# Deep Learning for Computer Vision

## Fall 2022

https://cool.ntu.edu.tw/courses/189345 (NTU COOL)

http://vllab.ee.ntu.edu.tw/dlcv.html (Public website)

Yu-Chiang Frank Wang 王鈺強, Professor

Dept. Electrical Engineering, National Taiwan University

2022/11/1

# What to Cover Today…

- **Recurrent Neural Network** & **Transformer**
  - Attention in RNN
  - *Attention is All You Need*: Transformer
  - Transformer for Visual Analysis
    - Visual Classification
    - Semantic Segmentation & More

- **Vision & Language**
  - Image Captioning
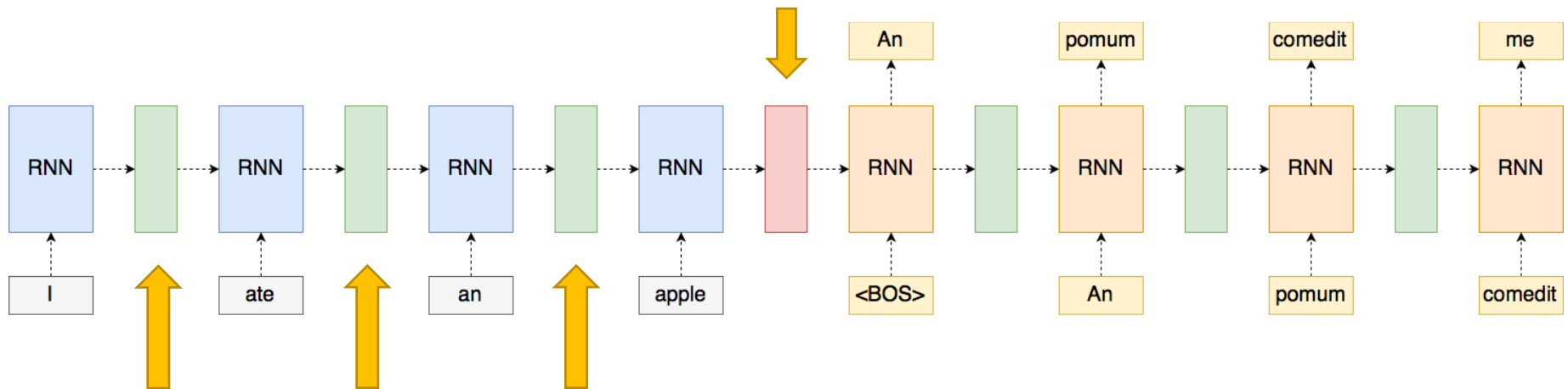  - Text-to-Image Synthesis



"a corgi wearing a bow tie and a birthday hat"



*Teddy bears shopping for groceries in the style of ukiyo-e*
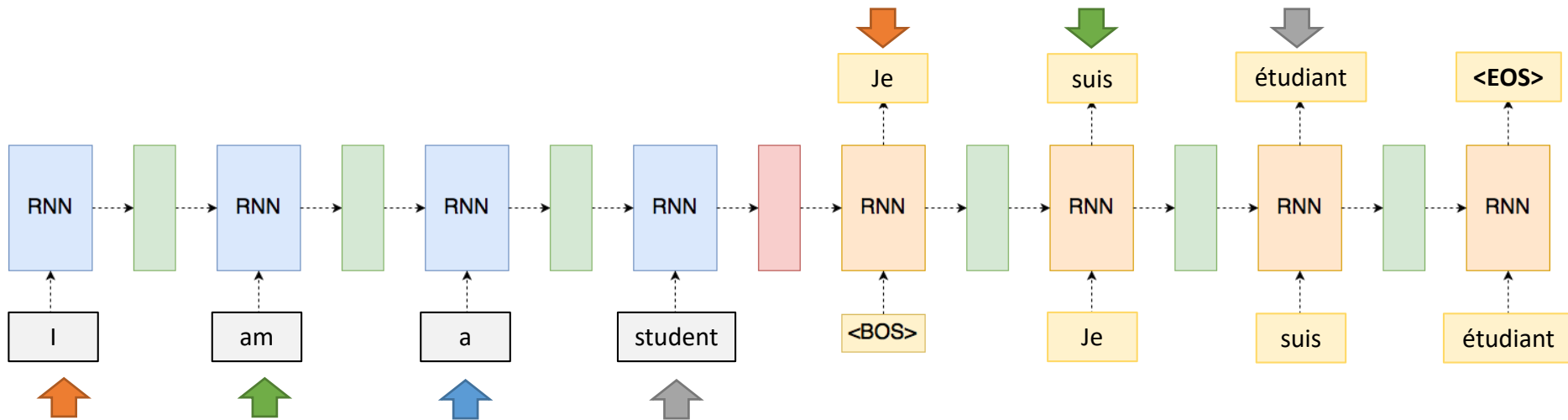
# What's the Potential Problem of RNN?

- Each hidden state vector extracts/carries information across time steps (some might be diluted downstream).

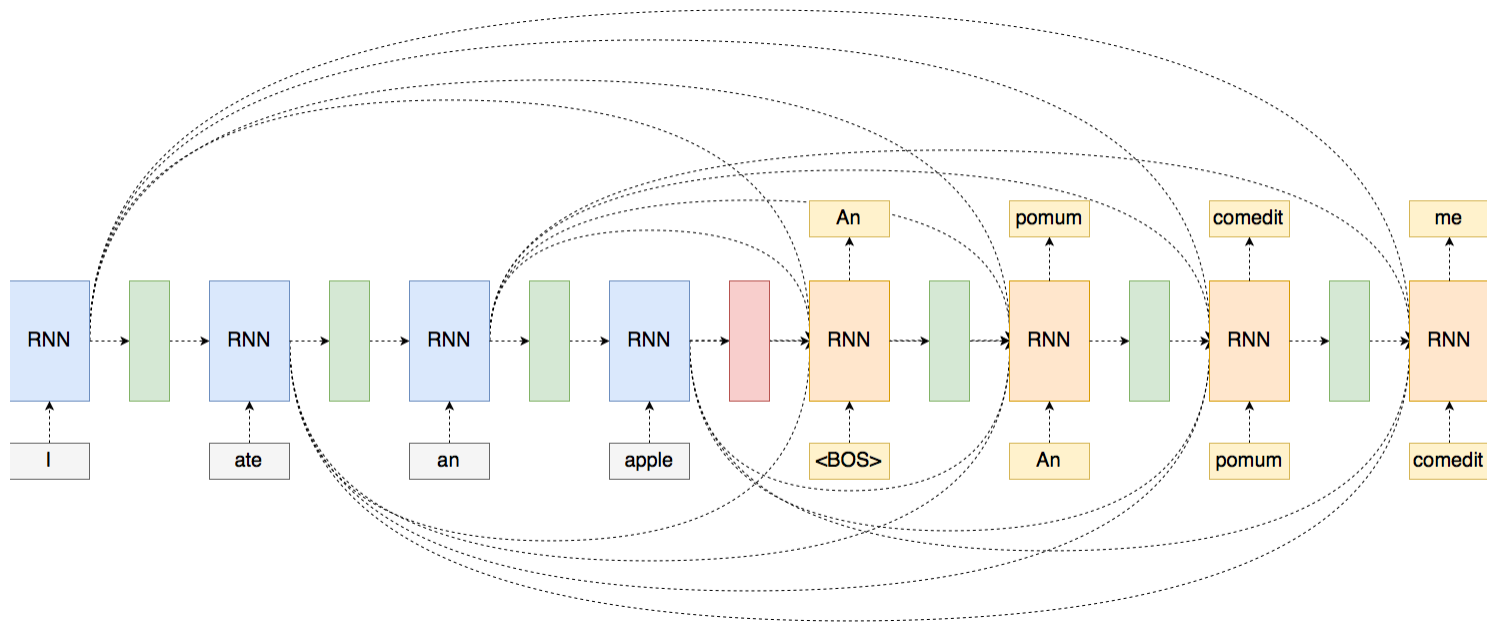- However, information of the entire input sequence is embedded into a single hidden state vector.

# What's the Potential Problem? (cont'd)

- Outputs at different time steps have particular meanings.
- However, synchrony between **input** and **output seqs** is not required.
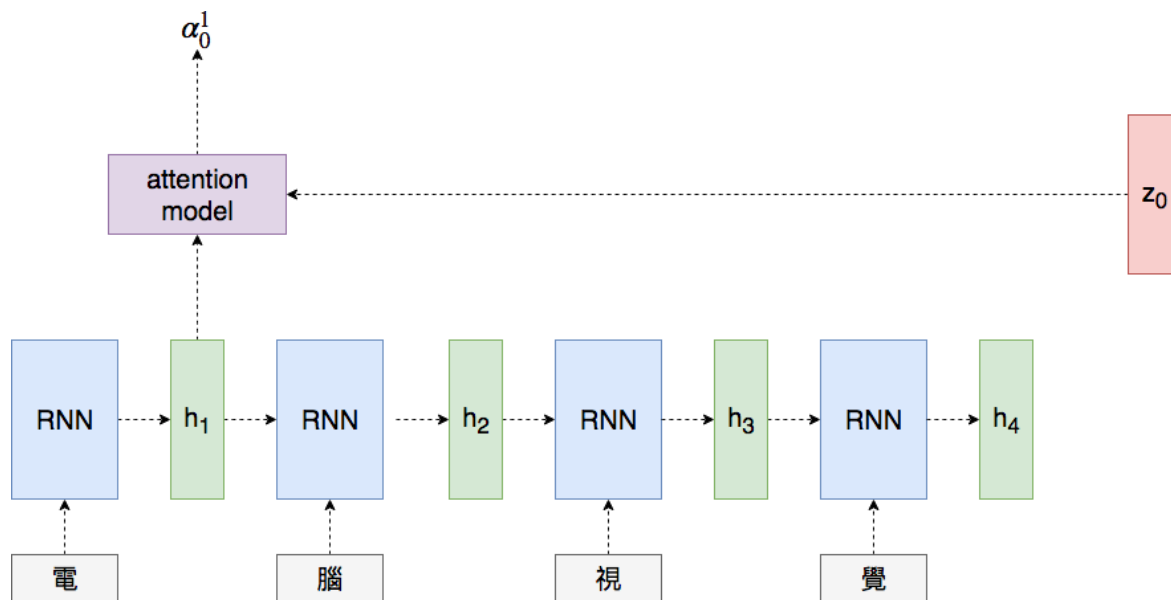
# What's the Potential Problem? (cont'd)

