# Incorporating External Information for Visual Question Answering

Jialin Wu

UT Austin

Committee members: Raymond Mooney, Gregory Durrett, David Harwath, Dhruv Batra, Roozbeh Mottaghi

# Visual Question Answering
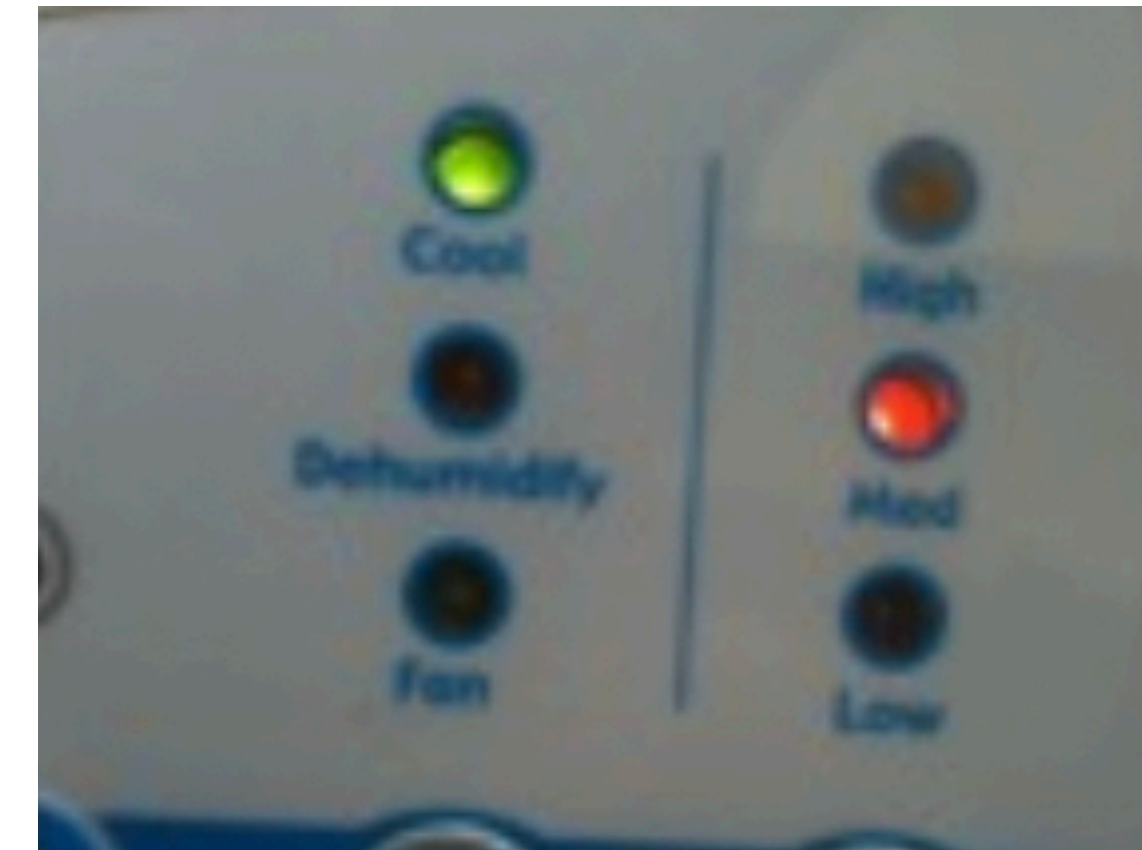
Multimodal and open-ended task

Helping the visually impaired

Does this boy have a full wetsuit on?

Is this air conditioner on fan, dehumidifier or air conditioning?



[Antol et al., CVPR 2015]



Image from CVPR 2022 tutorial



[Gurari et al., ECCV 2018]

# VQA Models

# Structural Visual Features

Question: Does this boy have a full wetsuit on?



Answer : Yes

Objects: man, wetsuit, wave, man, surf board, …

Attributes: man is young, wetsuit is full, wave is small, …

Relationships: wetsuit on man, man riding wave, …

[Anderson et al., CVPR2018; Wu et al., ACL2019; Lu et al., CVPR2020; Cho et al., ICML2021;… ]

[Li et al., ICCV2019; Hudson et al., CVPR2019;… ]

4

# Talk Outline

Core Contributions:

| Info | Image Captions | Human Explanations | Open Knowledge | Open Knowledge |
|---|---|---|---|---|
| Topic | General VQA | Fighting Imbalanced Training Distribution | Outside Knowledge VQA | Knowledge Retrieval |

Future Directions

# Talk Outline

Core Contributions:

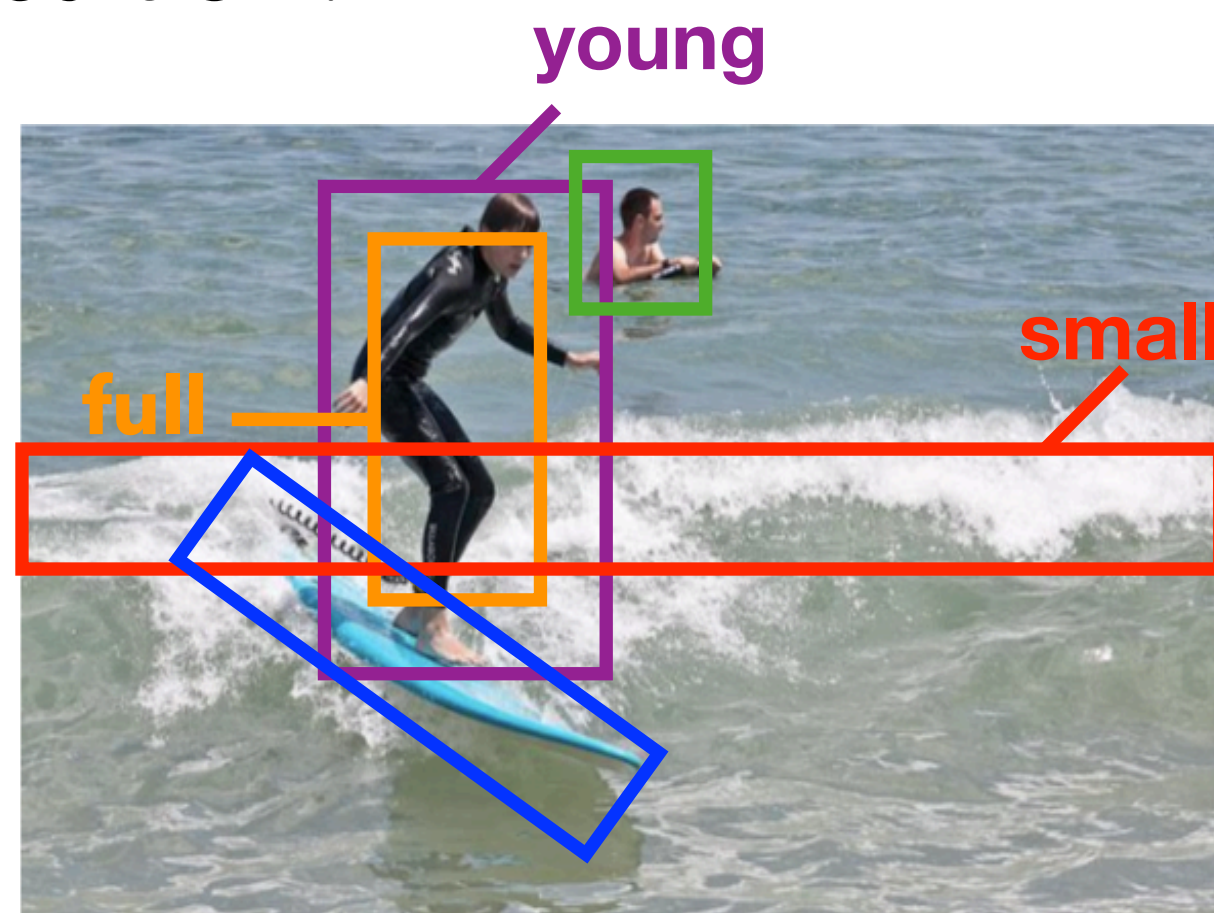| Info | Image Captions | Human Explanations | Open Knowledge | Open Knowledge |
|---|---|---|---|---|
| Topic | General VQA | Fighting Imbalanced Training Distribution | Outside Knowledge VQA | Knowledge Retrieval |

Future Directions

6

# Captions vs Structural Visual Features

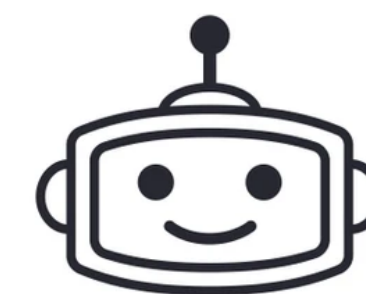Question: Does this boy have a full wetsuit on?



