

Visually Grounded Language Understanding and Generation

Jiasen Lu

Georgia Institute of Technology

A clean kitchen with a double sink, a black kettle on stove and a refrigerator.




A clean kitchen with a double **sink**, a black **kettle** on **stove** and a **refrigerator**.



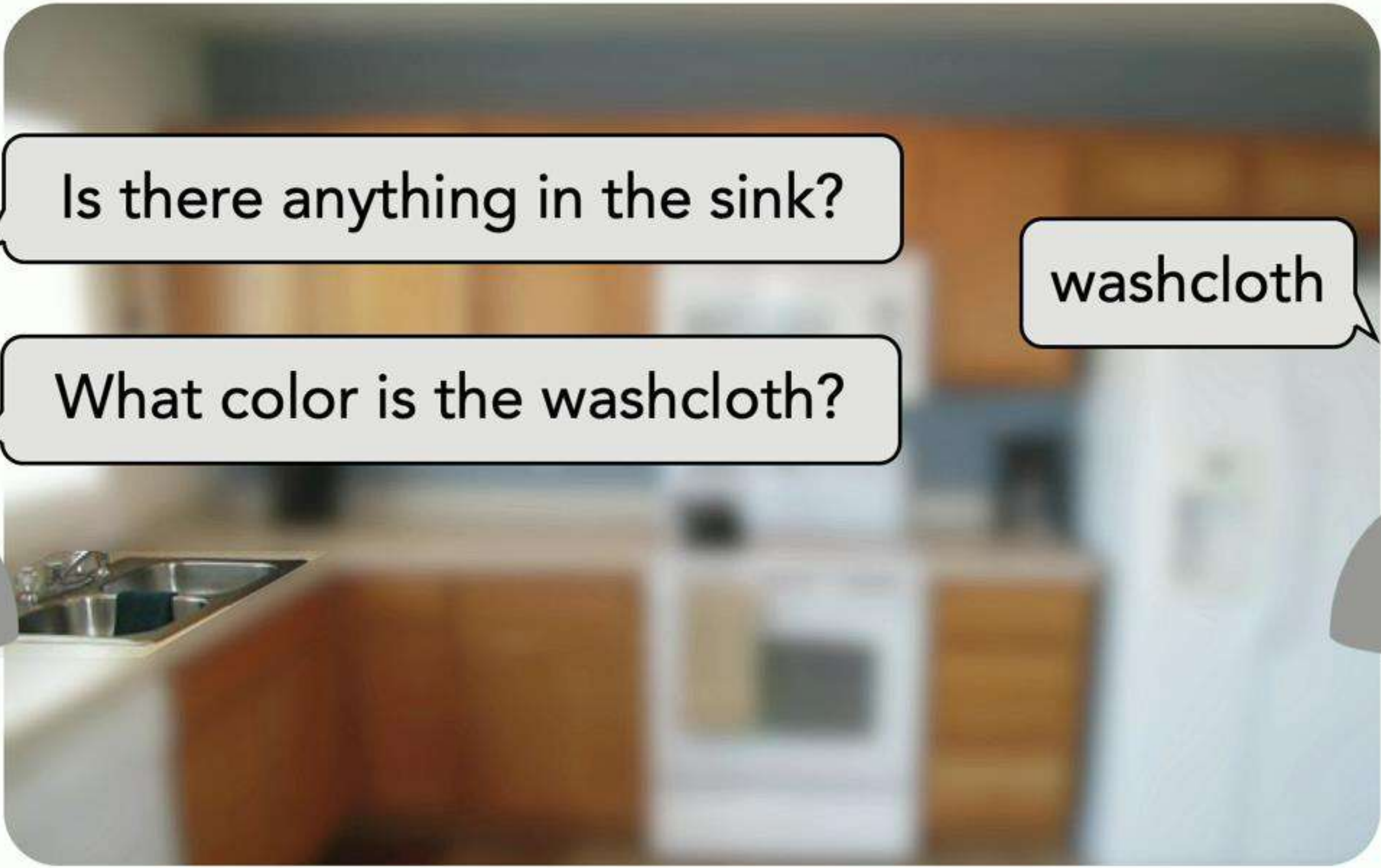
Is there anything in the sink?





Is there anything in the sink?

washcloth



Is there anything in the sink?

washcloth

What color is the washcloth?

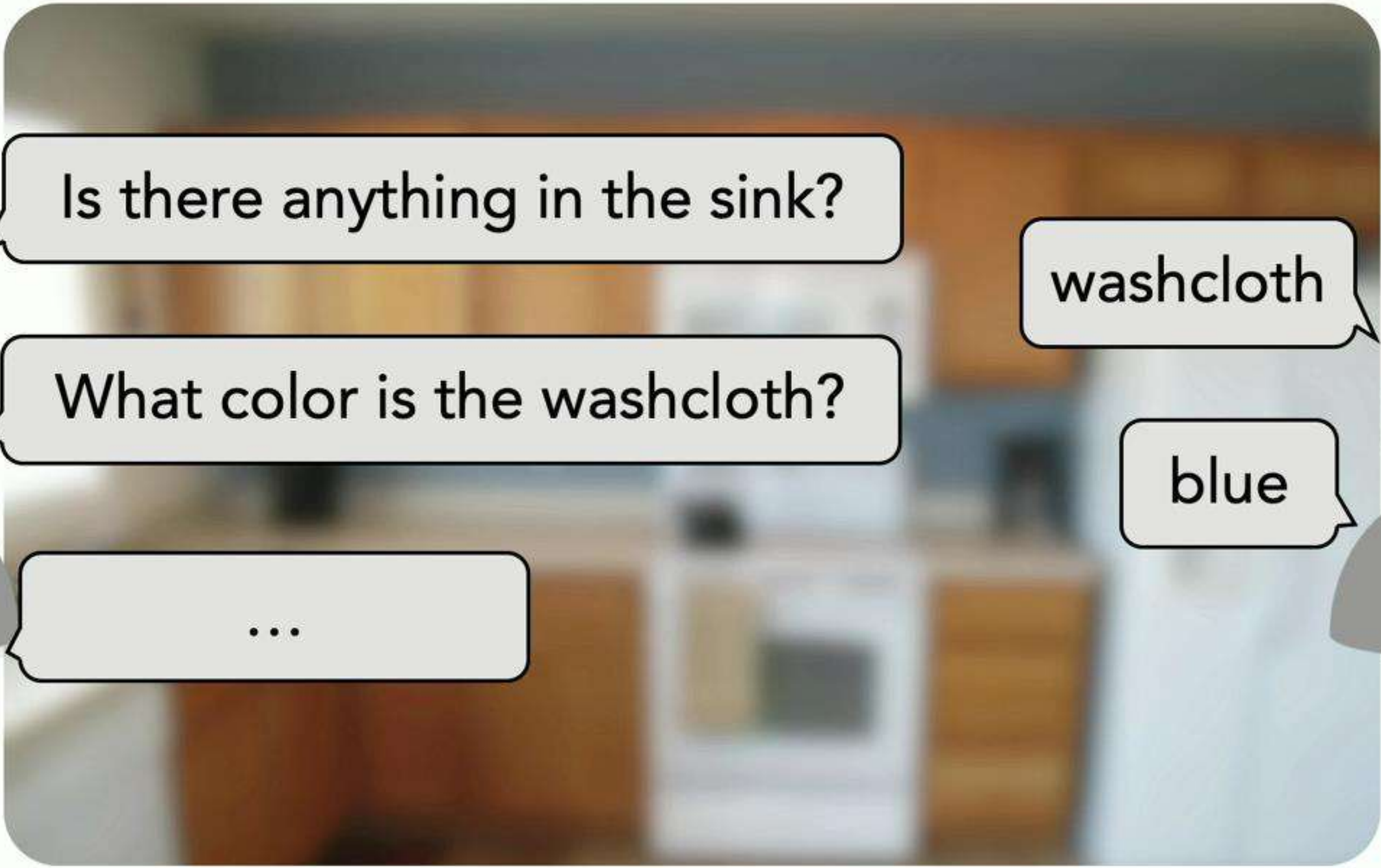


Is there anything in the sink?

washcloth

What color is the washcloth?

blue



Is there anything in the sink?

washcloth

What color is the washcloth?

blue

...



Is there anything in the sink?

washcloth

What color is the washcloth?

blue

...



Visual Dialog



The man at bat readies to swing at the pitch while the umpire looks on.



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



A large bus sitting next to a very tall building.



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

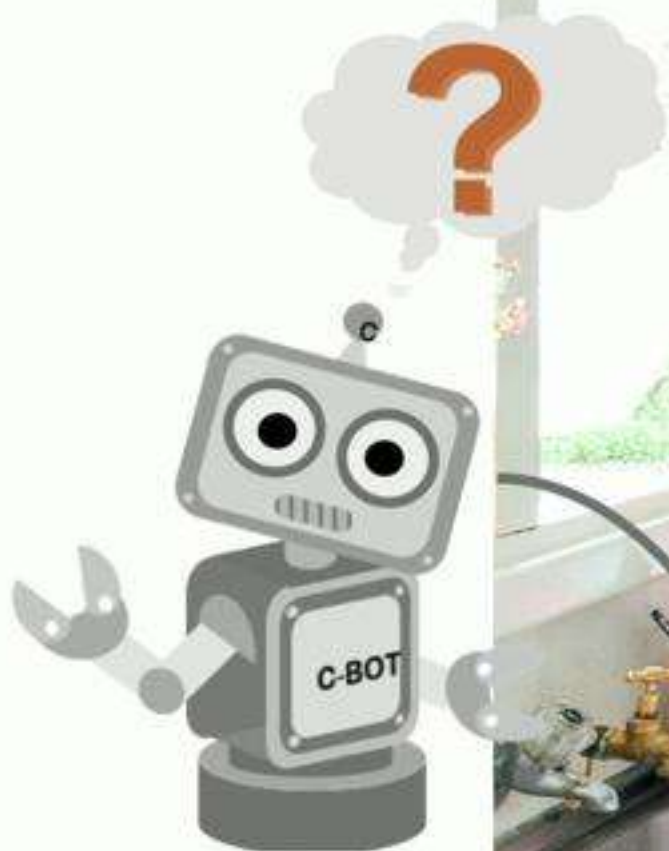


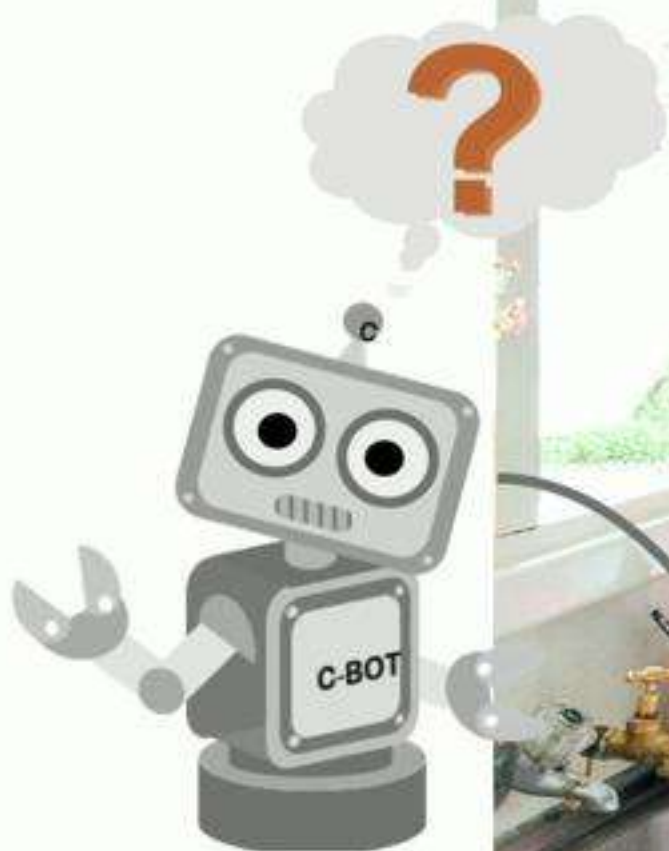
Visual Dialog

- Q: What is the gender of the one in the white shirt ?
A: She is a woman
Q: What is she doing ?
A: Playing a Wii game
Q: Is that a man to her right
A: No, it's a woman

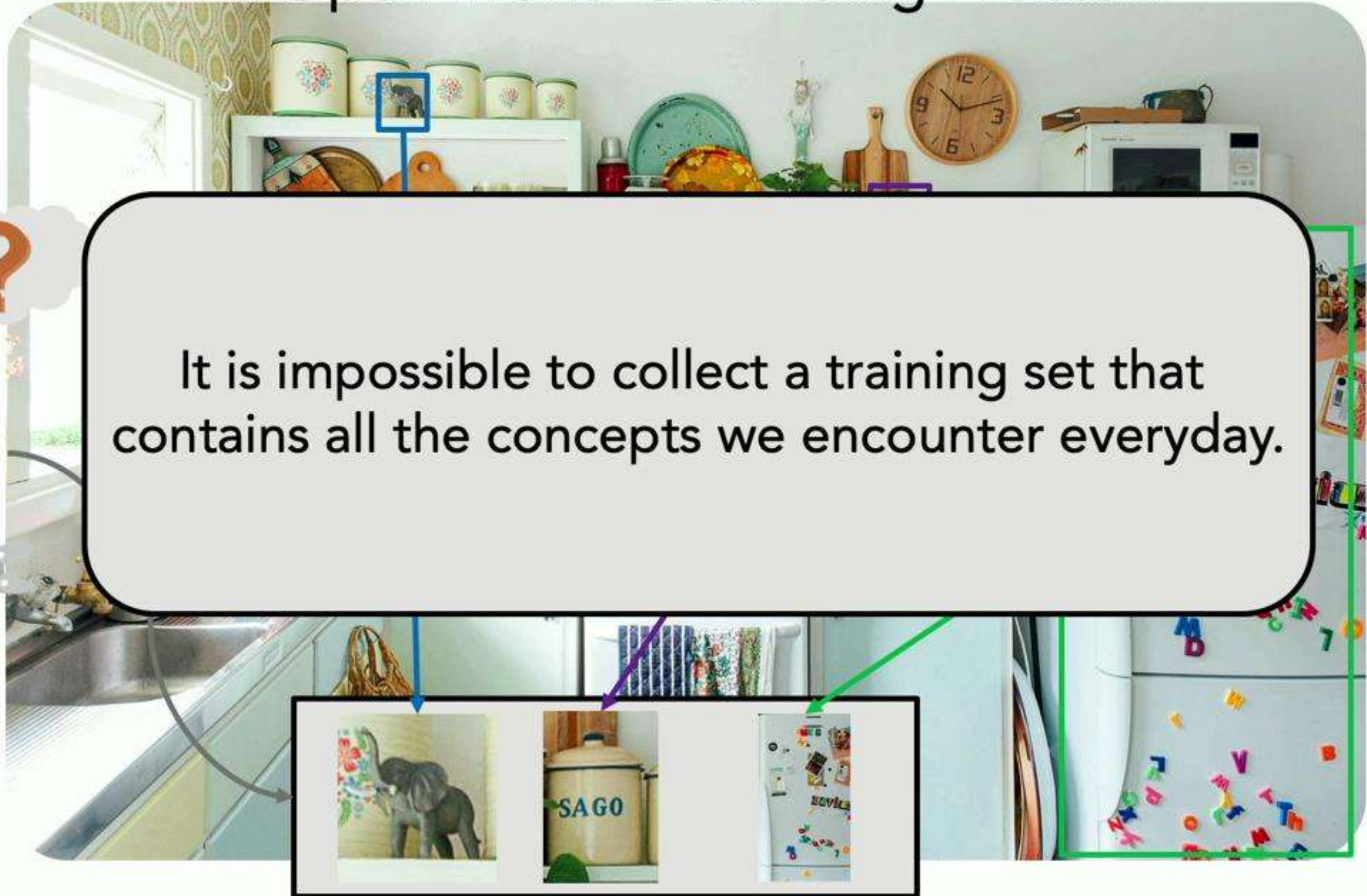


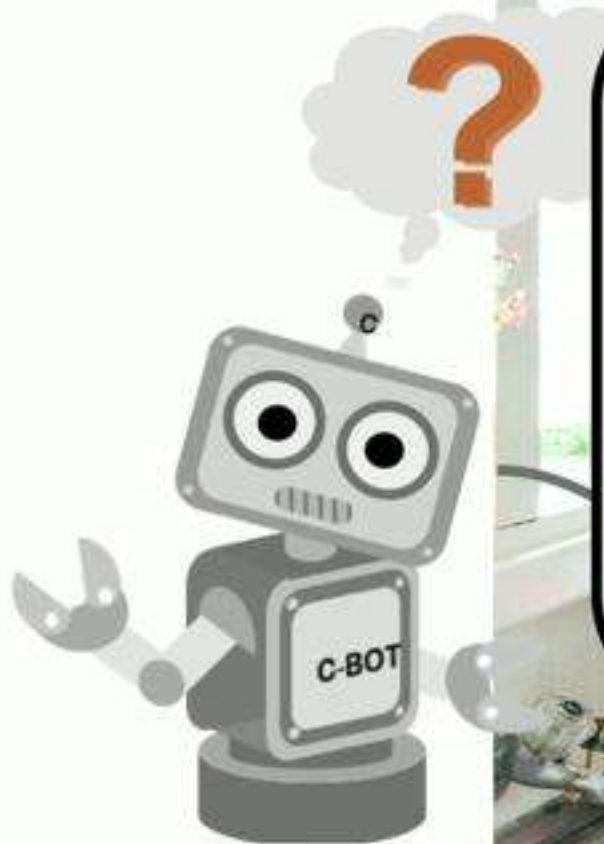






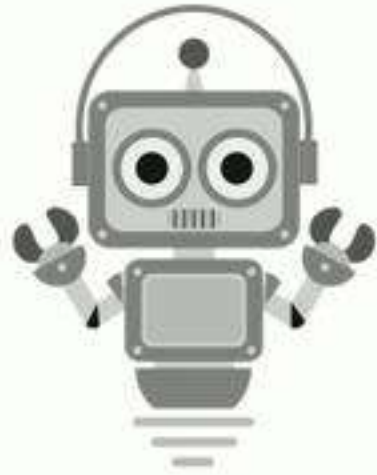
The Open-World Grounding Problem



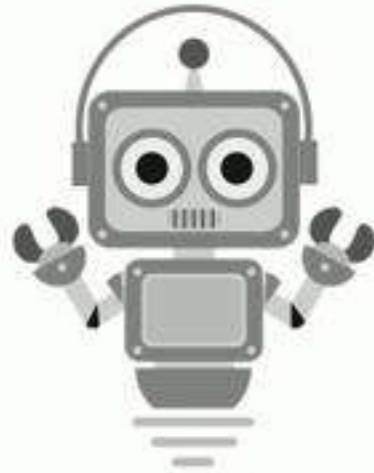


How to leverage **grounding** learned from **other sources** to improve multi-modal AI capabilities ?





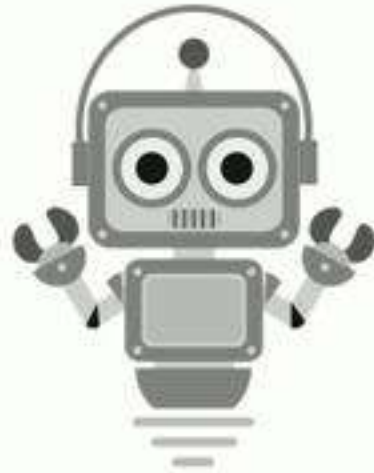
Virtual assistance



Virtual assistance



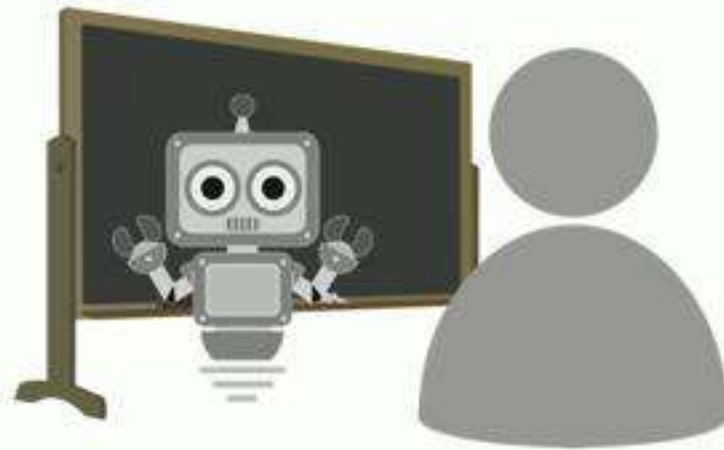
Assistive Technology



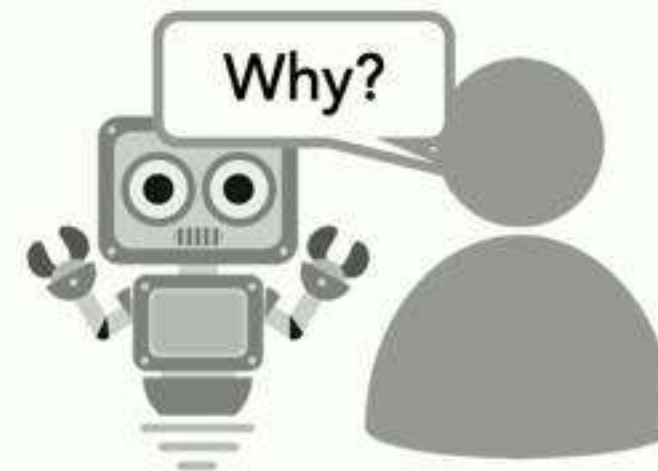
Virtual assistance



Assistive Technology



Machine Teaching

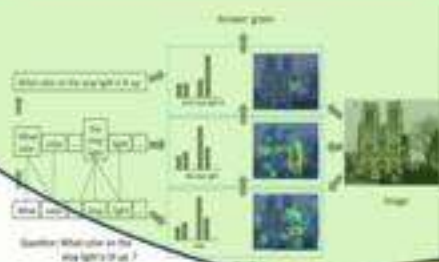


Explainable AI

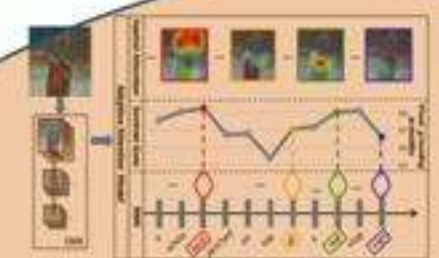


Visual Question Answering
[ICCV 2015, IJCV]

VQA

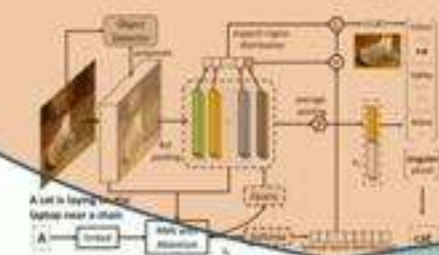


Hierarchical Question-Image
Co-Attention for VQA.
[NeurIPS 2016]

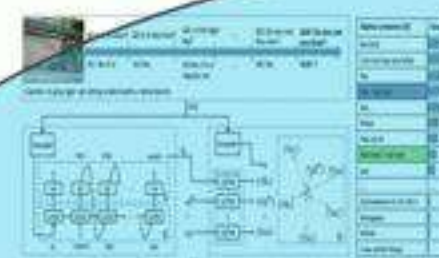


Knowing when to Look: adaptive
attention for Image captioning
[CVPR 2017]

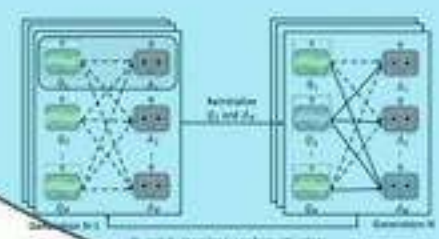
Image
Captioning



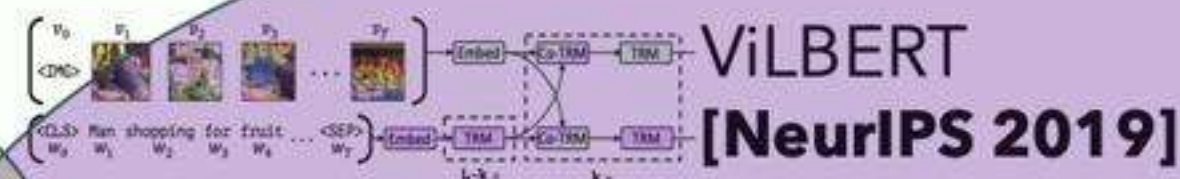
Neural Baby Talk
[CVPR 2018]



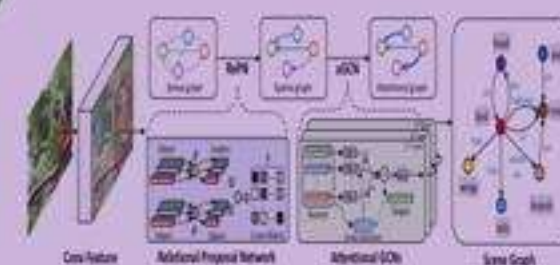
Discriminative Learning for visual
dialog
[NeurIPS 2018] (Visual) Dialog



Emergence of Compositional Language
with Deep Generational Transmission.
[In Submission]

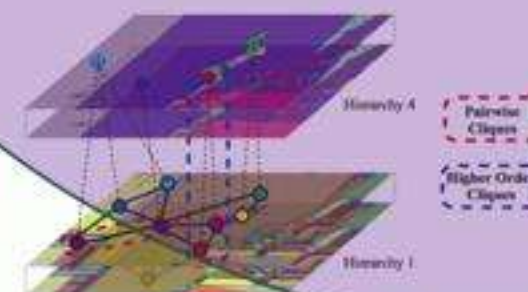


ViLBERT
[NeurIPS 2019]

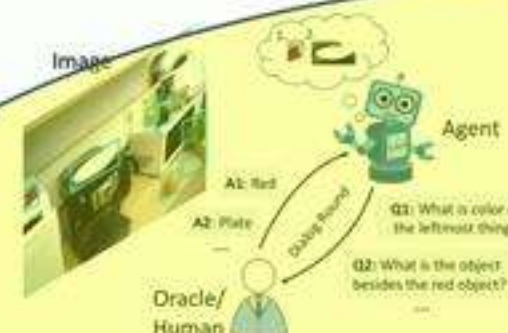


Graph RCNN
[ECCV 2018]

Image/Video
understanding



Human Action
Segmentation
[CVPR 2015]



Visual Curiosity
[CORL 2018]

Embodied Agent

Exit the bedroom and go towards the table. Go to the stairs on the left of the couch. Wait on the third step.



Self-Monitoring Navigation
Agent via Auxiliary Progress
Estimation.
[ICLR 2019]

This Talk



WALL-E

Neural Baby Talk

CVPR 2018



Mr. Potato Head

ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks

NeurIPS
2019



Jarvis

Recent & Future Work

Image Captioning

The fundamental capability to describe what is seen.

Input:



Image Captioning

The fundamental capability to describe what is seen.

Input:



Desired
output:

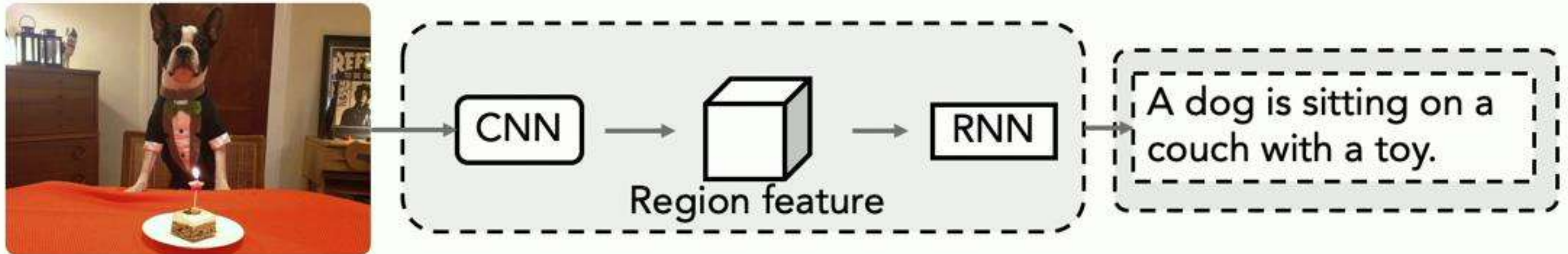
*A table that has a mug and
a plate with food on it.*



*A dog wearing goggles and
leathers sitting on a motorcycle.*

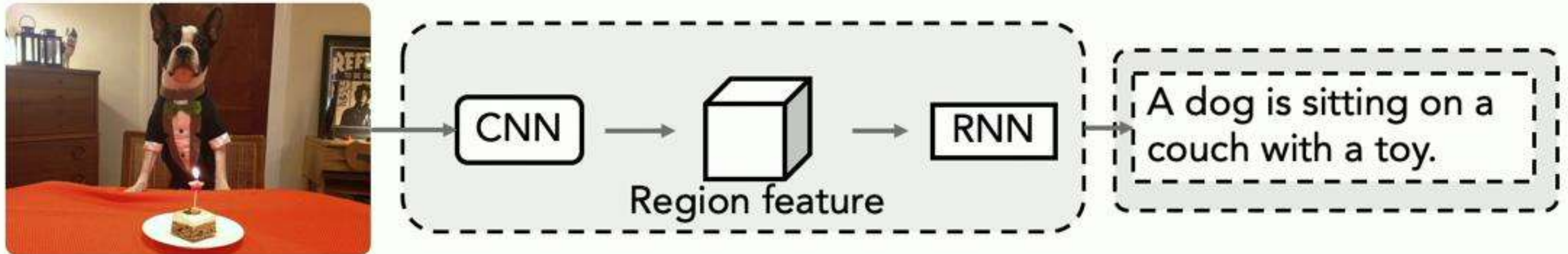
Related Work

SOTA neural image captioning system.



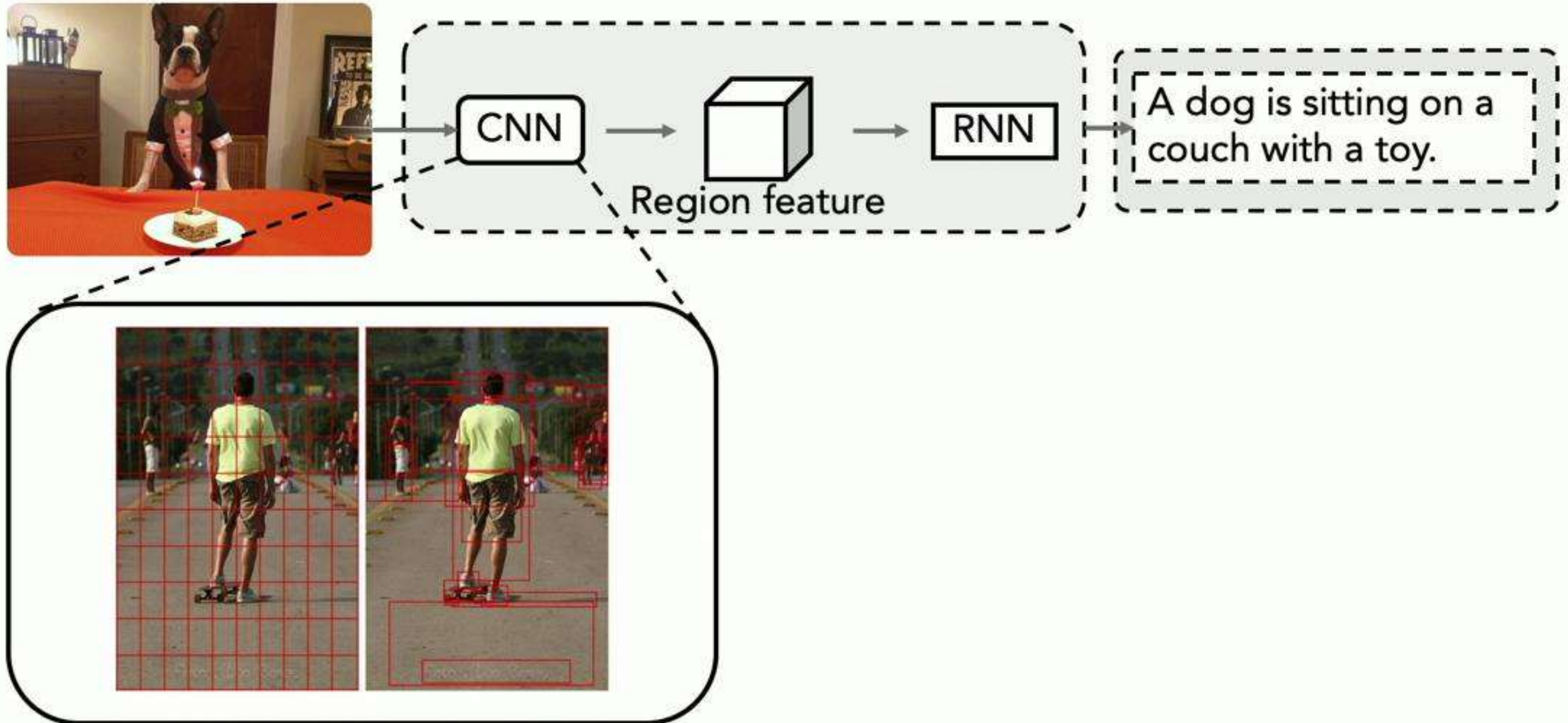
Related Work

Bottom up and Top down attention for image captioning



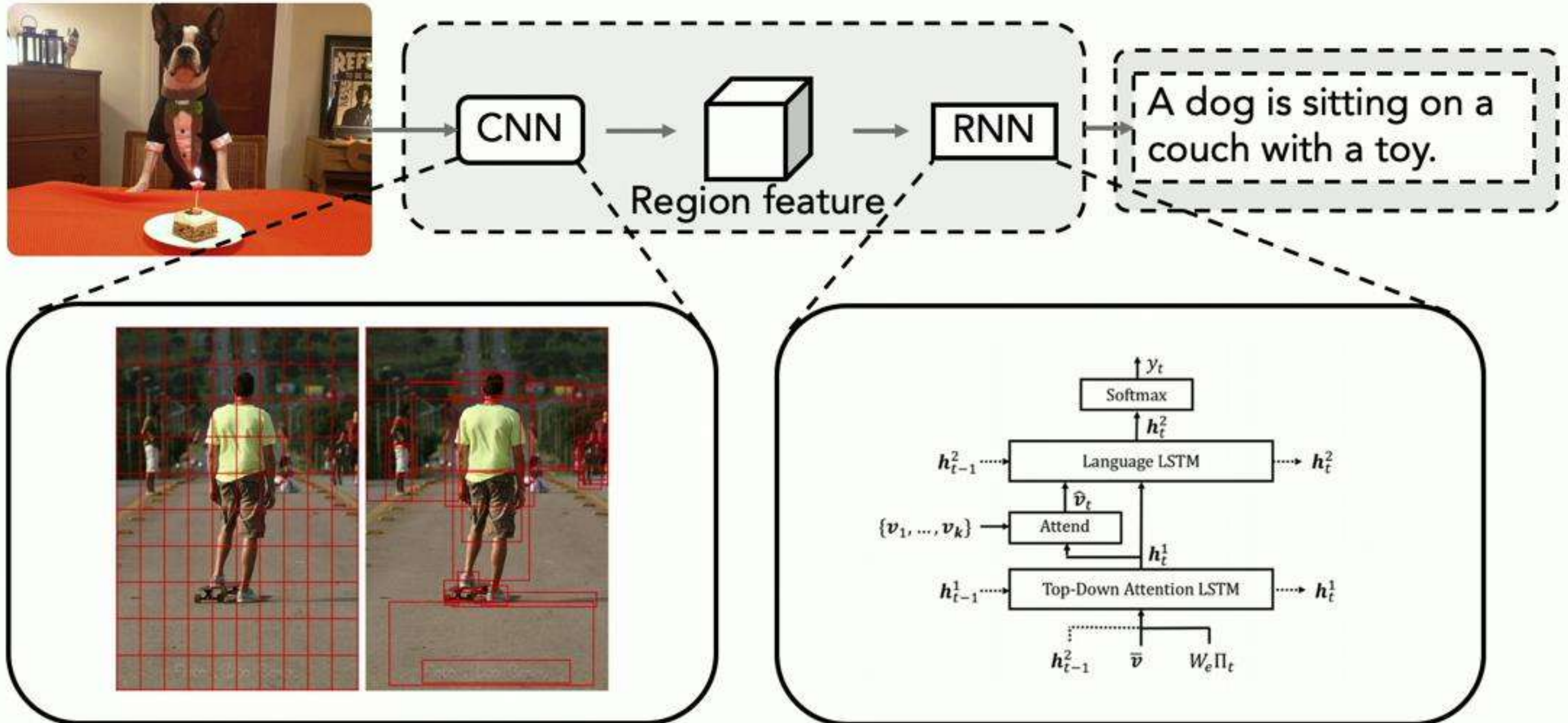
Related Work

Bottom up and Top down attention for image captioning



Related Work

Bottom up and Top down attention for image captioning



Related Work

SOTA neural image captioning system.



Two elephants and a baby elephant walking together.

Related Work

SOTA neural image captioning system.



Two elephants and a baby elephant walking together.



A cat is standing on a sign that says "UNK".

Related Work

SOTA neural image captioning system.



Two elephants and a baby elephant walking together.



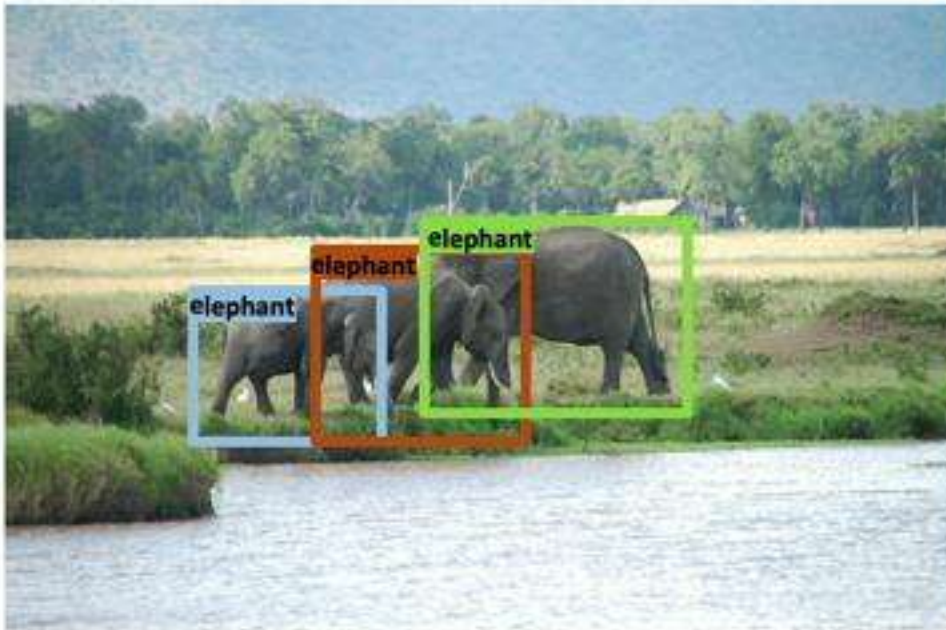
A cat is standing on a sign that says "UNK".



A man standing on a beach holding a surfboard.

Related Work

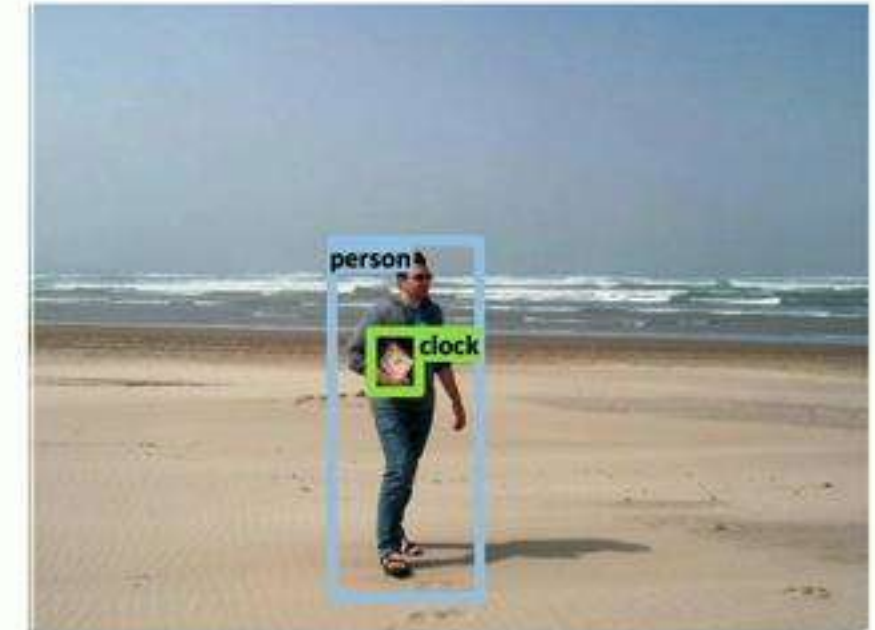
Object detection and OCR output



Two elephants and a baby elephant walking together.



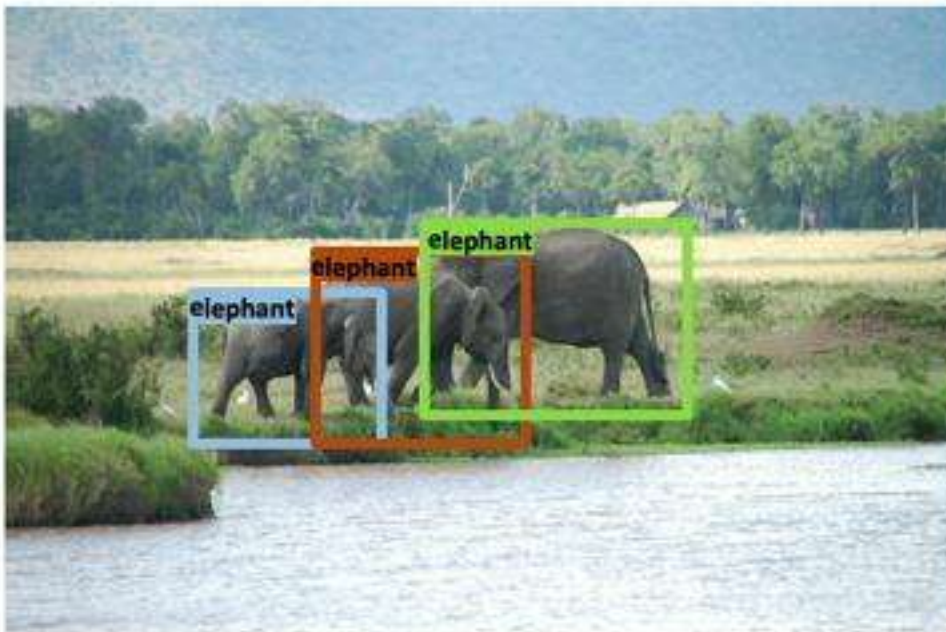
A cat is standing on a sign that says "UNK".



A man standing on a beach holding a surfboard.

Related Work

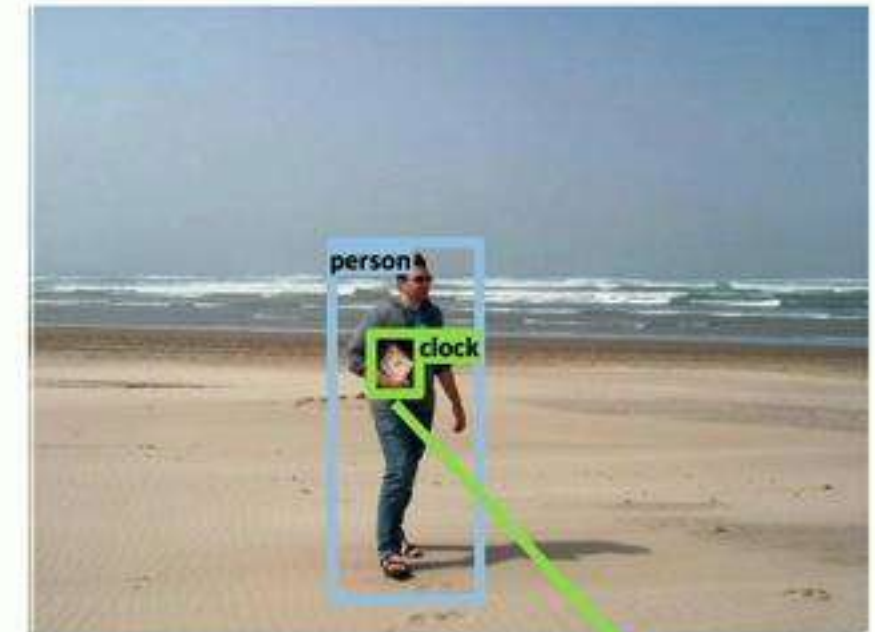
Learn from object detection and OCR.



Two elephants and a baby elephant walking together.



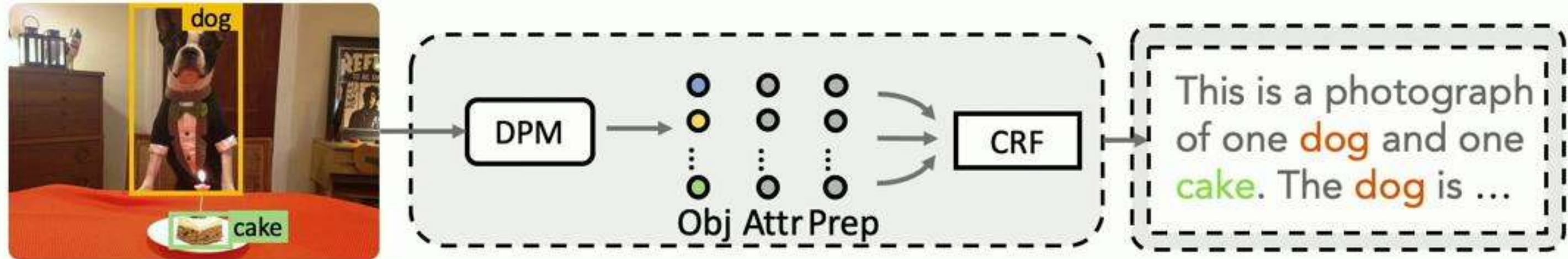
A cat is standing on a sign that says "~~UNK~~" - "Abundzu".



A man standing on a beach holding a ~~surfboard~~ - clock.

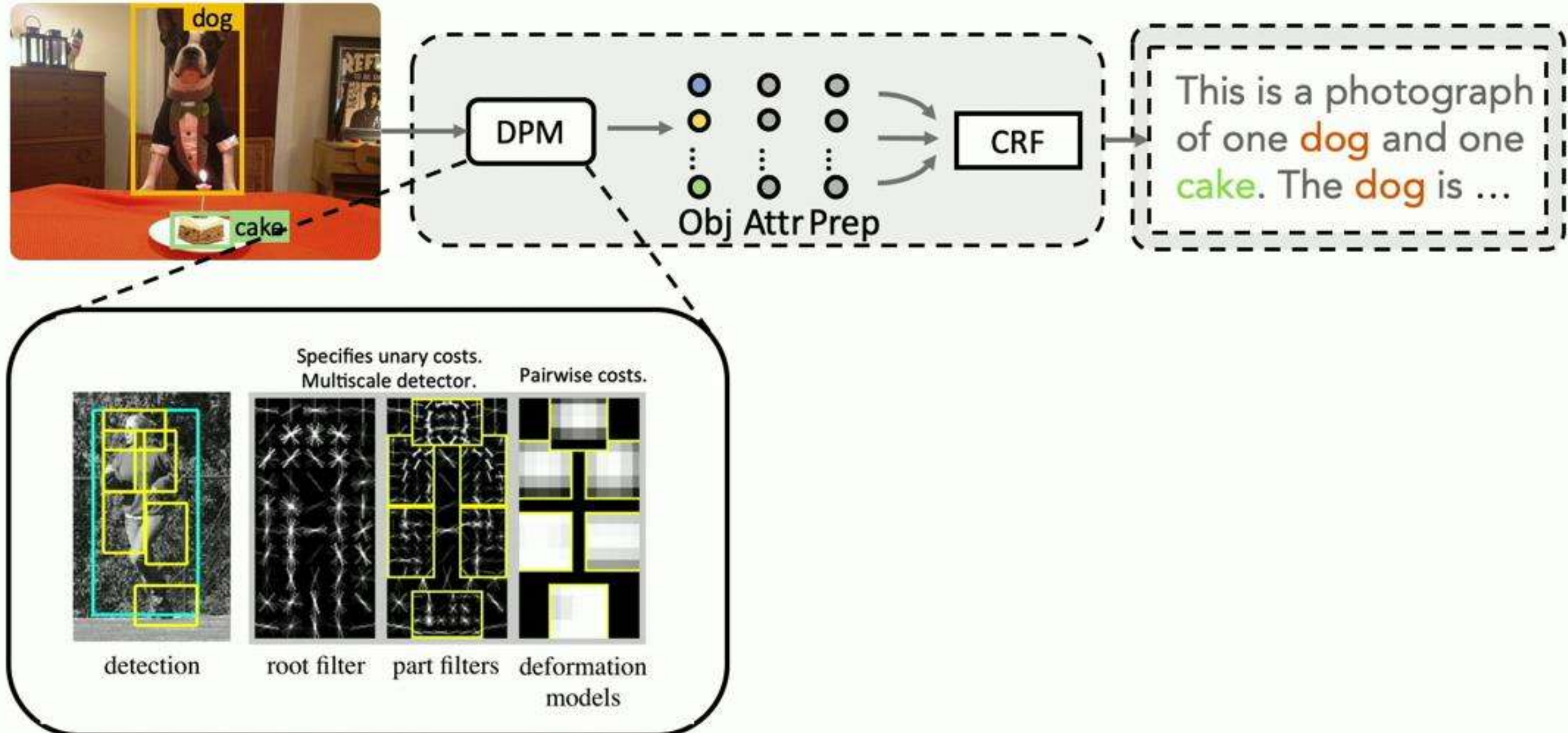
Related Work

Image captioning system before deep learning “revolution”.



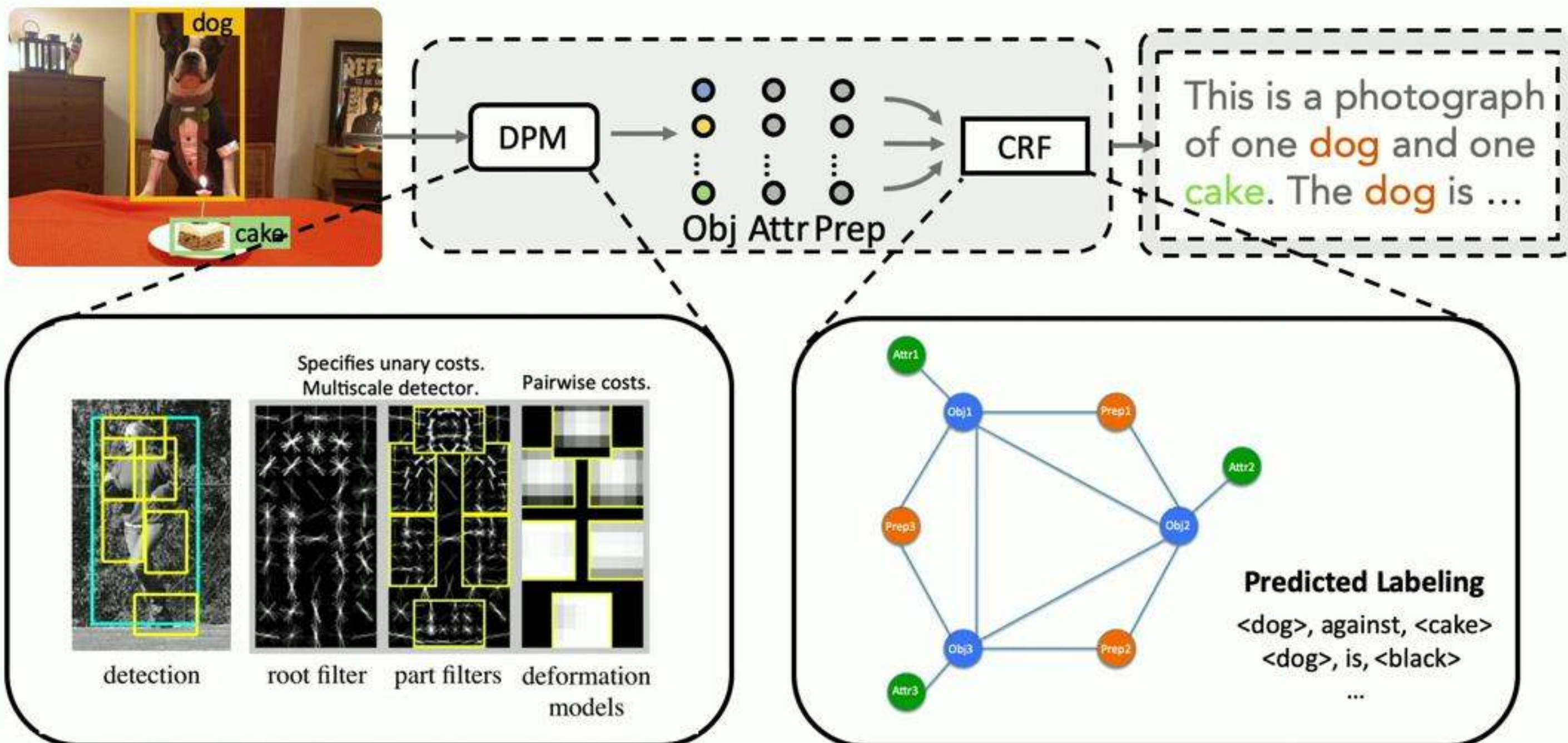
Related Work

Image captioning system before deep learning "revolution".



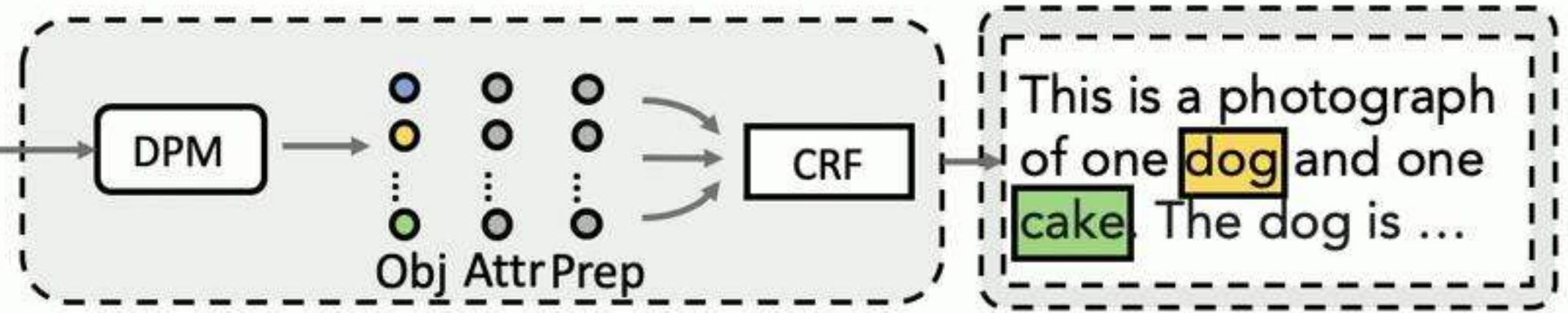
Related Work

Image captioning system before deep learning "revolution".



