

Experiments in Retrieval-Augmented Image Captioning

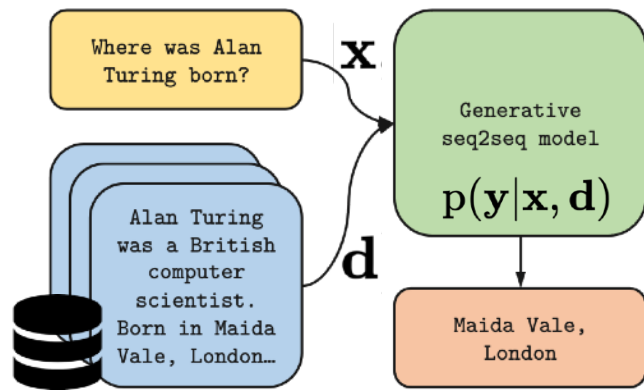


~~Rita Ramos~~ Desmond Elliott

Department of Computer Science
University of Copenhagen

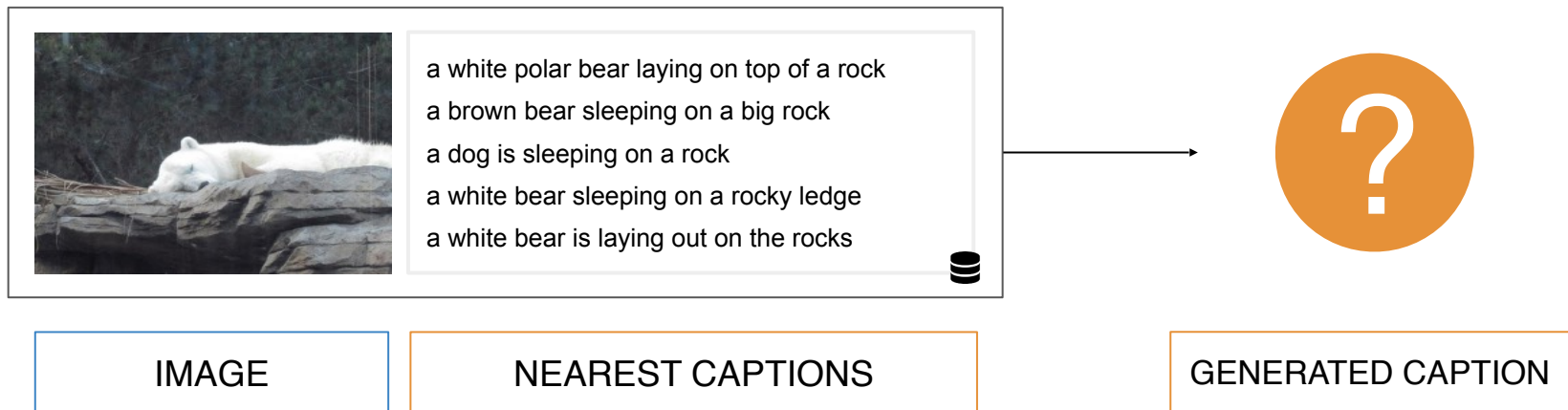
Retrieval Augmented Generation

- Combine the power of in-weights learning with in-context adaptation through retrieval augmentation
- Given a datastore of facts, knowledge, documents, etc.
 - Combine the most relevant items from the datastore (d) with the input (x) for your task