Answer : Yes

Objects: man, wetsuit, wave

Relationships: man in wetsuit

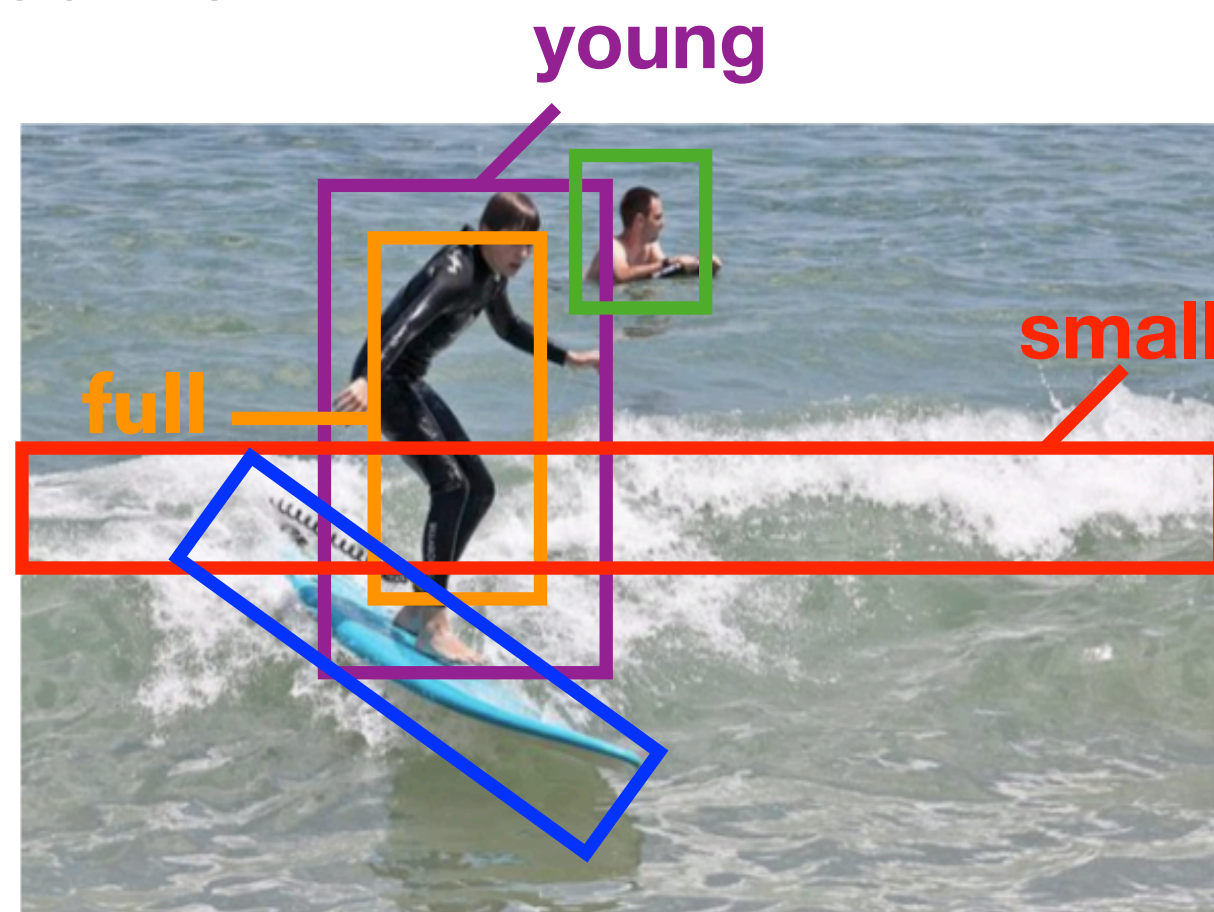Attributes: man is young, wetsuit is full, wave is small

Caption: A young surfer in a wetsuit surfs a small wave.

Succinct & Flexible

# Captions as Context

Question: Does this boy have a full wetsuit on?



Answer : Yes

A man on a blue surfboard on top of some rough water.

A young surfer in a wetsuit surfs a small wave.

A young man rides a surf board on a small wave while a man swims in the background.
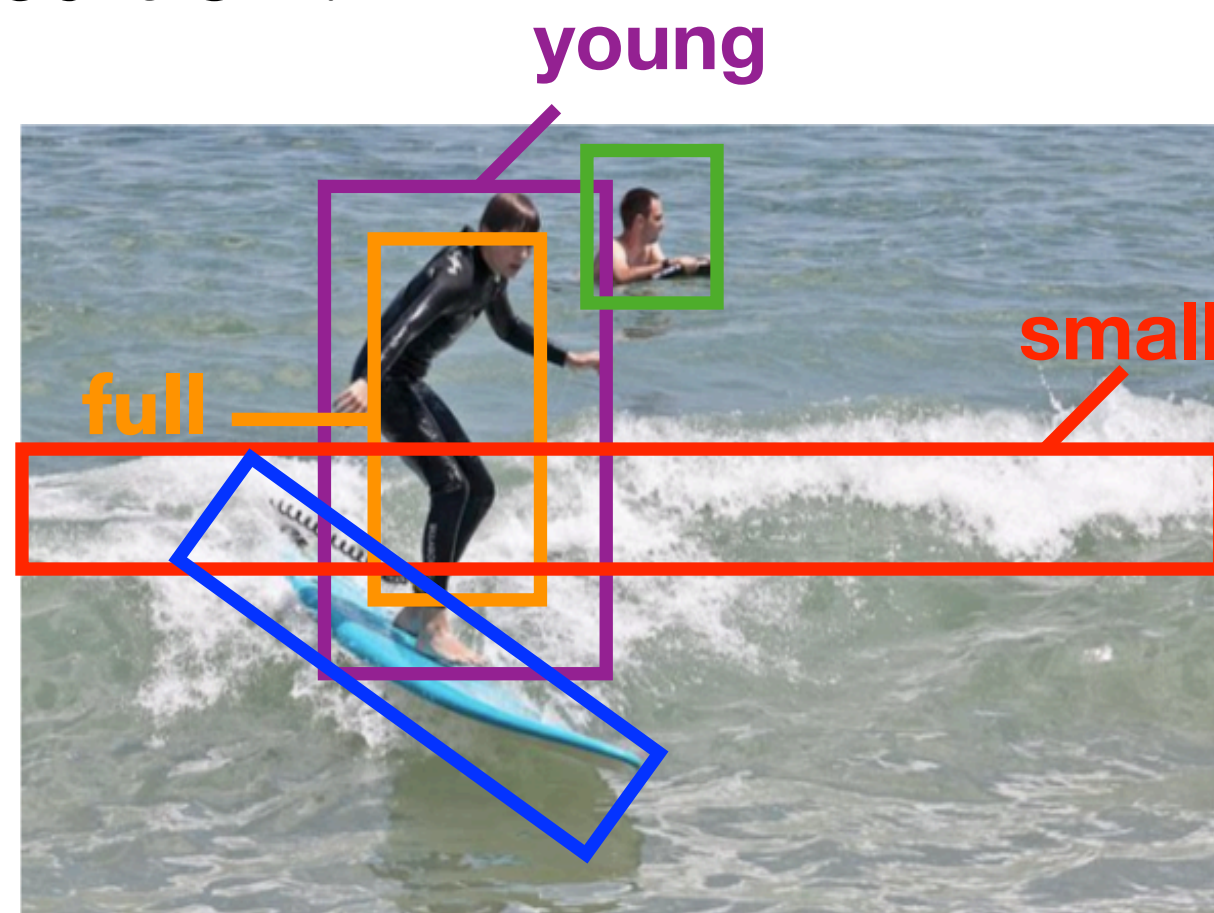
A young man is on his surf board with someone in the background.

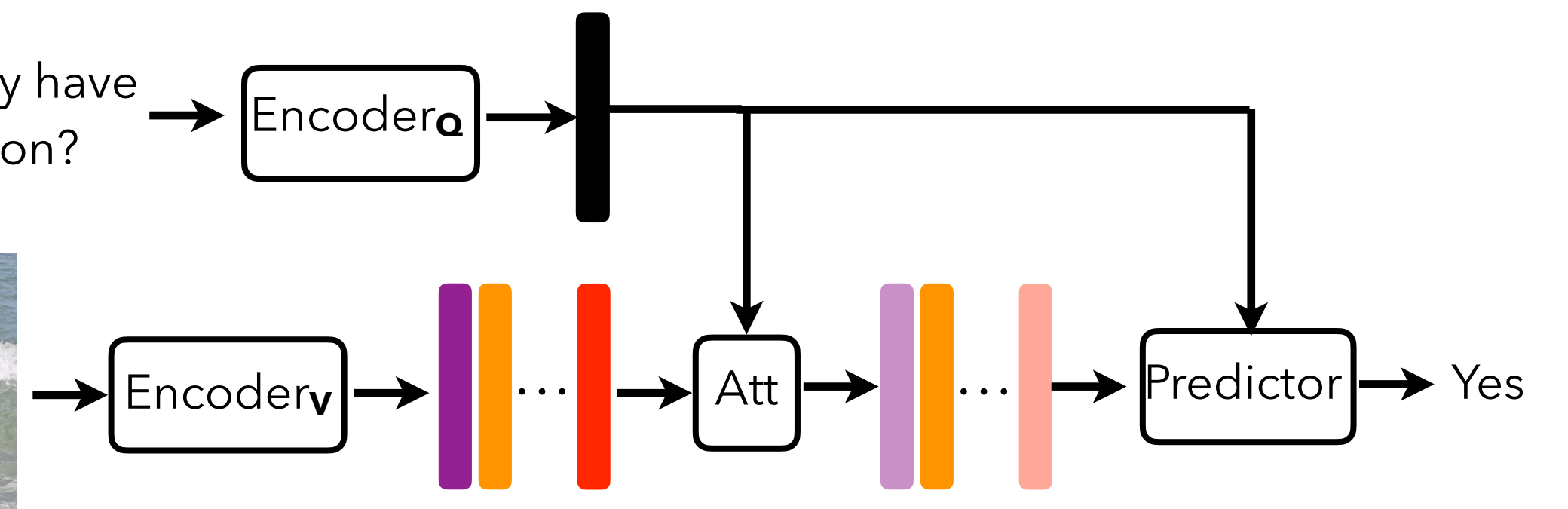A boy riding waves on his surf board in the ocean.

# Up-Down VQA Model

Question: Does this boy have a full wetsuit on?



Answer : Yes

[Anderson et al., CVPR2018]

9

# Relevant Captions



Question: Does this boy have a full wetsuit on?

Visual explanations

Answer : Yes

A young surfer in a **wetsuit** surfs a small wave.

Question: What color is the surfboard?

Visual explanations

Answer : Blue

A man on a **blue surfboard** on top of some rough water.

