connected graph is constructed



(a) Semantic Relation



(b) Spatial Relation



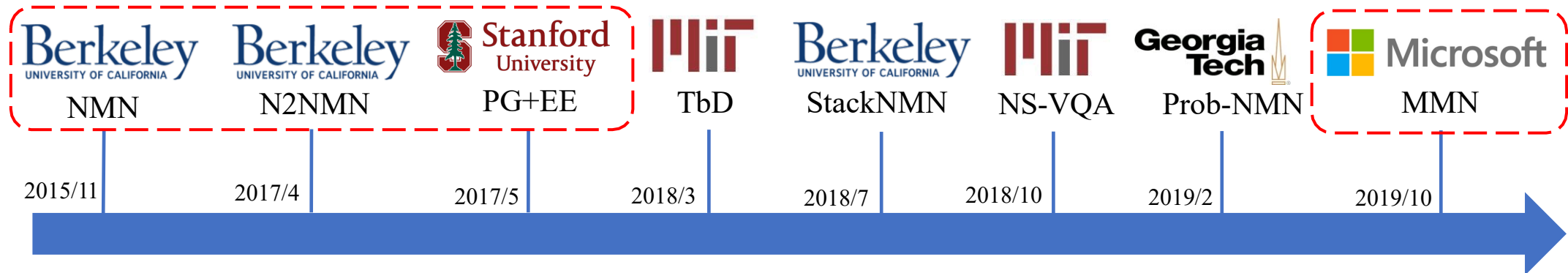
Q: Should the people be walking according to the light?
A:No

(c) Implicit Relation

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

Neural Module Network (NMN)

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



- 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

Consider a compositional model

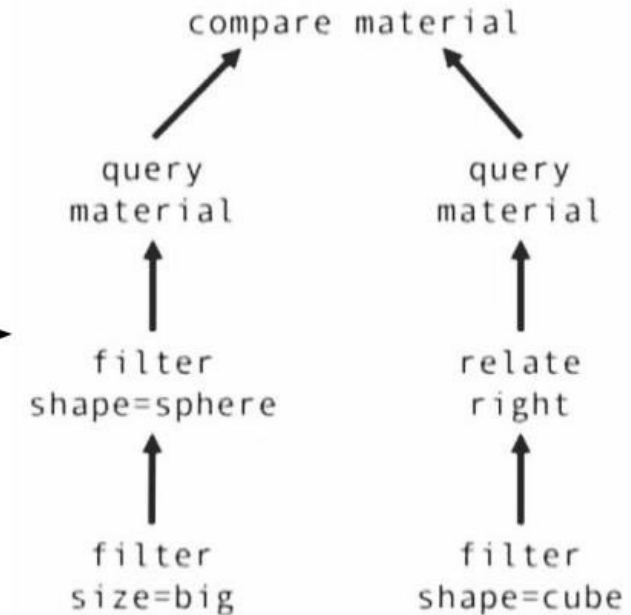
Q: How many spheres are the left of the big sphere and the same color as the small rubber cylinder?

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?

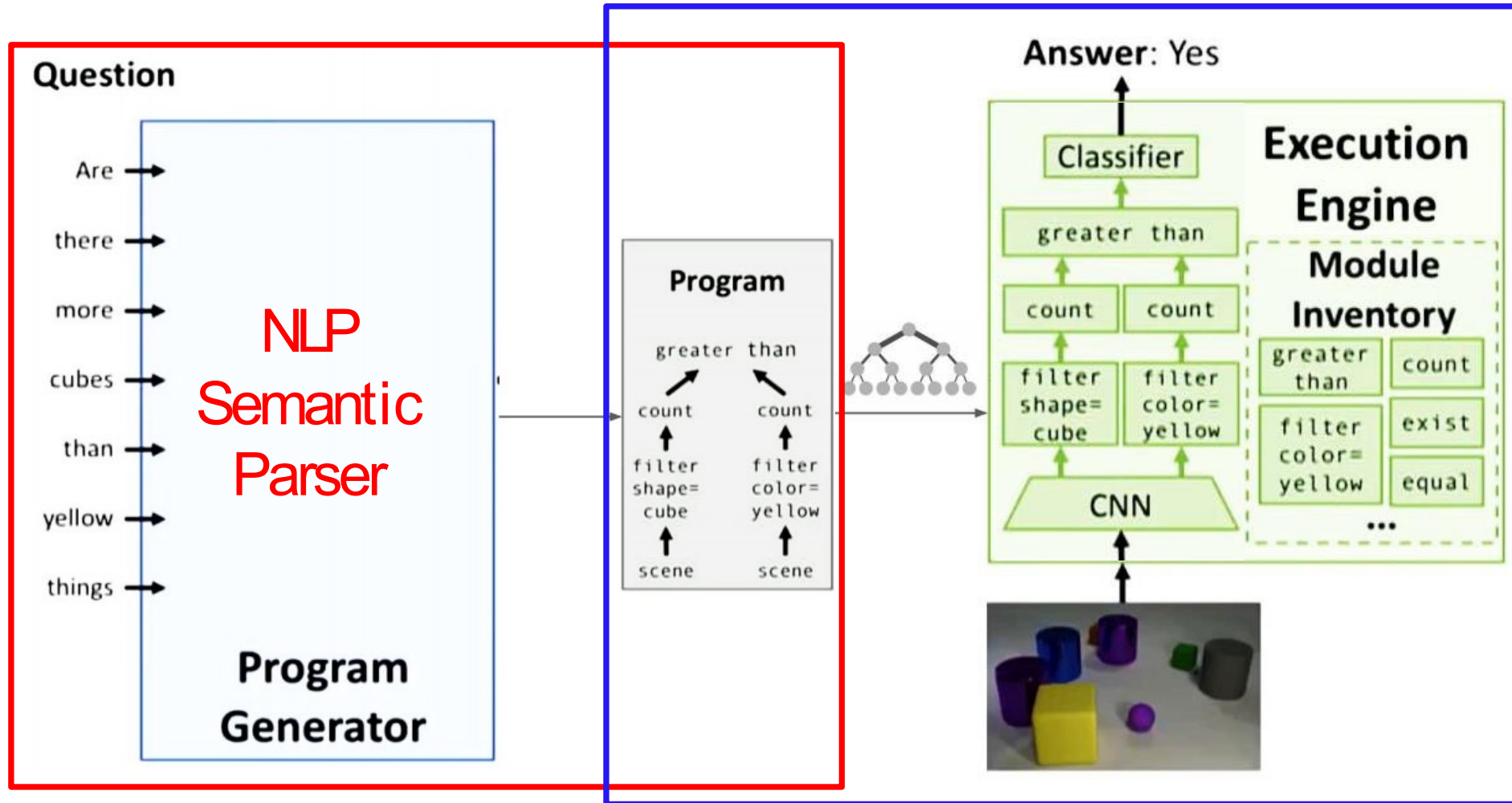
Common operations

Attributes identification
Counting objects
Comparisons
Spatial relationships
Logical operations