- Connecting every hidden state between encoder and decoder?
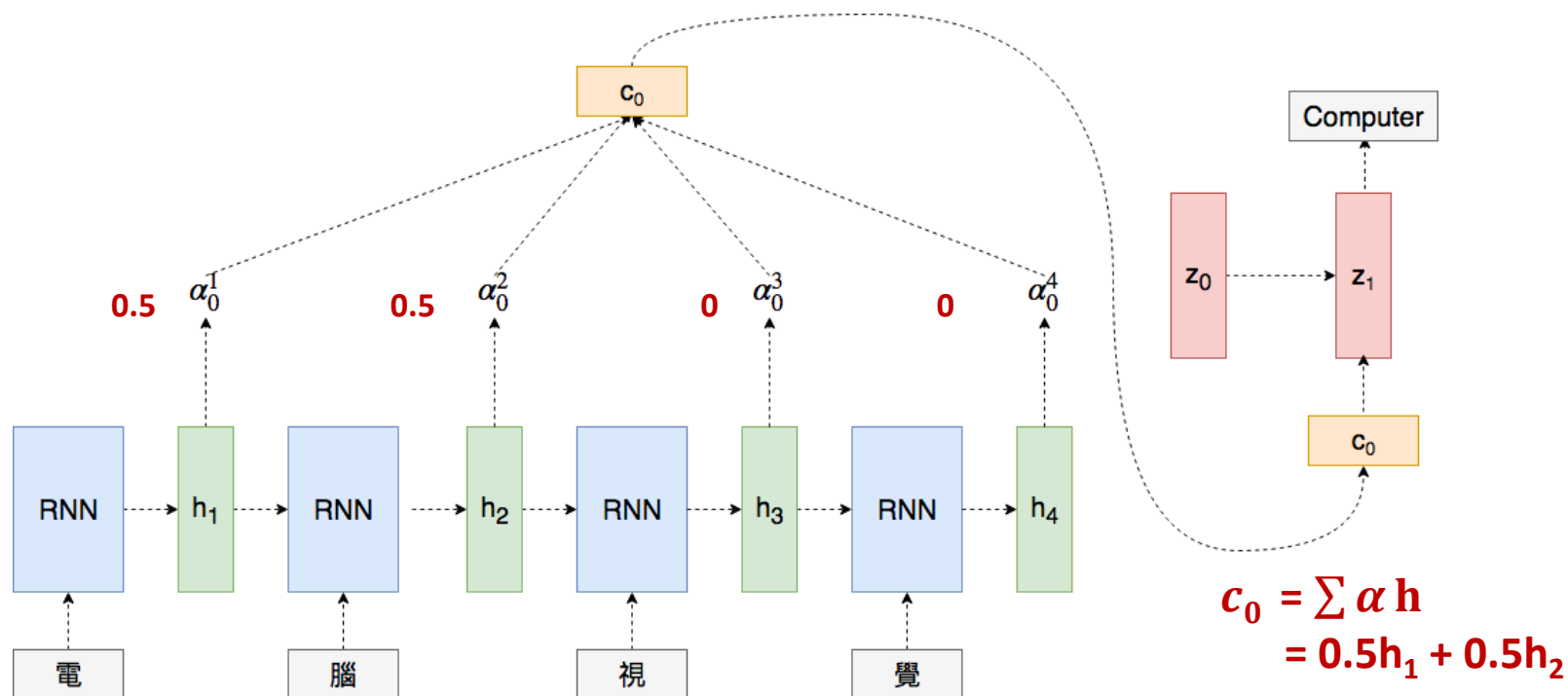


- Infeasible!
  - Both inputs and outputs are with varying sizes.
  - Overparameterized
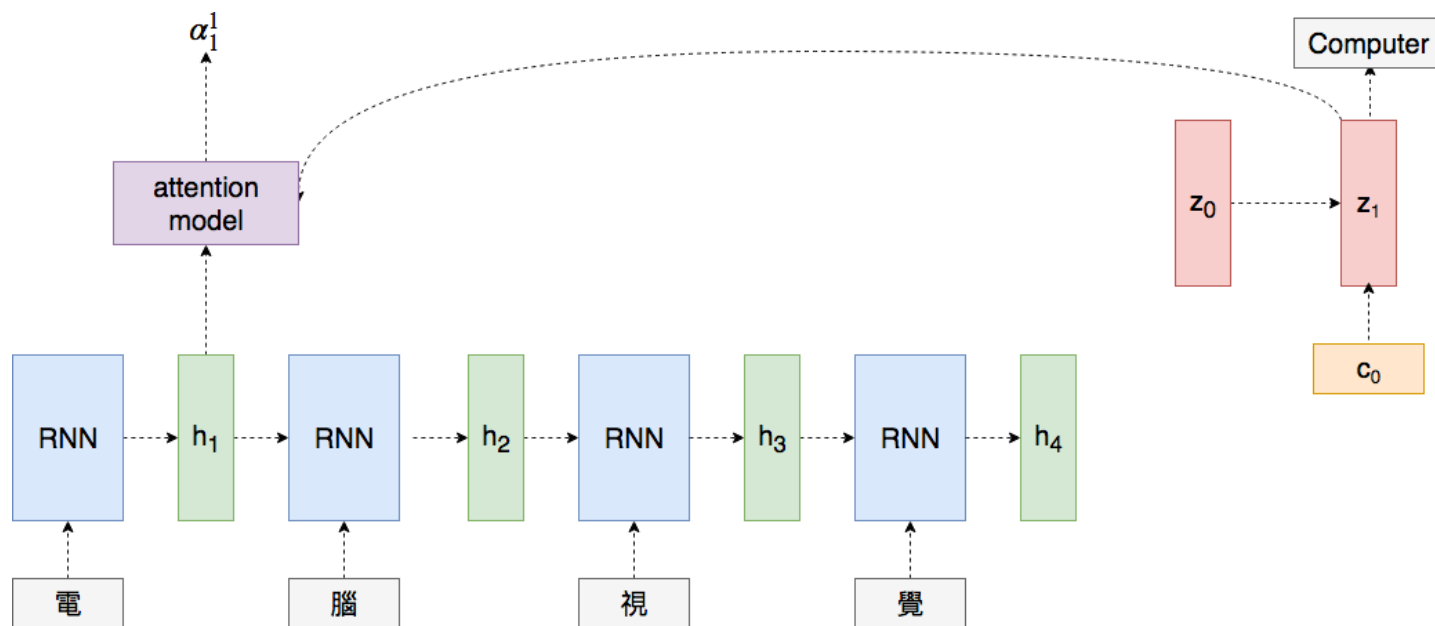
# Solution #1: Attention Model

- What should the attention model be?
  - A NN whose inputs are z and h while output is a scalar α, indicating the similarity between z and h.
- Most attention models are jointly learned with other parts of a network (e.g., classifier, regressor, etc.)
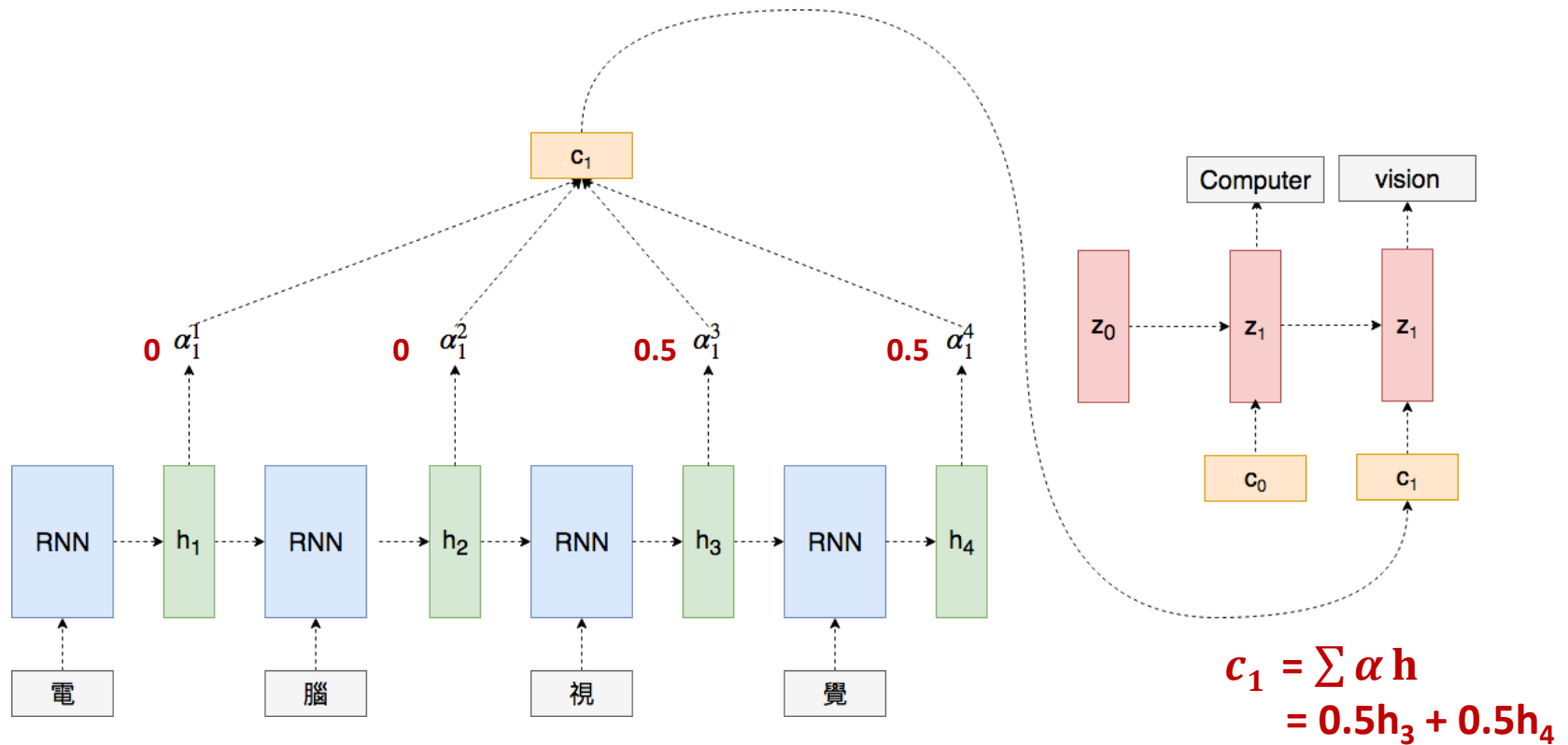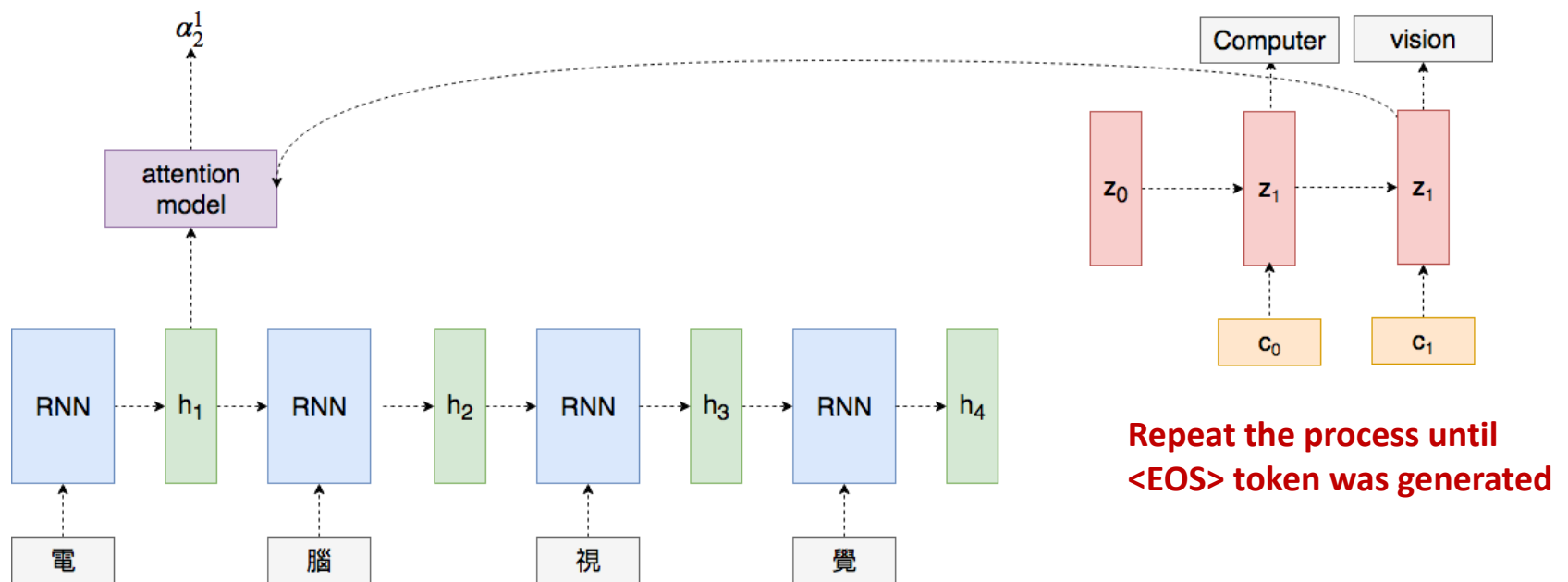  - Will see some examples later.