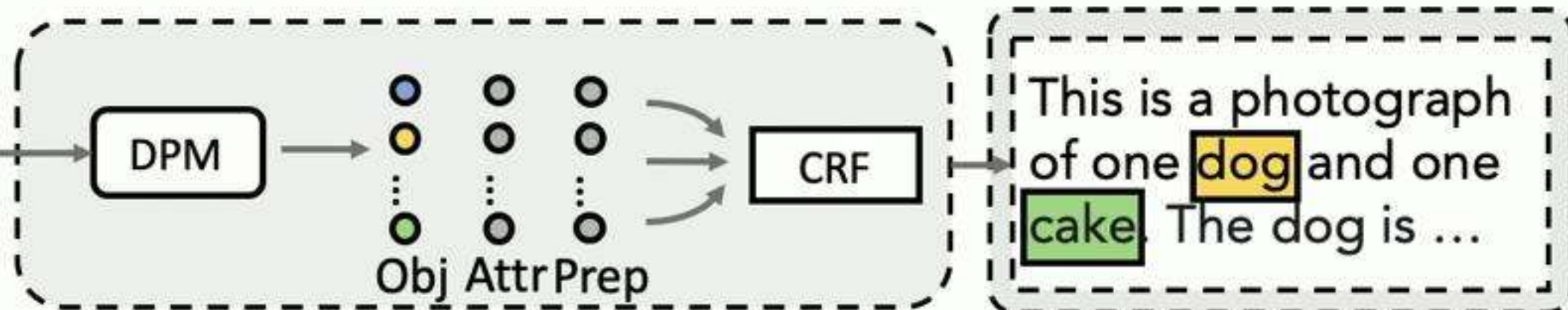
Motivation

More
Grounded

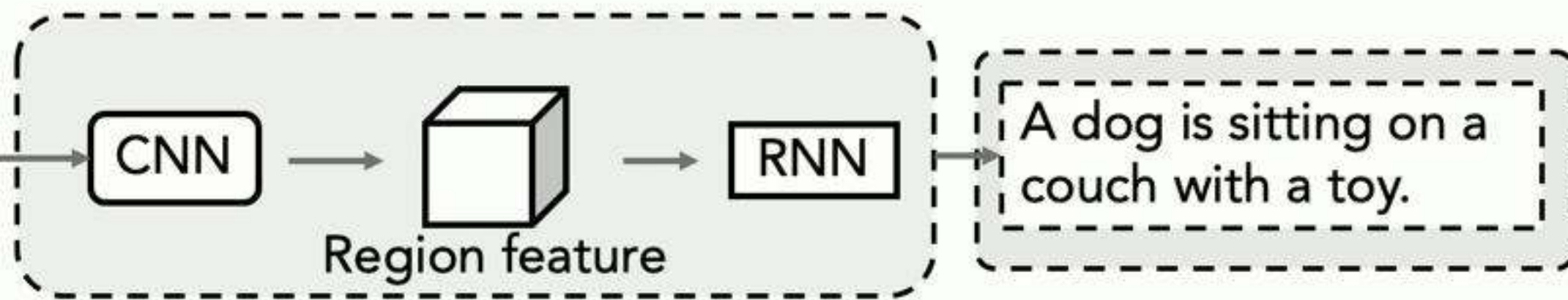


Motivation

More
Grounded

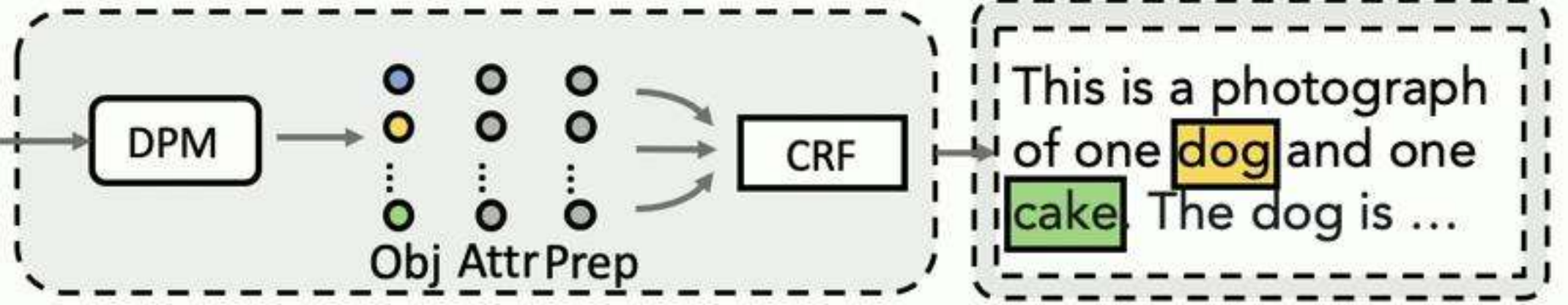


More
Natural



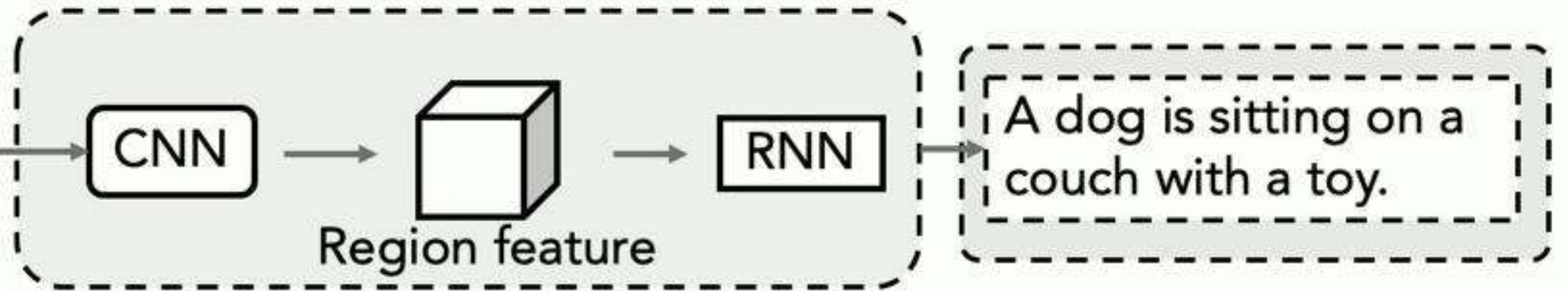
Motivation

More
Grounded



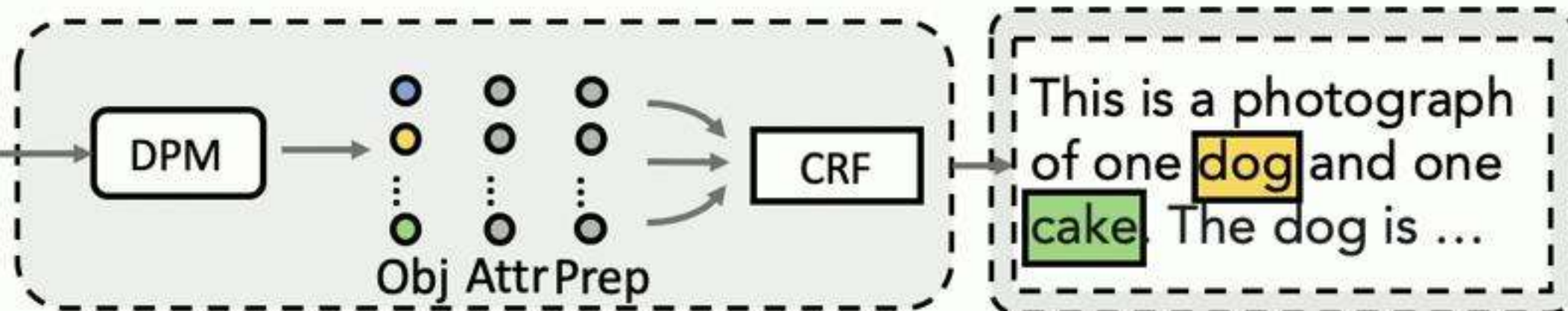
- Associate the named concept to the pixels / bounding box detections in the image.

More
Natural



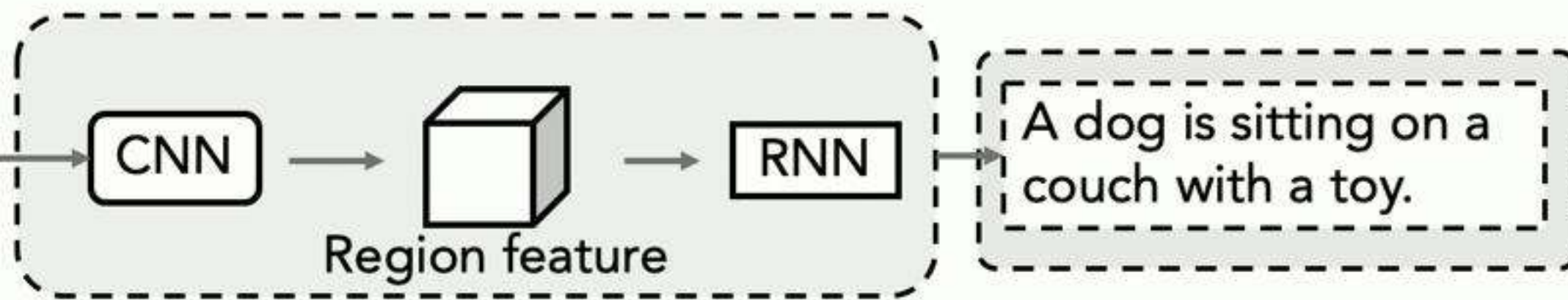
Motivation

More
Grounded

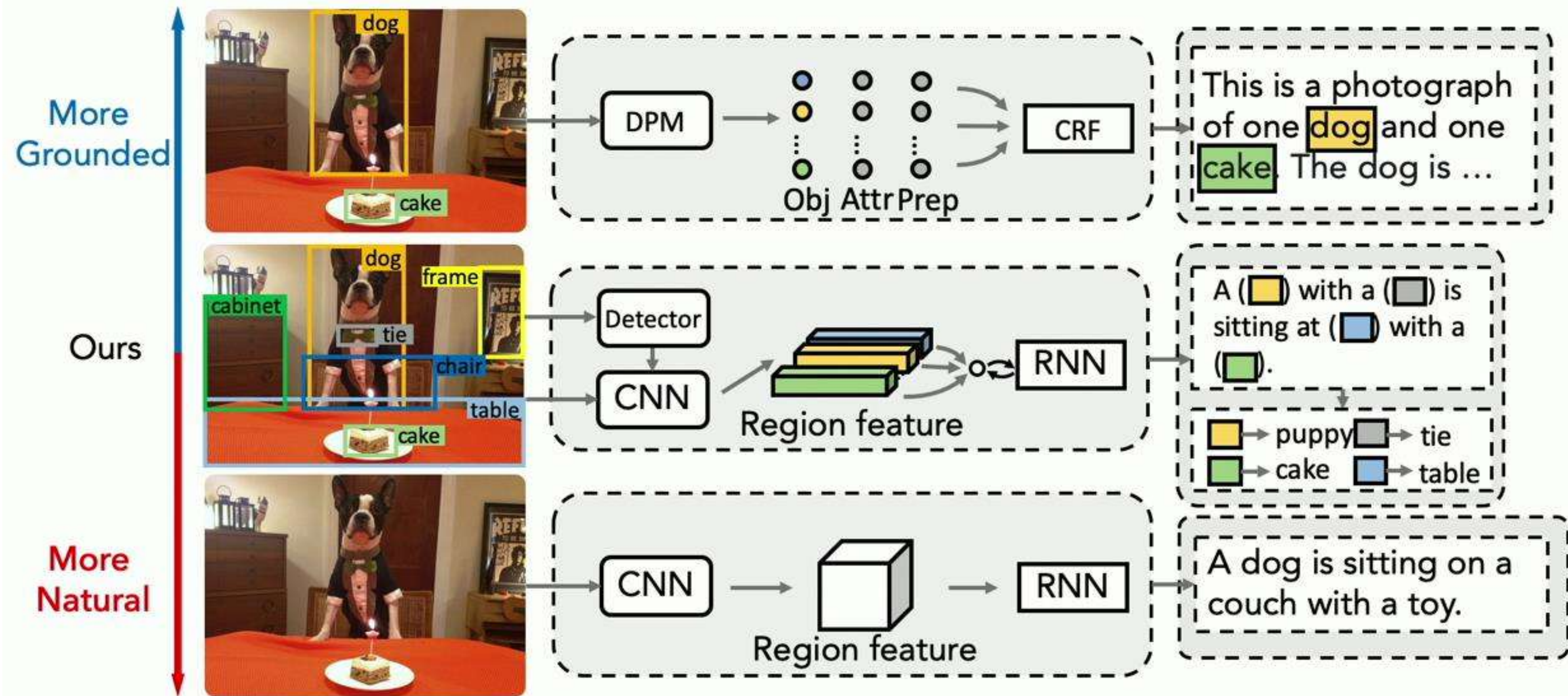


- Associate the named concept to the pixels / bounding box detections in the image.
- No human generated templated, language needs to be natural.

More
Natural

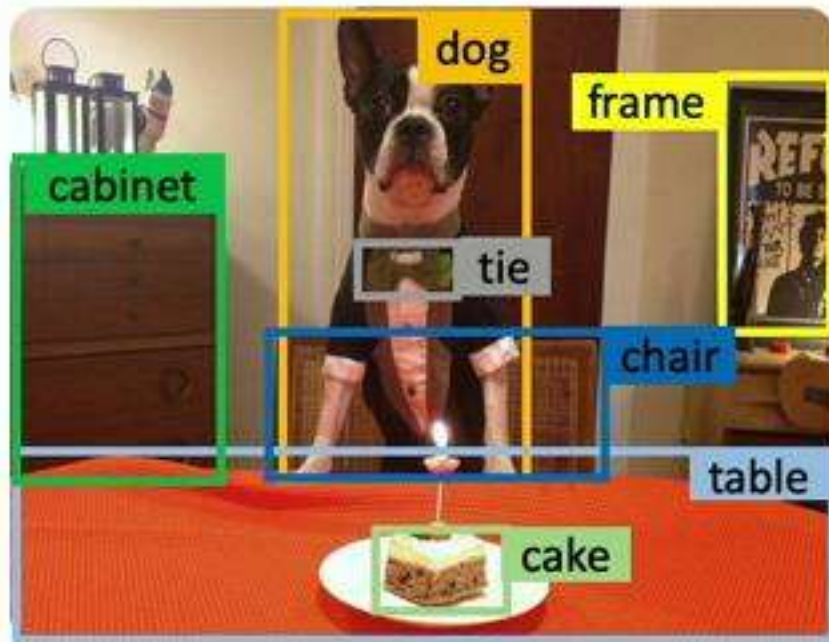


Motivation



Framework

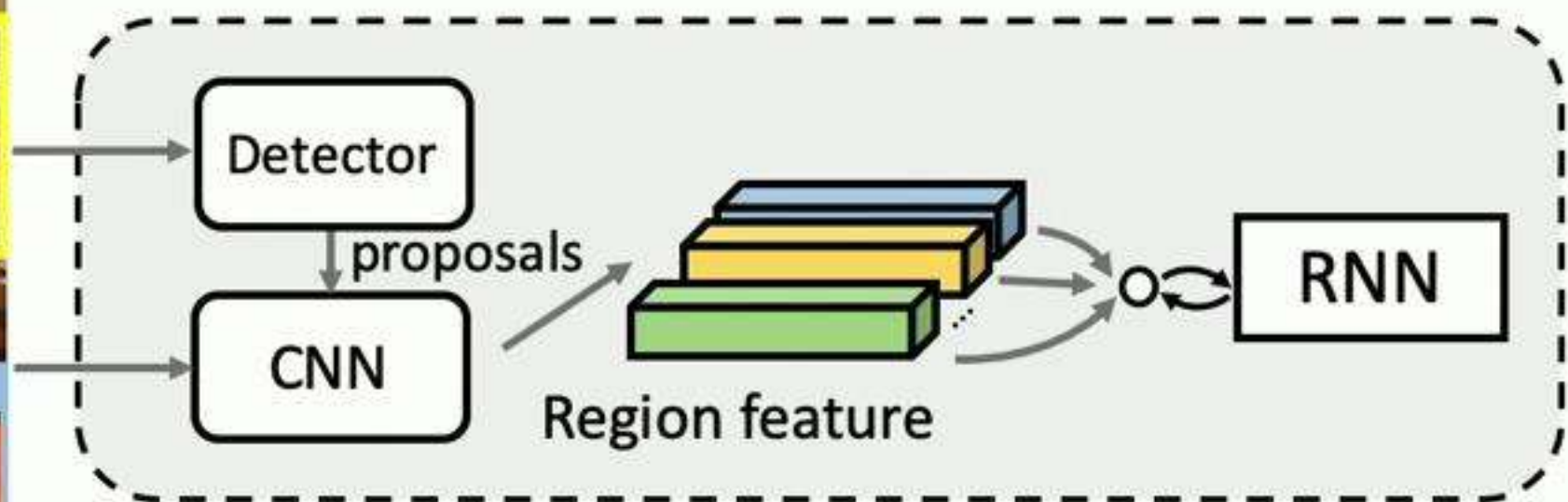
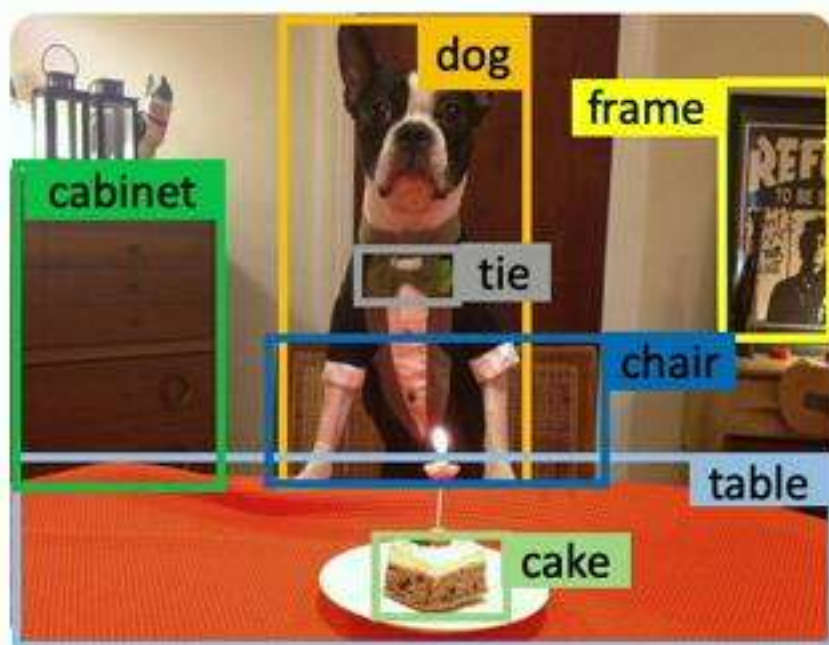
Neural Baby
Talk



Detector

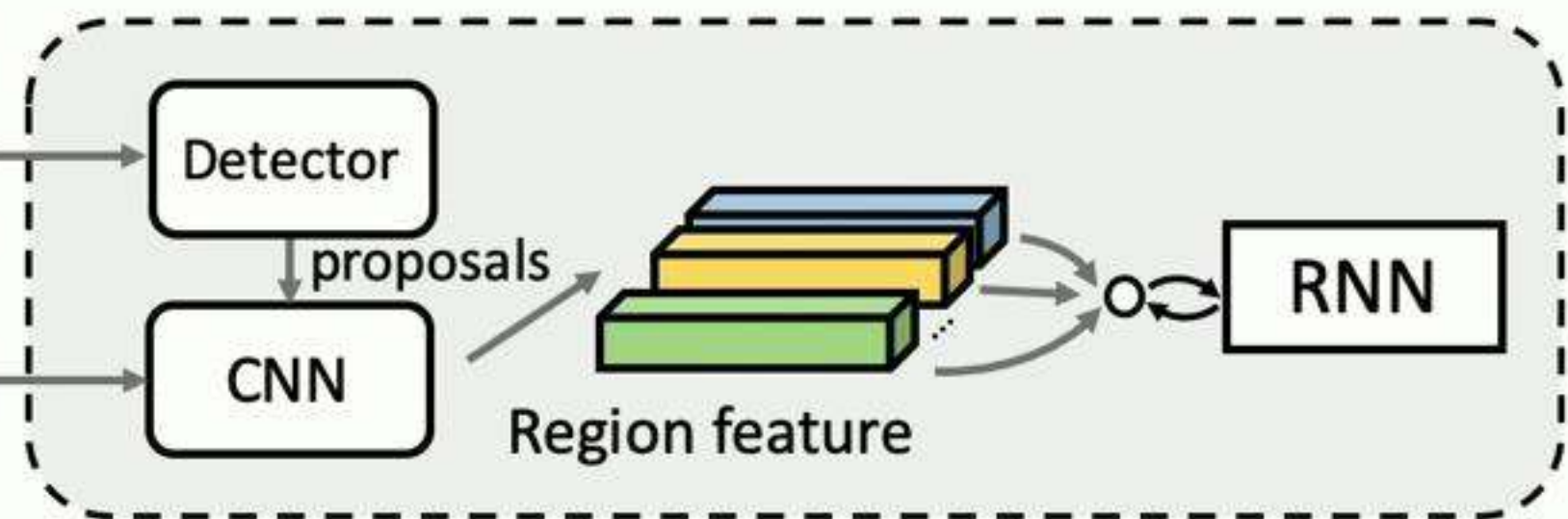
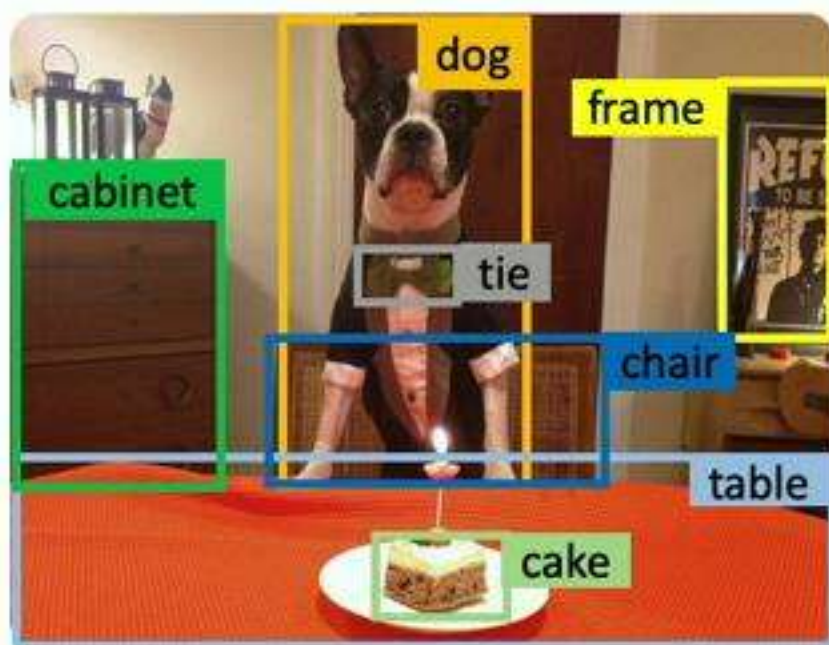
Framework

Neural Baby
Talk



Framework

Neural Baby
Talk

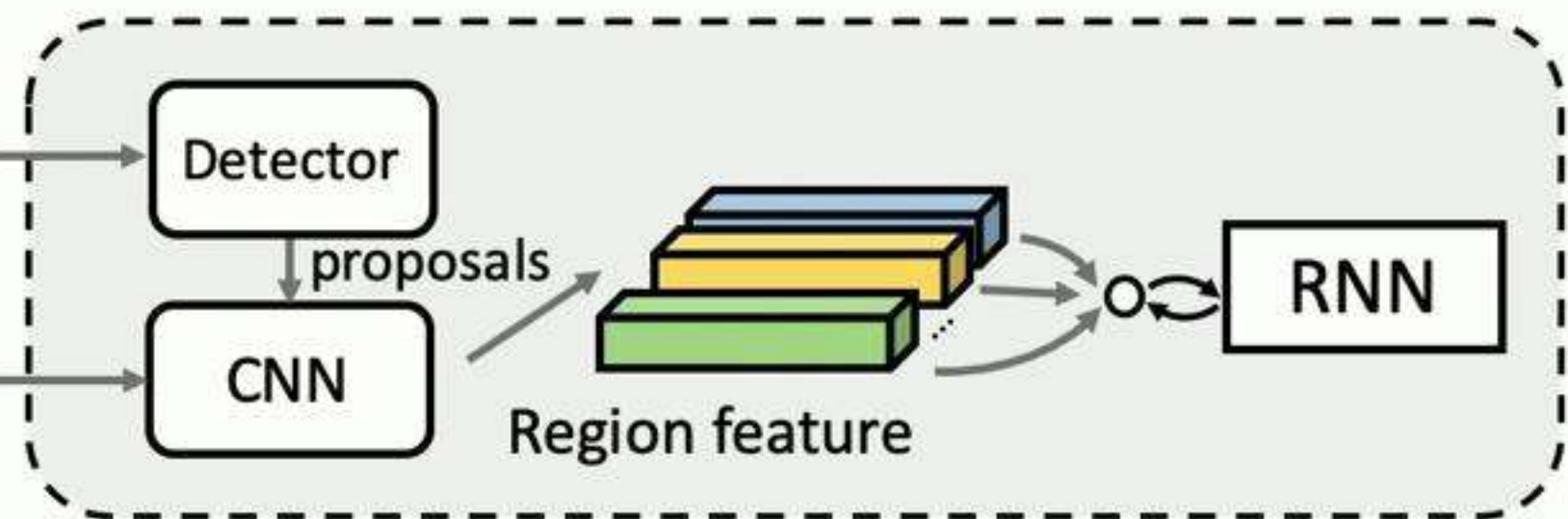
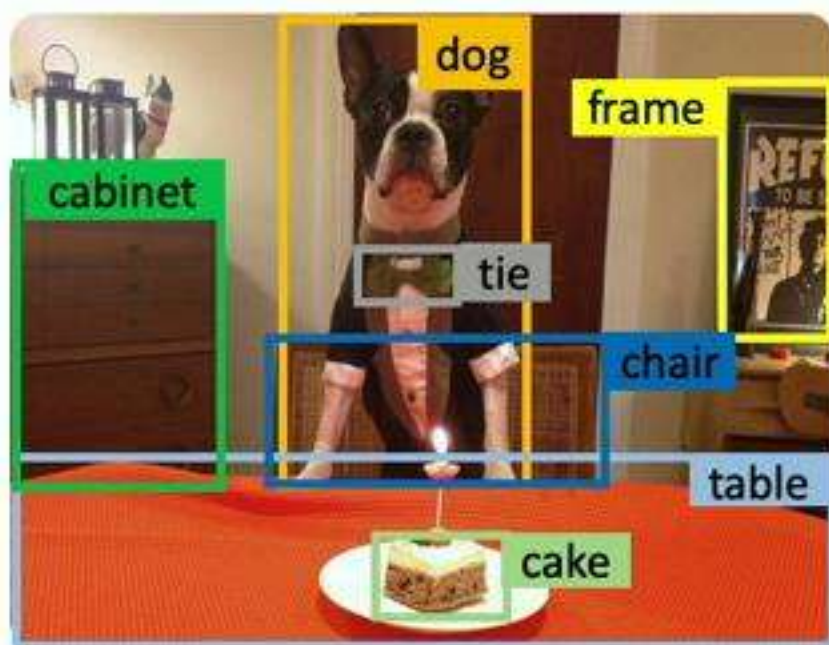


- Slotted caption template generation

A () with a () is sitting at () with a ().

Framework

Neural Baby Talk



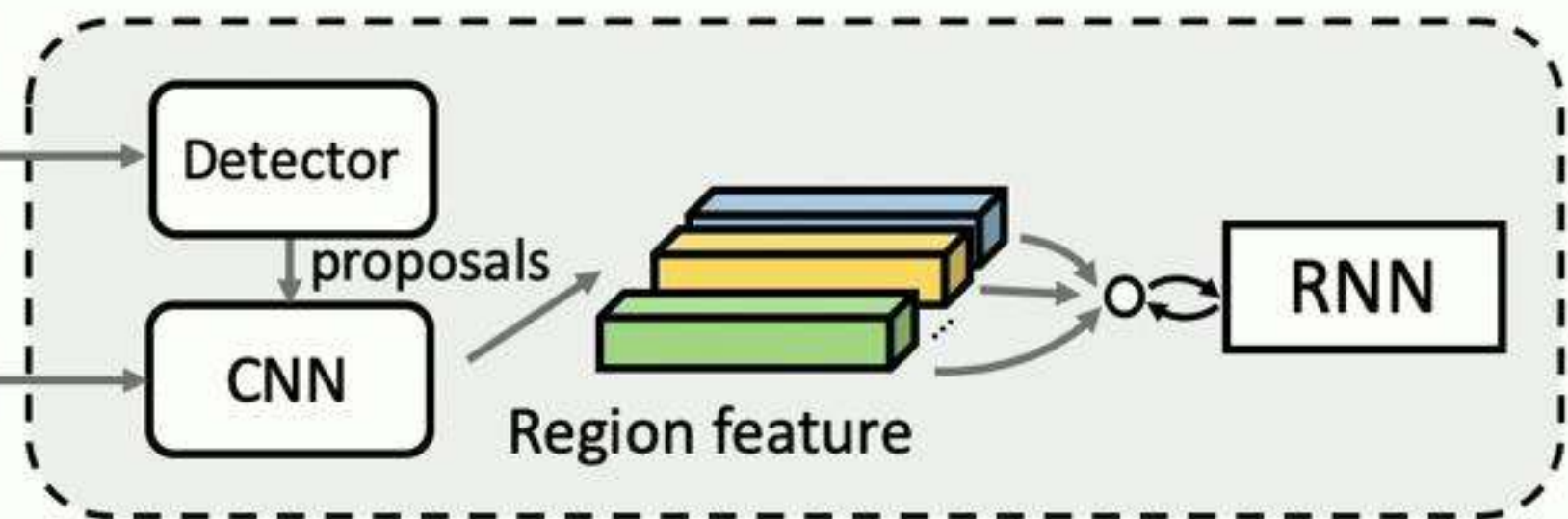
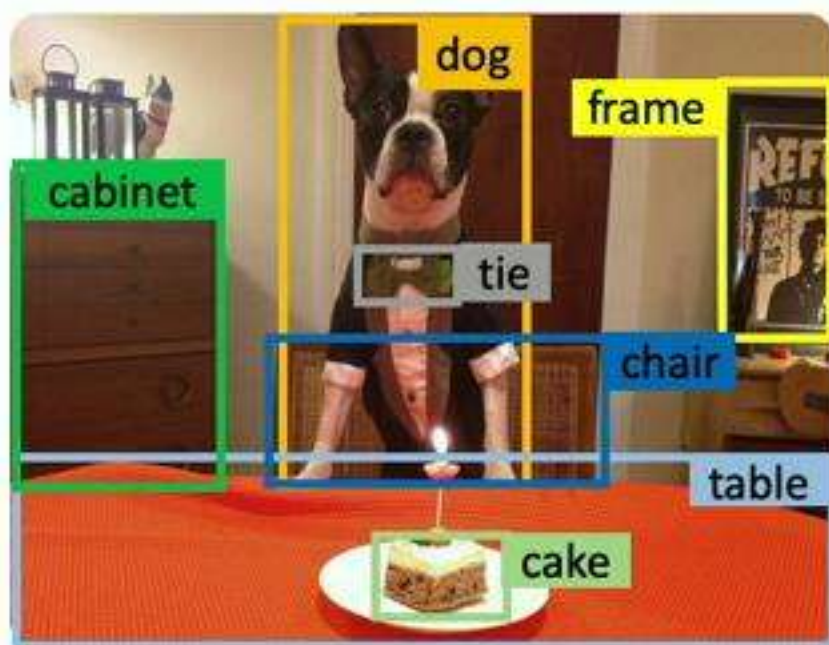
- Slotted caption template generation
- Filling in the slots

A () with a () is sitting at () with a ().

() → puppy () → tie
() → cake () → table

Framework

Neural Baby
Talk



- Slotted caption template generation
- Filling in the slots
- Final caption

A () with a () is sitting at () with a ().

→ puppy → tie
 → cake → table

A puppy with a tie is sitting at table with a cake.

Model

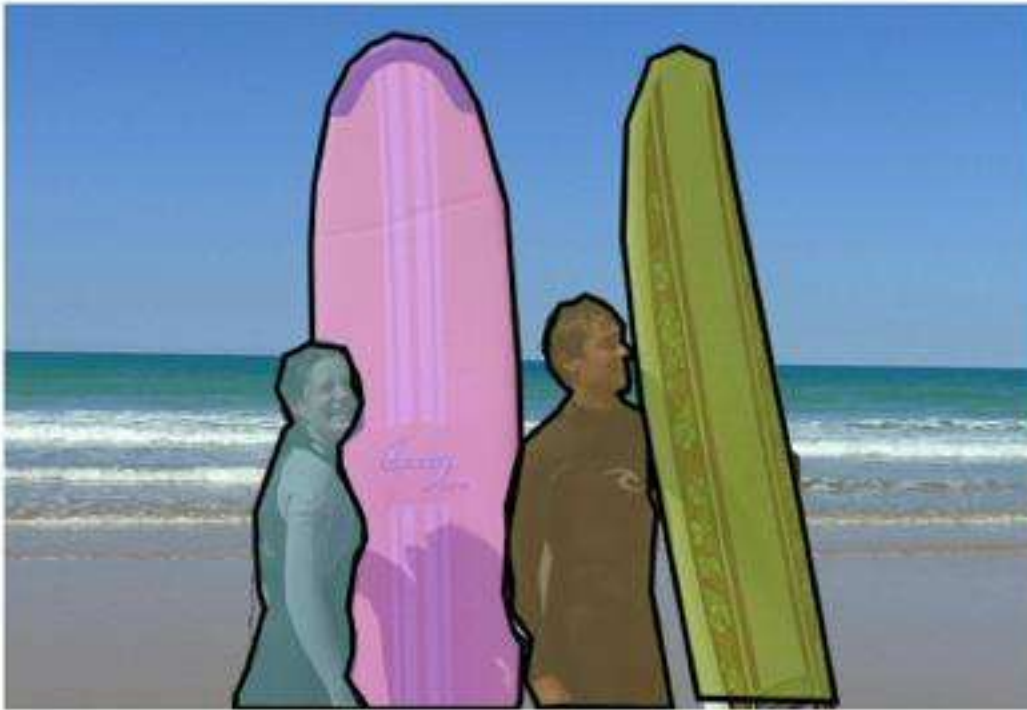
How to construct the caption template?



A young woman standing next to a man holding a surfboard.

Model

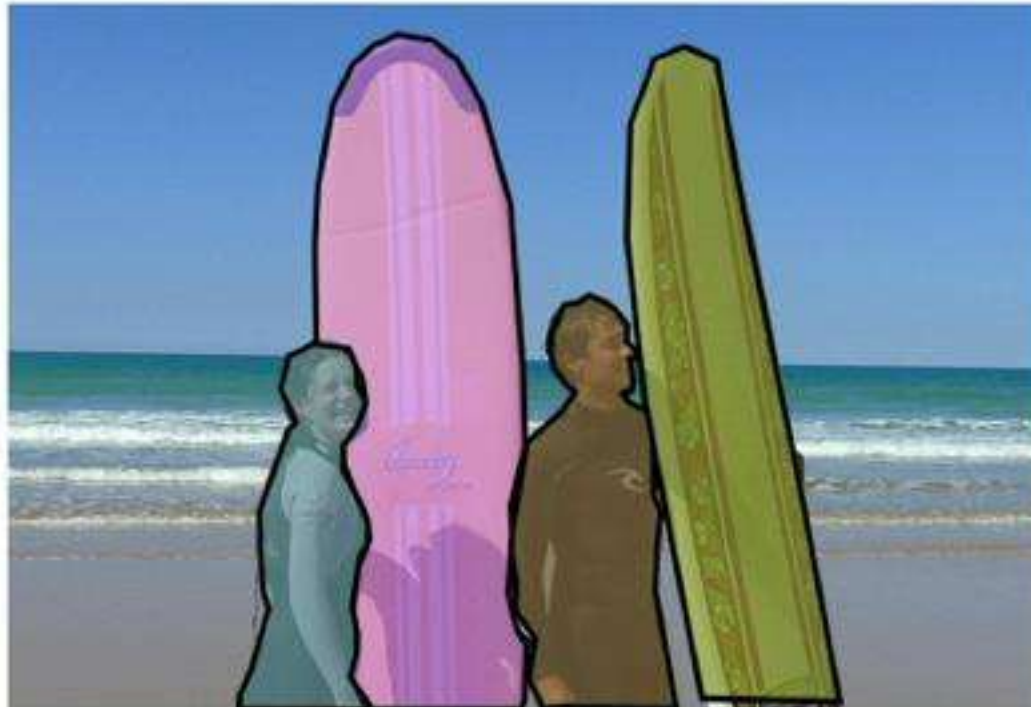
How to construct the caption template?



A young woman standing next to a man holding a surfboard.

Object label: `<person>`, `<person>`, `<surfboard>`, `<surfboard>`

How to construct the caption template?



A young woman standing next to a man holding a surfboard.

Object label: `<person>`, `<person>`, `<surfboard>`, `<surfboard>`

A young `woman` standing next to a `man` holding a `surfboard`.

Model

How to construct the caption template?



noun

A young woman standing next to a man holding a surfboard.

Object label: `<person>`, `<person>`, `<surfboard>`, `<surfboard>`

A young woman standing next to a man holding a surfboard.

A young standing next to a holding a .

Model

How to construct the caption template?



A young woman standing next to a man holding a surfboard.

Object label: `<person>`, `<person>`, `<surfboard>`, `<surfboard>`

A young `woman` standing next to a `man` holding a `surfboard`.

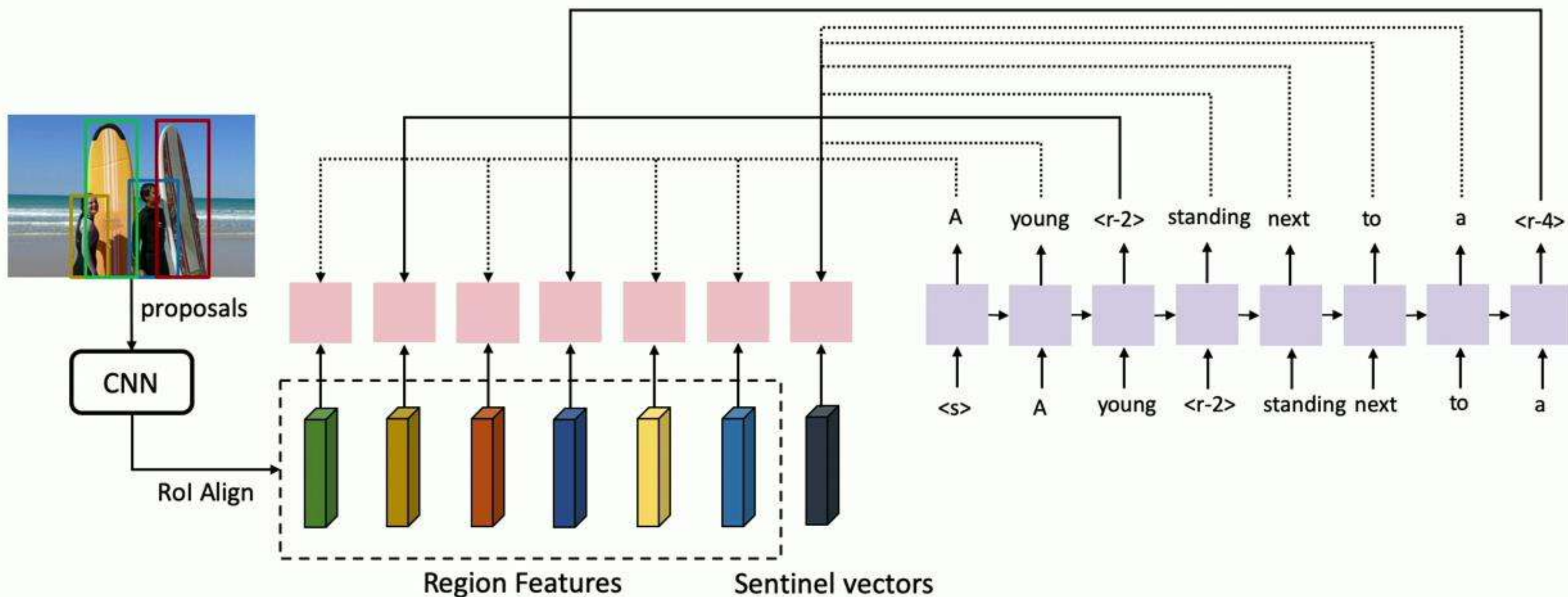
A young [redacted] standing next to a [redacted] holding a [redacted].

A [redacted] [redacted] [redacted] next to a [redacted] [redacted] a [redacted].

[redacted] noun [redacted] attribute [redacted] action

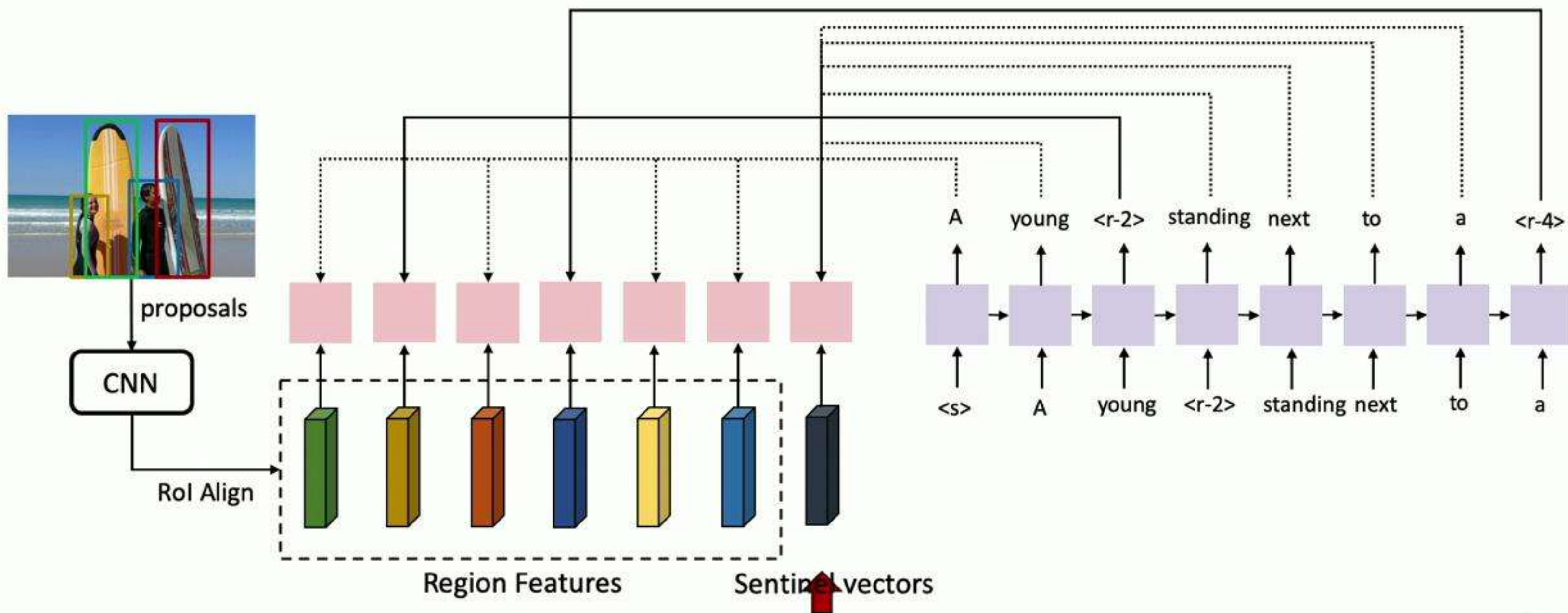
Model

Use pointer networks [Vinyals et.al 2015] with visual sentinel [Lu et.al 2017] to generate caption templates.



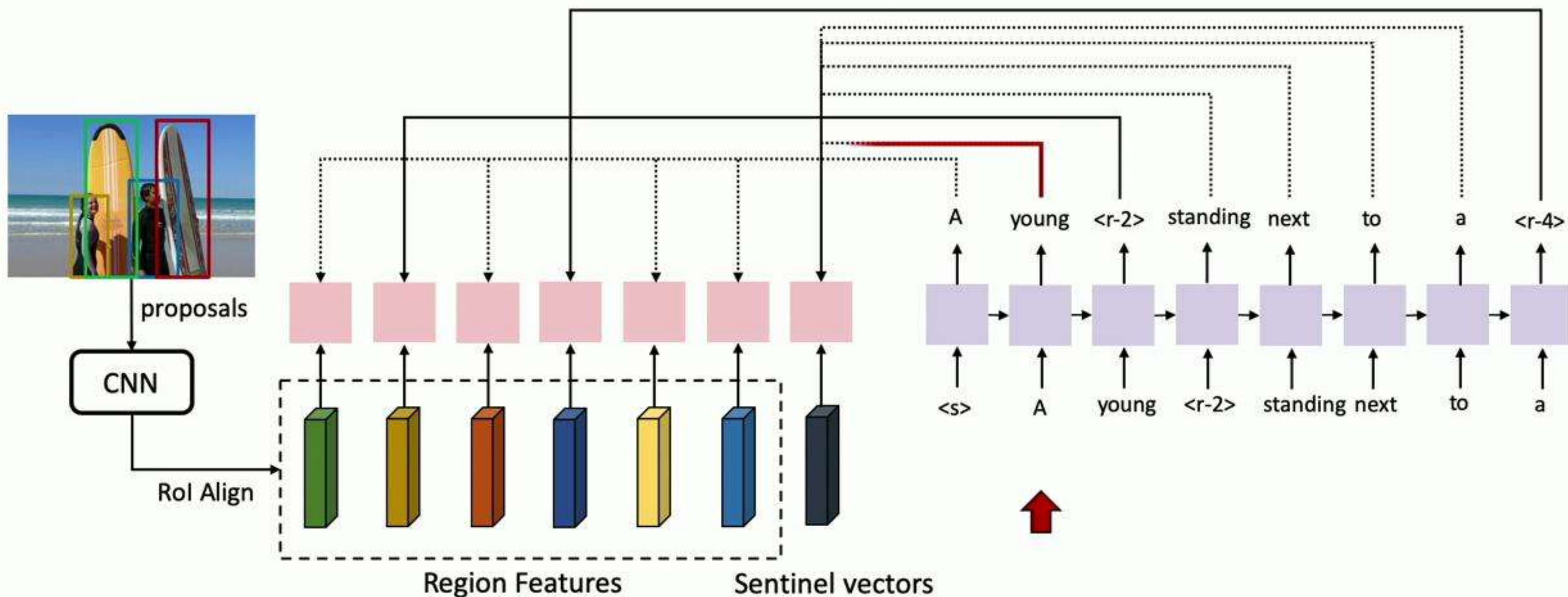
Model

Use pointer networks [Vinyals et.al 2015] with visual sentinel [Lu et.al 2017] to generate caption templates.



Model

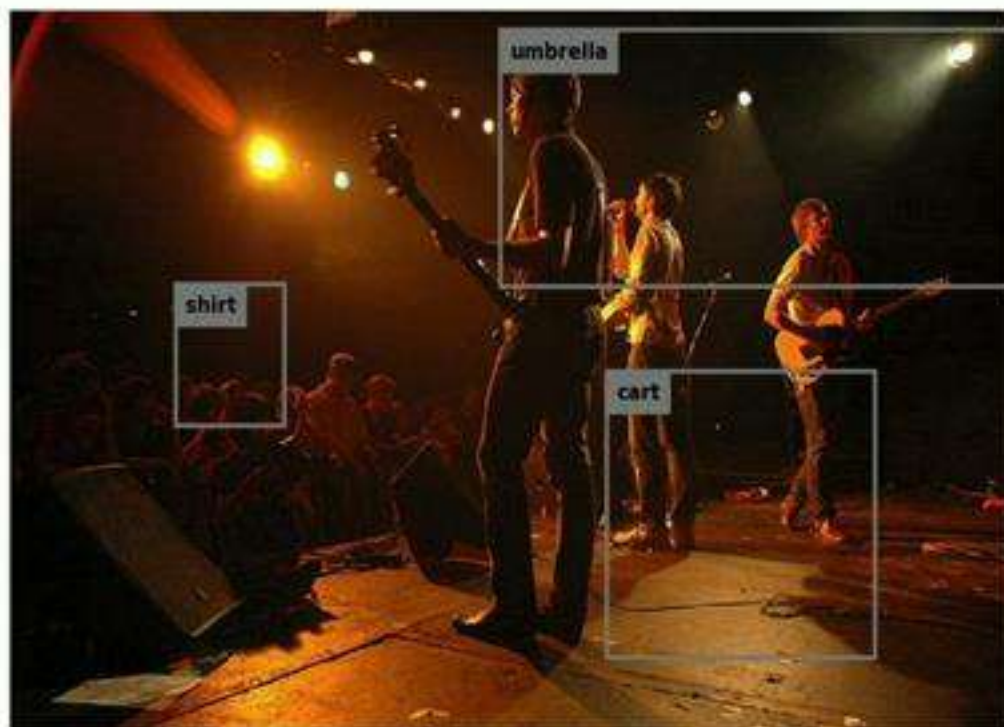
Use pointer networks [Vinyals et.al 2015] with visual sentinel [Lu et.al 2017] to generate caption templates.



How to handle inaccurate object detection?

Model

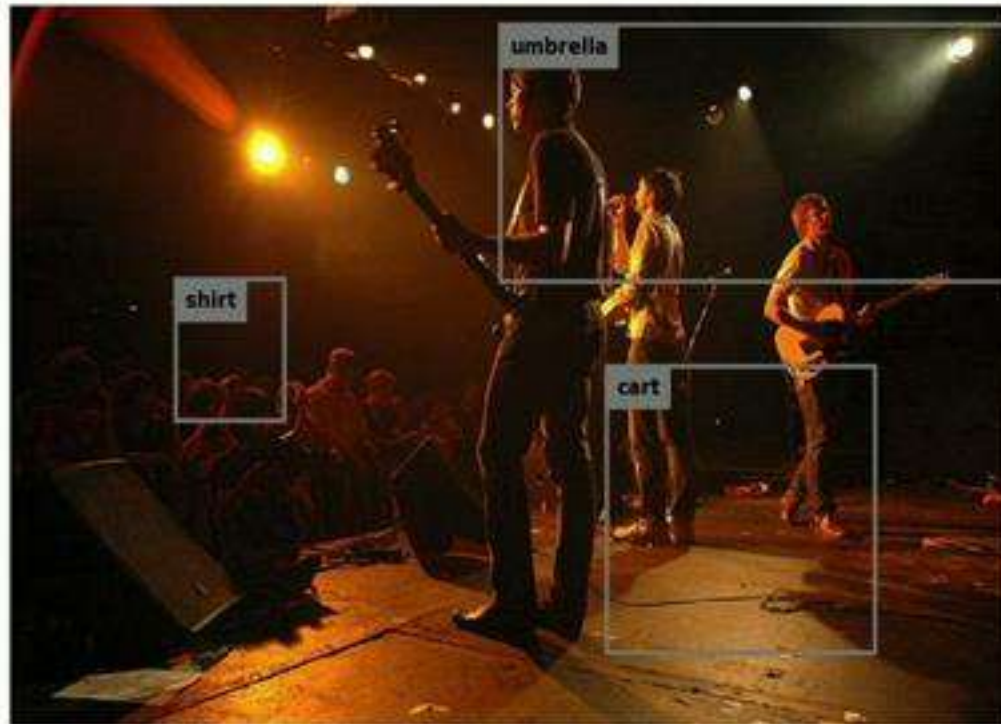
How to handle inaccurate object detection?



A man is playing guitar on the stage.

Model

How to handle inaccurate object detection?

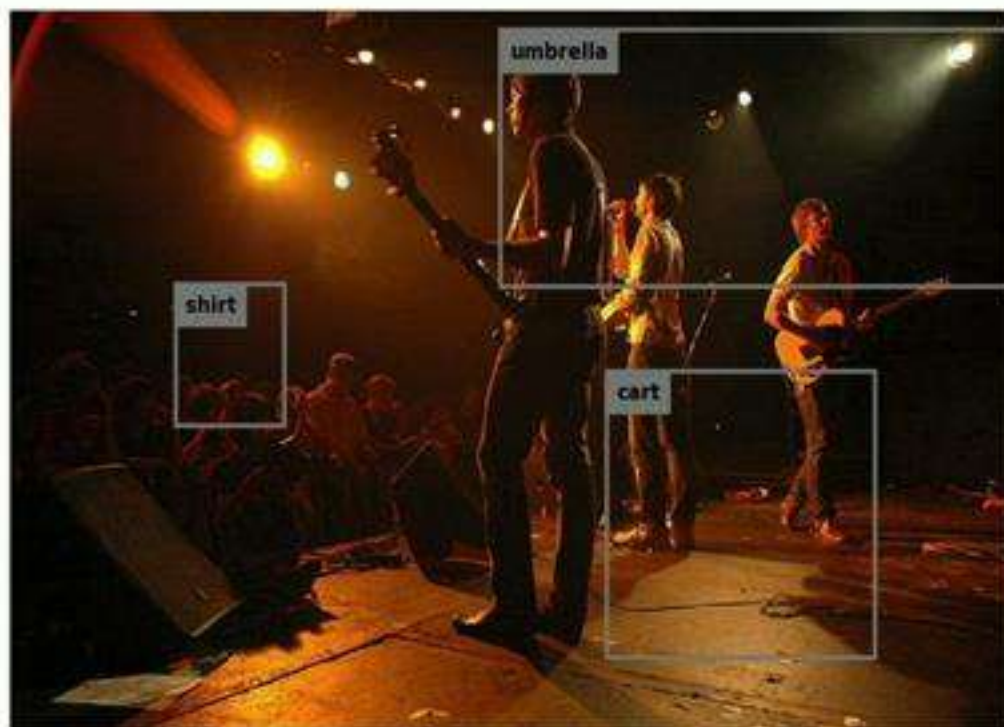


- Treat object labels as caption template.