# Combining Generated Captions



Question: Does this boy have a full wetsuit on?
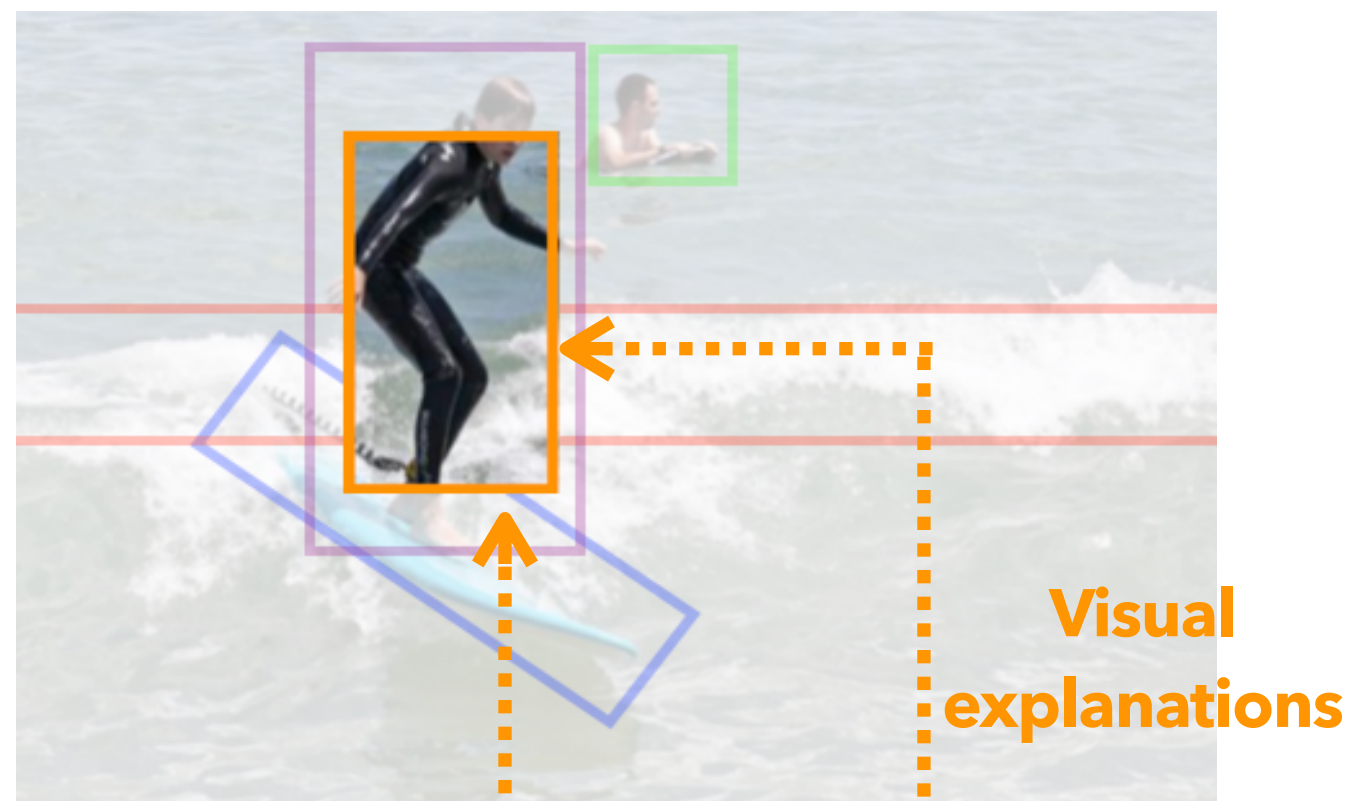
Answer : Yes

Q: Does this boy have a full wetsuit on?

GC1:A young surfer is riding a wave wearing black wetsuit.

GC2:A young man is surfing in wetsuit.

# Results on VQA v2



**VQA v2**[Goyal et al. CVPR2017]: a large-scale balanced dataset, 443k questions for training, 214k questions for validation and 447k questions for test

**Up-Down**[Anderson et al. CVPR 2018]: base model

**VQA-E**[Li et al. ECCV 2018]:Jointly generating explanations and VQA

**+IC**(ours):Using generated image captions

**+RC**(ours):Using generated question-relevant captions

**+HC**(ours):Using human captions

Using captions is an efficient way to provide visual context and question-relevant captions are better.

13

# Talk Outline

Core Contributions:

| Info | Image Captions | Human Explanations | Open Knowledge | Open Knowledge |
|------|----------------|--------------------|----------------|----------------|
| Topic | General VQA | **Fighting Imbalanced Training Distribution** | Outside Knowledge VQA | Knowledge Retrieval |

Future Directions

# VQA under Changing Priors



Question: What utensil is pictured?

Answer : Fork

No need to see the image, utensil means knife

Q: What utensil is pictured?

Encoder

VQA Model

Fork

Knife 90

Fork 10

[Agrawal et al., CVPR2018]

15

# Right for the Right Reasons

# Where to Focus

Question: What utensil is pictured?



Answer : Fork

There is a **fork** near the **cake**.

Textual Explanation

Visual Explanation

What **utensil** is pictured?

Question

# Learning to Focus

Question: What utensil is pictured?

At least one object is influential

Q: What utensil is pictured?

VQA Model → **Fork**

Answer : Fork

47.8

39.7

Up-Down

Condition 1

# Knowing Where to Focus

# What is the Focused One?

# Results on VQA-CP v2



**VQA-CP v2**[Agrawal et al. CVPR2018]: a reorganizing VQA v2 such that the answer distribution differs in the training and test splits.

**Up-Down**[Anderson et al. CVPR 2018]: base model

**AdvReg**[Ramakrishnan et al. NeurIPS 2018]: Adversarial Regularization on question only model

**RuBi**[Cadene et al. NeurIPS 2019]:Reducing Unimodal Biases by weight the output

**HINT**[Selvaraju et al. ICCV 2019]: Human Importance-aware Network Tuning

**+Q**(ours):Using question words

**+VE**(ours):Using visual explanation

**+TE**(ours):Using textual explanation

Human Explanations help VQA models be robust to unbalanced training distribution

21

# Talk Outline

Core Contributions:

| Info | Image Captions | Human Explanations | **Open Knowledge** | Open Knowledge |
|------|---------------|--------------------|--------------------|----------------|
| Topic | General VQA | Fighting Imbalanced Training Distribution | **Outside Knowledge VQA** | Knowledge Retrieval |

Future Directions

# Outside-Knowledge VQA

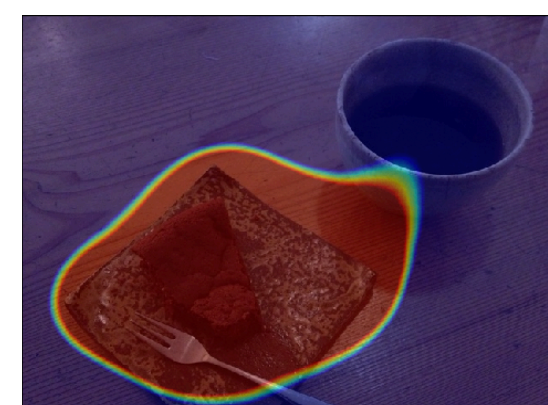Question: Which American president is associated with the stuffed animal seen here?



[Kenneth et al., CVPR2019]

# Outside-Knowledge VQA

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



*The name teddy bear comes from former United States President **Theodore Roosevelt**, who was often referred to as "Teddy"*

A: Teddy Roosevelt

[Kenneth et al., CVPR2019]

# (A) Factual Knowledge

Q: Which movie featured a man in this position telling his life story to strangers?





*Forrest Gump narrated his life's story at the northern edge of Chippewa Square in Savannah, Georgia, as he sat at a bus stop bench.*



A: Forrest Gump

# (B) Conceptual Knowledge



Q: Is this a healthy dish?



*Vegetarian food has property of healthy*
*Eating Vegetable has property of healthy*
*Beans are related to healthy*

**ConceptNet**
An open, multilingual knowledge graph

A: Yes

# (C) Visual Knowledge

Q: What breed of dog is the dog in this photo?

Golden Retriever

Google Images

A: Golden Retriever

# Is Question and Image Sufficient?

Q: What English city is famous for a tournament for the sport this man is playing?



Question + Image

> The modern game of tennis originated in Birmingham, England in the late 19th century as lawn tennis.

Question + Image + Answer (Wimbledon)

> **Wimbledon** is notable for the longest sponsorship in sports history due to its association with Slazenger who have supplied all tennis balls for the tournament since 1902.

# Fitting in the Question Context

Q: What **English city** is **famous for** a **tournament** for the **sport this man is playing**?



Question + Image

> The modern game of tennis originated in Birmingham, England in the late 19th century as lawn tennis.

Question + Image + Answer (Wimbledon)

> **Wimbledon** is **notable for** the longest sponsorship in sports history due to its association with Slazenger who have supplied all **tennis balls** for the **tournament** since 1902.

# Fitting in the Question Context

Q: What **English city** is **famous for** a **tournament** for the **sport this man is playing**?



**Wimbledon** is **notable for** the longest sponsorship in sports history due to its association with Slazenger who have supplied all **tennis balls** for the **tournament** since 1902.

Question + Image answer (Wimbledon)

# Multimodal Answer Validation (MAVEx) Outline

- **MAVEx Answer-Guided Knowledge Retrieval**

- MAVEx Answer Validation Model

- Experimental Results

# MAVEx Retrieval Outline

Parsing entities

↓

Retrieving knowledge

↓

Matching knowledge to entities

# Extracting Noun Phrases

Q: Which movie featured a man in this position telling his life story to strangers?

S1: Forrest Gump featured a man in this position…
S2: Speed featured a man in this position telling…

⋮





33

# Parsing Entities

Q: Which movie featured a man in this position telling his life story to strangers?





Queries

Movie

34

# Parsing Entities

Q: Which movie featured a man in this position telling his life story to strangers?