# Solution #1: Attention Model



$c_0 = \sum \alpha \, h$
$= 0.5h_1 + 0.5h_2$

# Solution: Attention Model

# Solution: Attention Model



$$c_1 = \sum \alpha\, h$$
$$= 0.5h_3 + 0.5h_4$$

# Solution: Attention Model



**Repeat the process until <EOS> token was generated**

# Selected Attention Models
# for Image-Based Applications

- Image Captioning
  - Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML '15

- Visual Question Answering
  - Zhu et al, "Visual7W: Grounded Question Answering in Images", CVPR '16

- Image Classification
  - Mnih et al, "Recurrent Models of Visual Attention", NIPS '14

# Image Captioning with Attention

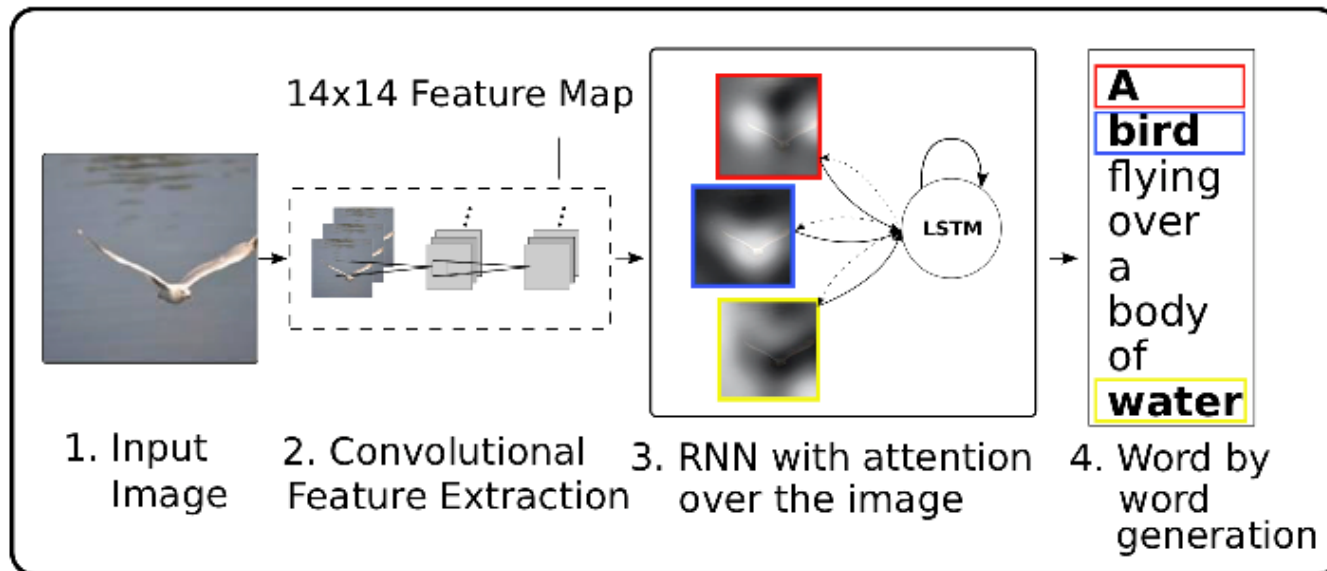- RNN focuses visual attention at different spatial locations when generating corresponding words during captioning.



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

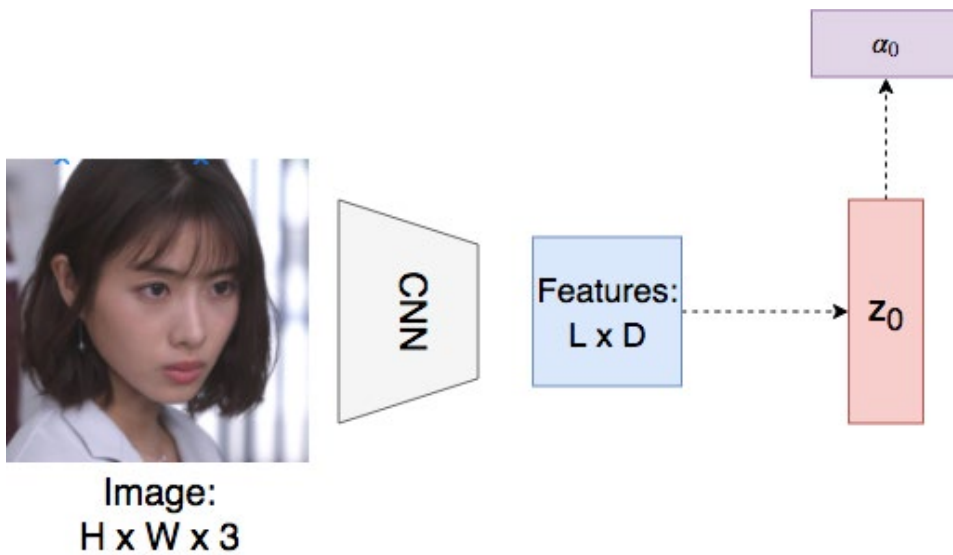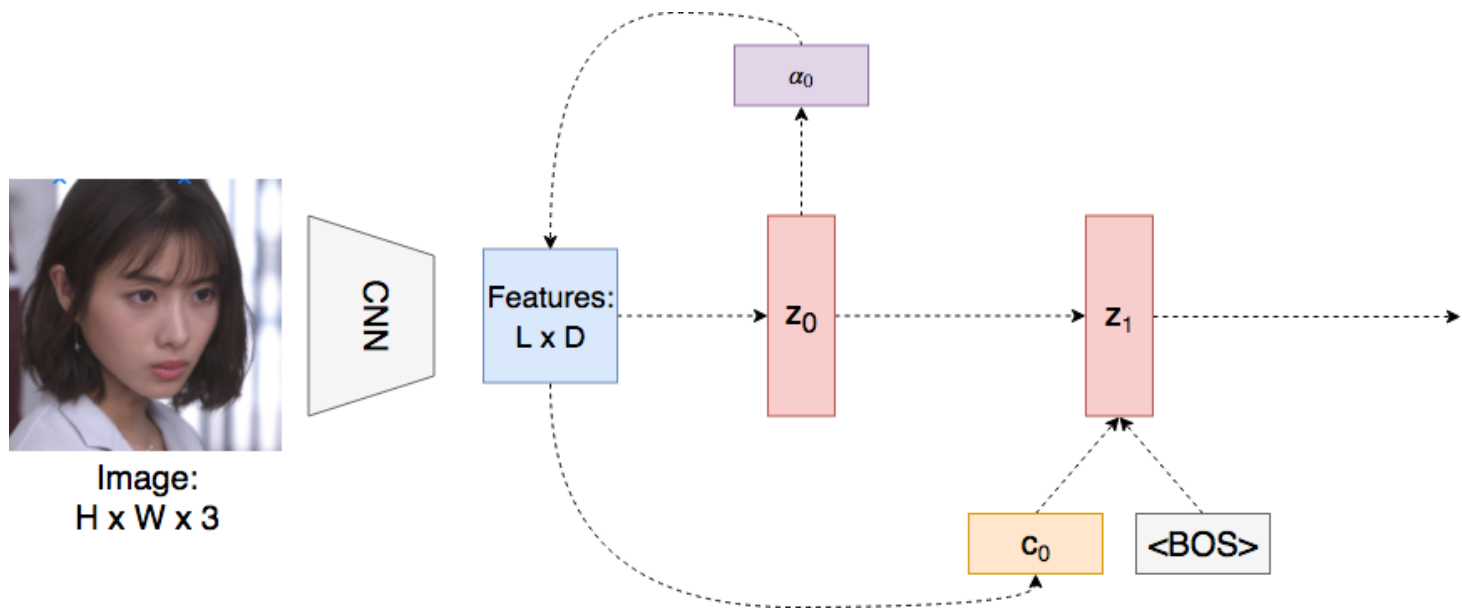# Image Captioning with Attention