Template: *A man is playing guitar on the stage.*

A man is playing guitar on the stage.

How to handle inaccurate object detection?



A man is playing guitar on the stage.

- Treat object labels as caption template.

Template: *A man is playing guitar on the stage.*

- Dynamically identify the visual words.

$\text{IoU}(\text{detected } \text{bbox}, \text{GT } \text{bbox}) > 0.5 \text{ and labels} == \text{words}$

Filling in the slots

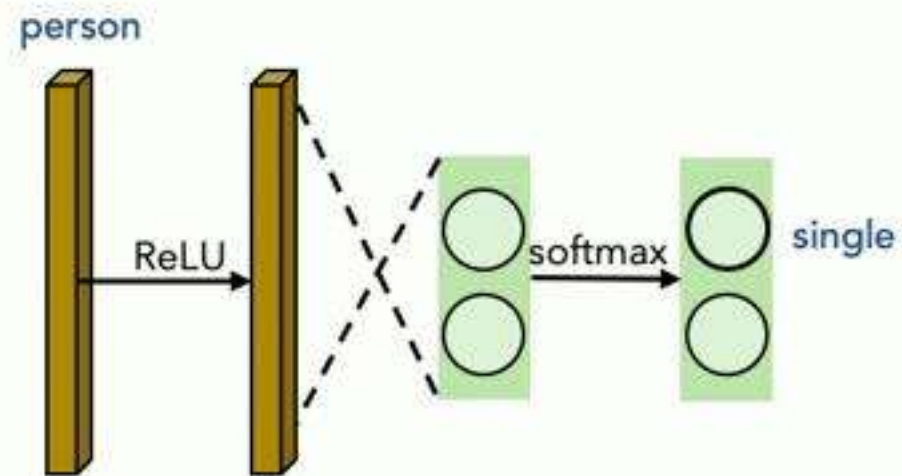
Input: A young <region-2> standing next to a <region-3> holding a <region-5>.
person, person, surfboard (coarse names from the object detector)

Model

Filling in the slots

Input: A young <region-2> standing next to a <region-3> holding a <region-5>.
person, person, surfboard (coarse names from the object detector)

1) Classify Plurality

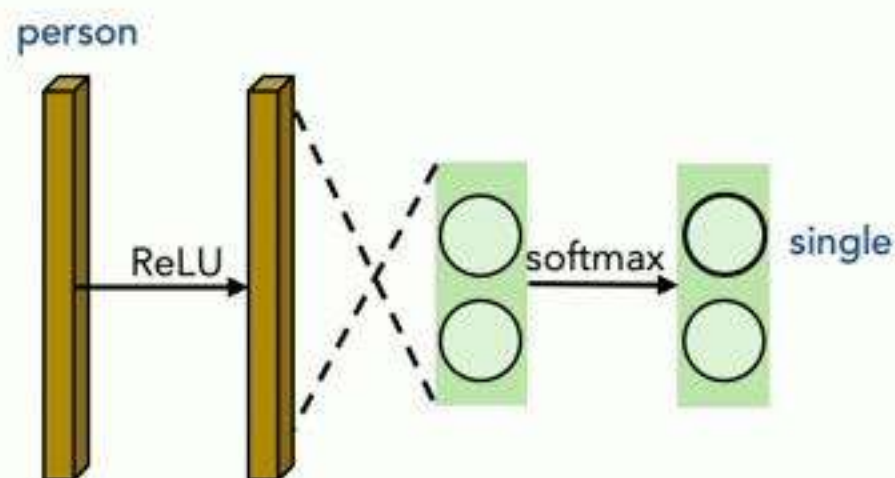


Model

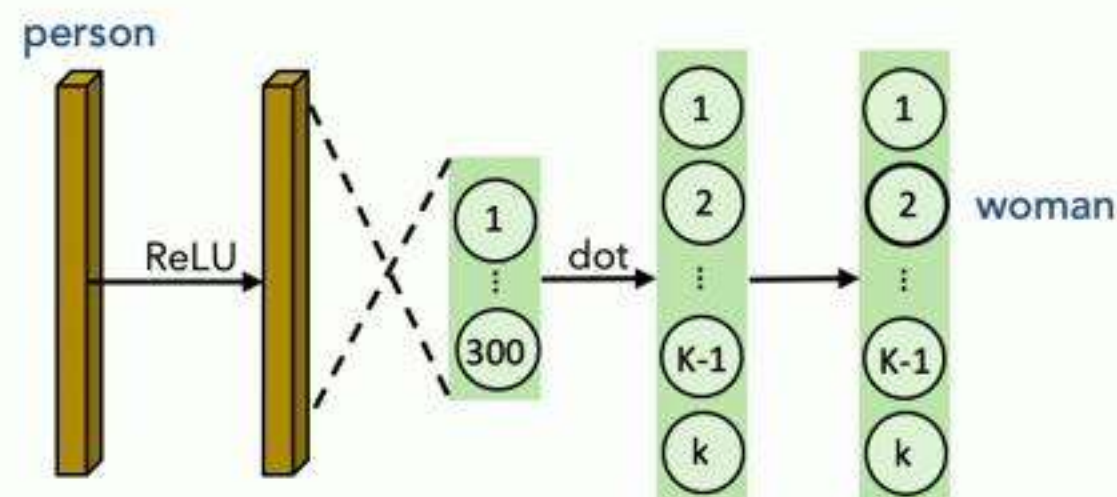
Filling in the slots

Input: A young <region-2> standing next to a <region-3> holding a <region-5>.
person, person, surfboard (coarse names form the object detector)

1) Classify Plurality



2) Determine Fine Grained Category

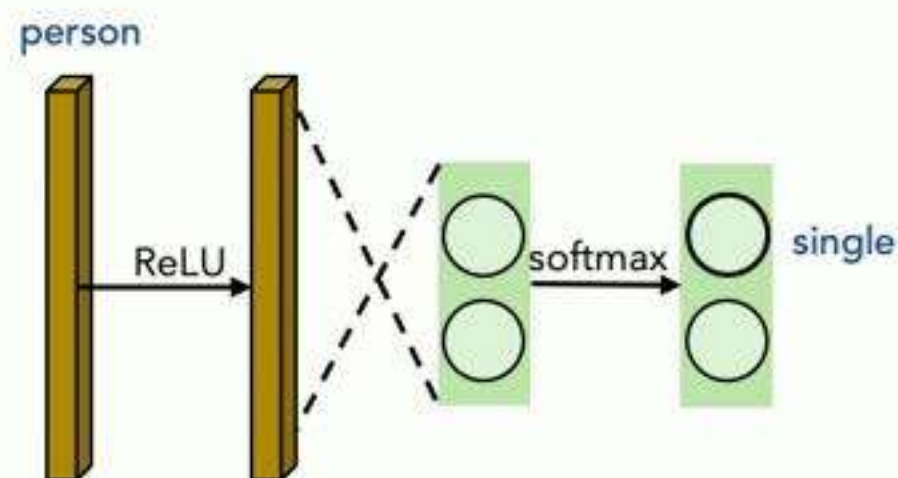


Model

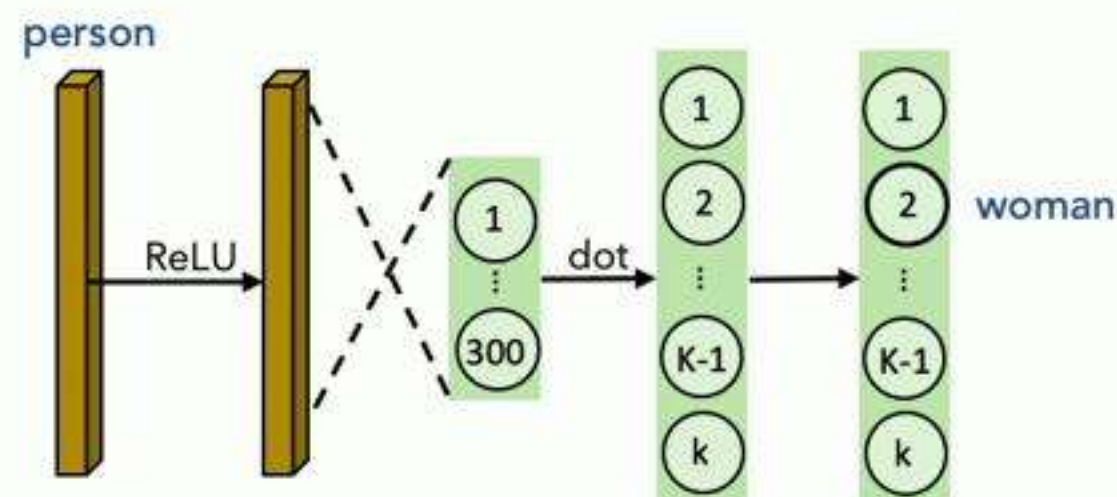
Filling in the slots

Input: A young <region-2> standing next to a <region-3> holding a <region-5>.
person, person, surfboard (coarse names form the object detector)

1) Classify Plurality



2) Determine Fine Grained Category



Output: A young woman standing next to a man holding a surfboard.

Objective

Minimize this cross-entropy loss

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | \mathbf{y}_{1:t-1}^*) 1_{(y_t^* = y^{txt})} + p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*) \left(\frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*) \right) 1_{(y_t^* = y^{vis})} \right)$$

Objective

Minimize this cross-entropy loss

Template word prob.

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | \mathbf{y}_{1:t-1}^*) 1_{(y_t^* = y^{txt})} + p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*) \left(\frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*) \right) 1_{(y_t^* = y^{vis})} \right)$$

Objective

Minimize this cross-entropy loss

Template word prob.

Refinement prob.

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | y_{1:t-1}^*) 1_{(y_t^* = y^{txt})} + p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*) \left(\frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*) \right) 1_{(y_t^* = y^{vis})} \right)$$

Objective

Minimize this cross-entropy loss

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(\begin{array}{c} \text{Template word prob.} \\ p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | \mathbf{y}_{1:t-1}^*) 1_{(y_t^* = y^{txt})} \end{array} + \begin{array}{c} \text{Refinement prob.} \\ p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*) \end{array} \left(\begin{array}{c} \text{target region prob.} \\ \frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*) 1_{(y_t^* = y^{vis})} \end{array} \right) \right)$$

Objective

Minimize this cross-entropy loss

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(\begin{array}{c} \text{Template word prob.} \\ p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | \mathbf{y}_{1:t-1}^*) 1_{(y_t^* = y^{txt})} \end{array} + \begin{array}{c} \text{Refinement prob.} \\ p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*) \end{array} \left(\begin{array}{c} \text{target region prob.} \\ \frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*) \end{array} \right) 1_{(y_t^* = y^{vis})} \right)$$

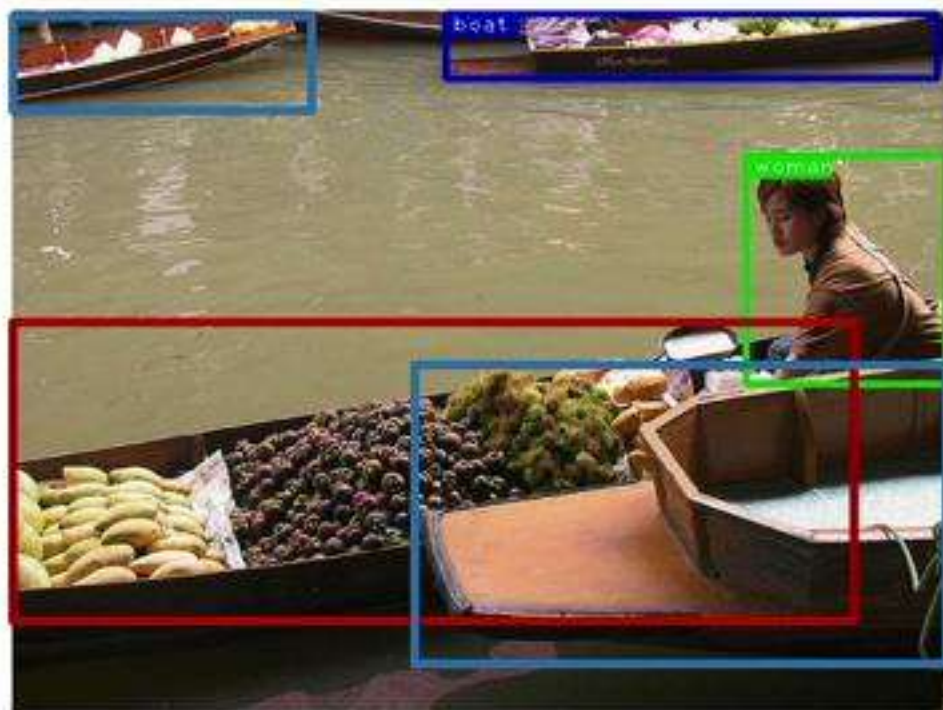
Co-reference when different kind of supervision exists.

Objective

Minimize this cross-entropy loss

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(\underbrace{p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | y_{1:t-1}^*) 1_{(y_t^* = y^{txt})}}_{\text{Template word prob.}} + \underbrace{p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*)}_{\text{Refinement prob.}} \left(\underbrace{\frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*)}_{\text{target region prob.}} \right) 1_{(y_t^* = y^{vis})} \right)$$

Co-reference when different kind of supervision exists.



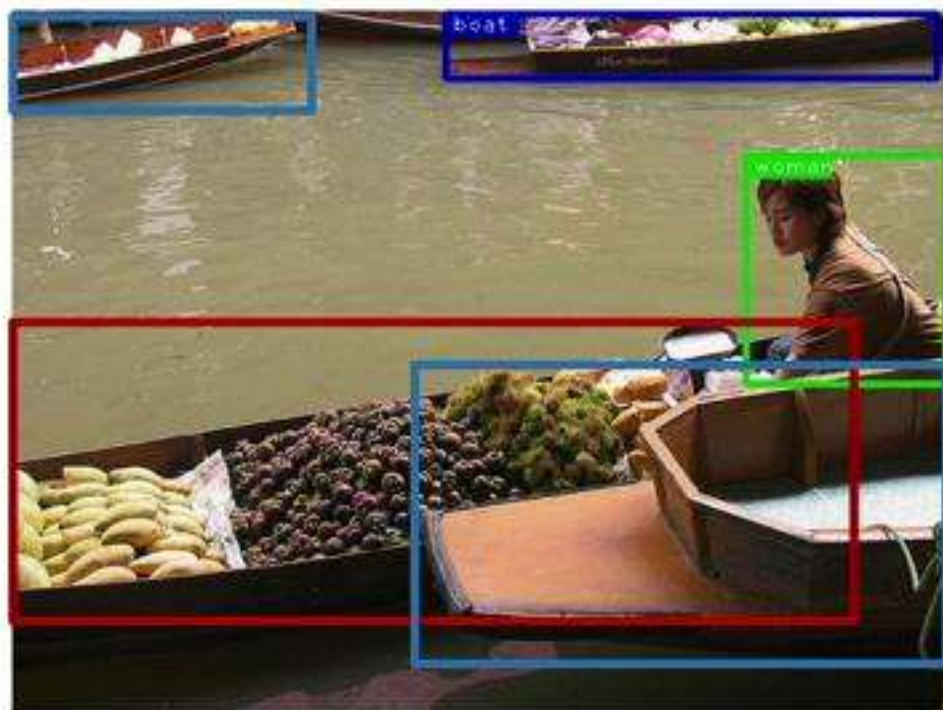
A young woman is sitting inside a boat.

Objective

Minimize this cross-entropy loss

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(\begin{array}{c} \text{Template word prob.} \\ p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | y_{1:t-1}^*) 1_{(y_t^* = y^{txt})} \end{array} + \begin{array}{c} \text{Refinement prob.} \\ p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*) \end{array} \left(\begin{array}{c} \text{target region prob.} \\ \frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*) \end{array} \right) 1_{(y_t^* = y^{vis})} \right)$$

Co-reference when different kind of supervision exists.



A young woman is sitting inside a boat.

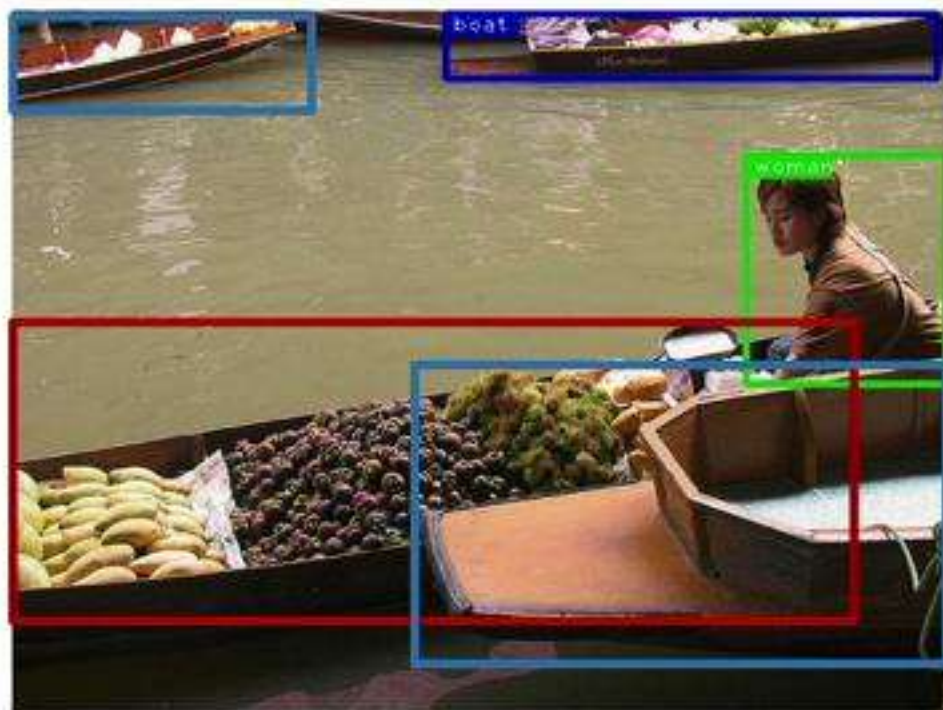
Problem: which boat is the caption referred to.

Objective

Minimize this cross-entropy loss

$$L(\boldsymbol{\theta}) = - \sum_{t=1}^T \log \left(\begin{array}{c} \text{Template word prob.} \\ p(y_t^* | \tilde{r}, y_{1:t-1}^*) p(\tilde{r} | y_{1:t-1}^*) 1_{(y_t^* = y^{txt})} \end{array} + \begin{array}{c} \text{Refinement prob.} \\ p(b_t^*, s_t^* | \mathbf{r}_t, y_{1:t-1}^*) \end{array} \left(\begin{array}{c} \text{target region prob.} \\ \frac{1}{m} \sum_{i=1}^m p(r_t^i | y_{1:t-1}^*) 1_{(y_t^* = y^{vis})} \end{array} \right) \right)$$

Co-reference when different kind of supervision exists.



A young woman is sitting inside a boat.

Problem: which boat is the caption referred to.

Solution: maximize the averaged target region prob.

Datasets

- Flickr30k: 31,783 images, 5 captions per image, 275,555 annotated bounding boxes.
- COCO: 164,062 images, 5 captions per image.

Object category to words

- For COCO dataset. (e.g., "person" mapping to ["child", "baker", ...])

Caption pre-processing

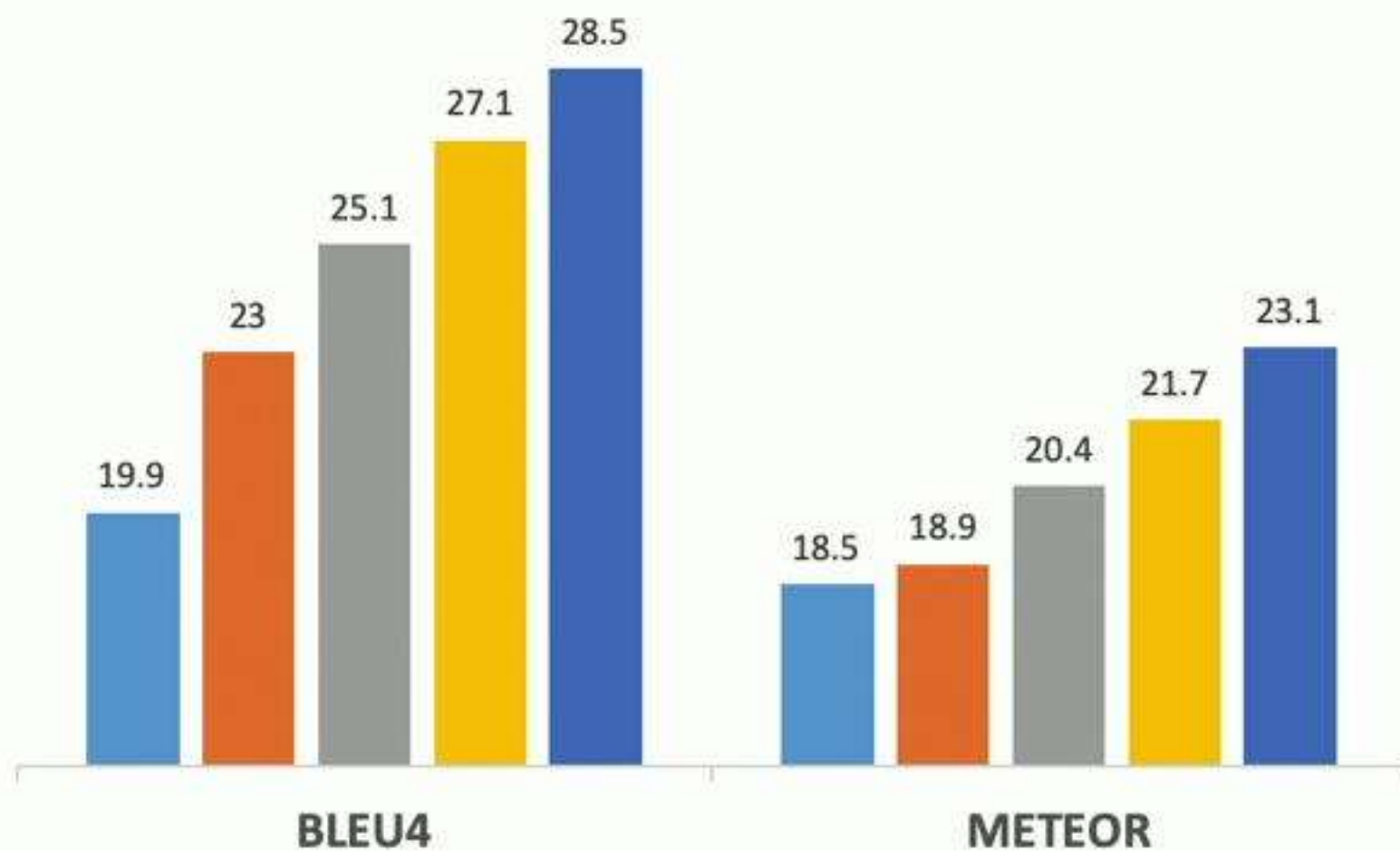
- Caption truncation (if > 16 words)
- Building vocabulary (9,587 words for COCO, 6,864 words for Flickr30k)

Results

Standard Image Captioning

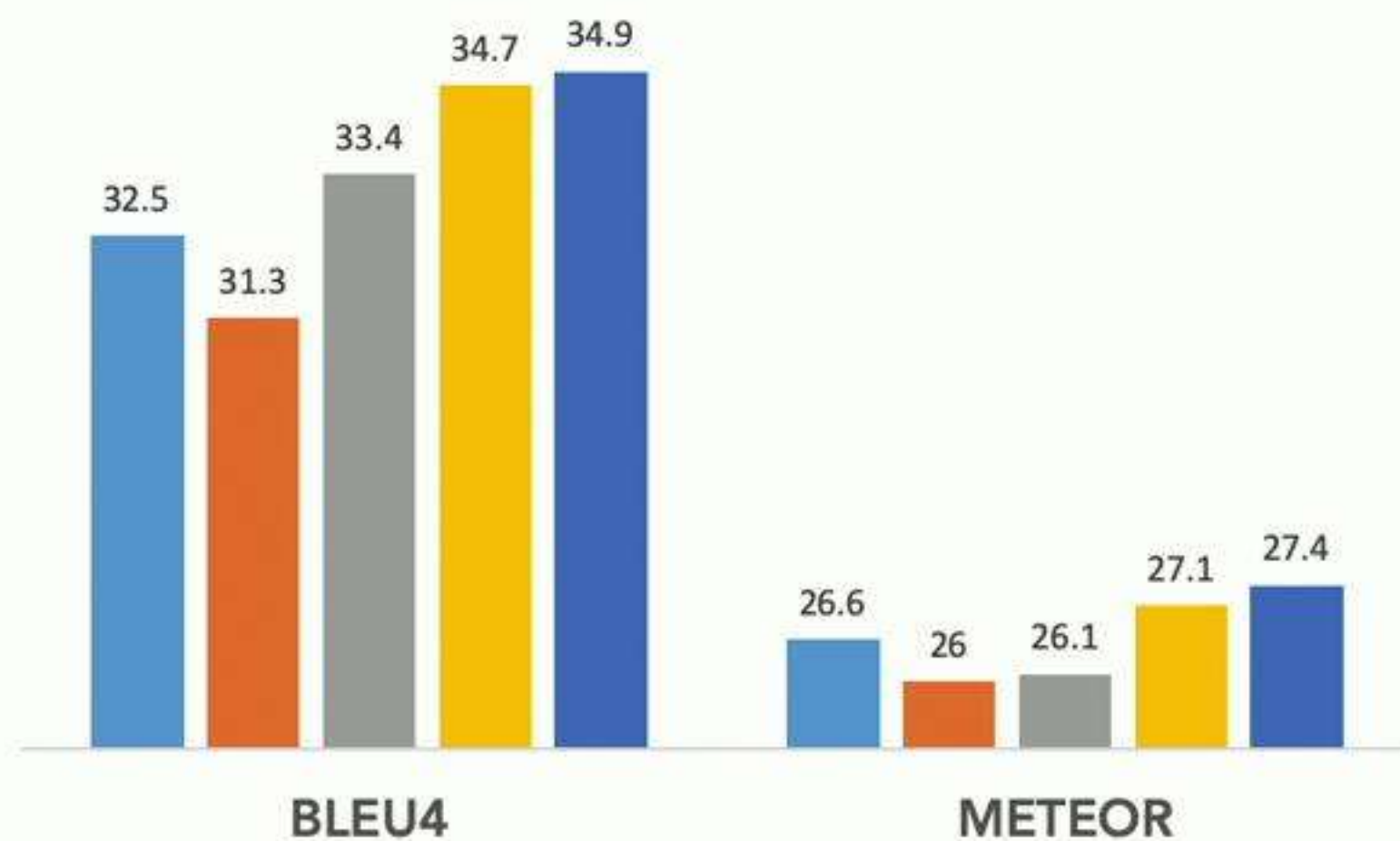
Flickr30k

■ Hard-Attention ■ ATT-FCN ■ Adaptive ■ NBT ■ NBT*



COCO

■ Adaptive ■ Att2in ■ Up-Down ■ NBT ■ NBT*

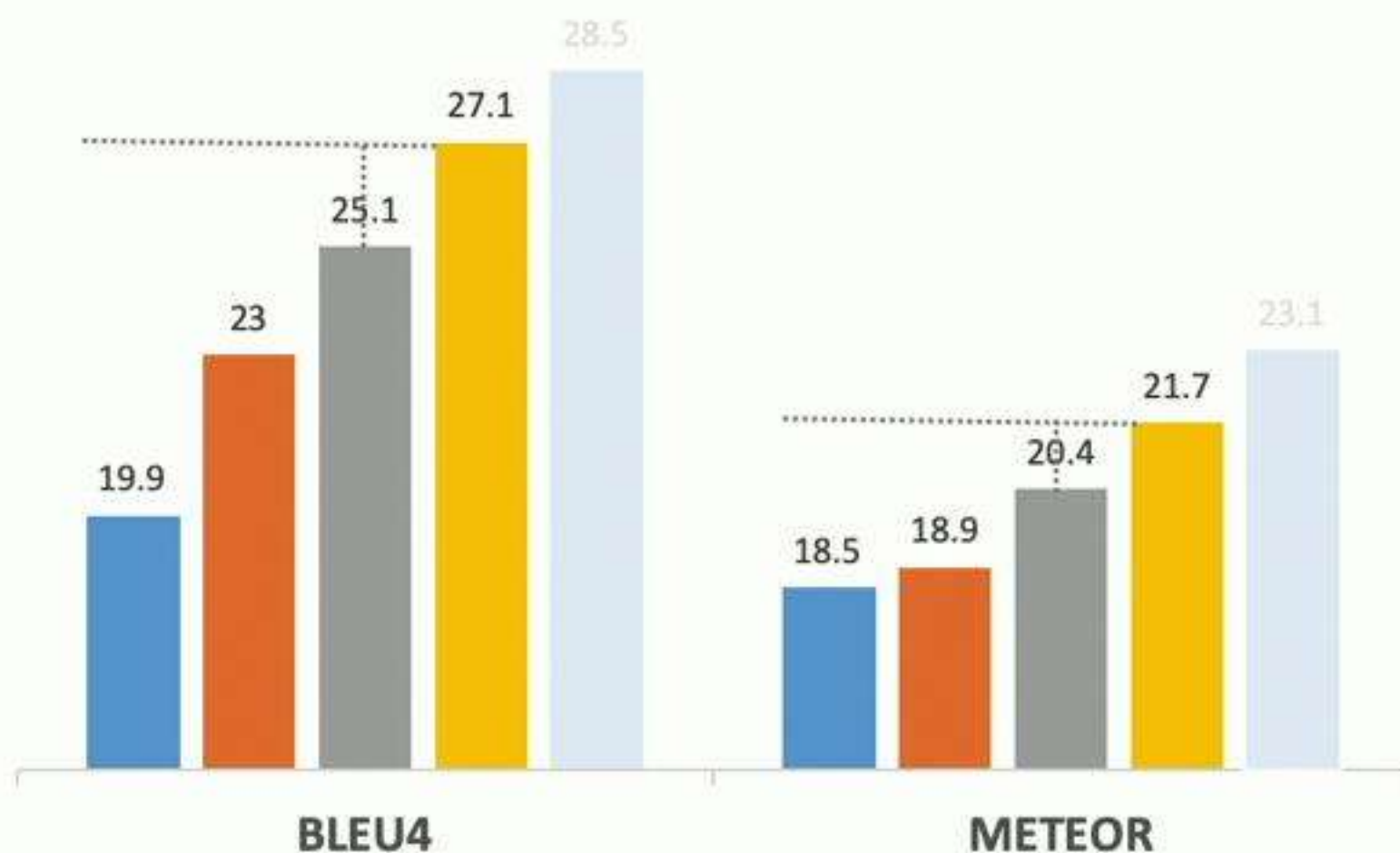


Results

Standard Image Captioning

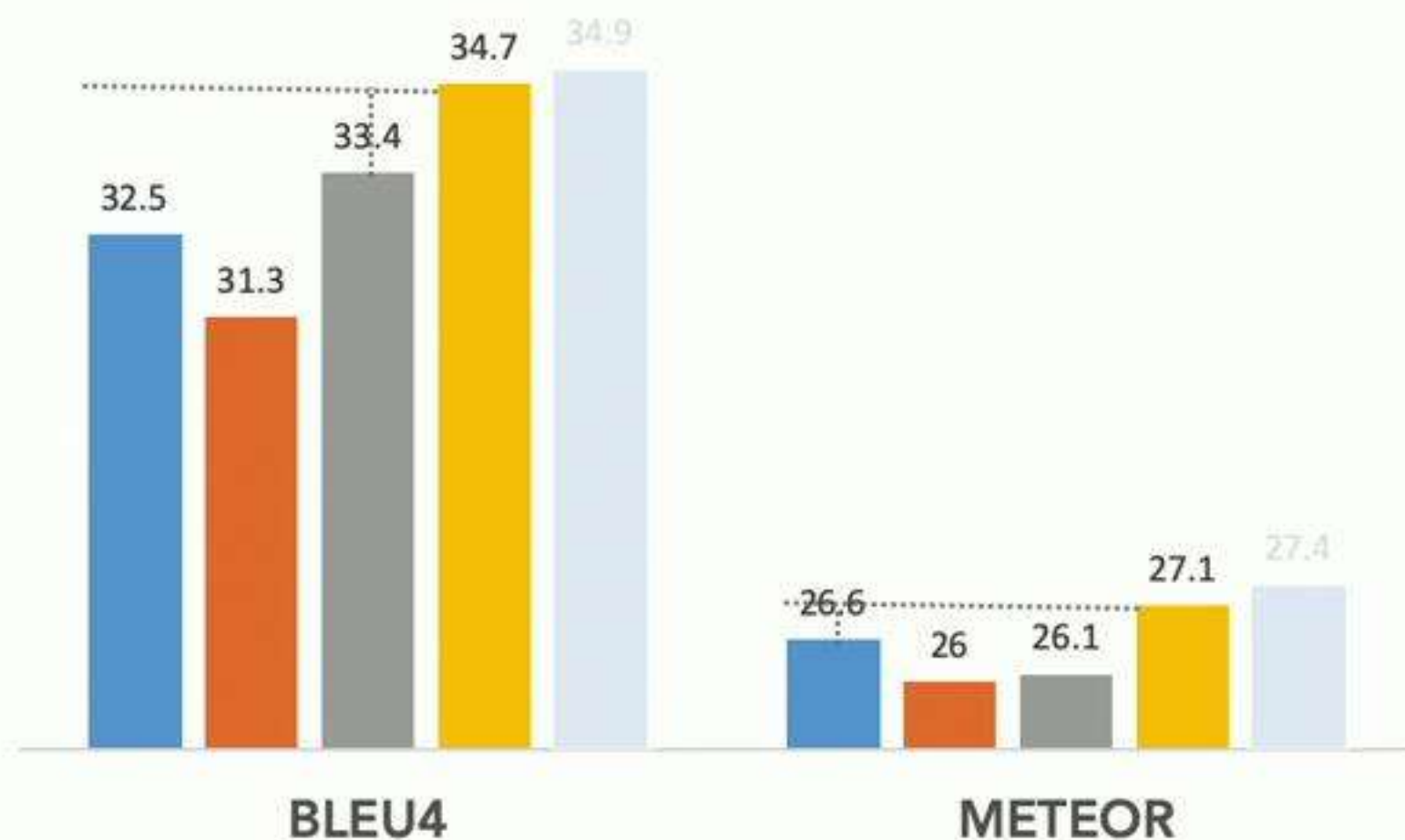
Flickr30k

■ Hard-Attention ■ ATT-FCN ■ Adaptive ■ NBT ■ NBT*



COCO

■ Adaptive ■ Att2in ■ Up-Down ■ NBT ■ NBT*

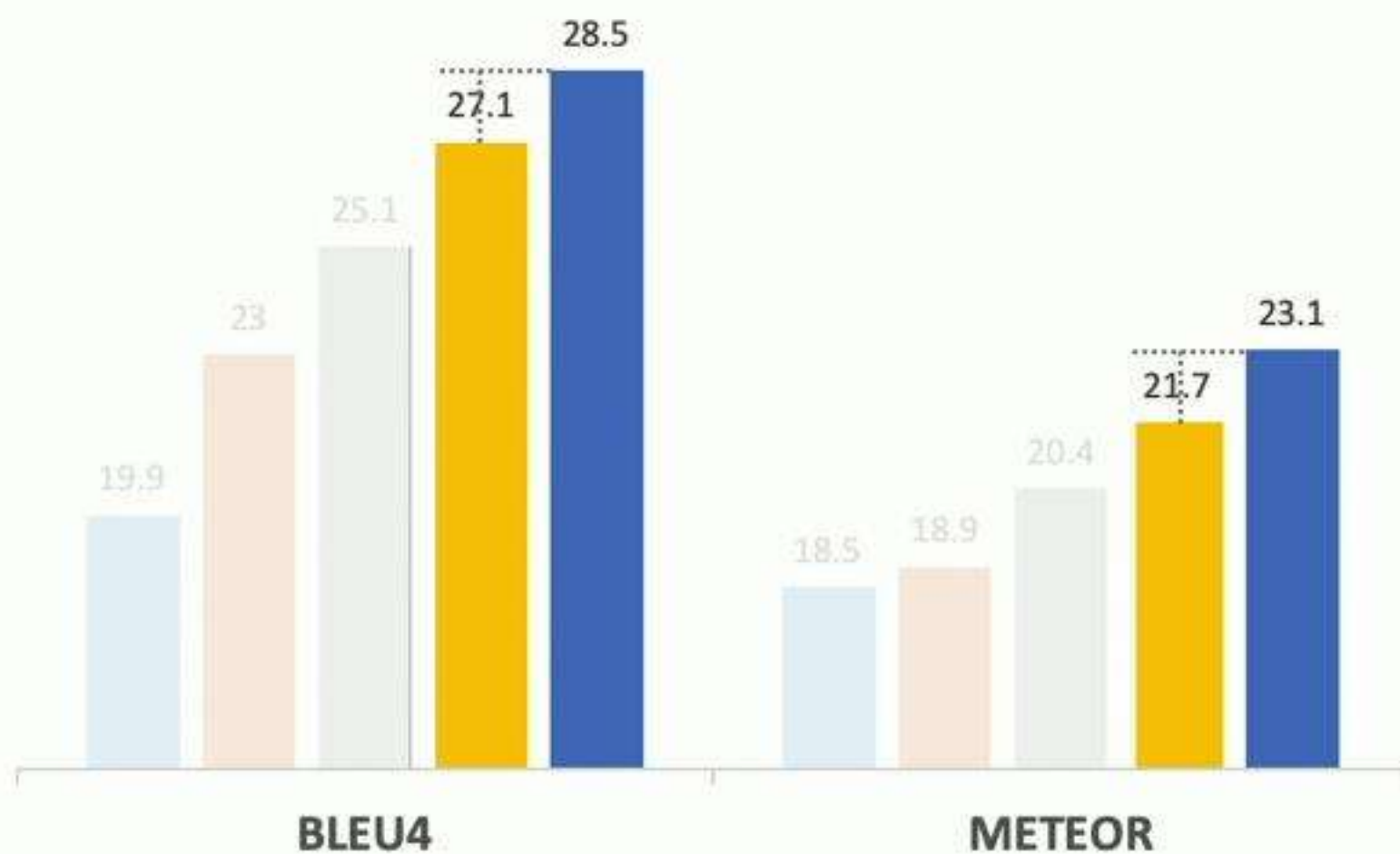


Results

Standard Image Captioning

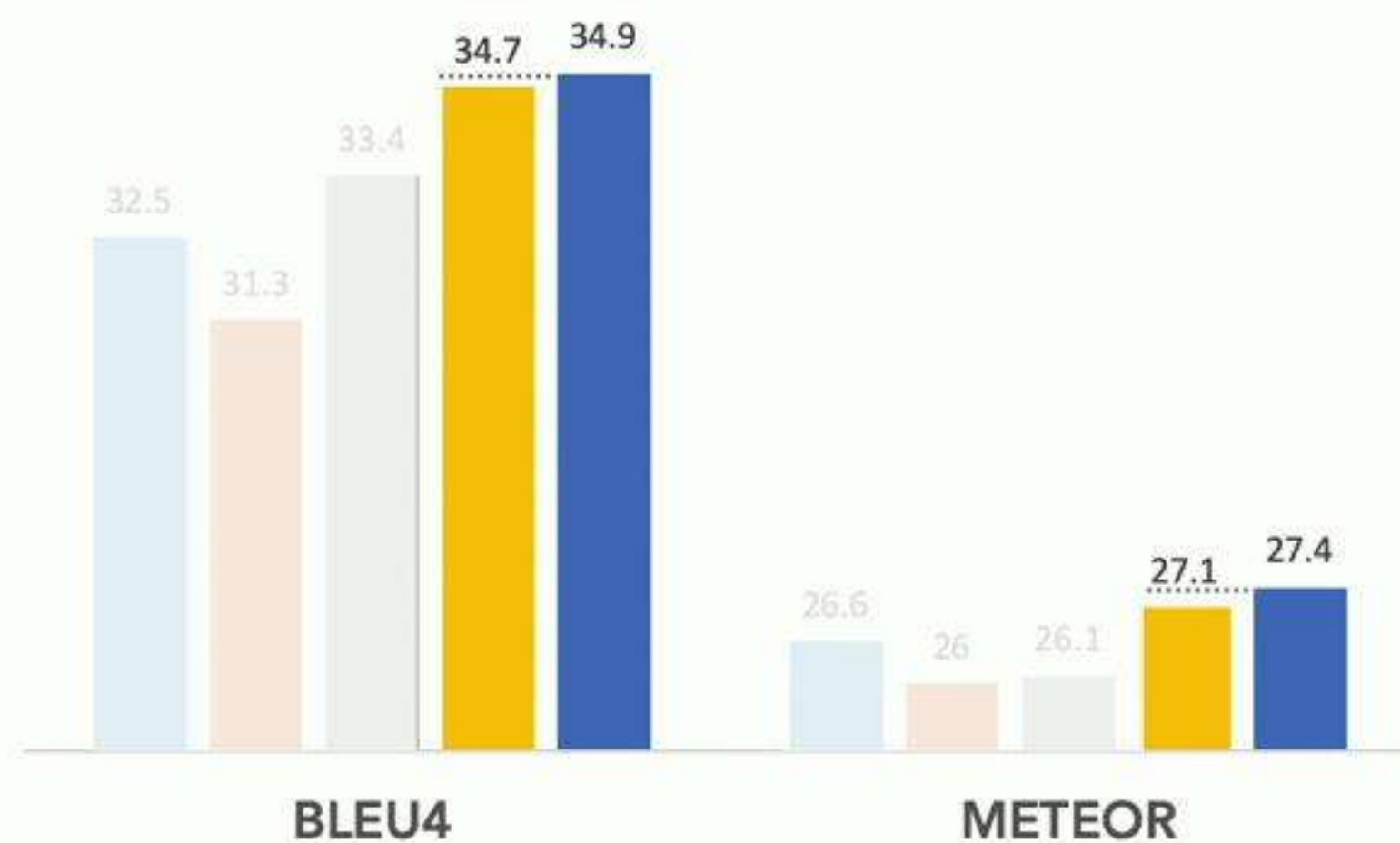
Flickr30k

■ Hard-Attention ■ ATT-FCN ■ Adaptive ■ NBT ■ NBT*



COCO

■ Adaptive ■ Att2in ■ Up-Down ■ NBT ■ NBT*

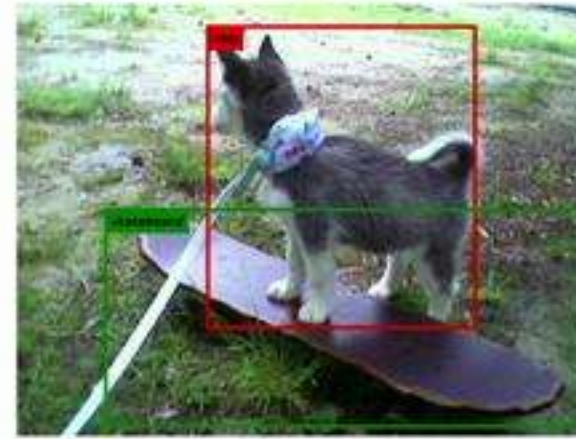
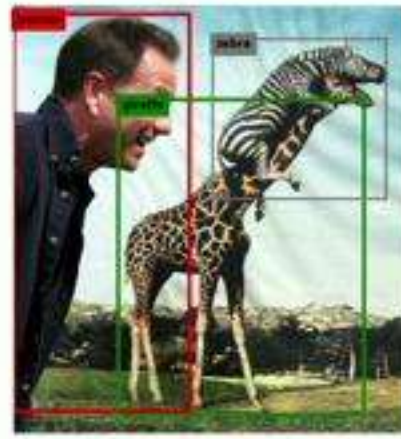
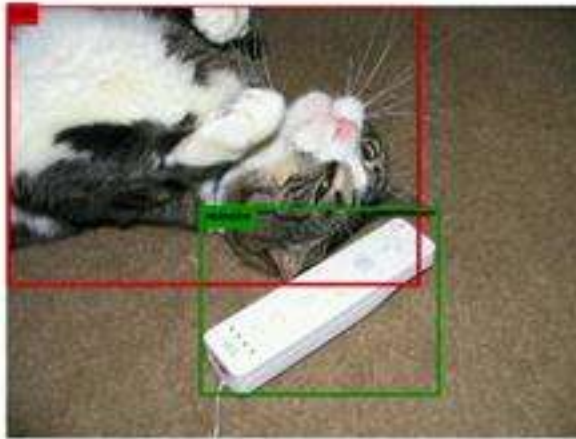


Robust Image Captioning

Results

Robust Image Captioning

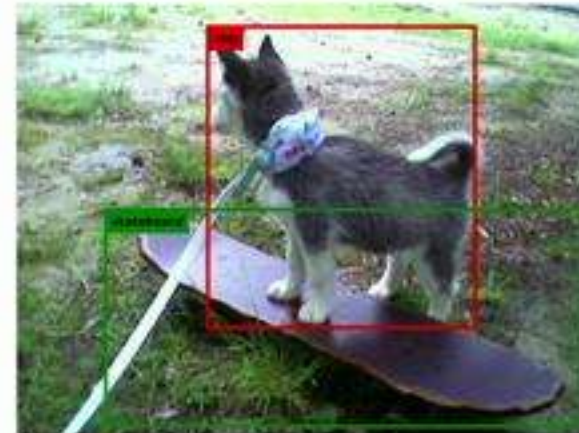
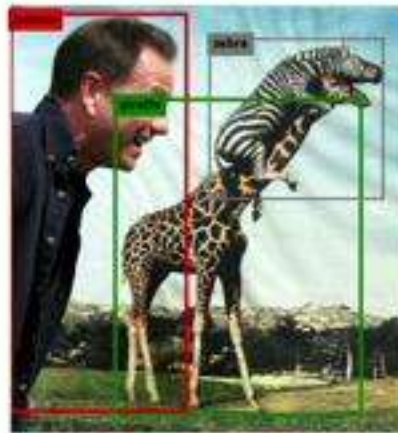
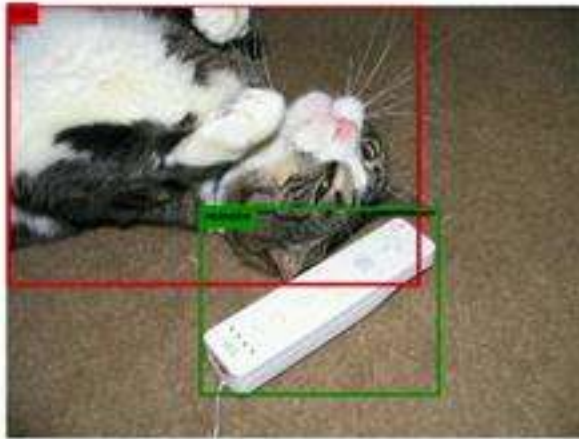
- To evaluate image captioning for novel scene compositions.



Results

Robust Image Captioning

- To evaluate image captioning for novel scene compositions.



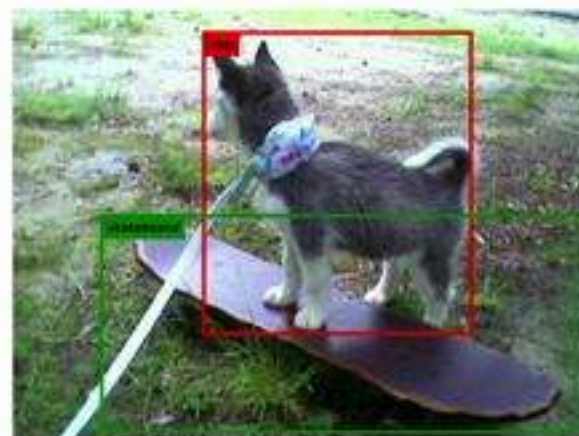
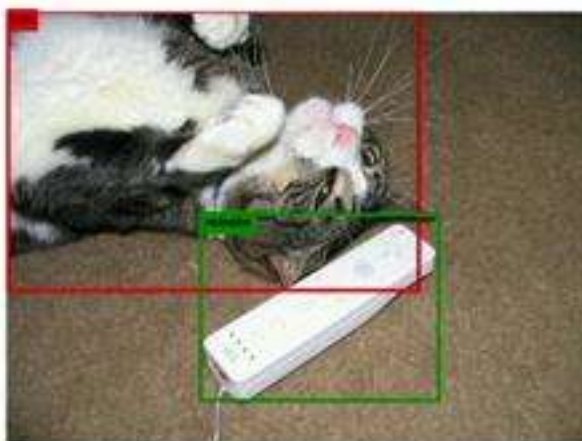
Robust-COCO split

- Distribution of co-occurring objects in train data is different from test data.
- Sufficient examples from each category in train set.
- Novel compositions (pairs) of categories in test set.

Results

Robust Image Captioning

- To evaluate image captioning for novel scene compositions.



Robust-COCO split

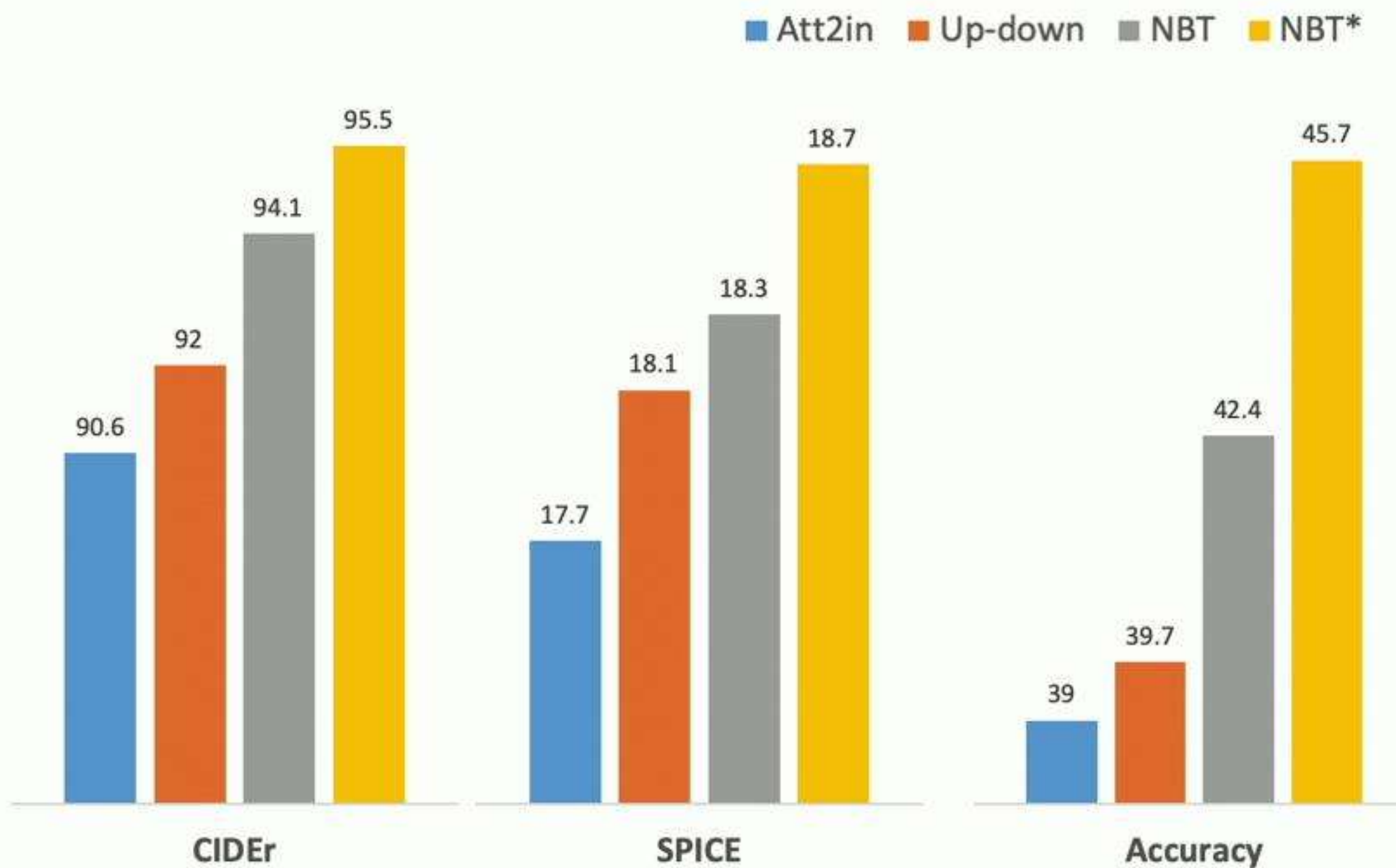
- Distribution of co-occurring objects in train data is different from test data.
- Sufficient examples from each category in train set.
- Novel compositions (pairs) of categories in test set.

Accuracy (new metric)

- Whether or not a generated caption includes the new object combination.
- 100% accuracy for at least one mention of the novel category pair.