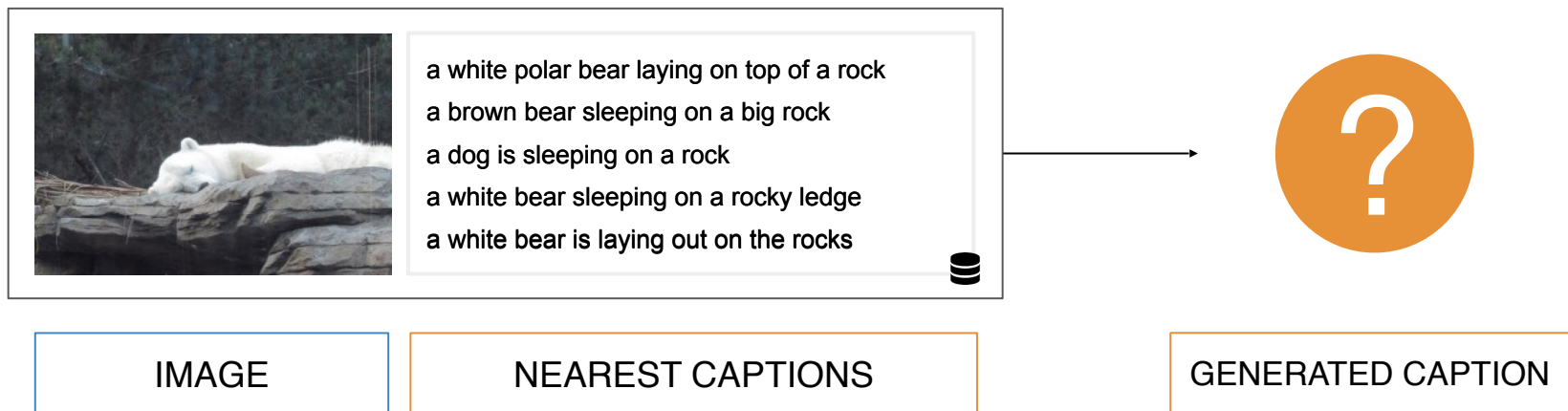
Multimodal Retrieval Augmentation

- Combine the most relevant items from the datastore with the input



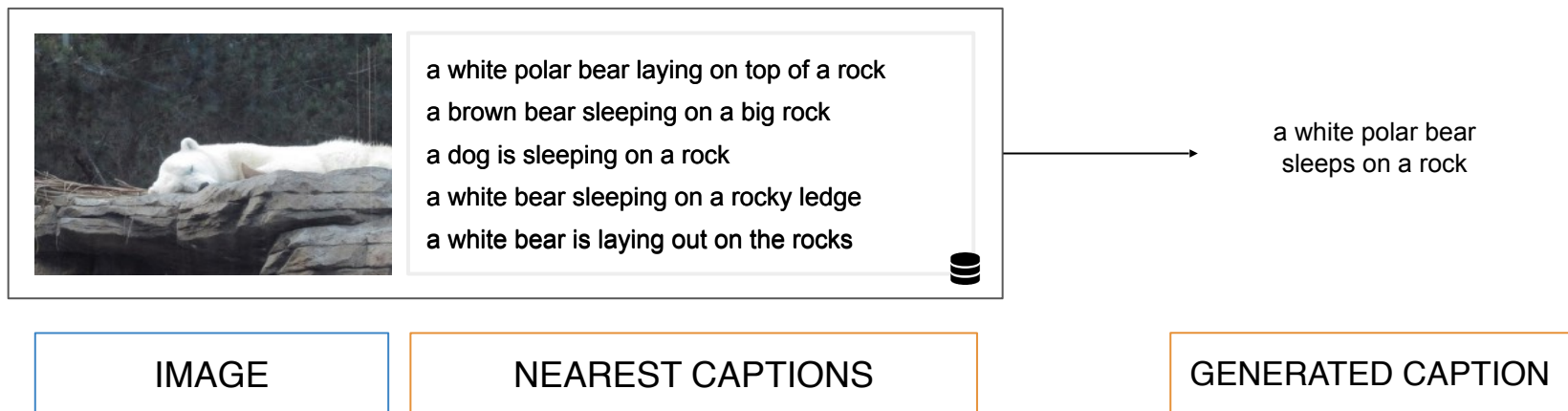
Multimodal Retrieval Augmentation

- Combine the most relevant items from the datastore with the input



Multimodal Retrieval Augmentation

- Combine the most relevant items from the datastore with the input



... at ACL 2025

... at ACL 2025

Towards Text-Image Interleaved Retrieval

Xin Zhang^{1,2*}, Ziqi Dai^{1*}, Yongqi Li², Yanzhao Zhang, Dingkun Long

Pengjun Xie, Meishan Zhang^{1†}, Jun Yu¹, Wenjie Li², Min Zhang¹

¹Harbin Institute of Technology, Shenzhen ²The Hong Kong Polytechnic University

{zhangxin2023, ziqi.dai}@stu.hit.edu.cn zhangmeishan@hit.edu.cn

Release at <https://github.com/vec-ai/wikiHow-TIIR>

... at ACL 2025

Towards Text-Image Interleaved Retrieval

Xin Zhang

Pengjie

¹Harbin Institute of Technology

{zhangxin, pengjie}

hust.edu.cn

WavRAG: Audio-Integrated Retrieval Augmented Generation for Spoken Dialogue Models

Yifu Chen^{1,†} **Shengpeng Ji**^{1,†} **Haoxiao Wang**^{1,†} **Ziqing Wang**³ **Siyu Chen**¹

Jinzheng He² **Jin Xu**² **Zhou Zhao**^{1,*}

¹ Zhejiang University ² Alibaba Group ³ Beijing University of Technology

[†] Equal contribution. ^{*} Corresponding author.

