

**Spring 2023**

# ADVANCED TOPICS IN COMPUTER VISION

---

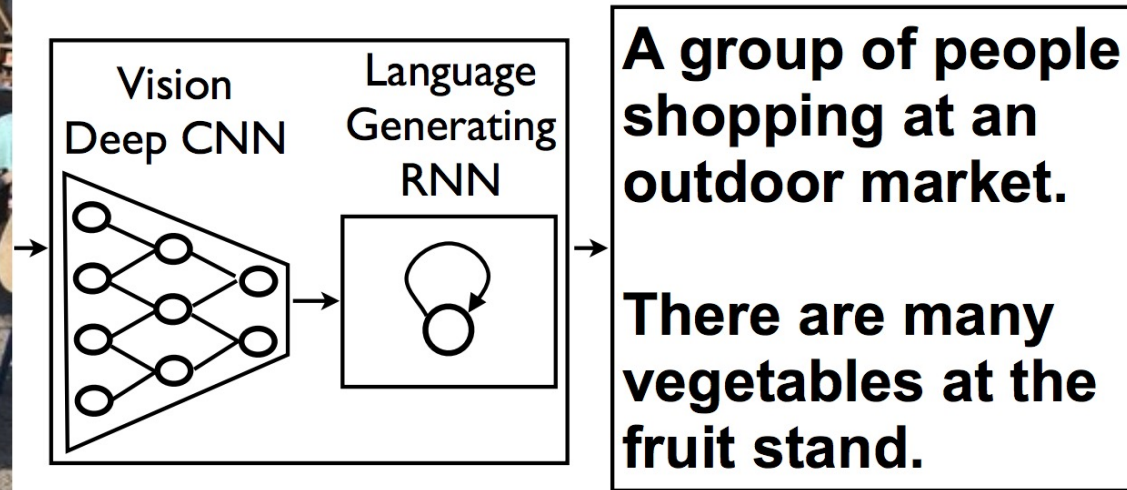
**Atlas Wang**

Assistant Professor, The University of Texas at Austin

**Visual Informatics Group@UT Austin**

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

# Vision + Language: Applications (1)



Visual Captioning: Vinyals et al. 2015

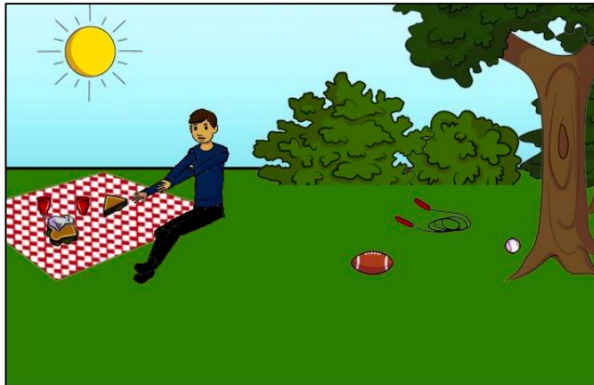
# Vision + Language: : Applications (2)



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)

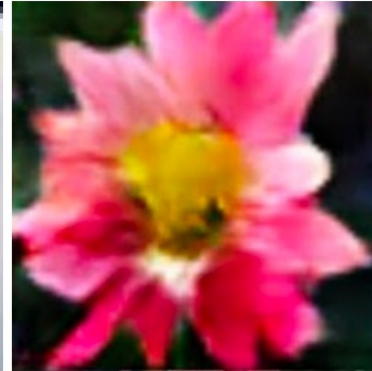
This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face



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

a man is eating a hot dog in a crowd



A fun video running visual captioning model real-time made by Kyle McDonald. Source: <https://vimeo.com/146492001>

# Image Captioning with CNN-LSTM

- Problem Formulation

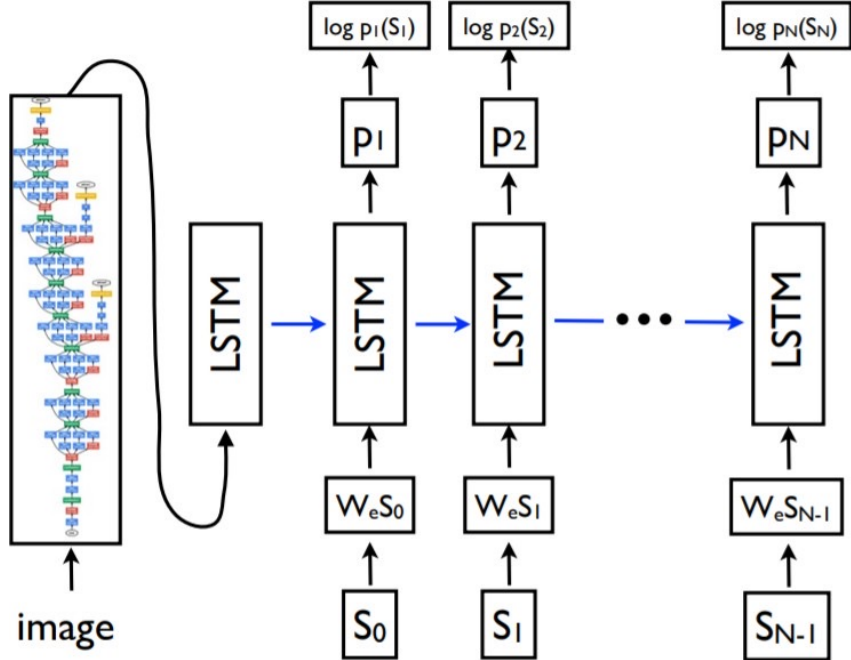
$$\theta^* = \arg \max_{\theta} \sum_{(I,S)} \log p(S|I; \theta)$$

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

- The Encoder-Decoder framework



“Show and Tell”

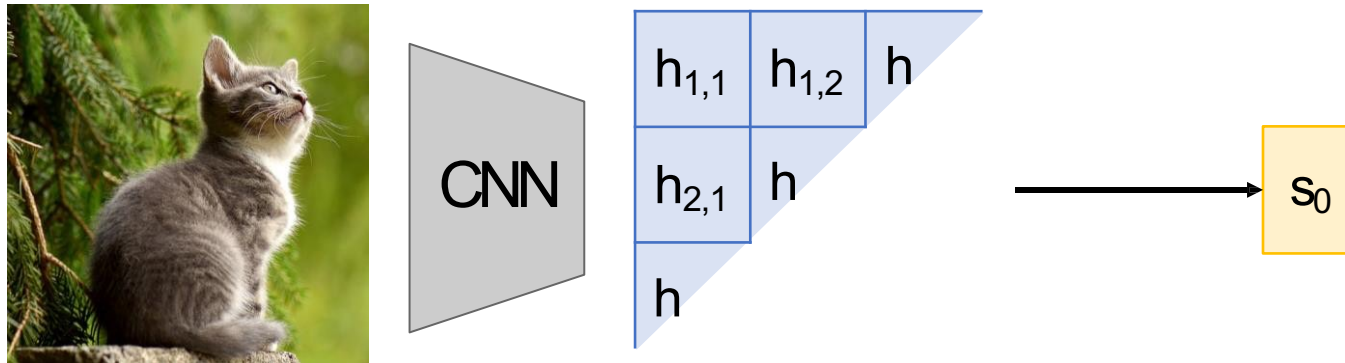




# Image Captioning with Soft Attention

- 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



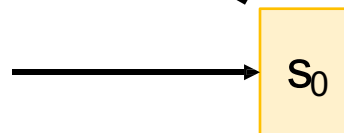
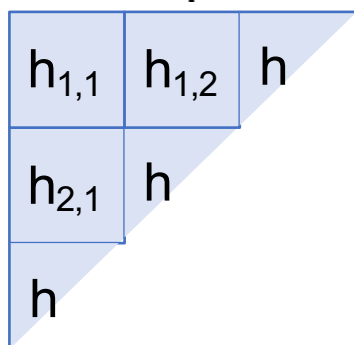
Use a CNN to compute a grid of features for an image

# Image Captioning with Soft Attention

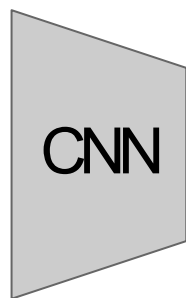
$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

Alignment scores

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

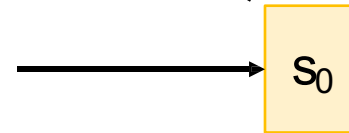
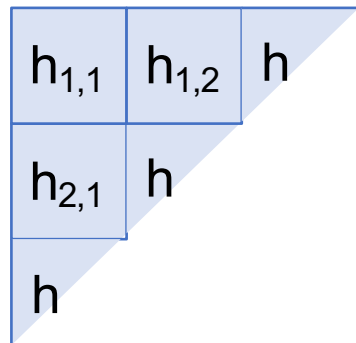
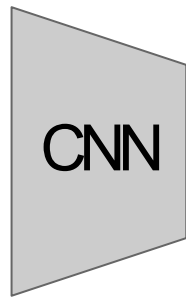
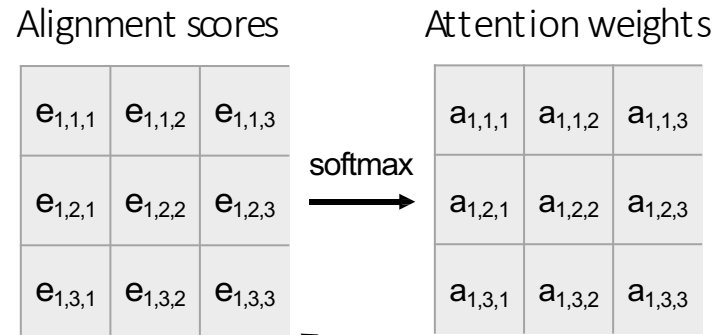


Use a CNN to compute a grid of features for an image



# Image Captioning with Soft Attention

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



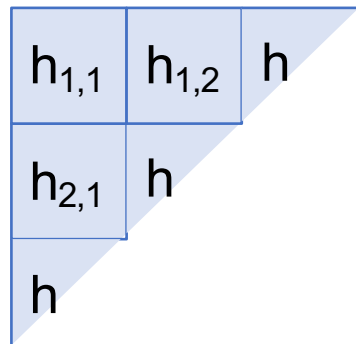
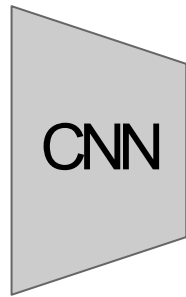
Use a CNN to compute a grid of features for an image

# Image Captioning with Soft Attention

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

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

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



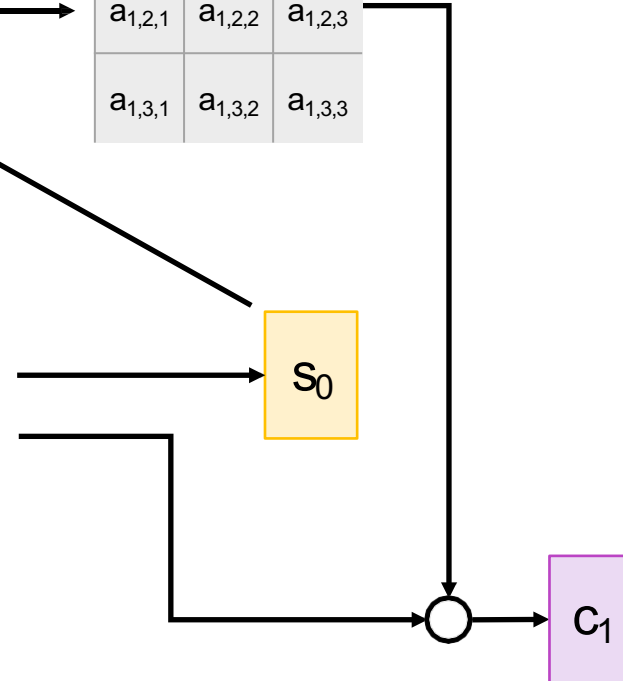
Alignment scores

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

softmax

Attention weights

$a_{1,1,1}$	$a_{1,1,2}$	$a_{1,1,3}$
$a_{1,2,1}$	$a_{1,2,2}$	$a_{1,2,3}$
$a_{1,3,1}$	$a_{1,3,2}$	$a_{1,3,3}$



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

# Image Captioning with Soft Attention

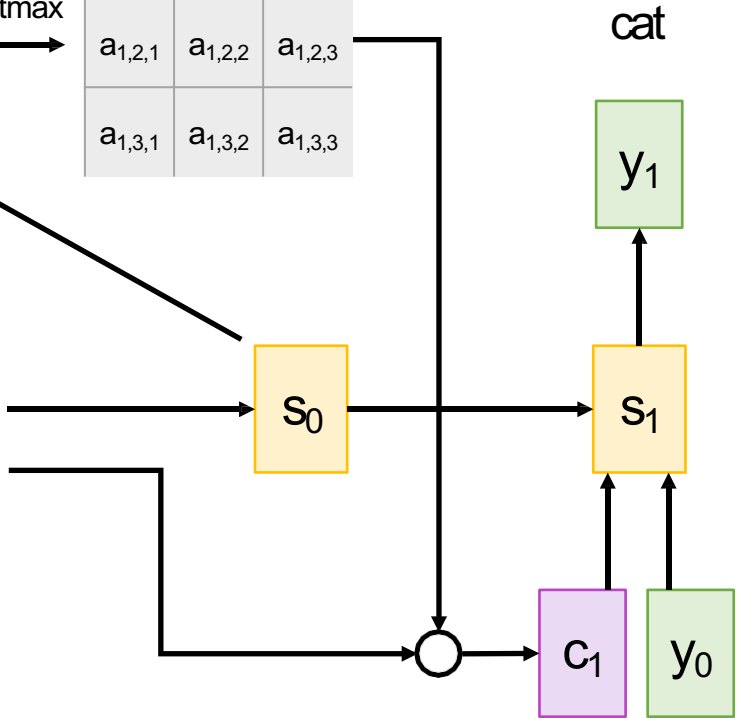
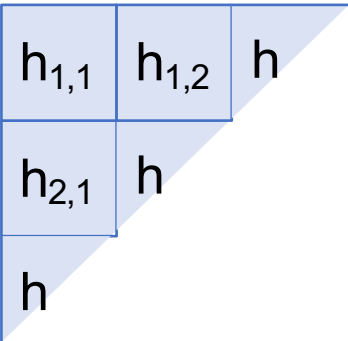
$$e_{t,i,j} = f_{att}(s_{t-1}, h_{i,j})$$

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

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



CNN



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

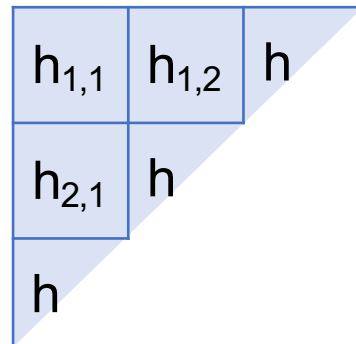
[START]

# Image Captioning with Soft Attention

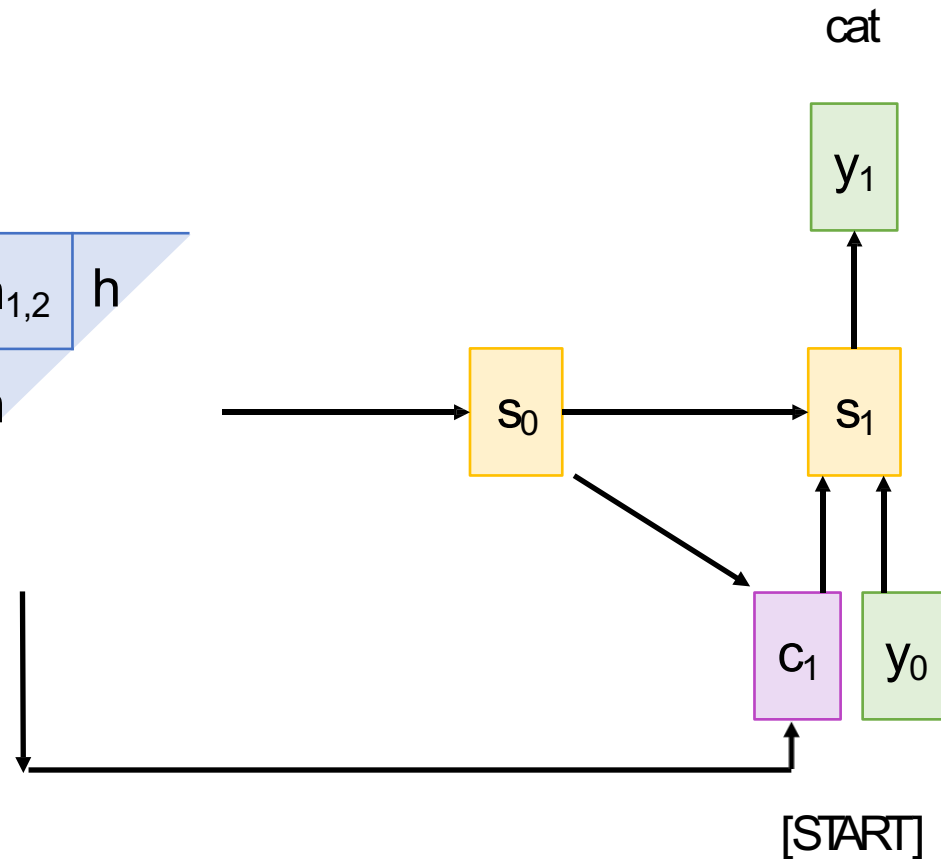
$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$
$$a_{t,:} = \text{softmax}(e_{t,:})$$
$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



CNN



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



# Image Captioning with Soft Attention

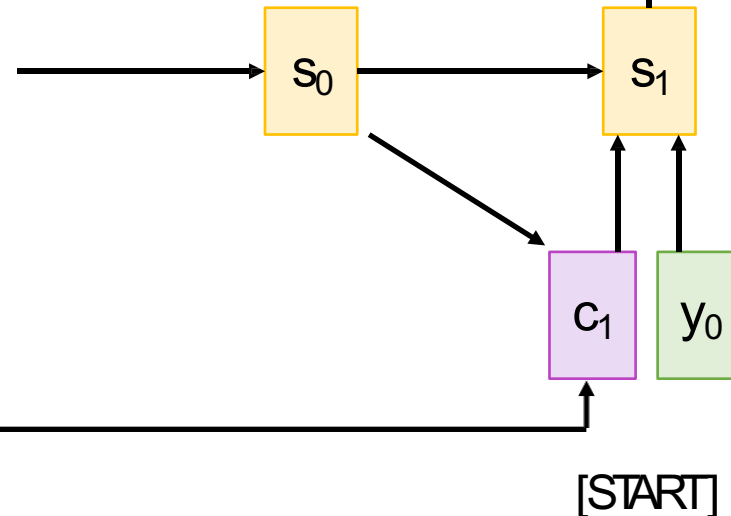
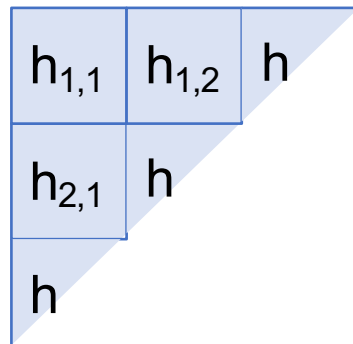
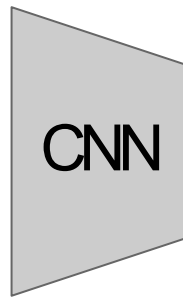
$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

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

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

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



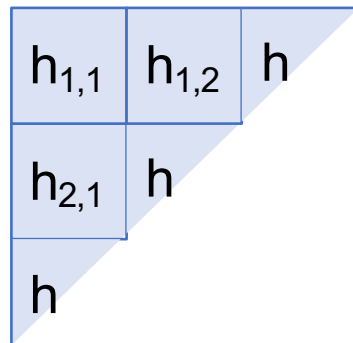
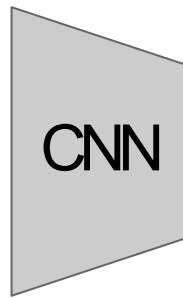
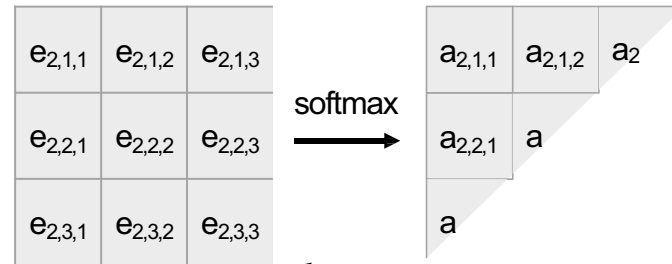
# Image Captioning with Soft Attention

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

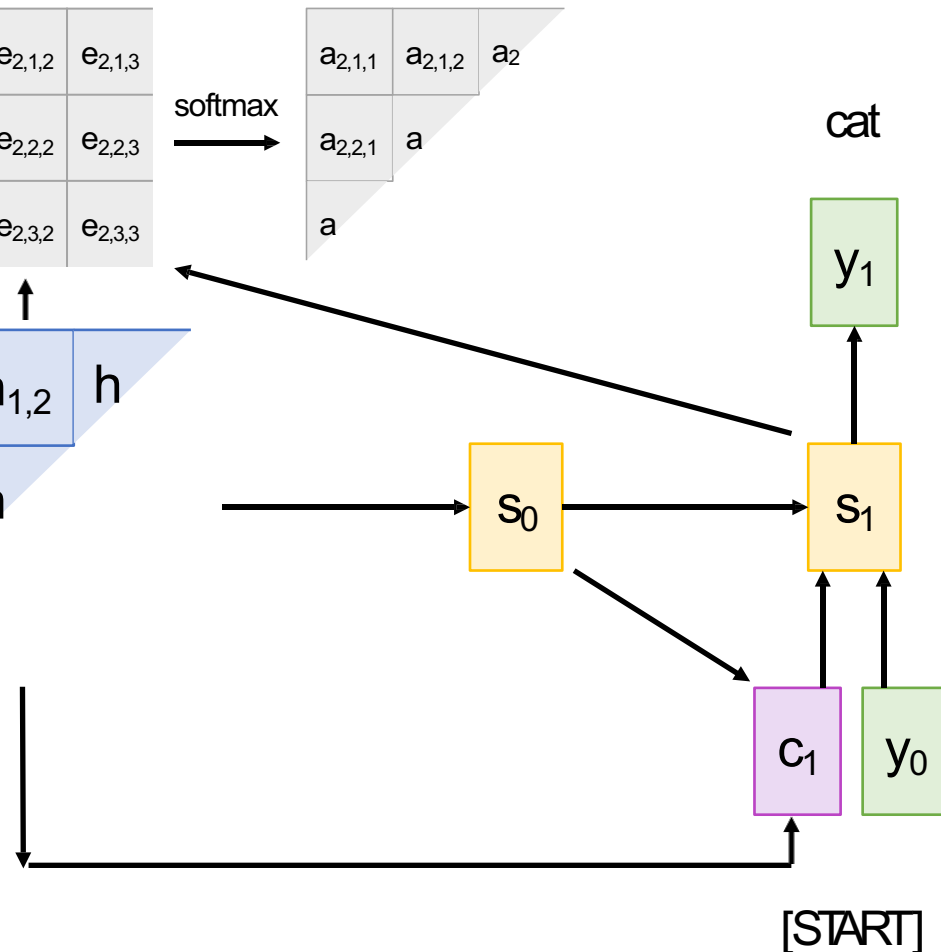
$$a_{t,:} = \text{softmax}(e_{t,:})$$

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

Alignment scores      Attention weights



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



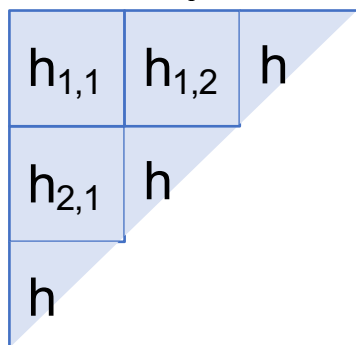
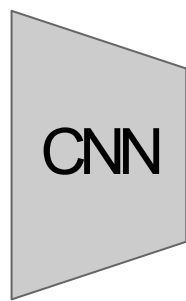
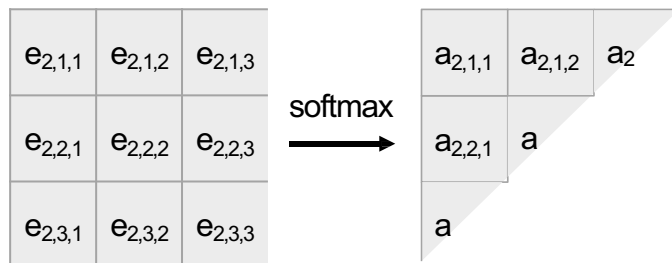
# Image Captioning with Soft Attention