# Image Captioning with Attention

**Distribution of attention over L locations**



**Weighted combination of features**
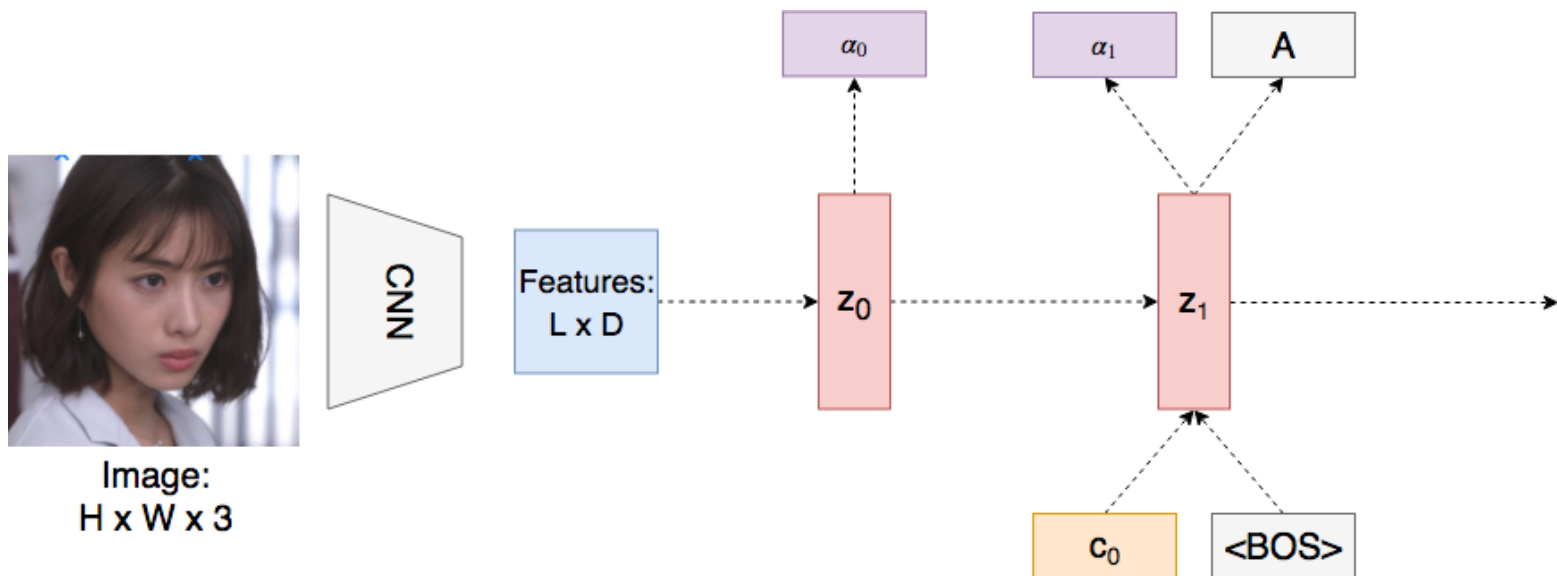
# Image Captioning with Attention
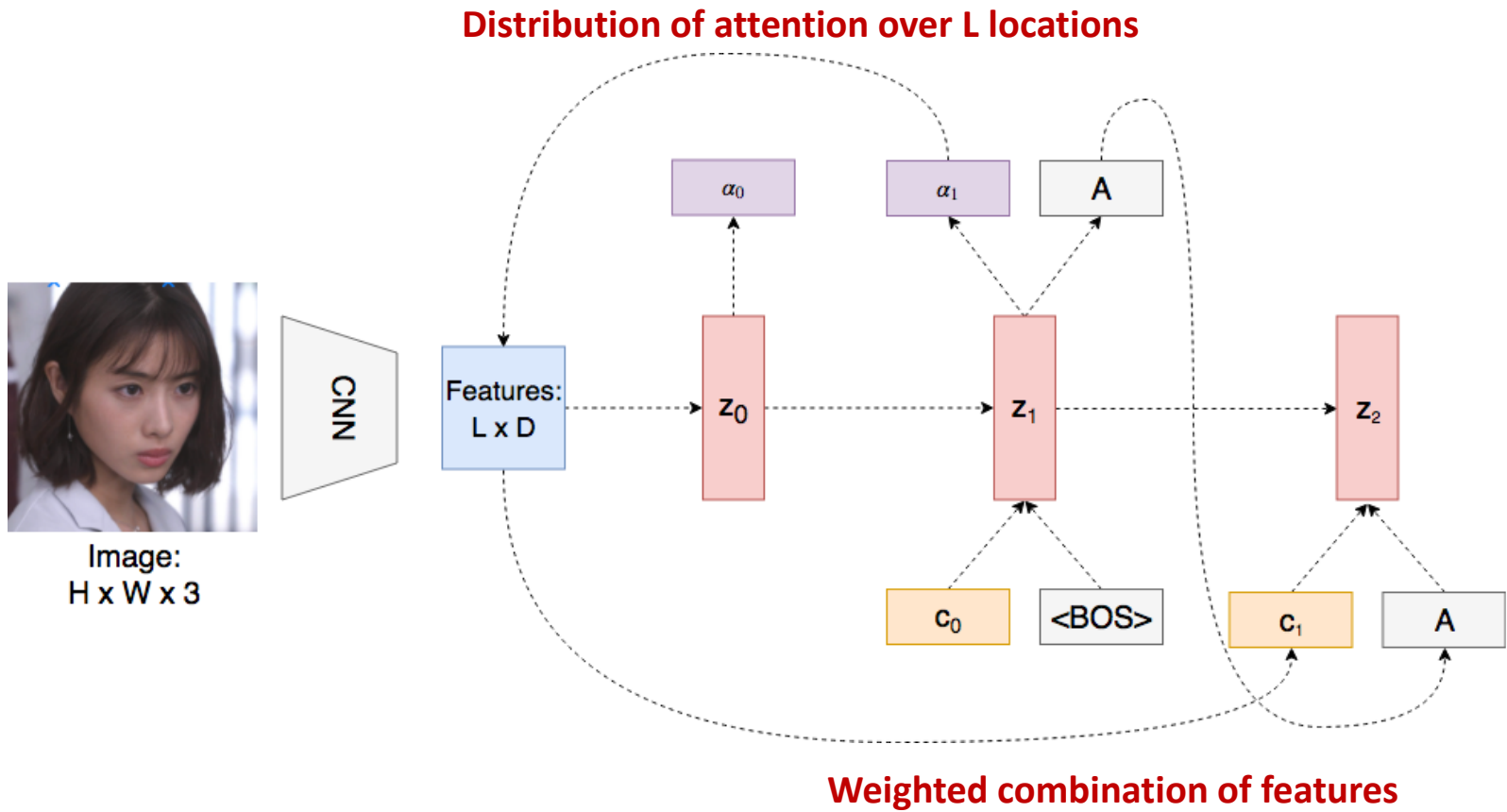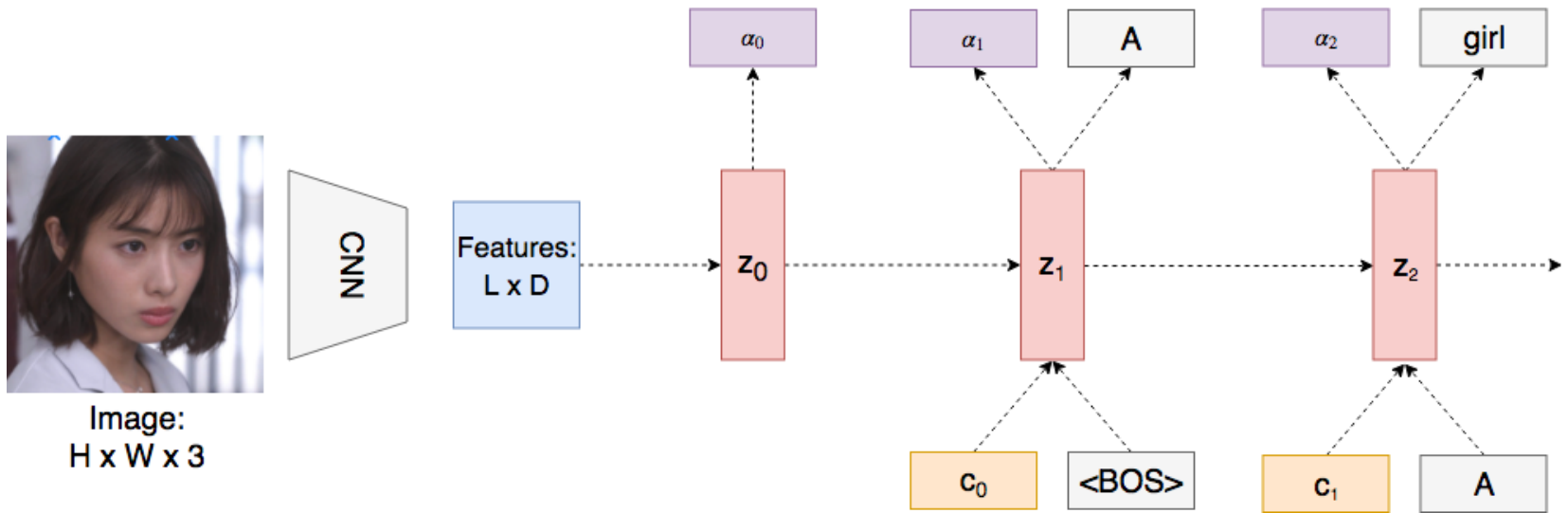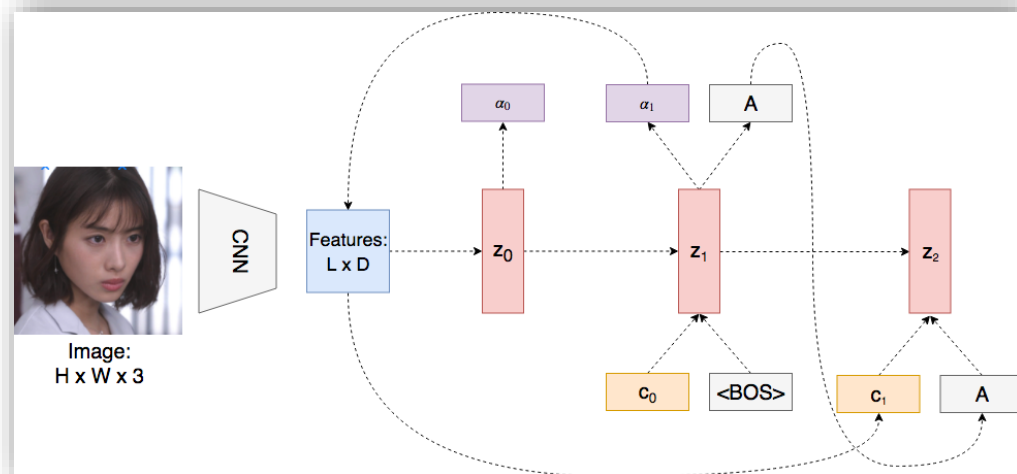
# Image Captioning with Attention