Eve106298@163.com

zhaozhou@zju.edu.cn

... at ACL 2025

Towards Text-Image

Maximal Matching Matters: Preventing Representation Collapse for Robust Cross-Modal Retrieval

**Xin Zhang
Pengjie**

¹Harbin Institute of Technology
{zhangxin, pengjie}

I

WavRAG: Audio

Hani Alomari
Virginia Tech
hani@vt.edu

Anushka Sivakumar
Virginia Tech
anushkas01@vt.edu

Andrew Zhang
Virginia Tech
azhang42@vt.edu

Chris Thomas
Virginia Tech
chris@cs.vt.edu

Yifu Chen^{1,†} **Shengpeng Ji** **Yuanhao Huang** **Ziyang Huang** **Yi Yu Chen**

Jinzheng He² **Jin Xu**² **Zhou Zhao**^{1,*}

¹ Zhejiang University ² Alibaba Group ³ Beijing University of Technology

[†] Equal contribution. ^{*} Corresponding author.

Eve106298@163.com

zhaozhou@zju.edu.cn

... at ACL 2025

Towards Text-Image

Maximal Matching Matters: Preventing Representation Collapse for Robust Cross-Modal Retrieval

Xin Zhang
Pengjie

¹Harbin Institute of Technology
{zhangxin, pengjie}

WavRAG: Audio Retrieval-Augmented Generation

Hani Alomari
Virginia Tech
hani@vt.edu

Anushka Sivakumar
Virginia Tech
anushkas01@vt.edu

Andrew Zhang
Virginia Tech
azhang42@vt.edu

Chris Thomas
Virginia Tech
chris@cs.vt.edu

Yifu Chen^{1,†} **Shengpeng Gao**² **Yiwei Wang**² **Ziyang Wang**² **Yi Yu Chen**²
Jinzheng He² **Jin Xu**² **Zhou Zhao**^{1,*}
¹ Zhejiang University ² Alibaba Group ³ Beijing University of Technology
[†] Equal contribution

Eve
zhao

VISA: Retrieval Augmented Generation with Visual Source Attribution

Xueguang Ma^{*,1} **Shengyao Zhuang**^{*,2,3} **Bevan Koopman**^{2,3}
Guido Zuccon³ **Wenhu Chen**¹ **Jimmy Lin**¹

¹University of Waterloo ²CSIRO ³University of Queensland

x93ma@uwaterloo.ca, s.zhuang@uq.edu.au

... at ACL 2025

Towards Text-Image

Xin Zhang
Pengjie

¹Harbin Institute of Technology
{zhangxin, pengjie}

WavRAG: Audio

Maximal Matching Matters: Preventing Representation Collapse for Robust Cross-Modal Retrieval

Hani Alomari
Virginia Tech
hani@vt.edu

Anushka Sivakumar
Virginia Tech
anushkas01@vt.edu

Andrew Zhang
Virginia Tech
azhang42@vt.edu

Chris Thomas
Virginia Tech
chris@cs.vt.edu

OMGM: Orchestrate Multiple Granularities and Modalities for Efficient Multimodal Retrieval

Wei Yang^{*}, Jingjing Fu[†], Rui Wang, Jinyu Wang, Lei Song, Jiang Bian
Microsoft Research Asia
wyang6621@gmail.com, {jifu, ruiwa, jinywan, lesong, jiabia}@microsoft.com

Generation with Visual Source Attribution

Yiyao Zhuang^{*,2,3}, Bevan Koopman^{2,3}
Wenhu Chen¹, Jimmy Lin¹

²CSIRO ³University of Queensland
o.ca, s.zhuang@uq.edu.au

... at ACL 2025

Towards Text-Image

Maximal Matching Matters: Preventing Representation Collapse for Robust Cross-Modal Retrieval

Xin Zhang
Pengjie

¹Harbin Institute of Technology
{zhangxi, pengjie}@hit.edu.cn

WavRAG: Audio-Visual Retrieval-Augmented Generation

Hani Alomari
Virginia Tech
hani@vt.edu

Anushka Sivakumar
Virginia Tech
anushkas01@vt.edu

Andrew Zhang
Virginia Tech
azhang42@vt.edu

Chris Thomas
Virginia Tech
chris@cs.vt.edu

OMCM: Orchestrate Multiple Granularities and Modalities for Efficient



Ask in Any Modality

A Comprehensive Survey on Multimodal Retrieval-Augmented Generation

Mohammad Mahdi Abootorabi^{†,✱}, Amirhosein Zobeiri[°], Mahdi Dehghani[‡], Mohammadali Mohammadkhani[§],
Bardia Mohammadi[✱], Omid Ghahroodi[†], Mahdieh Soleymani Baghshah^{§,✱}, Ehsaneddin Asgari^{✱,*}

[†]Qatar Computing Research Institute, [‡]Saarland University, [✱]Zuse School ELIZA, [°]University of Tehran,

[✱]Max Planck Institute for Software Systems, [‡]K.N. Toosi University of Technology, [§]Sharif University of Technology