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

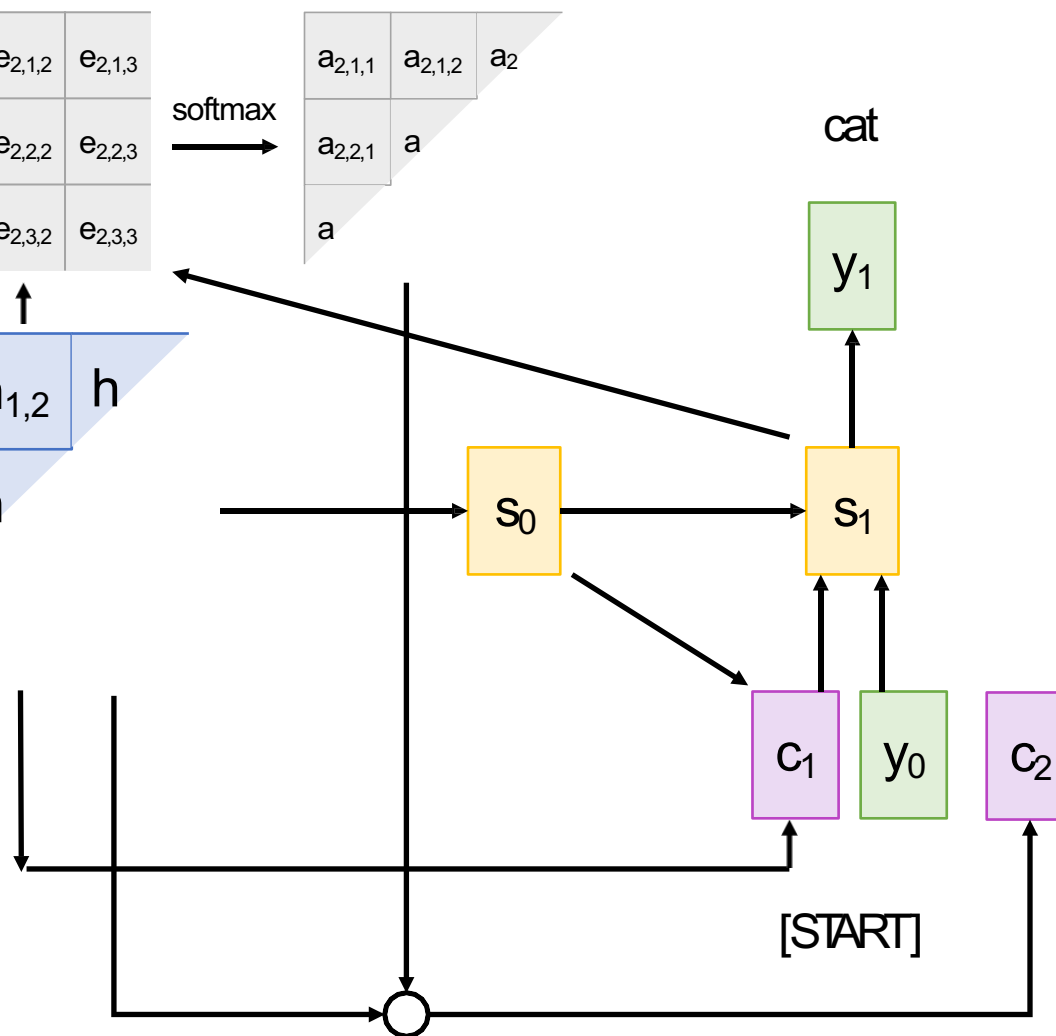
$$a_{t,:} = \text{softmax}(e_{t,:})$$

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

Alignment scores      Attention weights



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



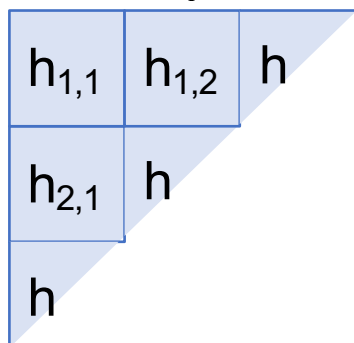
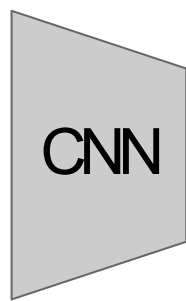
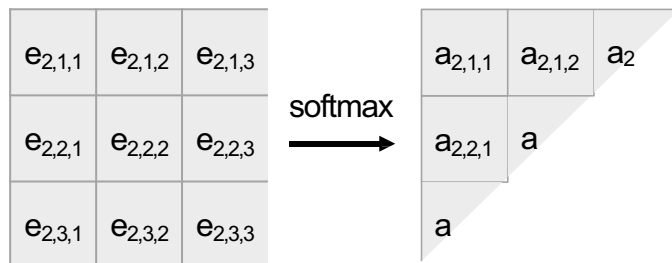
# Image Captioning with Soft Attention

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

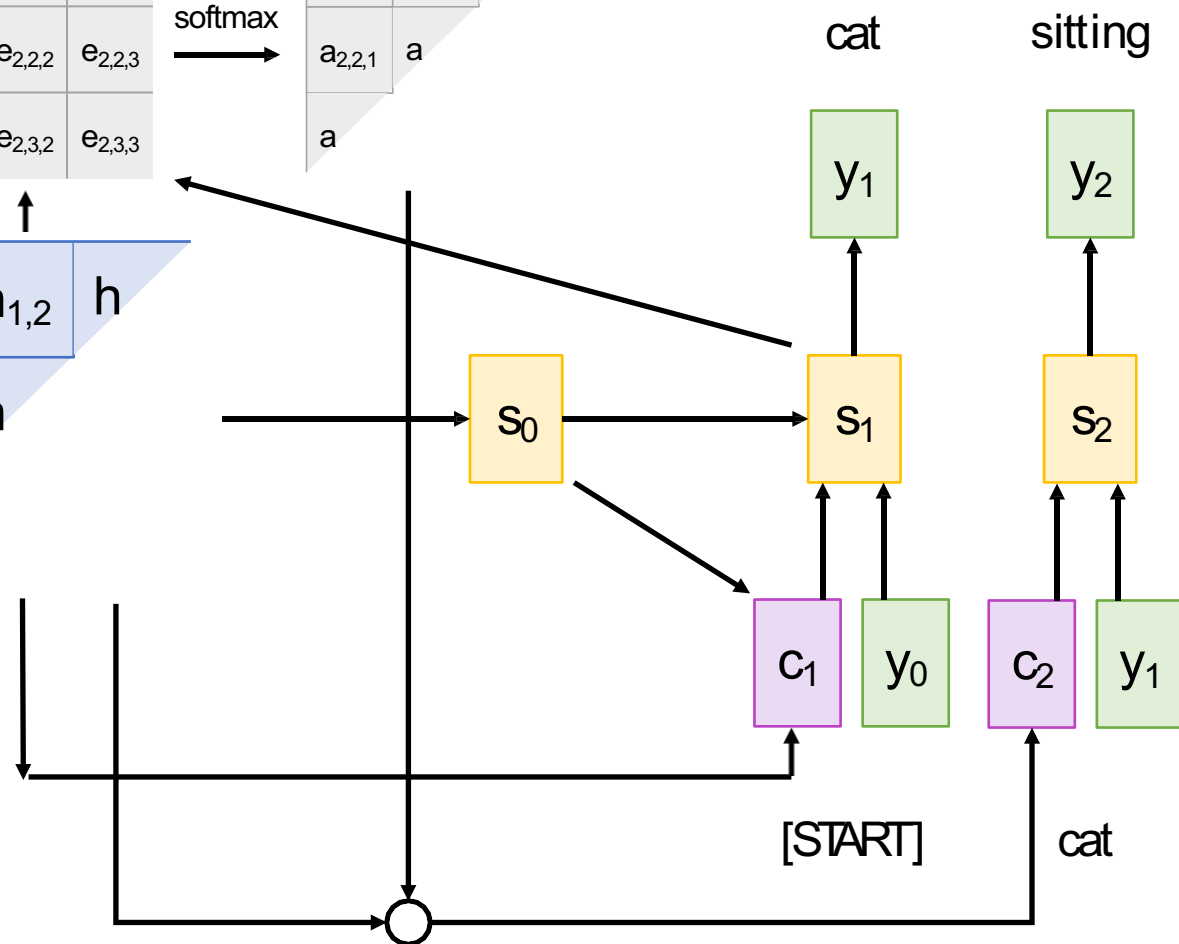
$$a_{t,:} = \text{softmax}(e_{t,:})$$

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

Alignment scores      Attention weights



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



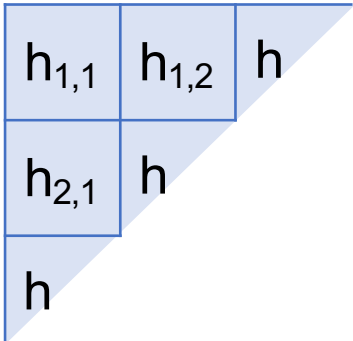
# Image Captioning with Soft Attention

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

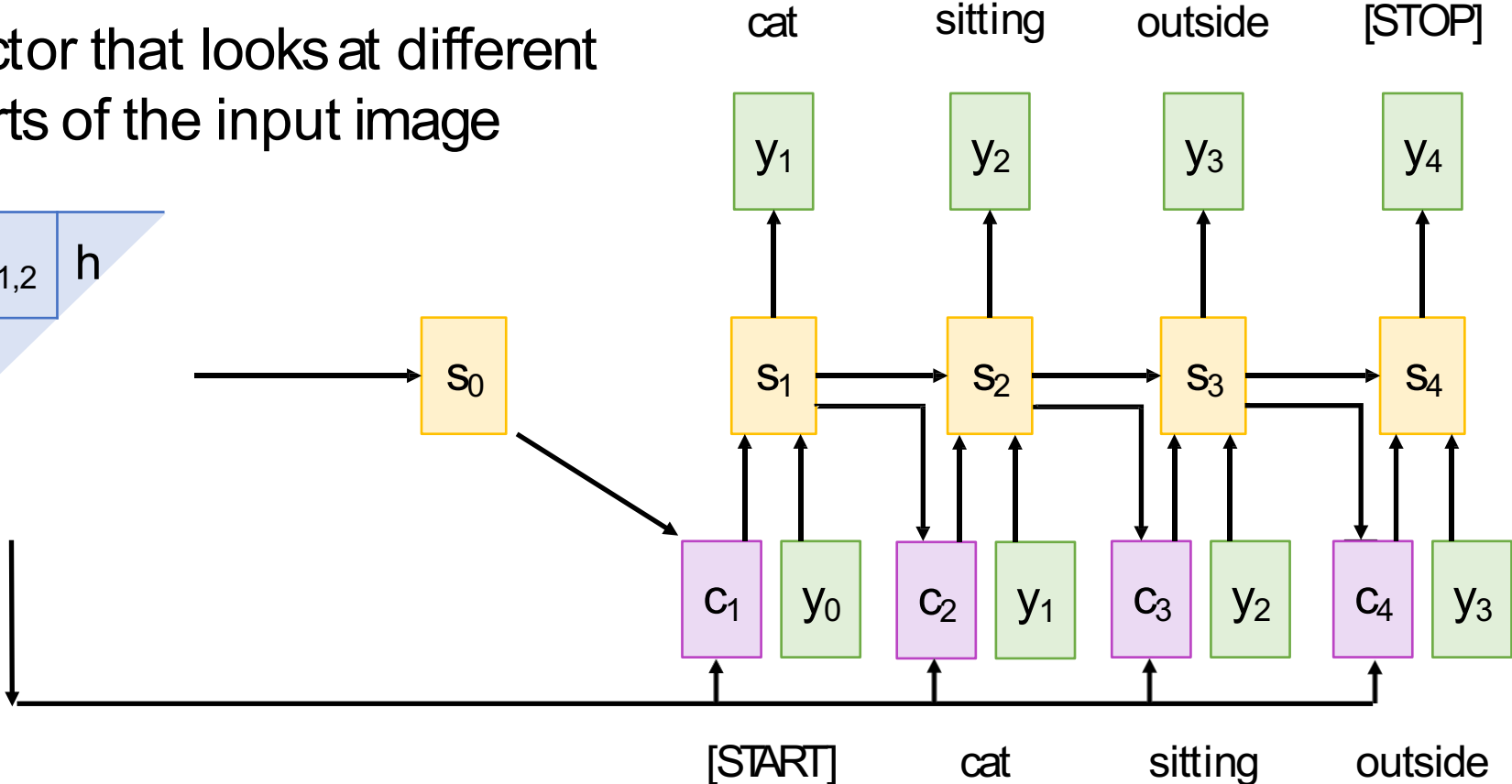
$$a_{t,:} = \text{softmax}(e_{t,:})$$

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

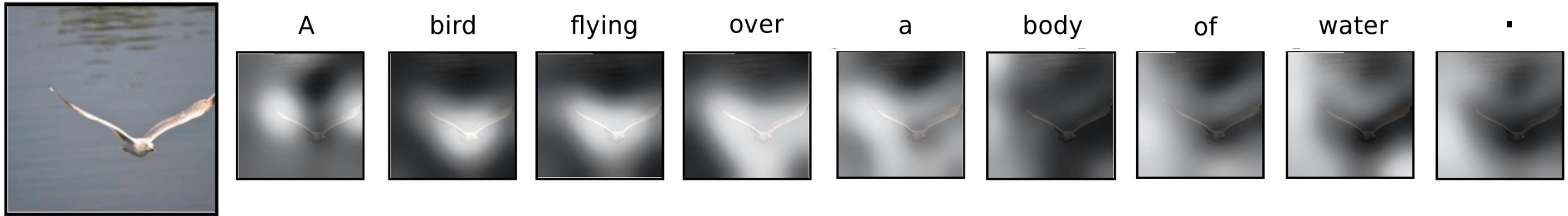
Each timestep of decoder uses a different context vector that looks at different parts of the input image



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



# Image Captioning with Soft Attention



# 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

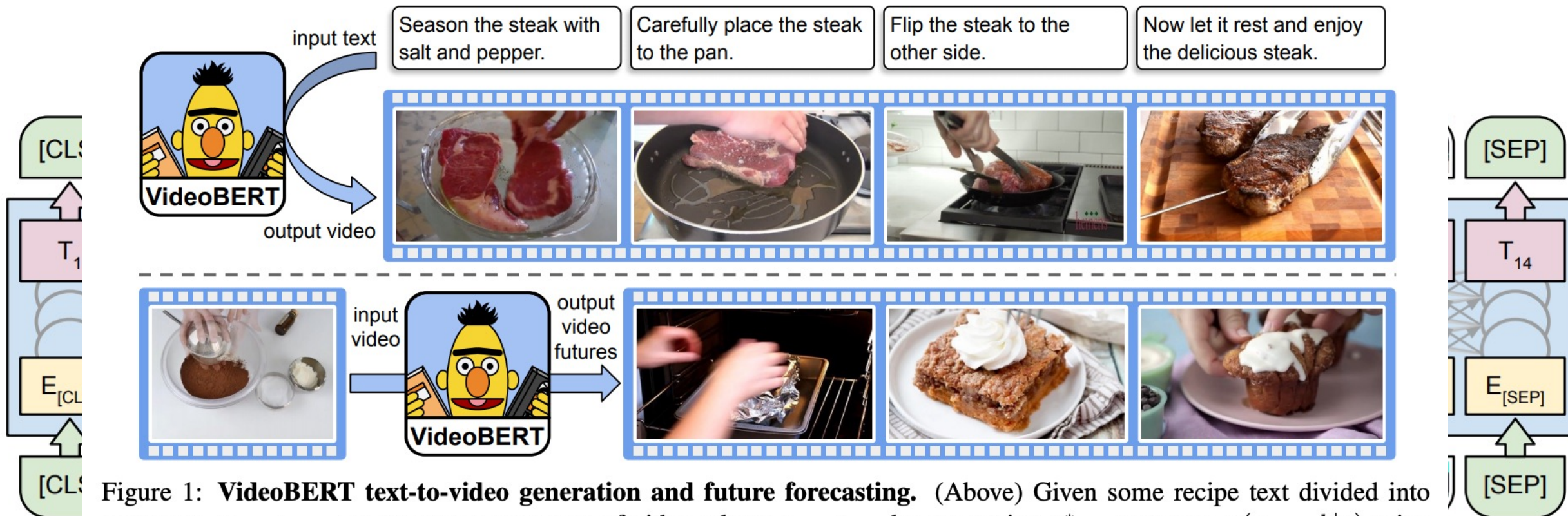


Figure 1: **VideoBERT text-to-video generation and future forecasting.** (Above) Given some recipe text divided into sentences,  $y = y_{1:T}$ , we generate a sequence of video tokens  $x = x_{1:T}$  by computing  $x_t^* = \arg \max_k p(x_t = k|y)$  using VideoBERT. (Below) Given a video token, we show 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

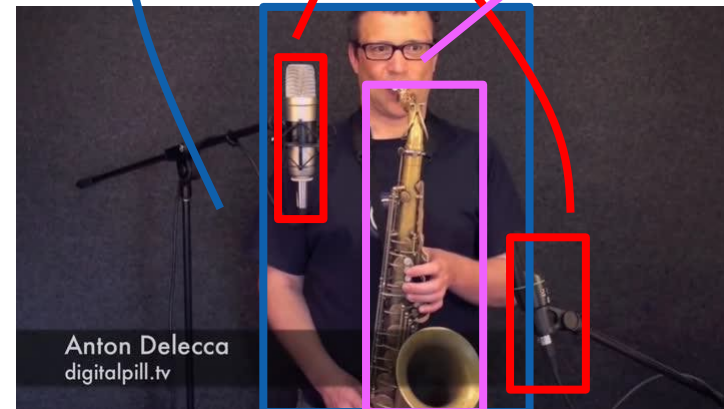
- Essentially, visual description + object grounding or detection
- To achieve better result interpretability, we need grounding!
  - Image domain: Neural Baby Talk, etc.
  - Video domain: Grounded Video Description, etc.
- Requires special dataset that has both description and bounding box



# Single-Frame Annotation



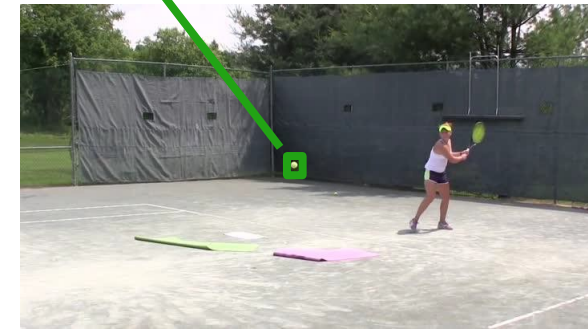
We see a man playing a saxophone  
in front of microphones.



# Multi-Frame Annotation

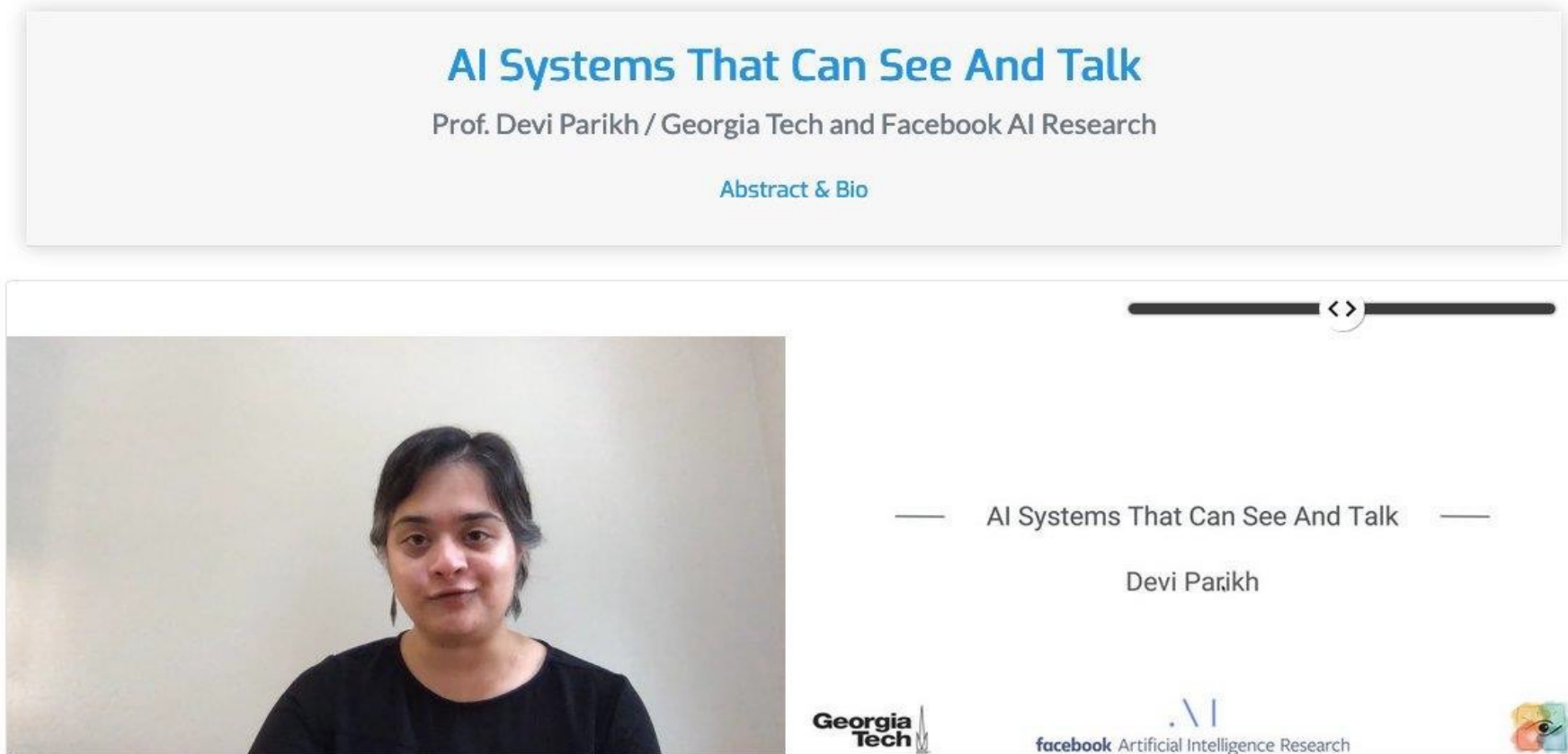


Two women are on a tennis court, showing the technique to posing and hitting the ball.