**Network architecture
corresponding to the
third question**

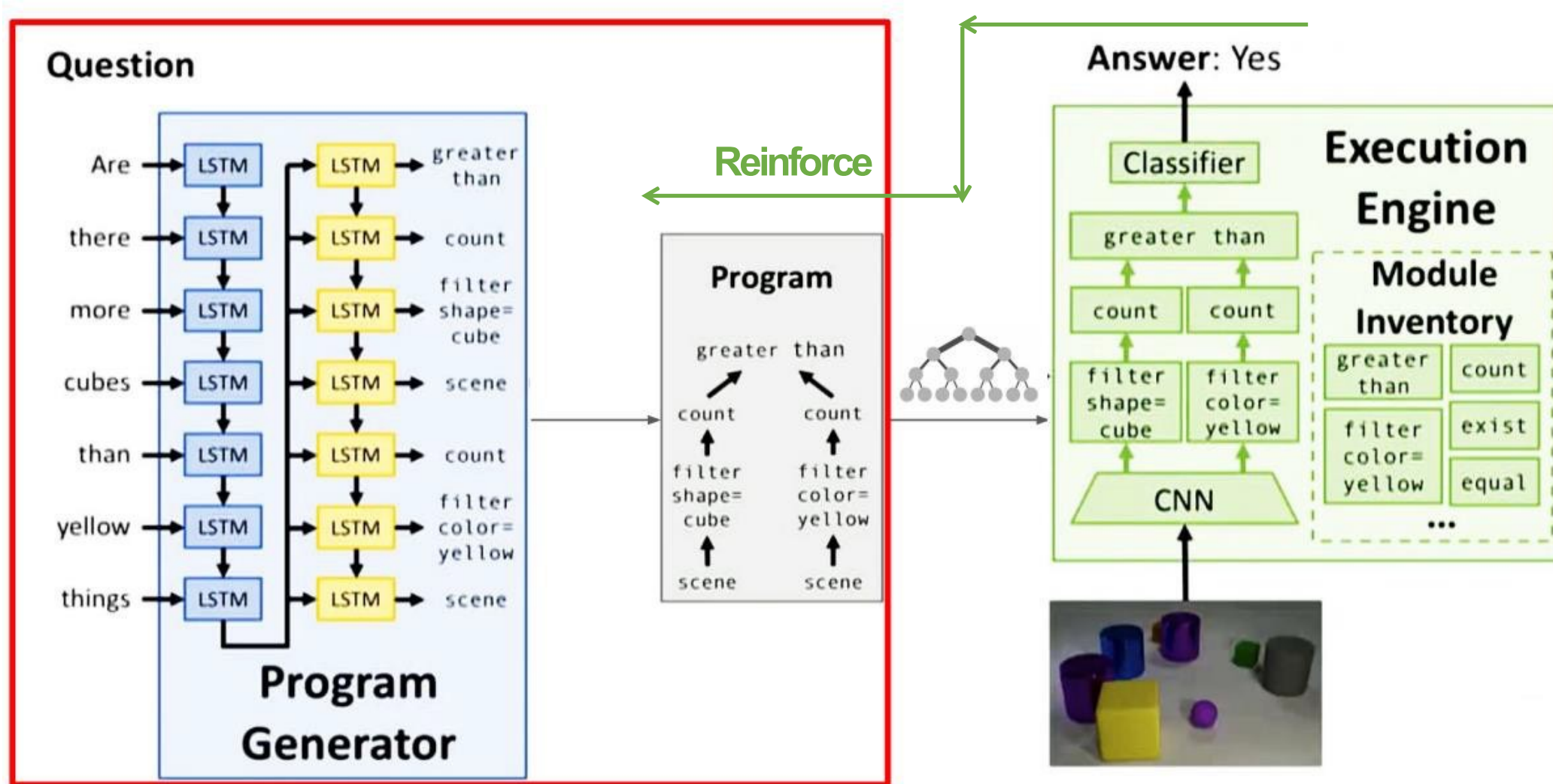
Overview of the NMN approach

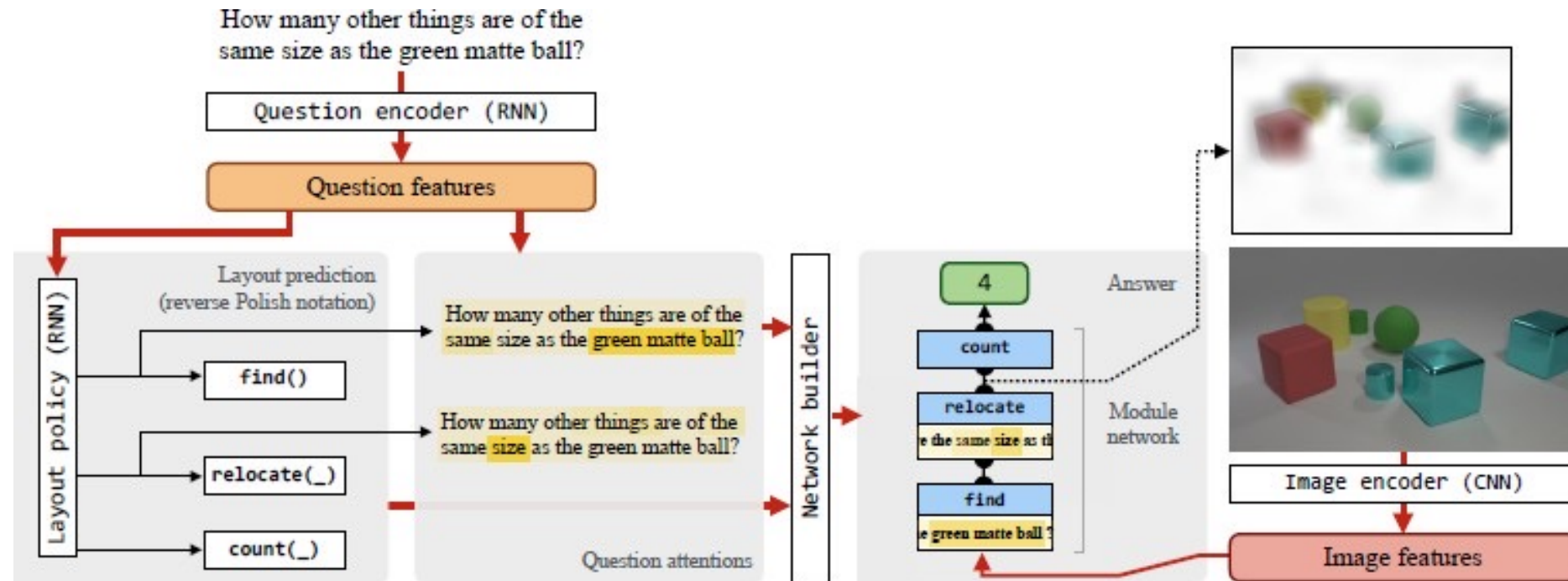
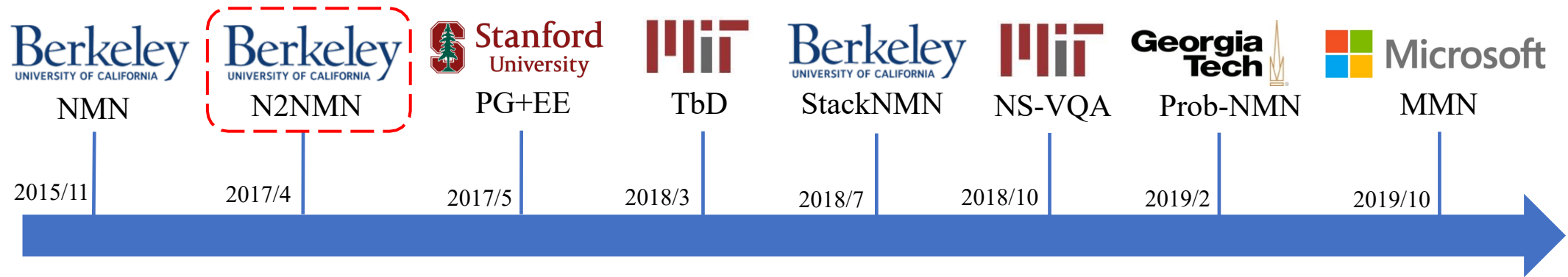


Uses some pre-trained parser


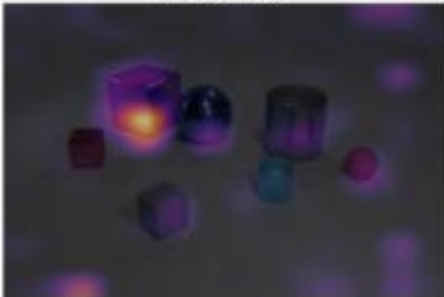
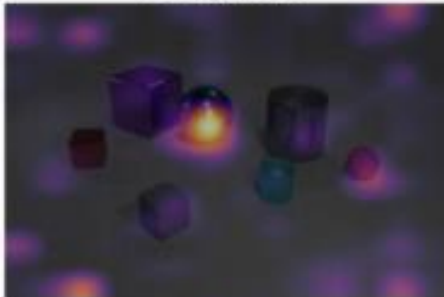