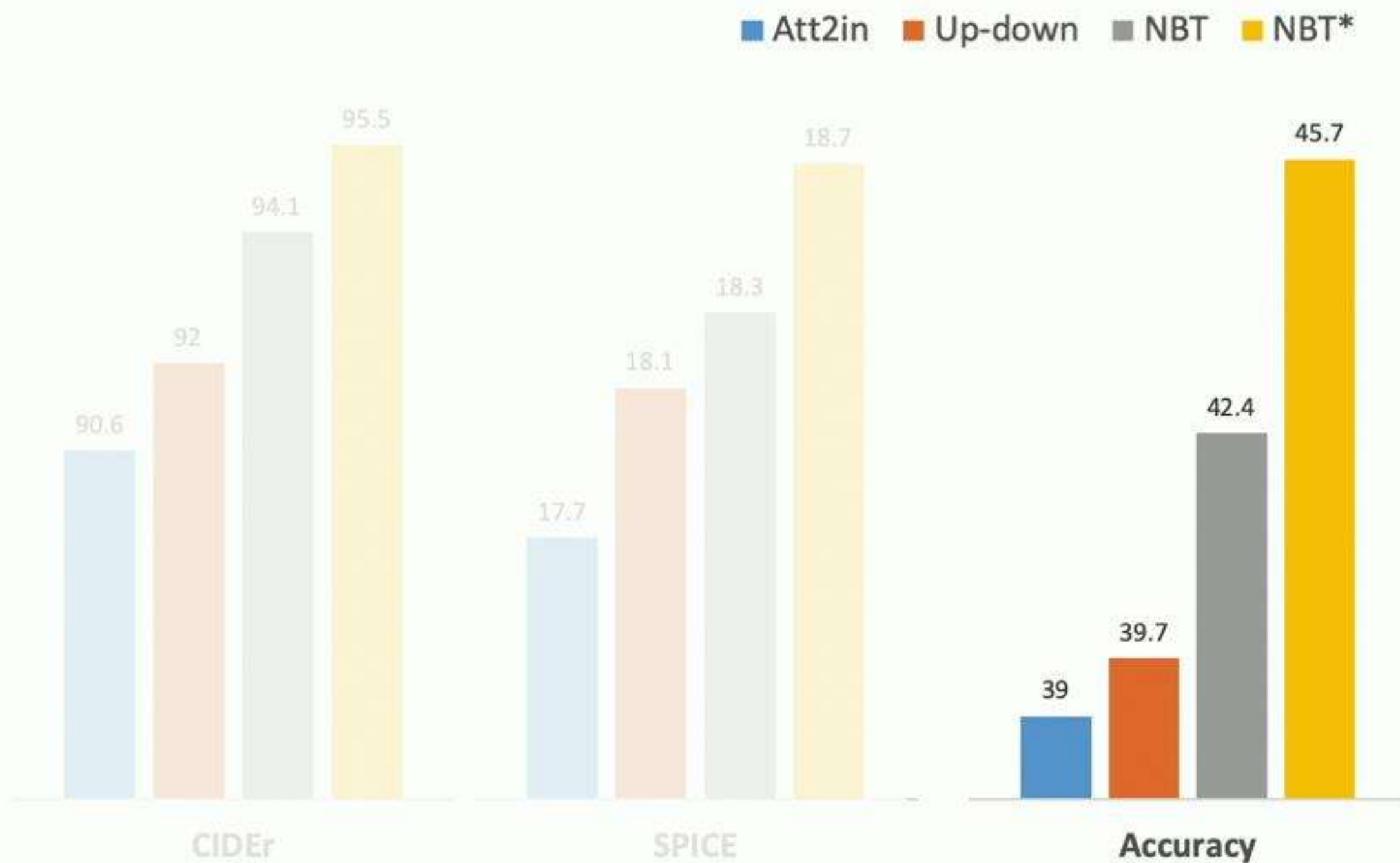
Results

Robust Image Captioning



Results

Robust Image Captioning



Novel Object Captioning [*Hendricks et.al.*]

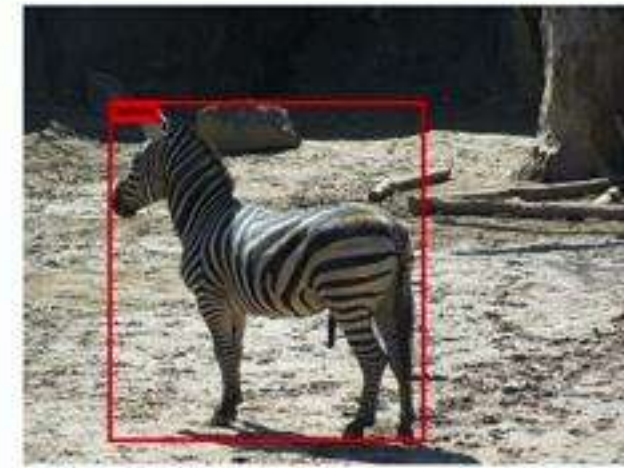
Novel Object Captioning [*Hendricks et.al.*]

- Describe image with novel objects.

Results

Novel Object Captioning [Hendricks et.al.]

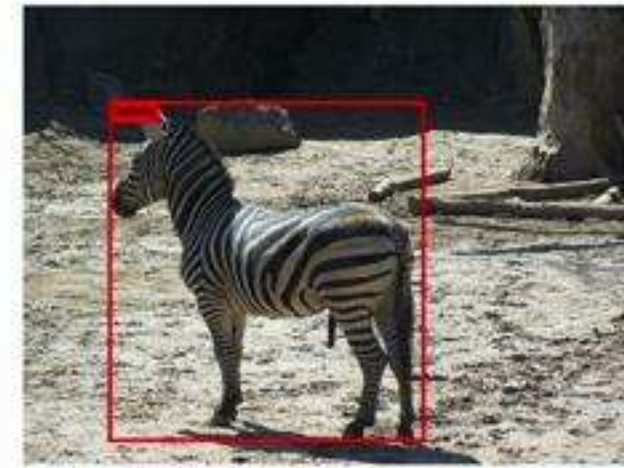
- Describe image with novel objects.
- Excludes all the pairs that contain at least one of the eight objects in COCO ("bottle", "bus", "couch", "microwave", "pizza", "racket", "suitcase", and "zebra")



Results

Novel Object Captioning [Hendricks et.al.]

- Describe image with novel objects.
- Excludes all the pairs that contain at least one of the eight objects in COCO ("bottle", "bus", "couch", "microwave", "pizza", "racket", "suitcase", and "zebra")

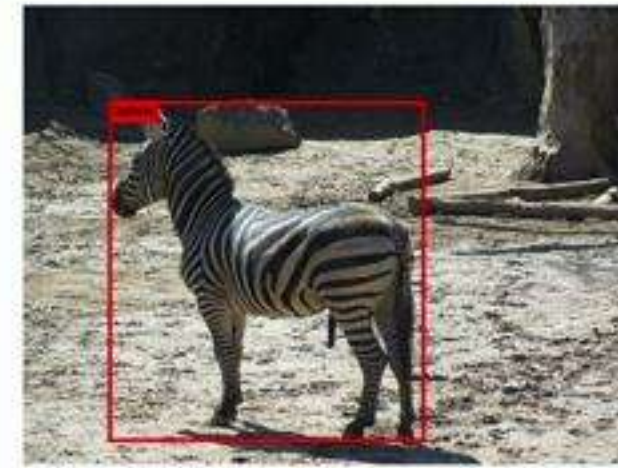


- Test set is split into in-domain and out-of-domain subsets.

Results

Novel Object Captioning [Hendricks et.al.]

- Describe image with novel objects.
- Excludes all the pairs that contain at least one of the eight objects in COCO ("bottle", "bus", "couch", "microwave", "pizza", "racket", "suitcase", and "zebra")

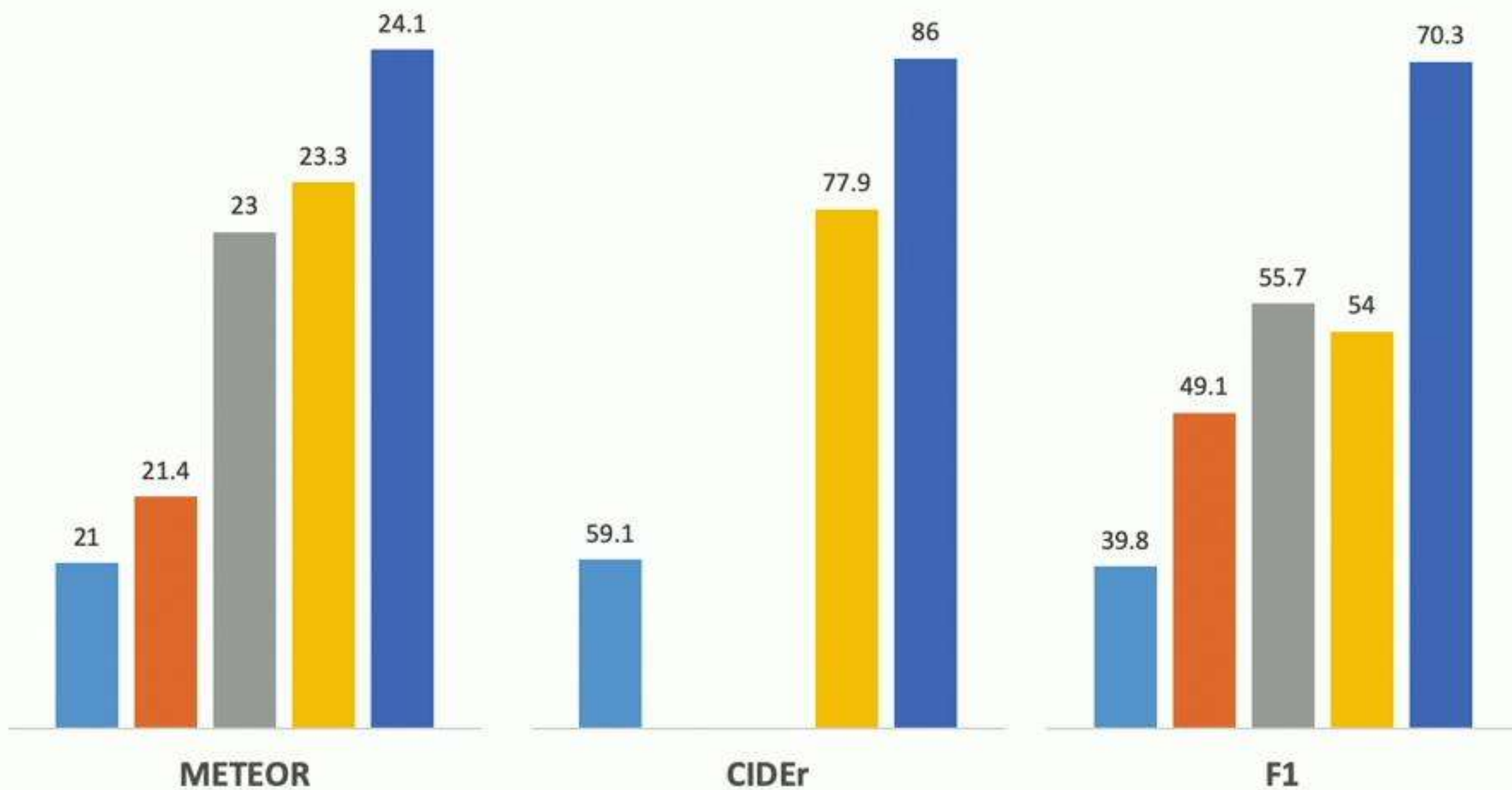


- Test set is split into in-domain and out-of-domain subsets.
- F1 score (metric)

Checks if the excluded object is correctly mentioned in the generated caption.

Results

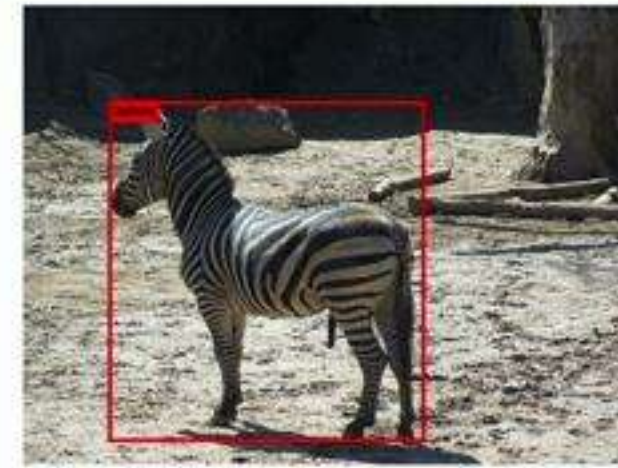
Novel Object Captioning



Results

Novel Object Captioning [Hendricks et.al.]

- Describe image with novel objects.
- Excludes all the pairs that contain at least one of the eight objects in COCO ("bottle", "bus", "couch", "microwave", "pizza", "racket", "suitcase", and "zebra")

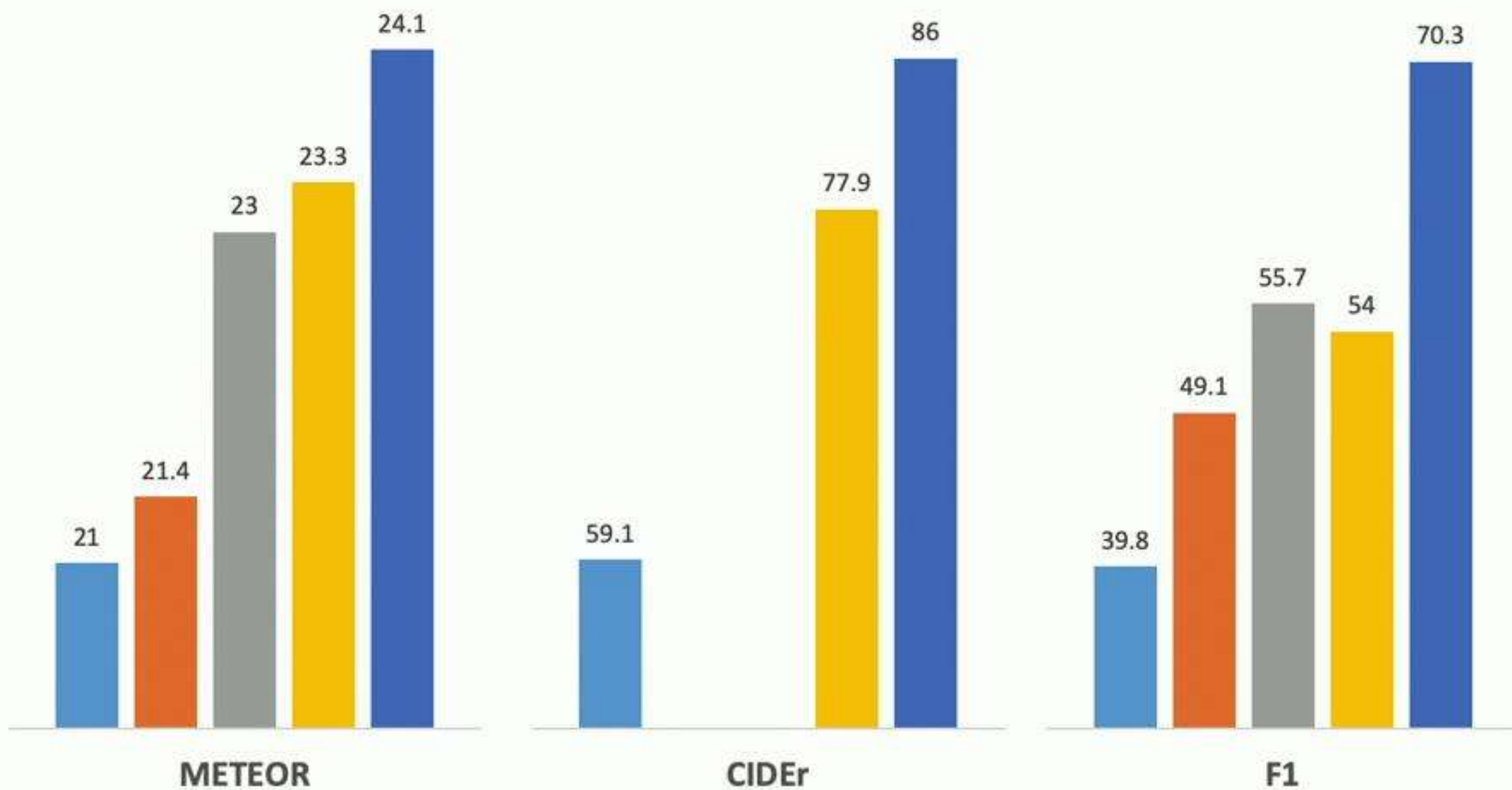


- Test set is split into in-domain and out-of-domain subsets.
- F1 score (metric)

Checks if the excluded object is correctly mentioned in the generated caption.

Results

Novel Object Captioning



Summary

Two-stage approach for image captioning

- generate hybrid template
- fills the slots with categories recognized by object detector

Summary

Two-stage approach for image captioning

- generate hybrid template
- fills the slots with categories recognized by object detector

Limitation

Summary

Two-stage approach for image captioning

- generate hybrid template
- fills the slots with categories recognized by object detector

Limitation

Only ground on noun words which found by object detector

- Extend to attributes, actions, OCR etc.

Summary

Two-stage approach for image captioning

- generate hybrid template
- fills the slots with categories recognized by object detector

Limitation

Only ground on noun words which found by object detector

- Extend to attributes, actions, OCR etc.

Limited by the image captioning dataset.

- Learn from other data sources -- ViLBERT.

This Talk



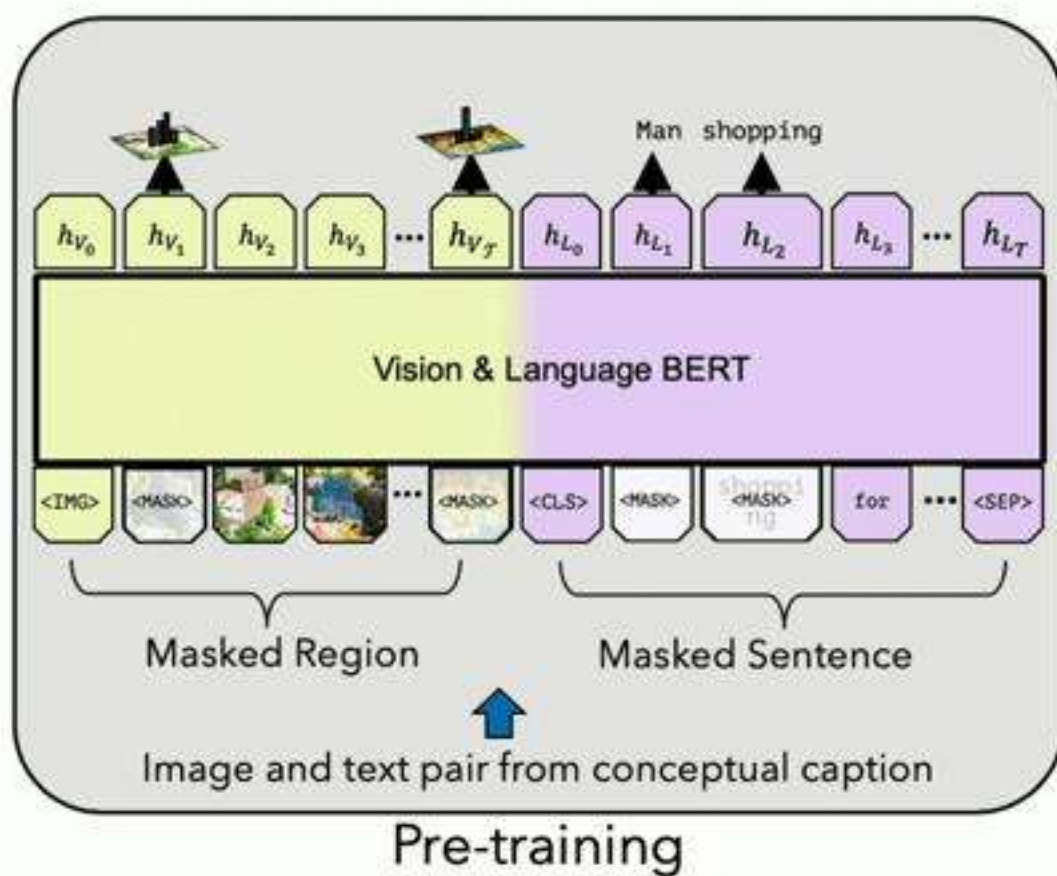
ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks

NeurIPS
2019



ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks

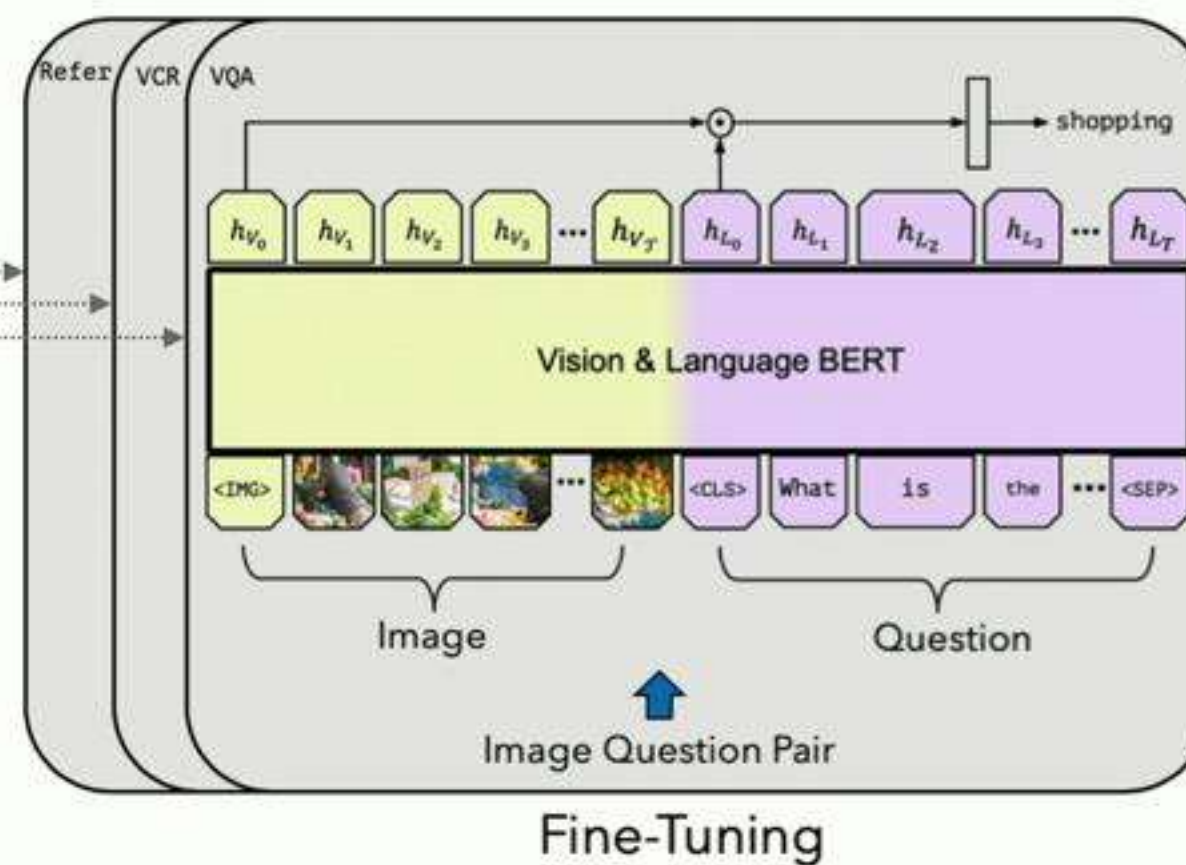
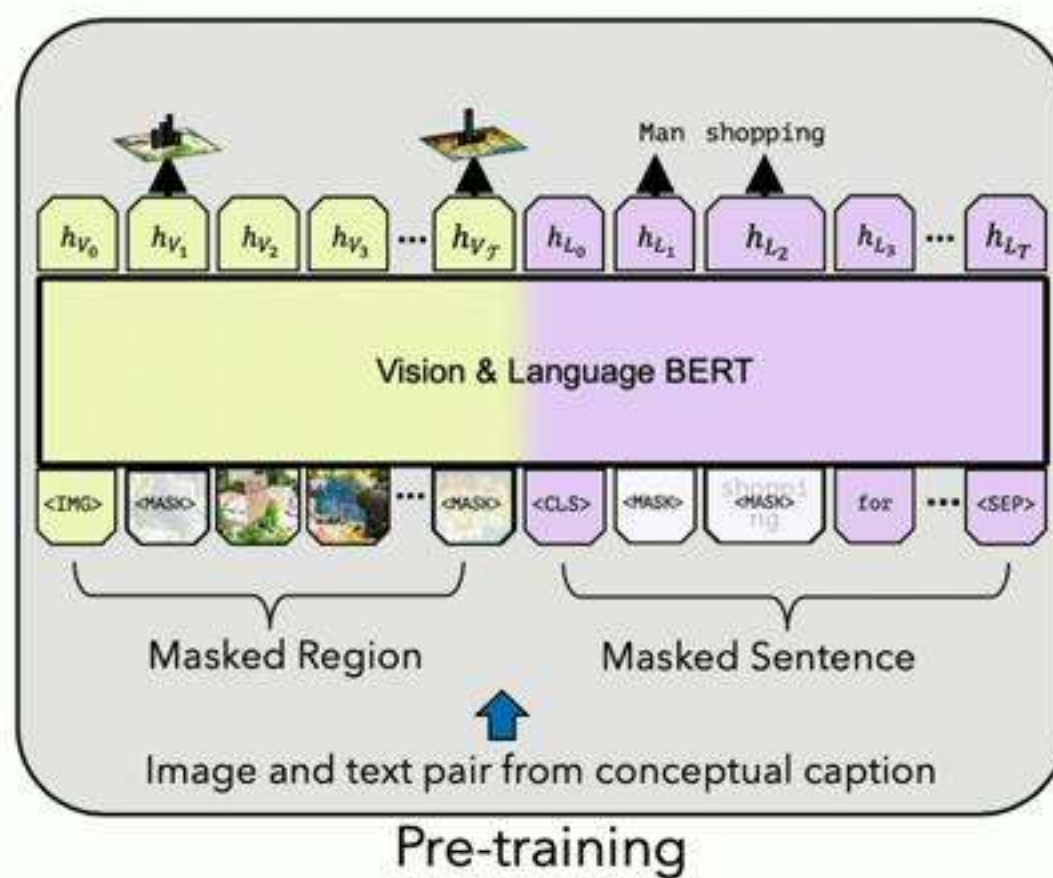
NeurIPS
2019



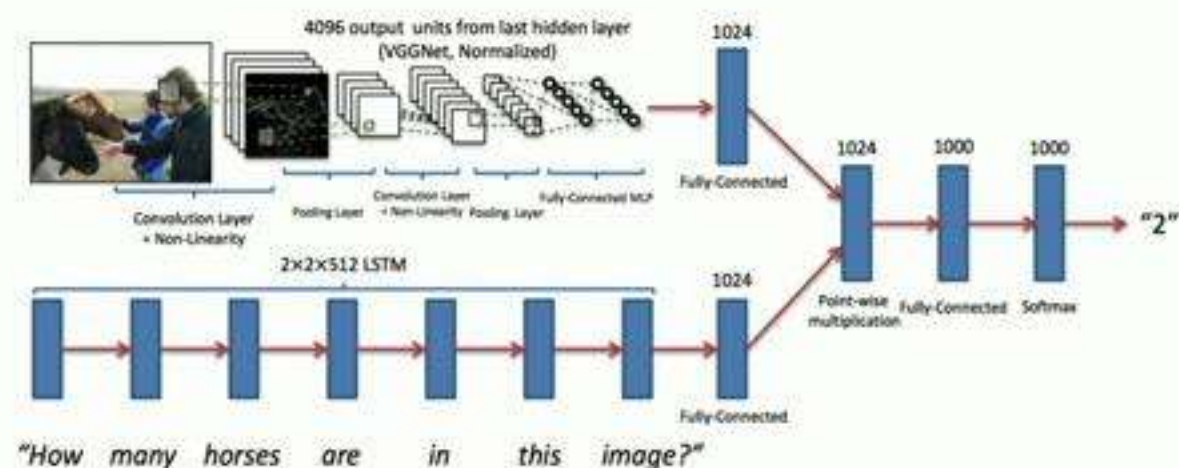


ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks

NeurIPS
2019



Vision and Language



Visual Question Answering

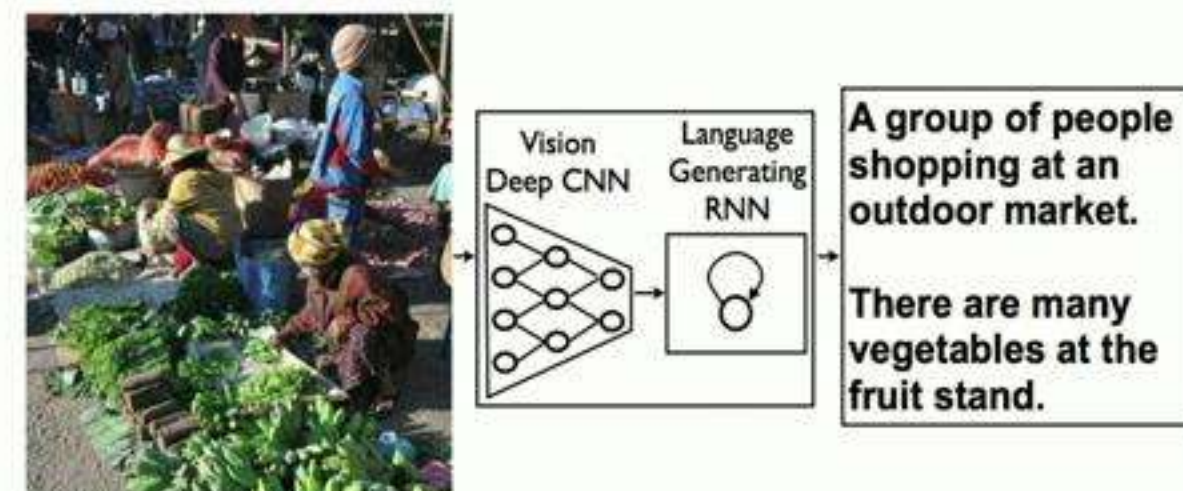
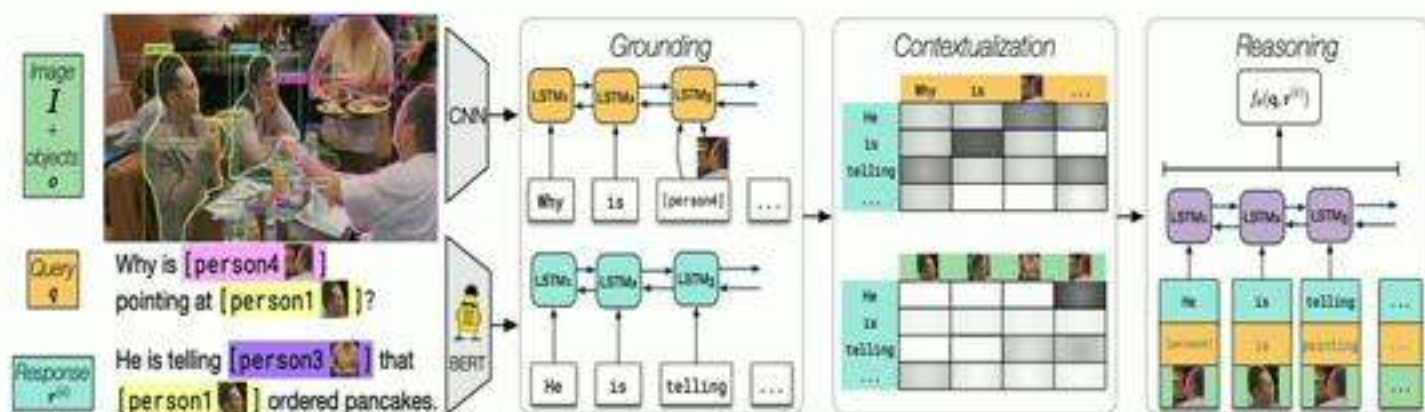
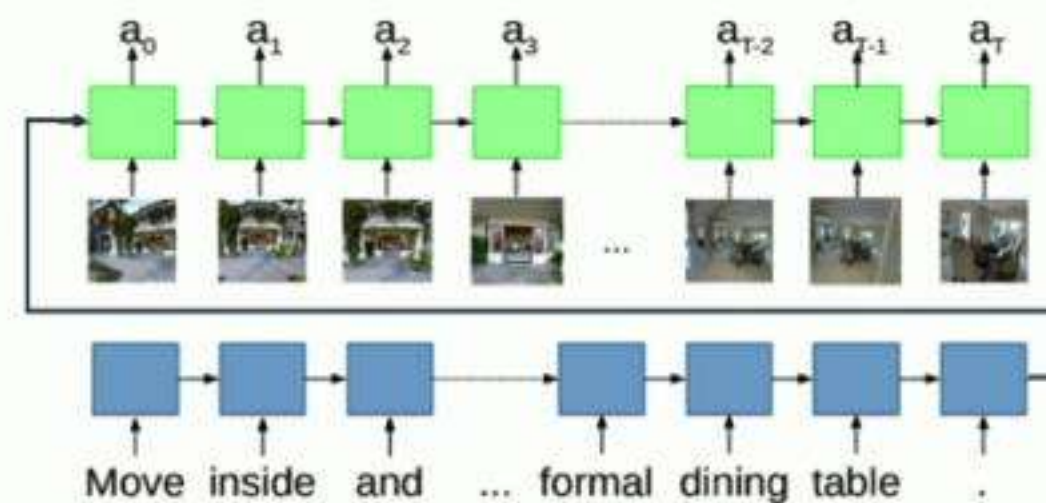


Image Captioning

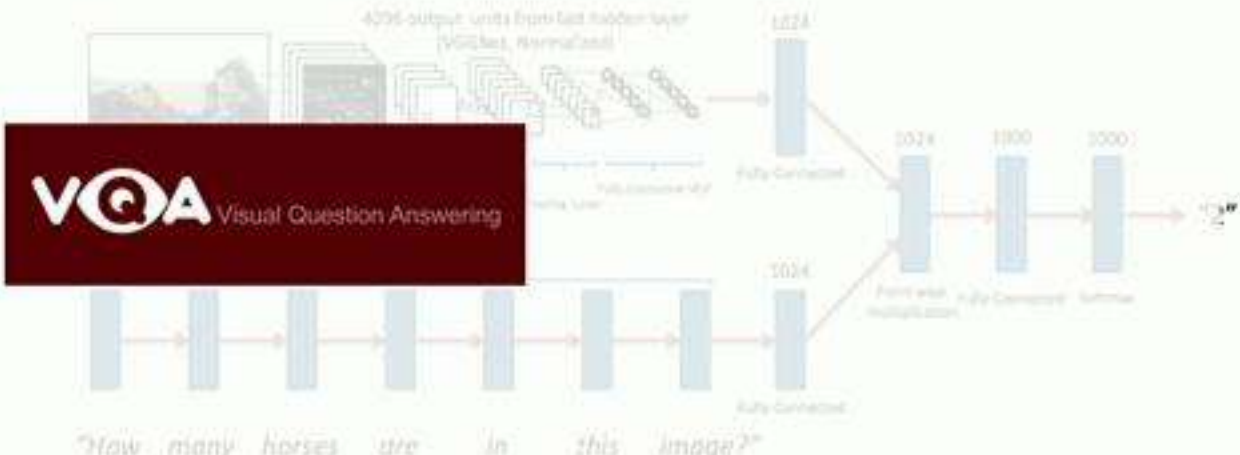


Visual Commonsense Reasoning

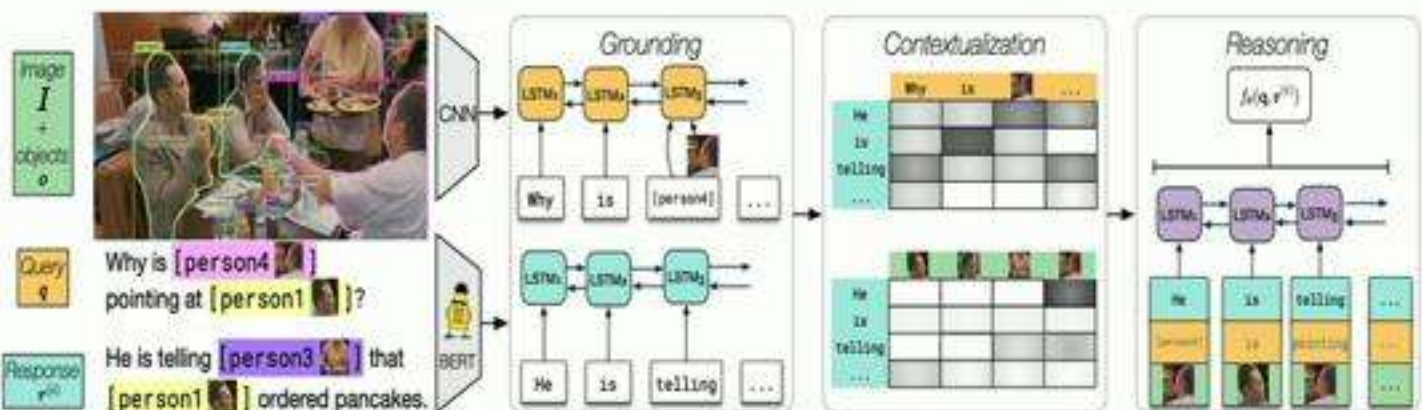


Vision and Language Navigation

Vision and Language



Visual Question Answering



Visual Commonsense Reasoning

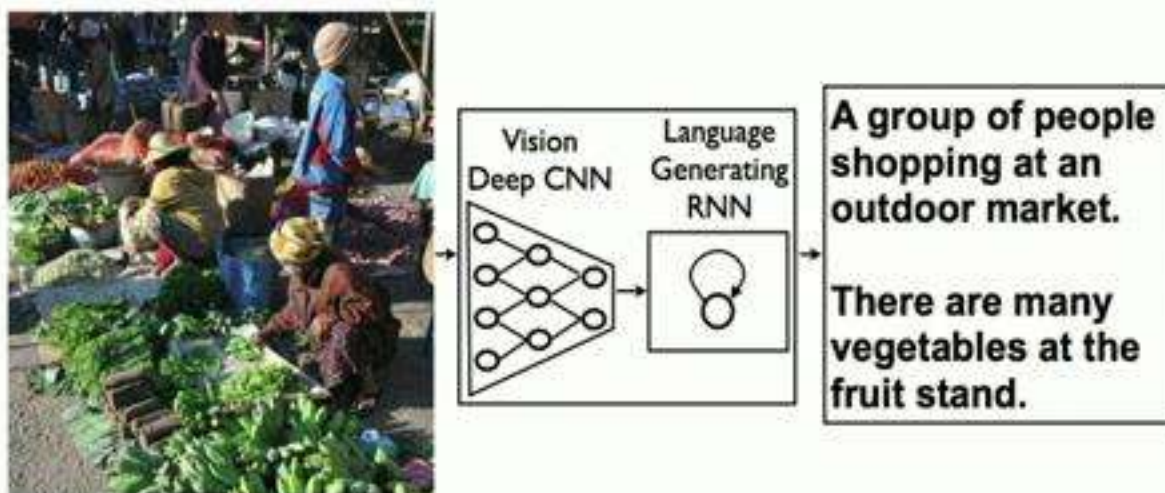
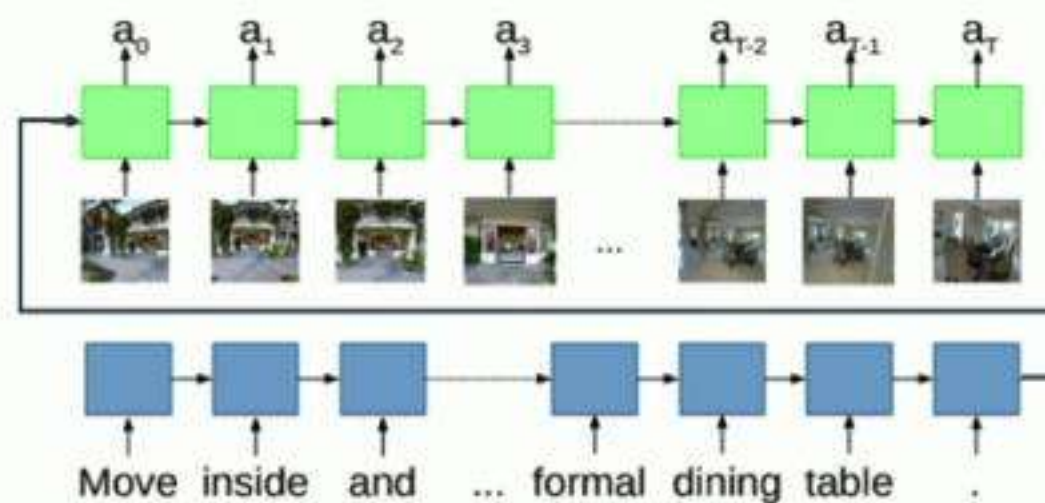
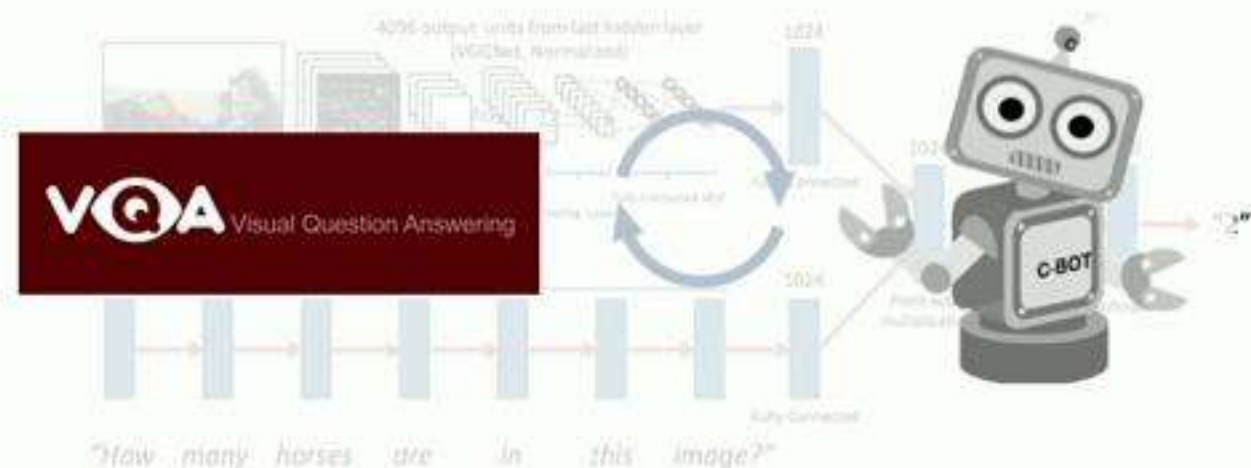


Image Captioning



Vision and Language Navigation

Vision and Language



Visual Question Answering

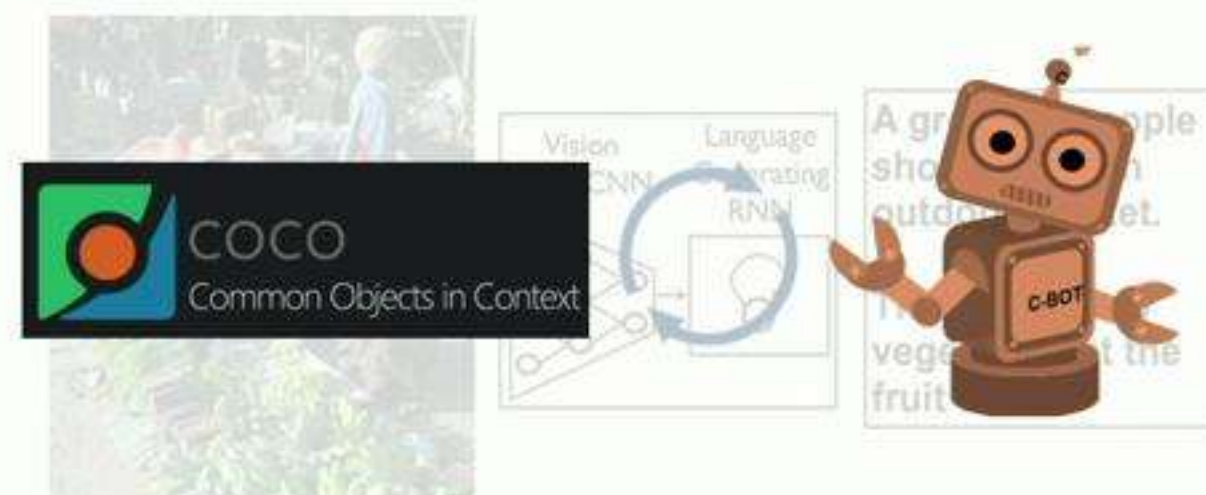
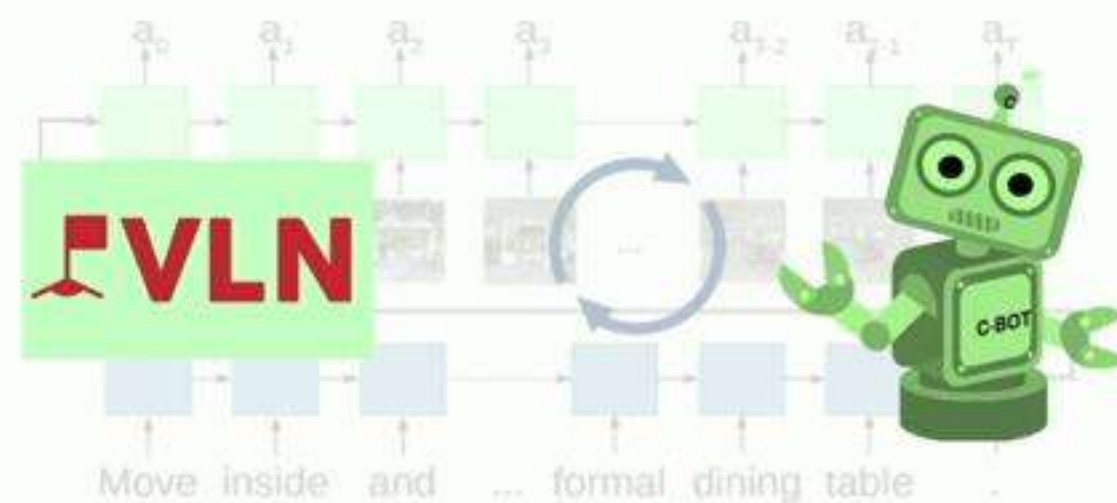


Image Captioning

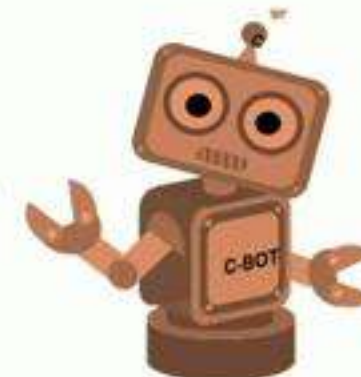
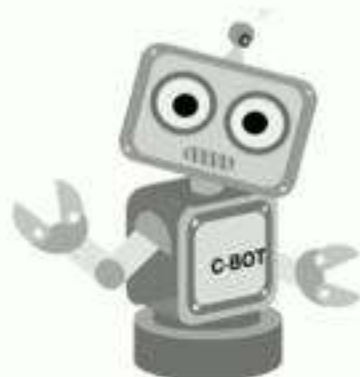


Visual Commonsense Reasoning

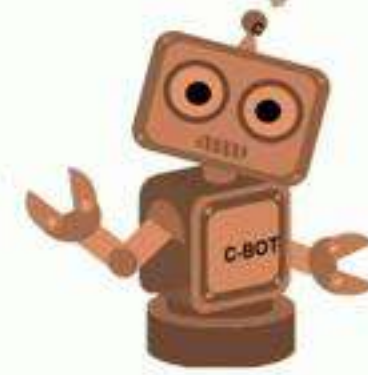
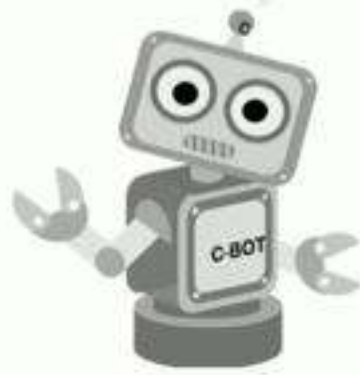


Vision and Language Navigation

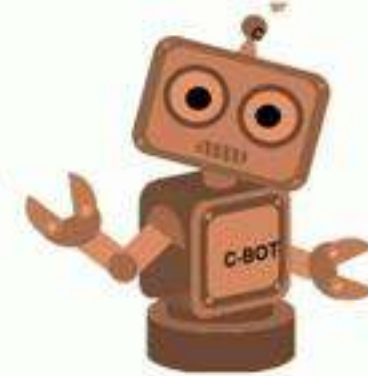
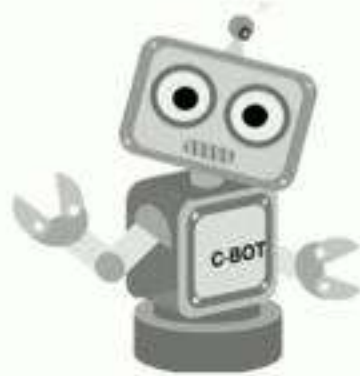
Visual Grounding



Visual Grounding

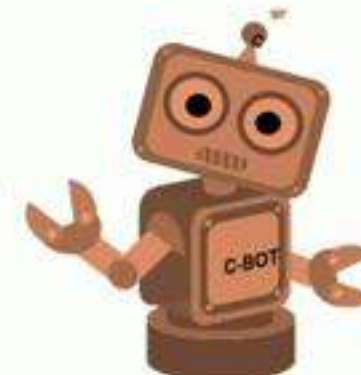
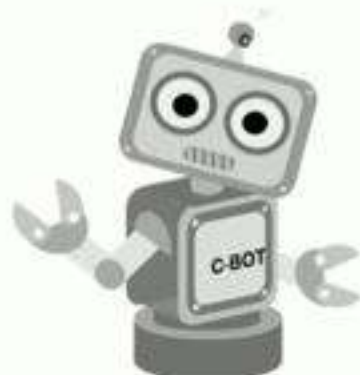


Visual Grounding



Q: What type of plant is this?

Visual Grounding

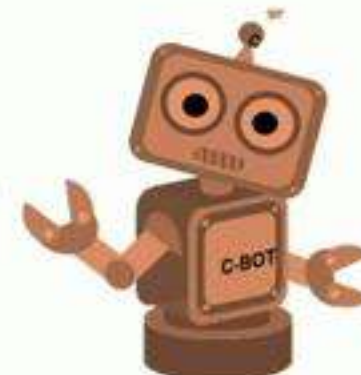
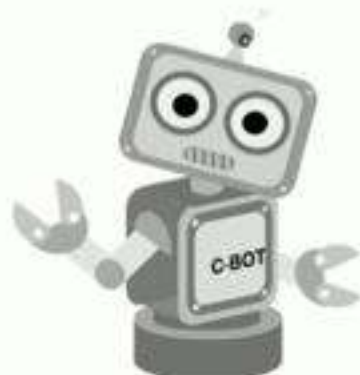


Q: *What type of plant is this?*

A: **Banana**

C: A bunch of red and yellow **flowers** on a branch.

Visual Grounding



Q: *What type of plant is this?*

A: **Banana**

C: A bunch of red and yellow **flowers** on a branch.

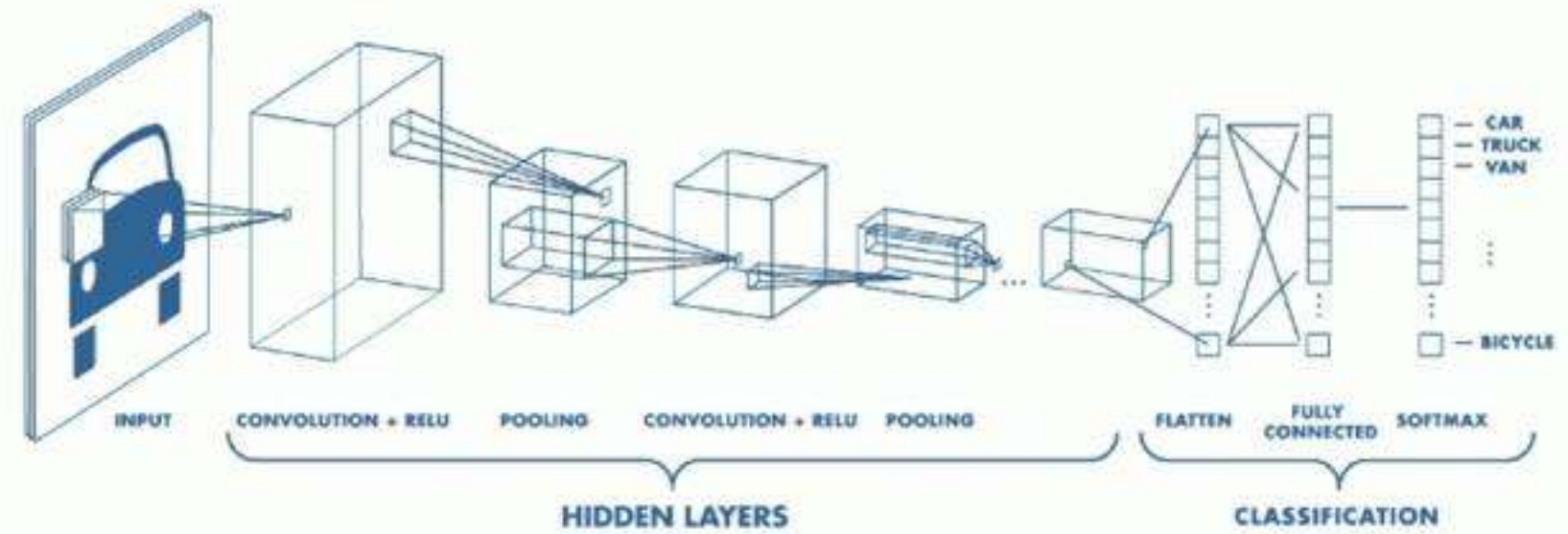
Common model for **visual grounding** and leverage them on a wide array of vision-and-language tasks

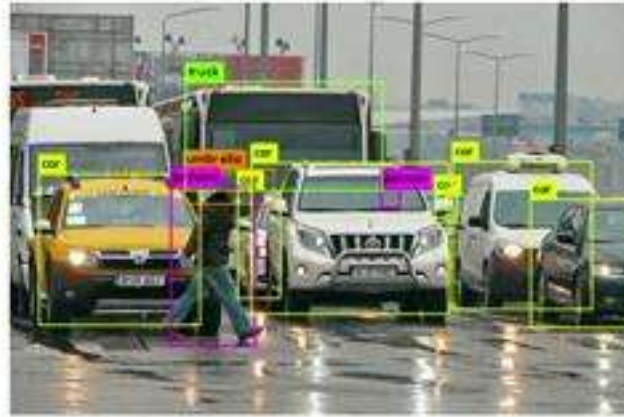
Pretrain-Transfer

Pretrain-Transfer



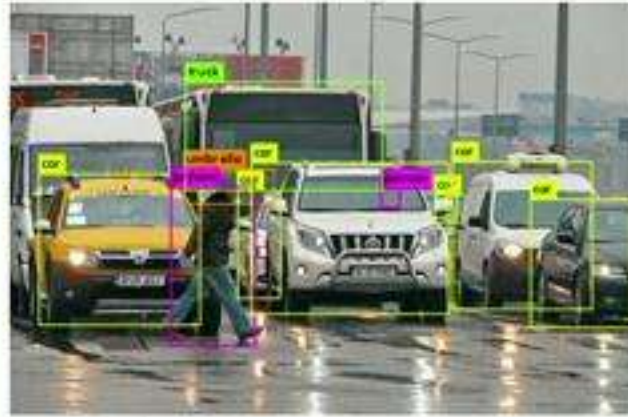
Pretrain-Transfer





Object Detection

Pretrain-Transfer

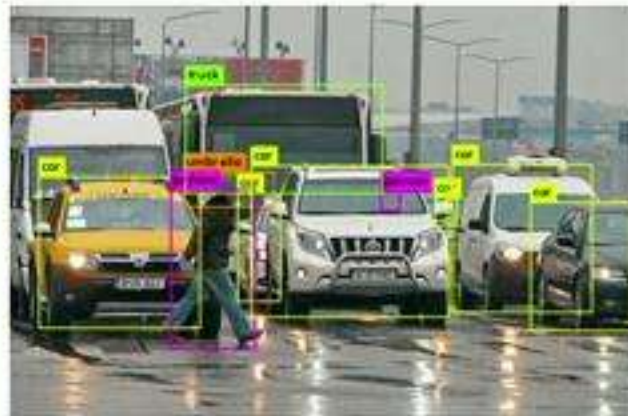


Object Detection



Semantic Segmentation

Pretrain-Transfer

The logo for ImageNet, featuring the word "IMAGENET" in a grey, sans-serif font. The letter "A" is replaced by a small, stylized icon of a green square, an orange square, and a red square arranged in a triangular pattern.