# Problem Overview (2): VQA and Visual Reasoning

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



The image shows a presentation slide with a video feed of Prof. Devi Parikh on the left. The slide content includes the title 'AI Systems That Can See And Talk', the presenter's name 'Devi Parikh', and logos for Georgia Tech, Facebook AI Research, and a colorful eye icon. A navigation bar at the top right features a double arrow icon.

**AI Systems That Can See And Talk**  
Prof. Devi Parikh / Georgia Tech and Facebook AI Research

[Abstract & Bio](#)

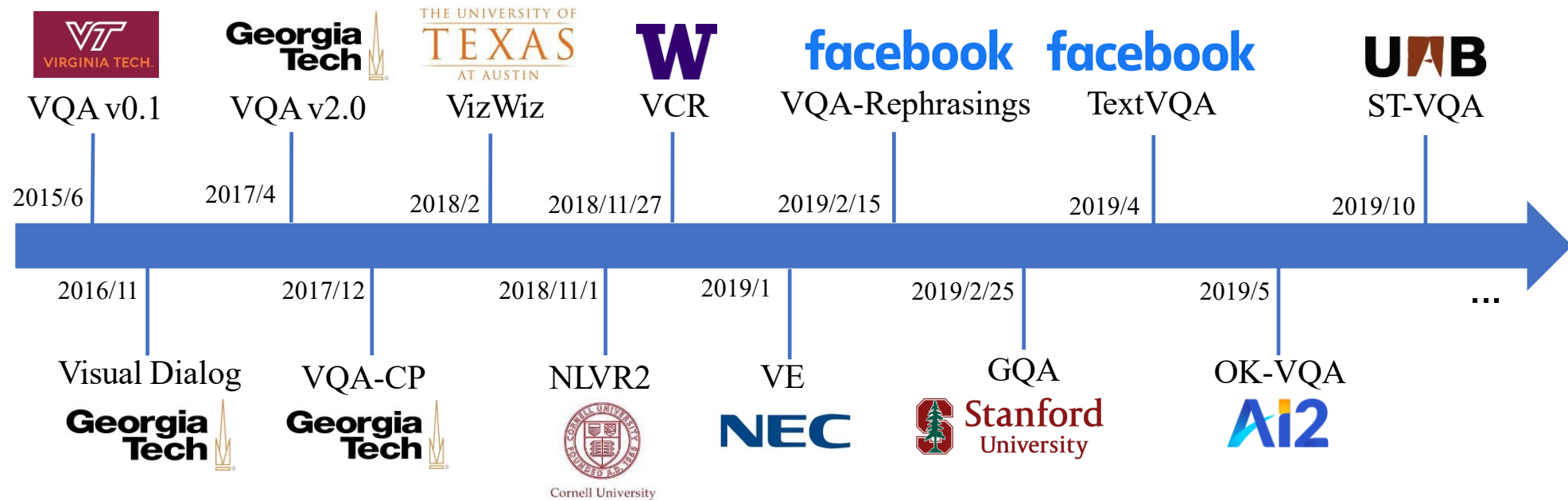
AI Systems That Can See And Talk

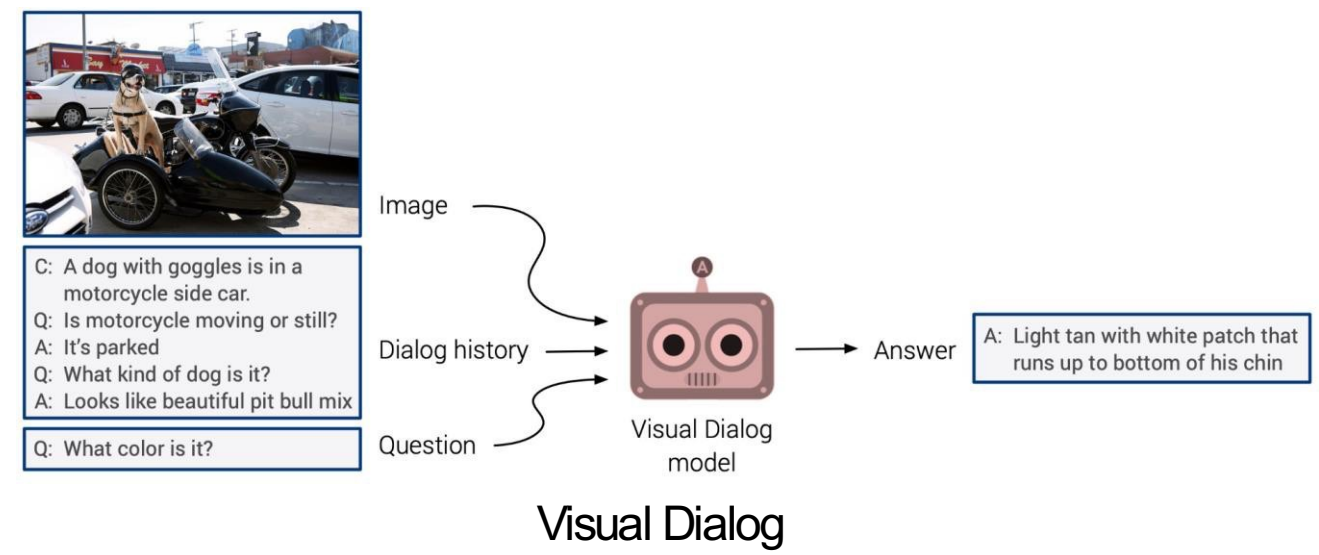
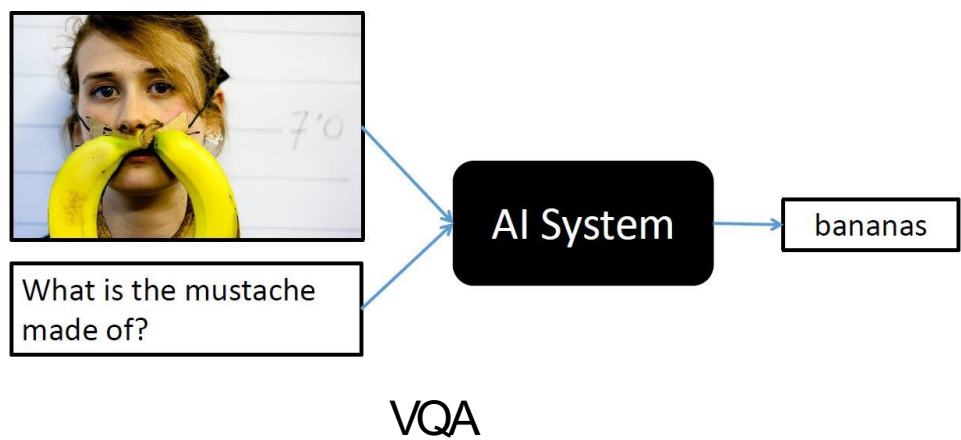
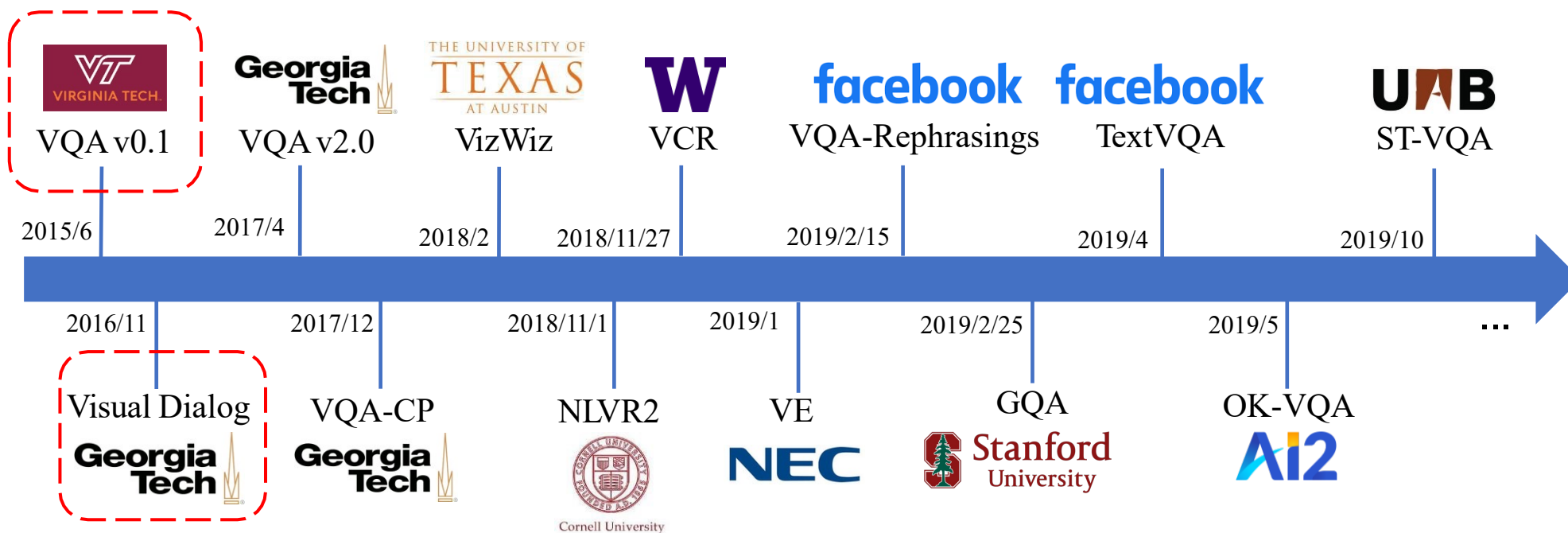
Devi Parikh

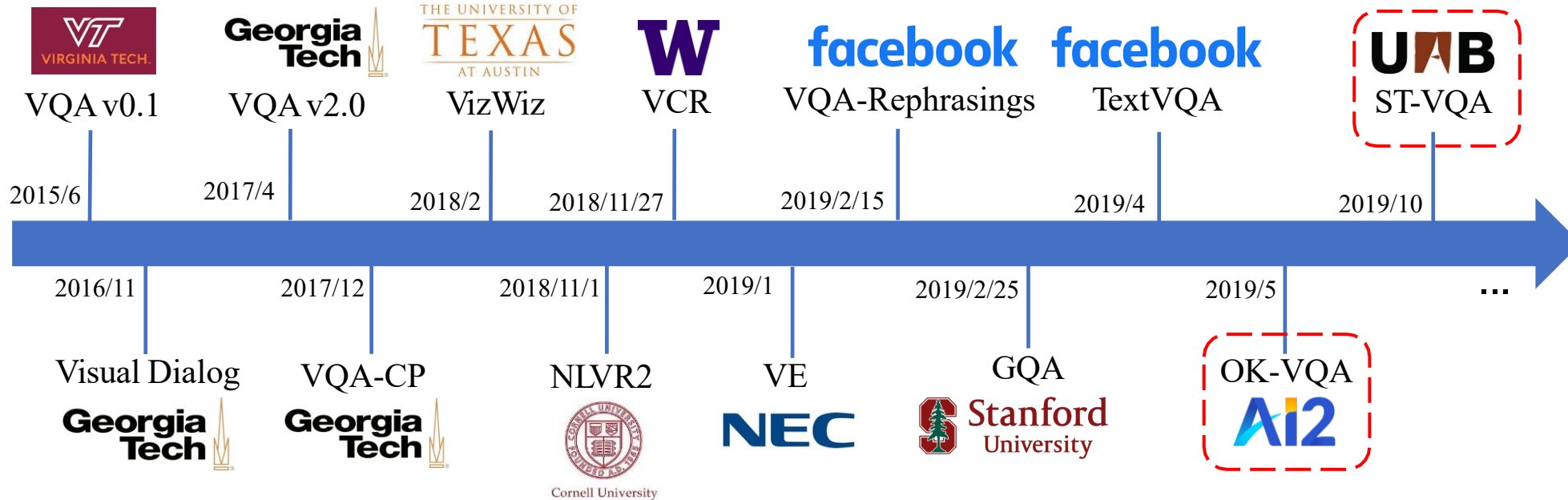
Georgia Tech | facebook Artificial Intelligence Research

# 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?

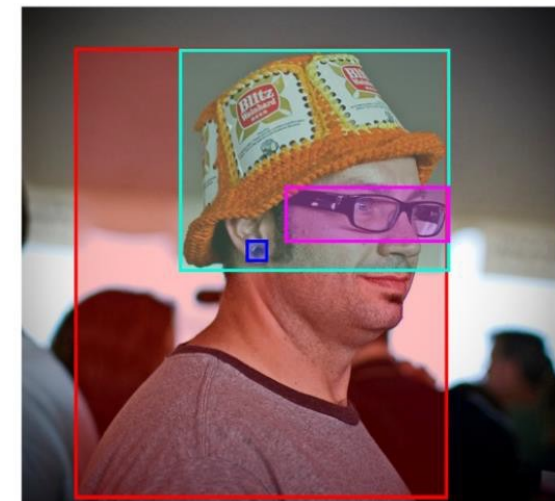
A: Stop

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

short deer



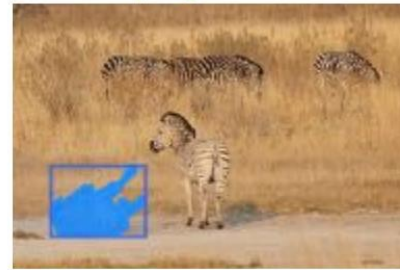
walking people



wipers on trains



zebra lying on savanna



black shirt



hatchback car



mark on chicken



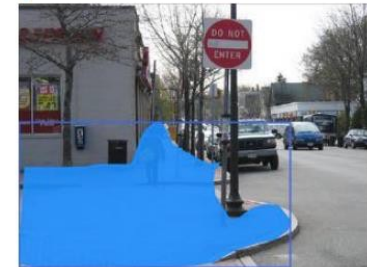
glass bottles



blonde hair



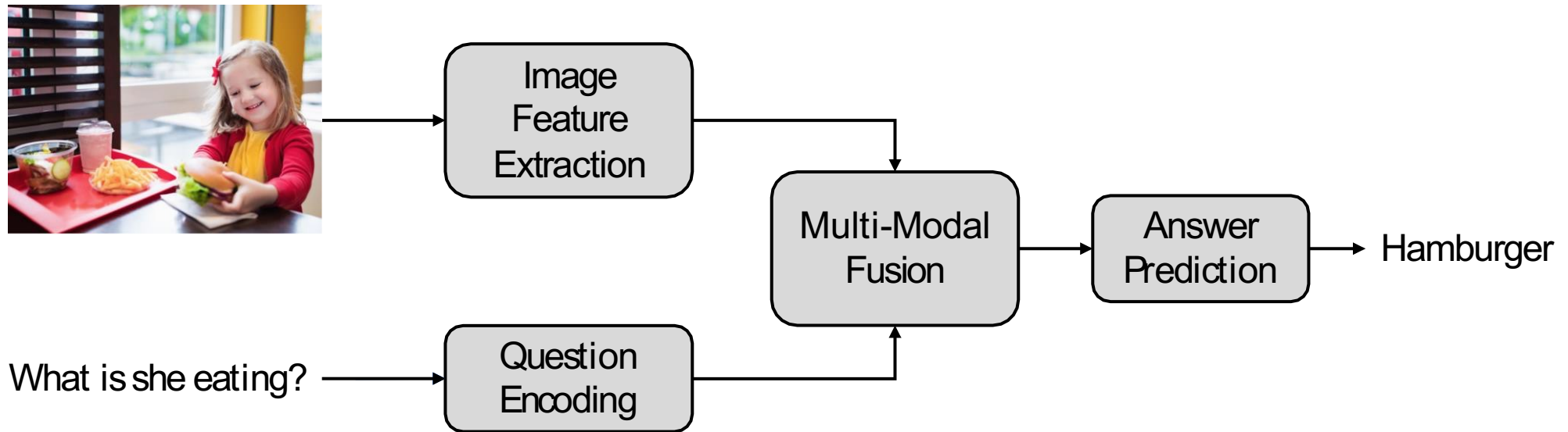
pedestrian crosswalk





# Approach Overview

- How a typical system looks like

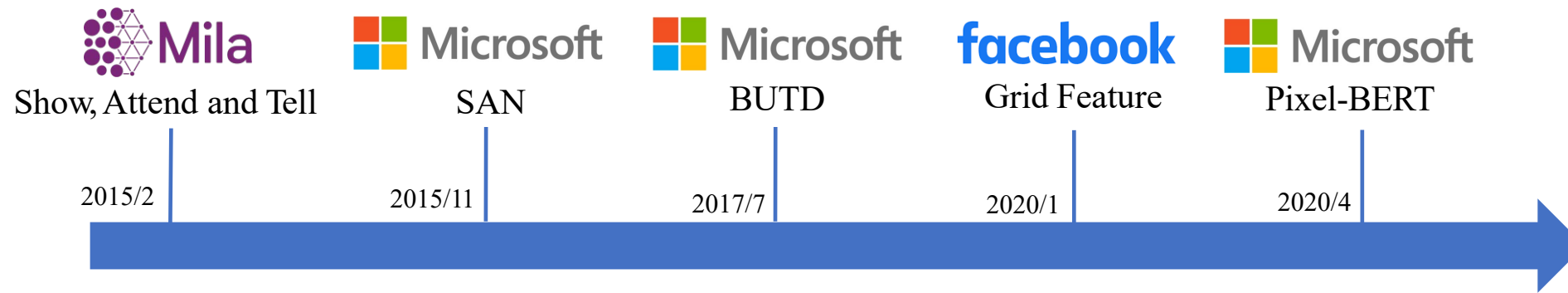


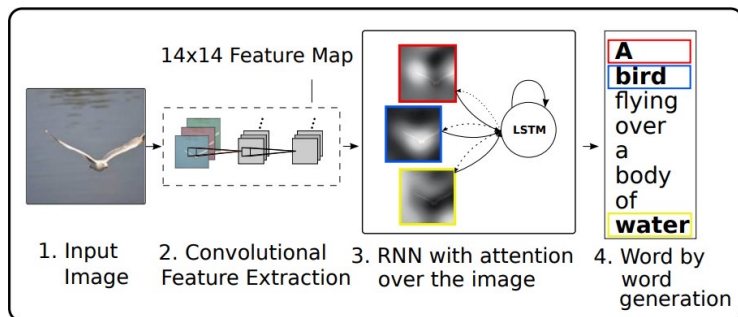
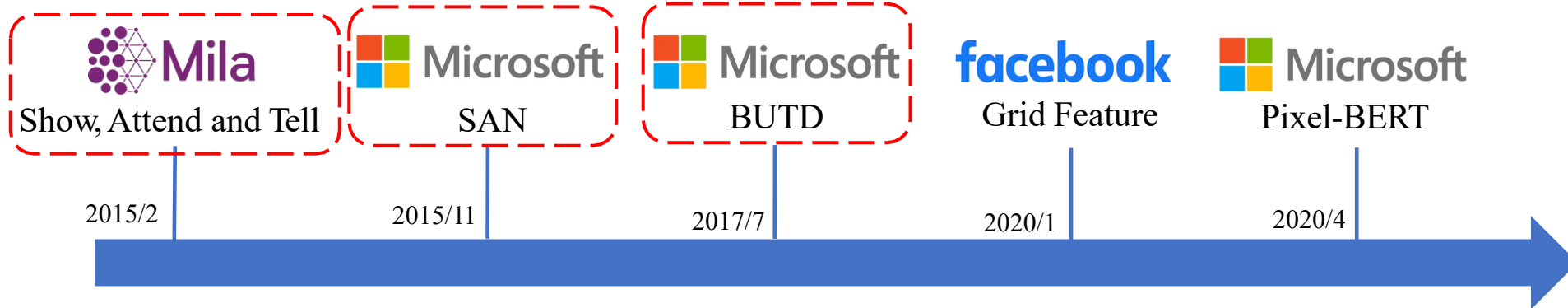
# Research Challenges & Opportunities

- 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

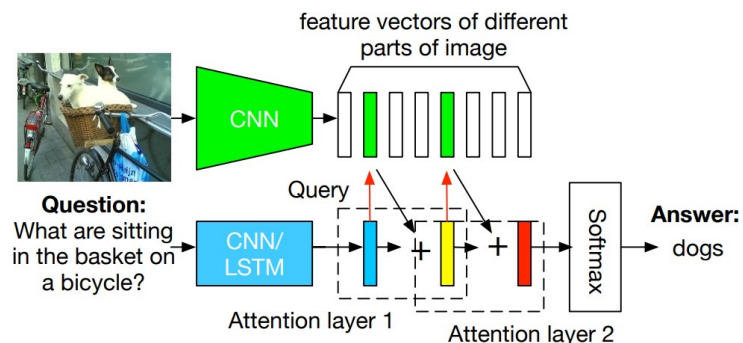
# Better Image Feature Preparation

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

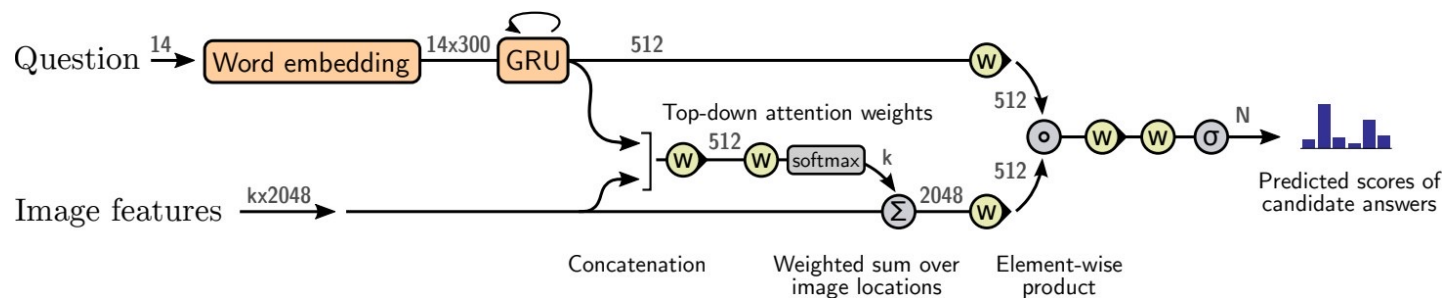




Show, Attend and Tell



Stacked Attention Network



2017 VQA Challenge Winner

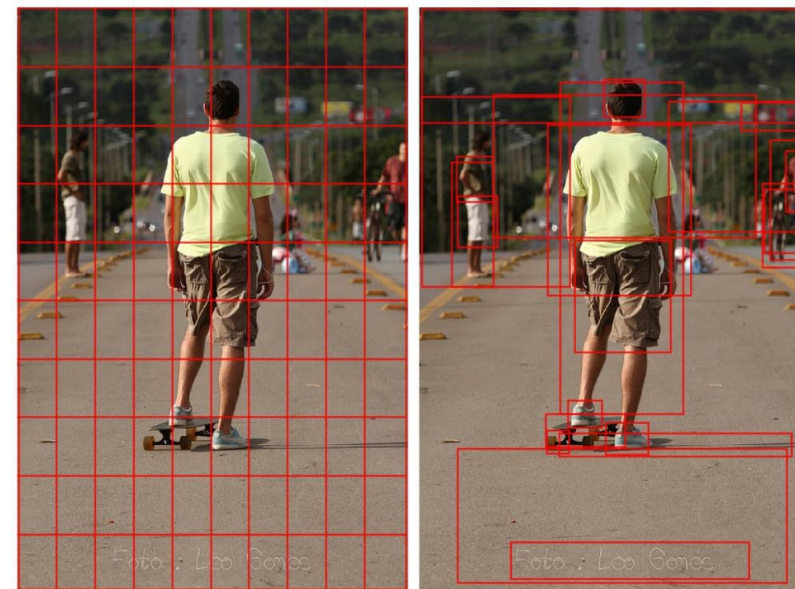


Figure 1. Typically, attention models operate on CNN features corresponding to a uniform grid of equally-sized image regions (left). Our approach enables attention to be calculated at the level of objects and other salient image regions (right).