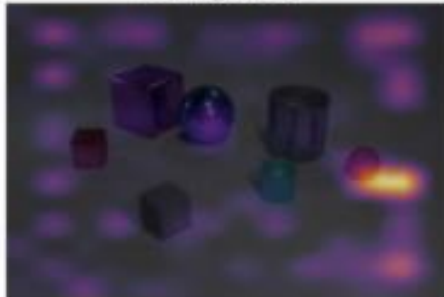
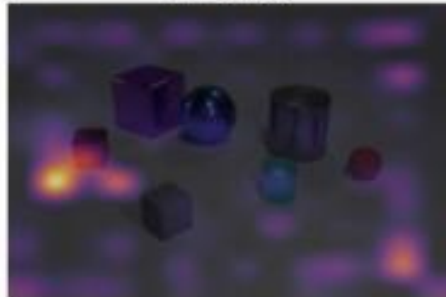
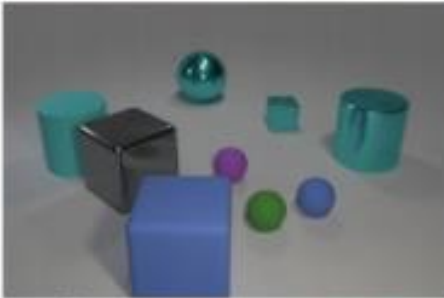
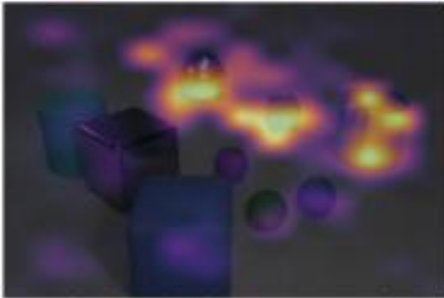
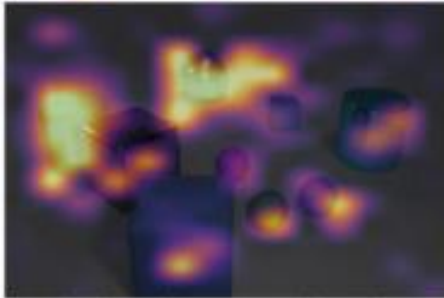
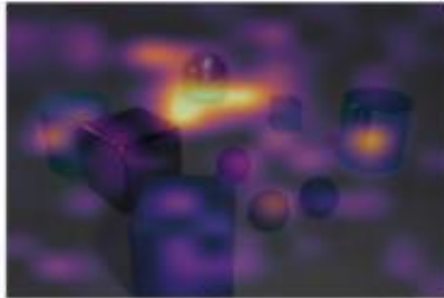
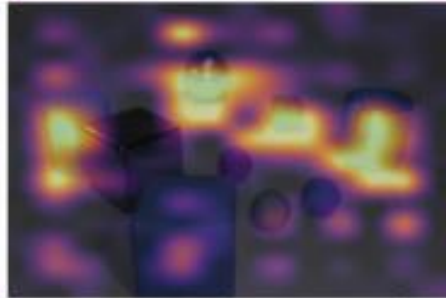
Trained separately