Object Detection

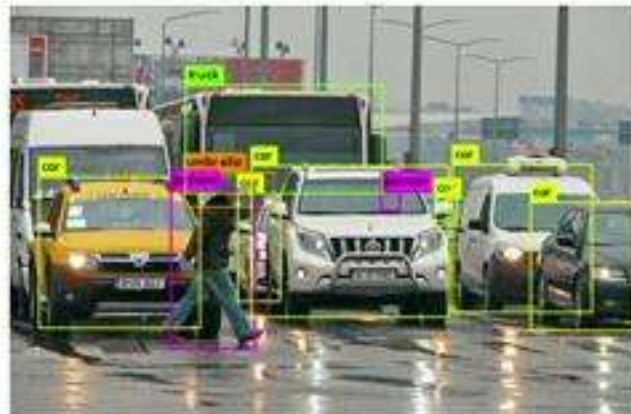


Semantic Segmentation



Pose Estimation

Pretrain-Transfer



Object Detection



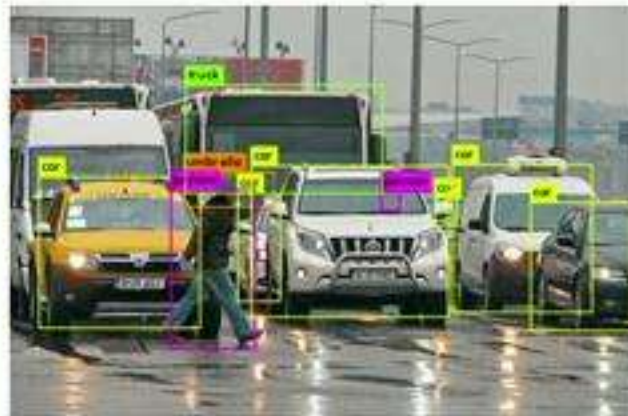
Semantic Segmentation



Pose Estimation



Pretrain-Transfer



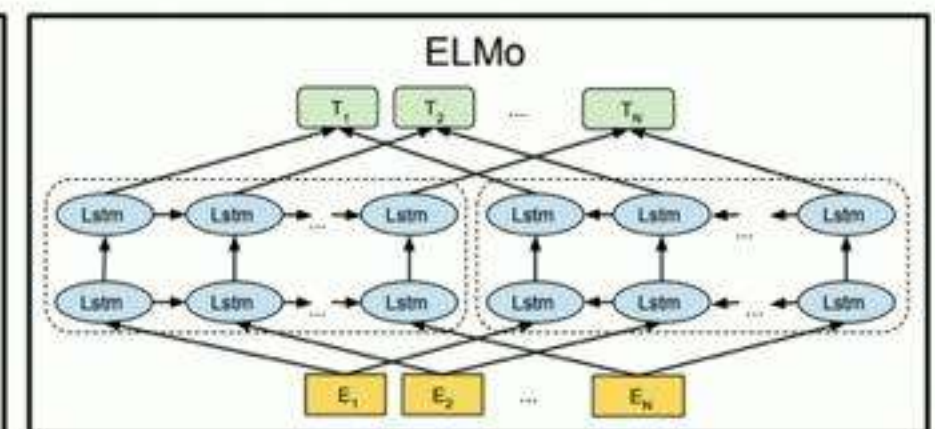
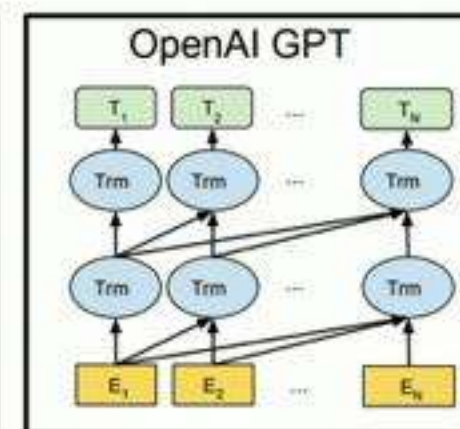
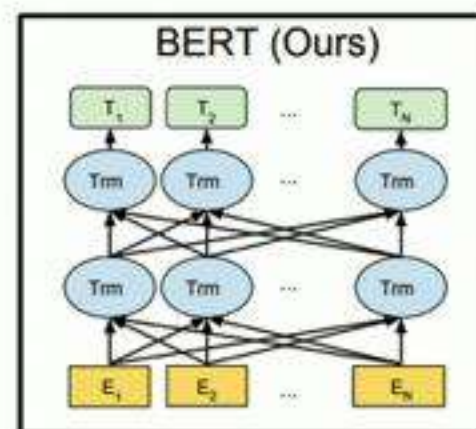
Object Detection



Semantic Segmentation



Pose Estimation



Pretrain-Transfer



Object Detection



Semantic Segmentation



Pose Estimation



Passage Segment

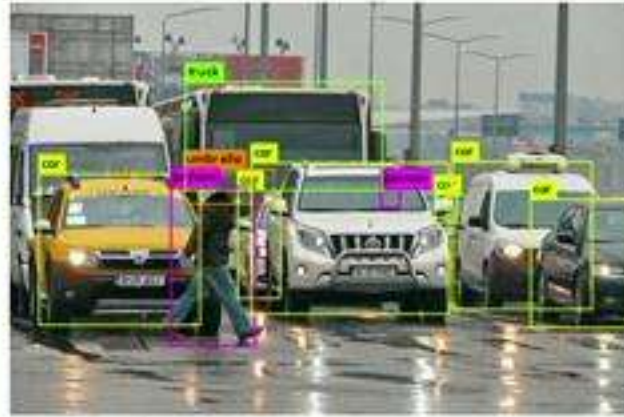
...The European Parliament and the Council of the European Union have powers of amendment and veto during the legislative process...

Question

Which governing bodies have veto power?

Question Answering

Pretrain-Transfer



Object Detection



Semantic Segmentation



Pose Estimation



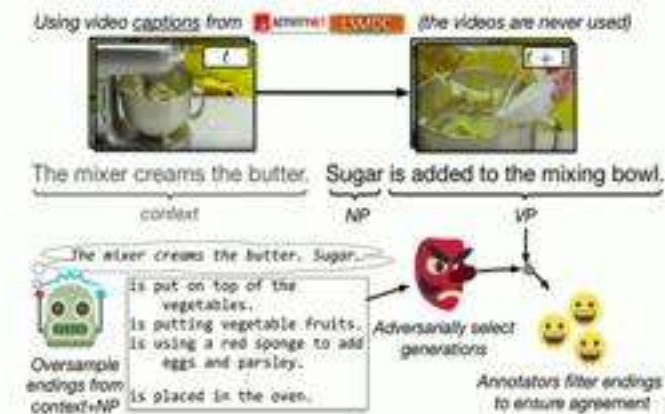
Passage Segment

...The European Parliament and the Council of the European Union have powers of amendment and veto during the legislative process...

Question

Which governing bodies have veto power?

Question Answering



Commonsense Inference

Pretrain-Transfer

IMAGENET



Object Detection



Semantic Segmentation



Pose Estimation

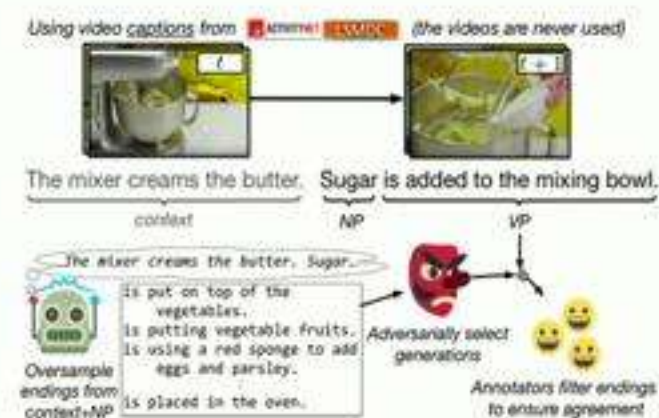
Passage Segment

...The European Parliament and the Council of the European Union have powers of amendment and veto during the legislative process...

Question

Which governing bodies have veto power?

Question Answering



Commonsense Inference

Sentiment	Tweets
Negative	@united is the worst. Nonrefundable First class tickets? Oh because when you select Global/FC their system auto selects economy w/upgrade. @united I will not be flying you again
Neutral	@VirginAmerica my drivers license is expired by a little over a month. Can I fly Friday morning using my expired license? @VirginAmerica any plans to start flying direct from DAL to LAS?
Positive	@VirginAmerica done! Thank you for the quick response, apparently faster than sitting on hold ;) @united I appreciate your efforts getting me home!

Sentiment Analysis

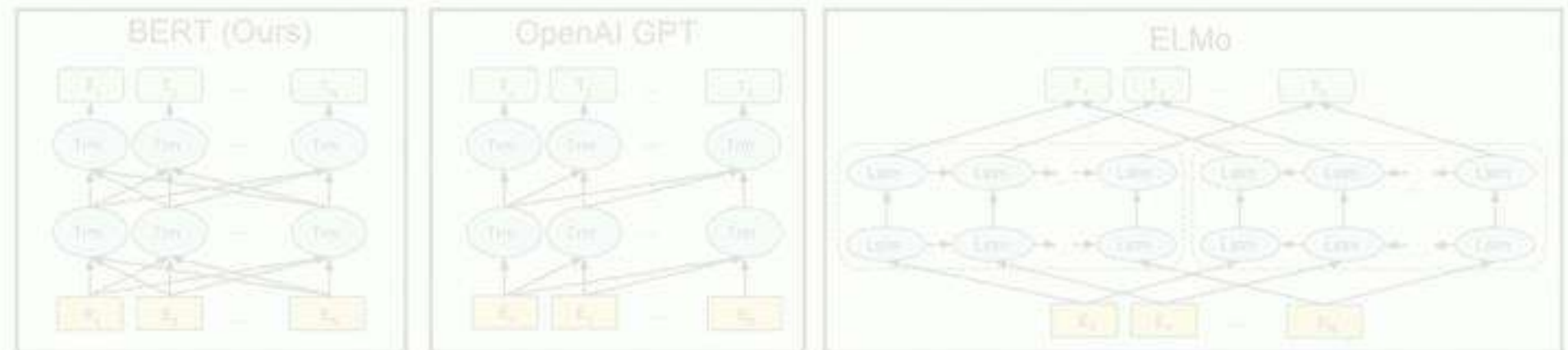
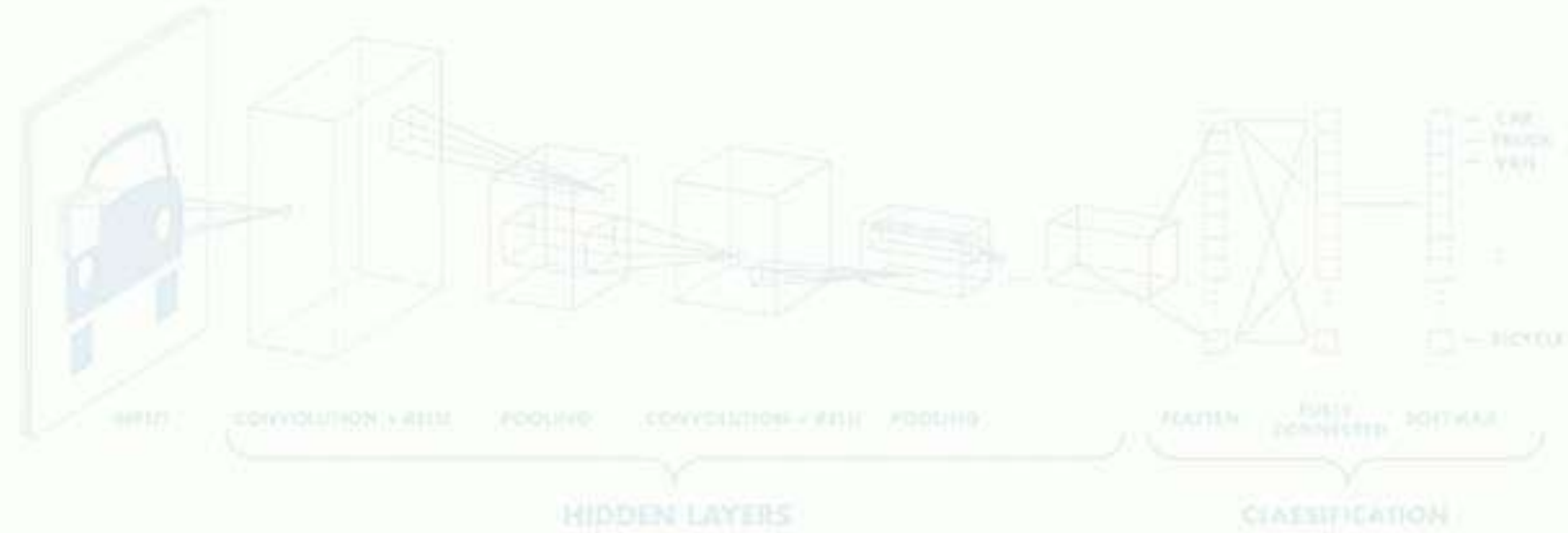
Pretrain-Transfer



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset



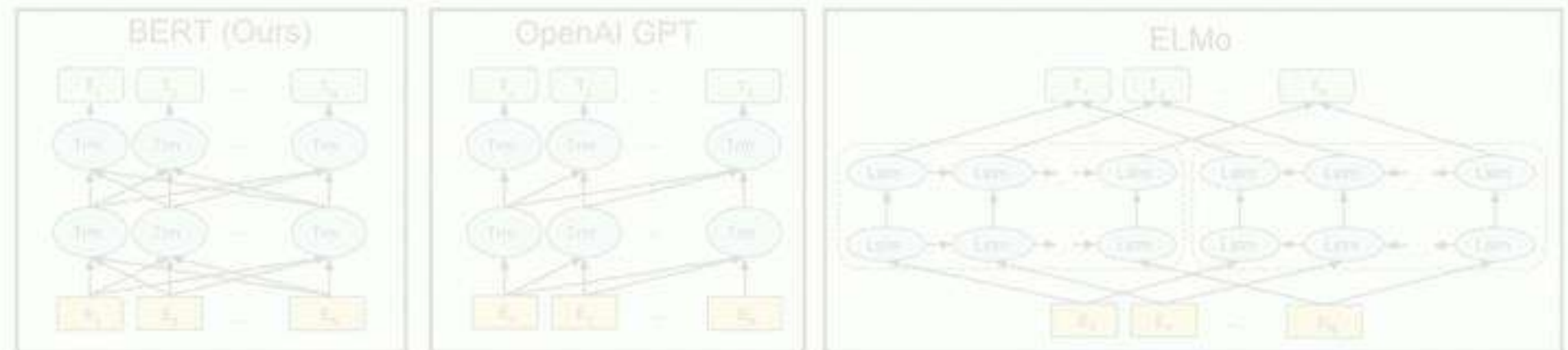
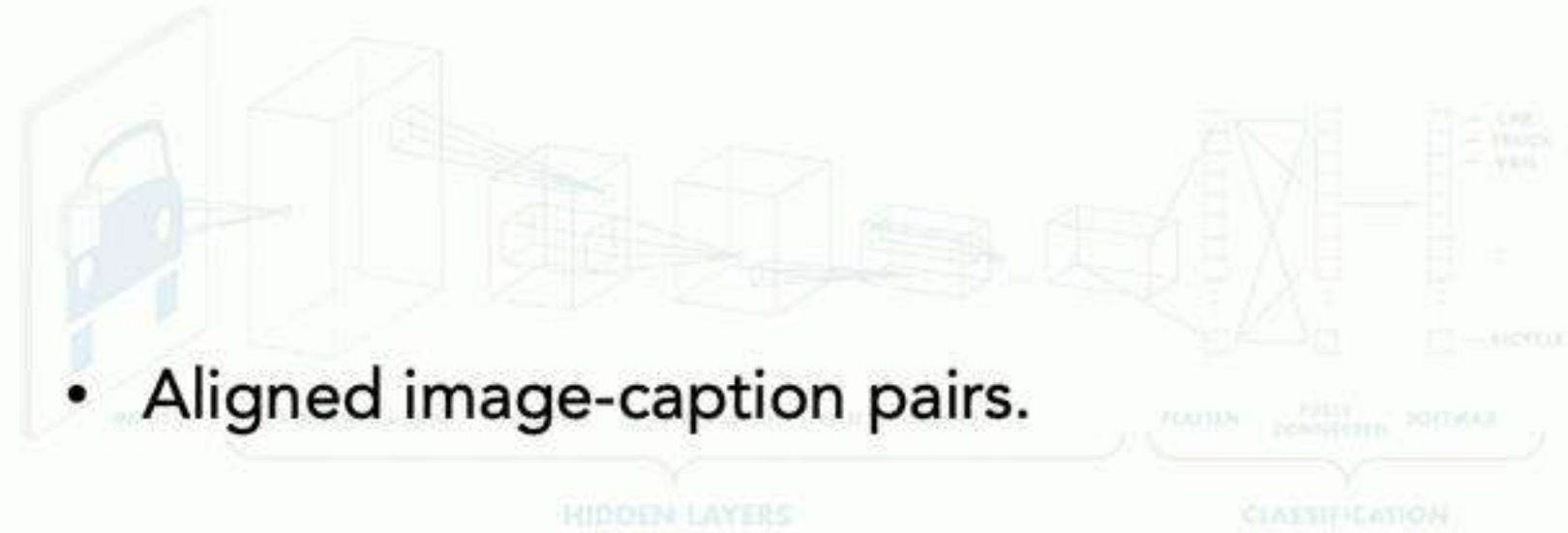
Pretrain-Transfer



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset



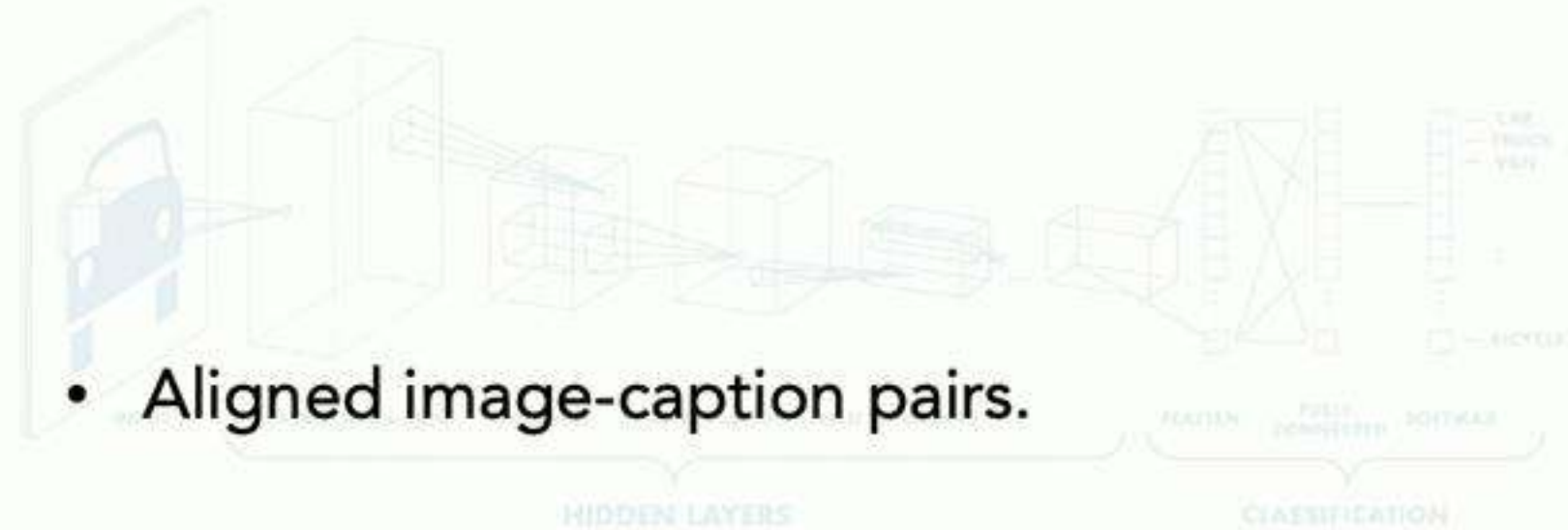
Pretrain-Transfer



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

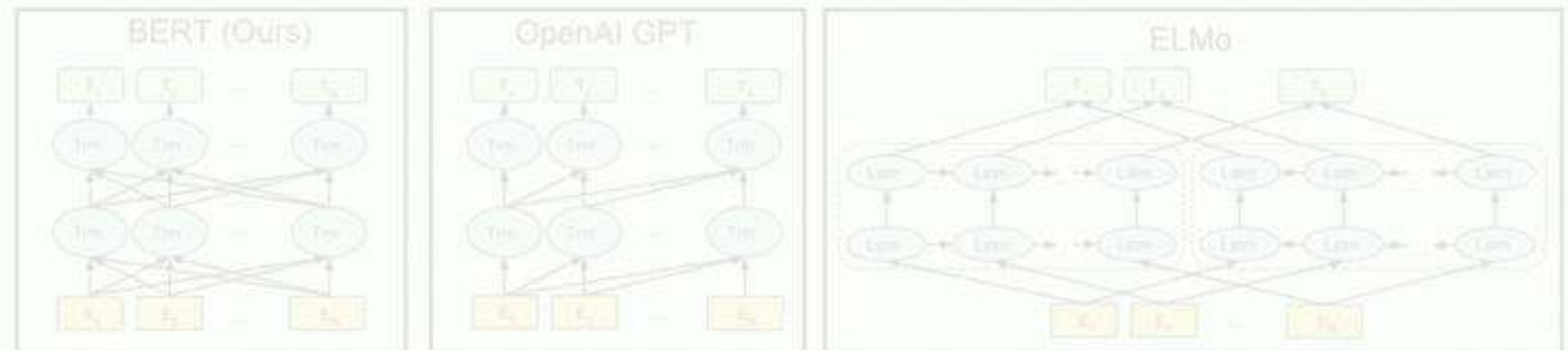
Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset



- Aligned image-caption pairs.

- Large Scale. (3.3 million)



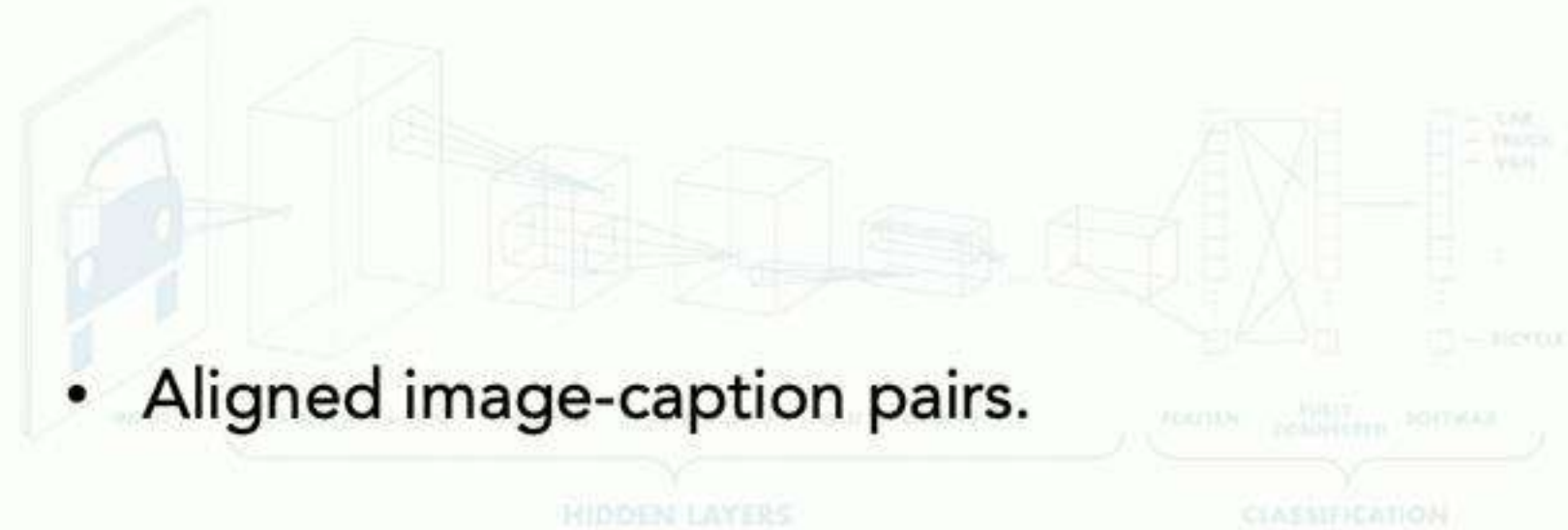
Pretrain-Transfer



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

Conceptual Captions: pop artist performs at the festival in a city.

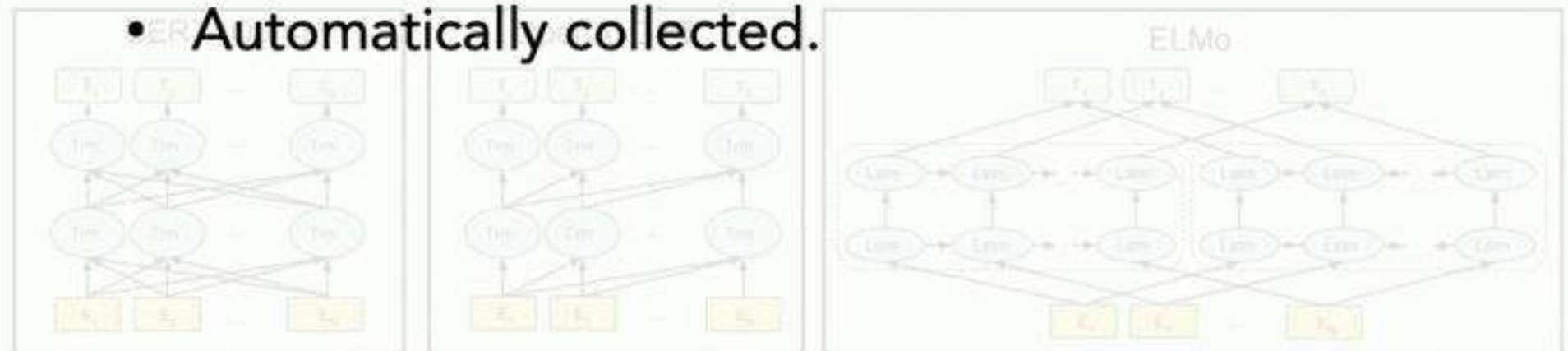
Conceptual Caption Dataset



- Aligned image-caption pairs.

- Large Scale. (3.3 million)

- Automatically collected.



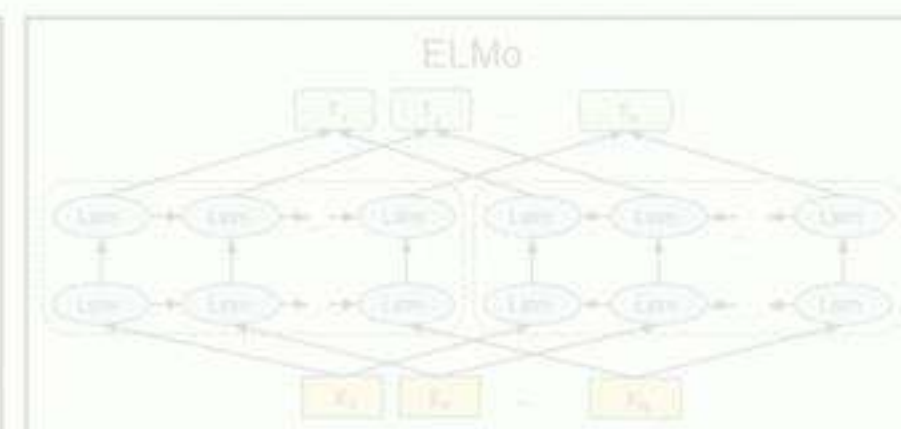
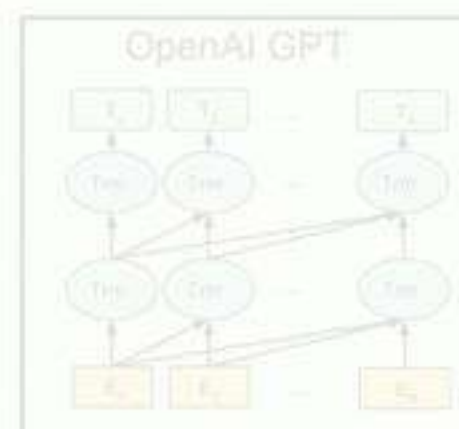
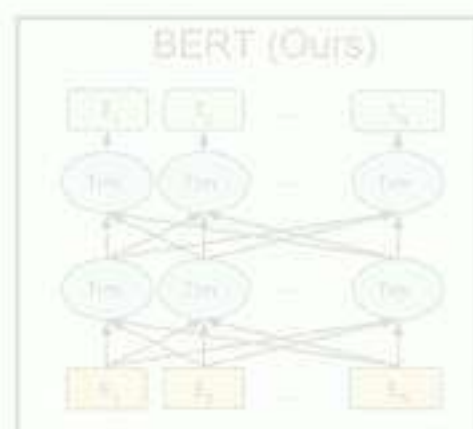
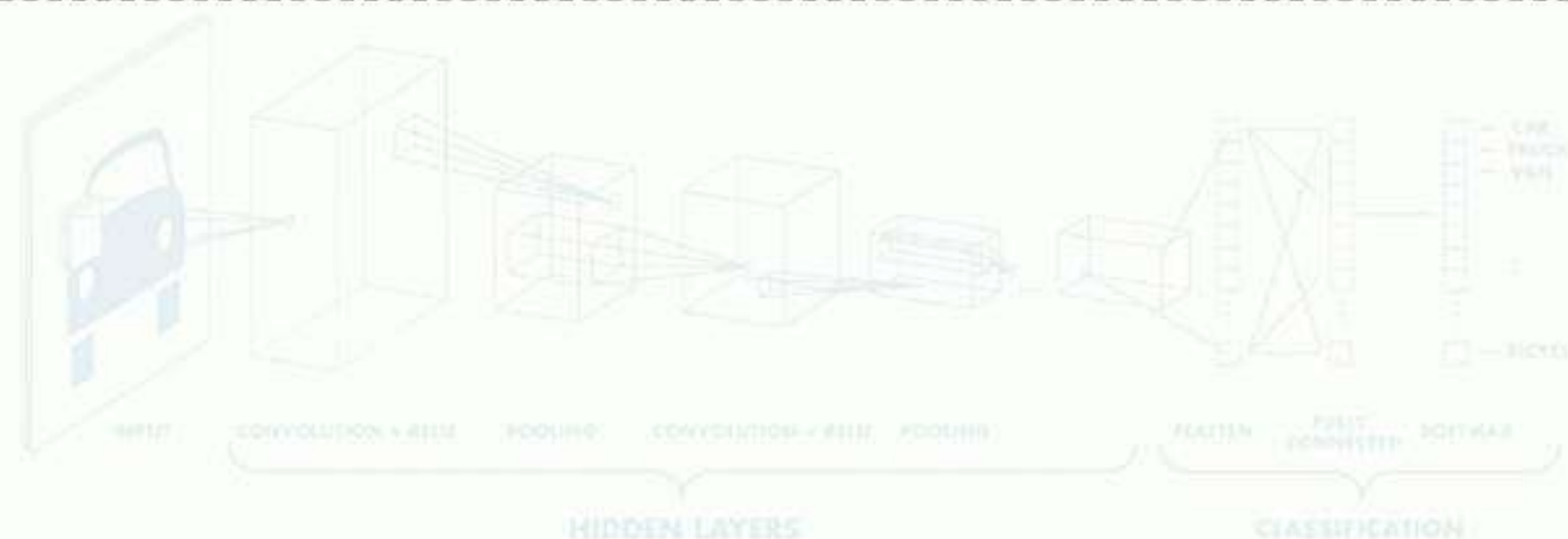
BERT



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset



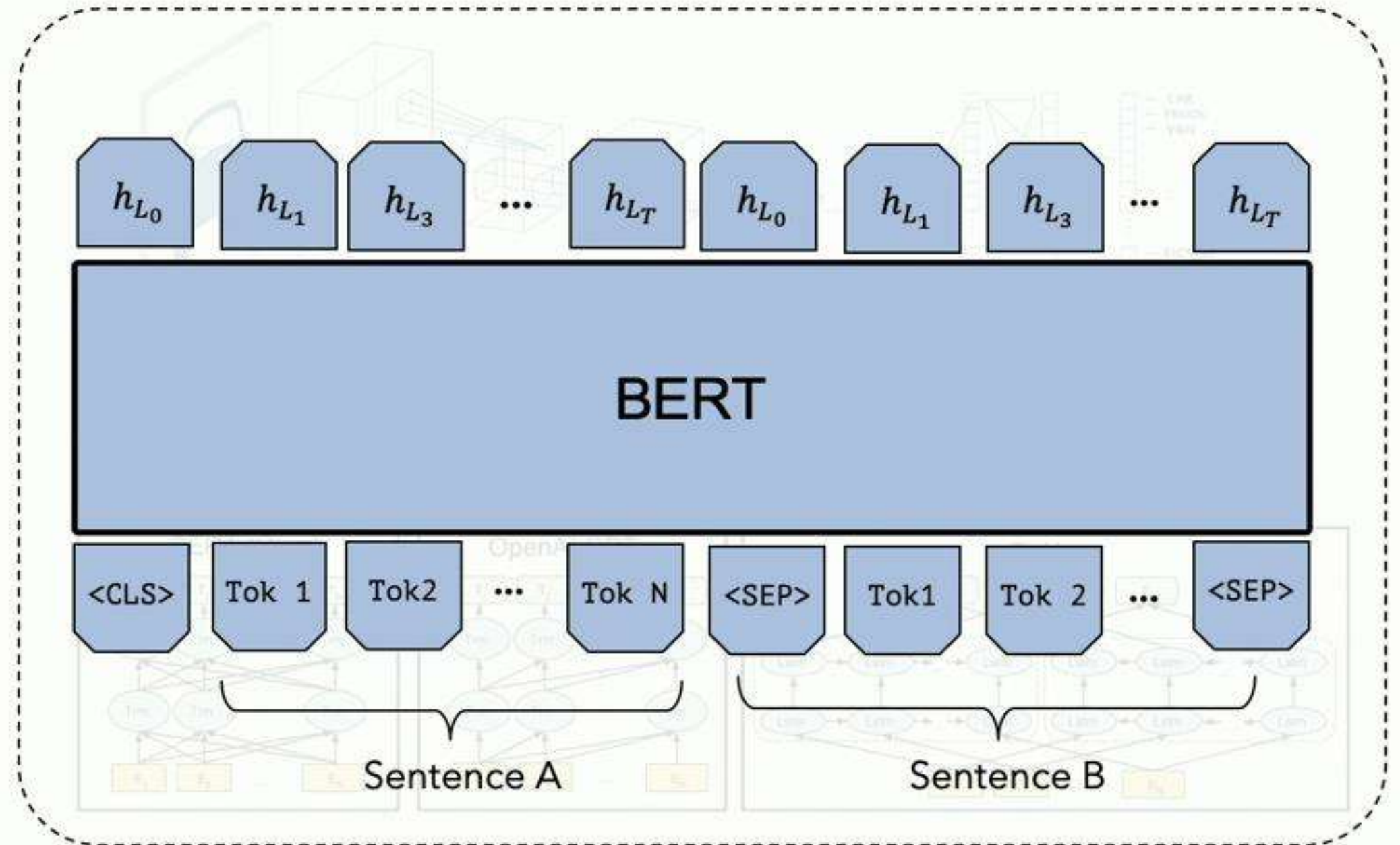
BERT



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset

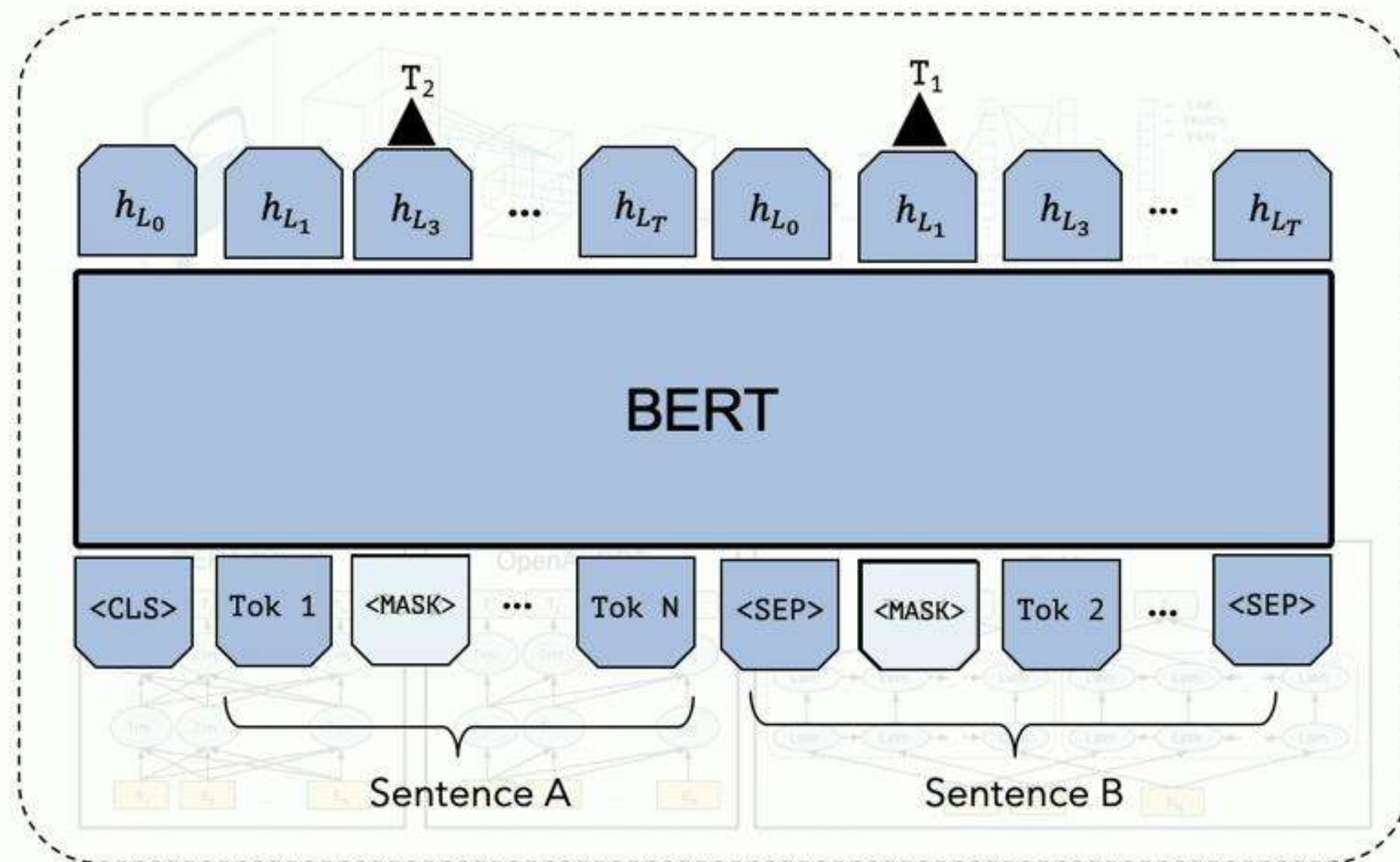




Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset



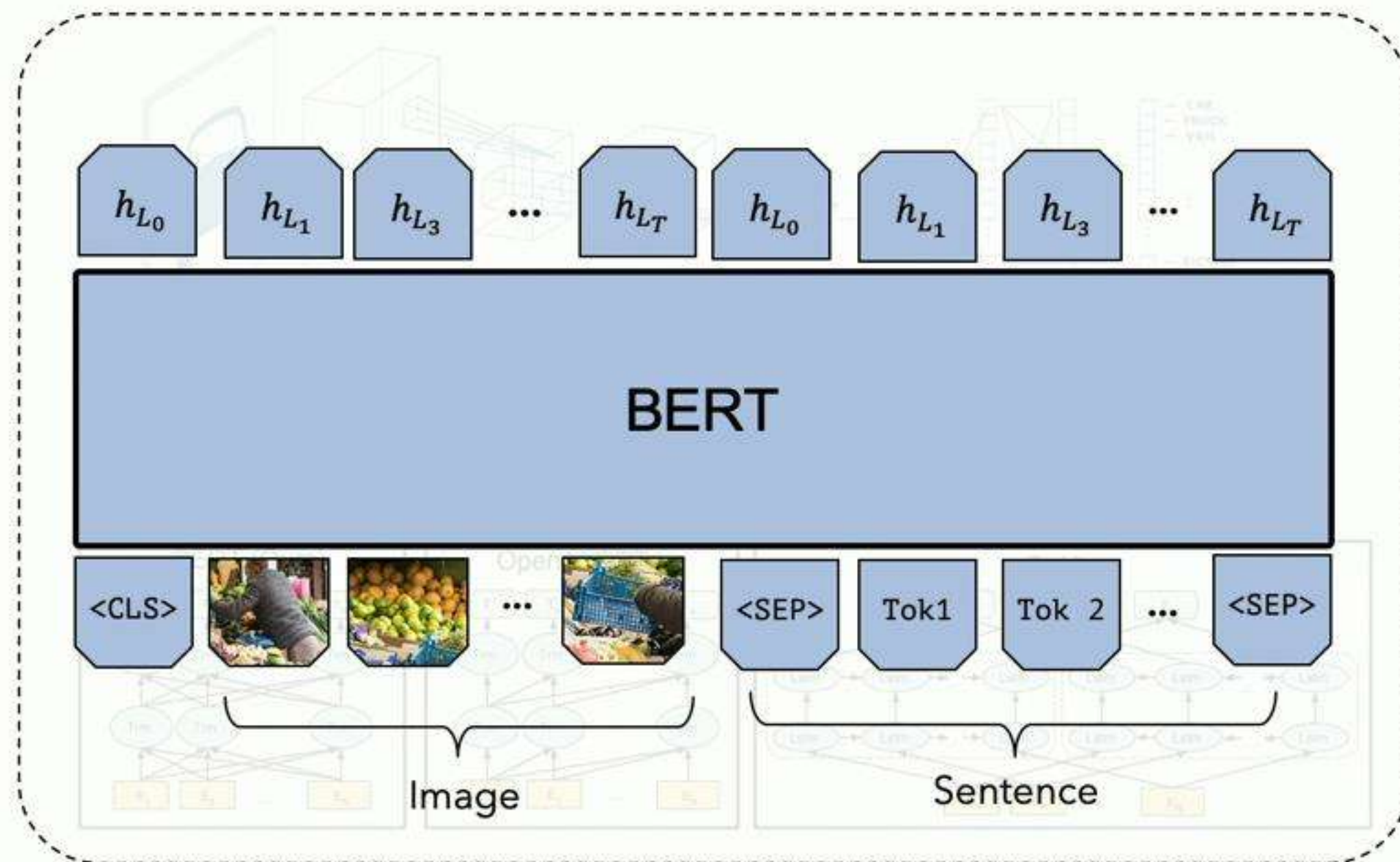
Single-Stream model



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset



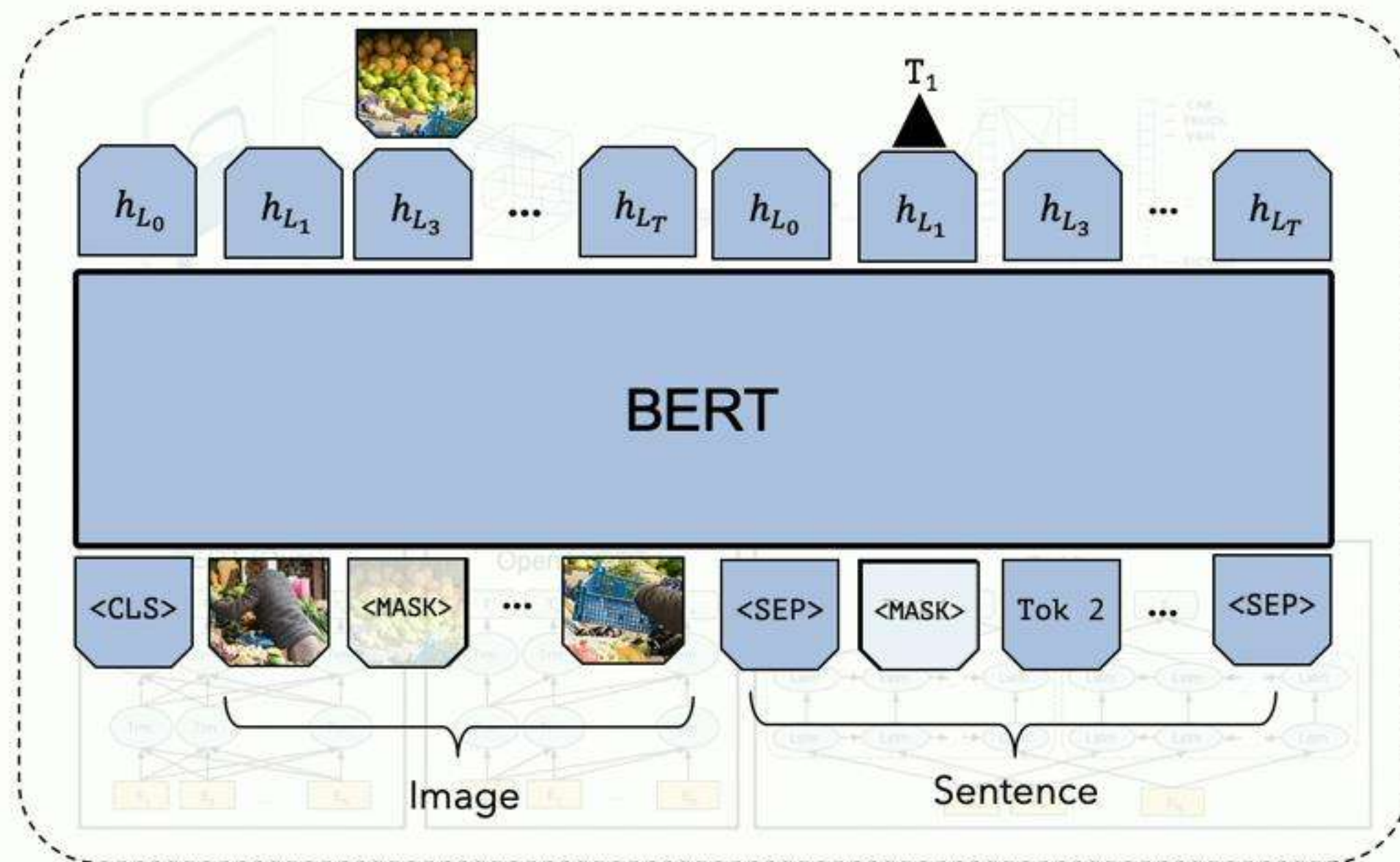
Single-Stream model



Alt-text: Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

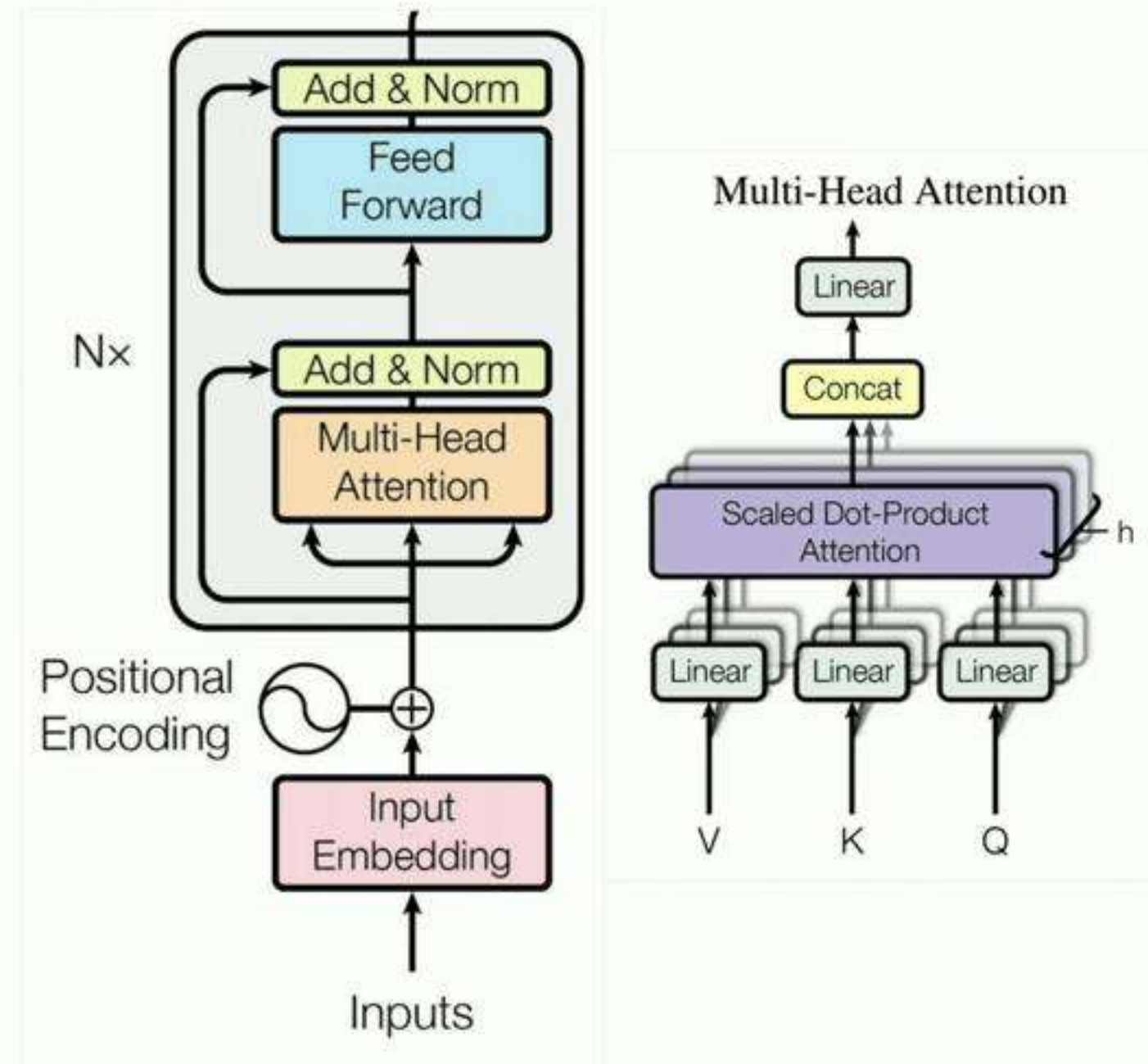
Conceptual Captions: pop artist performs at the festival in a city.