- 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

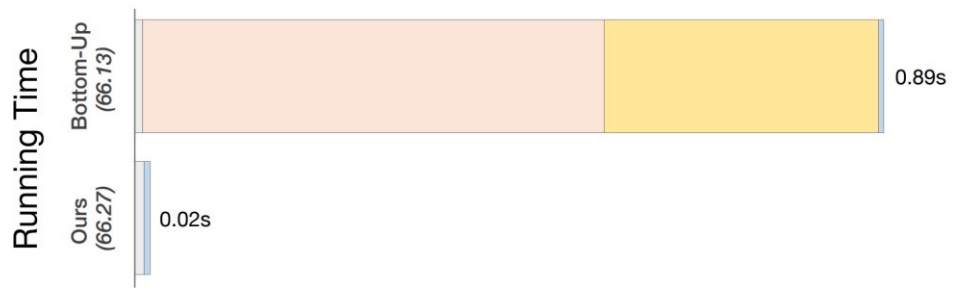
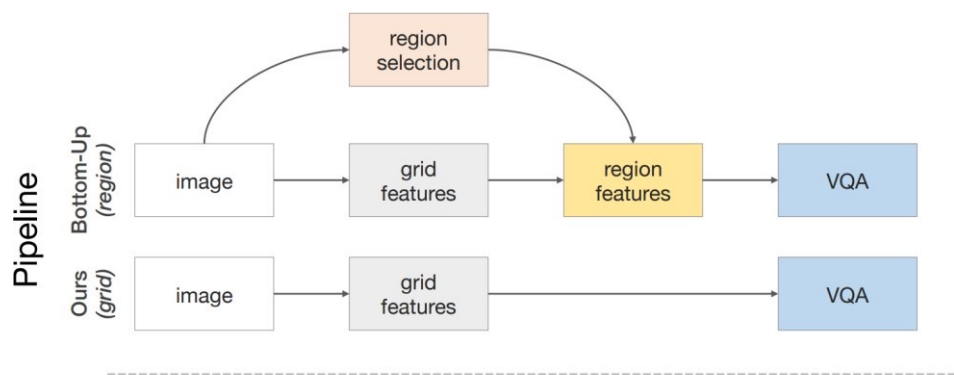
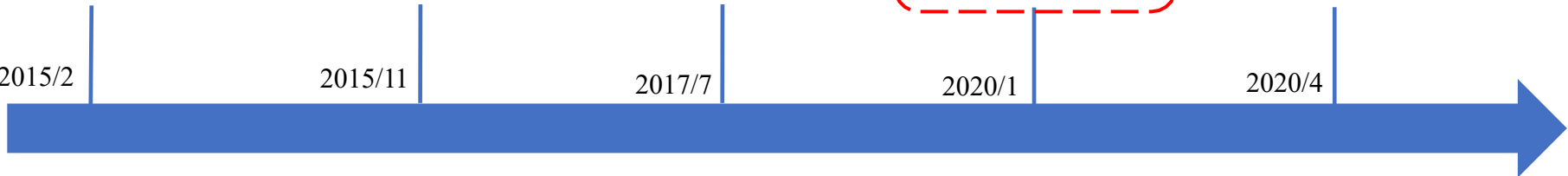
2017/7

 **facebook**  
Grid Feature

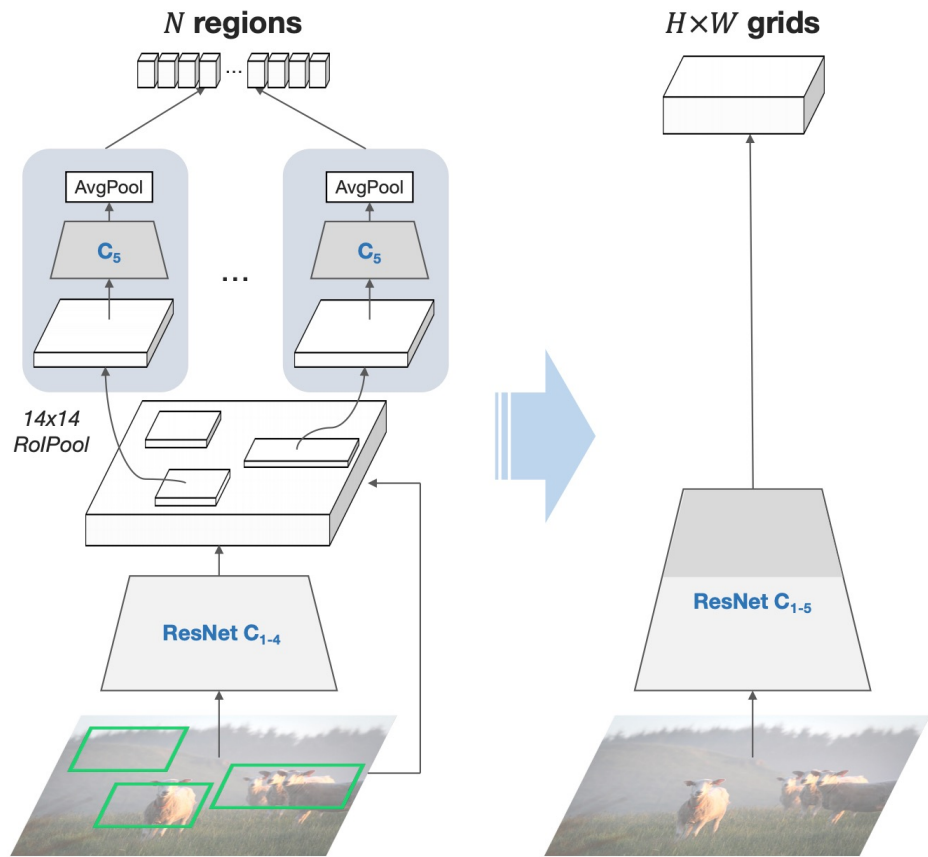
2020/1

 **Microsoft**  
Pixel-BERT

2020/4



**In Defense of Grid Features for VQA**



 **Mila**  
Show, Attend and Tell

2015/2

 **Microsoft**

SAN

2015/11

 **Microsoft**

BUTD

2017/7

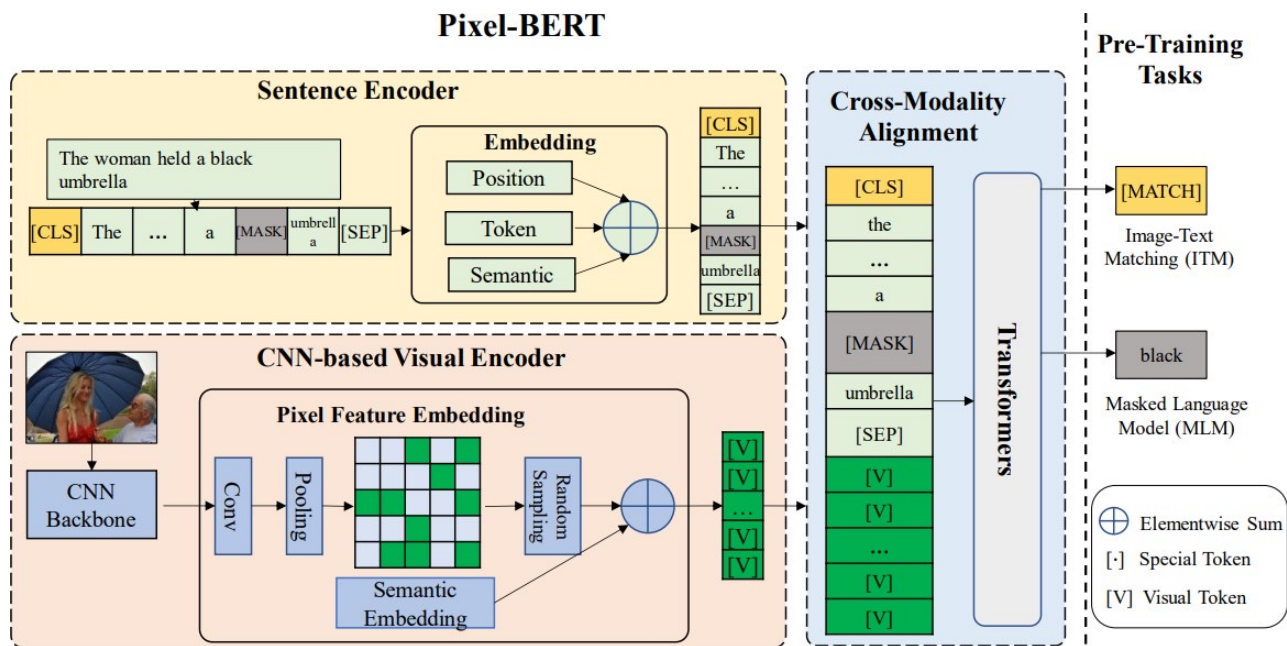
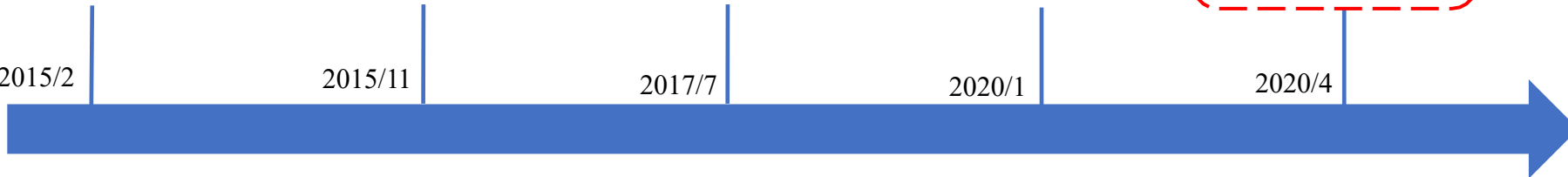
**facebook**

Grid Feature

2020/1

 **Microsoft**  
Pixel-BERT

2020/4



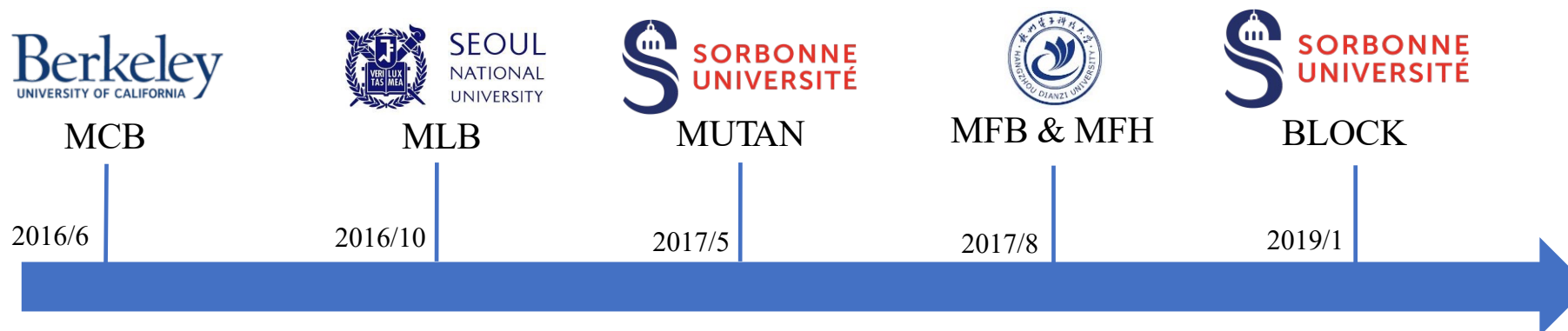
Model	test-dev	test-std
MUTAN[5]	60.17	-
BUTD[2]	65.32	65.67
ViLBERT[21]	70.55	70.92
VisualBERT[19]	70.80	71.00
VLBERT[29]	71.79	72.22
LXMERT[33]	72.42	72.54
UNITER[6]	72.27	72.46
Pixel-BERT (r50)	71.35	71.42
Pixel-BERT (x152)	<b>74.45</b>	<b>74.55</b>

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

[1] Pixel-BERT: Aligning Image Pixels with Text by Deep Multi-Modal Transformers, 2020

# 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





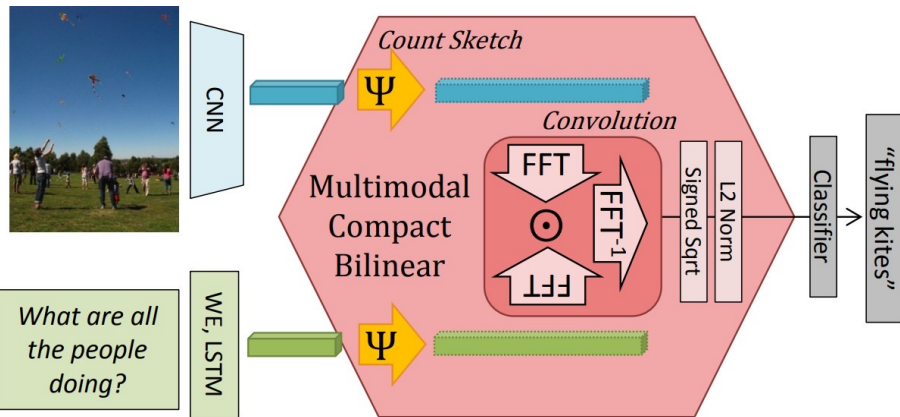
2016/6

2016/10

2017/5

2017/8

2019/1



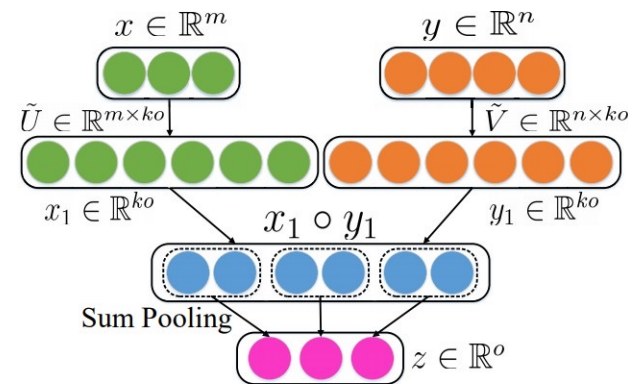
### Multimodal Compact Bilinear Pooling

*2016 VQA Challenge Winner*

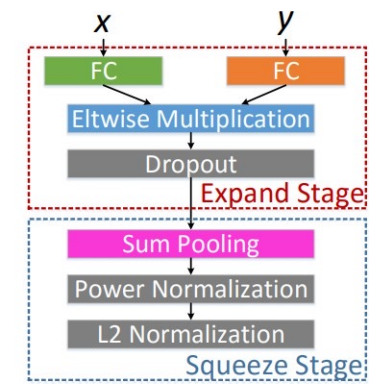
However, the feature after FFT is very high dimensional.

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

### Multimodal Low-rank Bilinear Pooling



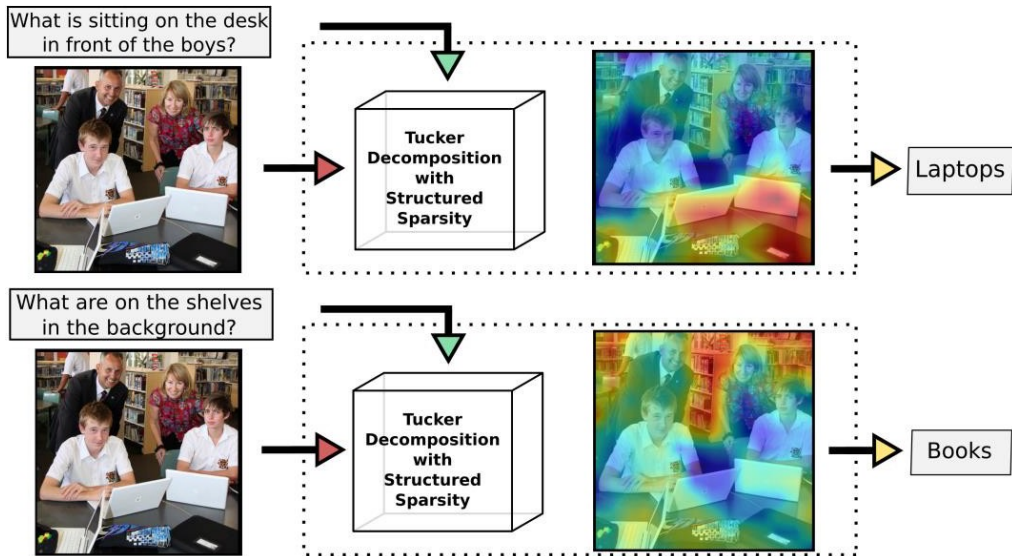
(a) Multi-modal Factorized Bilinear Pooling



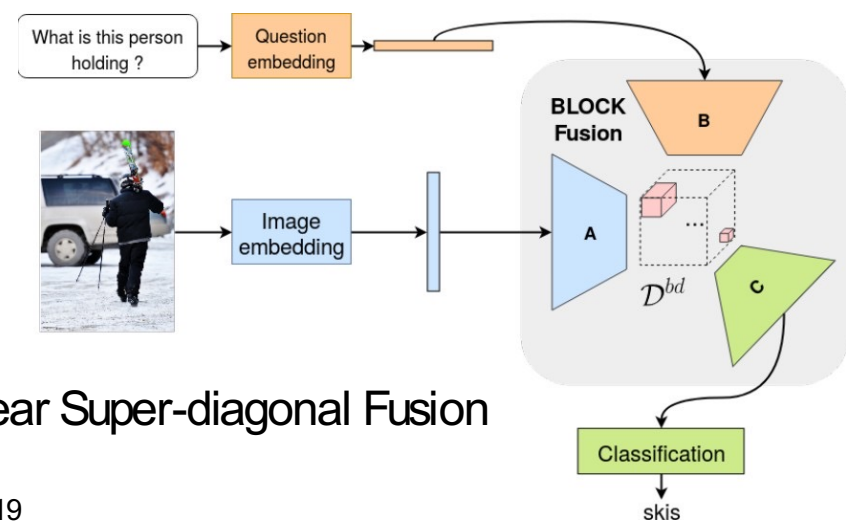
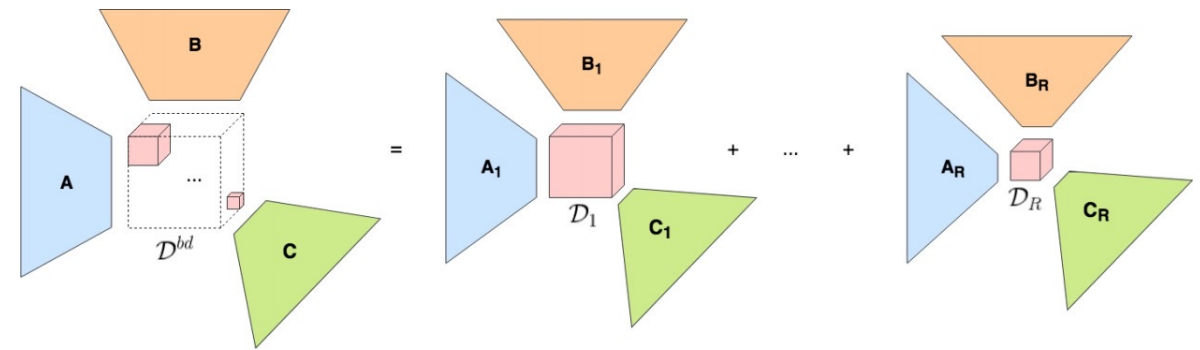
(b) MFB module

- 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

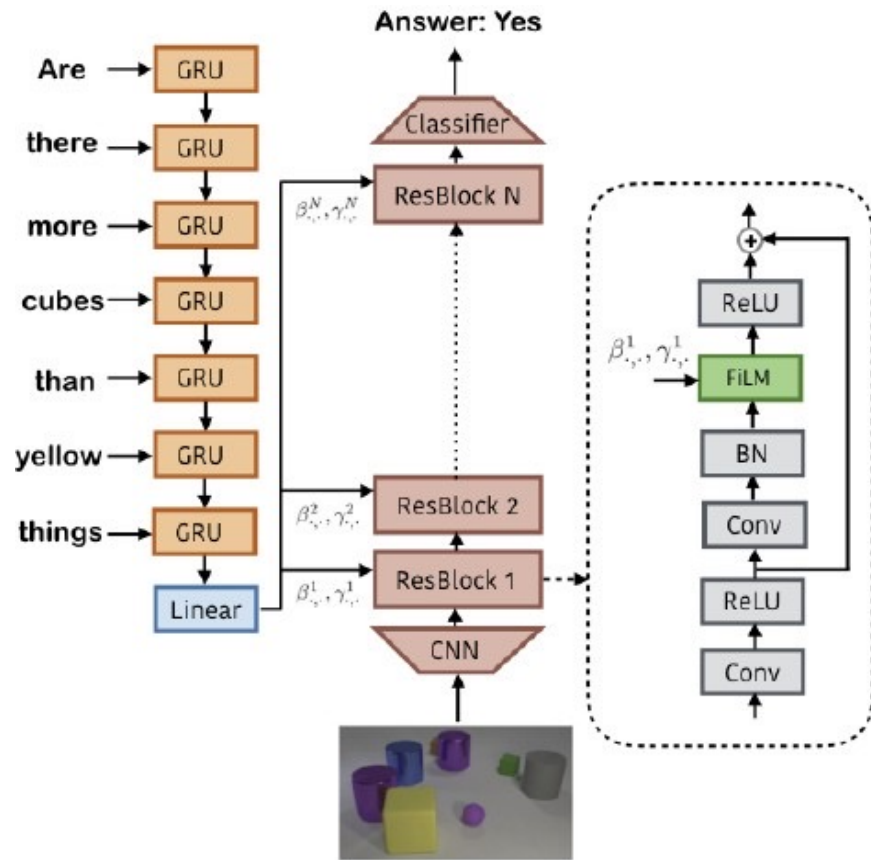


Bilinear Super-diagonal Fusion

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

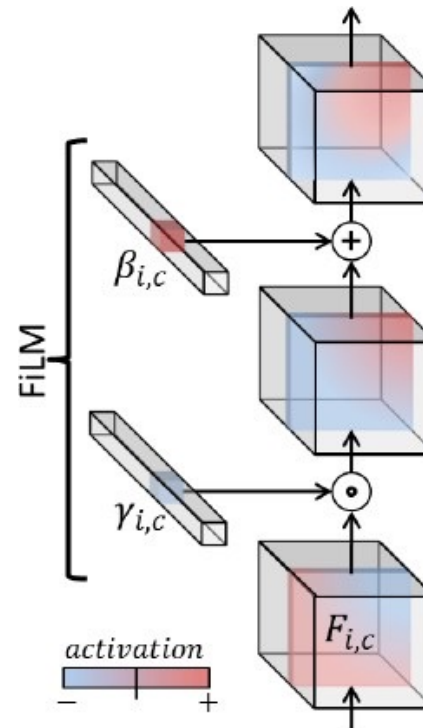
# FiLM: Feature-wise Linear Modulation



$$\gamma_{i,c} = f_c(\mathbf{x}_i) \quad \beta_{i,c} = h_c(\mathbf{x}_i),$$

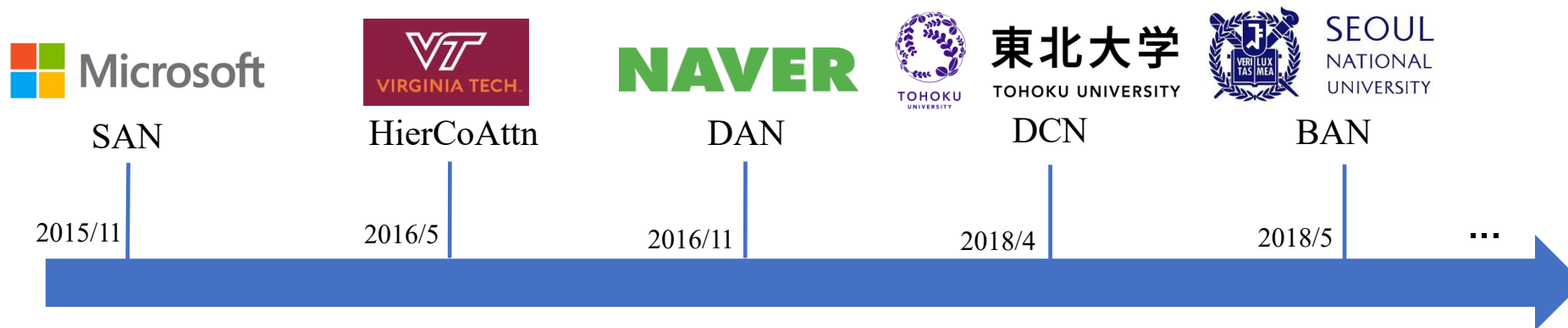
$$FiLM(\mathbf{F}_{i,c} | \gamma_{i,c}, \beta_{i,c}) = \gamma_{i,c} \mathbf{F}_{i,c} + \beta_{i,c}.$$