# Image Captioning with Attention



**Repeat the process until <EOS> token was generated**

# Image Captioning with Attention

- Visualization



A woman is throwing a frisbee in a park.

A dog is standing on a hardwood floor.

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

A little girl sitting on a bed with a teddy bear.

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

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

# Selected Attention Models
# for Image-Based Applications

- Image Captioning
  - Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML '15
- Visual Question Answering
  - Zhu et al, "Visual7W: Grounded Question Answering in Images", CVPR '16
- Image Classification
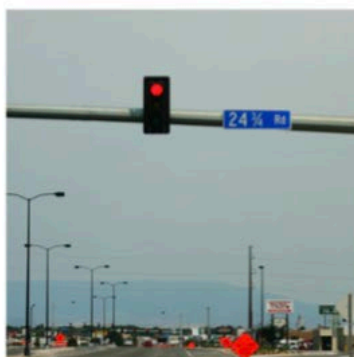  - Mnih et al, "Recurrent Models of Visual Attention", NIPS '14

# Visual Question Answering
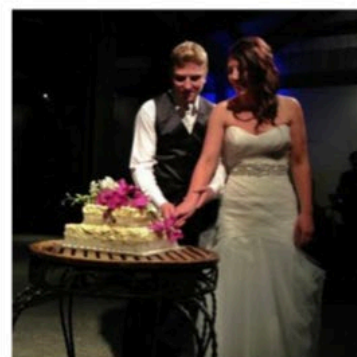
- Examples of multiple-choice QA & pointing QA



Q: What endangered animal is featured on the truck?

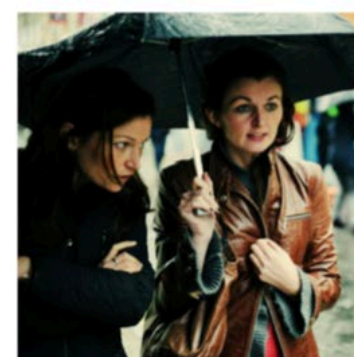A: **A bald eagle.**
A: A sparrow.
A: A humming bird.
A: A raven.

Q: Where will the driver go if turning right?

A: **Onto 24 ¾ Rd.**
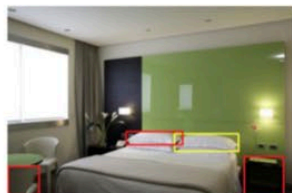A: Onto 25 ¾ Rd.
A: Onto 23 ¾ Rd.
A: Onto Main Street.

Q: When was the picture taken?

A: **During a wedding.**
A: During a bar mitzvah.
A: During a funeral.
A: During a Sunday church service.

Q: Who is under the umbrella?

A: **Two women.**
A: A child.
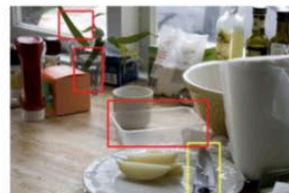A: An old man.
A: A husband and a wife.

Q: Which pillow is farther from the window?

Q: Which step leads to the tub?

Q: Which is the small computer in the corner?
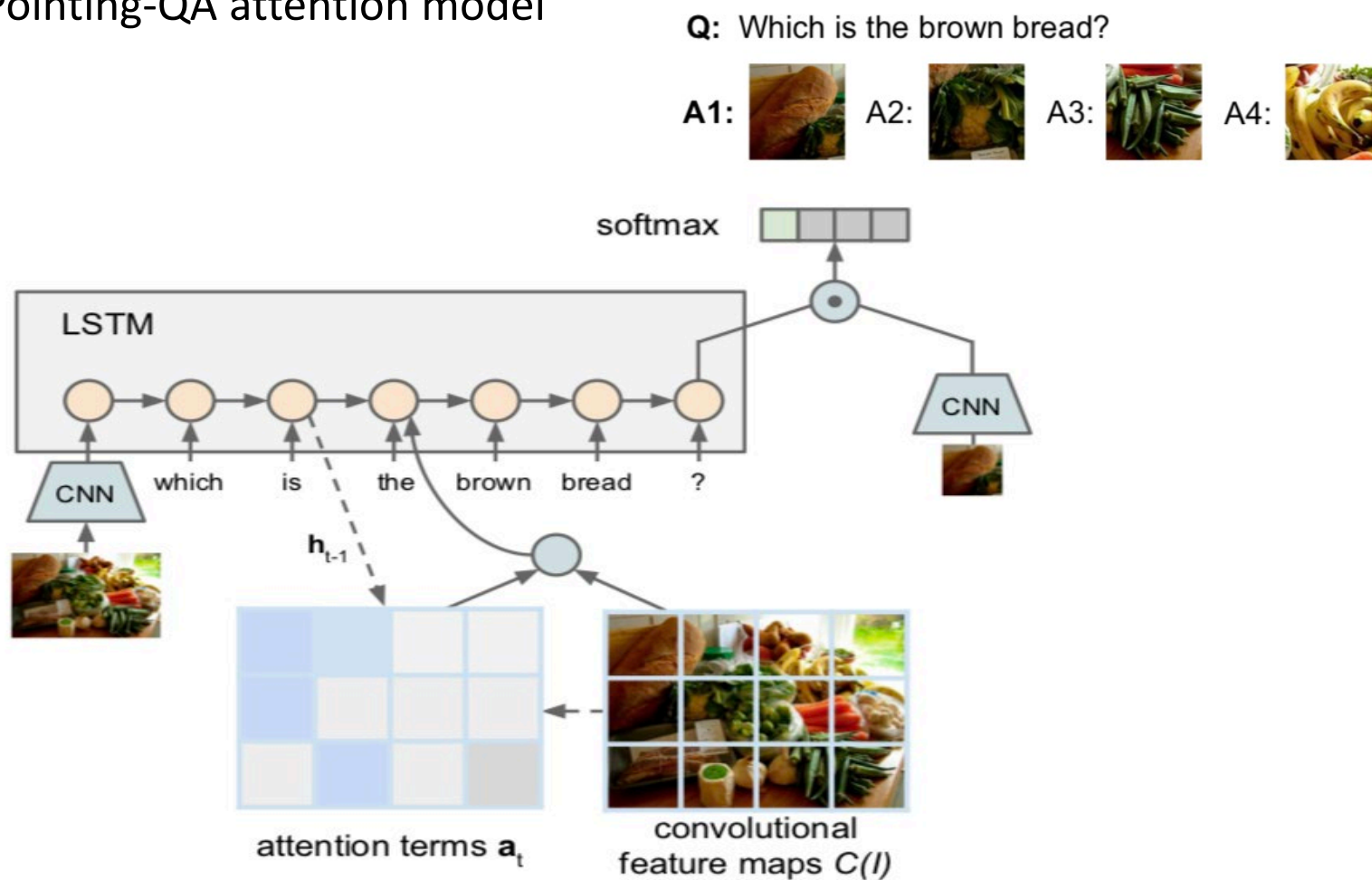
Q: Which item is used to cut items?

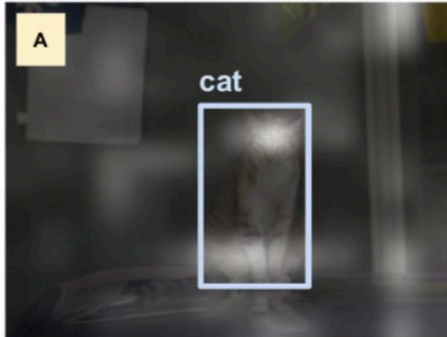Q: Which doughnut has multicolored sprinkles?

Q: Which man is wearing the red tie?

20

# Visual Question Answering with Attention

- Pointing-QA attention model

# Visual Question Answering with Attention (cont'd)



**A & B** (answers related to physicla objs):
The peaks of the attention maps reside in the bounding boxes of the target objects.

**C & D** (answers w/ non-physical objs):
The bottom two examples show QA pairs with answers not explicitly containing objects. The attention heat maps are scattered around the image grids.