Queries

Movie
Man
Sitting man



35

# Parsing Entities

Q: Which movie featured a man in this position telling his life story to strangers?



Queries

Movie
Man
Sitting man
Life story
Stranger



36

# Parsing Entities

S1: Forrest Gump featured a man in this position telling …
S2: Speed featured a man in …





Queries

Movie
Man
Sitting man
Life story
Stranger
Forrest Gump
Gump
Speed
⋮

# Retrieving Articles

Queries

Movie
Man
Sitting man
Life story
Stranger
**Forrest Gump**
Gump
Speed
⋮

*Forrest Gump is a 1994 American comedy-drama film directed by Robert Zemeckis and written by Eric Roth. …… The Gump family home set was built along the Combahee River near Yemassee, South Carolina, and the nearby land was used to film Curran's home as well as some of the Vietnam scenes. ……. Forrest Gump narrated his life's story at the northern edge of Chippewa Square in Savannah, Georgia, as he sat at a bus stop bench.*

WIKIPEDIA
The Free Encyclopedia

Article    Talk

*Forrest Gump*

From Wikipedia, the free encyclopedia

# Retrieving Sentences



BERTScore

*Forrest Gump is a 1994 American comedy-drama film directed by Robert Zemeckis and written by Eric Roth.* …… *The Gump family home set was built along the Combahee River near Yemassee, South Carolina, and the nearby land was used to film Curran's home as well as some of the Vietnam scenes.* ……. *Forrest Gump narrated his life's story at the northern edge of Chippewa Square in Savannah, Georgia, as he sat at a bus stop bench.*

S1: Forrest Gump featured a man in this position…
S2: Speed featured a man in this position telling…

39

# Matching Knowledge

Queries

Movie
Man
Sitting man
Life story
Stranger
Forrest Gump
Gump
Speed
⋮

A man is an adult male human.

The novel also features Gump as an astronaut, a professional wrestler, and a chess player.

Forrest Gump narrated his life's story at the …, as he sat at a bus stop bench

Forrest Gump is a 1994 American comedy-drama film …

Speed is a 1994 American action thriller film directed by Jan de Bont …

WIKIPEDIA
The Free Encyclopedia

# Matching Knowledge

Queries

Movie
Man
Sitting man
Life story
Stranger
Forrest Gump
Gump
Speed
⋮

WIKIPEDIA
The Free Encyclopedia

A **man** is an adult **male human**.

The novel also features **Gump** as an astronaut, a professional wrestler, and a chess player.

**Forrest Gump** narrated his **life's story** at the …, as he **sat** at a bus stop bench

**Forrest Gump** is a 1994 American comedy-drama **film** …

**Speed** is a 1994 American action thriller **film** directed by Jan de Bont …

# Matching Knowledge

Queries

Movie
Man
Sitting man
Life story
Stranger
Forrest Gump
Gump
Speed
⋮

**ConceptNet**
An open, multilingual knowledge graph

A **Gentleman** is at a **movie**

**Forrest Gump** is a **film**

**Story** of my **life** is related to tell me about it

**Strangers** is related to people

**Speed** is related to **film**

# Matching Visual Knowledge

Q: Which movie featured a man in this position telling his life story to strangers?

# Matching Visual Knowledge

S1: Forrest Gump featured a man in this position telling …
S2: Speed featured a man in …

# MAVEx Retrieval Results

Q: Which movie featured a man in this position telling his life story to strangers?
S1: Forrest Gump featured a man in this position telling …
S2: Speed featured a man in …



**Queries**
Movie
Man
Sitting man
Life story
Stranger
Forrest Gump
Gump
Speed

WIKIPEDIA
The Free Encyclopedia

A **man** is an adult **male human**.

The novel also features **Gump** as an astronaut, a professional wrestler, and a chess player.

**Forrest Gump** narrated his **life's story** at the …, as he **sat** at a bus stop bench

**Forrest Gump** is a 1994 American comedy-drama **film** …

**Speed** is a 1994 American action thriller **film** directed by Jan de Bont …

**Queries**
Movie
Man
Sitting man
Life story
Stranger
Forrest Gump
Gump
Speed

ConceptNet
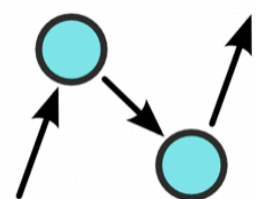An open, multilingual knowledge graph

A **Gentleman** is at a **movie**

**Forrest Gump** is a **film**

**Story** of my **life** is related to tell me about it

**Strangers** is related to people

**Speed** is related to **film**

ViLBERT

Google Images

45

# Multimodal Answer Validation

- MAVEx Answer-Guided Knowledge Retrieval

- **MAVEx Answer Validation Model**

- Experimental Results

46

# Answer Validation

Q: Which movie featured a man in this position telling his life story to strangers?



Answer1:
Forrest Gump
Answer2:
Speed

Queries
   Movie
   Man
   Sitting man
   Life story
   Stranger
   Forrest Gump
   Gump
   Speed

**WIKIPEDIA**
The Free Encyclopedia

A **man** is an adult **male human**.

The novel also features **Gump** as an astronaut, a professional wrestler, and a chess player.