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

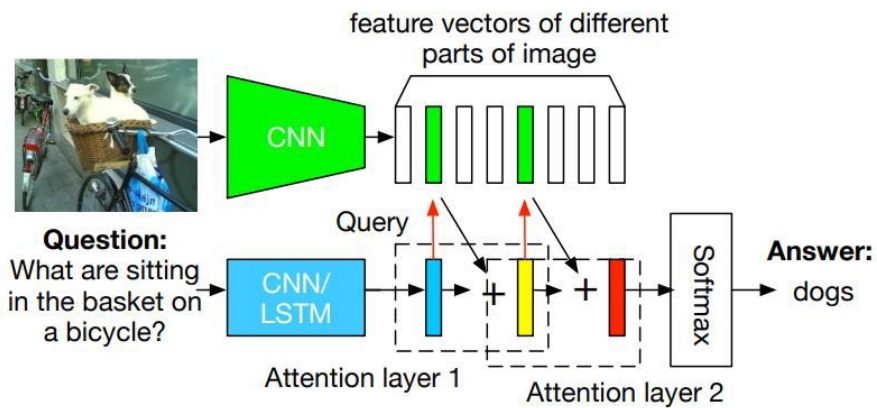


SEOUL  
NATIONAL  
UNIVERSITY

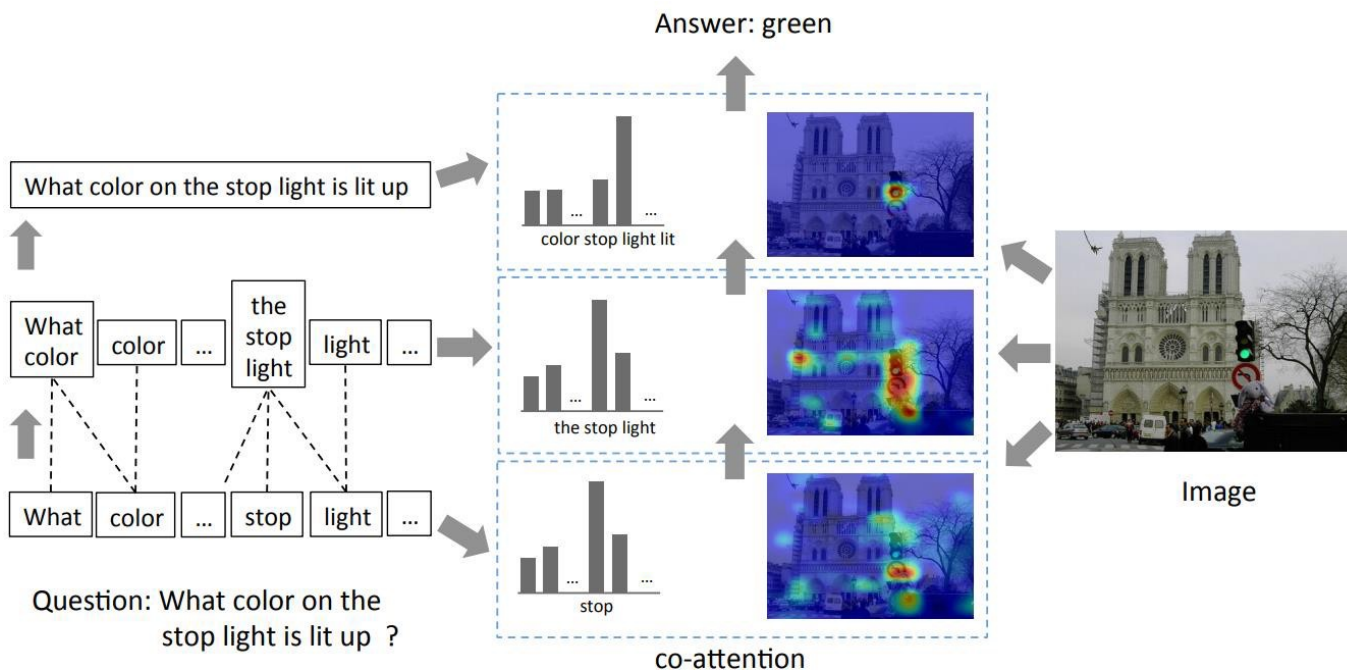
BAN

2018/5

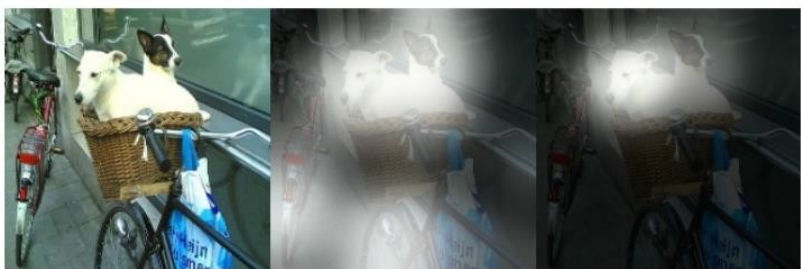
...



(a) Stacked Attention Network for Image QA



## Parallel Co-attention and Alternative Co-attention



(b) Visualization of the learned multiple attention layers.

- 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

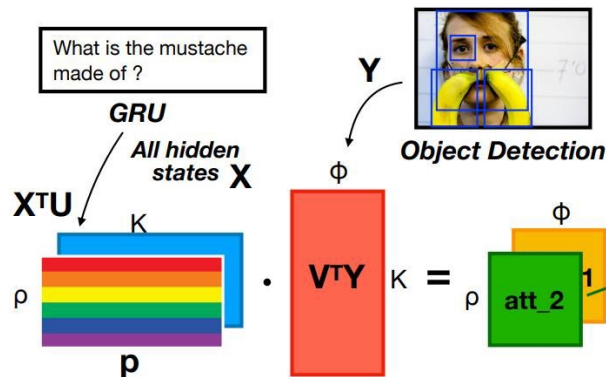
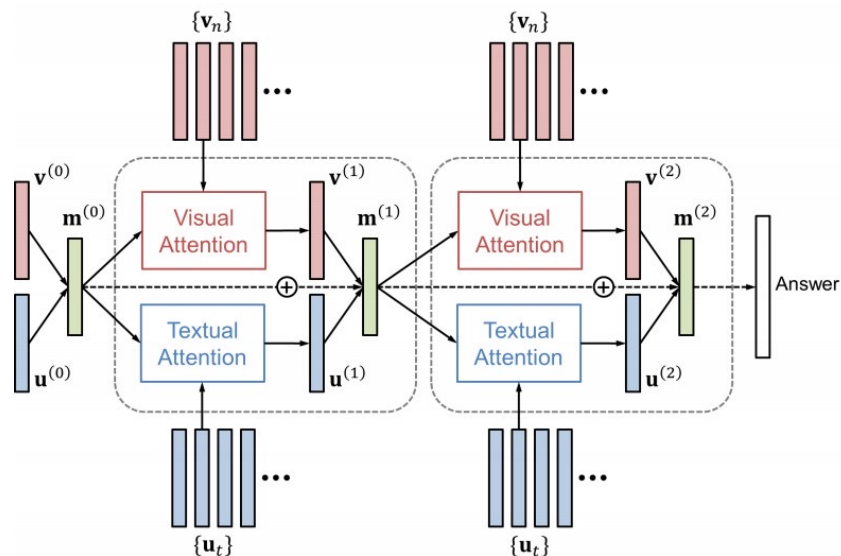


SEOUL NATIONAL UNIVERSITY

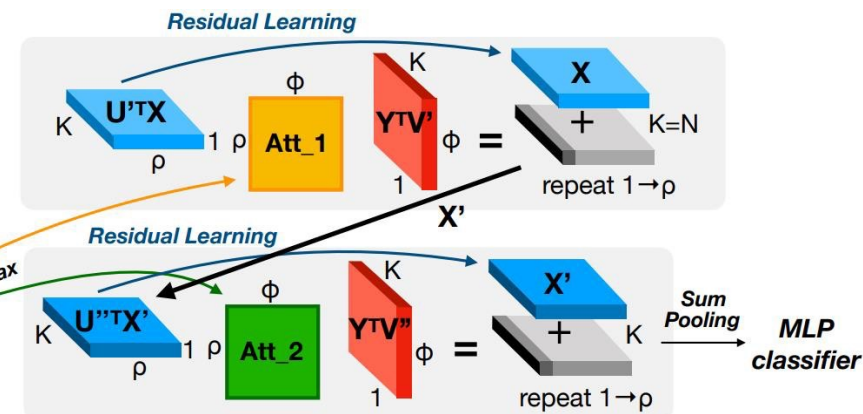
BAN

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

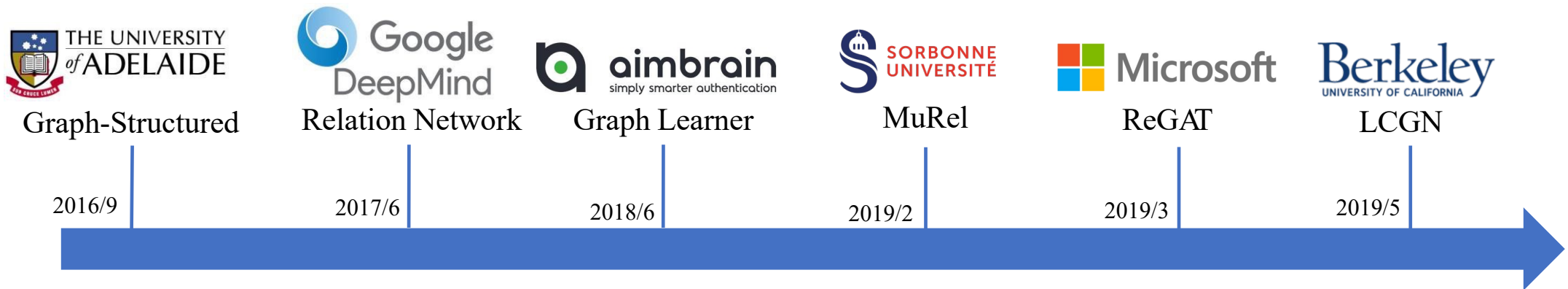
DCN: Dense Co-attention Network

1 Stacked Attention Networks for Image Question Answering, CVPR2016

2 Improved Fusion of Visual and Language Representations by Dense Symmetric Co-Attention for Visual Question Answering, CVPR2018

# 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



2017/6



Graph Learner

2018/6



MuRel

2019/2



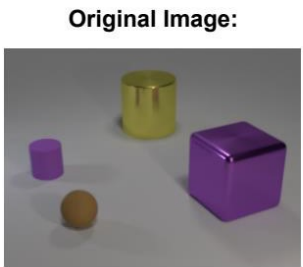
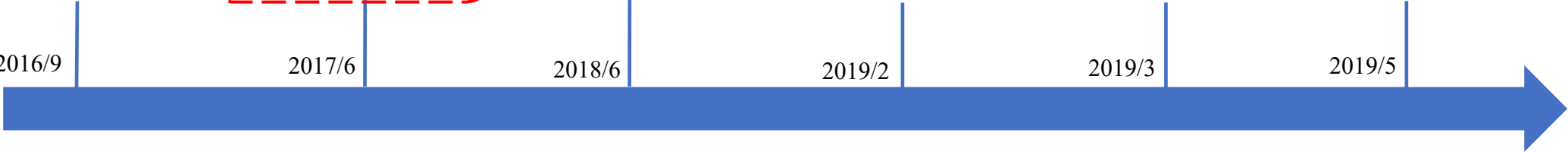
ReGAT

2019/3



LCGN

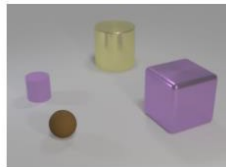
2019/5



Original Image:

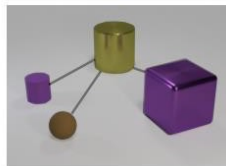
Non-relational question:

What is the size of the brown sphere?

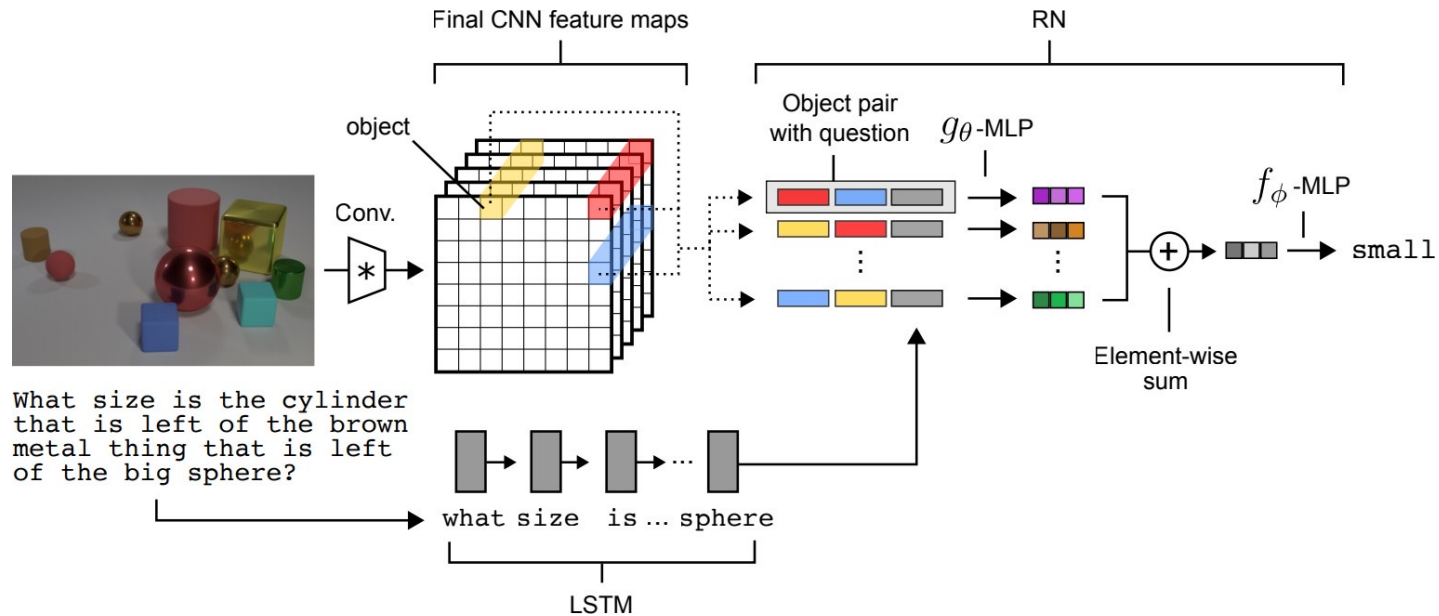


Relational question:

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



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



Relational Network: A fully-connected graph is constructed



Graph-Structured

2016/9



Relation Network

2017/6



Graph Learner

2018/6



MuRel

2019/2

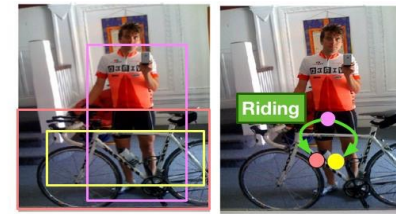
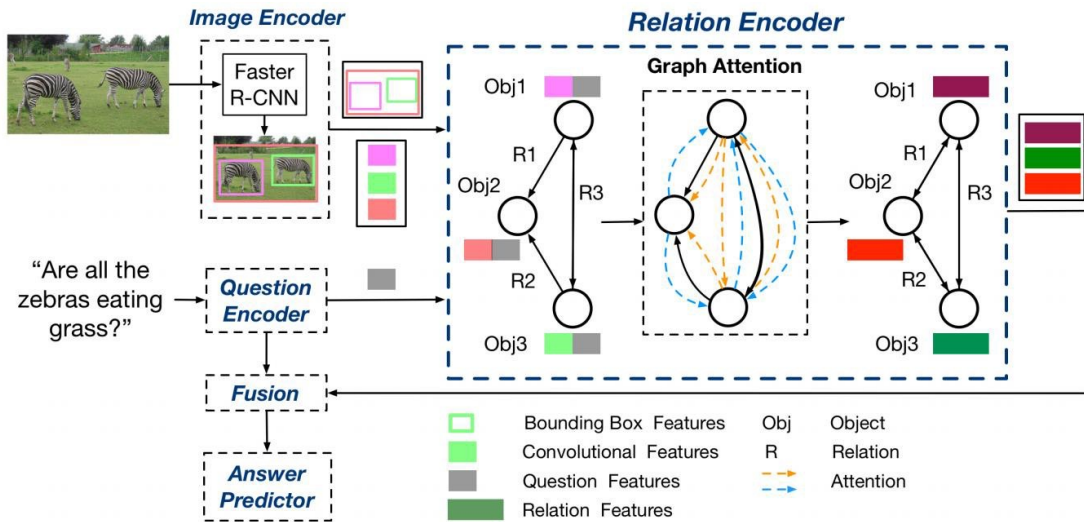


2019/3

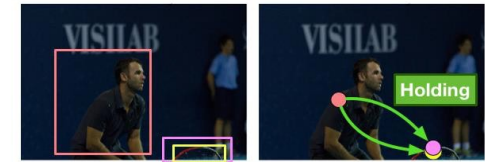


LCGN

2019/5

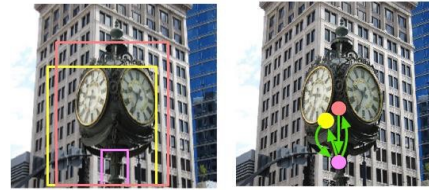


Q: Is this the typical fashion for riding this bike?  
A: Yes

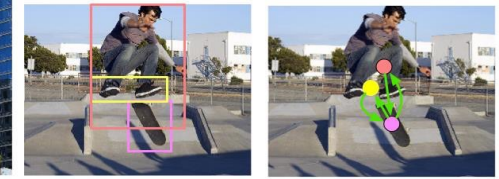


Q: What is he holding?  
A: Tennis Racket

(a) Semantic Relation



Q: What's the clock attached to?  
A: Pole



Q: Are his feet touching the skateboard?  
A: No

(b) Spatial Relation



Q: Where is the vase?  
A: On the table



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

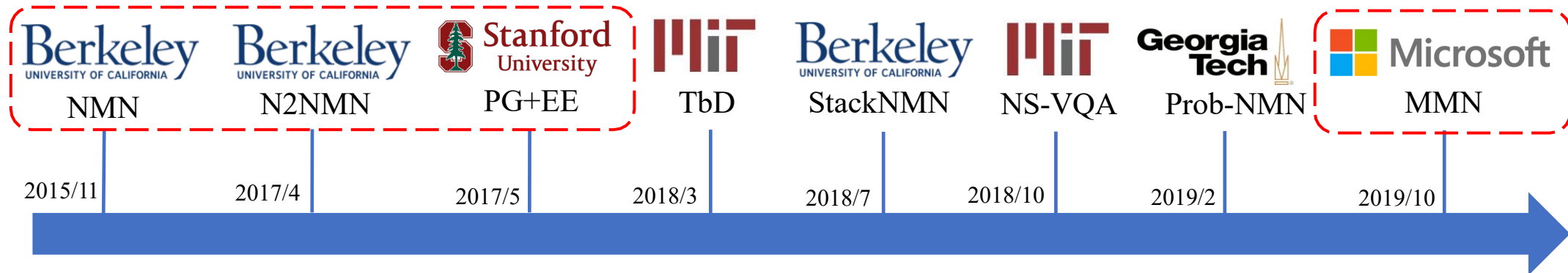
(c) Implicit Relation

- **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

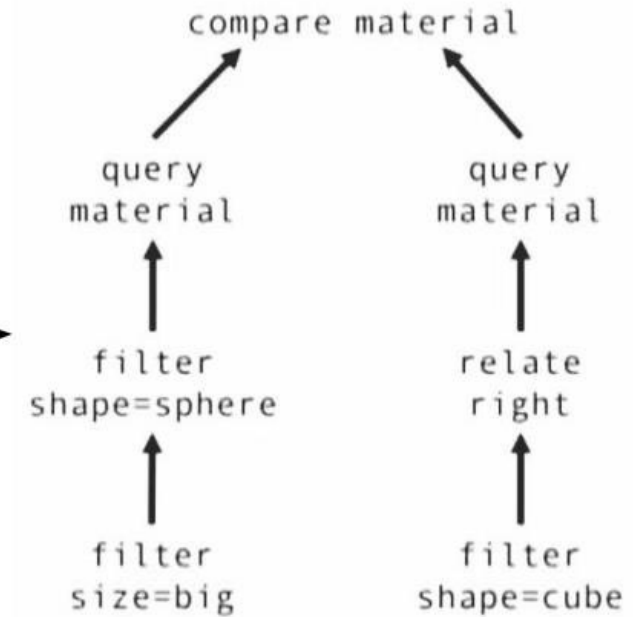
Q: How many spheres are the left of the big sphere and the same color as the small rubber cylinder?

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?