# Selected Attention Models
# for Image-Based Applications

- Image Captioning
  - Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML '15
- Visual Question Answering
  - Zhu et al, "Visual7W: Grounded Question Answering in Images", CVPR '16
- Image Classification
  - Mnih et al, "Recurrent Models of Visual Attention", NIPS '14

# Glimpse Sensor & Glimpse Network

**Glimpse sensor**: extracts a retina-like representation centered at $l_{t-1}$ that contains multiple resolution patches.
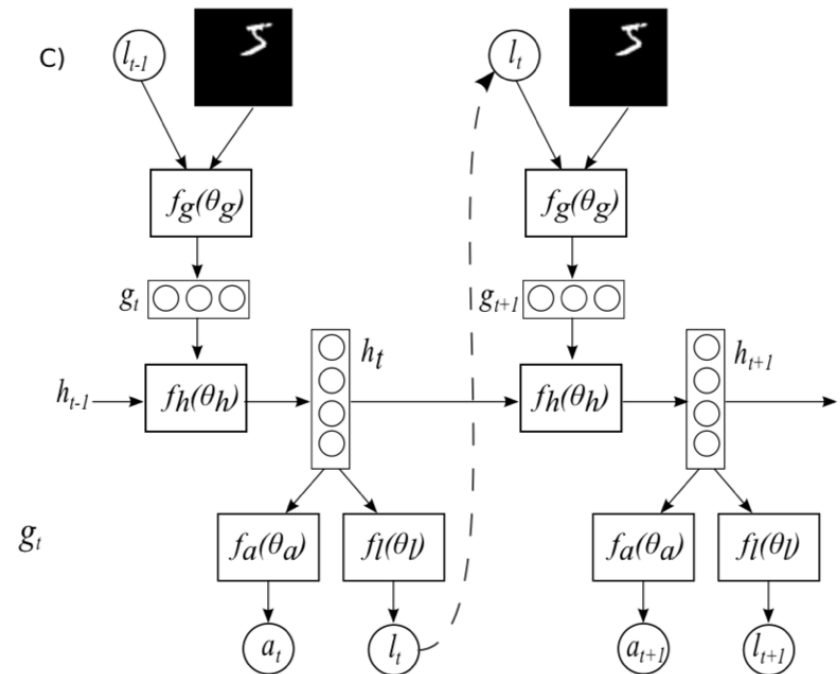


**Glimpse network**: given location $l_{t-1}$ and image $x_t$, use the glimpse sensor to extract retina representation, which is mapped into a joint hidden space.

**RNN-based model architecture**: the core network takes the glimpse representation as input with the hidden state vector from the prevision step, and outputs the new hidden state resulting in **location** and **action** networks to predict the next location to attend and the associated action.

Mnih et al, "Recurrent Models of Visual Attention", NIPS '14

# Architecture: RNN with Attention Models

# Example Results



- Original MNIST

- Translated MNIST

# Example: Actual Glimpse Path

# What to Cover Today…

- **Recurrent Neural Network** & **Transformer**
  - Attention in RNN
  - *Attention is All You Need*: Transformer
  - Transformer for Visual Analysis
    - Visual Classification
    - Semantic Segmentation & More

- **Vision & Language**
  - Image Captioning
  - Text-to-Image Synthesis



"a corgi wearing a bow tie and a birthday hat"



*Teddy bears shopping for groceries in the style of ukiyo-e*

# RNN with Attention is Good, But..

- Attention in a pre-defined sequential order

- Information loss due to long sequences…

OUTPUT | I am a student

INPUT | Je suis étudiant

# RNN with Attention is Good, But..

- Connecting every hidden state between encoder and decoder?



- Infeasible!
  - Both inputs and outputs are with varying sizes.
  - Overparameterized

# Solution #2: Transformer

- "Attention is all you need", NeurIPS 2017

- More details available at:
  http://jalammar.github.io/illustrated-transformer/

# Transformer

- "Attention is all you need", NeurIPS 2017
- Self-attention for text translation

# Self-Attention (1/5)

- Query **q**, key **k**, value **v** vectors are learned from each input **x**

$$q_i = W^Q x_i$$
$$k_i = W^K x_i$$
$$v_i = W^V x_i$$

# Self-Attention (2/5)

- Relation between each input is modeled by inner-product of query **q** and key **k**.

$$a_{1,i} = \frac{q_1 \cdot k_i}{\sqrt{d}}, \text{ where } a \in R, q, k \in R^d$$

# Self-Attention (3/5)

- SoftMax is applied:

$$0 \leq \hat{a}_i = e^{a_i} / \sum_{j}^{N} e^{a_j} \leq 1 \text{ , for I =1, ..., N}$$



| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

$\hat{a}_{1,1}$ $\hat{a}_{1,2}$ $\hat{a}_{1,N}$

**SoftMax**

$a_{1,1}$ $a_{1,2}$ $a_{1,N}$

$q_1$ $k_1$ $v_1$ $q_2$ $k_2$ $v_2$ $q_N$ $k_N$ $v_N$

$x_1$ $x_2$ ... $x_N$

# Self-Attention (4/5)

- Value vectors **v** are aggregated
  with attention weight $\hat{a}$ , i.e., $y_1 = \sum_i^N \hat{a}_i \cdot v_i$



| Input | | Thinking | | Machines |
|---|---|---|---|---|
| Embedding | | $x_1$ | | $x_2$ |
| Queries | | $q_1$ | | $q_2$ |
| Keys | | $k_1$ | | $k_2$ |
| Values | | $v_1$ | | $v_2$ |
| Score | | $q_1 \cdot k_1 = 112$ | | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | | 14 | | 12 |
| Softmax | | 0.88 | | 0.12 |
| Softmax X Value | | $v_1$ | | $v_2$ |
| Sum | | $z_1$ | | $z_2$ |

# Self-Attention (5/5)

- All $y_i$ can be computed in parallel

- $y_i$ considers $x_1 \sim x_N$, modeling long-distance dependencies.

- Global feature can be obtained by average-pooling over $y_1 \sim y_N$

# Self-Attention: Implementation

- Input sequence can be represented as a N x $d_{in}$ matrix

- * denotes matrix multiplication

$x_i \in R^{d_{in}}$

Input matrix

# Self-Attention: Implementation

- Output matrix Y
- All operations are **matrix multiplication**, can be parallelized on GPU.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

# Multi-Head Self-Attention (1/4)

- Perform self-attention at different subspaces, implying attention over different input feature types (e.g., representations, modalities, positions, etc.)

# Multi-Head Self-Attention (2/4)

- Perform self-attention at different subspaces, implying attention over different input types

- See example below



Attention weights
of Head 1

Attention weights
of Head 2

# Multi-Head Self-Attention (3/4)

- A 2-head example, output of two heads are concatenated.

# Multi-Head Self-Attention (4/4)

- A 2-head example, output of two heads are concatenated.

# The Residuals

- A residual connection followed by layer normalization

# The Decoder in Transformer

- Design similar to that of encoder,
  except the 1$^{st}$ decoder takes additional inputs (of predicted word embeddings).

# The Decoder in Transformer

- Design similar to that of encoder,
  except the 1st decoder takes additional inputs (of predicted word embeddings).

# Recap: Transformer

- "Attention is all you need", NeurIPS 2017

- We didn't cover positional encoding (particularly for language translation)

- More info available at:
  http://jalammar.github.io/illustrated-transformer/

# What to Cover Today...

- **Recurrent Neural Network & Transformer**
  - Attention in RNN
  - *Attention is All You Need*: Transformer
  - Transformer for Visual Analysis
    - Visual Classification
    - Semantic Segmentation & More

- **Vision & Language**
  - Image Captioning
  - Text-to-Image Synthesis



"a corgi wearing a bow tie and a birthday hat"



*Teddy bears shopping for groceries in the style of ukiyo-e*

# Vision Transformer

- "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR, 2021. (Google Research)

- Partition the input image into a **patch sequence**

- An additional **token** (*) is appended to perform attention on patches

- Both the "*" token and positional embeddings (denoted by 0, 1, 2 …) are **trainable vectors**



Vision Transformer (ViT)

# Query-Key-Value Attention in ViT

- Assume that the input is partitioned into 4 patches and the feature dimension is 3, that is, P=4 and D=3

- Note that there are (P+1) rows since we have an additional token



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

$P$: Patch Number
$D$: Dimension
$W_{Q,K,V}$: Learnable Matrices
$Q, K, V$: Query, Key, Value
$A$: Attention Matrix

# Query-Key-Value Attention in ViT

- By performing attention, the input sequence X (of length P+1) is "transformed" into another sequence Y with the same length

- That is why it is called "**Transformer**" and how it is a **seq2seq** model



$P$: Patch Number
$D$: Dimension
$W_{Q,K,V}$: Learnable Matrices
$Q, K, V$: Query, Key, Value
$A$: Attention Matrix

# Query-Key-Value Attention in ViT

- In standard vision transformer, we only take the **first output token** of the output sequence (the **first row** of Y) for classification purposes
- This corresponds to the output when **token "0"** serves as query



$P$: Patch Number
$D$: Dimension
$W_{Q,K,V}$: Learnable Matrices
$Q, K, V$: Query, Key, Value
$A$: Attention Matrix

52

# Visualization

- To visualize the attention maps, we take the attention scores from the **first row** of A (when token "0" serves as query)

- Note the first element is excluded, and thus there are **P scores** corresponding to the P image patches



$P$: Patch Number
$D$: Dimension
$W_{Q,K,V}$: Learnable Matrices
$Q, K, V$: Query, Key, Value
$A$: Attention Matrix

# Example Visualization

# PS-ViT

- Vision Transformers with Progressive Sampling
- Progressively select important patches by shifting patch centers



Yue et al. "Vision transformer with progressive sampling." ICCV 2021

# PS-ViT (cont'd)

# Example Visualization

# Transformer for Semantic Segmentation

- Segmentation via attention



Strudel et al. "Segmenter: Transformer for Semantic Segmentation." ICCV 2021

# Transformer for Semantic Segmentation

- Using different class tokens ("Tree", "Sidewalk", "Person", …) as queries



Segmenter

Input Image

Patch embedding
Position embedding
Patch encoding

Tree
Sidewalk
Person

Flatten and Project

Transformer Encoder

Mask Transformer

Scalar Product

Class Masks

Upsample and Argmax

Segmentation Map

# Example Visualization



(a) Patch size $32 \times 32$     (b) Patch size $16 \times 16$     (c) Patch size $8 \times 8$     (d) Ground Truth

# Self-Supervised Learning (SSL) for Transformer

- Learning discriminative representations from **unlabeled** data
- Create self-supervised tasks via **data augmentation**

Chen et al. "A simple framework for contrastive learning of visual representations." ICML 2020

# Self-Supervised Transformer

- Vision Transformer + **SSL**
- Maximize the similarity between the augmented version and itself
- Avoid collapse with **student-teacher** network



$$\theta_T \leftarrow \tau\theta_T + (1-\tau)\theta_S$$

Vision Transformer

Caron et al. "Emerging properties in self-supervised vision transformers." ICCV 2021