Inferring and Executing Programs



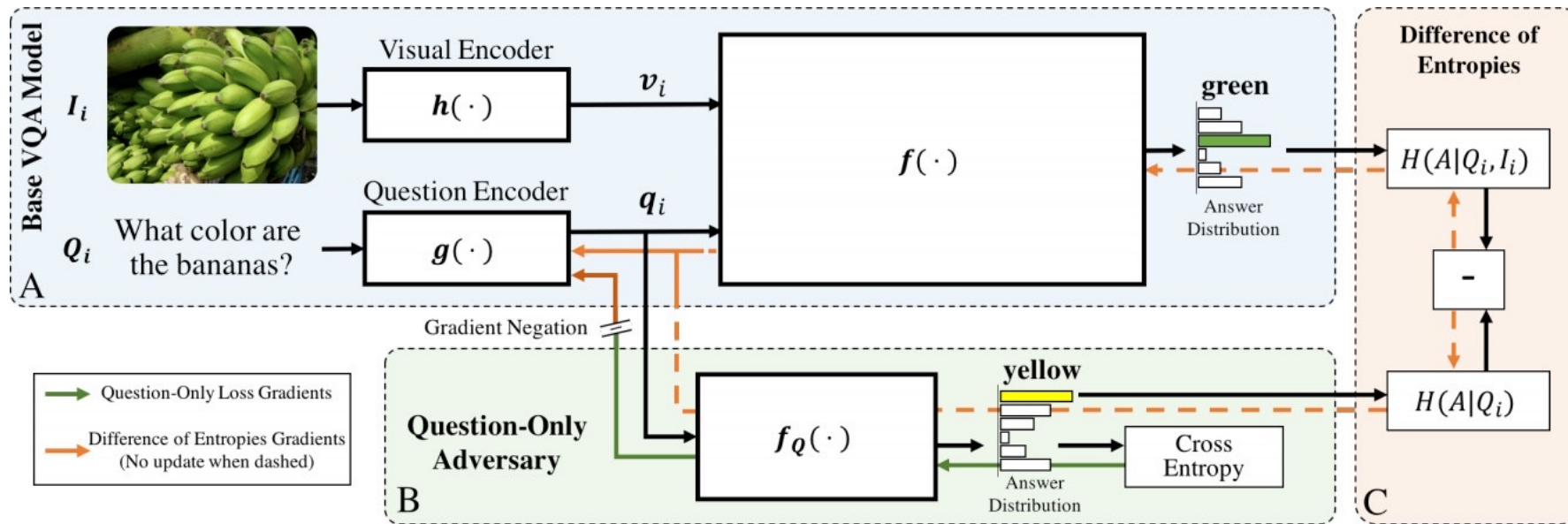


What do the modules learn?

<p>Q: What shape is the...</p>	<p>...<u>purple</u> thing?</p>	<p>...<u>blue</u> thing?</p>	<p>...<u>red</u> thing <u>right</u> of the <u>blue</u> thing?</p>	<p>...<u>red</u> thing <u>left</u> of the <u>blue</u> thing?</p>
	<p>A: cube</p>	<p>A: sphere</p>	<p>A: sphere</p>	<p>A: cube</p>
				
				
<p>Q: How many <u>cyan</u> things are...</p>	<p>...<u>right</u> of the <u>gray</u> cube?</p>	<p>...<u>left</u> of the <u>small</u> cube?</p>	<p>...right of the gray cube <u>and</u> left of the small cube?</p>	<p>...right of the gray cube <u>or</u> left of the small cube?</p>
	<p>A: 3</p>	<p>A: 2</p>	<p>A: 1</p>	<p>A: 4</p>

Robust VQA: an example

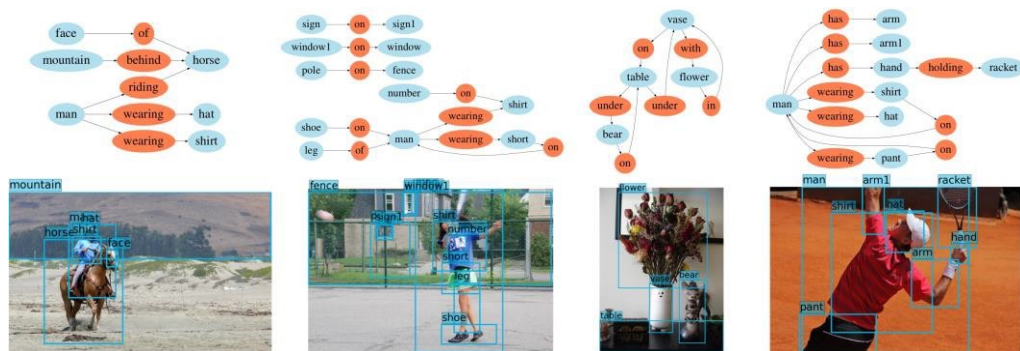
- Overcoming language prior with adversarial regularization



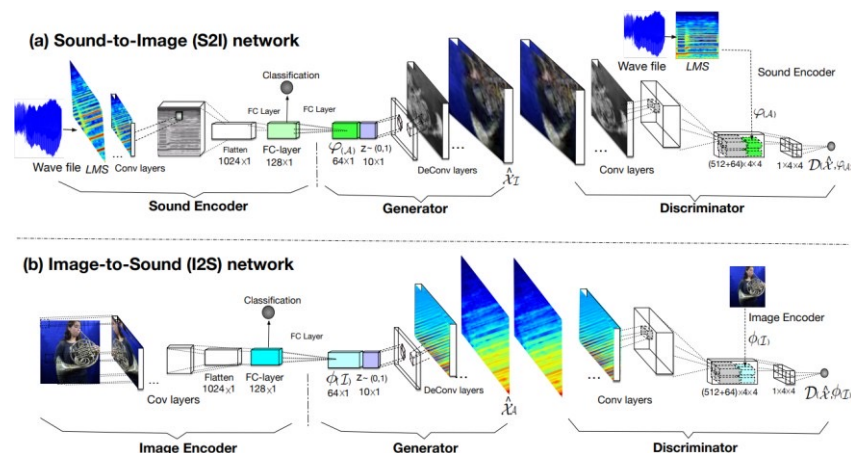
Problem Overview (3): Text-to-Image Generation

- Text-to-Image Synthesis
 - StackGAN, AttnGAN, TAGAN, ObjGAN ...
- Text-to-Video Synthesis
 - GAN-based methods, VAE-based methods, StoryGAN ...
- Dialogue-based Image Synthesis
 - ChatPainter, CoDraw, SeqAttnGAN ...