Correspondence: soleymani@sharif.edu and easgari@hbku.edu.qa

<https://multimodalrag.github.io>

Generation with Visual Source Attribution

Yiyao Zhuang^{*,2,3} Bevan Koopman^{2,3}
Wenhu Chen¹ Jimmy Lin¹

²CSIRO ³University of Queensland

o.ca, s.zhuang@uq.edu.au

... at ACL 2025

Towards Text-Image

Xin Zhang
Pengjie

¹Harbin Institute of Technology
{zhangxin, pengjie}

WavRAG: Audio

Maximal Matching Matters: Preventing Representation Collapse for Robust Cross-Modal Retrieval

Hani Alomari
Virginia Tech
hani@vt.edu

Anushka Sivakumar
Virginia Tech
anushkas01@vt.edu

Andrew Zhang
Virginia Tech
azhang42@vt.edu

Chris Thomas
Virginia Tech
chris@cs.vt.edu

MegaPairs: Massive Data Synthesis for Universal Multimodal Retrieval

Junjie Zhou^{1,2*}, Yongping Xiong^{1*}, Zheng Liu^{2,5†}, Ze Liu^{3,2}, Shitao Xiao²,
Yueze Wang², Bo Zhao⁴, Chen Jason Zhang⁵, Defu Lian^{3†}

¹ State Key Laboratory of Networking and Switching Technology,
Beijing University of Posts and Telecommunications

² Beijing Academy of Artificial Intelligence

³ University of Science and Technology of China

⁴ Shanghai Jiao Tong University, ⁵ The Hong Kong Polytechnic University
{junjiebu, zhengliu1026}@gmail.com, liandefu@ustc.edu.cn

A Comprehensive Survey of

Mohammad Mahdi Abootorabi^{♣†}, An
Bardia Mohammadi[♣], Omid Ghahmipour

[†]Qatar Computing Research Institute

[♣]Max Planck Institute for Software Systems, ^{*}K.N. Toosi University of Technology, [~]Sharif University of Technology

Correspondence: soleymani@sharif.edu and easgari@hbku.edu.qa

<https://multimodalrag.github.io>

Source Attribution

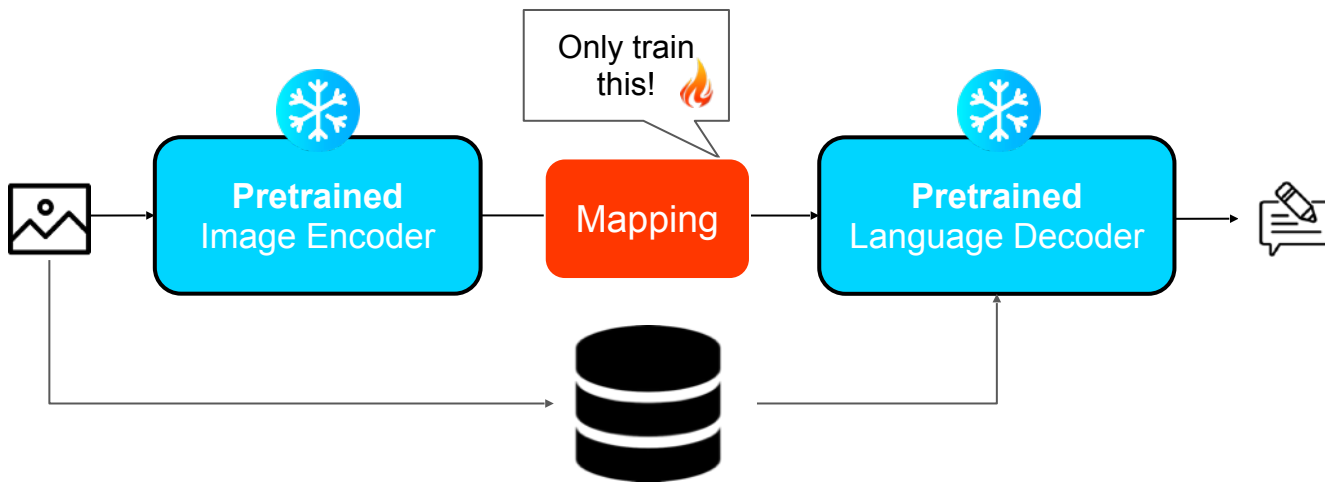
soleymani^{2,3}
¹

Queensland

o.ca, s.zhuang@uq.edu.au

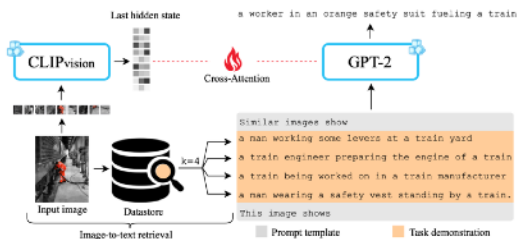
Our Motivation

- Train lightweight image captioning models using frozen backbones
 - CLIPCap (Mokady et al. 2021), I-Tuning (Luo et al. 2023)
- ... and using retrieval augmentation to *assist* the decoder