Conceptual Caption Dataset



Transformer

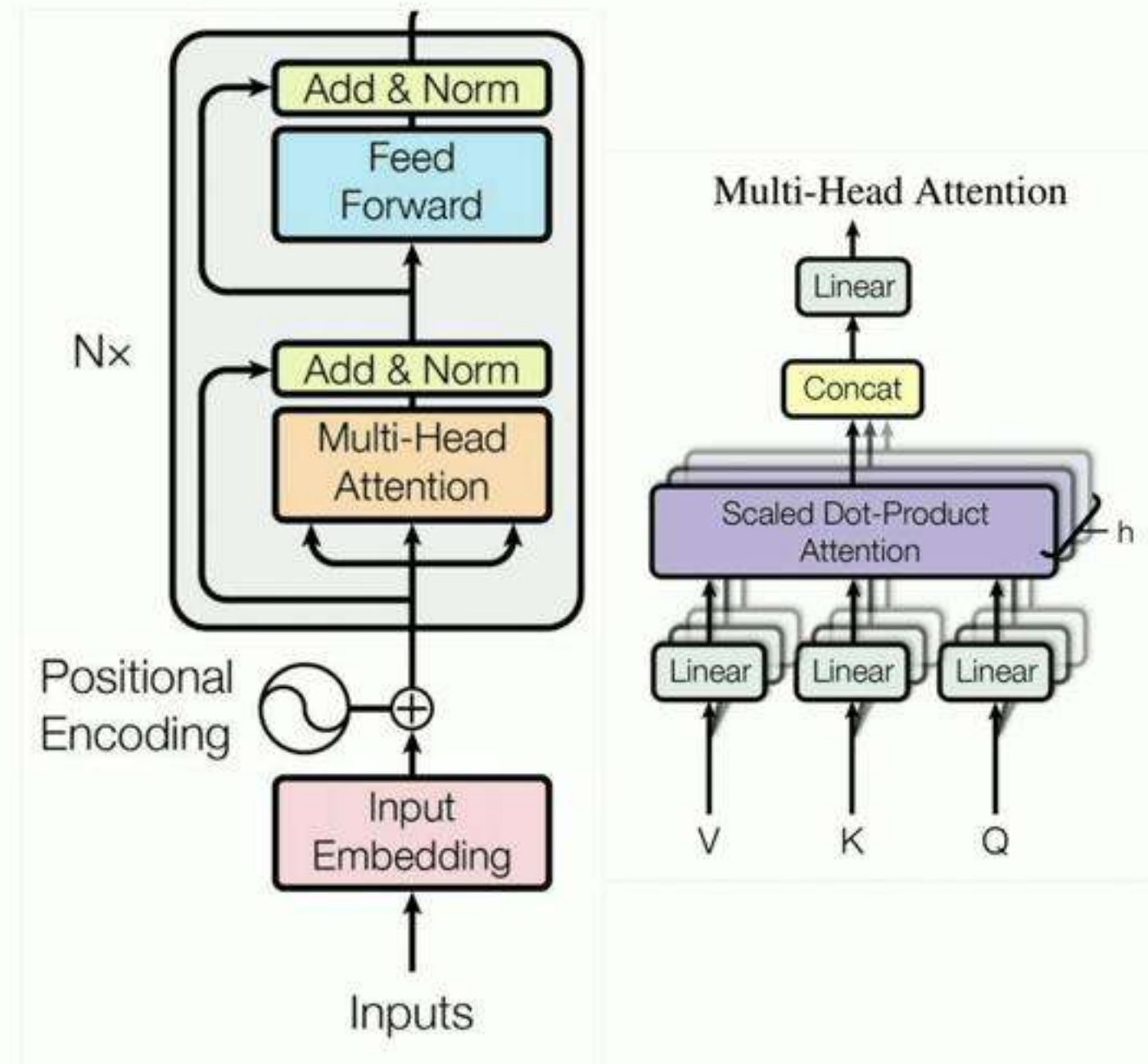
Transformer encoder



Transformer

Transformer encoder

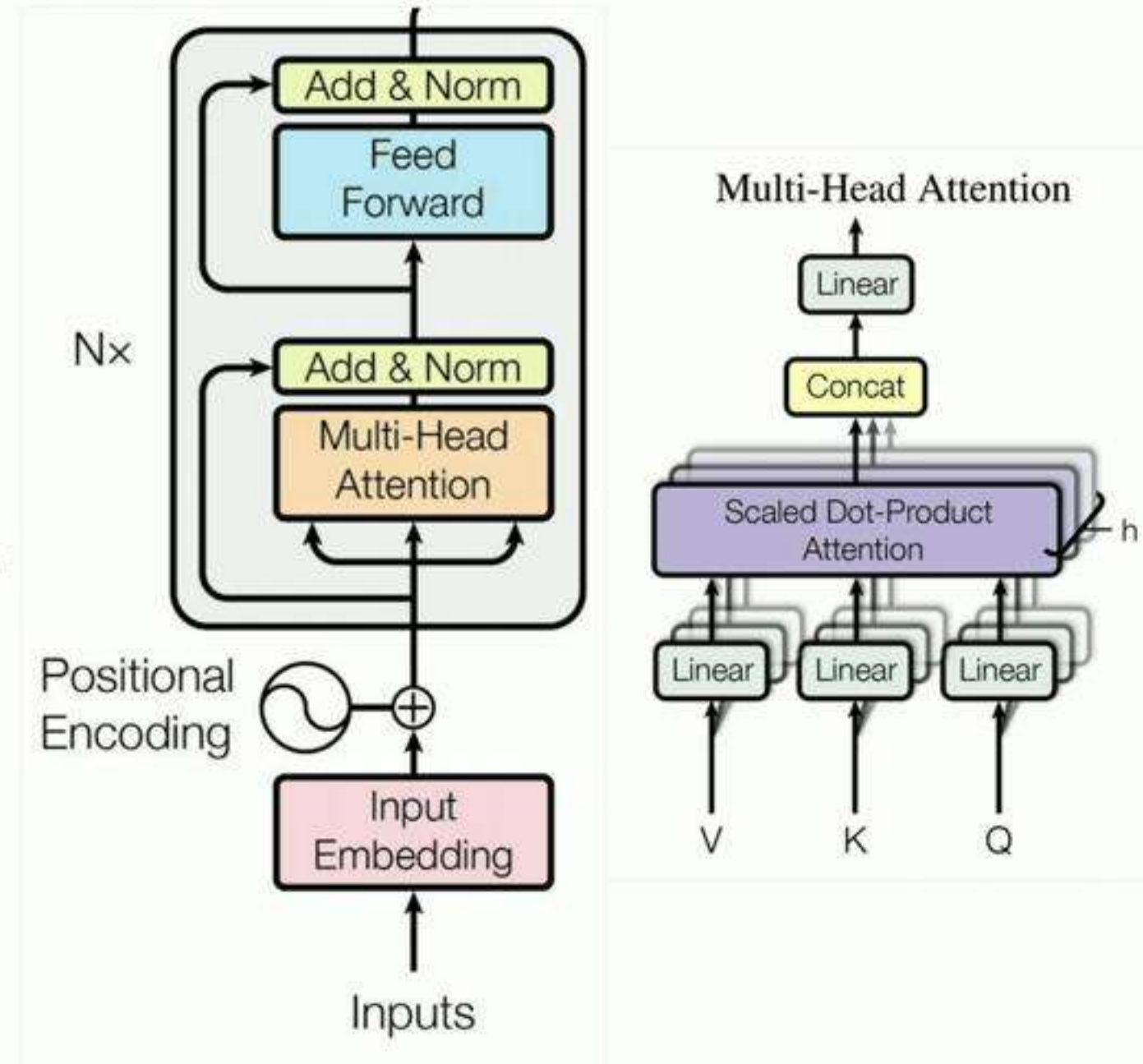
- Multi-headed self attention
 - Model context



Transformer

Transformer encoder

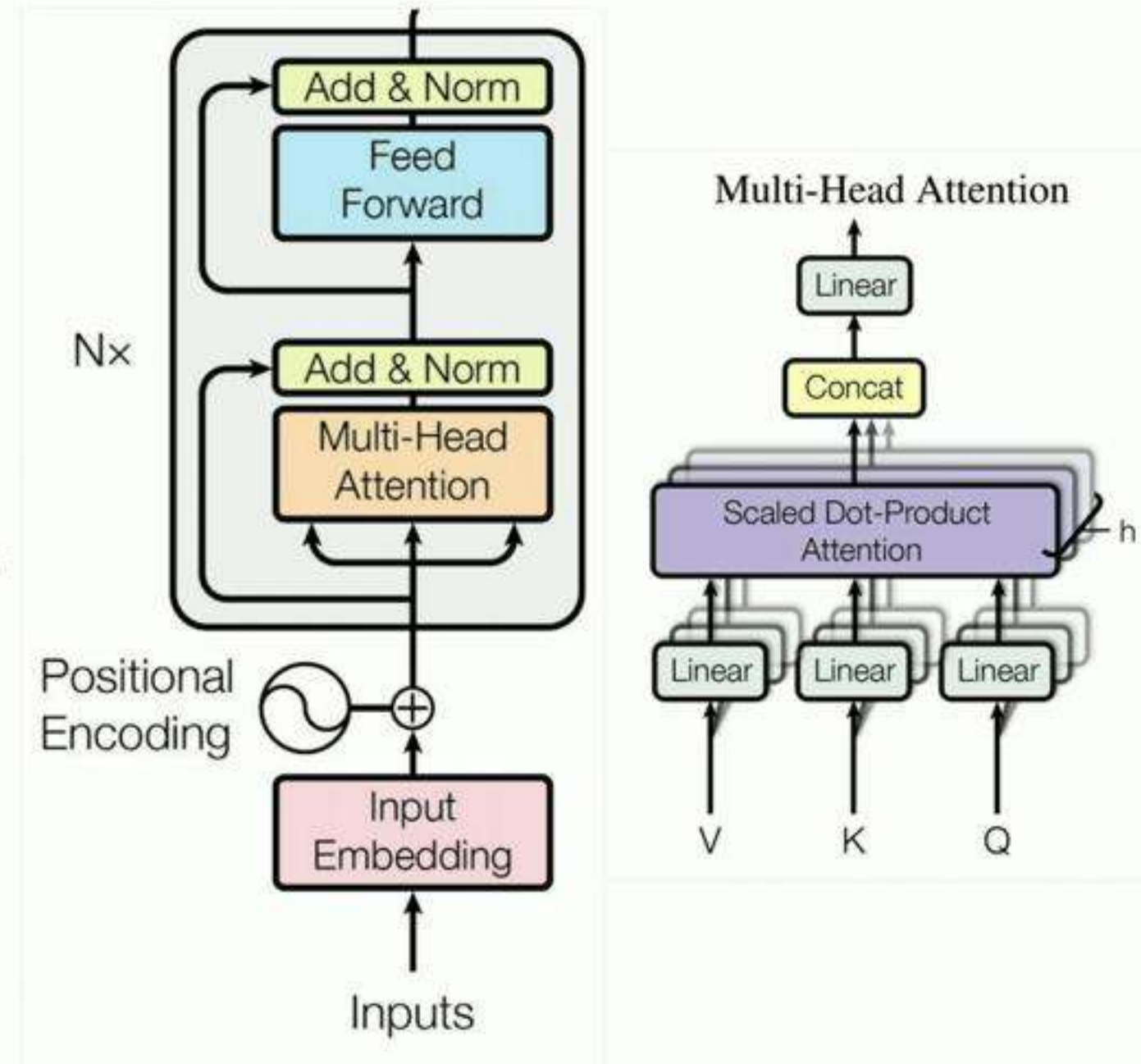
- Multi-headed self attention
 - Model context
- Feed-forward layers
 - Computes non-linear hierarchical features



Transformer

Transformer encoder

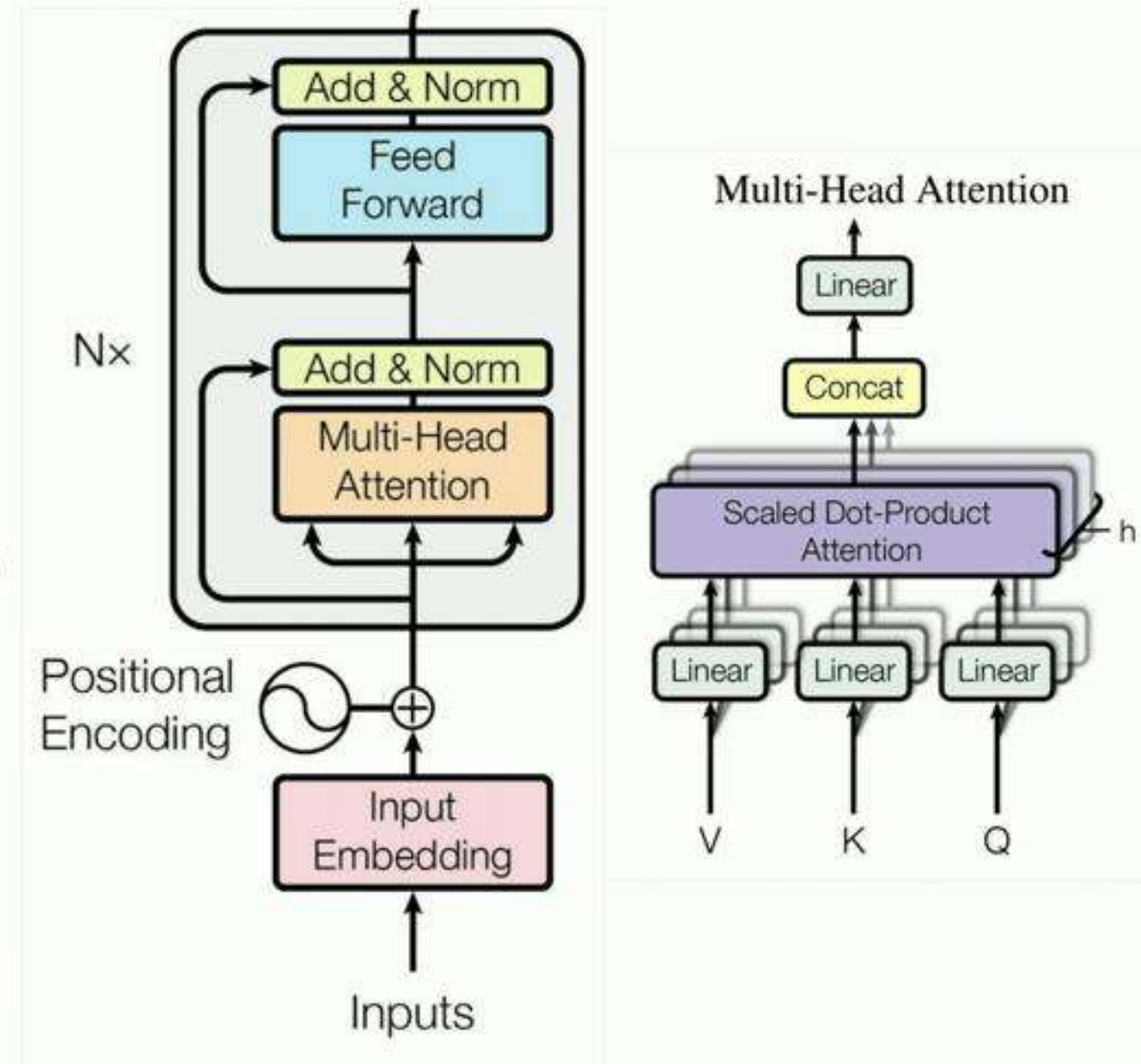
- Multi-headed self attention
 - Model context
- Feed-forward layers
 - Computes non-linear hierarchical features
- Layer norm and residuals
 - Makes training deep networks healthy



Transformer

Transformer encoder

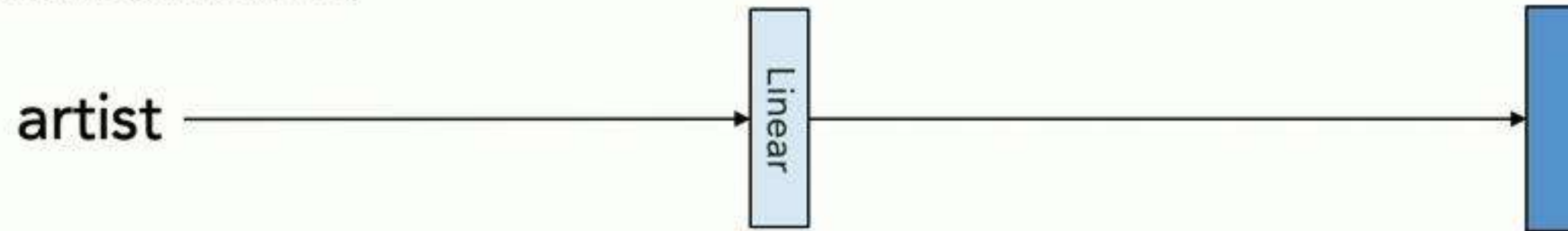
- Multi-headed self attention
 - Model context
- Feed-forward layers
 - Computes non-linear hierarchical features
- Layer norm and residuals
 - Makes training deep networks healthy
- Positional embeddings
 - Allows model to learn relative positioning



Problem: Different modalities may requires different level of abstractions.

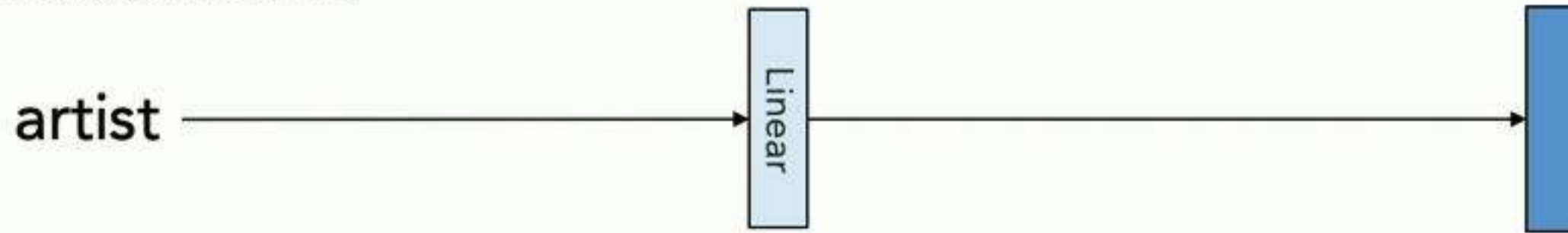
Problem: Different modalities may requires different level of abstractions.

- Linguistic stream:

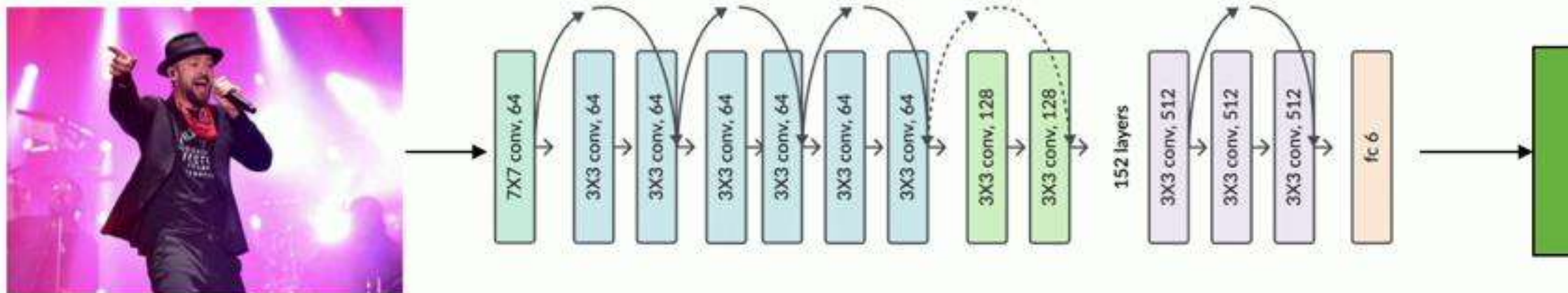


Problem: Different modalities may requires different level of abstractions.

- Linguistic stream:

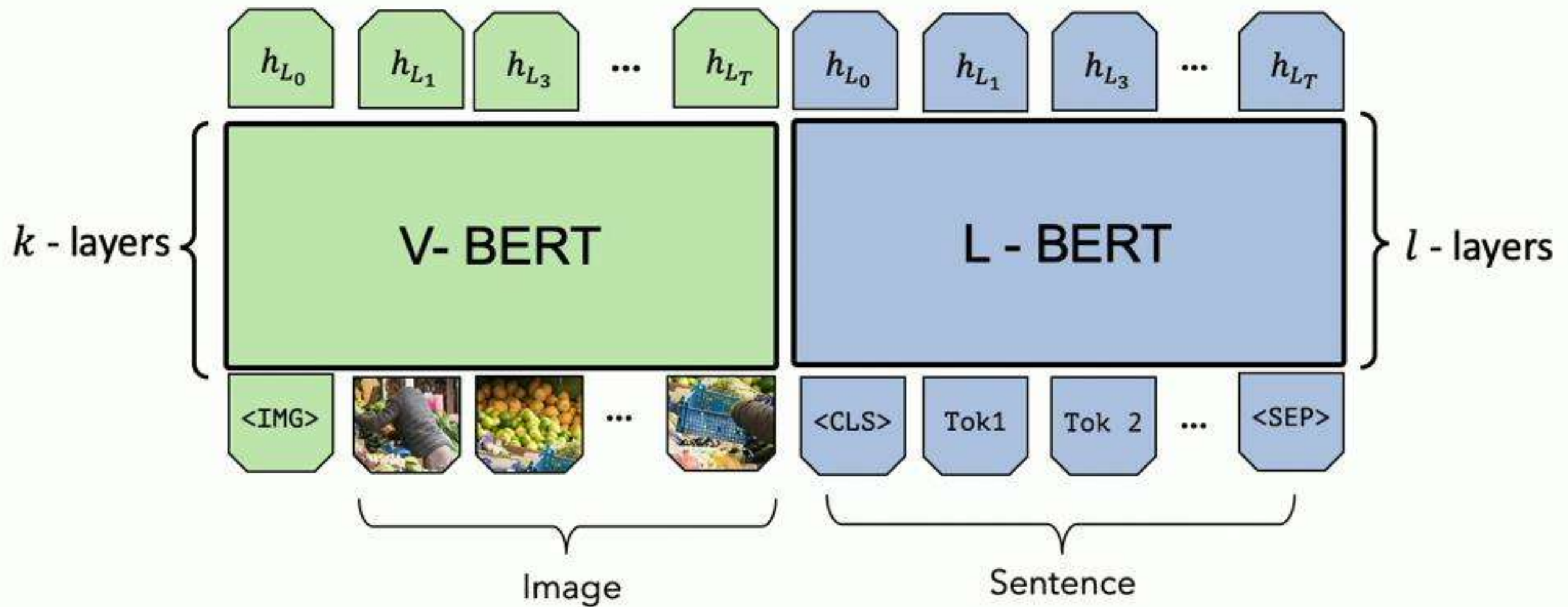


- Visual stream:



Solution: two-stream model which process visual and linguistic separately.

Solution: two-stream model which process visual and linguistic separately.



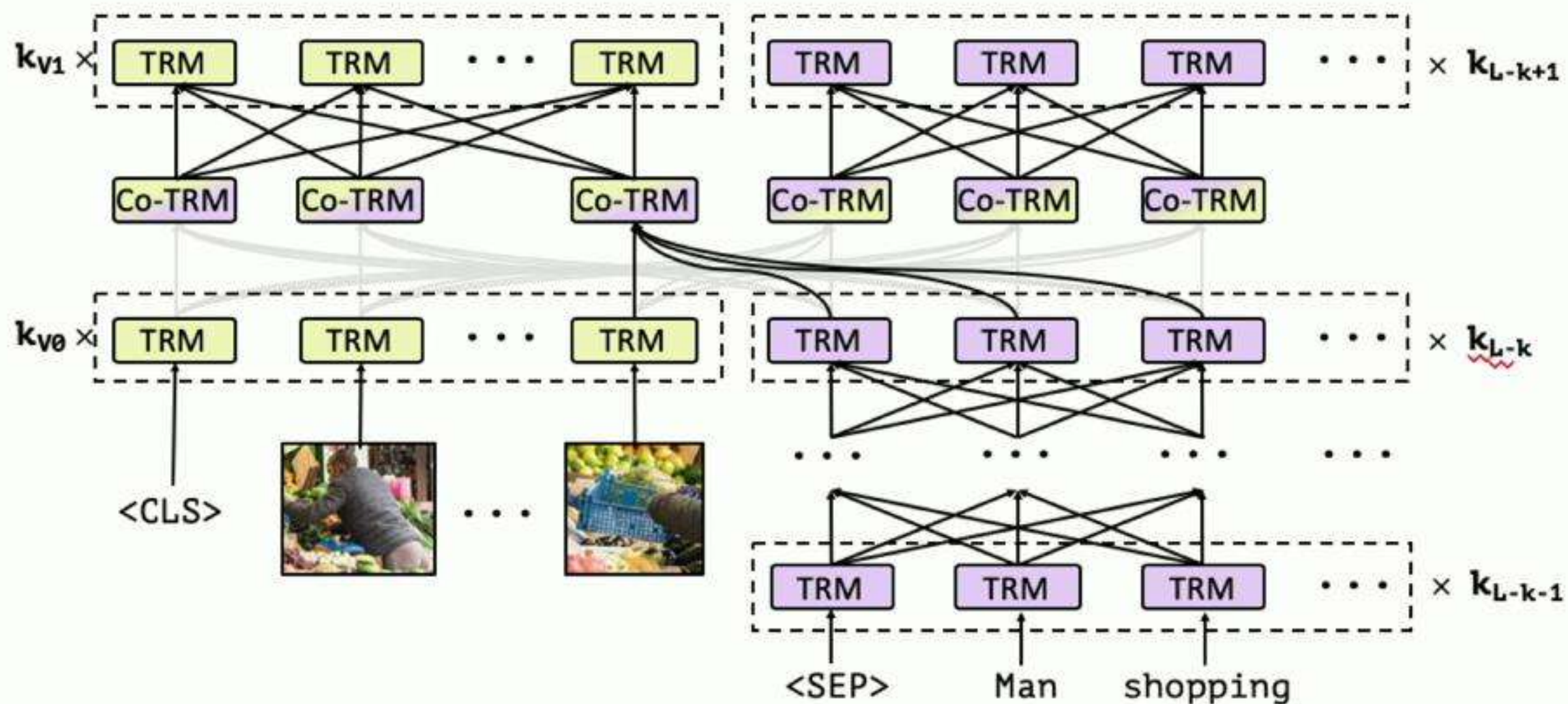
Problem: how to fuse two different modality?

Problem: how to fuse two different modality?

Solution: use **co-attention** [Lu.et.al 2016] to fuse information between different source.

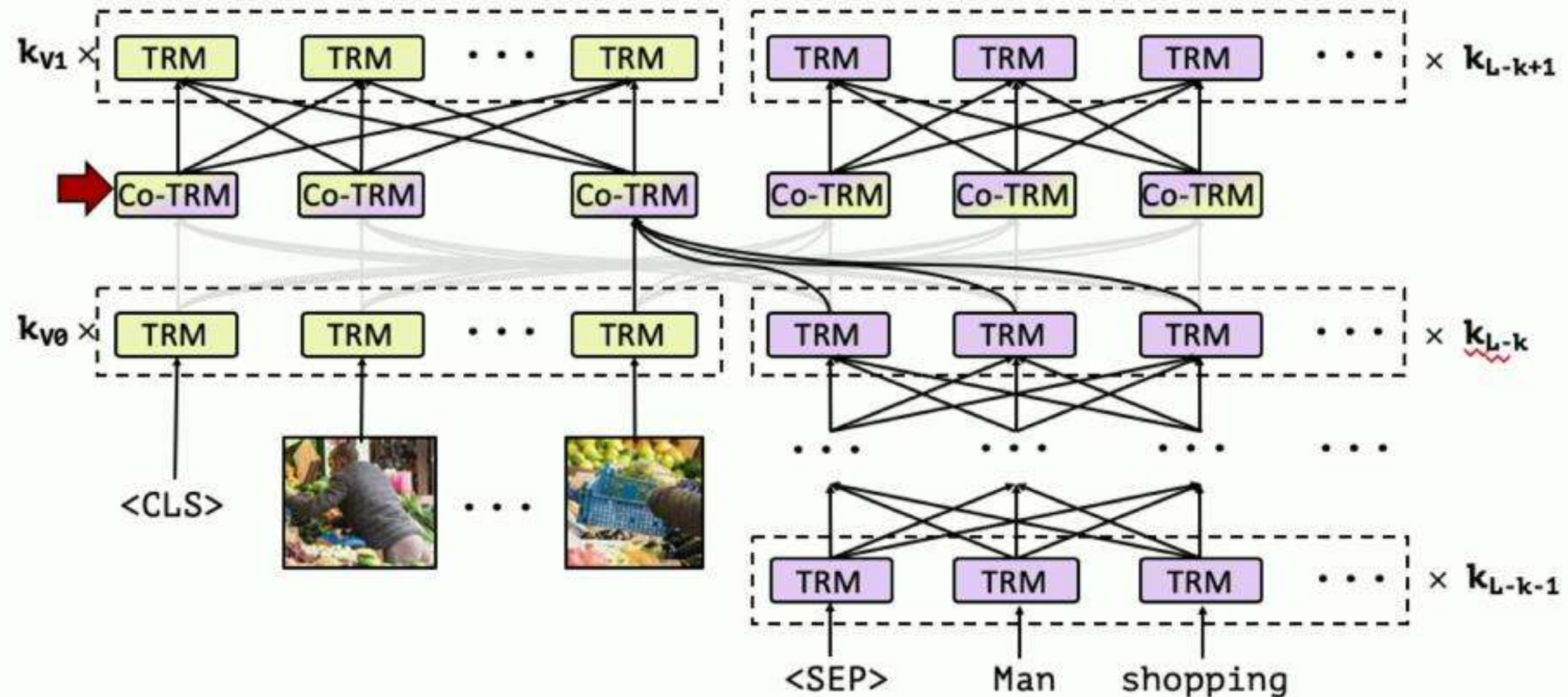
Problem: how to fuse two different modality?

Solution: use **co-attention** [Lu.et.al 2016] to fuse information between different source.

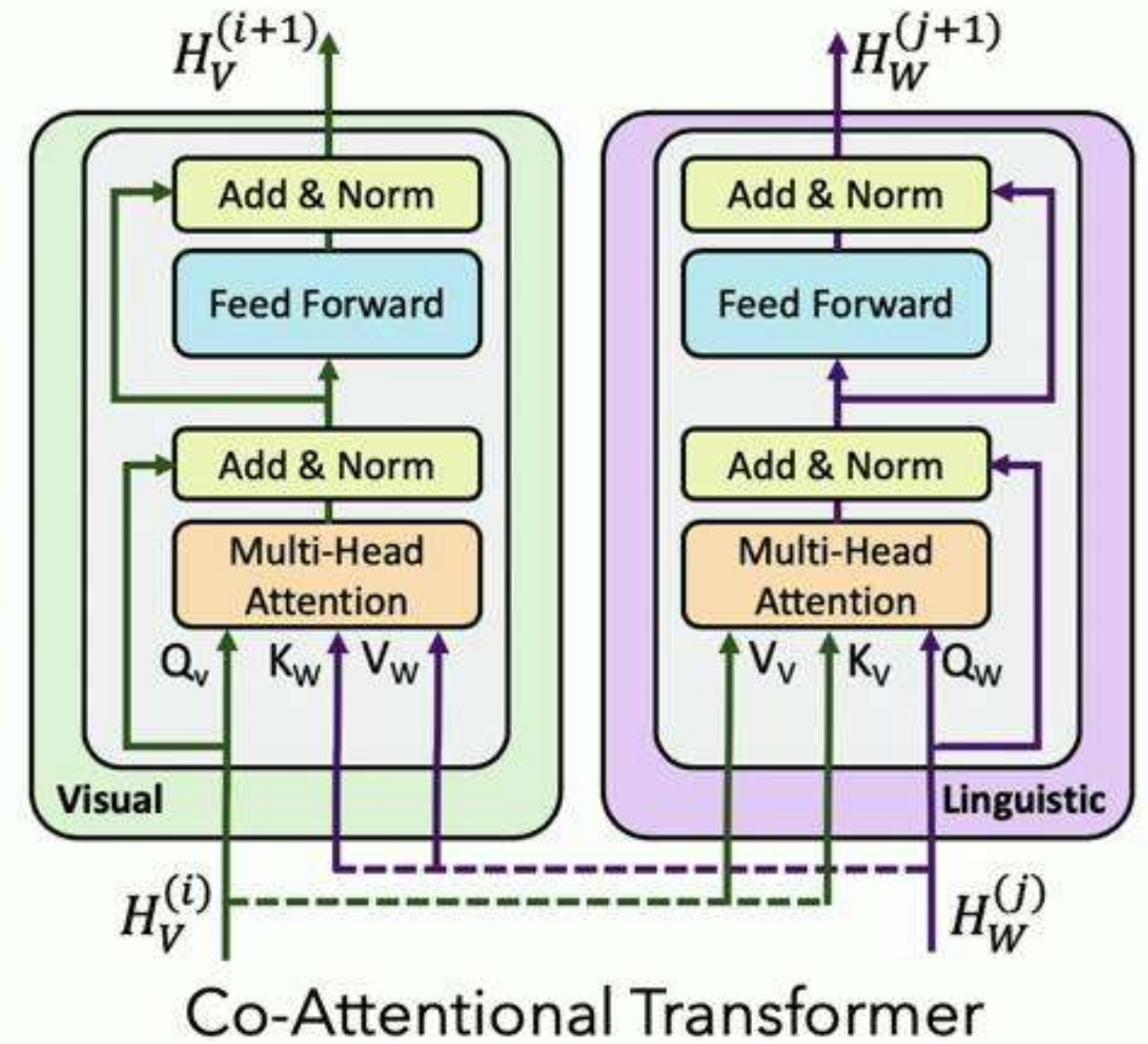


Problem: how to fuse two different modality?

Solution: use **co-attention** [Lu.et.al 2016] to fuse information between different source.

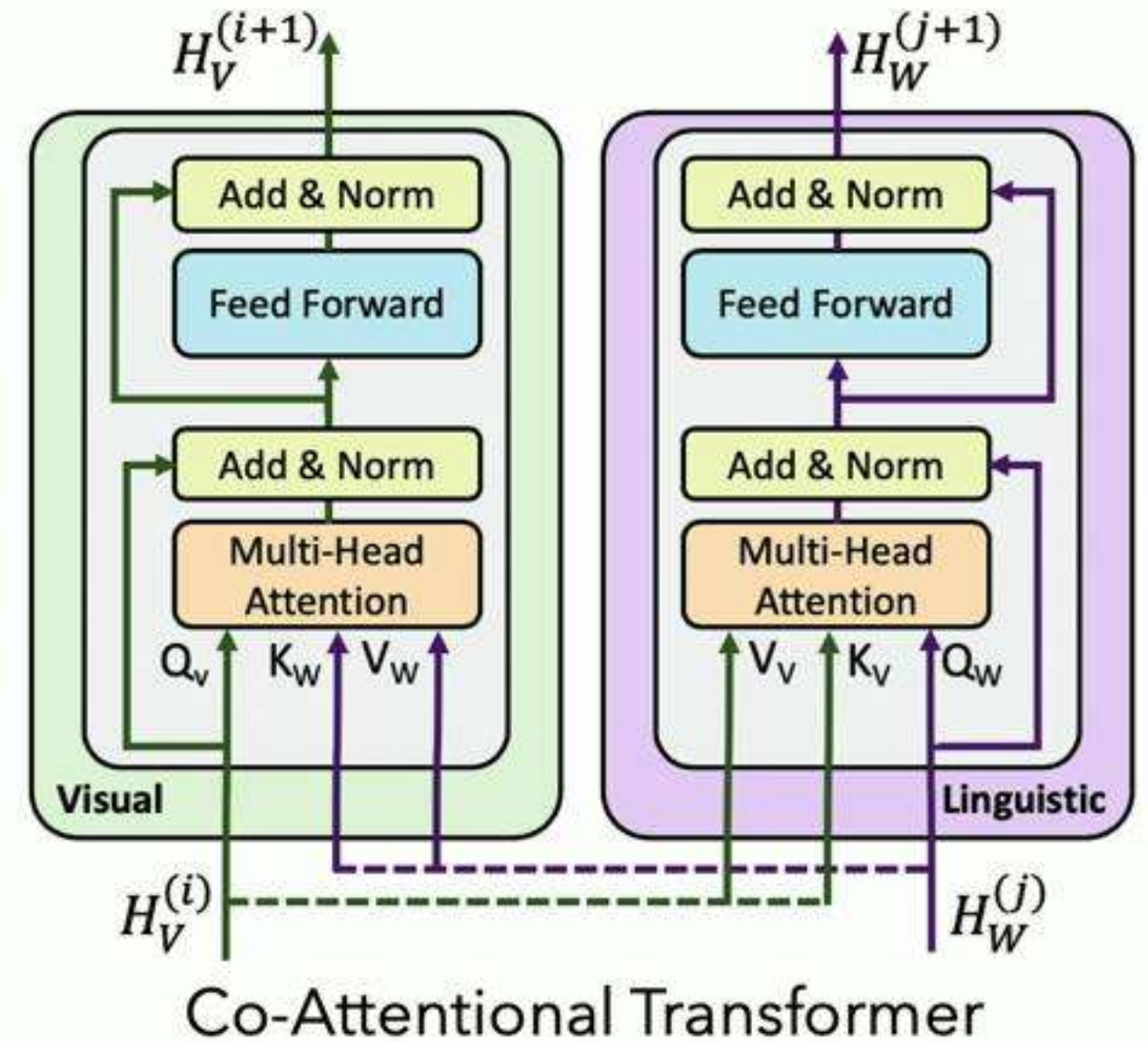


Co-Attentional Transformer



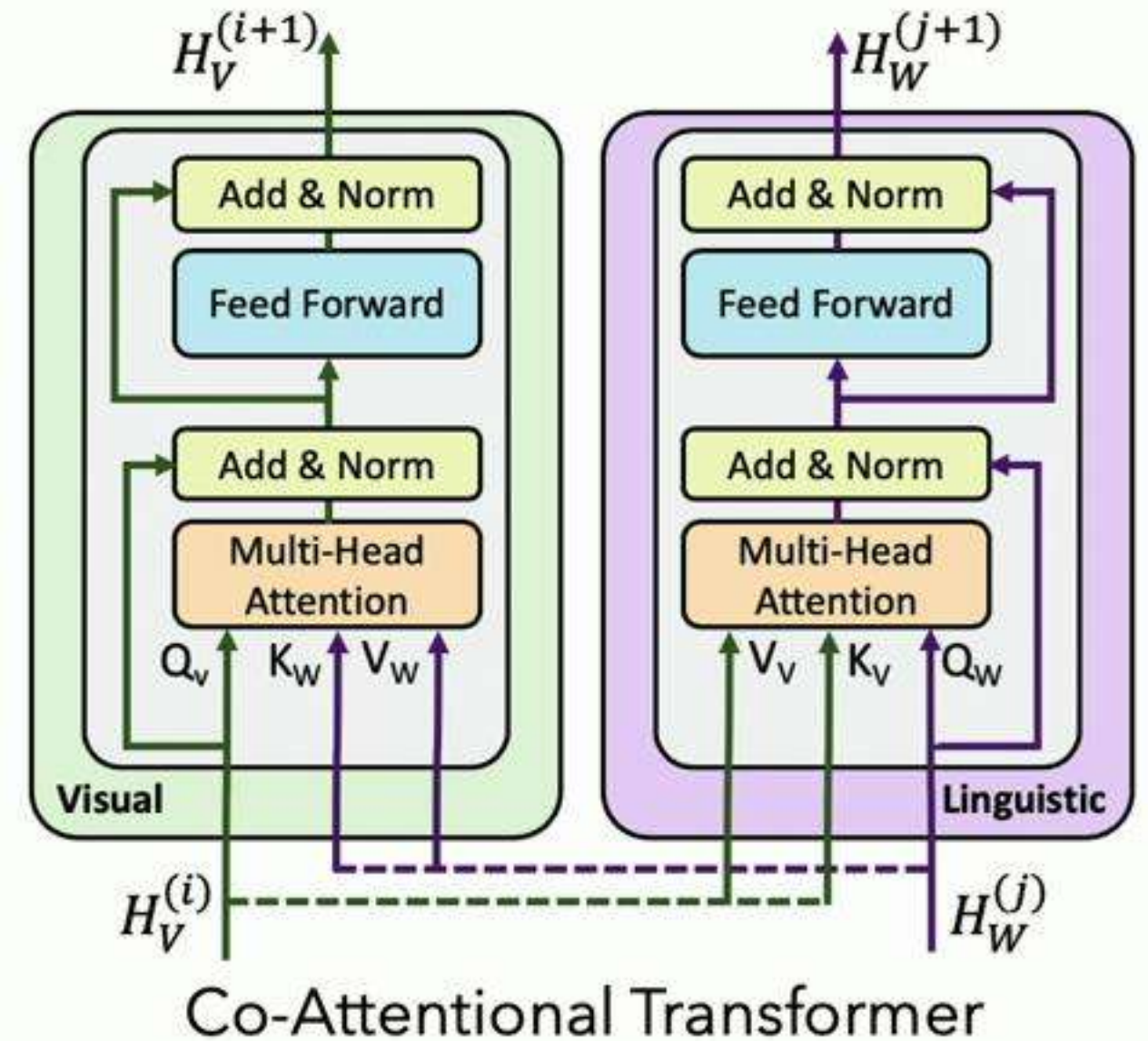
Co-Attentional Transformer

- Transformer encoder with query from another modality.



Co-Attentional Transformer

- Transformer encoder with query from another modality.
- Aggregate information with residual add operation.



Pre-training Objective

Pre-training Objective

Masked multi-modal modelling

Pre-training Objective

Masked multi-modal modelling

- Follows masked LM in BERT.

Masked multi-modal modelling

- Follows masked LM in BERT.
- 15% of the words or image regions to predict.

Masked multi-modal modelling

- Follows masked LM in BERT.
- 15% of the words or image regions to predict.
- Linguistic stream:
 - 80% of the time, replace with *[MASK]*.
 - 10% of the time, replace random word.
 - 10% of the time, keep same.

Masked multi-modal modelling

- Follows masked LM in BERT.
- 15% of the words or image regions to predict.
- Linguistic stream:
 - 80% of the time, replace with *[MASK]*.
 - 10% of the time, replace random word.
 - 10% of the time, keep same.
- Visual stream:
 - 80% of the time, replace with zero vector.

Masked multi-modal modelling

- Follows masked LM in BERT.
- 15% of the words or image regions to predict.
- Linguistic stream:
 - 80% of the time, replace with *[MASK]*.
 - 10% of the time, replace random word.
 - 10% of the time, keep same.
- Visual stream:
 - 80% of the time, replace with zero vector.

Multi-modal alignment prediction

Image Representation

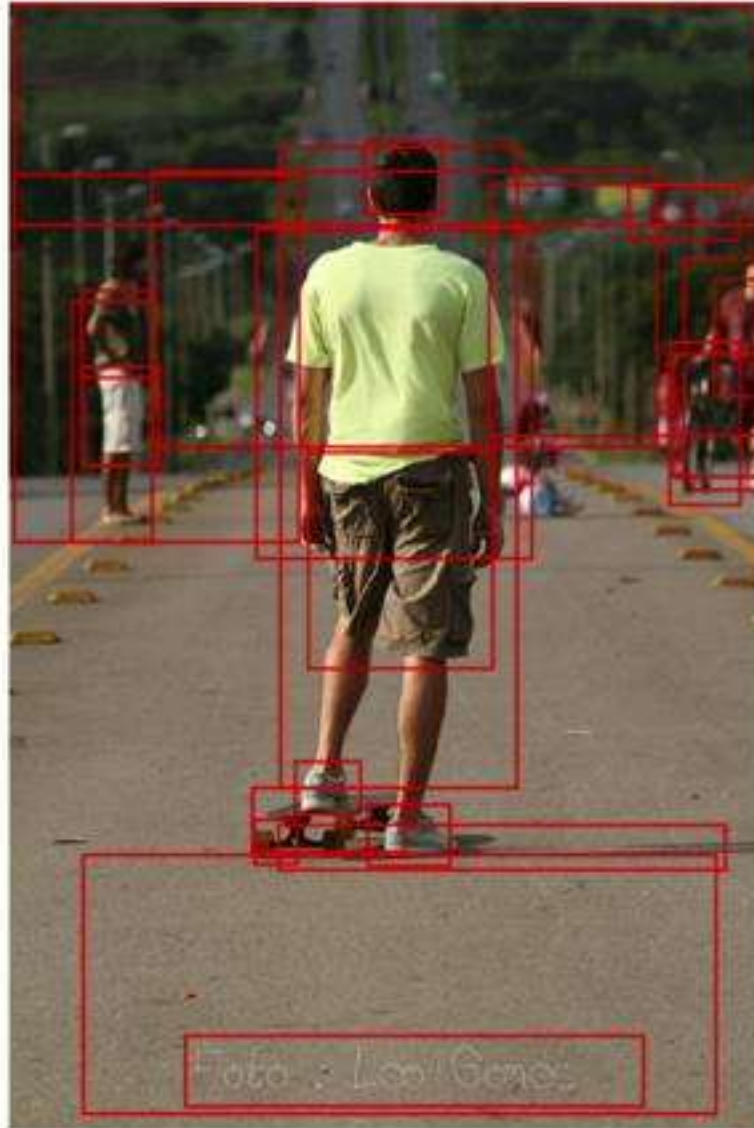
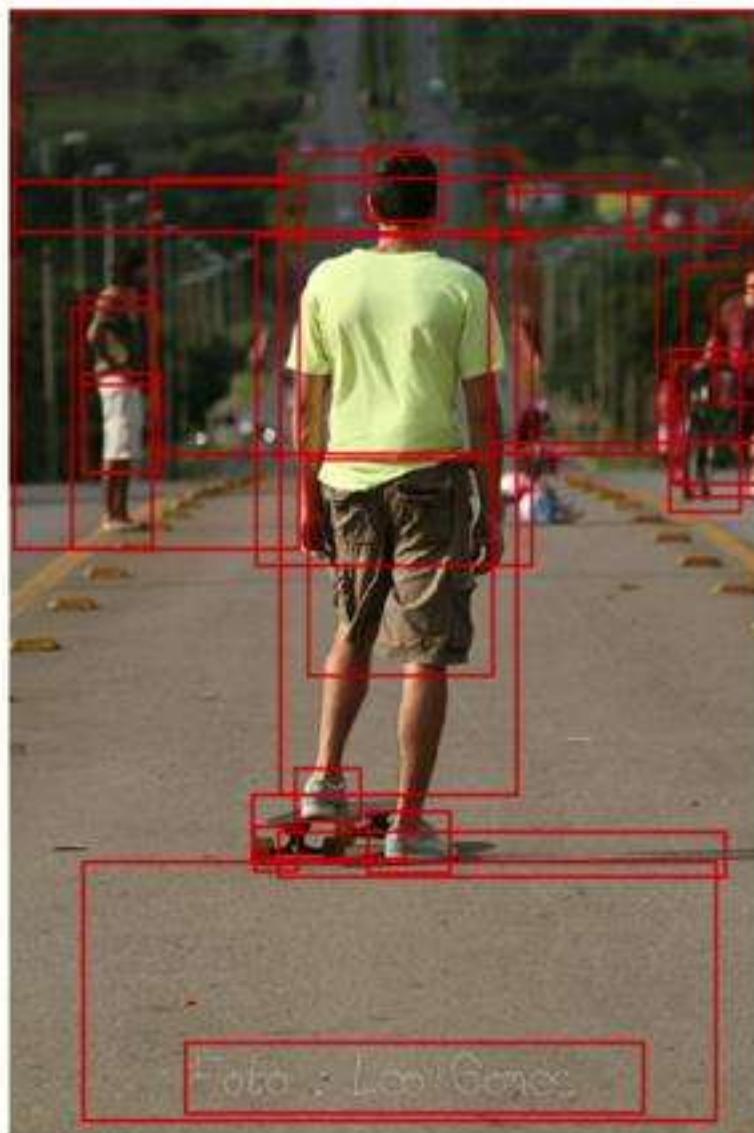
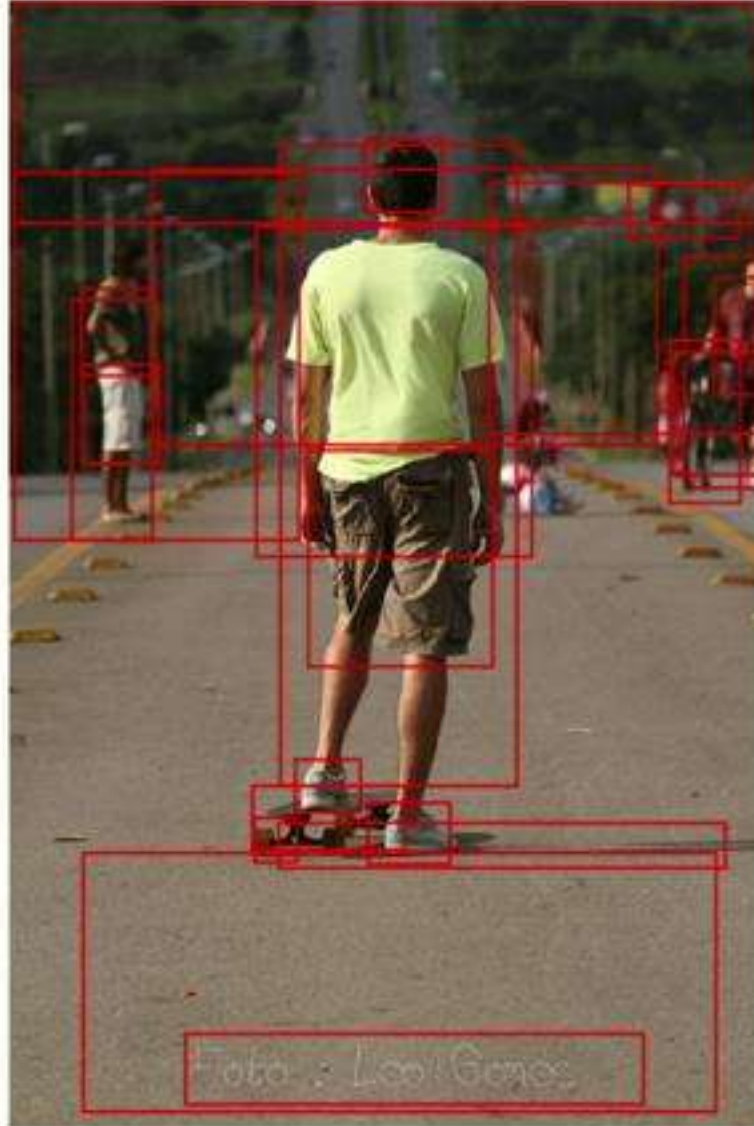


Image Representation



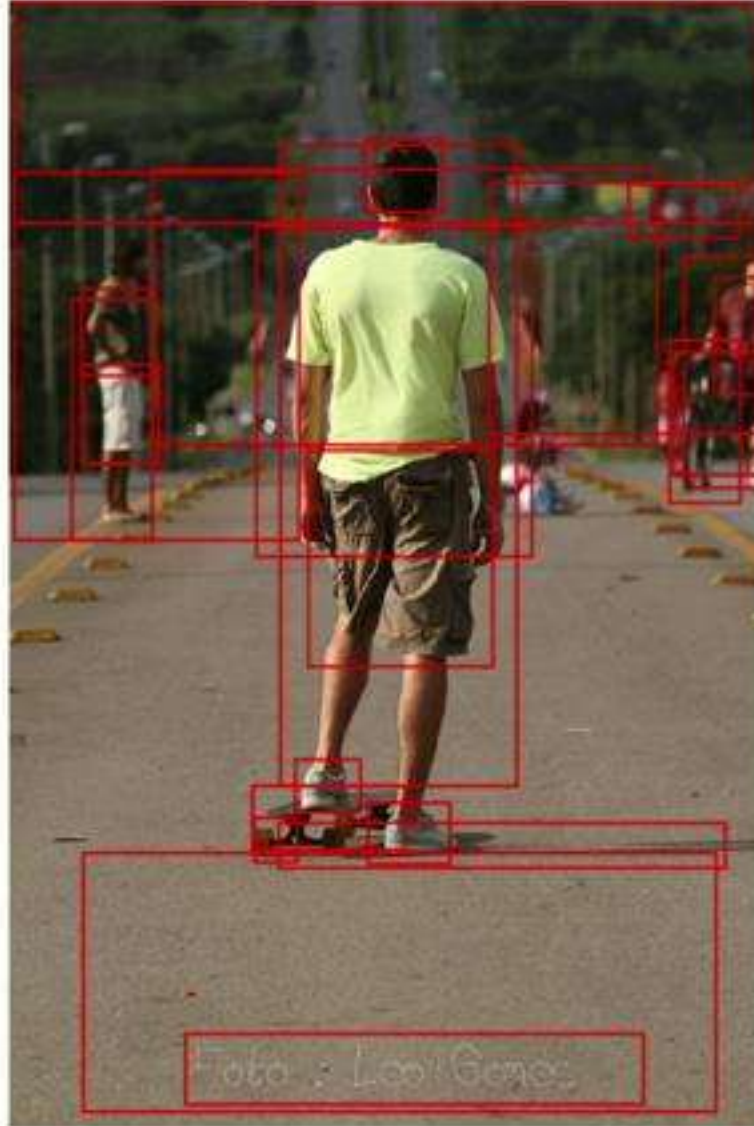
- Faster R-CNN with Res101 backbone.

Image Representation



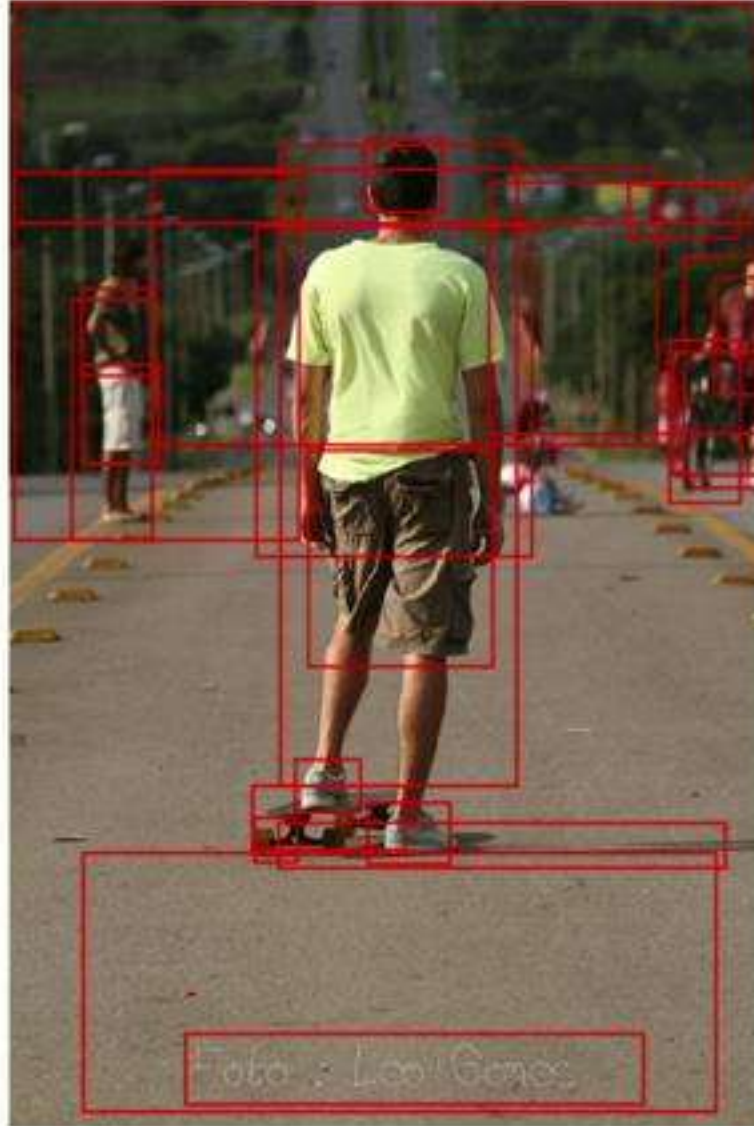
- Faster R-CNN with Res101 backbone.
- Trained on Visual Genome dataset.

Image Representation



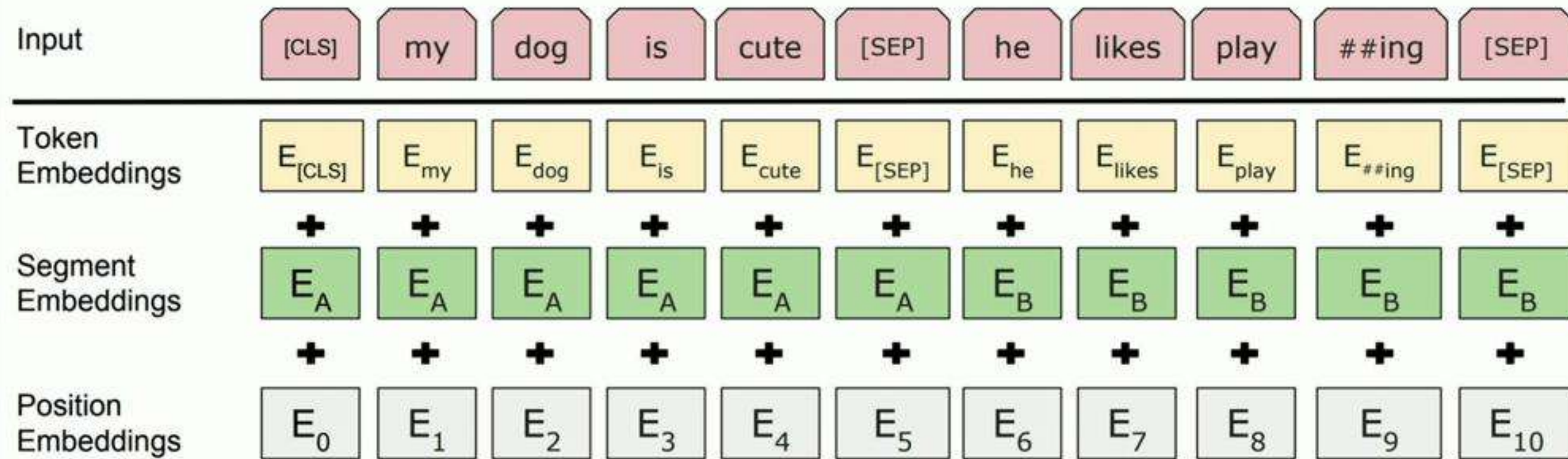
- Faster R-CNN with Res101 backbone.
- Trained on Visual Genome dataset.
- 1600 detection classes.