**Forrest Gump** narrated his **life's story** at the …, as he **sat** at a bus stop bench

**Forrest Gump** is a 1994 American comedy-drama **film** …

**Speed** is a 1994 American action thriller **film** directed by Jan de Bont …

Answer-agnostic features

Answer specific features

Answer2: Speed

Answer1: Forrest Gump

Knowledge Module (Visual) — MLP — MLP

Knowledge Module (Concept) — MLP — MLP

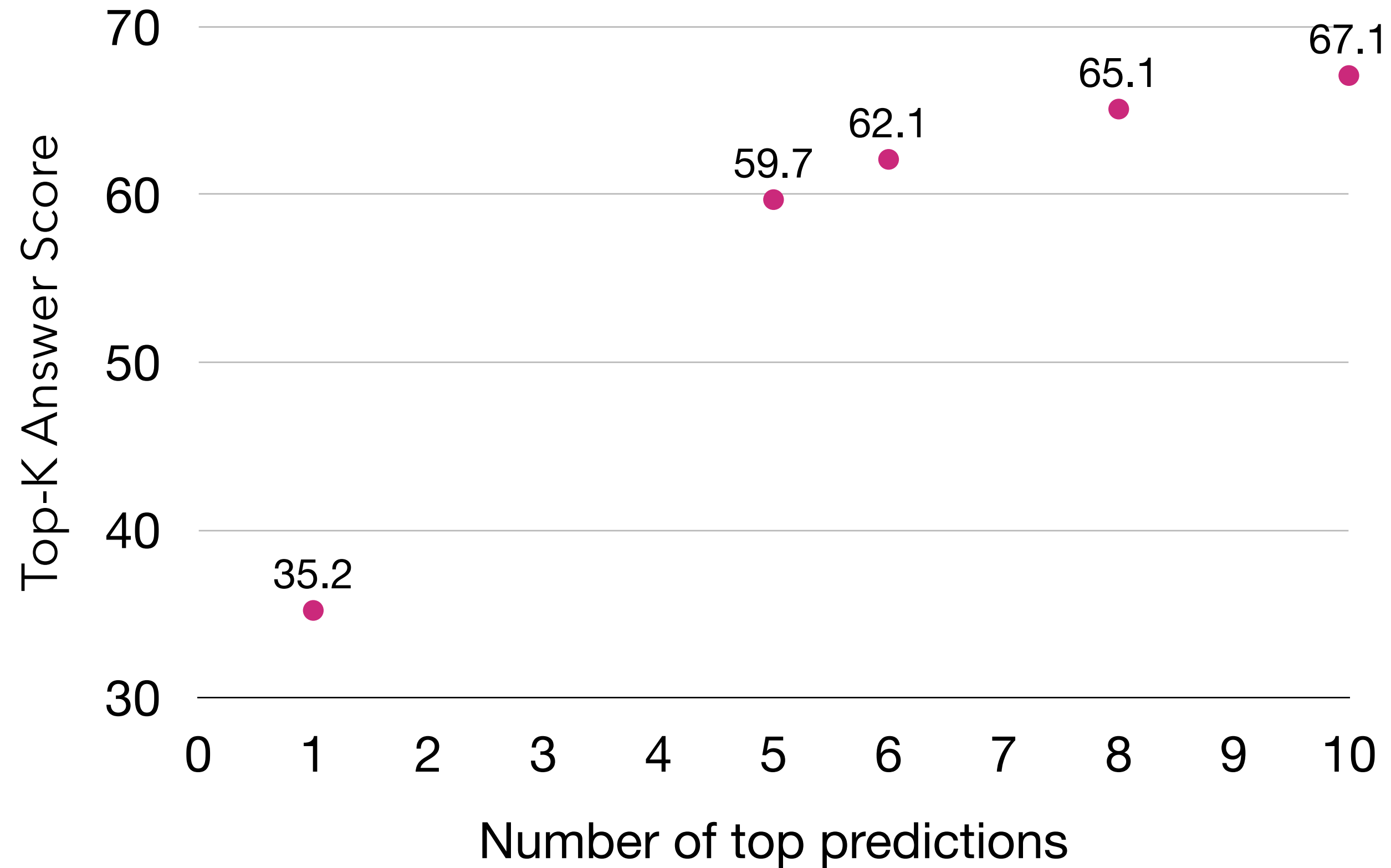Knowledge Module (Wikipedia) — MLP — MLP

47

# Multimodal Answer Validation

- MAVEx Answer-Guided Knowledge Retrieval

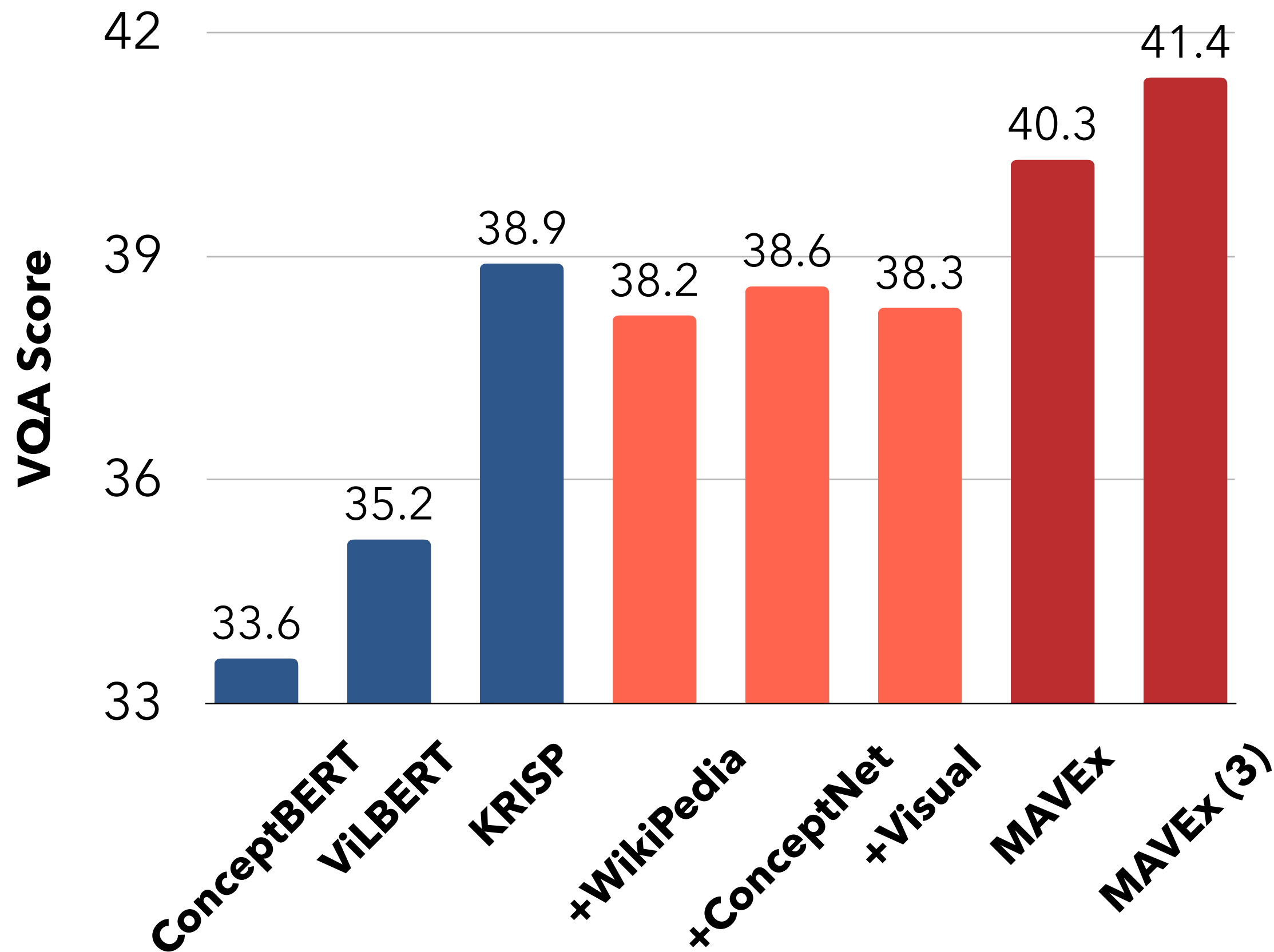- MAVEx Answer Validation Model

- **Experimental Results**

# Identifying Answer candidates

VQA systems are good candidate generators

# Performance on OKVQA v1.1

**OKVQA v1.1**[Kenneth et al. CVPR2019]: a large scale outside knowledge VQA dataset, containing 9k questions for training and 5k questions for test

**ConceptBert**[Gardères et al. EMNLP 2020]:Bert on ConceptNet

**ViLBERT**[Lu et al. CVPR 2020]: base multimodal transformer

**KRISP**[Kenneth et al. CVPR 2021]: Knowledge Reasoning with Implicit and Symbolic rePresentations

**+WikiPedia**(ours):Only use Wikipedia knowledge

**+ConceptNet**(ours):Only use ConceptNet knowledge

**+Visual**(ours):Only use Google Image knowledge

**MAVEx**(ours):Our full model

**MAVEx(3)**(ours): an ensemble of 3 models

Chart — VQA Score by model:
- ConceptBERT: 33.6
- ViLBERT: 35.2
- KRISP: 38.9
- +WikiPedia: 38.2
- +ConceptNet: 38.6
- +Visual: 38.3
- MAVEx: 40.3
- MAVEx (3): 41.4

OK-VQA benefits from answer validation with multiple knowledge sources

50

# Talk Outline

Core Contributions:

| Info | Image Captions | Human Explanations | Open Knowledge | **Open Knowledge** |
|------|----------------|---------------------|-----------------|---------------------|
| Topic | General VQA | Fighting Imbalanced Training Distribution | Outside Knowledge VQA | **Knowledge Retrieval** |

Future Directions

# Rethinking on Knowledge Retrieval

Hard term-matching based retriever (TF-IDF, BM25, MAVEx Retriever)



**Queries**
Movie
Man
Sitting man
Life story
Stranger
Forrest Gump
Gump
Speed

WIKIPEDIA
The Free Encyclopedia

A **man** is an adult **male human**.

The novel also features **Gump** as an astronaut, a professional wrestler, and a chess player.

**Forrest Gump** narrated his **life's story** at the …, as he **sat** at a bus stop bench

**Forrest Gump** is a 1994 American comedy-drama **film** …

**Speed** is a 1994 American action thriller **film** directed by Jan de Bont …

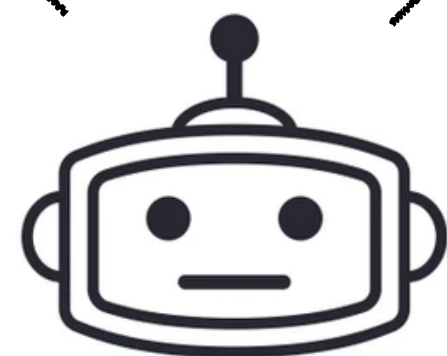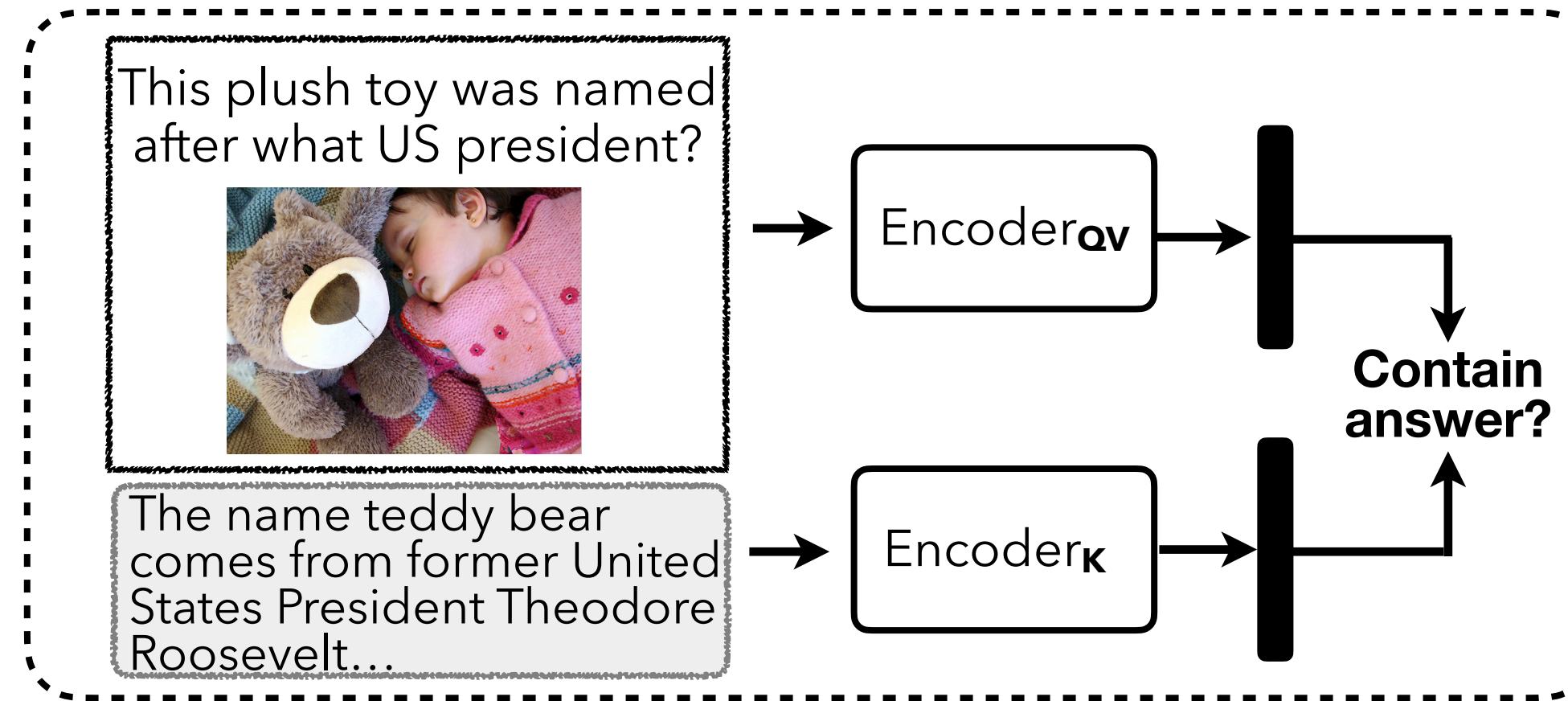Pros: Clear and explicit context matching
Cons: Bad at tackling synonyms and understanding semantic

# Dense Passage Retrieval

Q: This plush toy was named after what US president?