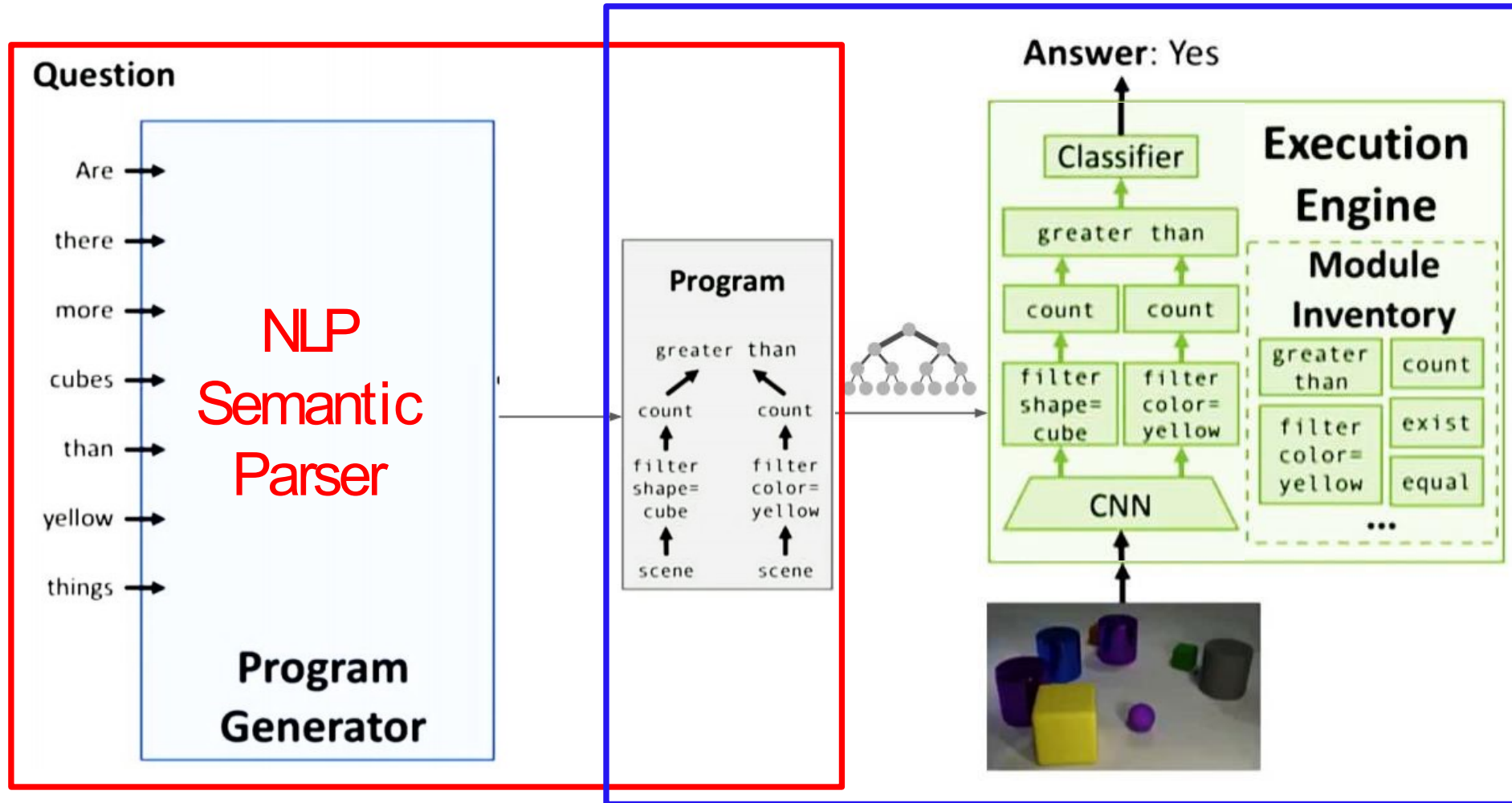
## Common operations

Attributes identification  
Counting objects  
Comparisons  
Spatial relationships  
Logical operations



**Network architecture  
corresponding to the  
third question**

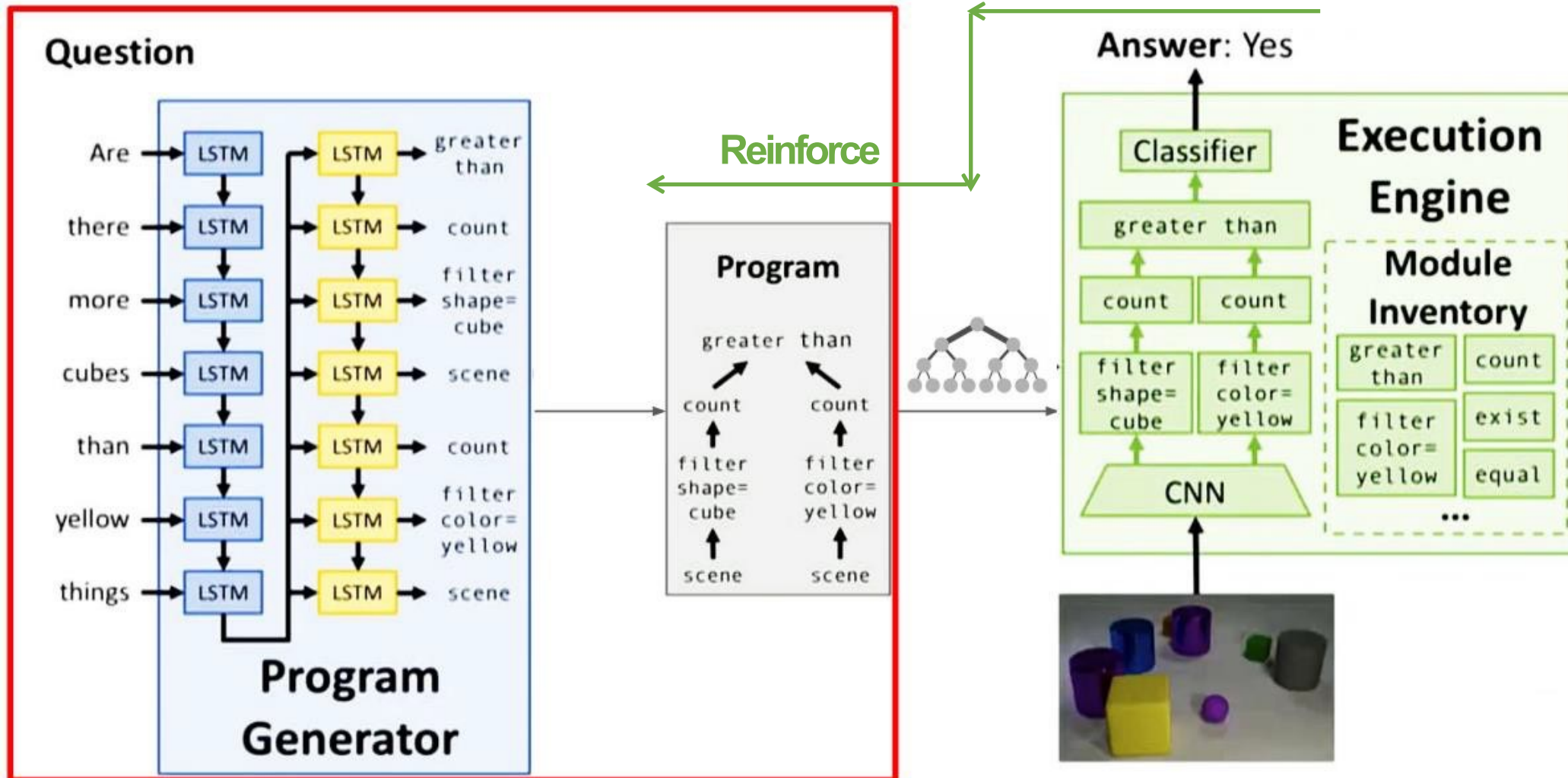
# Overview of the NMN approach

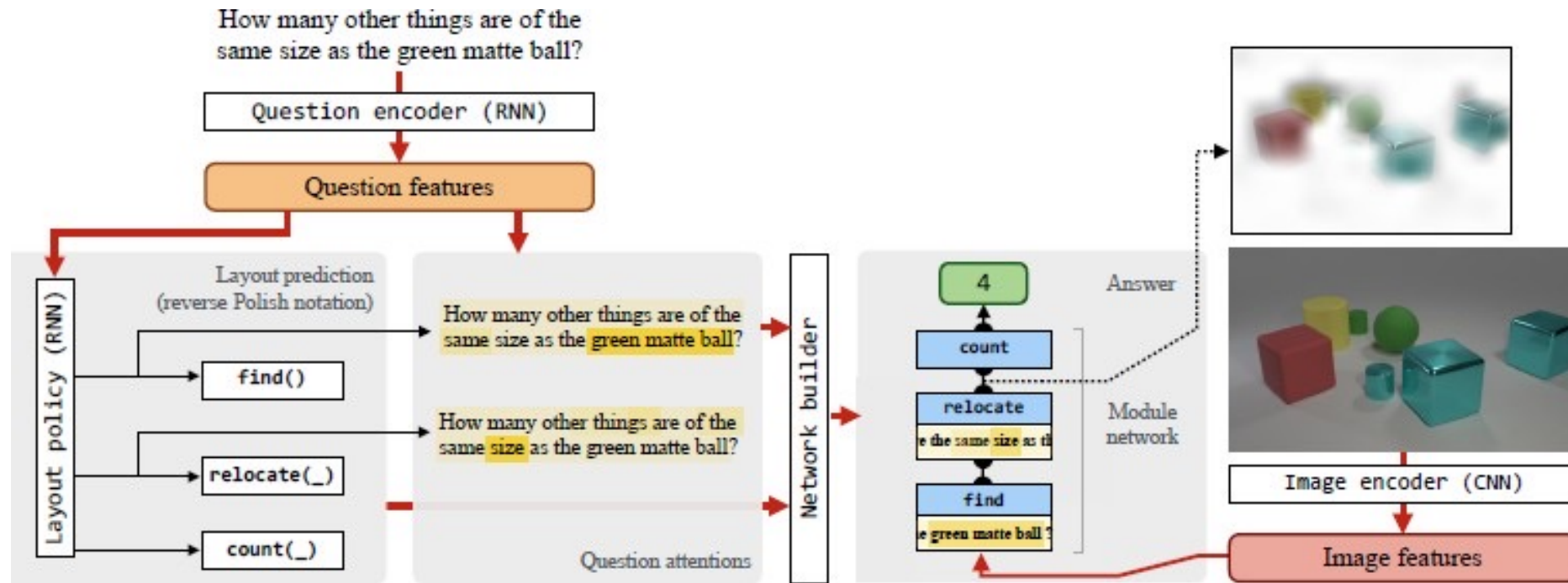
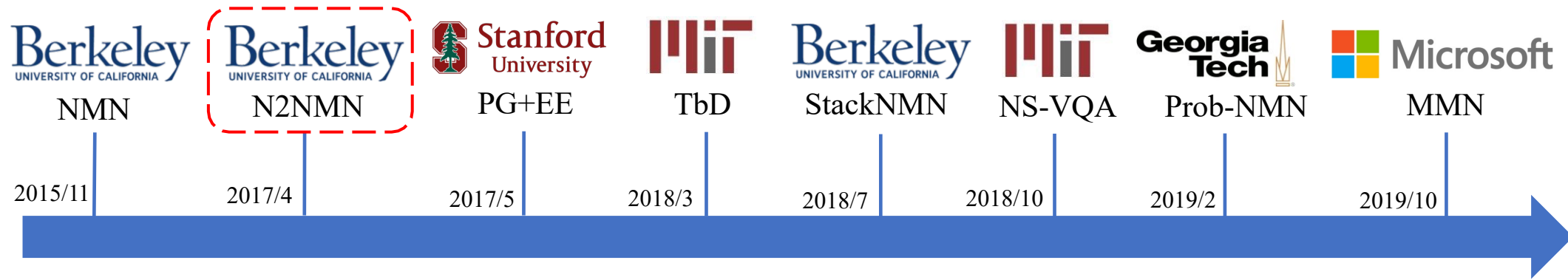


Uses some pre-trained parser

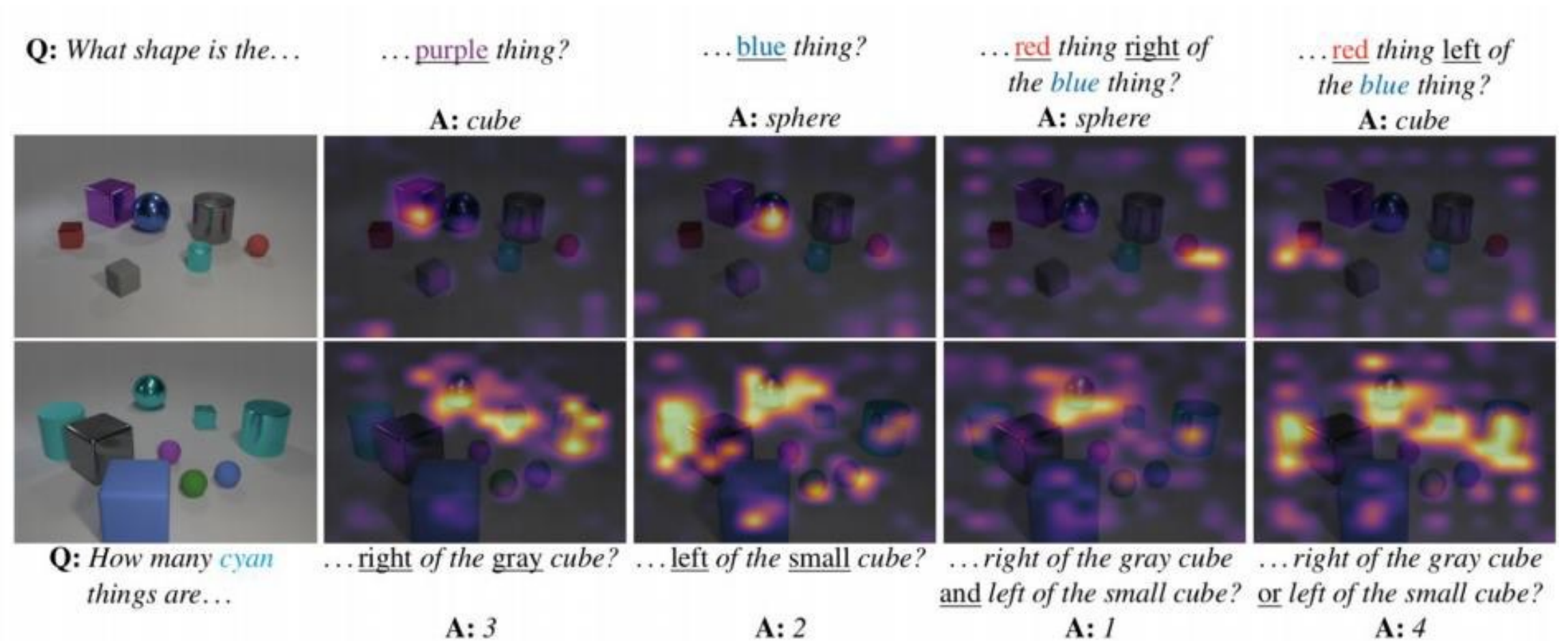
Trained separately

# Inferring and Executing Programs





# What do the modules learn?

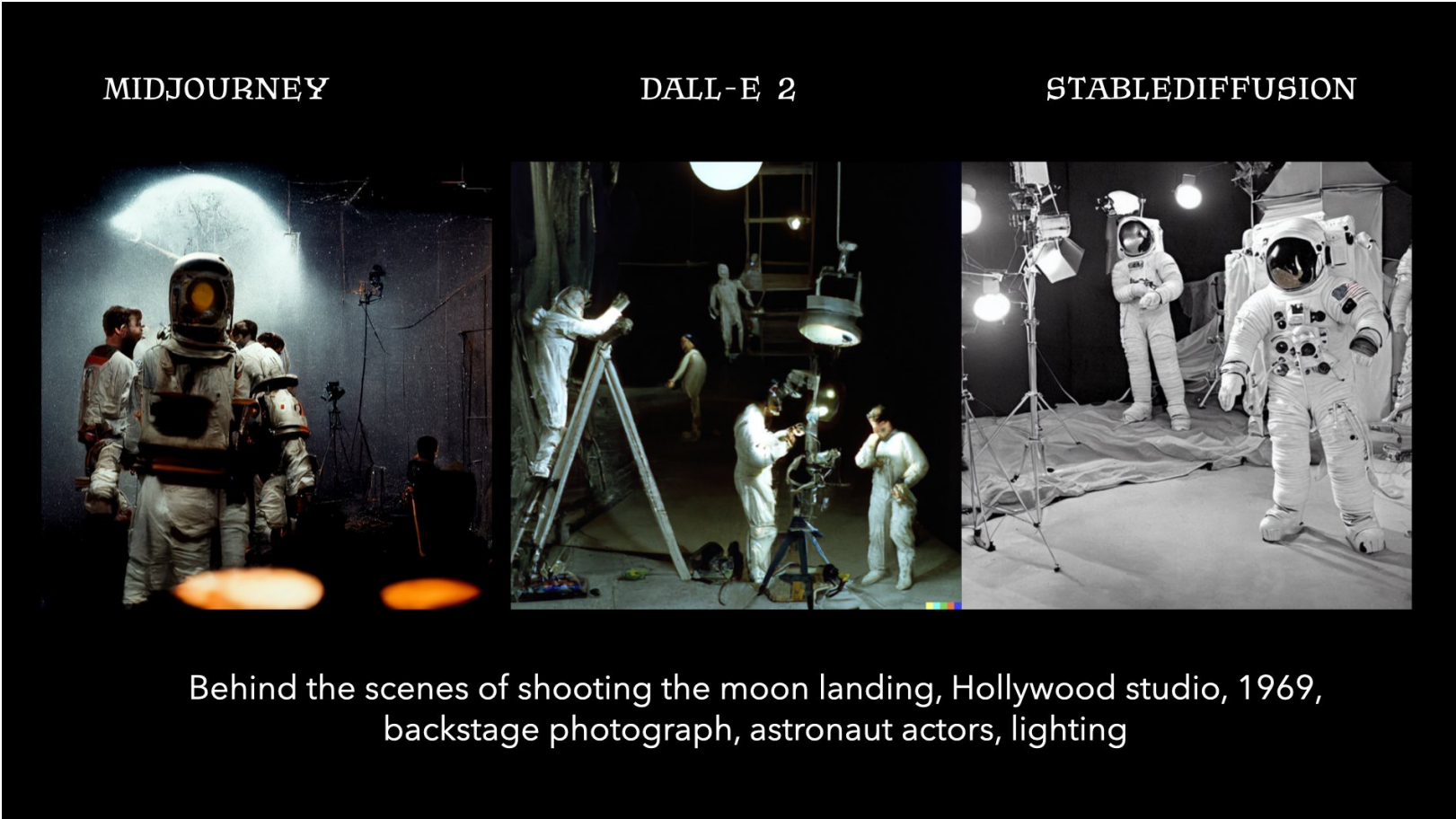


# Now, a much more ambitious task (and a \$\$\$\$ market): Text2Image!

MIDJOURNEY

DALL-E 2

STABLEDIFFUSION



Behind the scenes of shooting the moon landing, Hollywood studio, 1969, backstage photograph, astronaut actors, lighting

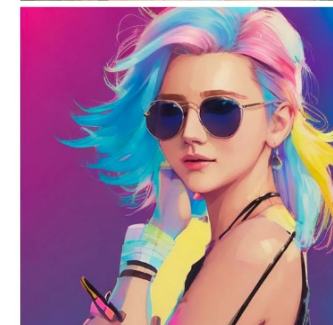
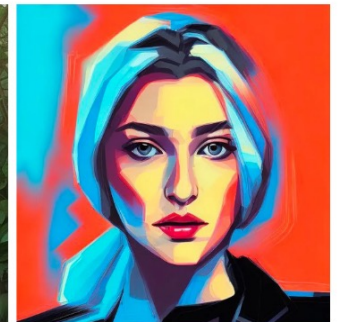
The New York Times

IT HAPPENED ONLINE

## How Is Everyone Making Those A.I. Selfies?

Images generated with Lensa AI are all over social media, but at what cost?

Give this article



<https://cvpr2022-tutorial-diffusion-models.github.io/> (lots of slides borrowed hereinafter)

Lensa AI, a popular iPhone app, uses your selfies and artificial intelligence to create portraits in a variety of styles. Lensa AI

## DALL·E 2

“a teddy bear on a skateboard in times square”



[“Hierarchical Text-Conditional Image Generation with CLIP Latents”](#)  
Ramesh et al., 2022

## Imagen

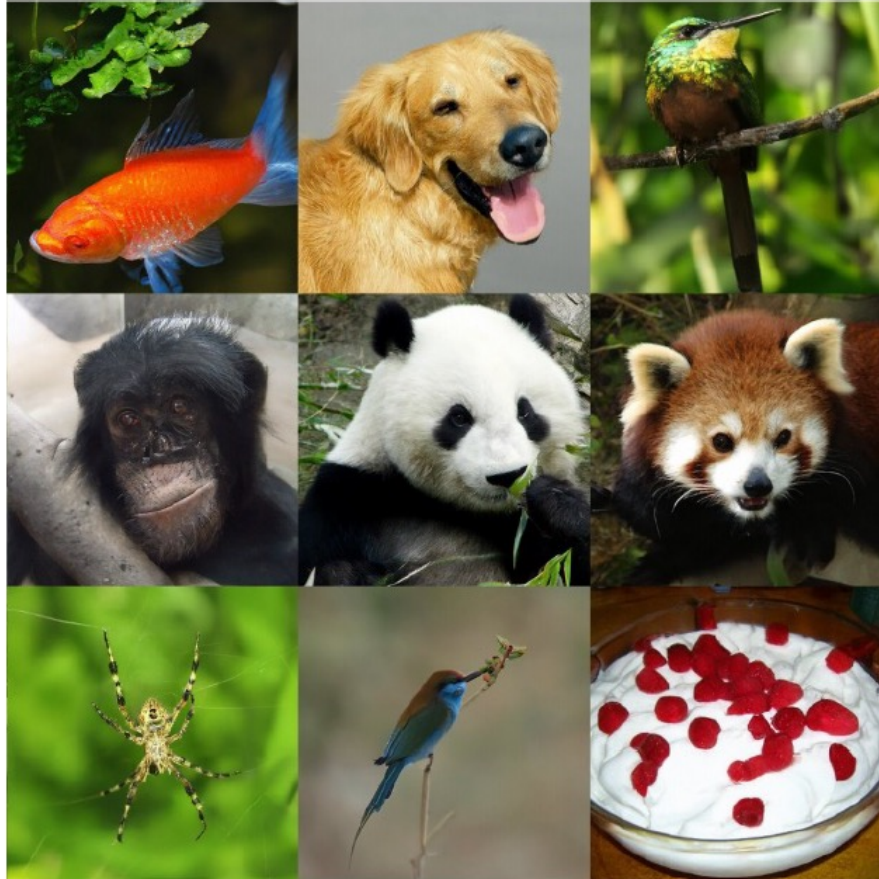
A group of teddy bears in suit in a corporate office celebrating the birthday of their friend. There is a pizza cake on the desk.