# Qualitative & Quantitative Results



| Method | Arch. | Param. | im/s | Linear | $k$-NN |
|---|---|---|---|---|---|
| Supervised | RN50 | 23 | 1237 | 79.3 | 79.3 |
| SCLR [12] | RN50 | 23 | 1237 | 69.1 | 60.7 |
| MoCov2 [15] | RN50 | 23 | 1237 | 71.1 | 61.9 |
| InfoMin [67] | RN50 | 23 | 1237 | 73.0 | 65.3 |
| BarlowT [81] | RN50 | 23 | 1237 | 73.2 | 66.0 |
| OBoW [27] | RN50 | 23 | 1237 | 73.8 | 61.9 |
| BYOL [30] | RN50 | 23 | 1237 | 74.4 | 64.8 |
| DCv2 [10] | RN50 | 23 | 1237 | 75.2 | 67.1 |
| SwAV [10] | RN50 | 23 | 1237 | **75.3** | 65.7 |
| DINO | RN50 | 23 | 1237 | **75.3** | **67.5** |

# What to Cover Today…

- **Recurrent Neural Network & Transformer**
  - Attention in RNN
  - *Attention is All You Need*: Transformer
  - Transformer for Visual Analysis
    - Visual Classification
    - Semantic Segmentation & More

- **Vision & Language**
  - Image Captioning
  - Text-to-Image Synthesis





"a corgi wearing a bow tie and a birthday hat"

*Teddy bears shopping for groceries in the style of ukiyo-e*

# *A picture is worth a thousand words...*
# Is it that simple?



- Thing
- Airplane
- Flying airplane in blue sky
- A Lufthansa MD-11 cargo plane in blue sky flying over mountainous terrain

# Vision + Language → ?

- Image Captioning

- Image Manipulation/Completion

- Composed Image Retrieval

- Visual Question Answering (VQA) and many more…

# Image Captioning



Applications: semantics understanding, image-text retrieval, medical AI, etc.

# Image Captioning (cont'd)

- Training a captioning model requires a large amount of image-caption data pairs

- Image captioning in the wild:
  - Describing images with novel content during inference
  - For example, COCO dataset has 80 object categories.
    How to generalize captioning models to Open Image (w/ 600 classes)?

- Domain-specific image captioning:
  - From general-purpose captioning to task-oriented captioning



COCO (80 classes)

Two pug **dogs** sitting on a **bench** at the beach.

A **child** is sitting on a **couch** and holding an **umbrella**.



Open Images (600 classes)

goat    artichoke    accordion

dolphin    waffle    balloon

# Image Captioning *in the Wild*

- **Novel Object Captioning (NOC)**
    - Training with captioned and uncaptioned data
      captioned data: labeled image data with captions (e.g., COCO)
      uncaptioned data: only labels of novel classes available (e.g., Open Images)



**We have captioning data**



**Data with labels for novel objects
but w/o captions**

# Novel Object Captioning

- **VIVO**: **Vi**sual **Vo**cabulary Pre-Training for Novel Object Caption Captioning (AAAI'21)
    - Pre-training a cross-modality Transformer for vision & language tasks
    - Pre-training really matters, since it's been observed in
        - Computer Vision (e.g., models pre-trained on ImageNet)
        - Natural Language Processing (e.g., BERT pre-trained on Wikipedia)

Object detection,
semantic segmentation, etc.

Question answering,
Sentence classification, etc.

# Novel Object Captioning (cont'd)

- **VIVO**: **Vi**sual **Vo**cabulary Pre-Training for Novel Object Caption Captioning
  - Pre-training: uncaptioned image data containing novel class labels
  - Fine-tuning: (a limited amount of) image data with class labels & descriptions



(a) Pre-training: learn visual vocabulary

(b) Fine-tuning: learn sentence description

# Novel Object Captioning (cont'd)

- **VIVO**: **Vi**sual **Vo**cabulary Pre-Training for Novel Object Caption Captioning
  - Pre-training: uncaptioned image data containing novel class labels
  - Fine-tuning: (a limited amount of) image data with class labels & descriptions
  - Inference:
    - Inputs: image (with region features & tags) & [CLS]
    - Output: caption



**(c) Inference: novel object captioning**

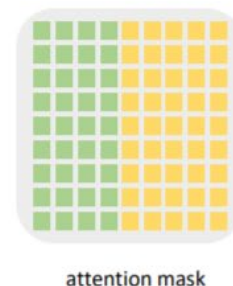# Novel Object Captioning (cont'd)

- VIVO: Visual Vocabulary Pre-Training for Novel Object Caption Captioning
  - Properly aligned image and text data for captioning

# Beyond Image Captioning: Unified Vision & Language Model

- **Oscar**: **O**bject-**S**emanti**c**s **A**ligned **P**re-training for Vision-Language Tasks (ECCV'20)
  - Training data:
    triplets of caption-tag-region
  - Objectives:
    1. Masked token loss for words & tags
    2. Contrastive loss tags and others
  - Fine-tuning:
    5 vision & language tasks (VQA, image-text retrieval, image captioning, NOC, etc.)

# Semantics-Aligned Pre-training for V+L Tasks

- **Oscar**: **O**bject-**S**emanti**c**s **A**ligned P**r**e-training for Vision-Language Tasks (ECCV'20)
  - Training:
    - Inputs: triplets of caption-tag-region
    - Objectives: Masked token loss for words & tags + Contrastive loss tags and others
  - Fine-tuning:
    5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)



75

# Semantics-Aligned Pre-training for V+L Tasks (cont'd)

- **Oscar**: **O**bject-**S**emanti**c**s **A**ligned **P**re-training for Vision-Language Tasks (ECCV'20)
  - Training:
    - Inputs: triplets of word-tag-region
    - Objectives: Masked token loss for words & tags + Contrastive loss tags and others
  - Fine-tuning:
    - 5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)

# Semantics-Aligned Pre-training for V+L Tasks (cont'd)

- **Oscar**: **O**bject-**S**emantics **A**ligned **Pr**e-training for Vision-Language Tasks (ECCV'20)
  - Fine-tuning:
    5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)
  - Take **image captioning** as an example
    - Training: triplets of image regions features + object tags + captions as inputs;
      caption tokens with full attention on image regions/tags but not the other way around
    - Inference: image regions, tags and **[CLS]** as inputs,
      with **[MASK]** tokens sequentially added/predicted

*Holding an apple* ⟷ 🍎 or 📱

- **Oscar** (cont'd)
  - Fine-tuning:
    5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)
  - Take **image-text retrieval** as an example
    - Training: aligned/mis-aligned image-text pairs as positive/negative input pairs,
      with **[CLS]** for binary classification (1/0)
    - Inference: for either image or text retrieval,
      calculate <u>classification score</u> of **[CLS]** for the input query

# Image Change Captioning

- **Goal: Caption the difference(s) between input images**
  - Inputs: images with difference(s) + ground truth caption for the difference(s)
  - For image pair with one change



- For image pair with multiple changes (Yue et al., ICCV'21)



Change captions

Caption 1: The large gray rubber sphere has disappeared. (delete)

Caption 2: There is no longer a large cyan metal cube. (delete)

Caption 3: The large brown metal sphere was moved from its original location. (move)

Caption 4: The small yellow rubber cylinder was replaced by a small red rubber sphere. (replace)
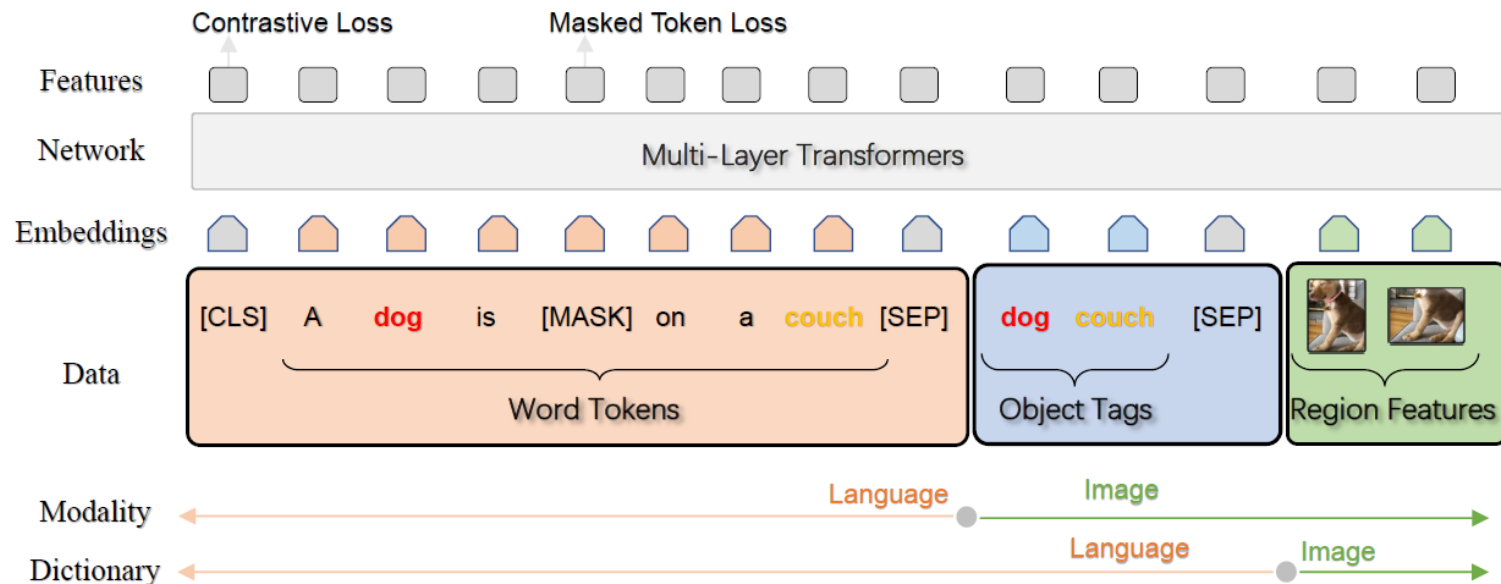
# Image Change Captioning

- **Goal: Caption the difference(s) between input images**
  - Inputs: images with difference(s) + ground truth caption for the difference(s)
  - For image pair with one change