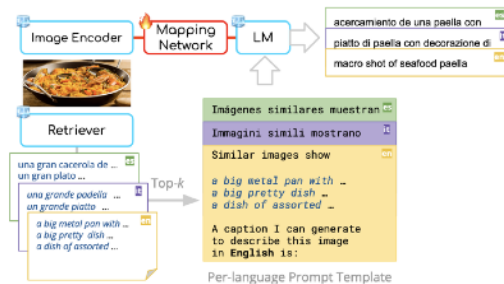


Overview

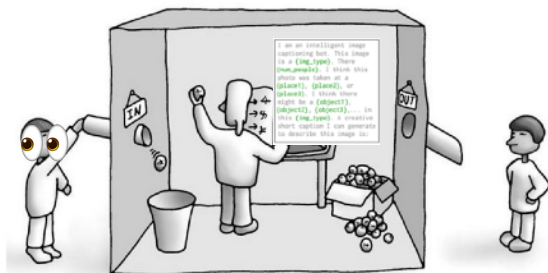
1. Lightweight RAG Captioning



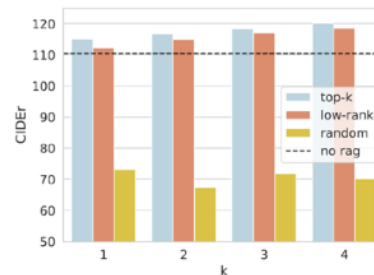
2. Lightweight Multilingual Training



3. Image-blind captioning



4. Understanding Multimodal RAG



SmallCap: Lightweight Image Captioning Prompted with Retrieval Augmentation

CVPR 2023



R. Ramos



B. Martins



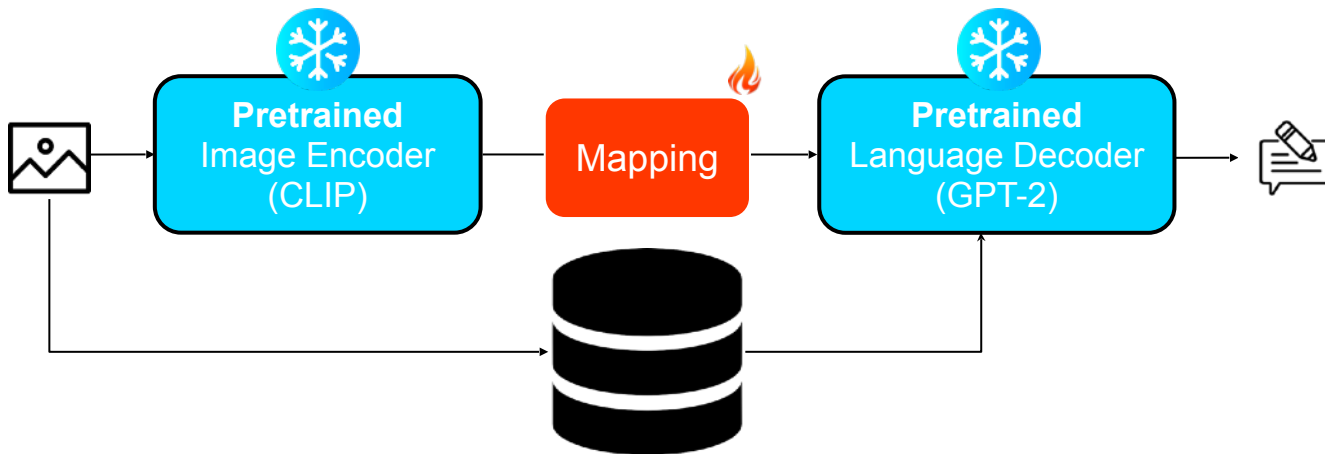
D. Elliott



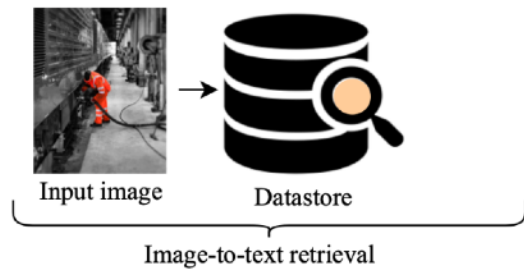
Y. Kementchedjheva

Lightweight Training through Retrieval

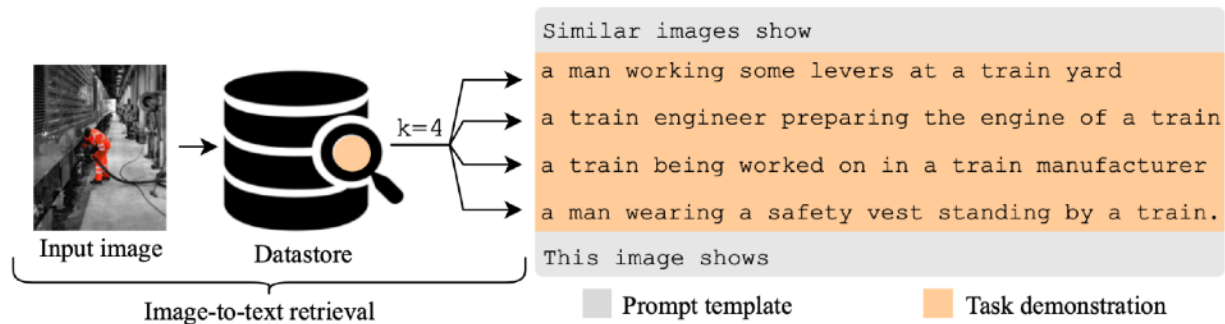
- Given the success of retrieval augmented generation, can we extend this to multimodality with a lightweight training paradigm?