Image Representation

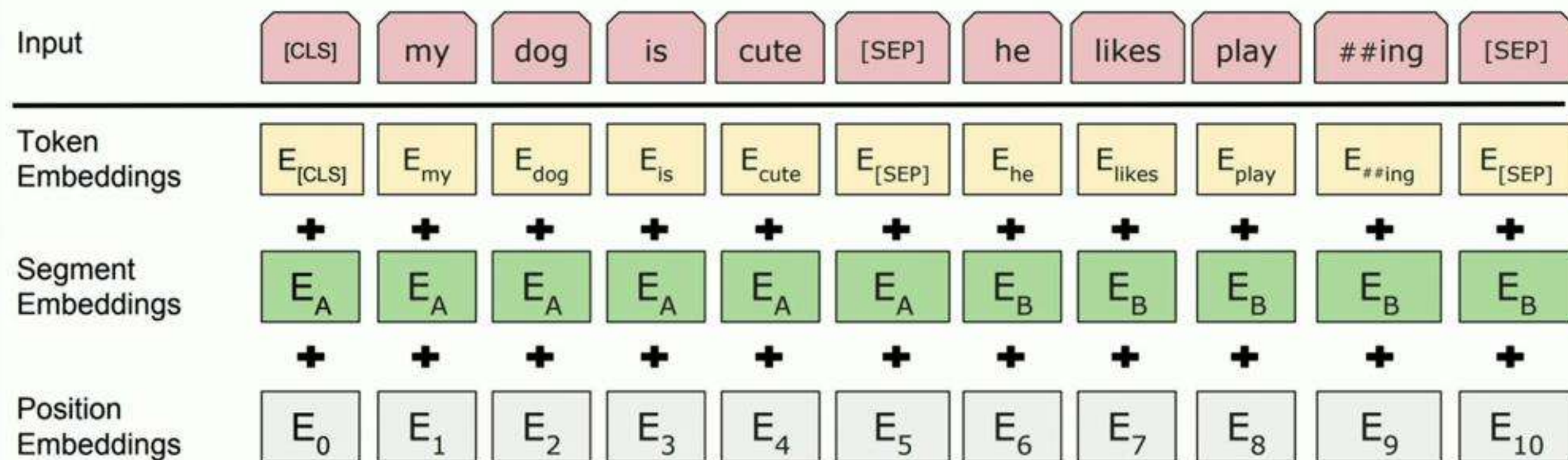


- Faster R-CNN with Res101 backbone.
- Trained on Visual Genome dataset.
- 1600 detection classes.
- Sum of region embeddings + location embeddings

Text Representation

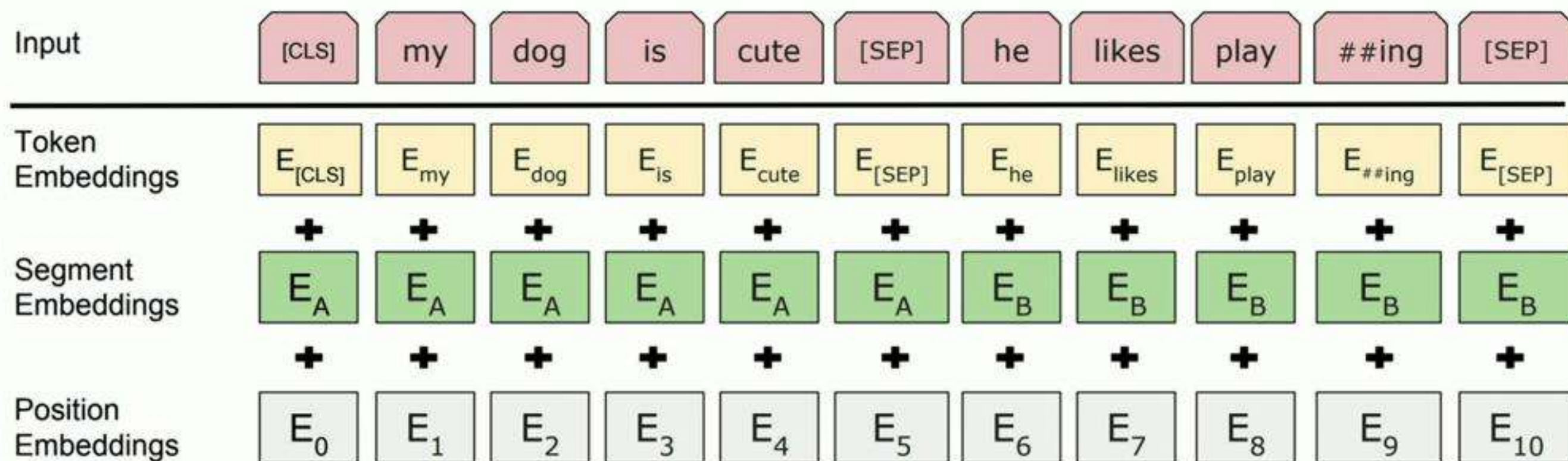


Text Representation



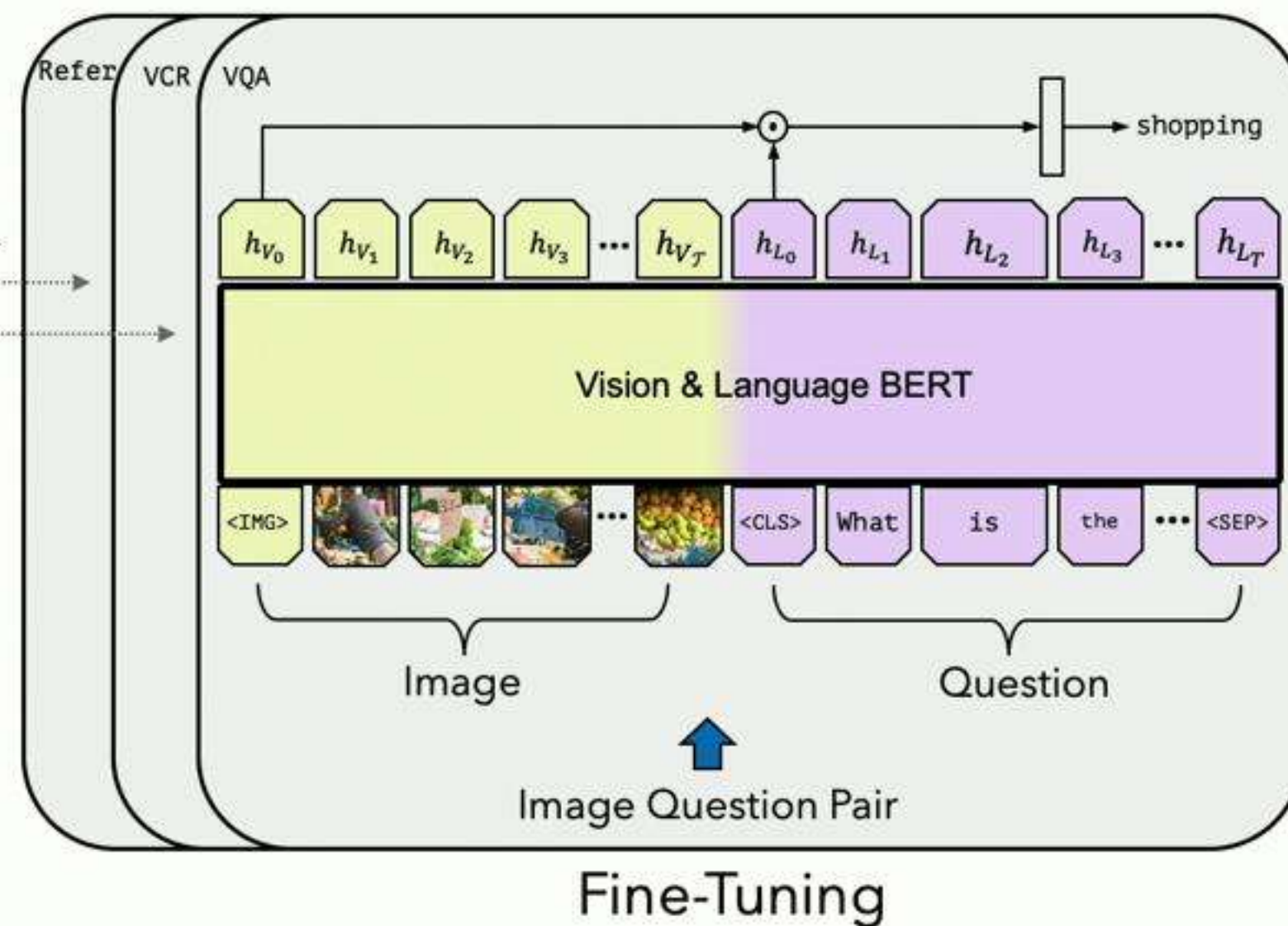
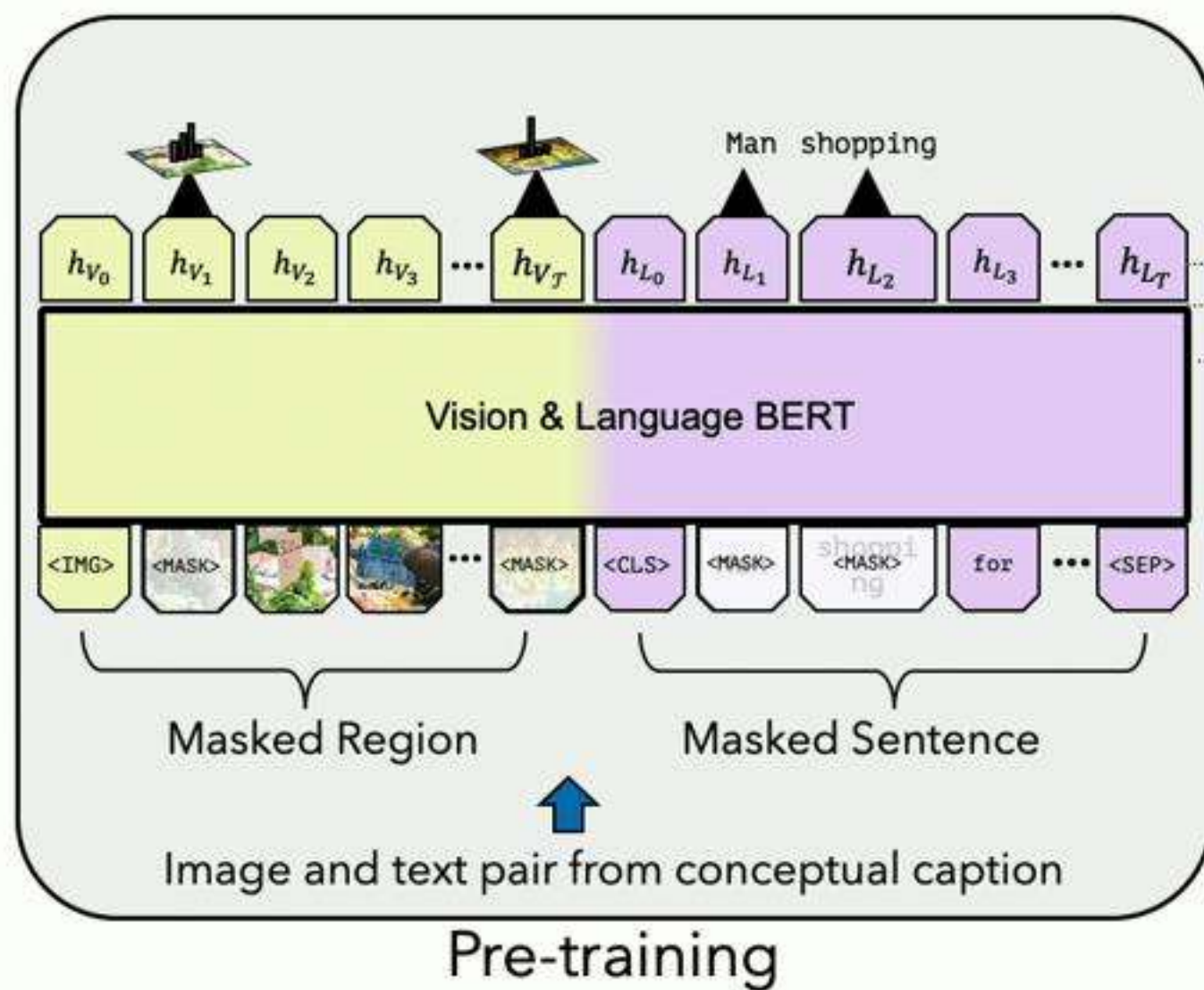
- Use 30,000 WordPiece vocabulary on input.

Text Representation



- Use 30,000 WordPiece vocabulary on input.
- Each token is sum of three embeddings

Fine-tuning Procedure



Tasks

Tasks



Is there something to cut the vegetables with?

VQA

Tasks



Is there something to cut the vegetables with?

VQA



Why is [person4] pointing at [person1]?

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

VCR Q→A

Rationale: a) is correct because...

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

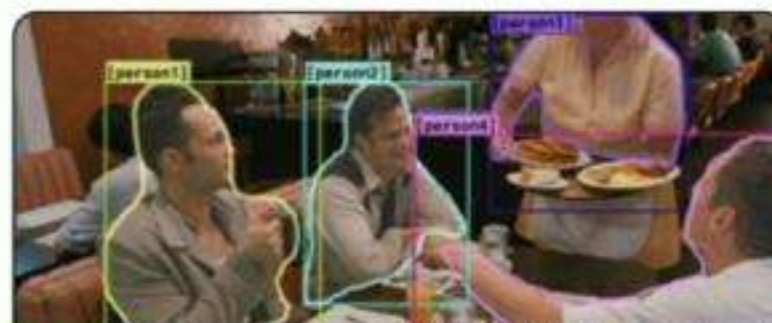
VCR QA→R

Tasks



Is there something to cut the vegetables with?

VQA



Why is [person4] pointing at [person1]?

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

VCR Q→A

Rationale: a) is correct because...

- a) [person1] has the pancakes in front of him.
- b) [person4] is taking everyone's order and asked for clarification.
- c) [person0] is looking at the pancakes both she and [person2] are smiling slightly.
- d) [person3] is delivering food to the table, and she might not know whose order is whose.

VCR QA→R



Guy in yellow dribbling ball

Referring Expressions

Tasks



Is there something to cut the vegetables with?

VQA



Why is [person4] pointing at [person1]?

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

VCR Q→A

Rationale: a) is correct because...

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

VCR QA→R



Guy in yellow dribbling ball

Referring Expressions

A large bus sitting next to a very tall building.

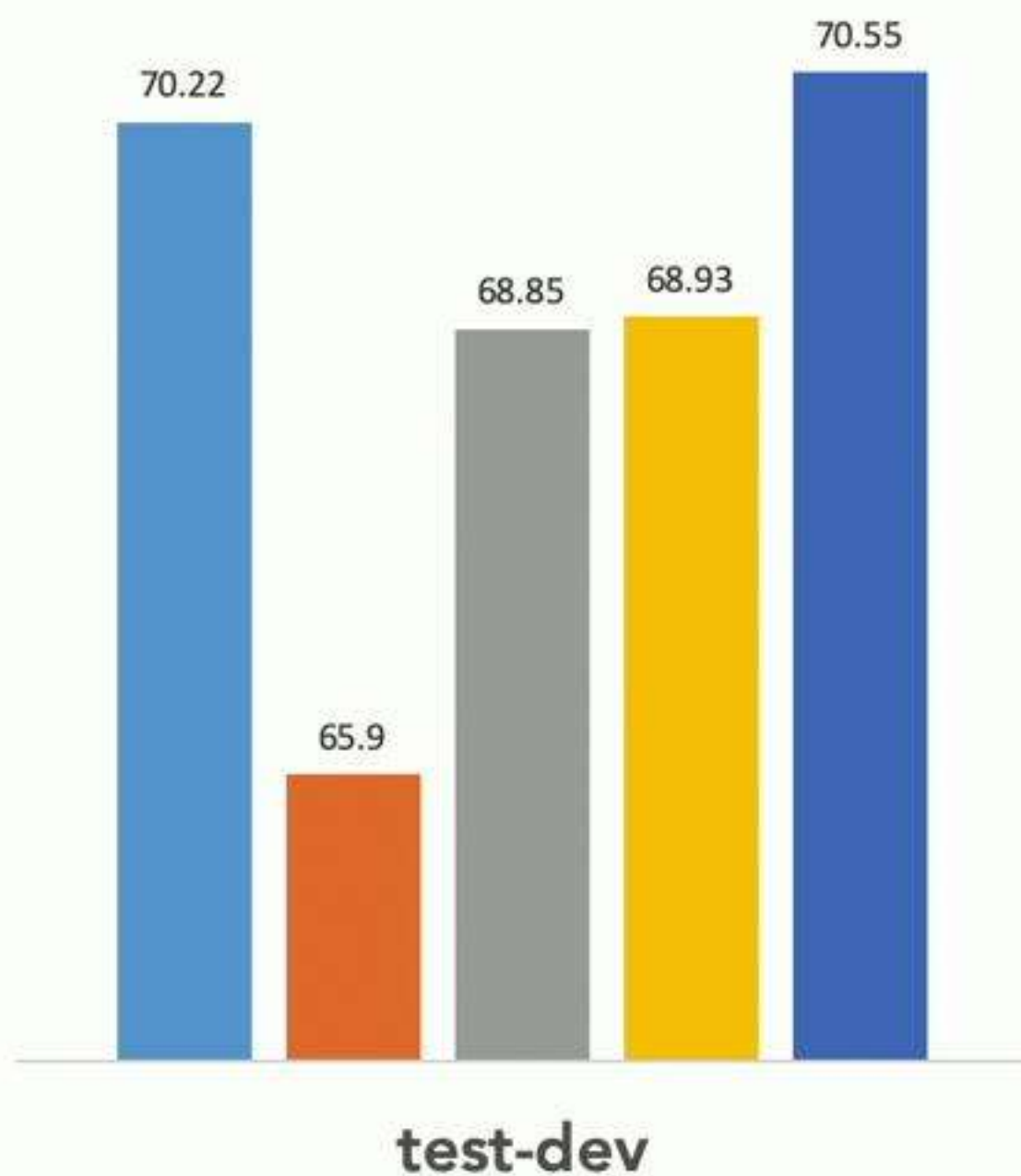


Caption-Based Image Retrieval

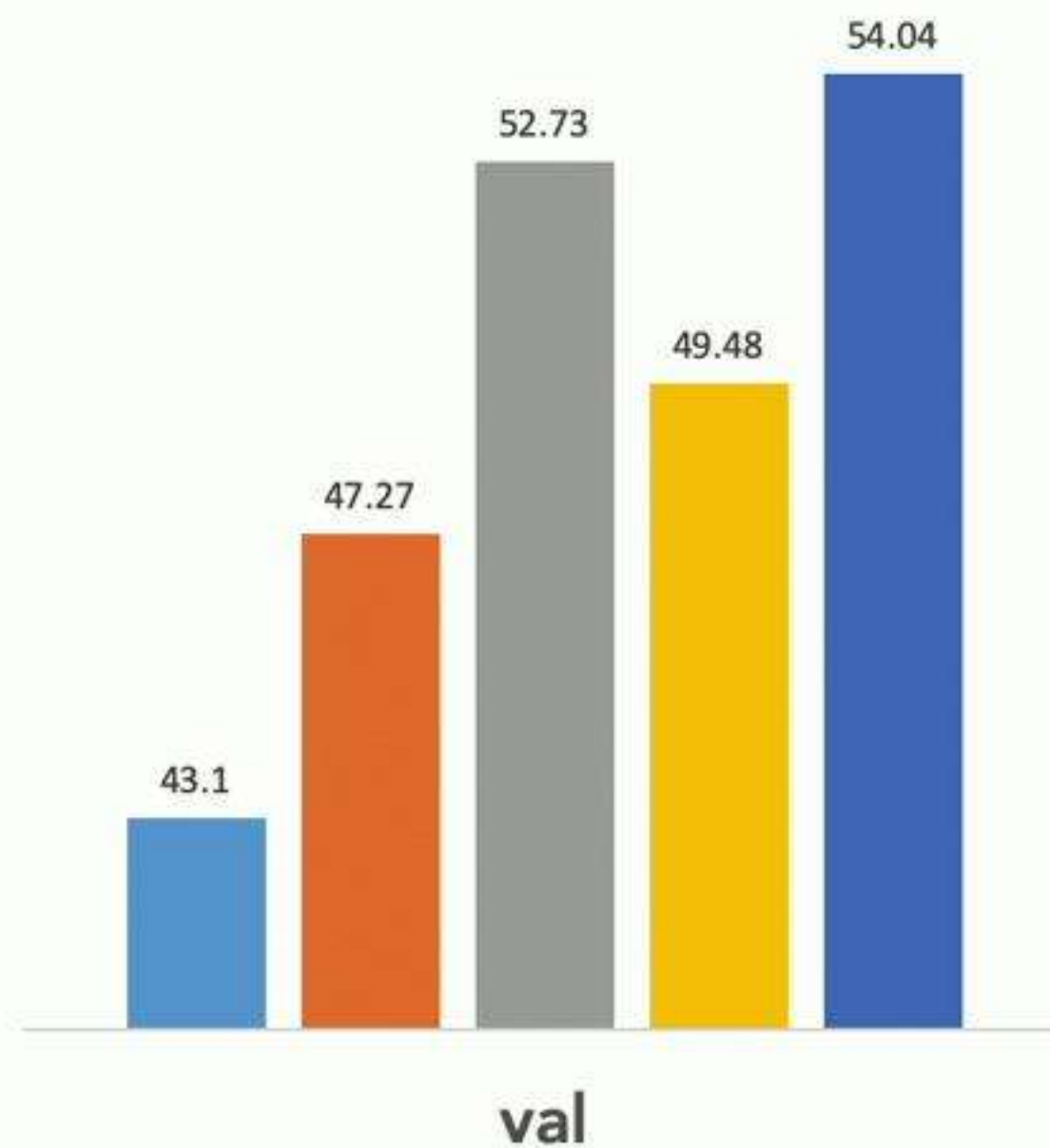
Results

■ SOTA ■ Single-Stream[†] ■ Single-Stream ■ ViLBERT[†] ■ ViLBERT

VQA



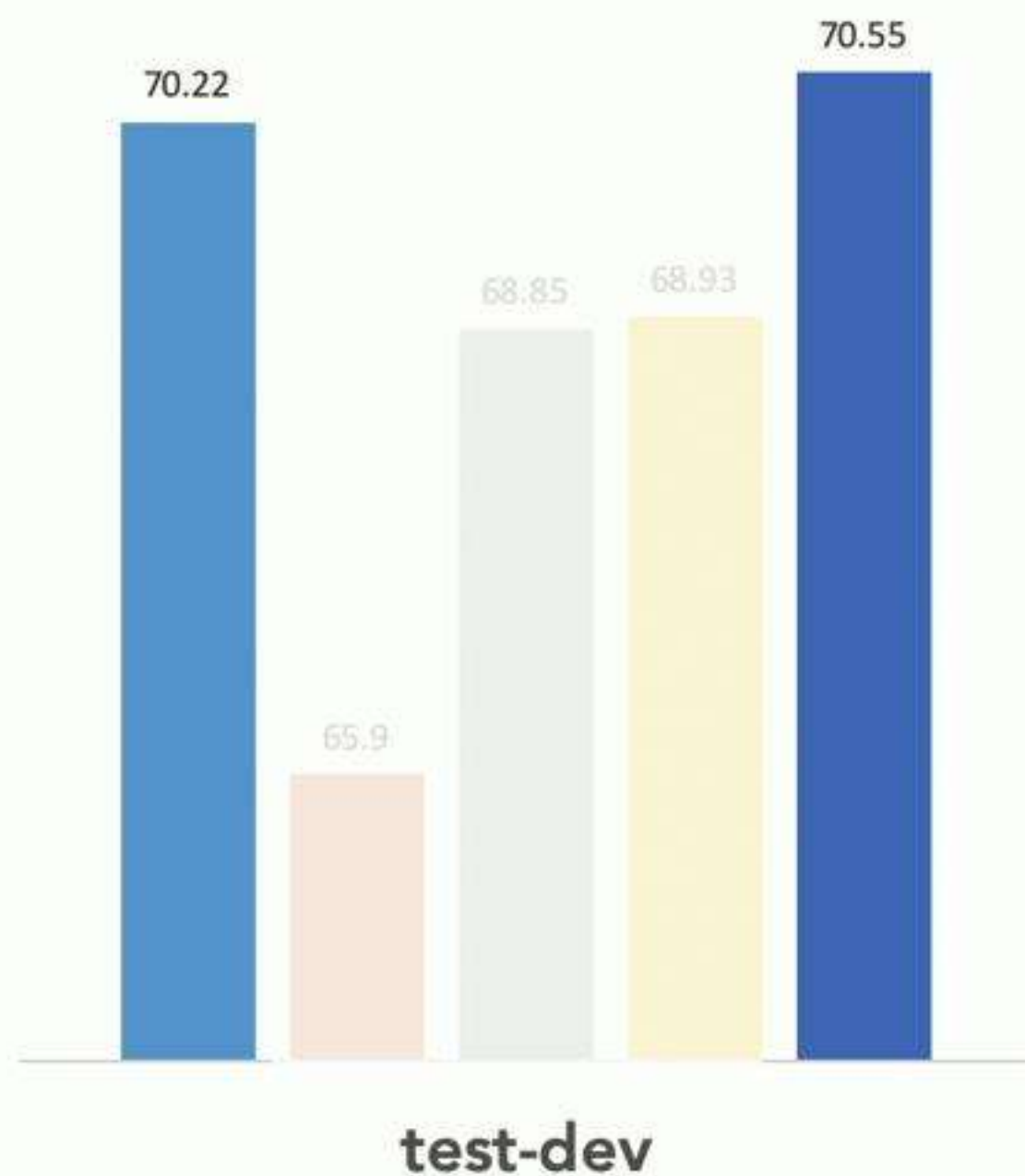
VCR Q->A



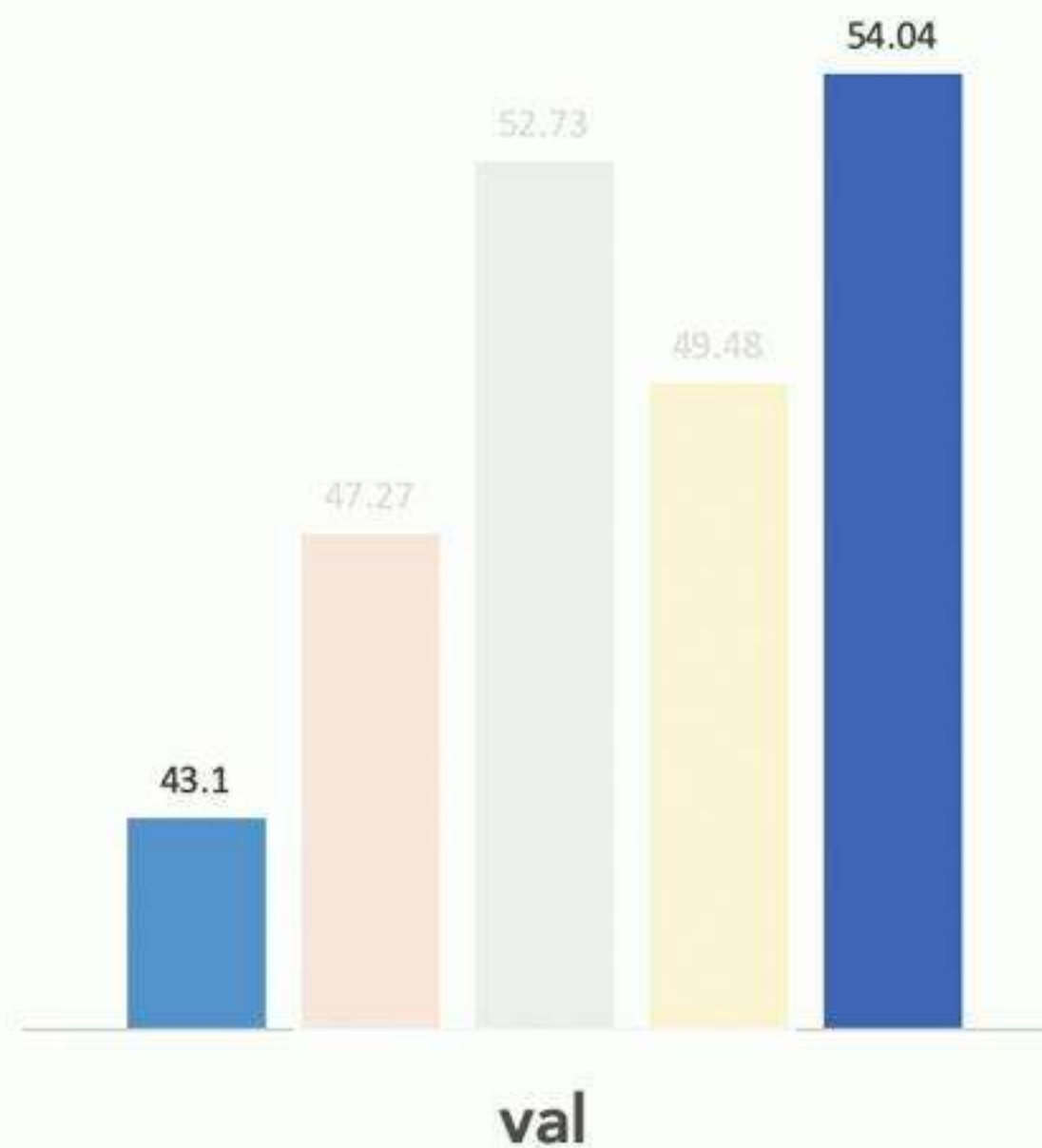
Results

■ SOTA ■ Single-Stream† ■ Single-Stream ■ ViLBERT† ■ ViLBERT

VQA



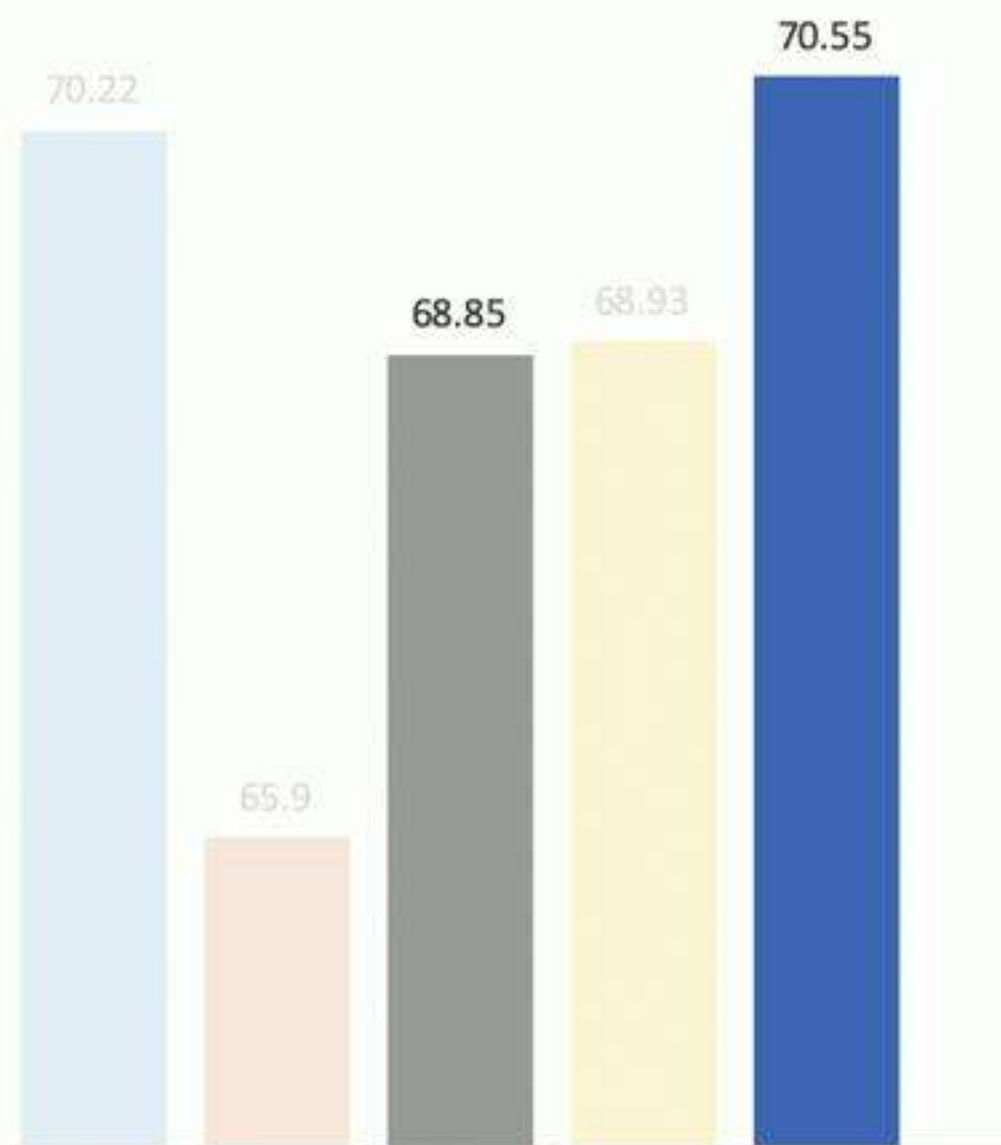
VCR Q->A



Results

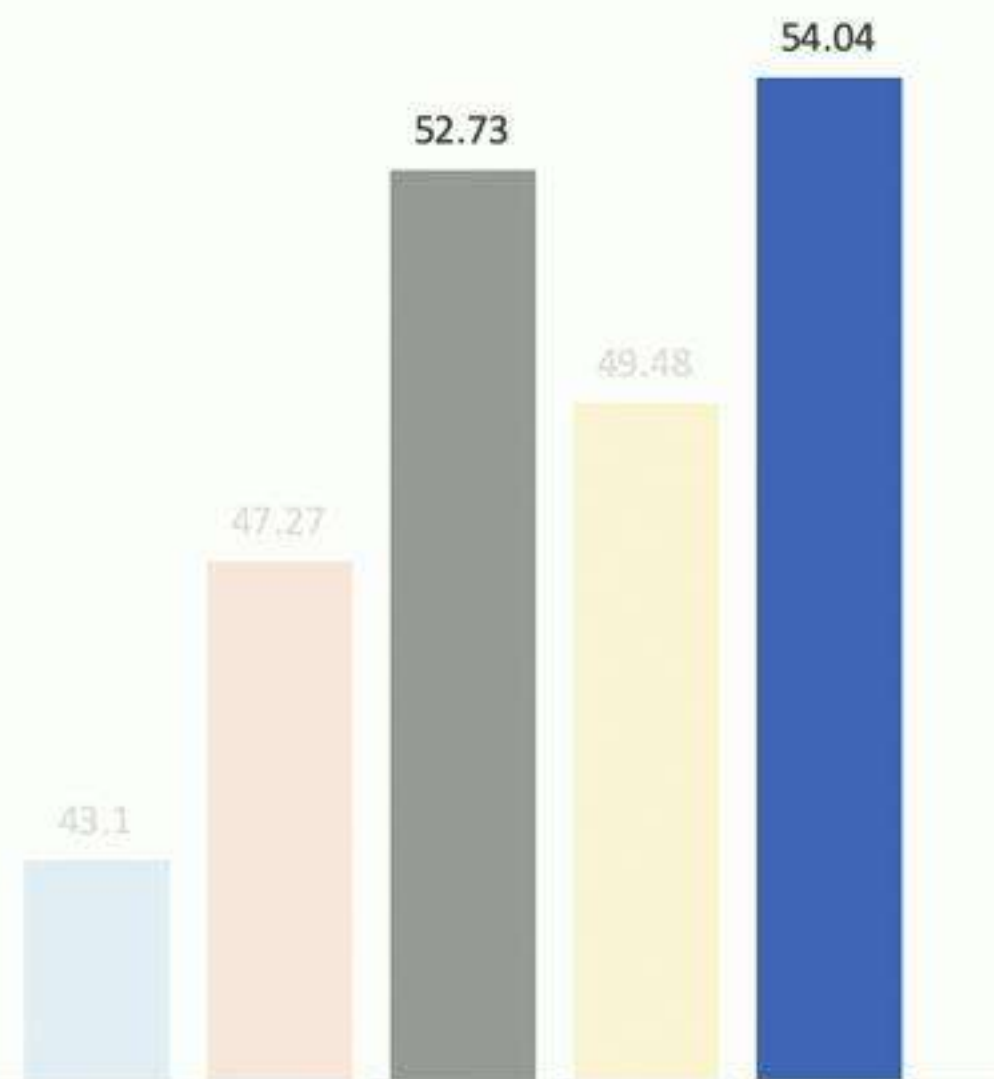
■ SOTA ■ Single-Stream[†] ■ Single-Stream ■ ViLBERT[†] ■ ViLBERT

VQA



test-dev

VCR Q->A

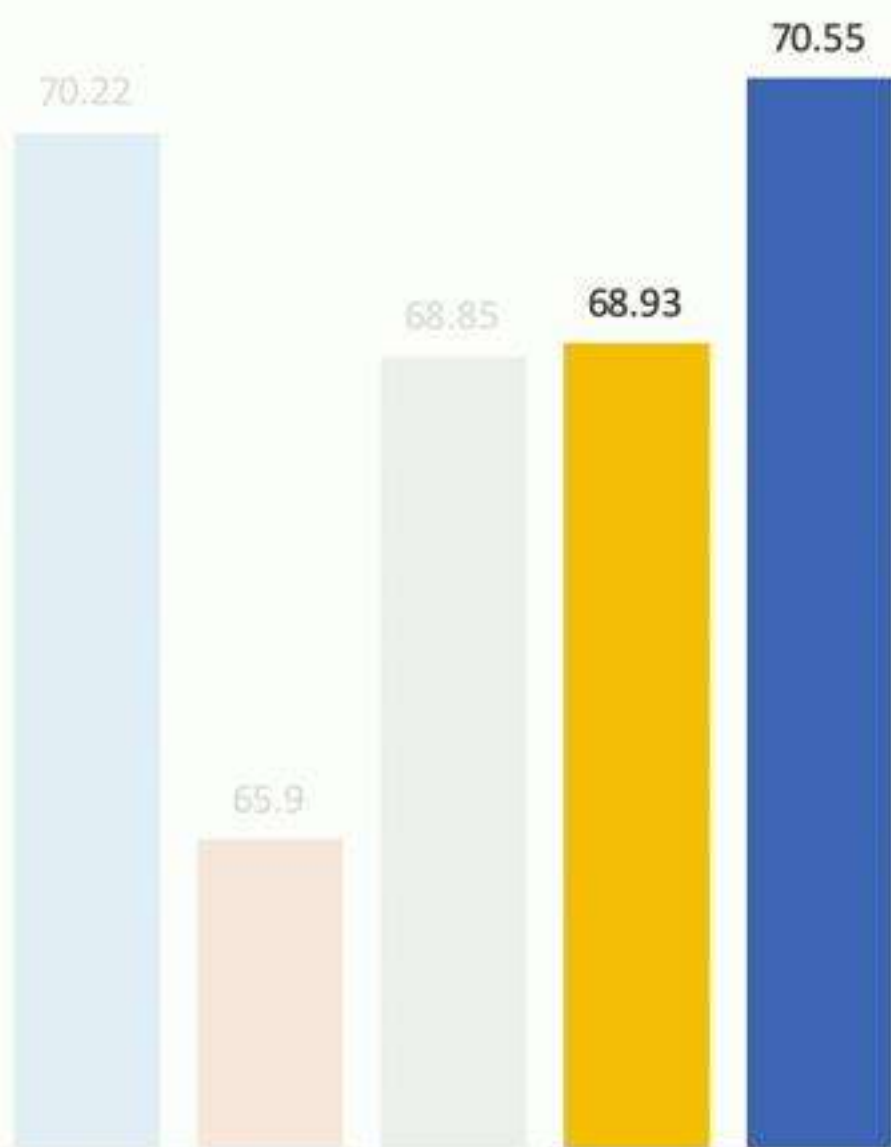


val

Results

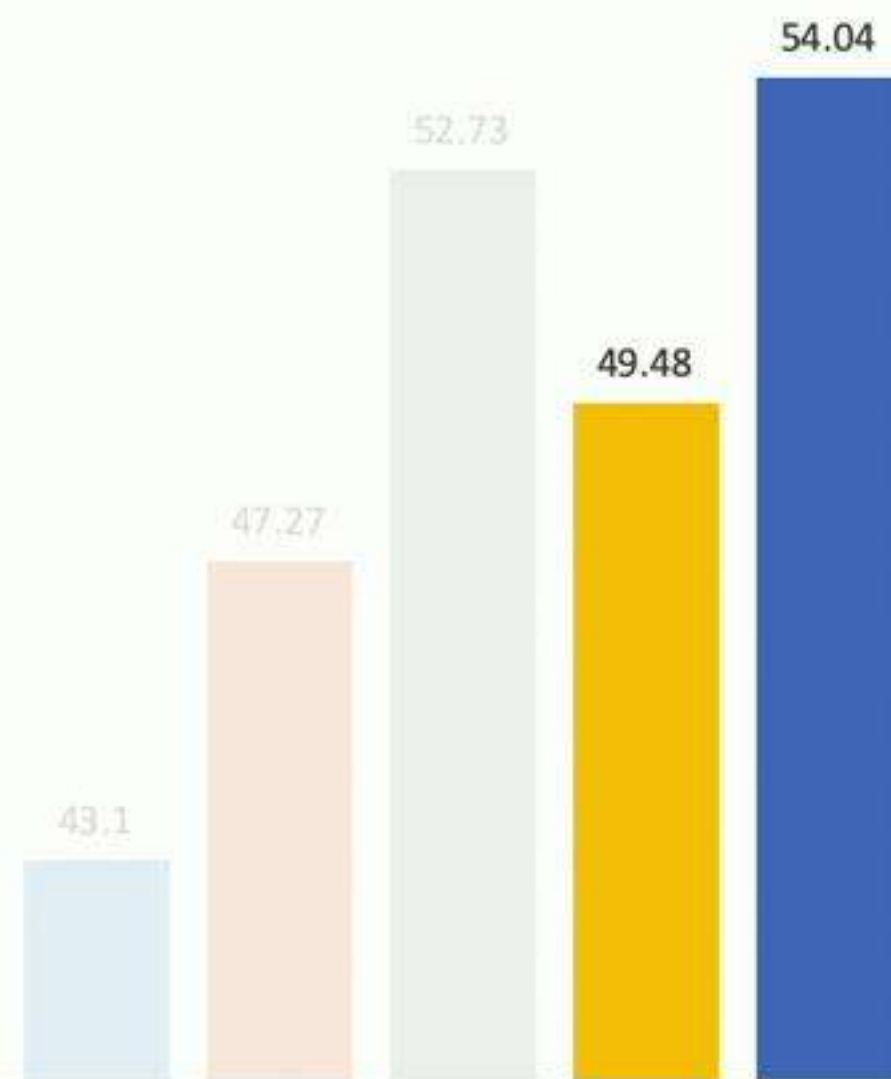
■ SOTA ■ Single-Stream† ■ Single-Stream ■ ViLBERT† ■ ViLBERT

VQA



test-dev

VCR Q->A



val

Results

■ SOTA ■ Single-Stream[†] ■ Single-Stream ■ ViLBERT[†] ■ ViLBERT

RefCOCO+

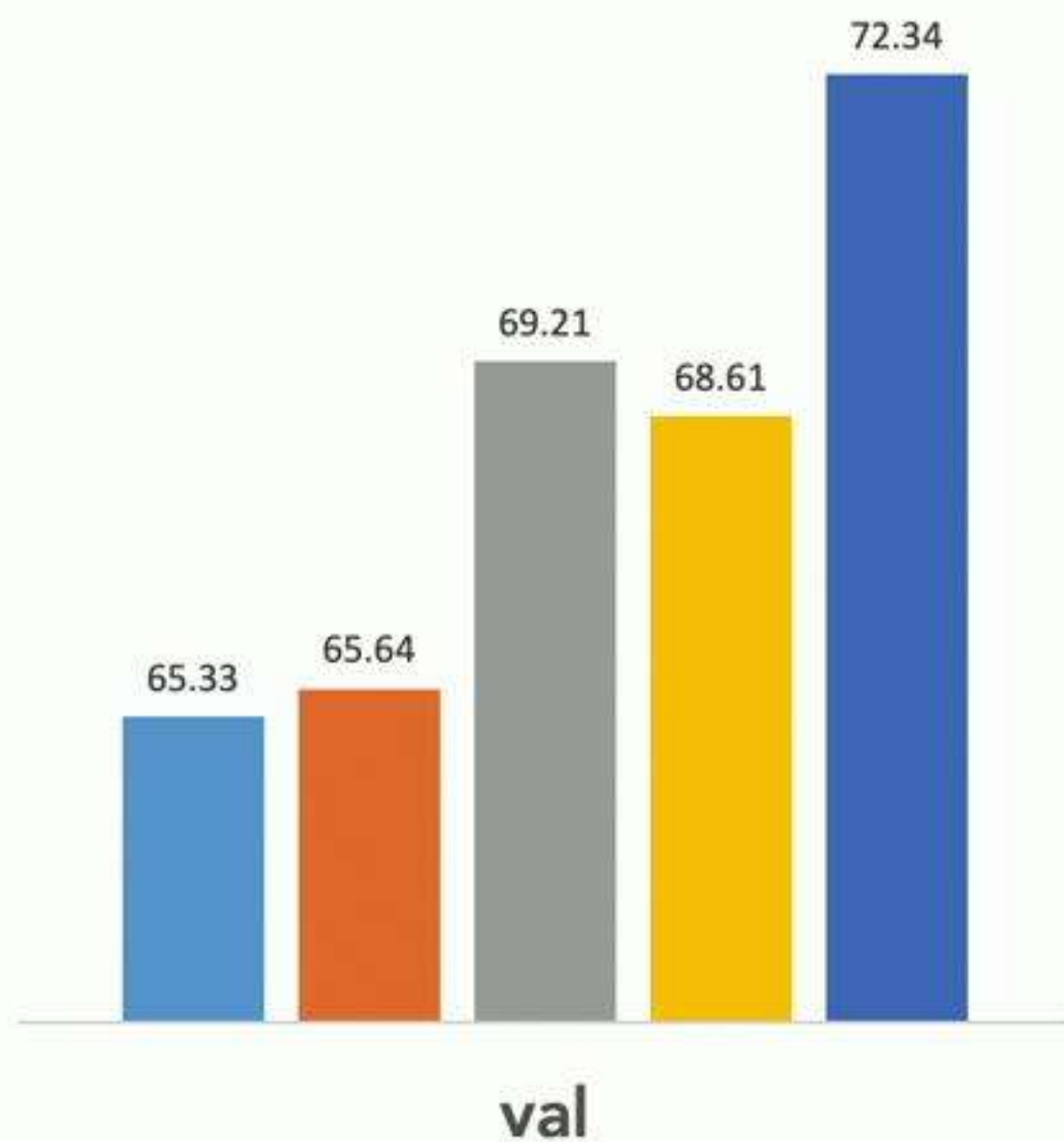
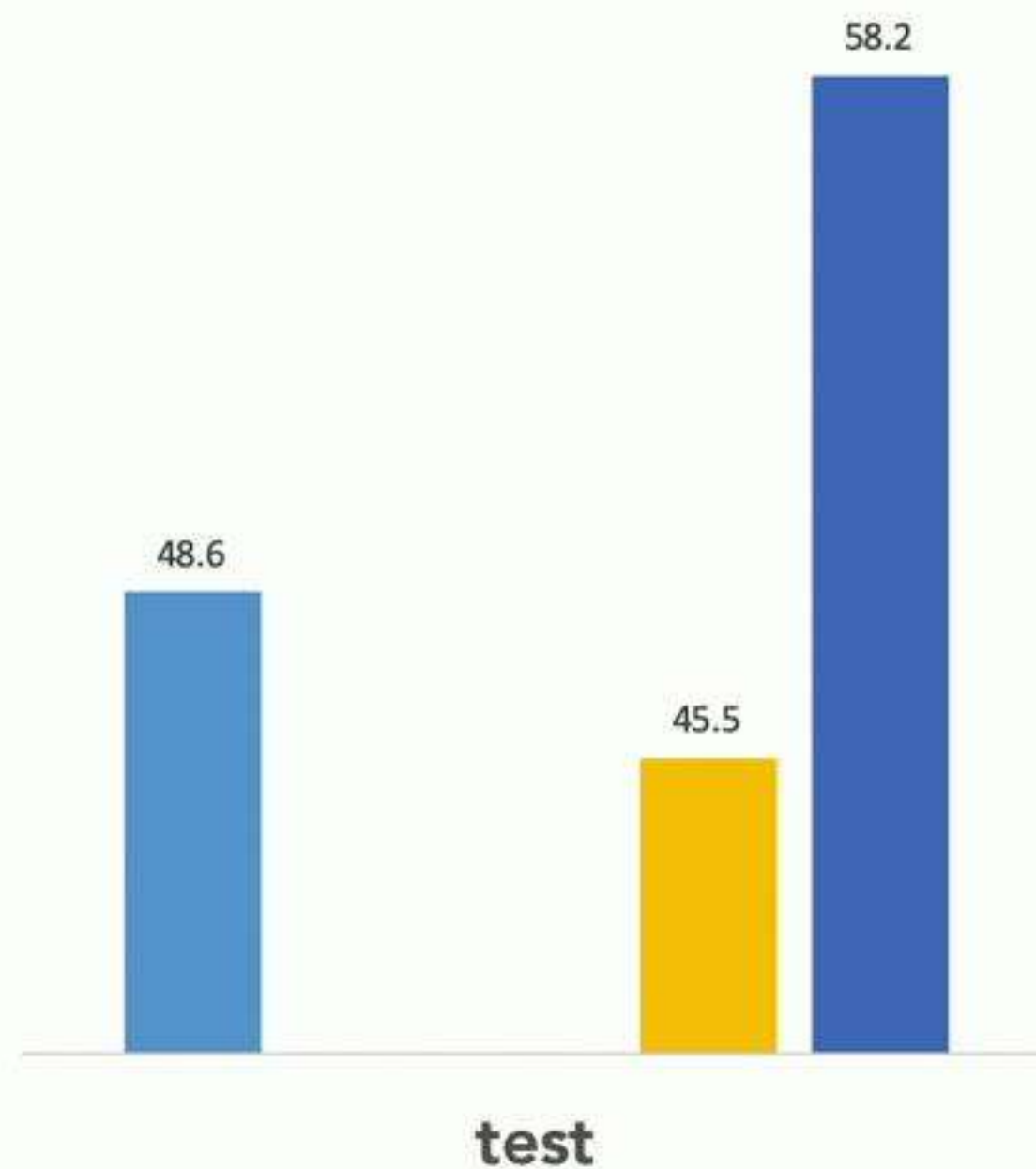