A: Teddy Roosevelt



This plush toy was named after what US president?

Encoder$_{QV}$

The name teddy bear comes from former United States President Theodore Roosevelt...

Encoder$_K$

Contain answer?

**Dense Passage Retriever**

Pros: Good at tackling synonyms and understanding semantic
Cons: Black-box context matching

[Qu et al., SIGIR2021; Luo et al., EMNLP2021]

# Objectives of Knowledge Retrieval

Q: This plush toy was named after what **US president**?



A: Teddy Roosevelt

**Informative**

The name teddy bear comes from former United States President **Theodore Roosevelt**, who was often referred to as "Teddy". …

# Objectives of Knowledge Retrieval

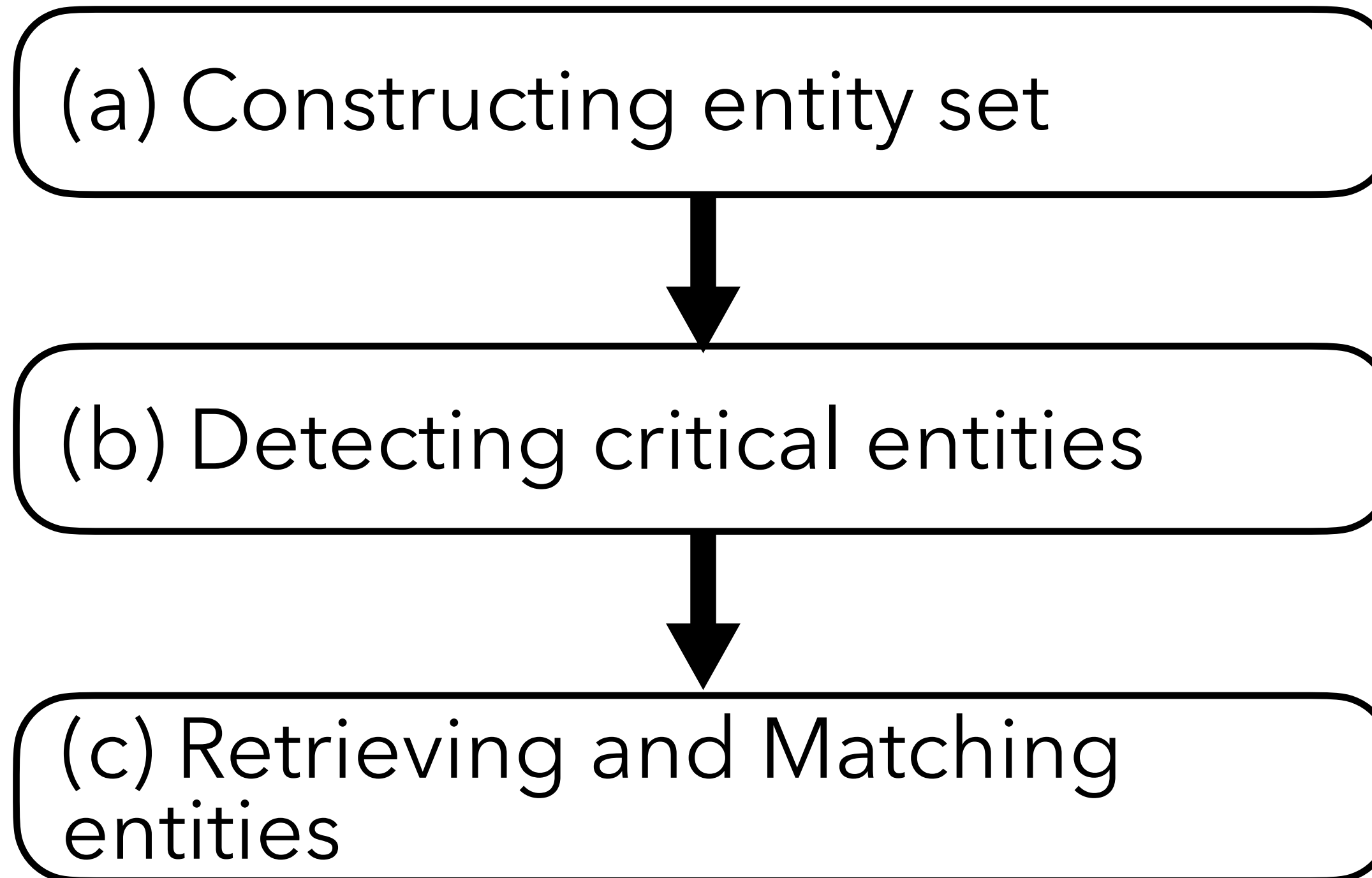Q: This **plush toy** was named after what **US president**?



A: Teddy Roosevelt

**Informative**
**Fitting in the query context**
**(Focusing on Critical Entity)**

The name **teddy bear** comes from former United States President **Theodore Roosevelt**, who was often referred to as "Teddy". …

# Entity-Focused Retrieval (EnFoRe)

(a) Constructing entity set

(b) Detecting critical entities

(c) Retrieving and Matching entities

# Entities Set Construction

Q: This plush toy was named after what US president?



A: Teddy Roosevelt

**Question-based Entities:**
Phrases in question: plush toy, president, …
Answer candidates: Roosevelt, …

# Entities Set Construction

Q: This plush toy was named after what US president?



A: Teddy Roosevelt

**Question-based Entities:**
Phrases in question: plush toy, president, …
Answer candidates: Roosevelt, …
Sub-questions: Teddy bear, …

What is this plush toy?
What is the US president?
…

# Entities Set Construction

Q: This plush toy was named after what US president?



A: Teddy Roosevelt

**Question-based Entities:**
Phrases in question: plush toy, president, …
Answer candidates: Roosevelt, …
Sub-questions: Teddy bear, …
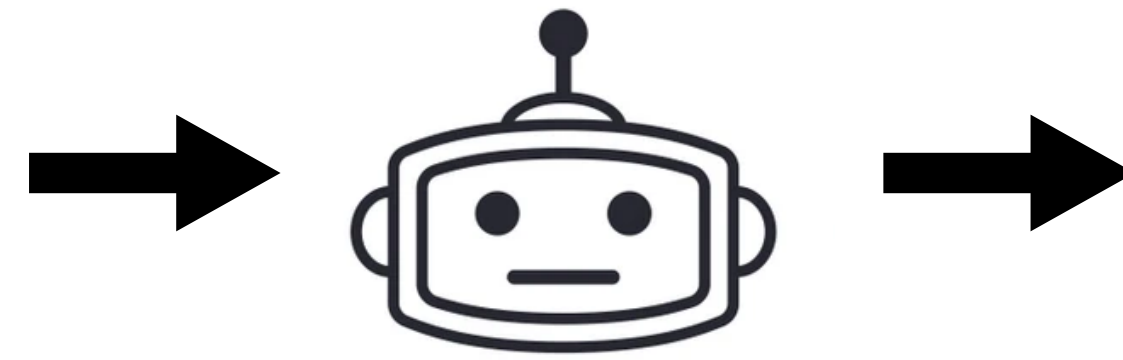
**Image-based Entities:**
Azure Tagging: Teddy bear, …

Tagging the common content, detecting brands and OCR.

# Entities Set Construction

Q: This plush toy was named after what US president?



A: Teddy Roosevelt

**Question-based Entities:**
Phrases in question: plush toy, president, …
Answer candidates: Roosevelt, …
Sub-questions: Teddy bear, …

**Image-based Entities:**
Azure Tagging: Teddy bear, …
Clip Tagging: Teddy bear, …

Clip model + WikiData

# Entities Set Construction

Q: This plush toy was named after what US president?



A: Teddy Roosevelt

**Question-based Entities:**
Phrases in question: plush toy, president, …
Answer candidates: Roosevelt, …
Sub-questions: Teddy bear, …

**Image-based Entities:**
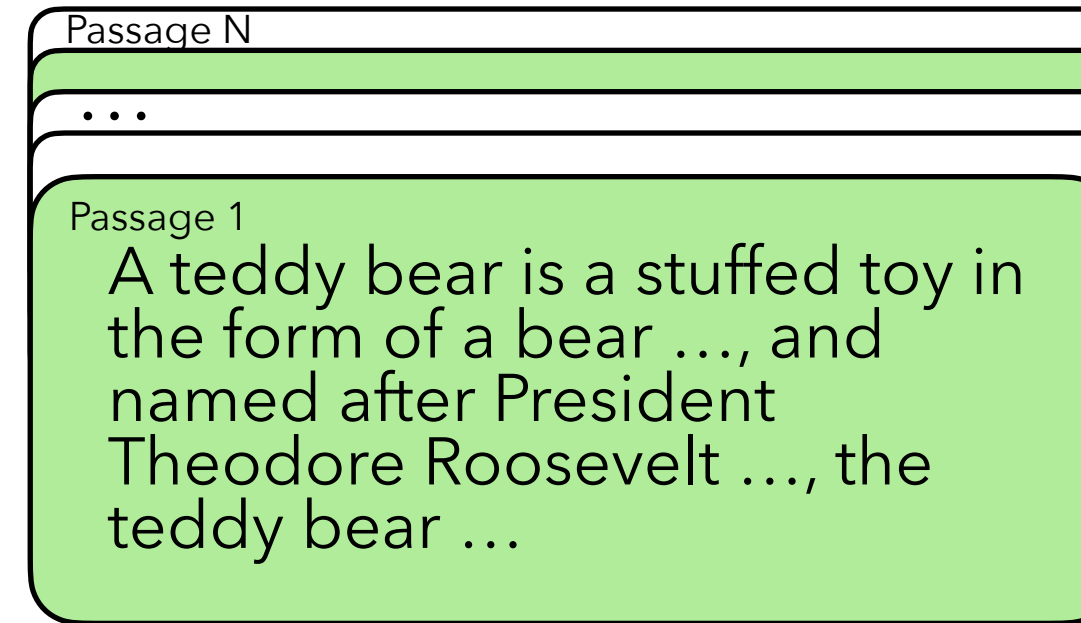Azure Tagging: Teddy bear, …
Clip Tagging: Teddy bear, …
Phrases in captions: Infant, …

There is an **infant** holding a toy.