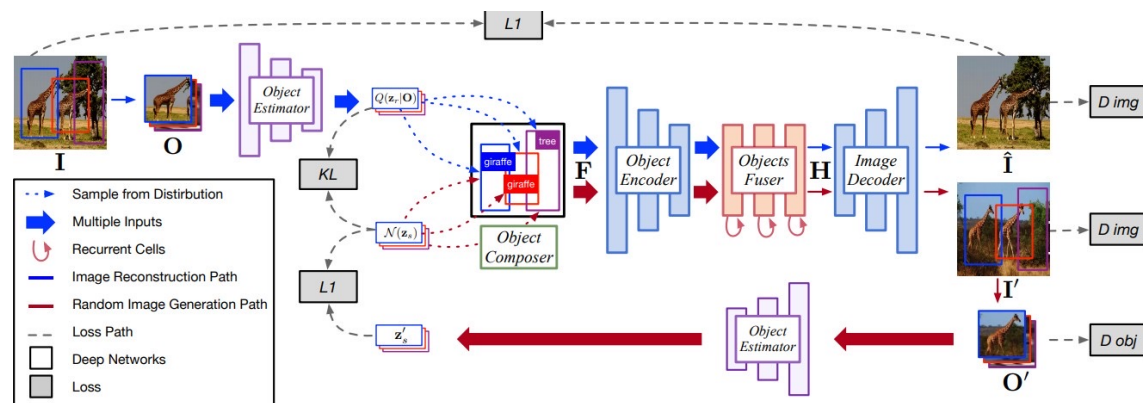
Conditional Image Synthesis



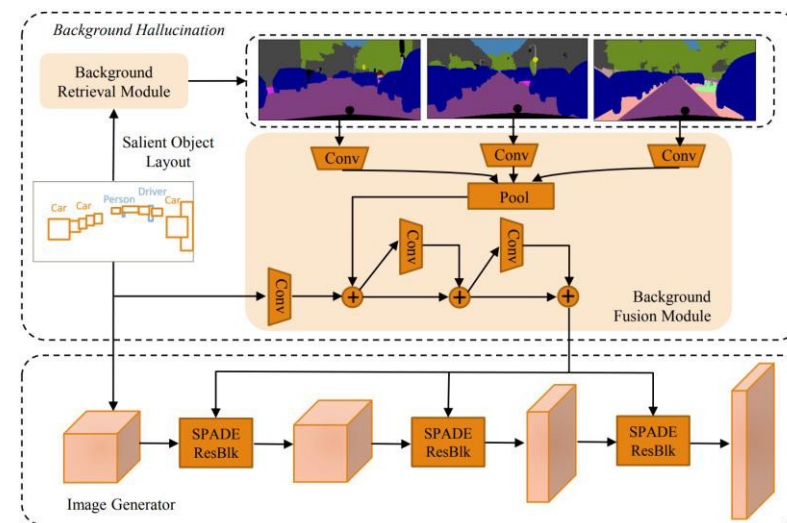
SceneGraph2img [Johnson et al., 2018]



Audio2img [Chen et al., 2019]

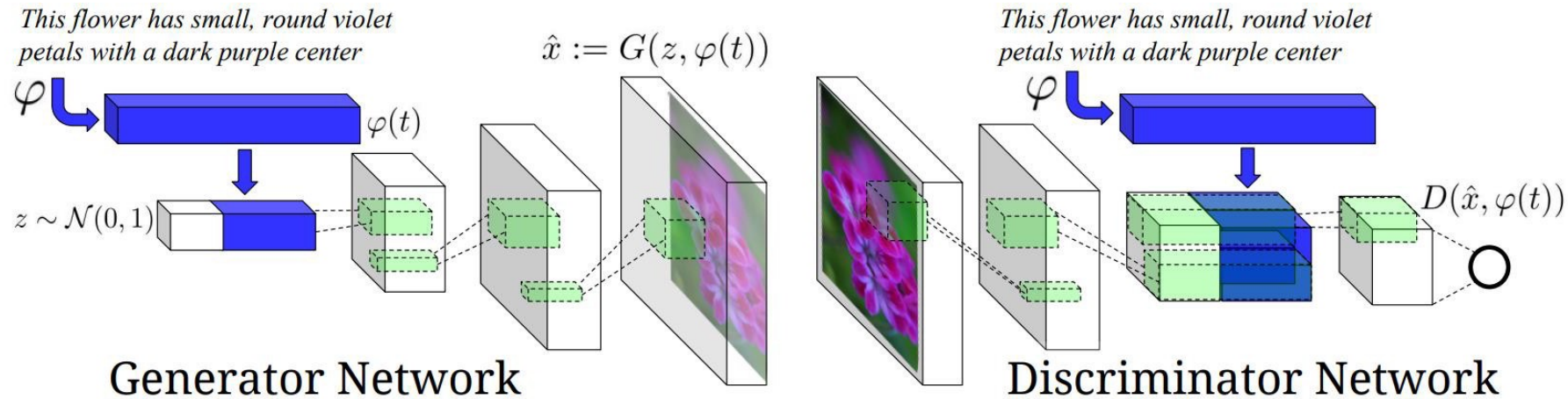
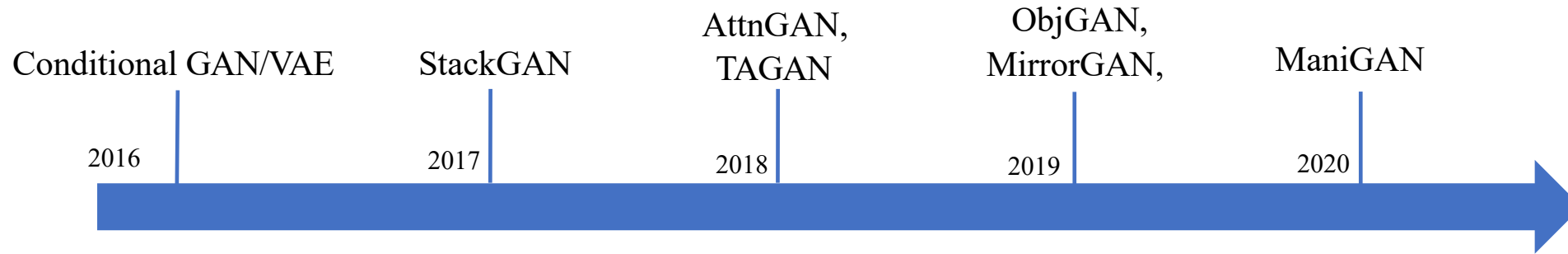