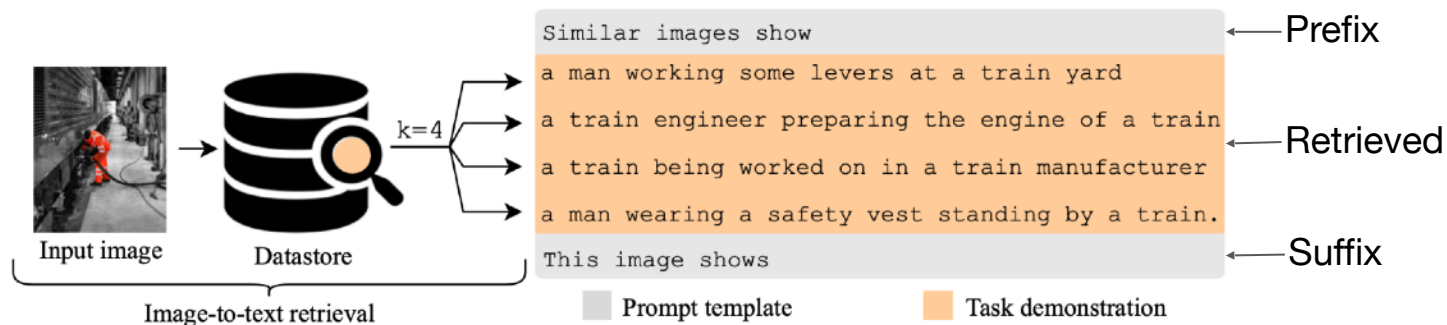
SmallCap Model



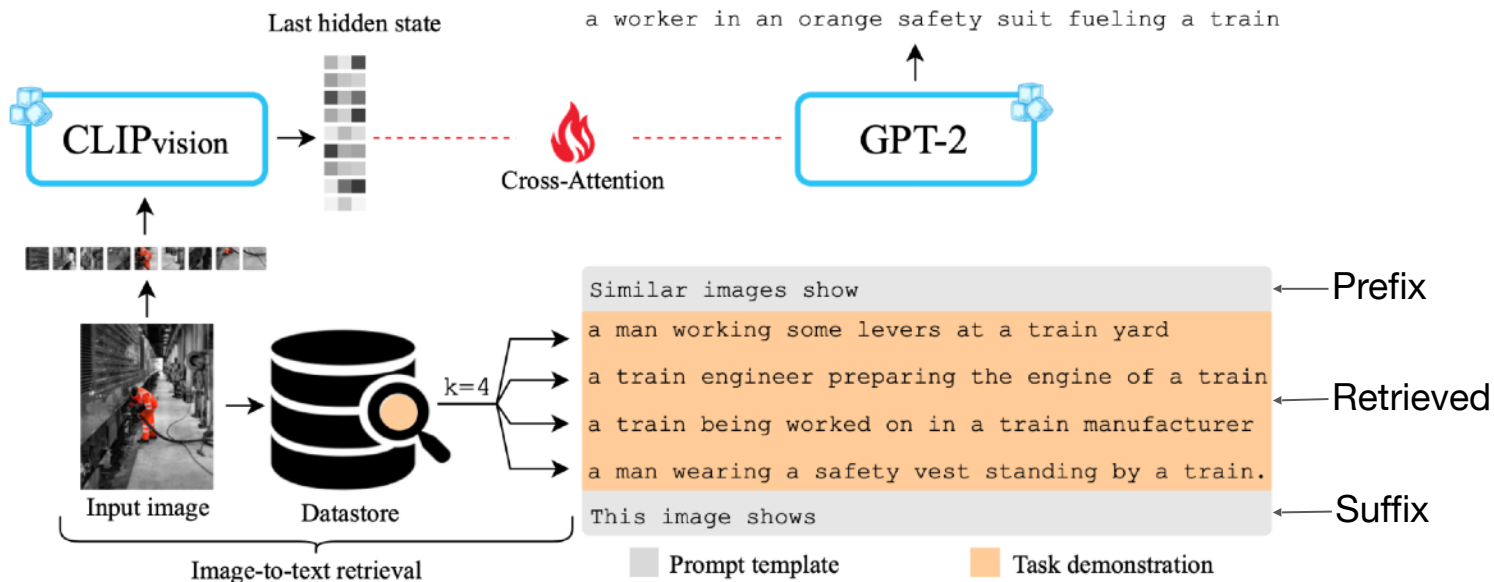
SmallCap Model



SmallCap Model

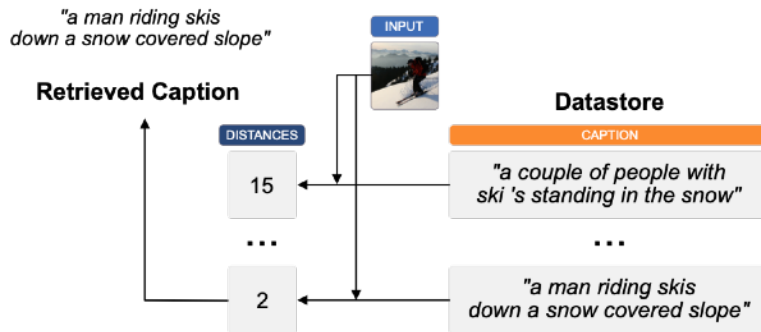


SmallCap Model



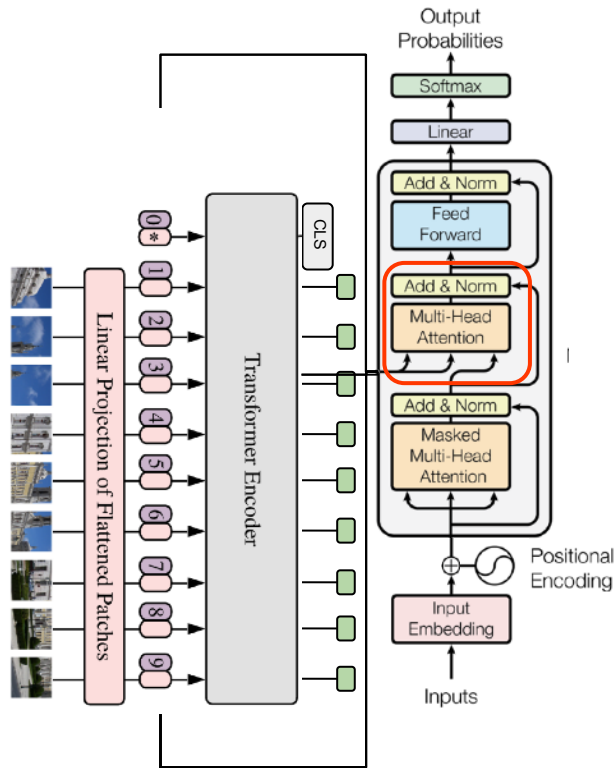
Retrieval System

- Build a datastore with high-dimensional dense vectors
 - FAISS: Facebook AI Similarity Search for nearest-neighbor search
 - Captions of images represented with CLIP embeddings
- Retrieve k nearest-neighbours captions from datastore
 - Image embedding compared against datastore caption vectors



Trained Cross-Attention Layers

- Autoregressive Transformer LMs only contain a multi-head **self-attention mechanism**
- We insert a randomly initialized **cross-attention mechanism** to attend to the visual encoder output embeddings



Experimental Setup

- Pretrained CLIP-ViT-B/32 and GPT/OPT backbone models
- Randomly initialize the cross-attention layer
- Train only on COCO in only 8 hours on 1 x 40GB NVIDIA A100 GPU