  - E.g., Robust Image Change Captioning, Dong et al., ICCV'19

# What to Cover Today…

- **Recurrent Neural Network** & **Transformer**
  - Attention in RNN
  - *Attention is All You Need*: Transformer
  - Transformer for Visual Analysis
    - Visual Classification
    - Semantic Segmentation & More

- **Vision & Language**
  - Image Captioning
  - Text-to-Image Synthesis





"a corgi wearing a bow tie and a birthday hat"

*Teddy bears shopping for groceries in the style of ukiyo-e*

# Image Manipulation

- Text-to-Image Synthesis & Manipulation
  - Task #1: Text-to-image generation
    - Produce images based on their descriptions
    - Training: image-caption pairs
    - Recent works: Show & Tell (CVPR'15), StackGAN (ICCV'17), DALL-E (OpenAI)
    - Example:

*Teddy bears shopping for groceries in the style of ukiyo-e*

DALL-E

- Text-to-Image Synthesis & Manipulation (cont'd)
  - Text-to-image generation
  - Task #2: Image manipulation by text instruction
    - Allow users to edit an image with complex instructions (e.g., add, remove, etc.)
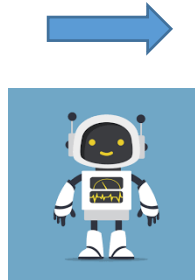    - Training: reference image & instruction as inputs; target image as output
    - E.g., GeNeVa-GAN (ICCV'19), TIM-GAN (MM'21)
  - Task #3: Text/caption-guided image manipulation
    - Edit image regions to match image descriptions
    - Training: image-caption pairs
    - E.g., GLIDE (OpenAI'21), Tedi-GAN (CVPR'21), ManiTrans (CVPR'22)



*make middle-left small gray object large*

**Fig. 1 Example of image manipulation by text instruction**



*A yellow tower.*

**Fig. 2 Example of text (caption)-guided image manipulation**

# Challenges in Text-Guided Image Manipulation

- ## Localization
    - Needs to identify objects in an image, locate the target location or objects of interest
    - Requires image understanding (with both semantics & spatial info)

- ## Manipulation
    - Needs to understand the input caption/instruction for manipulating images
    - Preserves object interaction and style to alleviate possible mismatch after manipulation

| Input | Localization | Manipulation |
|-------|--------------|--------------|

*a fire in the background*

# Text-Guided Image Manipulation (cont'd)

- Remarks & Opportunities
  - Not easy to collect training data with full supervision
  - Large-scale V&L pre-training models available (e.g., CLIP)
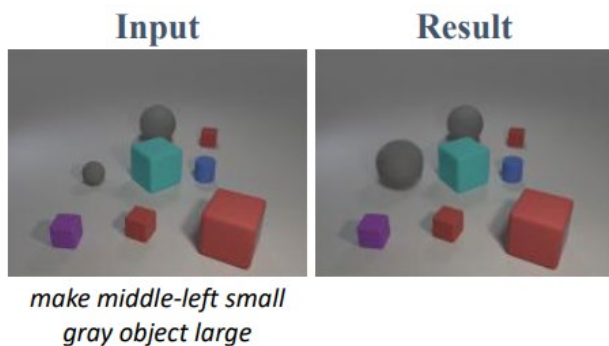  - **Task #2** (manipulate by instruction) vs**. Task #3** (manipulate by text guidance)

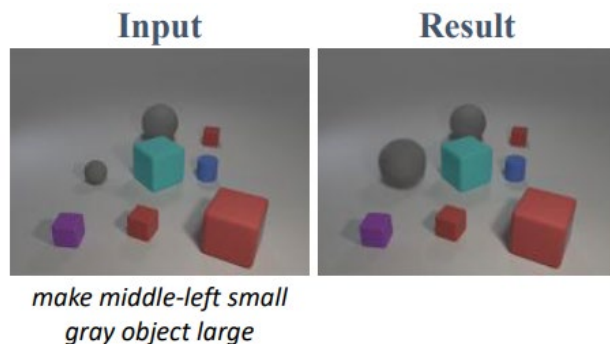

Fig. 1 Example of image manipulation by text instruction



Fig. 2 Example of text (caption)-guided image manipulation

  - Can scale up to industrial level with paired training data available

# Selected Work on Text-Guided Image Manipulation

- GLIDE
  - Developed by OpenAI in 2021
  - Training:
    - Image-caption pairs and randomly generated masks
    - Learns to recover the missing part based on the caption
  - Testing: image, caption, and mask annotated by user
  - Later extended by a recent CVPR'22 work (DiffusionCLIP) for semantics improvements



BEFORE  AFTER



"a corgi wearing a bow tie and a birthday hat"

"only one cloud in the sky today"

# Composed Image Retrieval

- Goal
  - Given a reference image and its modification text (i.e., a cross-modal query), retrieve the target image from the database
  - Very different from image-text or text-image retrieval!



+ I want to change it to longer sleeves and yellow in color.
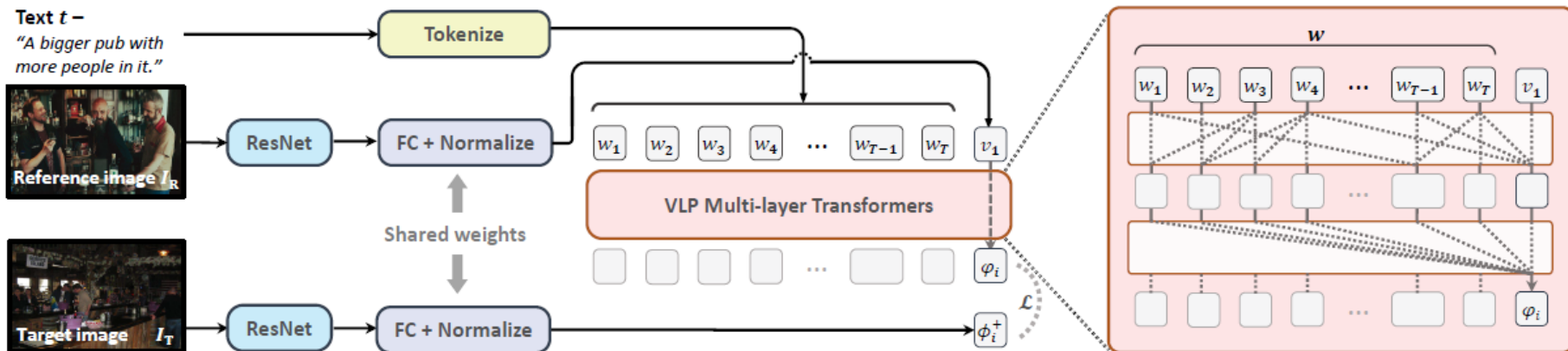
Reference Image            Modification Text                              Target Image

# Composed Image Retrieval with Pre-trained V&L Models

- Composed Image Retrieval using Pretrained LANguage Transformers (CIRPLANT)

  - Extract image features by a pre-trained ResNet

  - Aggregate information from modification text and reference image by a pre-trained **OSCAR**

  - Instead of use of output token [CLS], the derived output image feature ϕ is used for retrieval



Image Retrieval on Real-life Images with Pre-trained Vision-and-Language Models. *Zheyuan Liu et al*. ICCV 2021

# Retrieval with Text-Explicit Matching & Implicit Similarity

- **A**ttention-based **R**etrieval with

  **T**ext-**E**xplicit **M**atching and **I**mplicit **S**imilarity (**ARTEMIS**)

  - Image search with free-form text modifier

  - Cross-modal learning and visual retrieval

    - Text-guided attention is introduced ARTEMIS
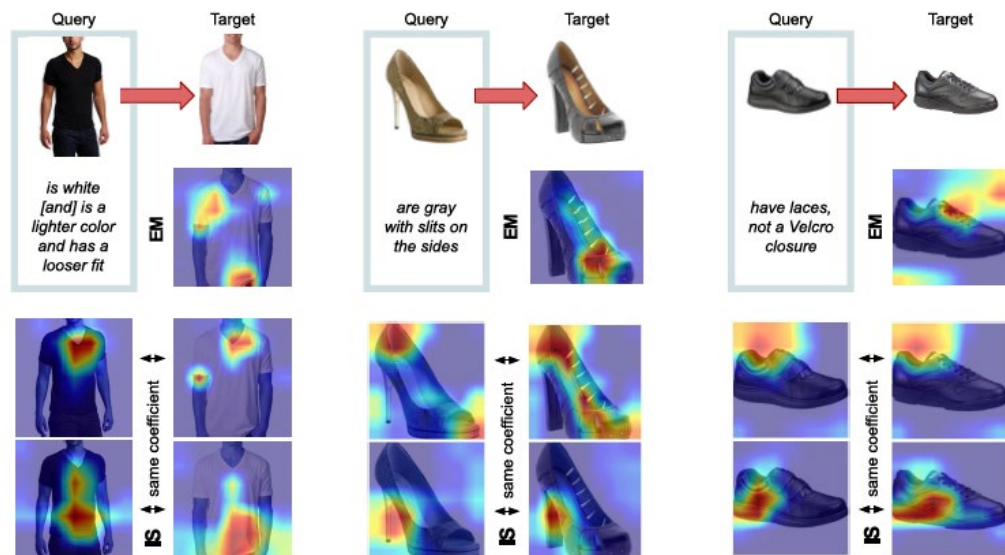
- **Attention-based Retrieval with**

  **Text-Explicit Matching and Implicit Similarity (ARTEMIS) (cont'd)**

  - **Implicit Similarity (IS):**

    attention mechanism focusing on what's not mentioned by text and should be preserved

  - **Explicit Matching (EM):**

    attention mechanism focusing on what's mentioned by text and should be changed.



ARTEMIS: Attention-based Retrieval with Text-Explicit Matching & Implicit Similarity. *Ginger Delmas et al.* ICLR 2022

- **Attention-based Retrieval with**
  **Text-Explicit Matching and Implicit Similarity (ARTEMIS) (cont'd)**
  - Example Results & Extension

# What to Cover Today...

- **Recurrent Neural Network & Transformer**
  - Attention in RNN
  - *Attention is All You Need*: Transformer
  - Transformer for Visual Analysis
    - Visual Classification
    - Semantic Segmentation & More

- **Vision & Language**
  - Image Captioning
  - Text-to-Image Synthesis

- **HW #3 is out!**



"a corgi wearing a bow tie and a birthday hat"



*Teddy bears shopping for groceries in the style of ukiyo-e*