Layout2img [Zhao et al., 2019]



BachGAN [Li et al., 2020]

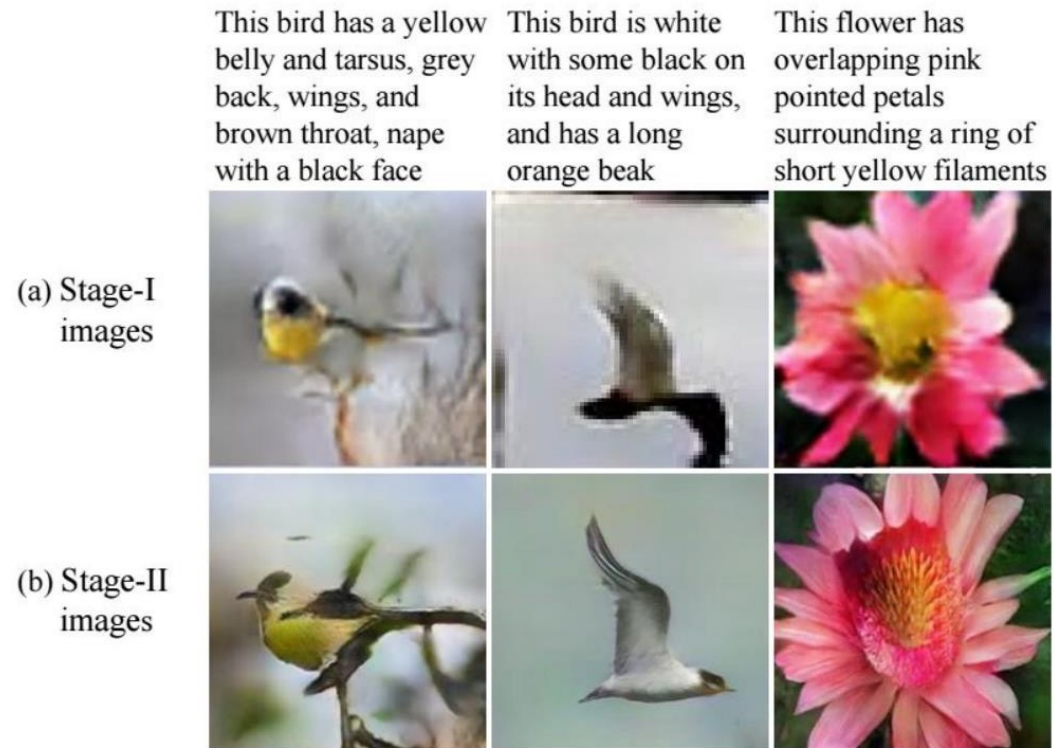
Text-to-Image Synthesis



StackGAN

- Stage 1.
 - Generates 64x64 images
 - Structural information
 - Low detail
- Stage 2.
 - Requires Stage 1. output
 - Upsamples to 256x256
 - Higher detail, photorealistic

Both stages take in the same conditioned textual input



AttnGAN

- Paying attentions to the relevant words in the natural language description
- Capture both both the global sentence level information and the fine-grained word level information



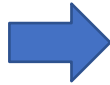
Text-to-Video Synthesis

- StoryGAN: Short story (sequence of sentences) → Sequence of images

Image Generation

Story Visualization

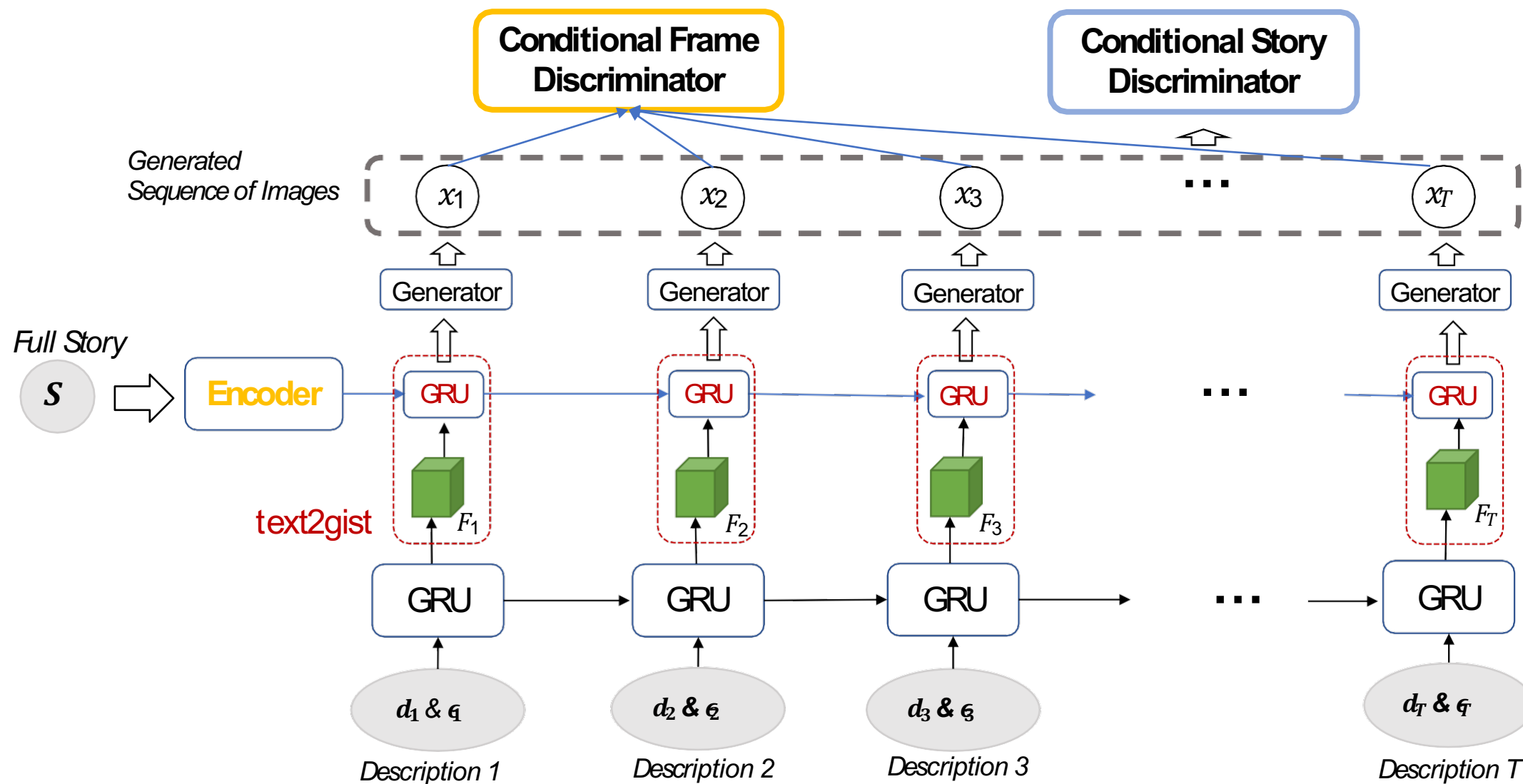
“A small yellow bird with a black crown and beak.”