61

# Oracle entities detection
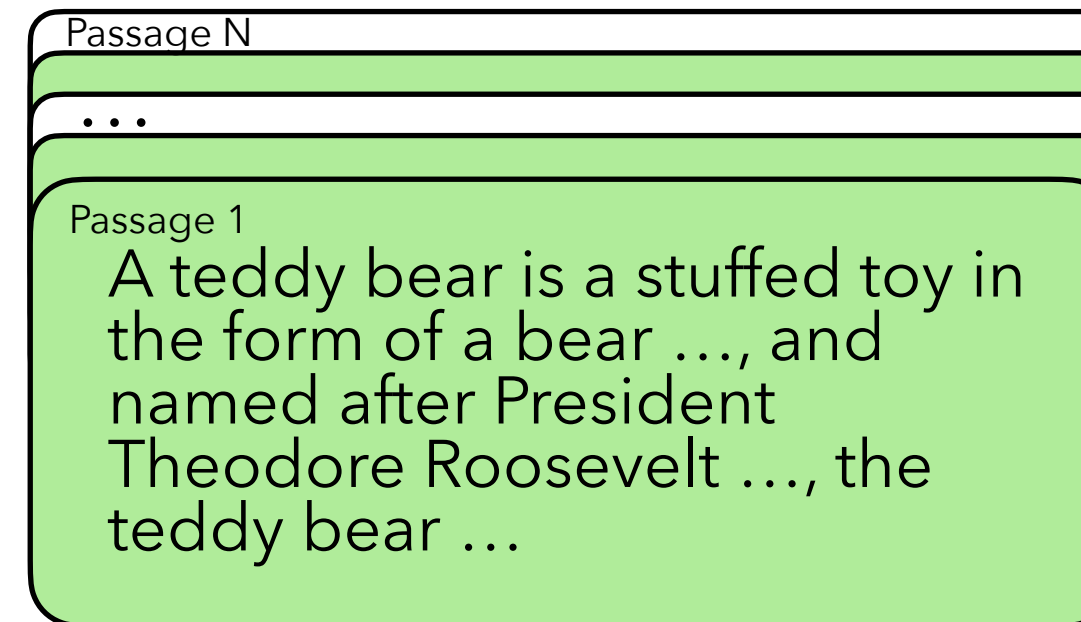
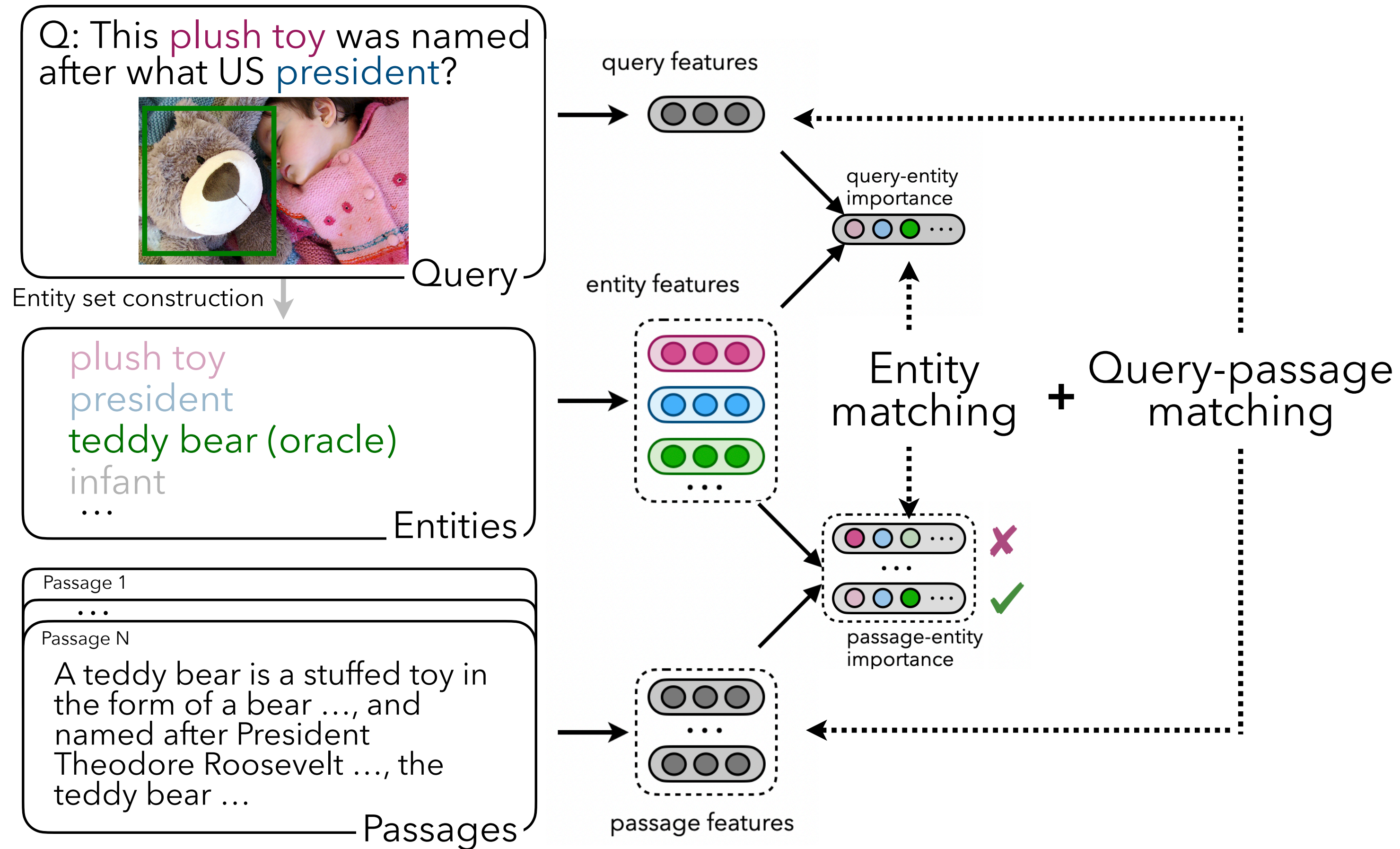This plush toy was named after what US president?

**Sparse Passage Retriever**

This plush toy was named after what US president?
Teddy Bear

Passage N

...

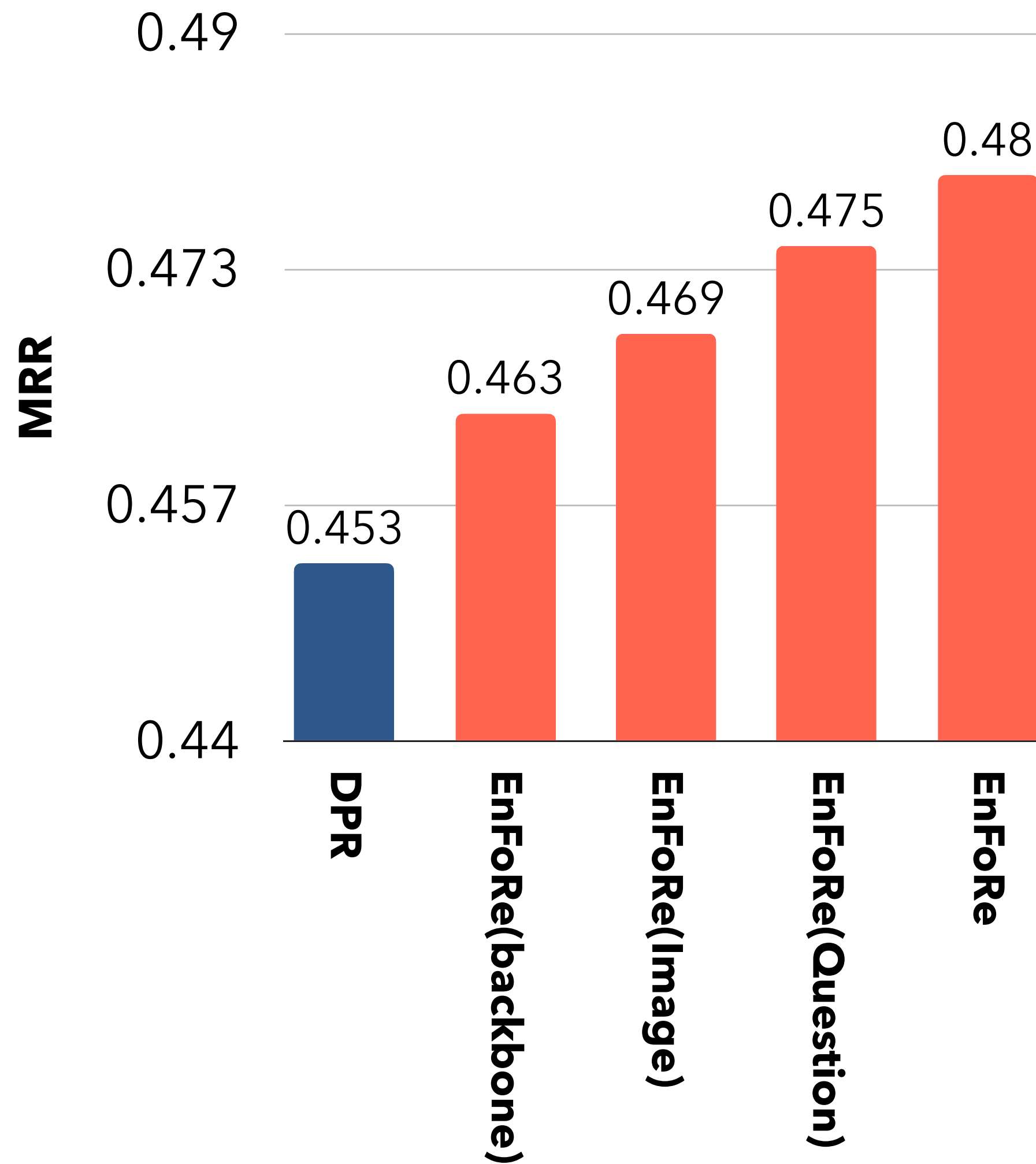Passage 1
A teddy bear is a stuffed toy in the form of a bear ..., and named after President Theodore Roosevelt ..., the teddy bear ...