[“Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding”](#), Saharia et al., 2022



# The Workhorse: *Diffusion Models*



[“Diffusion Models Beat GANs on Image Synthesis”](#)  
Dhariwal & Nichol, OpenAI, 2021

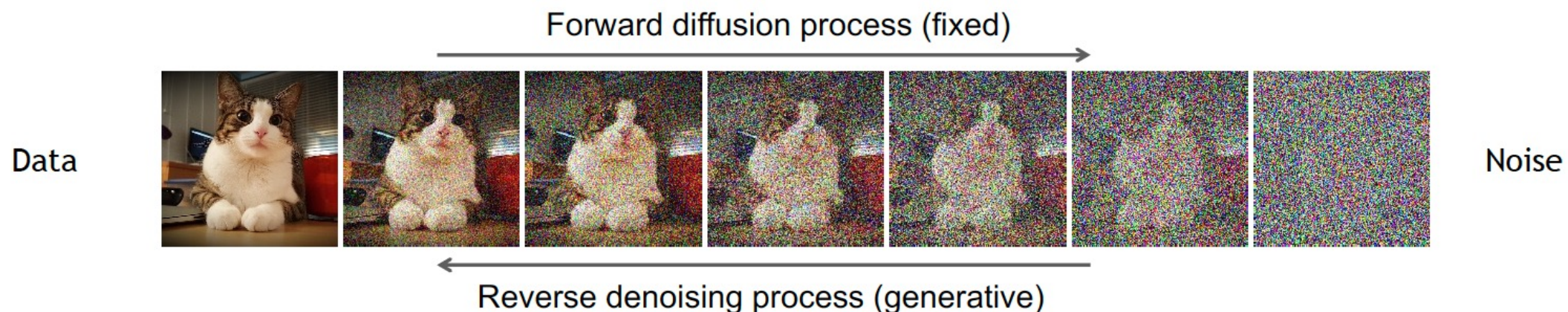


[“Cascaded Diffusion Models for High Fidelity Image Generation”](#)  
Ho et al., Google, 2021

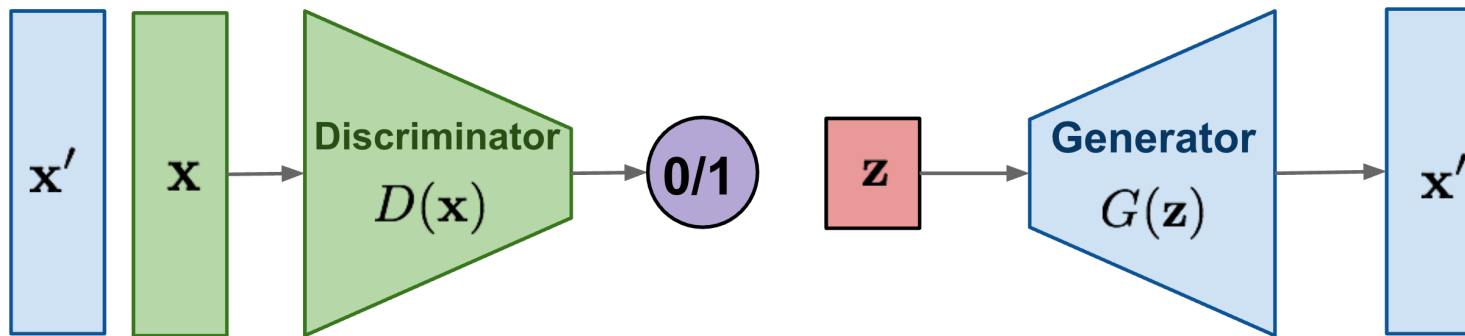
# Learning to generate by denoising

Denoising diffusion models consist of two processes:

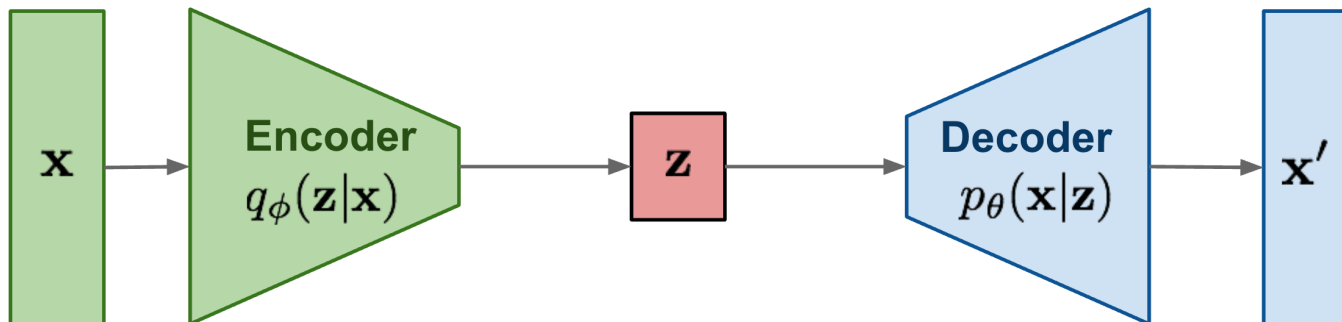
- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising



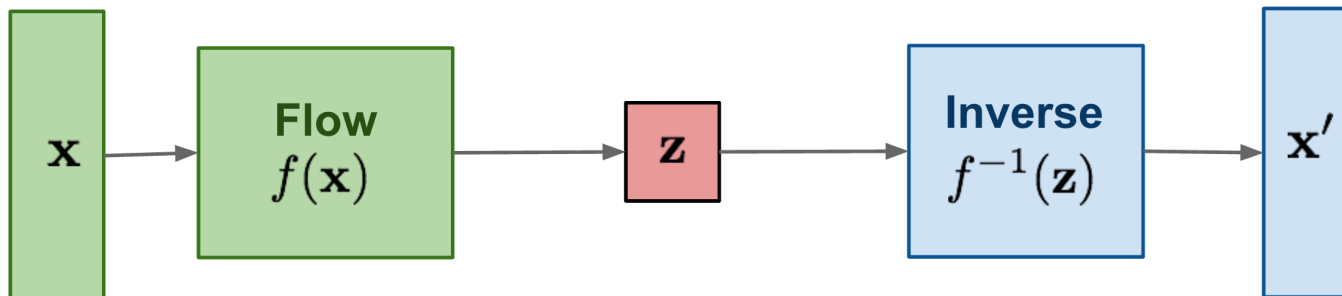
**GAN:** Adversarial training



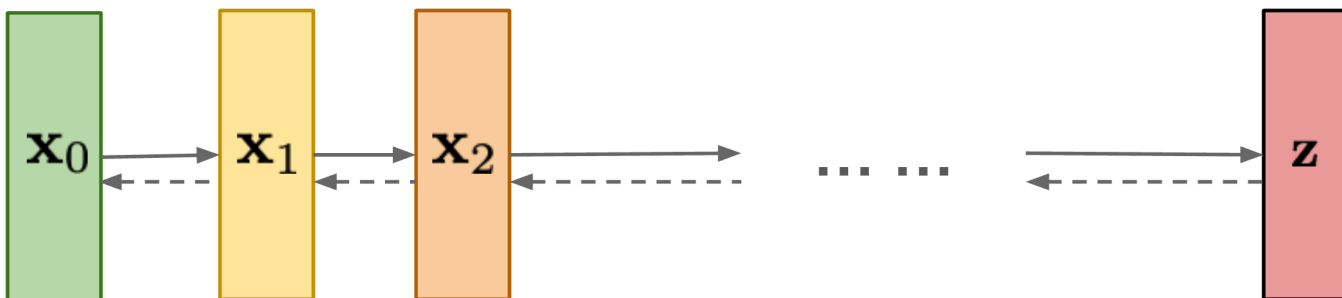
**VAE:** maximize variational lower bound



**Flow-based models:** Invertible transform of distributions

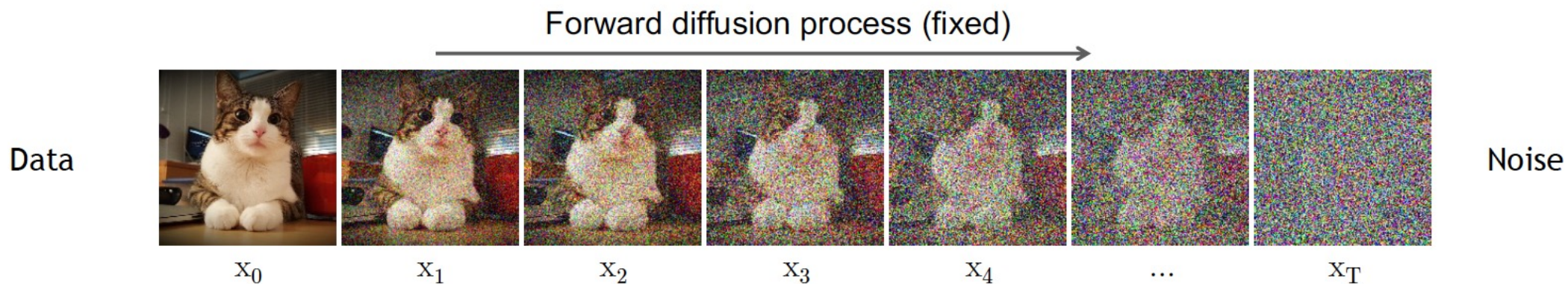


**Diffusion models:** Gradually add Gaussian noise and then reverse

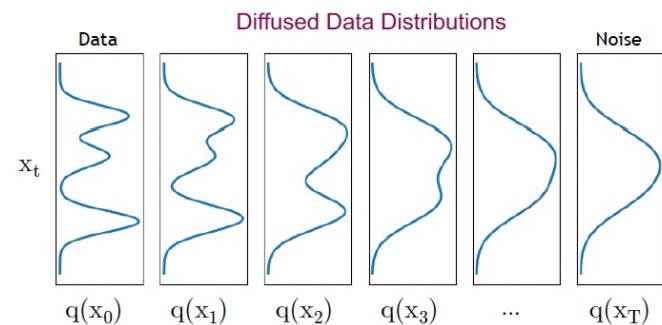


# Forward Diffusion Process

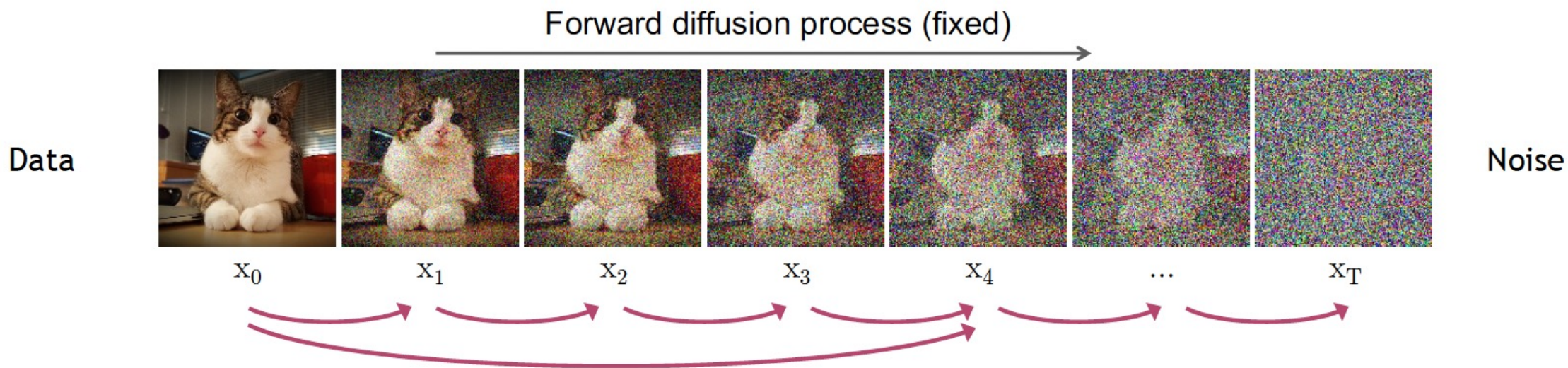
The formal definition of the forward process in T steps:



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad \rightarrow \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$



# Sampling at arbitrary time step with “reparameterization trick”



Define  $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$     $\longrightarrow$     $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$    (Diffusion Kernel)

For sampling:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$    where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

*The diffusion kernel is Gaussian convolution.*

$\beta_t$  values schedule (i.e., the noise schedule) is designed such that  $\bar{\alpha}_T \rightarrow 0$  and  $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

# Generative Learning by Denoising

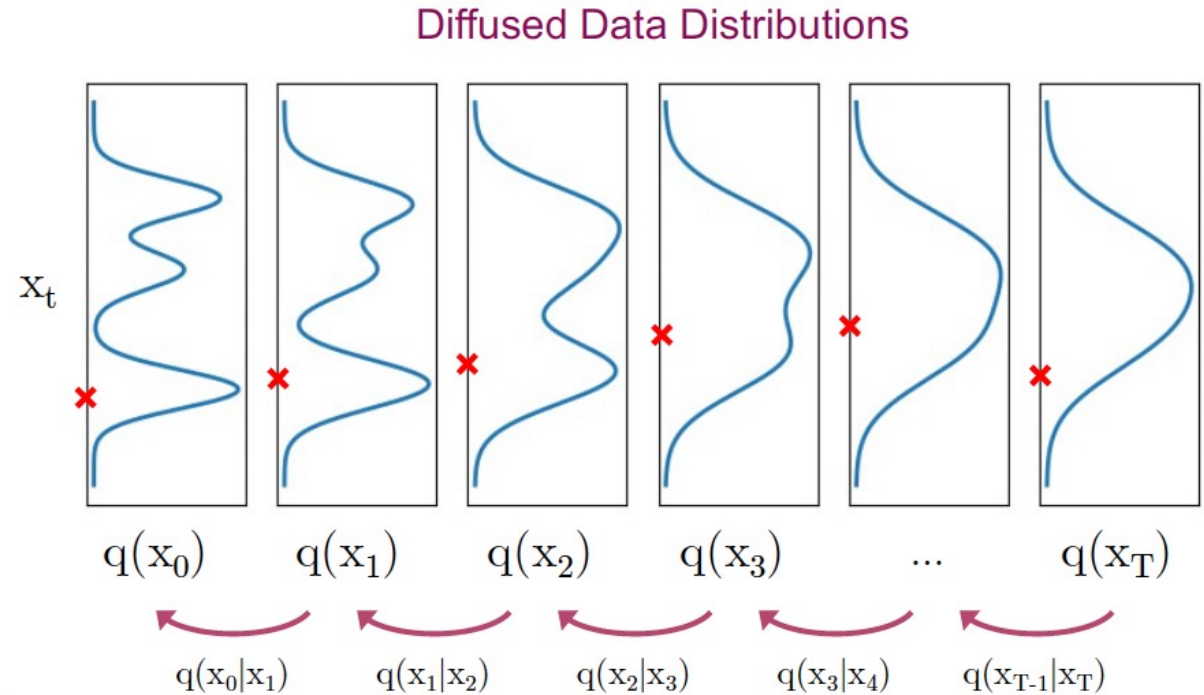
Recall, that the diffusion parameters are designed such that  $q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

**Generation:**

Sample  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Iteratively sample  $\mathbf{x}_{t-1} \sim \underbrace{q(\mathbf{x}_{t-1}|\mathbf{x}_t)}$

True Denoising Dist.

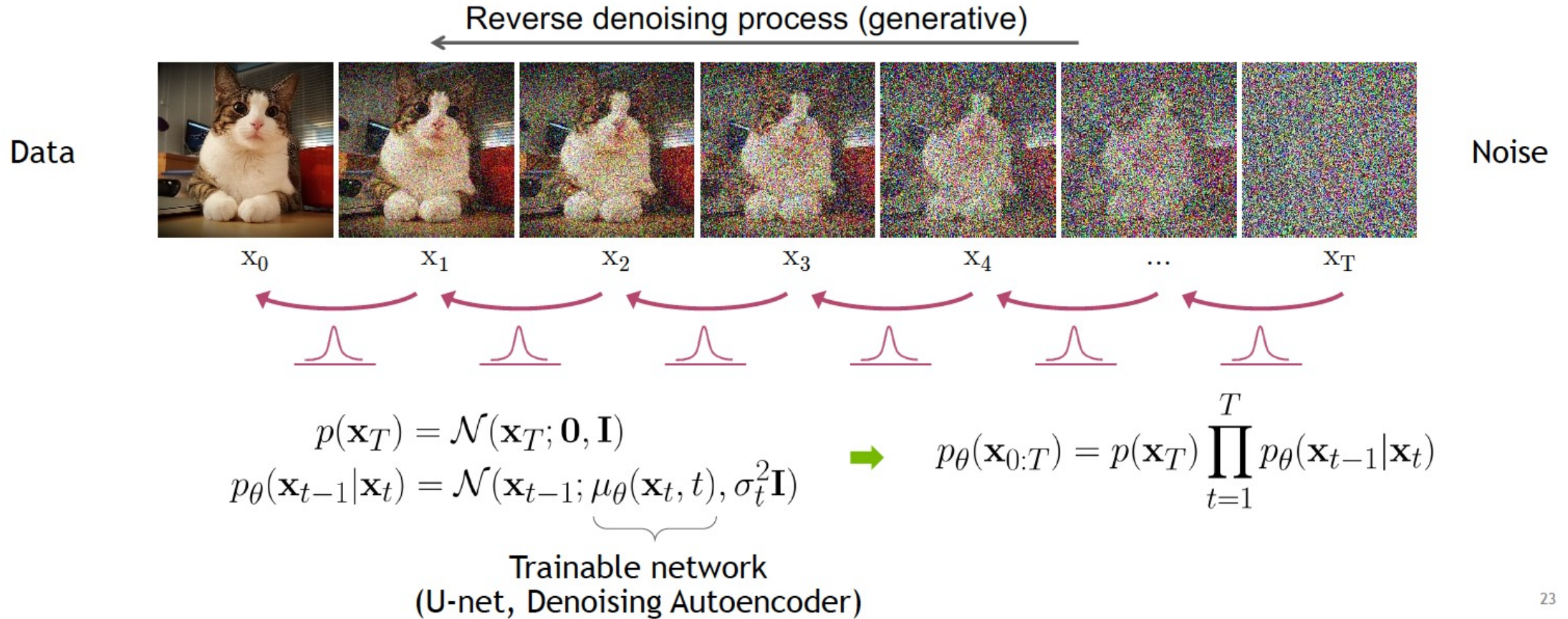


In general,  $q(\mathbf{x}_{t-1}|\mathbf{x}_t) \propto q(\mathbf{x}_{t-1})q(\mathbf{x}_t|\mathbf{x}_{t-1})$  is intractable.

Can we approximate  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ ? Yes, we can use a **Normal distribution** if  $\beta_t$  is small in each forward diffusion step.

# Reverse Denoising Process

Formal definition of forward and reverse processes in T steps:



# Denoising diffusion probabilistic models (DDPM)

---

## Algorithm 1 Training

---

- 1: **repeat**
  - 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
  - 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
  - 4:  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 5: Take gradient descent step on  
$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2$$
  - 6: **until** converged
- 

---

## Algorithm 2 Sampling

---

- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2: **for**  $t = T, \dots, 1$  **do**
  - 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
  - 5: **end for**
  - 6: **return**  $\mathbf{x}_0$
- 

- Denoising Diffusion models can be considered as a **special form of hierarchical VAEs**.

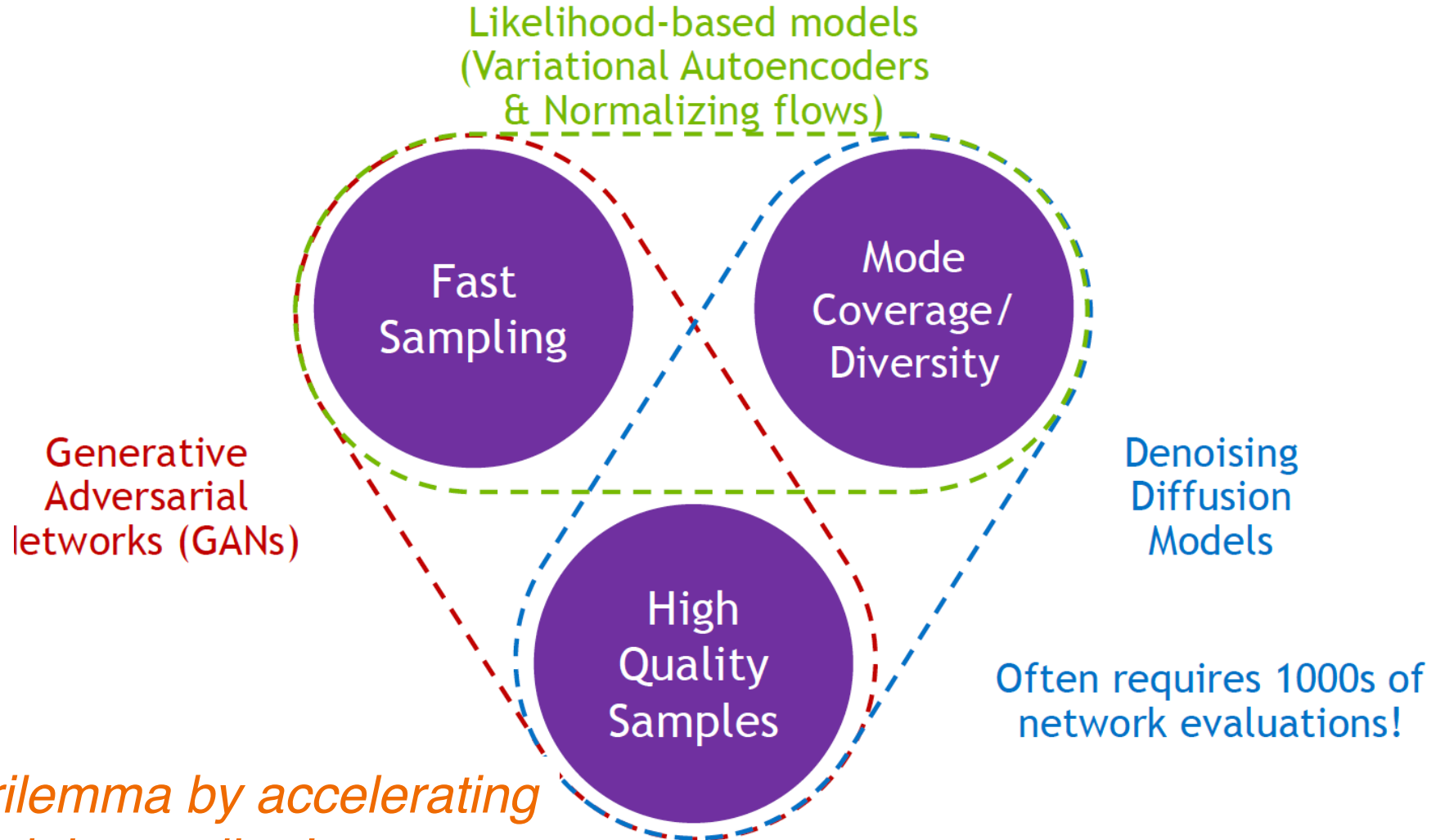
- The model is trained with some reweighting of the variational bound

However, in diffusion models:

- The encoder is fixed
- The latent variables always have the same dimension as the data (no “bottleneck”)
- Denoising model is shared across different timesteps

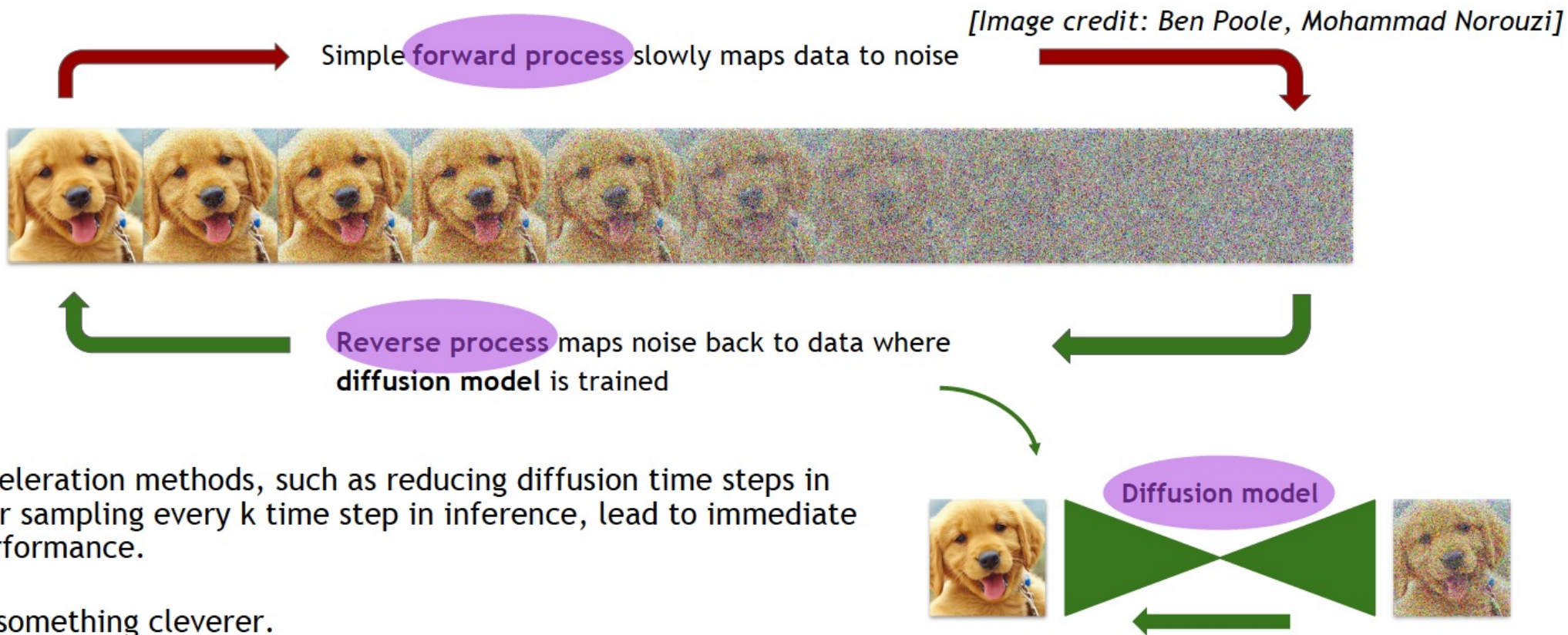


# The generative learning trilemma

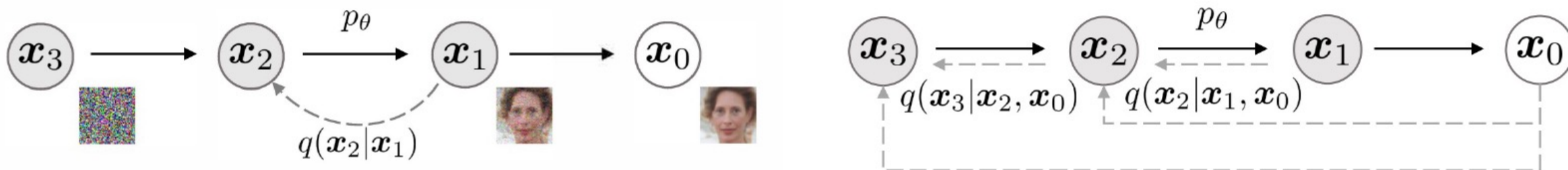


*Tackle the trilemma by accelerating diffusion model sampling!*

# How to accelerate diffusion models?



# From DDPM to DDIM: *Denoising diffusion implicit models*



## Main Idea

Design a family of non-Markovian diffusion processes and corresponding reverse processes.