“Pororo and Crong fishing together. Crong is looking at the bucket. Pororo has a fish on his fishing rod.”



StoryGAN



Precise Generation on CLEVR Dataset

- Given attributes of objects, generate the image

StoryGAN Ground Truth StackGAN

"Small purple rubber sphere, position is 1.4, -0.7."



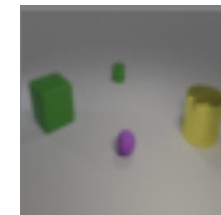
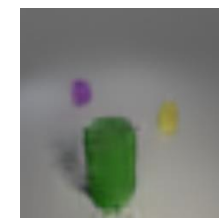
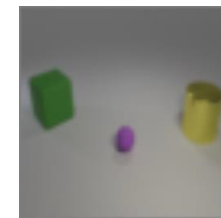
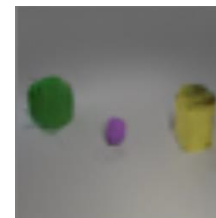
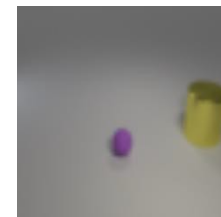
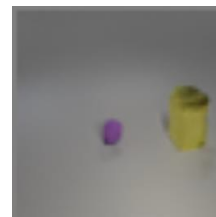
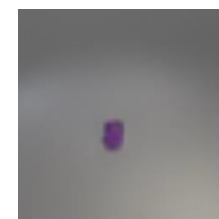
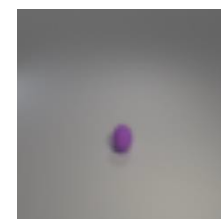
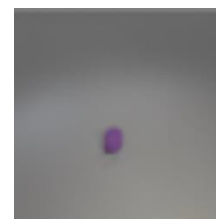
"Large yellow metallic cylinder, position is 2.1, 2.6."



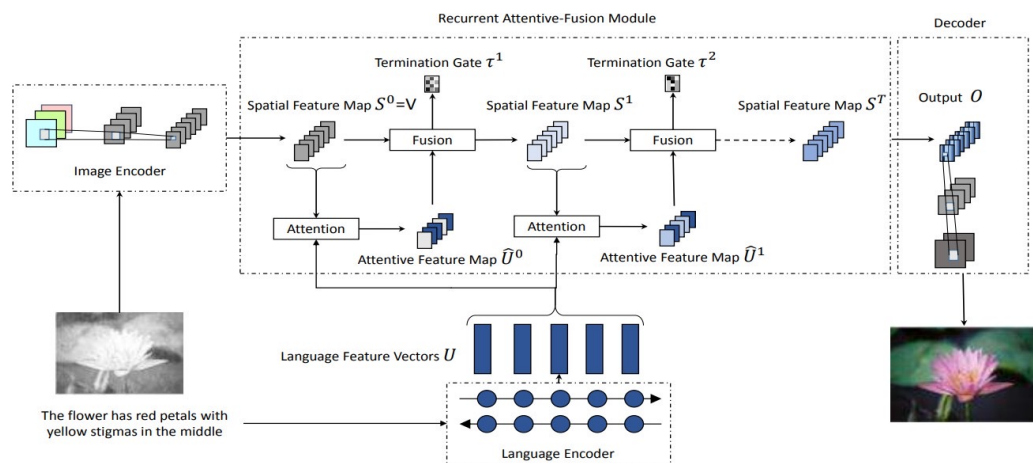
"Large green rubber cube, position is -2.0, -1.2."



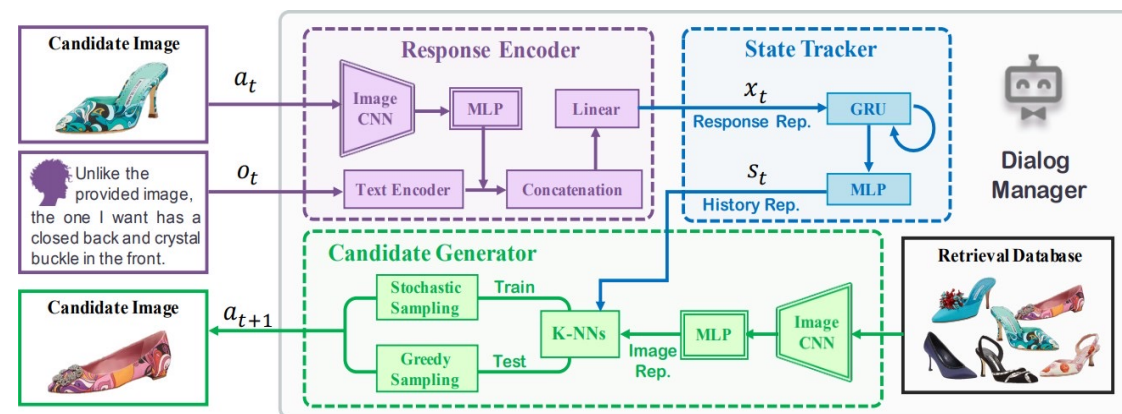
"Small green rubber cylinder, position is -2.5, 1.6."