$1/1 + 1/4$
$=1.25$

$\wedge$

Passage N

...

Passage 1
A teddy bear is a stuffed toy in the form of a bear ..., and named after President Theodore Roosevelt ..., the teddy bear ...

$1/1 + 1/2 + 1/4$
$=1.75$

62

# EnFoRe Model Overview

# EnFoRe Retrieval Results

**MRR**

0.49
0.473
0.457
0.44

0.453 (DPR)
0.463 (EnFoRe(backbone))
0.469 (EnFoRe(Image))
0.475 (EnFoRe(Question))
0.48 (EnFoRe)

DPR · EnFoRe(backbone) · EnFoRe(Image) · EnFoRe(Question) · EnFoRe

**DPR**[Qu et al. SIGIR 2021]:OKVQA version of DPR

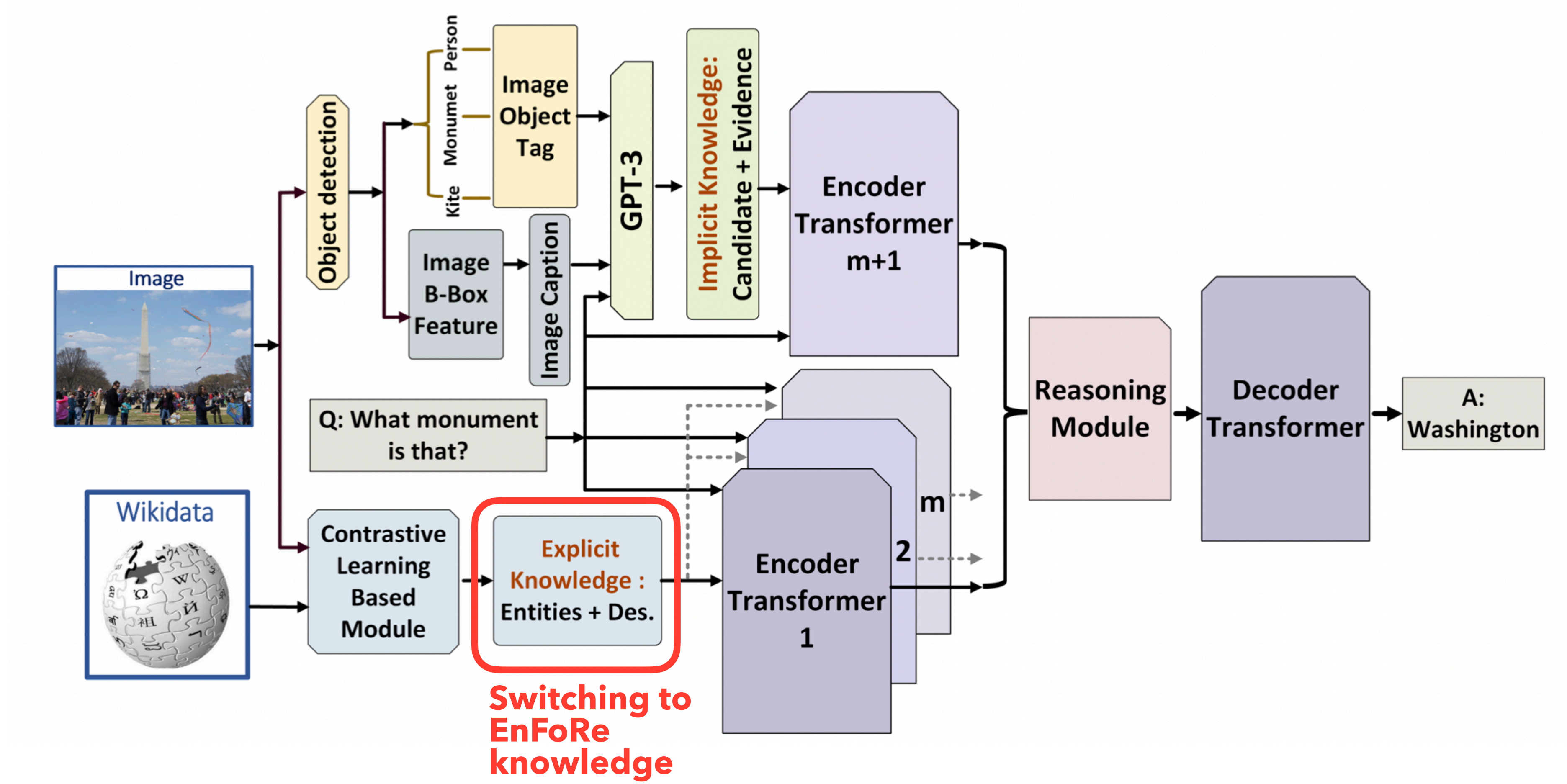**EnFoRe (backbone):** EnFoRe without entity matching, evaluating representation learning

**EnFoRe (Image):** EnFoRe re-ranking with Image-based entities

**EnFoRe (Question):** EnFoRe re-ranking with question-based entities

**EnFoRe:** Full EnFoRe model

64

# KAT Model

# VQA Results using EnFoRe Knowledge

**VQA Score**

| Model | Score |
|---|---|
| MAVEx | 0.382 |
| MAVEx + EnFoRe | 0.397 |
| KAT-base | 0.506 |
| KAT-base + EnFoRe | 0.522 |
| KAT-full | 0.544 |
| KAT-full + EnFoRe | 0.552 |

Y-axis: 0.37, 0.433, 0.497, 0.56

**OKVQA v1.1**[Kenneth et al. CVPR2019]: a large scale outside knowledge VQA dataset, containing 9k questions for training and 5k questions for test

**MAVEx**[Wu et al. AAAI 2022]:MAVEx with Wikipedia knowledge

**MAVEx + EnFoRe:** MAVEx with EnFoRe knowledge

**KAT-base:** Single T5-base Fusion-in-Decoder model

**KAT-base + EnFoFe:** KAT-base with EnFoRe knowledge

**KAT-full:** Ensembles of T5-large Fusion-in-Decoder model

**KAT-full + EnFoFe:** KAT-full with EnFoRe knowledge

The Results of KAT related model are evaluated by an unofficial metric that that takes the max over 1.0 and number of annotators agreements divided by 3.

# Human Evaluation

Which set of highlighted entities and background sentences contains the most supportive evidence to the answer?

Q: What is the main benefits of healthy properties these fruit contain? A: potassium



**EnFoRe wins 61.8%**

*Sent1: Naturally, the fruit of the Hassaku orange is a good source of vitamin C, folic acid, potassium and fiber.*
*Sent2: …*
*Sent3: …*

*Sent1: Fruit vegetable is a fruit commonly referred to as a vegetable because they are savory (not sweet).*
*Sent2: …*
*Sent3: …*

# Qualitative Results

Q: Who is famous for allegedly doing this in a lightning storm?
Prediction: Benjamin Franklin

*The electricity attraction from a lightning storm was done by **Benjamin Franklin** himself in the **kite** experiment that he talked about in a letter to Collinson dated October 19.*

Top Entities: **kite** flying, human-lifting **kite**, **kite**.

# Talk Outline

Core Contributions:

| Info | Image Captions | Human Explanations | Open Knowledge | Open Knowledge |
|------|----------------|--------------------|----------------|----------------|
| Topic | General VQA | Fighting Imbalanced Training Distribution | Outside Knowledge VQA | Knowledge Retrieval |

**Future Directions**

# Future Directions
## Logic-Aware Knowledge Retrieval

Q: What healthy properties do these **fruit** contain?
Prediction: **potassium**



Informative
**Fitting in the query context (Entity-based)**

*Naturally, the fruit of the Hassaku **orange** is a good source of vitamin C, folic acid, **potassium** and fiber.*

71

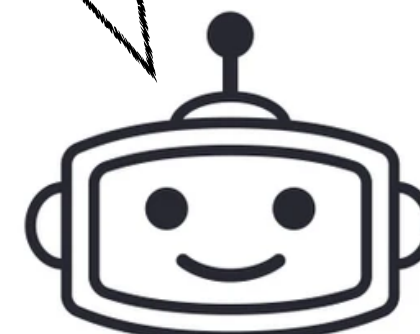# Step-by-step Knowledge Retrieval

Q: What is the main benefits of healthy properties these fruit contain?
Prediction: **potassium**



Breaking Down Questions

1. What fruits are the question refers to?
   **[grounding]**
2. What healthy properties do oranges have
   **[retrieval]**
3. What healthy properties do banana have
   **[retrieval]**
4. What is the shared healthy properties the
   **[reasoning]**
5. What is the main benefits of the healthy property
   **[retrieval]**

# Knowledge Sources

# Knowledge-Augmented Tasks

## Multimodal Pre-training

Coarsely labeled web images:
ConceptualCaptions (~3.3M)
YFCC100M(~100M Images)
JFT300M -3B (300M - 3B)

## Knowledge-Augmented Multimodal Pre-training



In the 15th and 16th centuries, Portuguese colonists started banana plantations in the Atlantic Islands, Brazil, and western Africa

## Open-vocabulary Object Detection



Common detection: Banana.

Open-vocabulary detection: raped banana, plantain, cavendish banana.

# Acknowledgements

## Thanks for you attention!