Image Retrieval

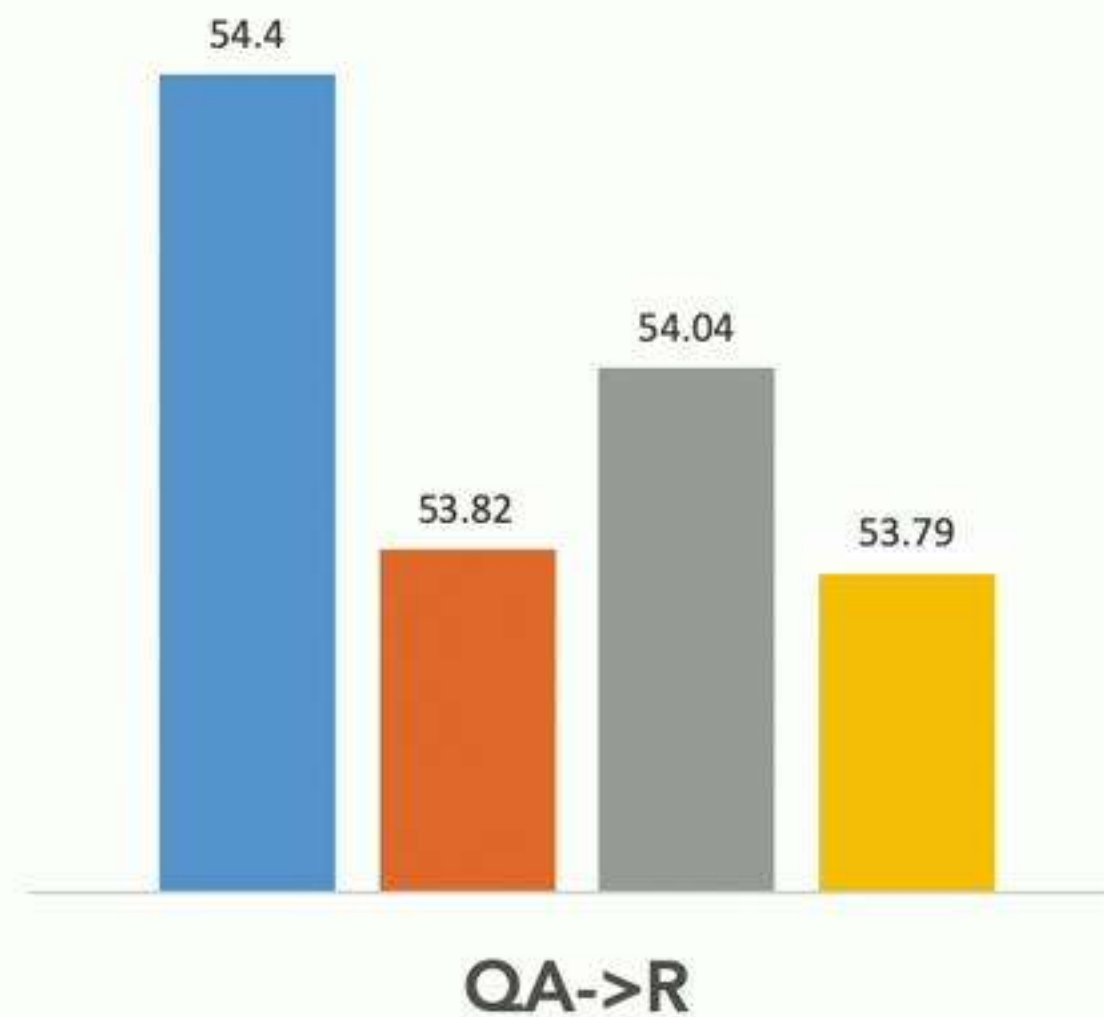
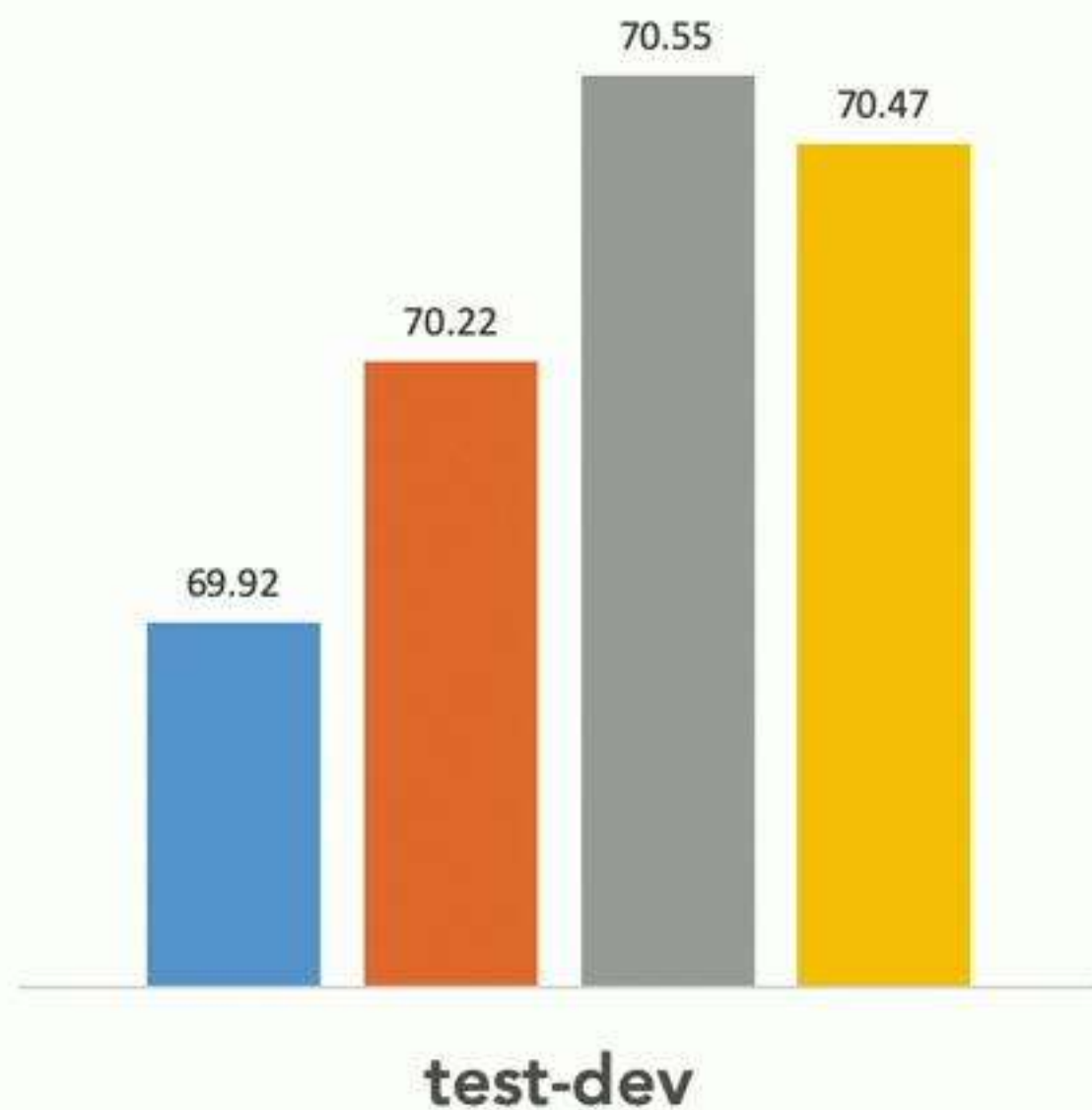


Ablations - Depth

■ 2-layer ■ 4-layer ■ 6-layer ■ 8-layer

VQA

VCR

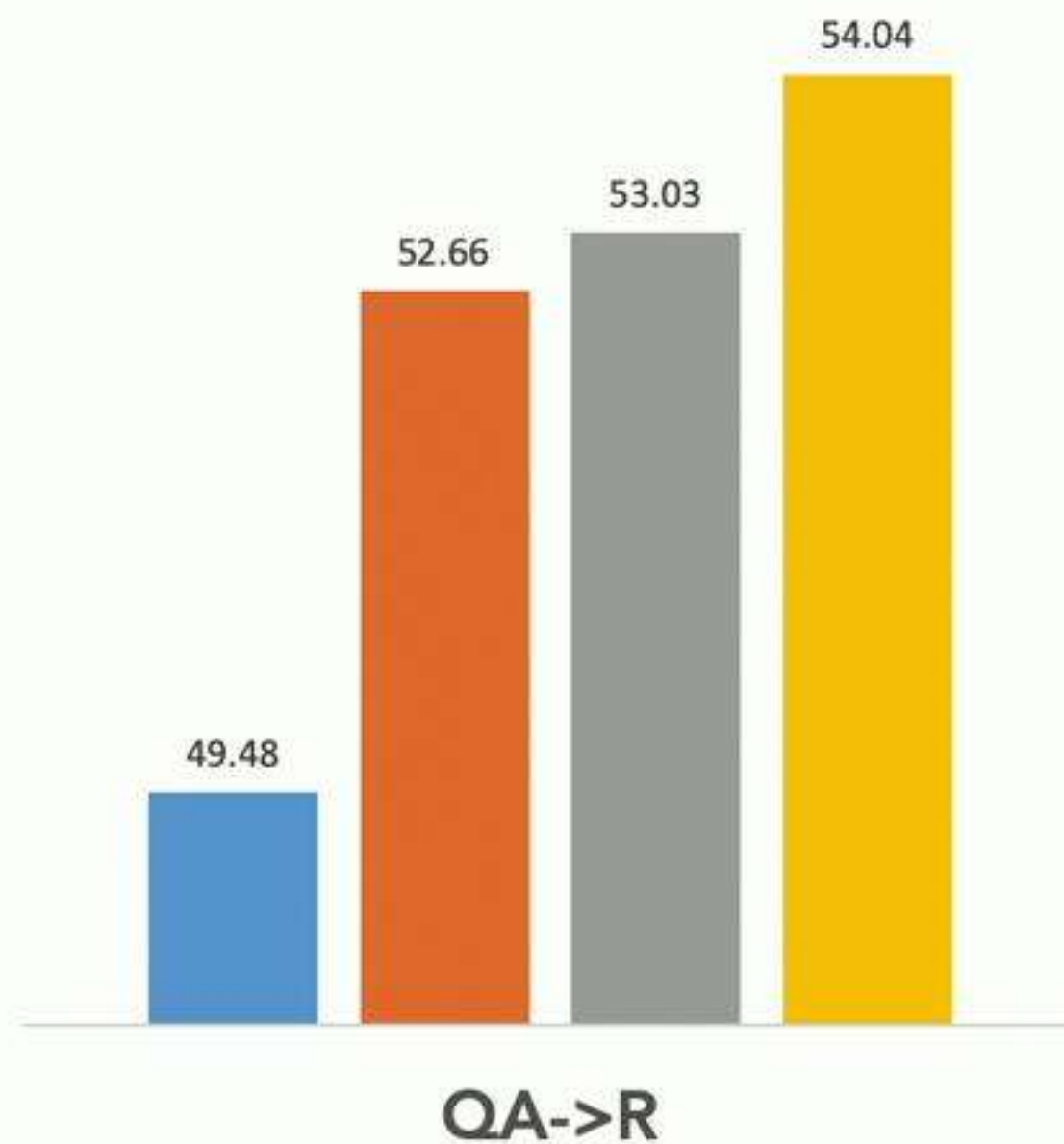
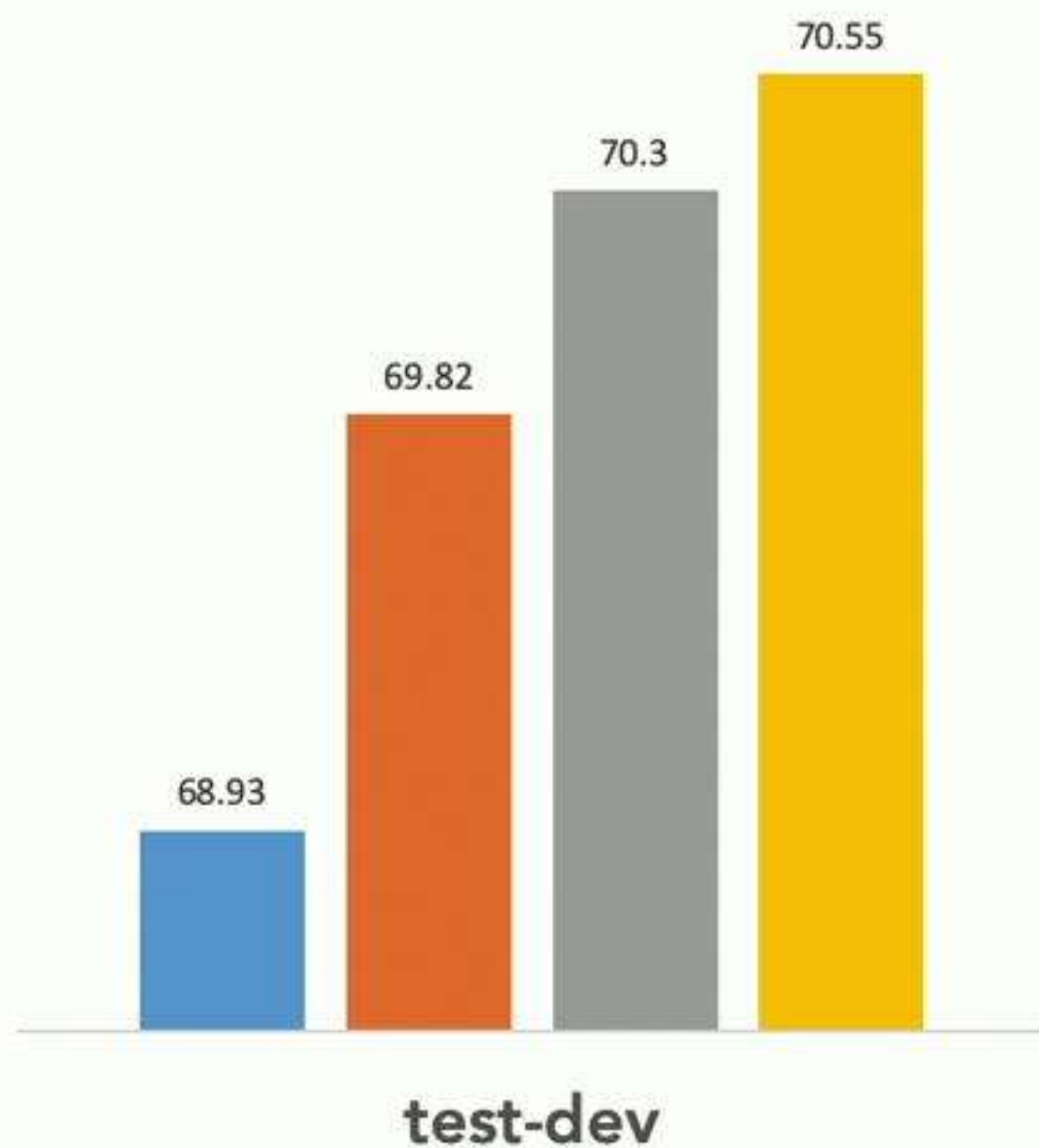


Ablations - Pretraining

■ 0% ■ 25% ■ 50% ■ 100%

VQA

VCR



Concurrent Work

	Method	Architecture	Visual Token	Pre-train Datasets	Pre-train Tasks	Downstream Tasks
Published Works	VideoBERT (Sun et al. 2019b)	single cross-modal Transformer	video frame	Cooking312K (Sun et al. 2019b)	1) sentence-image alignment 2) masked language modeling 3) masked visual-words prediction	1) zero-shot action classification 2) video captioning
Works Under Review / Just Got Accepted	CBT (Sun et al. 2019a)	two single-modal Transformer (vision & language respectively) + one cross-modal Transformer	video frame	Cooking312K (Sun et al. 2019b)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature regression	1) action anticipation 2) video captioning
	ViLBERT (Lu et al. 2019)	one single-modal Transformer (language) + one cross-modal Transformer (with restricted attention pattern)	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification	1) visual question answering 2) visual commonsense reasoning 3) grounding referring expressions 4) image retrieval 5) zero-shot image retrieval
	B2T2 (Alberti et al. 2019)	single cross-modal Transformer	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling	1) visual commonsense reasoning
	LXMERT (Hao Tan, 2019)	two single-modal Transformer (vision & language respectively) + one cross-modal Transformer	image RoI	‡ COCO Caption + VG Caption + VG QA + VQA + GQA	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification 4) masked visual-feature regression 5) visual question answering	1) visual question answering 2) natural language visual reasoning
Works in Progress	VisualBERT (Li et al. 2019b)	single cross-modal Transformer	image RoI	COCO Caption (Chen et al. 2015)	1) sentence-image alignment 2) masked language modeling	1) visual question answering 2) visual commonsense reasoning 3) natural language visual reasoning 4) grounding phrases
	Unicoder-VL (Li et al. 2019a)	single cross-modal Transformer	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification	1) image-text retrieval 2) zero-shot image-text retrieval
	Our VL-BERT	single cross-modal Transformer	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification	1) visual question answering 2) visual commonsense reasoning 3) grounding referring expressions

‡ LXMERT is pre-trained on COCO Caption (Chen et al. 2015), VG Caption (Krishna et al. 2017), VG QA (Zhu et al. 2016), VQA (Antol et al. 2015) and GQA (Hudson & Manning 2019).

Concurrent Work

	Method	Architecture	Visual Token	Pre-train Datasets	Pre-train Tasks	Downstream Tasks
Published Works	VideoBERT (Sun et al. 2019b)	single cross-modal Transformer	video frame	Cooking312K (Sun et al. 2019b)	1) sentence-image alignment 2) masked language modeling 3) masked visual-words prediction	1) zero-shot action classification 2) video captioning
Works Under Review / Just Got Accepted	CBT (Sun et al. 2019a)	two single-modal Transformer (vision & language respectively) + one cross-modal Transformer	video frame	Cooking312K (Sun et al. 2019b)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature regression	1) action anticipation 2) video captioning
	ViLBERT (Lu et al. 2019)	one single-modal Transformer (language) + one cross-modal Transformer (with restricted attention pattern)	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification	1) visual question answering 2) visual commonsense reasoning 3) grounding referring expressions 4) image retrieval 5) zero-shot image retrieval
	B2T2 (Alberti et al. 2019)	single cross-modal Transformer	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling	1) visual commonsense reasoning
	LXMERT (Hao Tan, 2019)	two single-modal Transformer (vision & language respectively) + one cross-modal Transformer	image RoI	‡ COCO Caption + VG Caption + VG QA + VQA + GQA	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification 4) masked visual-feature regression 5) visual question answering	1) visual question answering 2) natural language visual reasoning
Works in Progress	VisualBERT (Li et al. 2019b)	single cross-modal Transformer	image RoI	COCO Caption (Chen et al. 2015)	1) sentence-image alignment 2) masked language modeling	1) visual question answering 2) visual commonsense reasoning 3) natural language visual reasoning 4) grounding phrases
	Unicoder-VL (Li et al. 2019a)	single cross-modal Transformer	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification	1) image-text retrieval 2) zero-shot image-text retrieval
	Our VL-BERT	single cross-modal Transformer	image RoI	Conceptual Captions (Sharma et al. 2018)	1) sentence-image alignment 2) masked language modeling 3) masked visual-feature classification	1) visual question answering 2) visual commonsense reasoning 3) grounding referring expressions

‡ LXMERT is pre-trained on COCO Caption (Chen et al. 2015), VG Caption (Krishna et al. 2017), VG QA (Zhu et al. 2016), VQA (Antol et al. 2015) and GQA (Hudson & Manning 2019).

Summary

Summary

Task-agnostic visiolinguistic representations pretraining for visual grounding

- Introduce *pretrain-transfer* to vision and language tasks.
- Achieve SOTA on multiple vision and language tasks.

Summary

Task-agnostic visiolinguistic representations pretraining for visual grounding

- Introduce *pretrain-transfer* to vision and language tasks.
- Achieve SOTA on multiple vision and language tasks.

Limitations

The model can still learn inconsistent grounding by task specific finetuning.

- Training multiple vision and language task together.



Recent & Future Work

Summary

Task-agnostic visiolinguistic representations pretraining for visual grounding

- Introduce *pretrain-transfer* to vision and language tasks.
- Achieve SOTA on multiple vision and language tasks.

Limitations

The model can still learn inconsistent grounding by task specific finetuning.

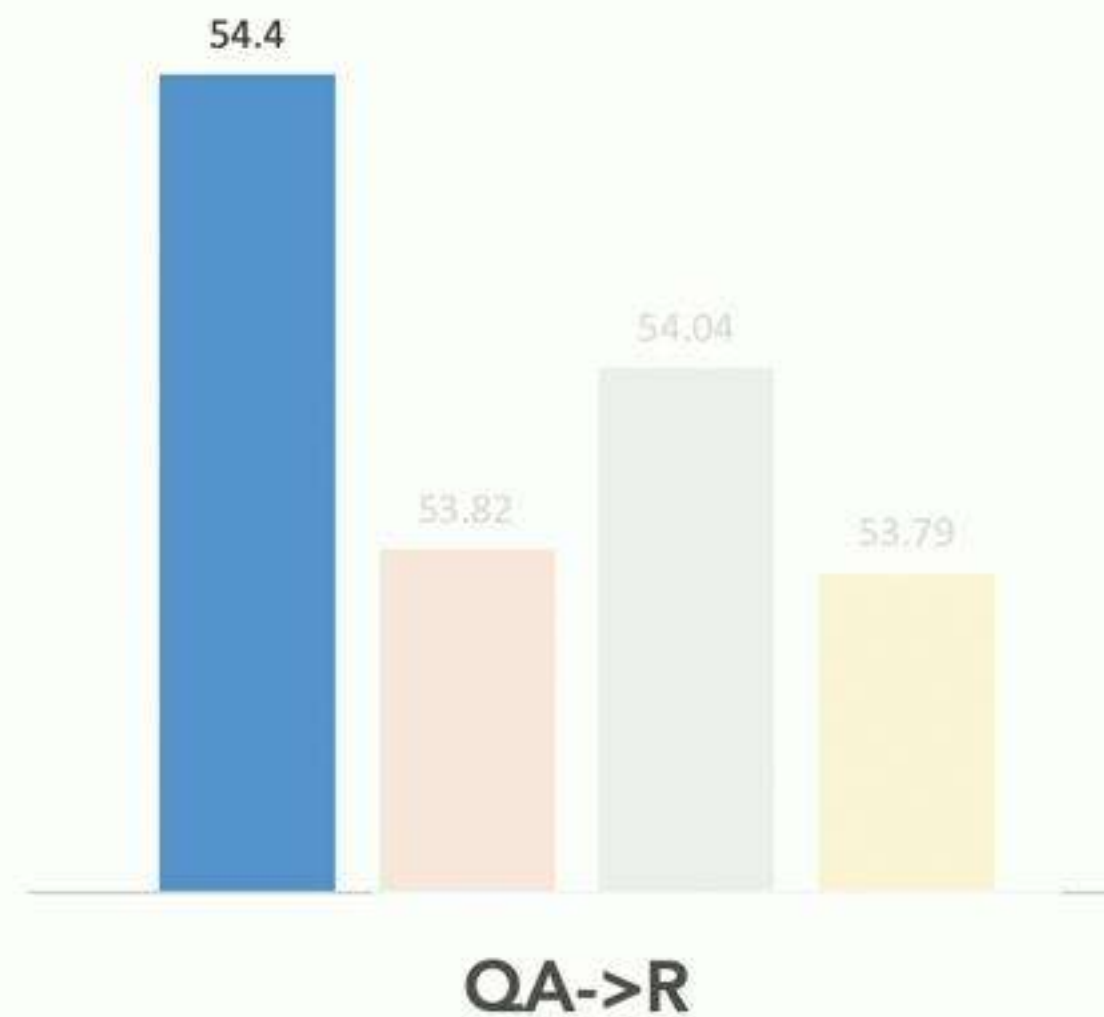
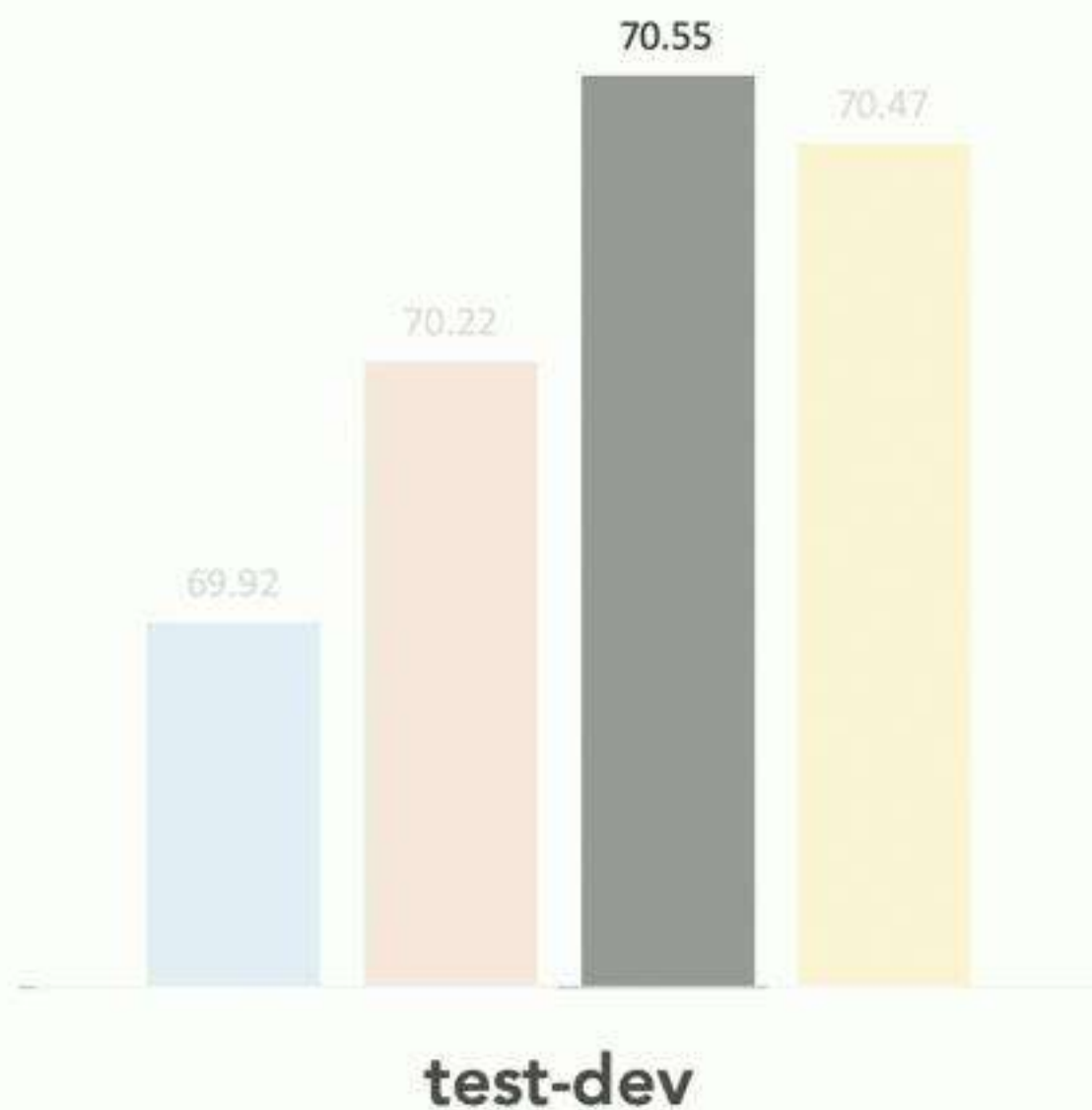
- Training multiple vision and language task together.

Ablations - Depth

■ 2-layer ■ 4-layer ■ 6-layer ■ 8-layer

VQA

VCR





Recent & Future Work

VISION AND LANGUAGE MULTI-TASK LEARNING

VISION AND LANGUAGE MULTI-TASK LEARNING

ViLBERT Problems:

- Inconsistent grounding by task specific finetuning.
- Four V&L tasks.
- Model is huge, overfitting.

VISION AND LANGUAGE MULTI-TASK LEARNING

ViLBERT Problems:

- Inconsistent grounding by task specific finetuning.
- Four V&L tasks.
- Model is huge, overfitting.

What we want:

- More tasks.
- Consistent Grounding.
- Explore the limit of the model.

VISION AND LANGUAGE MULTI-TASK LEARNING

ViLBERT Problems:

- Inconsistent grounding by task specific finetuning.
- Four V&L tasks.
- Model is huge, overfitting.

One Model for ALL V&L:

What we want:

- More tasks.
- Consistent Grounding.
- Explore the limit of the model.

VISION AND LANGUAGE MULTI-TASK LEARNING

ViLBERT Problems:

- Inconsistent grounding by task specific finetuning.
- Four V&L tasks.
- Model is huge, overfitting.

What we want:

- More tasks.
- Consistent Grounding.
- Explore the limit of the model.

One Model for ALL V&L:

VQA

- VQA
- Genome QA
- GQA

Image Description

- Caption based Retrieval (COCO)
- Caption based Retrieval (COCO)

Referring Expression

- Ref COCO
- Ref COCO+
- Ref COCOg
- Visual 7w
- GuessWhat

V&L Verification

- NLVR2
- Visual Entailment

VISION AND LANGUAGE MULTI-TASK LEARNING

One Model for ALL V&L:

VQA

- VQA
- Genome QA
- GQA

Image Description

- Caption based Retrieval (COCO)
- Caption based Retrieval (COCO)

Referring Expression

- Ref COCO
- Ref COCO+
- Ref COCOg
- Visual 7w
- GuessWhat

V&L Verification

- NLVR2
- Visual Entailment

VISION AND LANGUAGE MULTI-TASK LEARNING

One Model for ALL V&L:

VQA

- VQA
- Genome QA
- GQA

Image Description

- Caption based Retrieval (COCO)
- Caption based Retrieval (COCO)

Referring Expression

- Ref COCO
- Ref COCO+
- Ref COCOg
- Visual 7w
- GuessWhat

V&L Verification

- NLVR2
- Visual Entailment

- Benchmark ALL vision and language understanding tasks with ViLBERT.

VISION AND LANGUAGE MULTI-TASK LEARNING

One Model for ALL V&L:

VQA

- VQA
- Genome QA
- GQA

Image Description

- Caption based Retrieval (COCO)
- Caption based Retrieval (COCO)

Referring Expression

- Ref COCO
- Ref COCO+
- Ref COCOg
- Visual 7w
- GuessWhat

V&L Verification

- NLVR2
- Visual Entailment

- Benchmark ALL vision and language understanding tasks with ViLBERT.
- Study inter-connections within task group and between task group.

VISION AND LANGUAGE MULTI-TASK LEARNING

One Model for ALL V&L:

VQA

- VQA
- Genome QA
- GQA

Image Description

- Caption based Retrieval (COCO)
- Caption based Retrieval (COCO)

Referring Expression

- Ref COCO
- Ref COCO+
- Ref COCOg
- Visual 7w
- GuessWhat

V&L Verification

- NLVR2
- Visual Entailment

- Benchmark ALL vision and language understanding tasks with ViLBERT.
- Study inter-connections within task group and between task group.
- Explainable AI : Use other task outputs to provide explanations.

DIALOG WITHOUT DIALOG

1:



2:



3:



4:



Ongoing Work

DIALOG WITHOUT DIALOG

1:



2:



3:



4:



DIALOG WITHOUT DIALOG

1:



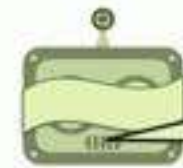
2:



3:



4:



Is the child laying or sitting?

DIALOG WITHOUT DIALOG

1:



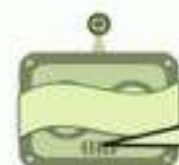
2:



3:



4:



Is the child laying or sitting?

laying



DIALOG WITHOUT DIALOG

1:



2:



3:

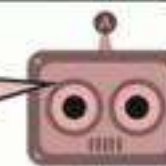


4:



Is the child laying or sitting?

laying



P: 4

DIALOG WITHOUT DIALOG

1:



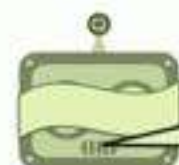
2:



3:

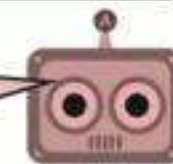


4:

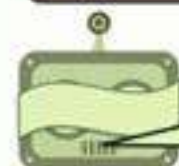


Is the child laying or sitting?

laying



P: 4



What color is the blanket?

DIALOG WITHOUT DIALOG

1:



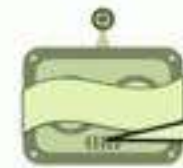
2:



3:

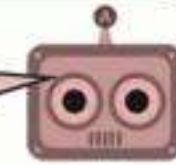


4:

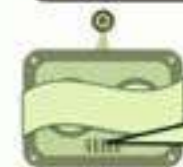


Is the child laying or sitting?

laying

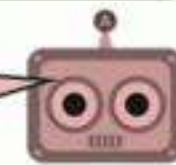


P: 4



What color is the blanket?

red



DIALOG WITHOUT DIALOG

1:



2:



3:

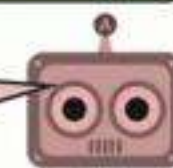


4:

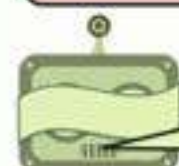


Is the child laying or sitting?

laying

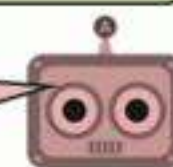


P: 4



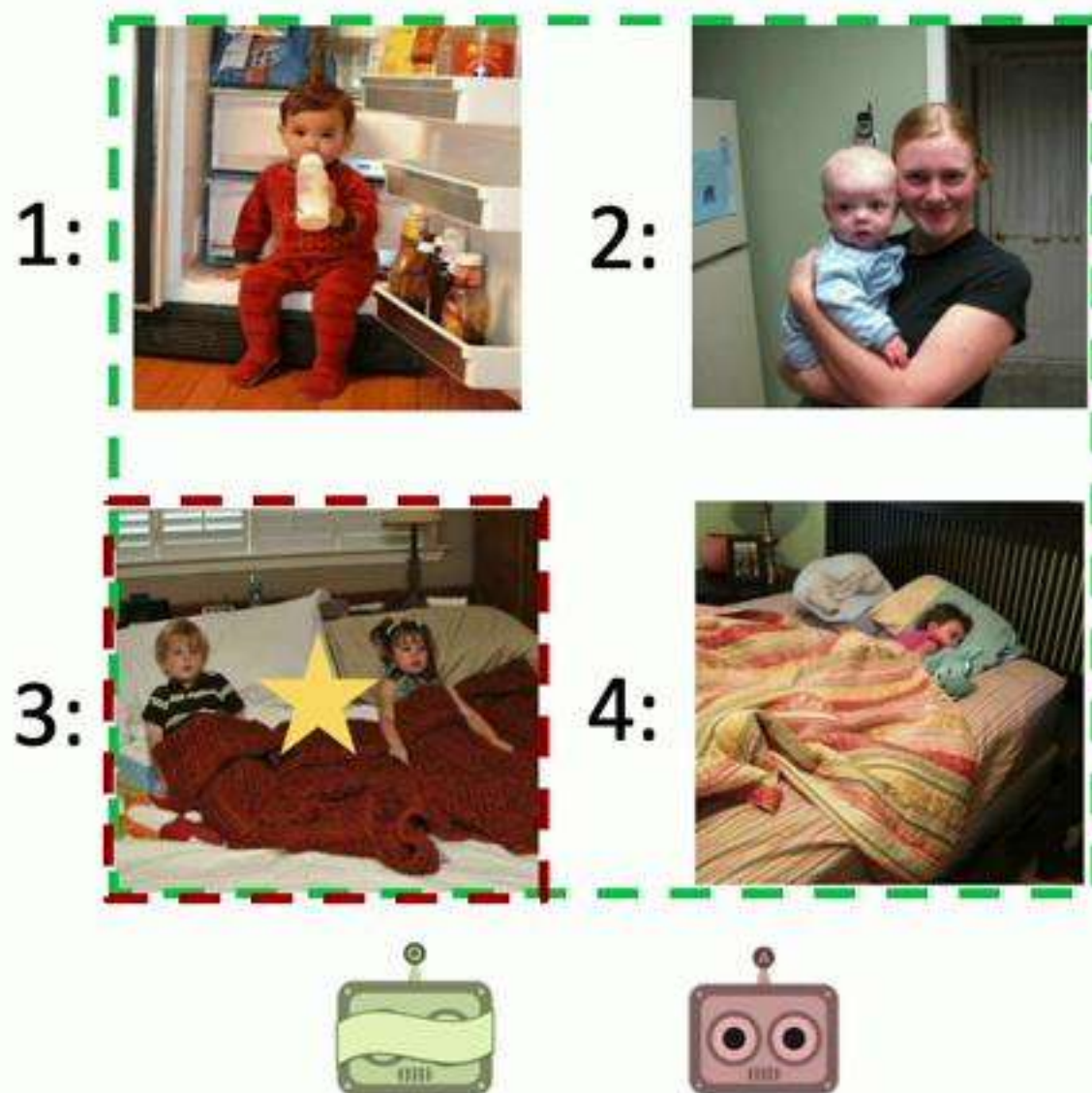
What color is the blanket?

red

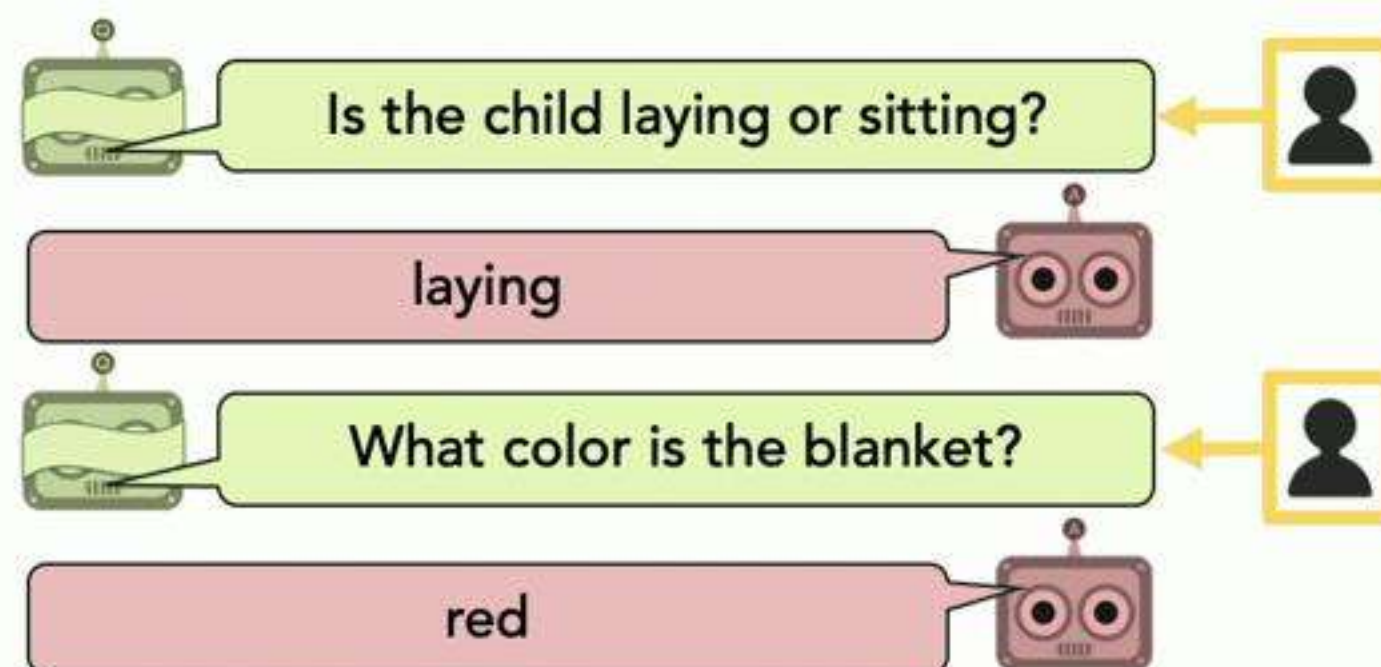


P: 3

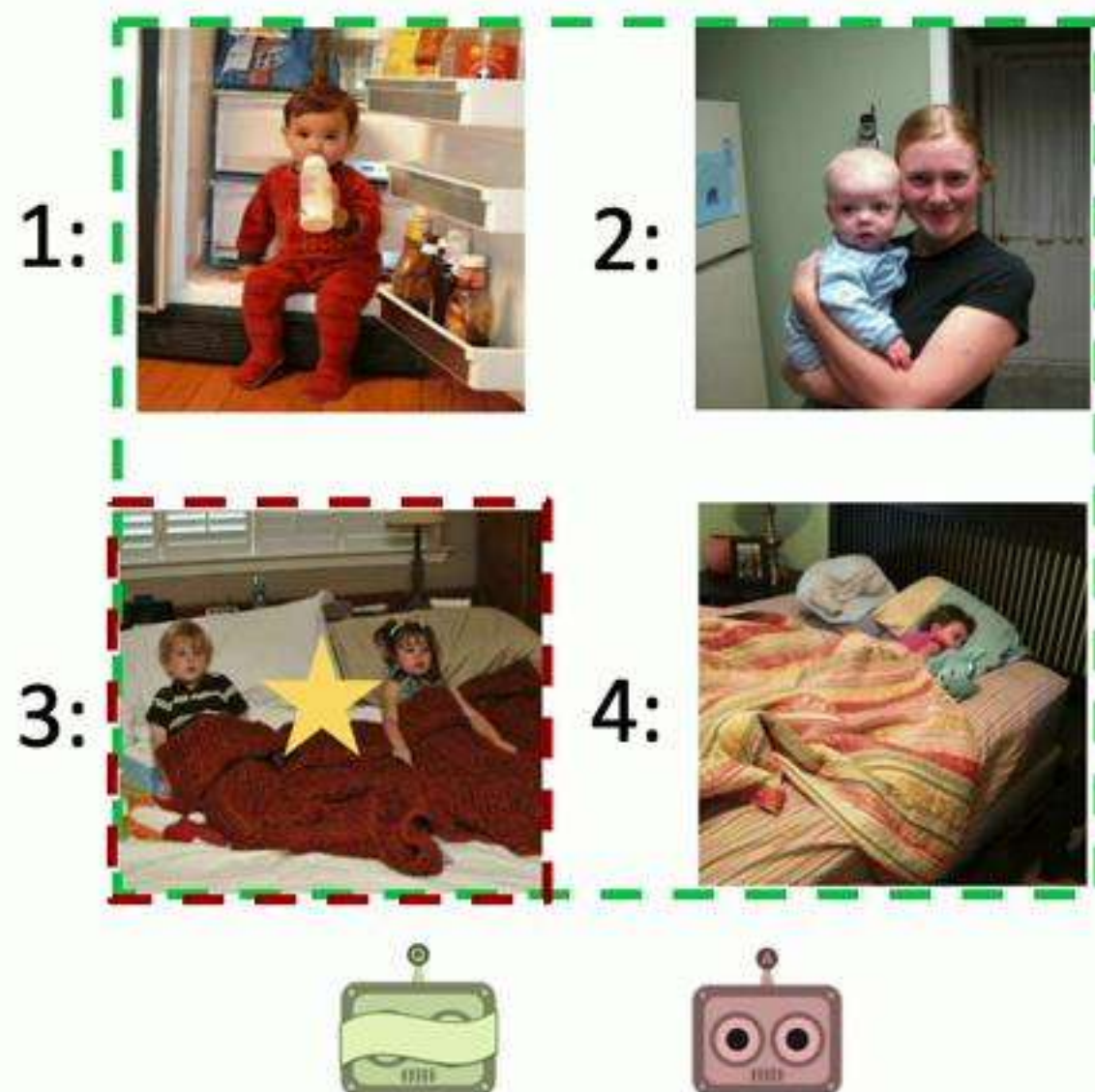
DIALOG WITHOUT DIALOG



How does Q-Bot know what to say?



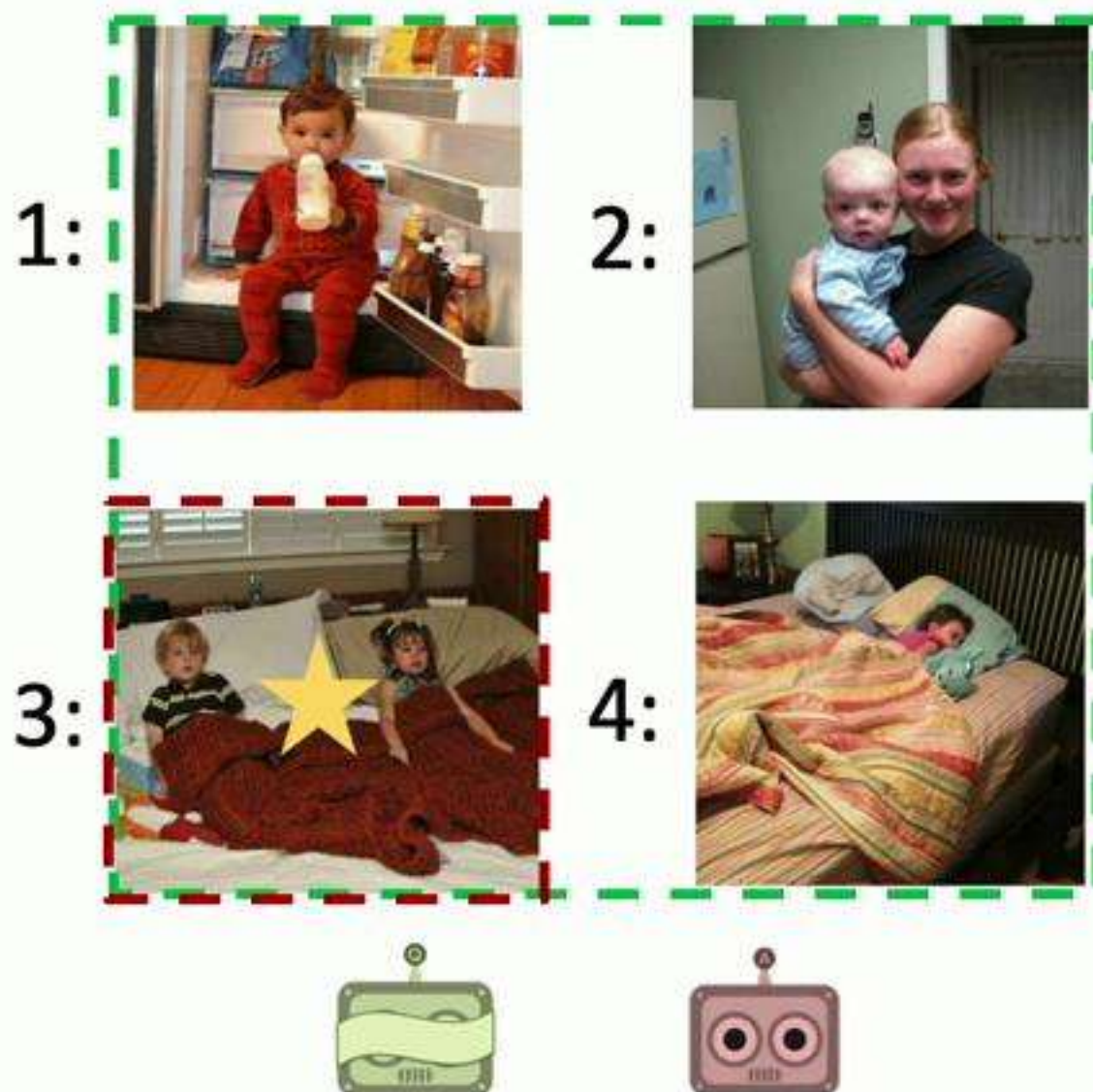
DIALOG WITHOUT DIALOG



How does Q-Bot know what to say?



DIALOG WITHOUT DIALOG



How does Q-Bot know what to say?



Goals (with sparse reward):

1. Language fluency
2. Task performance

Ongoing Work

Language Fluency from VQA

Who is wearing glasses?

man



woman



Is the umbrella upside down?

yes



no



Learn how to ask from VQA

Language Fluency from VQA

Who is wearing glasses?
man woman



Is the umbrella upside down?
yes no



Learn how to ask from VQA

Key Insight: LM with discrete latent action space

Interpolation

1. how many beds ?

10. where is he looking ?

Language Fluency from VQA

Who is wearing glasses?
man woman



Is the umbrella upside down?
yes no



Learn how to ask from VQA

Key Insight: LM with discrete latent action space

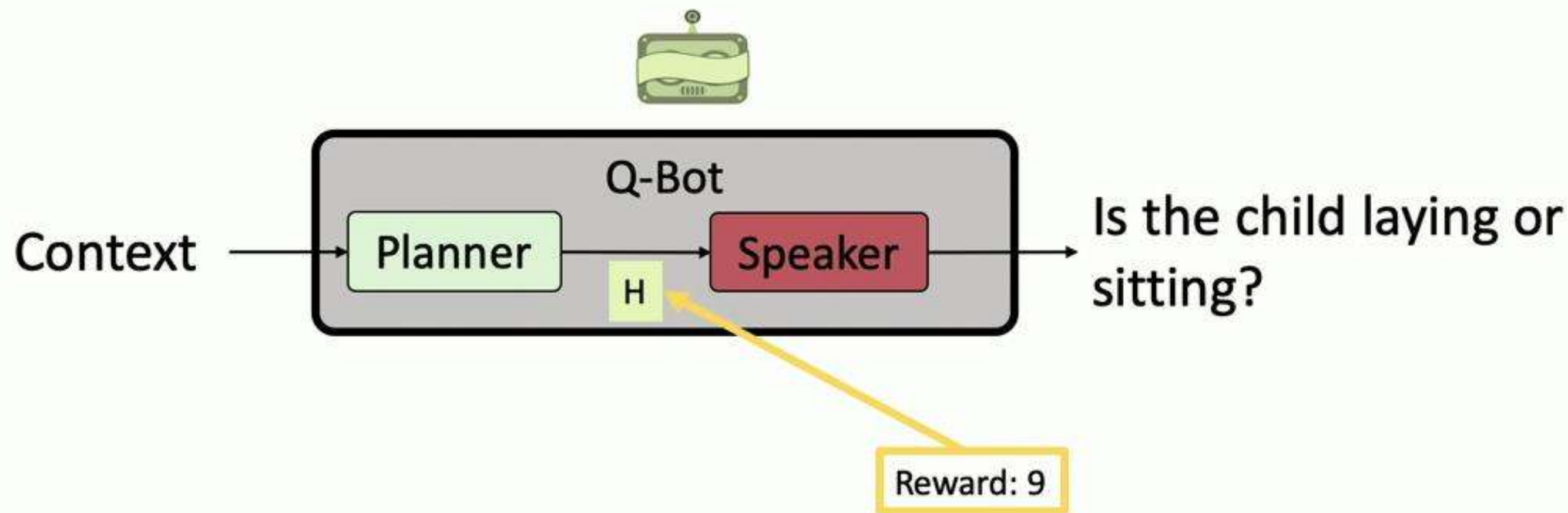
Interpolation

1. how many beds ?
2. how many cats ?
3. how many dogs ?
4. where is the dog ?
5. where is the man ?
6. where is this man ?
7. where is this woman ?
8. where is this ?
9. where is he ?
10. where is he looking ?

Ongoing Work

Task Performance via Fine-Tuning without Forgetting Language

Key Insight: Reinforce discrete latent action space.



Ongoing Work

Initial Results

Initial Results



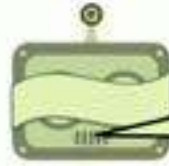
target



Initial Results



target



How many pets are in the photo?

Initial Results

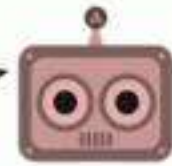


target



How many pets are in the photo?

1



Initial Results

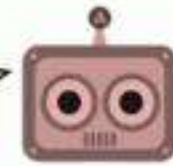


target



How many pets are in the photo?

1

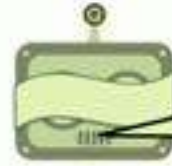


P: 1

Initial Results

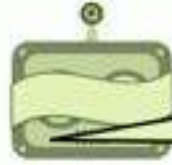
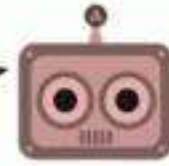


target



How many pets are in the photo?

1



What color is underbelly beneath the animal?

P: 1

Initial Results

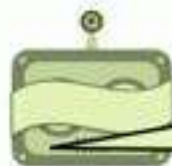
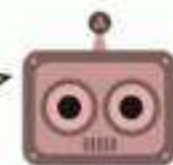


target



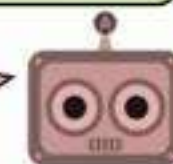
How many pets are in the photo?

1



What color is underbelly beneath the animal?

brown



P: 1

Initial Results

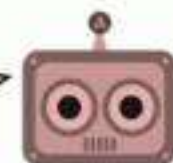


target



How many pets are in the photo?

1

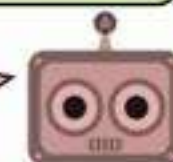


P: 1



What color is underbelly beneath the animal?

brown



P: 1

Initial Results

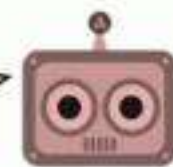


target



How many pets are in the photo?

1



P: 1



What color is underbelly beneath the animal?

brown



P: 1



What color is the zebra's fur ?

Initial Results

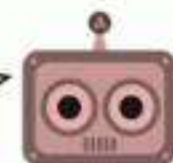


target

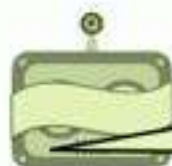


How many pets are in the photo?

1

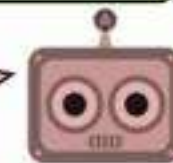


P: 1



What color is underbelly beneath the animal?

brown

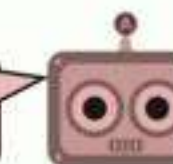


P: 1



What color is the zebra's fur ?

Black and white



Initial Results

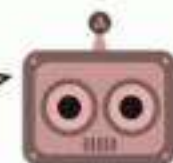


target

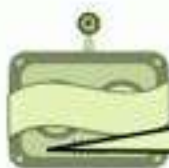


How many pets are in the photo?

1

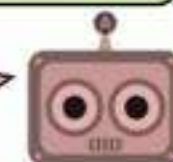


P: 1



What color is underbelly beneath the animal?

brown

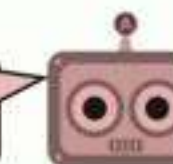


P: 1



What color is the zebra's fur ?

Black and white



P: 2

Initial Results

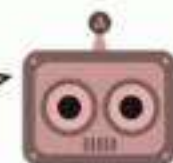


target

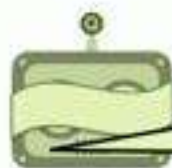


How many pets are in the photo?

1

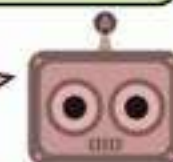


P: 1



What color is underbelly beneath the animal?

brown

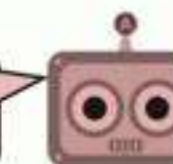


P: 1



What color is the zebra's fur ?

Black and white



P: 2



Is the animal looking at the camera?

Multi-modality representation learning in vision, language, sound and action.

Multi-modality representation learning in vision, language, sound and action.

- Environment around us is not un-modal

Multi-modality representation learning in vision, language, sound and action.

- Environment around us is not un-modal

Table:

Multi-modality representation learning in vision, language, sound and action.

- Environment around us is not un-modal

Table:



Multi-modality representation learning in vision, language, sound and action.

- Environment around us is not un-modal

Table:



Multi-modality representation learning in vision, language, sound and action.

- Environment around us is not un-modal

Table:



Providing a level surface on which objects may be placed

Multi-modality representation learning in vision, language, sound and action.

- Environment around us is not un-modal

Table:



Providing a level surface on which objects may be placed

- Learn a representation that connect all the modalities.
 - More complete understanding of the environment.

Emerging of goal orientated dialog with natural language.

- Dialog is the most natural way to communicate.

Emerging of goal orientated dialog with natural language.

- Dialog is the most natural way to communicate.
 - Impossible to collect dialog data for all the dialog tasks.

Emerging of goal orientated dialog with natural language.

- Dialog is the most natural way to communicate.
 - Impossible to collect dialog data for all the dialog tasks.
- Emerging of language (machine code) is not useful.

Emerging of goal orientated dialog with natural language.

- Dialog is the most natural way to communicate.
 - Impossible to collect dialog data for all the dialog tasks.
- Emerging of language (machine code) is not useful.
 - People can not understand them

Emerging of goal orientated dialog with natural language.

- Dialog is the most natural way to communicate.
 - Impossible to collect dialog data for all the dialog tasks.
- Emerging of language (machine code) is not useful.
 - People can not understand them
- Emerging of strategy (language that human can understand)

Common sense abstraction and causal reasoning.

- Intrinsic “motivation” behind the “scene”.

Common sense abstraction and causal reasoning.

- Intrinsic “motivation” behind the “scene”.
- Emerging of the commonsense and reasoning of the “motivation”

Common sense abstraction and causal reasoning.

- Intrinsic “motivation” behind the “scene”.
- Emerging of the commonsense and reasoning of the “motivation”
- More human-like and interpretable agent which helps people in the daily life.

End

QUESTIONS?