Dialogue-based Image Synthesis



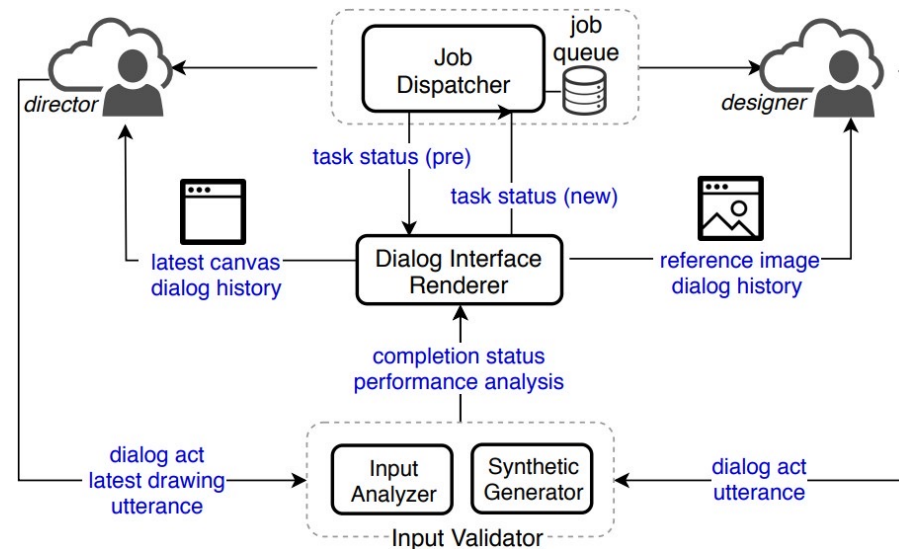
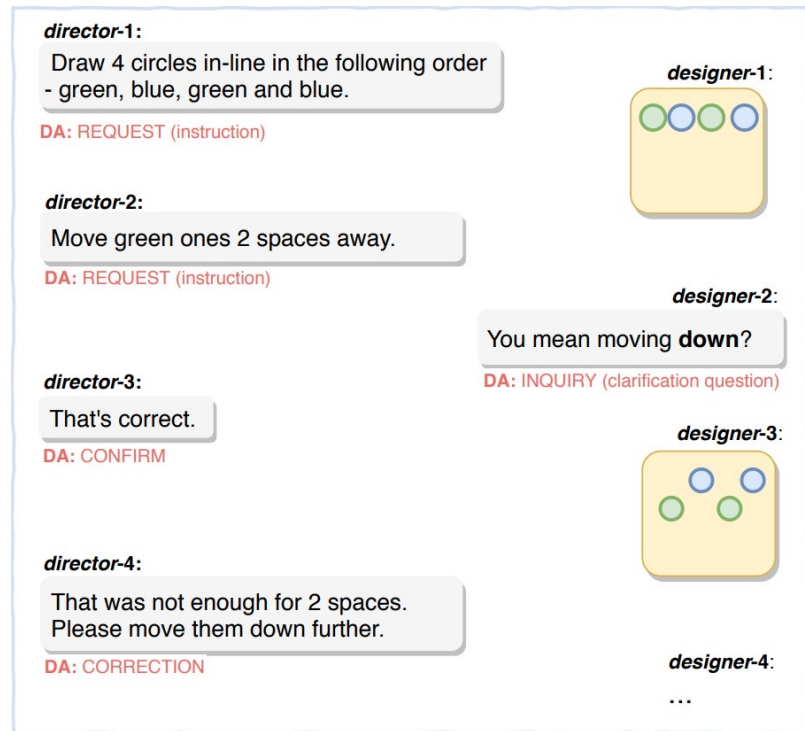
Text-based image editing
[Chen et al., 2018]



Dialogue-based image retrieval
[Guo et al., 2018]

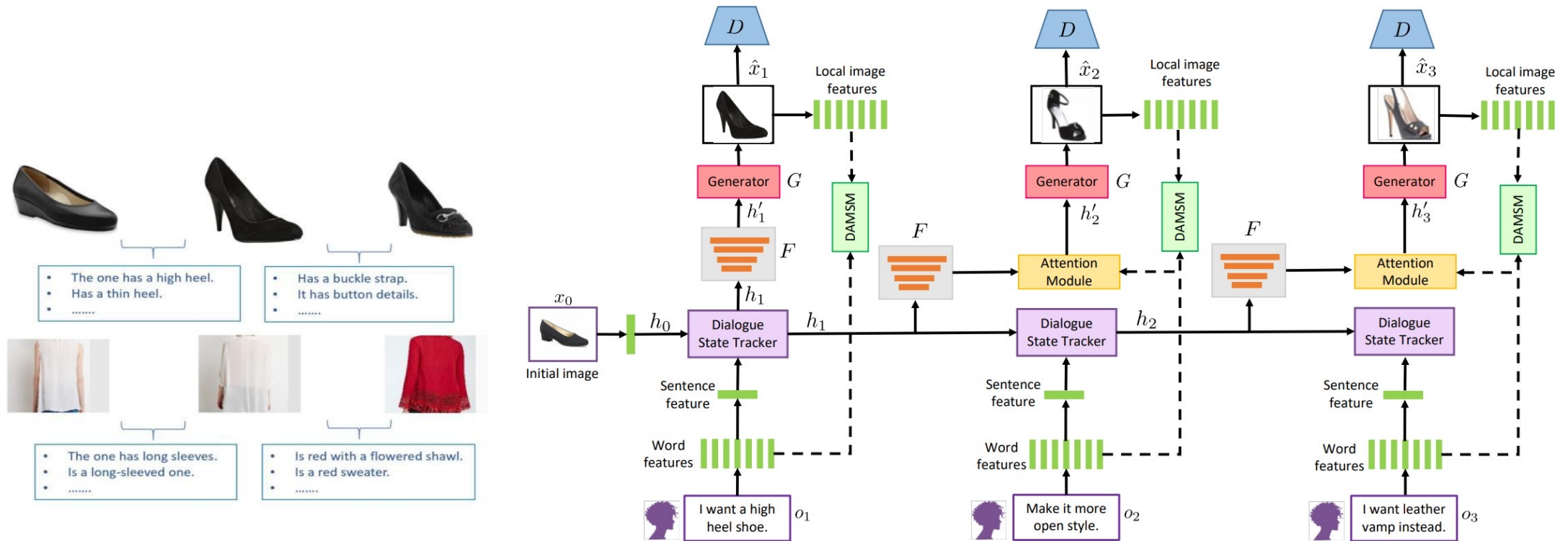
Chat-crowd

- A Dialog-based Platform for Visual Layout Composition



SeqAttnGAN

- Two new datasets: Zap-Seq and DeepFashion-Seq
- Extended from AttnGAN using sequential attention





The University of Texas at Austin
**Electrical and Computer
Engineering**
Cockrell School of Engineering