The process is designed such that the model can be optimized by the same surrogate objective as the original diffusion model.

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[ \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2 \right]$$

Therefore, can take a pretrained diffusion model but with more choices of sampling procedure.

# From DDPM to DDIM: *Denoising diffusion implicit models*

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t-1}}\hat{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \tilde{\sigma}_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_0}{\sqrt{1 - \bar{\alpha}_t}}, \tilde{\sigma}_t^2 \mathbf{I}\right)$$

- ... often using its **deterministic form**:  $\tilde{\sigma}_t^2 = 0, \forall t$
- With DDIM, it is possible to train the diffusion model up to any arbitrary number of forward steps but only **sample from a subset of steps in the generative process**

During generation, we only sample a subset of  $S$  diffusion steps  $\{\tau_1, \dots, \tau_S\}$  and the inference process becomes:

$$q_{\sigma, \tau}(\mathbf{x}_{\tau_{i-1}}|\mathbf{x}_{\tau_t}, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{\tau_{i-1}}; \sqrt{\bar{\alpha}_{\tau_{i-1}}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{\tau_{i-1}} - \sigma_t^2} \frac{\mathbf{x}_{\tau_t} - \sqrt{\bar{\alpha}_{\tau_t}}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_{\tau_t}}}, \sigma_t^2 \mathbf{I})$$

# Conditional Generation

Reverse process:  $p_\theta(\mathbf{x}_{0:T}|\mathbf{c}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c}), \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c}) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t, \mathbf{c}), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t, \mathbf{c}))$

Variational upper bound:  $L_\theta(\mathbf{x}_0|\mathbf{c}) = \mathbb{E}_q \left[ L_T(\mathbf{x}_0) + \sum_{t>1} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})) - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1, \mathbf{c}) \right].$

## Incorporate conditions into U-Net

- Scalar conditioning: encode scalar as a vector embedding, simple spatial addition or adaptive group normalization layers.
- Image conditioning: channel-wise concatenation of the conditional image.
- Text conditioning: single vector embedding - spatial addition or adaptive group norm / a seq of vector embeddings - cross-attention.

# Classifier guidance: Guiding Sampling using Using the gradient of a trained classifier

---

**Algorithm 1** Classifier guided diffusion sampling, given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , classifier  $p_\phi(y|x_t)$ , and gradient scale  $s$ .

---

Input: class label  $y$ , gradient scale  $s$

$x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$

**for all**  $t$  from  $T$  to 1 **do**

$\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$

$x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_\phi(y|x_t), \Sigma)$

**end for**

**return**  $x_0$

---

## Main Idea

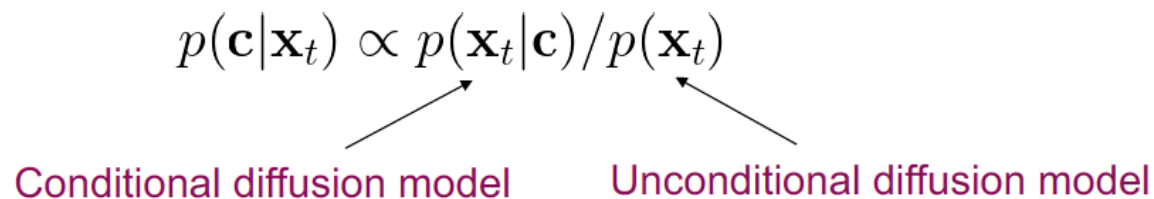
For class-conditional modeling of  $p(\mathbf{x}_t|\mathbf{c})$ , train an extra classifier  $p(\mathbf{c}|\mathbf{x}_t)$

Mix its gradient with the diffusion/score model during sampling

Sample with a modified score:  $\nabla_{\mathbf{x}_t} [\log p(\mathbf{x}_t|\mathbf{c}) + \omega \log p(\mathbf{c}|\mathbf{x}_t)]$

# Classifier-free guidance: Implicit trick via Bayesian rule

- Instead of training an additional classifier, get an “implicit classifier” by jointly training a conditional and unconditional diffusion model:

$$p(\mathbf{c}|\mathbf{x}_t) \propto p(\mathbf{x}_t|\mathbf{c})/p(\mathbf{x}_t)$$


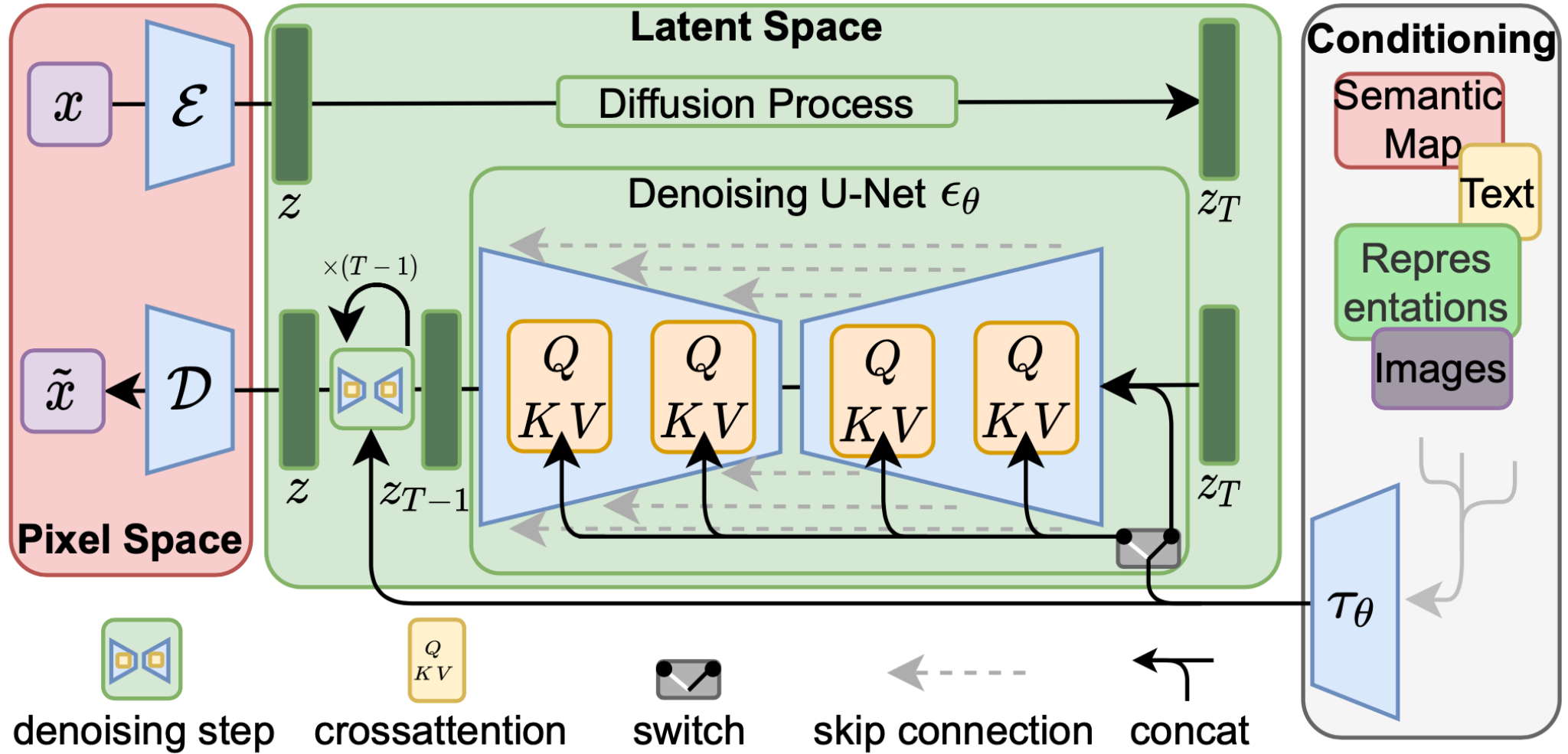
Conditional diffusion model      Unconditional diffusion model

- In practice,  $p(\mathbf{x}_t|\mathbf{c})$  and  $p(\mathbf{x}_t)$  by randomly dropping the condition of the diffusion model at certain chance.
- The modified score with this implicit classifier included is:

$$\begin{aligned}\nabla_{\mathbf{x}_t}[\log p(\mathbf{x}_t|\mathbf{c}) + \omega \log p(\mathbf{c}|\mathbf{x}_t)] &= \nabla_{\mathbf{x}_t}[\log p(\mathbf{x}_t|\mathbf{c}) + \omega(\log p(\mathbf{x}_t|\mathbf{c}) - \log p(\mathbf{x}_t))] \\ &= \nabla_{\mathbf{x}_t}[(1 + \omega) \log p(\mathbf{x}_t|\mathbf{c}) - \omega \log p(\mathbf{x}_t)]\end{aligned}$$

# Latent Diffusion Model (CVPR'22): Important Jump toward High-Resolution!

*DDIM sampler  
+ classifier-free  
guidance +  
many other  
tweaks ...*





# Latent Diffusion Model (CVPR'22): Important Jump toward High-Resolution!



```
python scripts/txt2img.py --prompt "a sunset behind a mountain range, vector image" --ddim_eta 1.0 --n_samples 1 --n_iter 1 --H 384 --W 1024 --scale 5.0
```

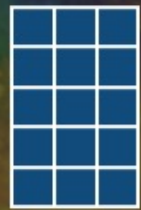


# Stable Diffusion

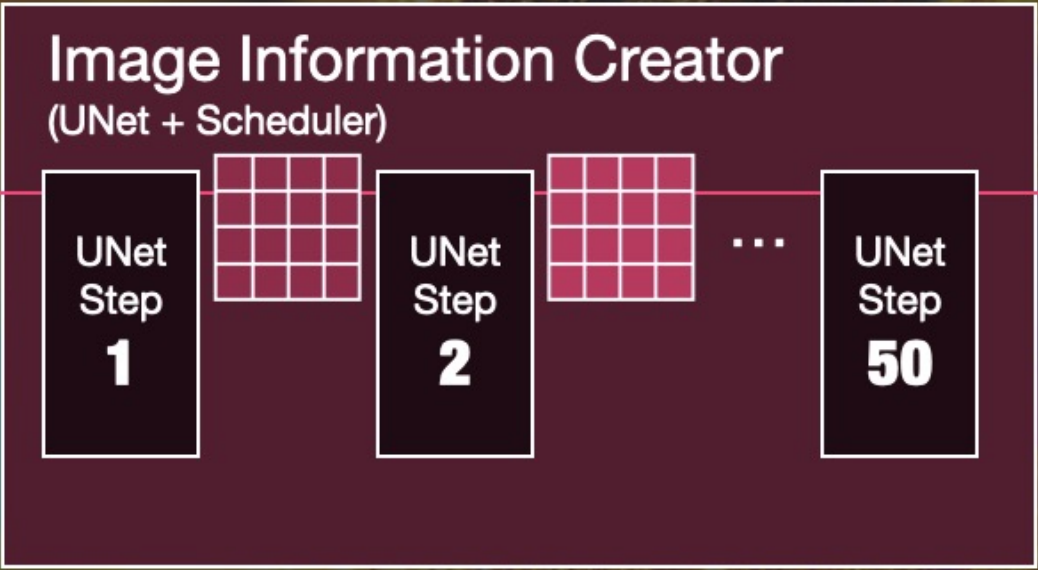
paradise  
cosmic  
beach

77 tokens

**Text Encoder**  
(CLIPText)



Token embeddings



Random image information tensor

Processed image information tensor

Diffusion

**Image Decoder**  
(Autoencoder decoder)

Generated image



# Personalizing Your Diffusion: DreamBooth

## Input

Images (~3-5) +  
subject's class name

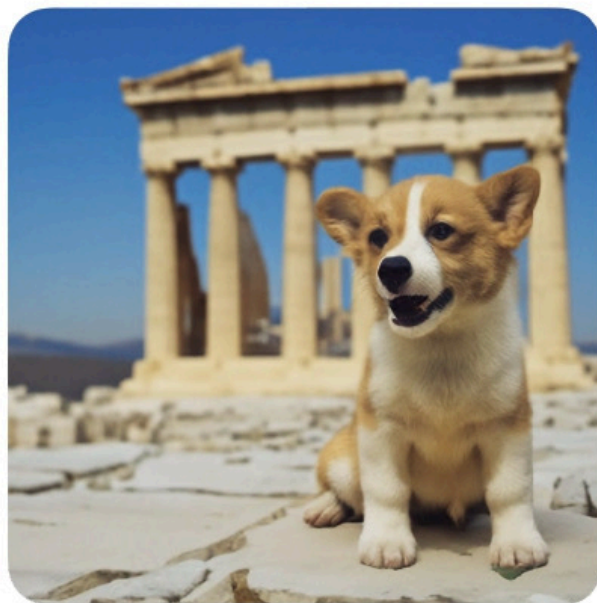


Input images

## Fine-Tuning

## Output

Unique  
identifier



*in the Acropolis*



*swimming*



*sleeping*

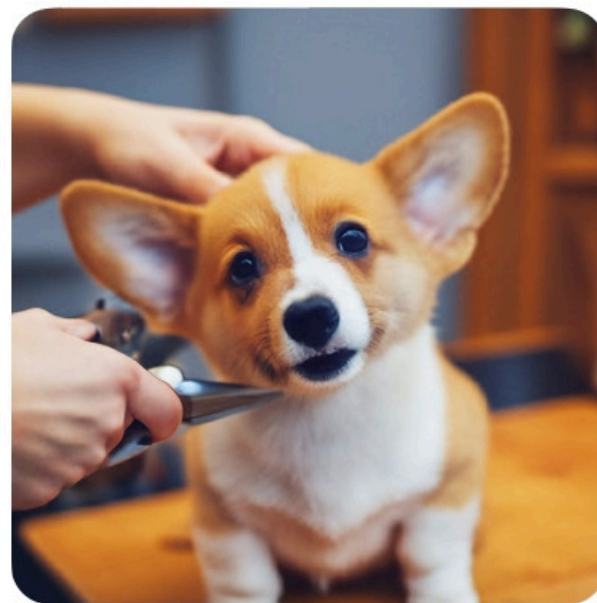


*in a doghouse*



*in a bucket*

## Inference



*getting a haircut*

# Personalizing Your Diffusion: Text Inversion



Input samples  $\xrightarrow{\text{invert}}$  " $S_*$ "



"An oil painting of  $S_*$ "



"App icon of  $S_*$ "



"Elmo sitting in the same pose as  $S_*$ "



"Crochet  $S_*$ "



Input samples  $\xrightarrow{\text{invert}}$  " $S_*$ "



"Painting of two  $S_*$  fishing on a boat"



"A  $S_*$  backpack"



"Banksy art of  $S_*$ "

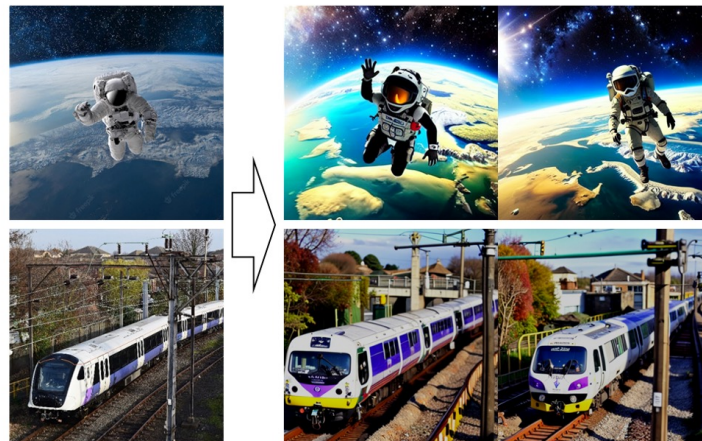


"A  $S_*$  themed lunchbox"

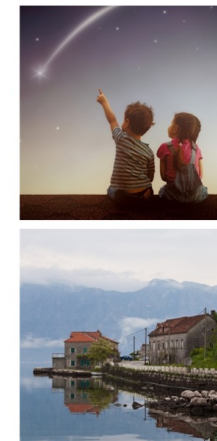
# Versatile Diffusion: All in One!



(a) Text-to-Image

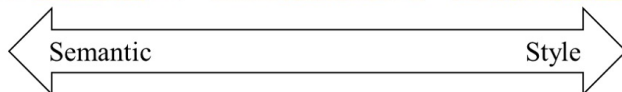


(b) Image-Variation



- There are stars that a child is watching about.
- Two young girls and a boy standing near a star.
- Two young girls are watching a star.
- Kids standing for their stars.
- Houses on the lake with boats and trees beside there with the mountains on the background.
- House, mountain, boat, somewhere near lake
- House on the cliff near the lake.
- Houses on the lake with the trees.

(c) Image-to-Text



(d) Disentanglement



(e) Dual-Guided Generation

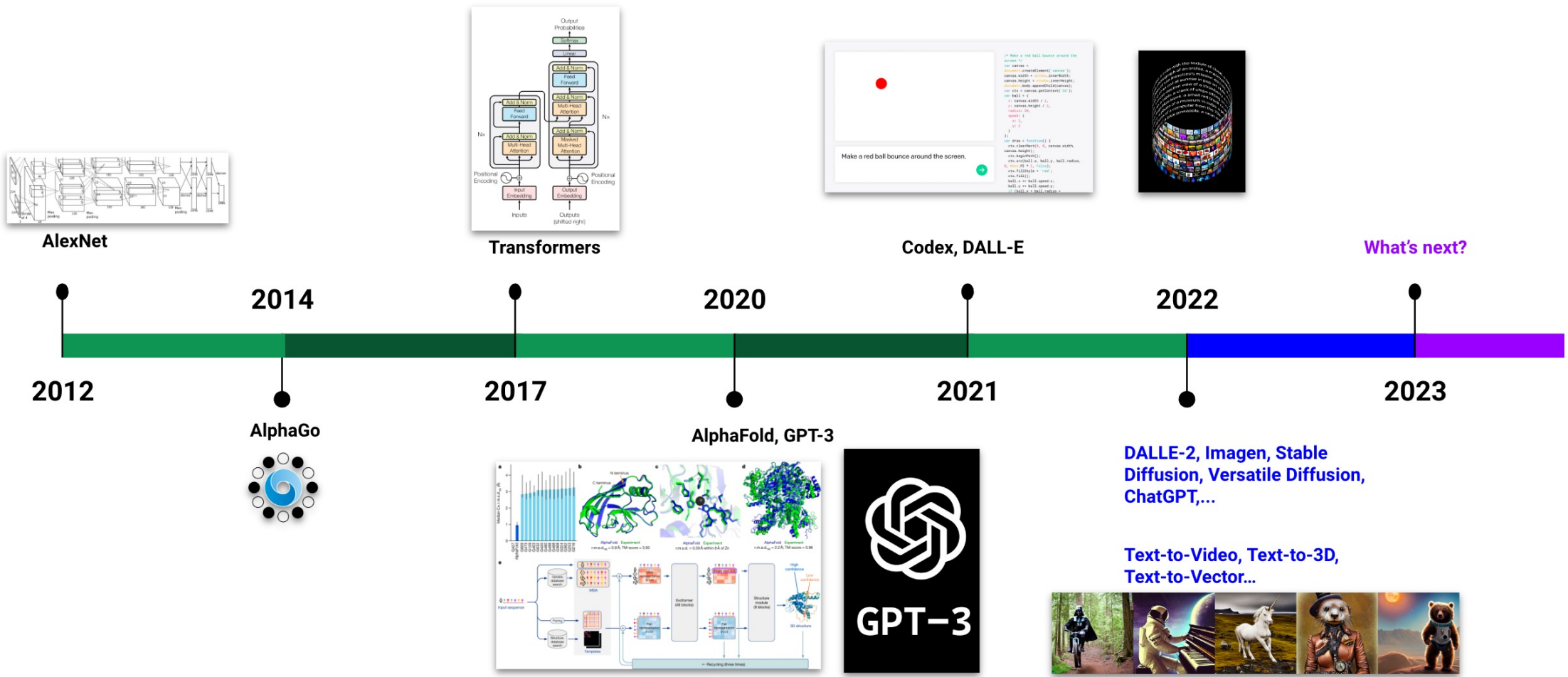


A house on a lake.

A house on a lake.  
tall castle

(f) Editable I2T2I

# Generative AI is revolutionizing the AI landscape right now..





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