Low-rank cross-attention

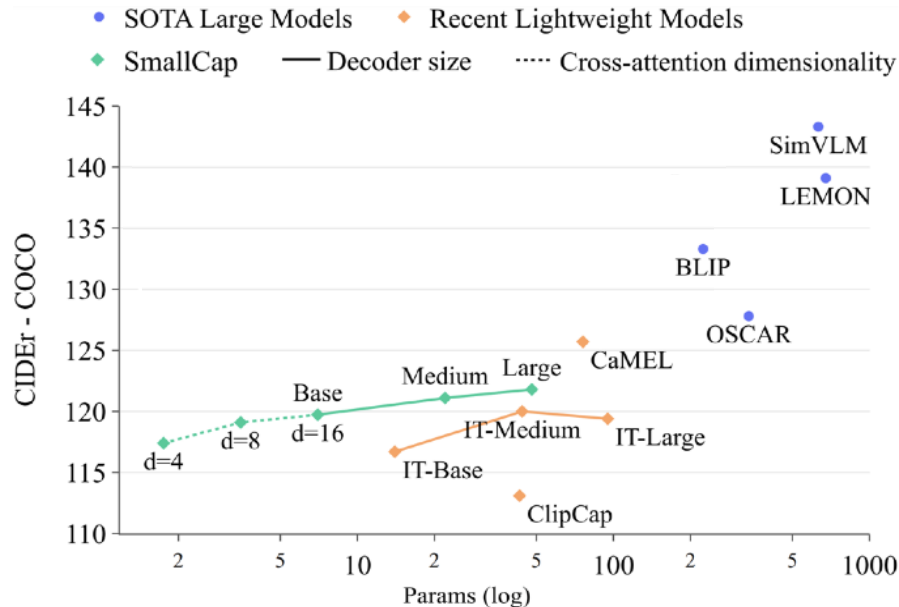
$$\text{Att}(\mathbf{QW}_i^Q, \mathbf{KW}_i^K, \mathbf{VW}_i^V)$$

$$\mathbf{W}_i^K, \mathbf{W}_i^Q, \mathbf{W}_i^V$$

$$\in \mathbb{R}^{d_{\text{encoder}} \times d}$$

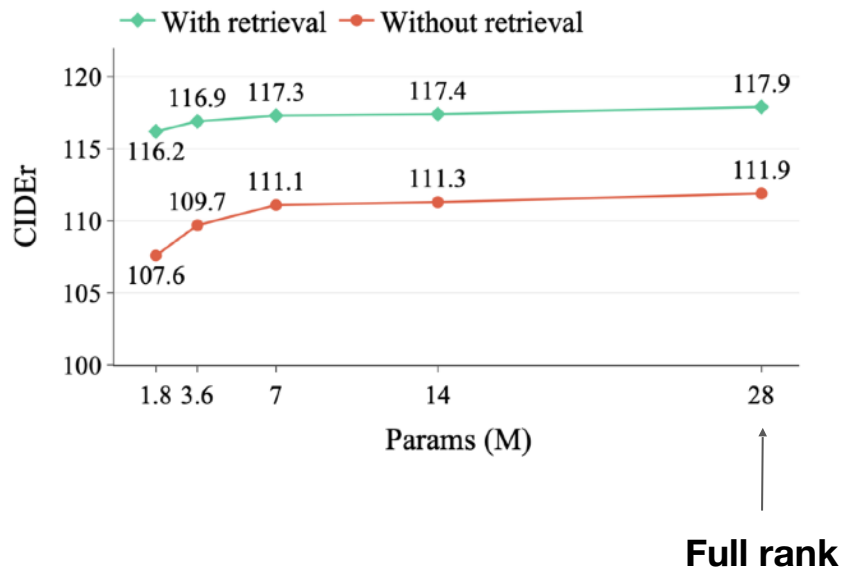
Attention rank	Params
d=64 (Full)	22M
d=16	7M
d=8	3.6M
d=4	1.8M

Results



- Outperform other lightweight approaches
- Effective with low-rank matrices: $4, 8, 16 \ll 64$
- Larger pretrained decoders further improve performance

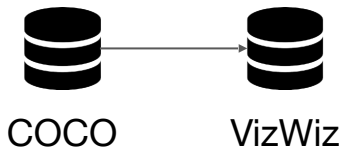
Importance of Retrieval Augmentation



- With retrieval:
 - Performance is stable across the range of cross-attention sizes
- Without retrieval:
 - Drop in performance
 - SmallCap model performance degrades at a higher rate

Training-Free Domain Transfer

- SmallCap was trained on COCO but we can easily swap the datastore



	Flickr30k	VizWiz	MSR-VTT
ClipCap	41.2	28.3	12.5
CaMEL	55.2	37.6	20.7
SmallCap	60.6	55.0	28.4

Qualitative Example from VizWiz



- some carrots potatoes garlic an onion and some chicken broth
- a selection of ingredients for soup includes carrots, meat, and prepackaged broth
- this is the makings of a meal with chicken and vegetables
- the meal has chicken, bread, and cole slaw

Generated caption:

a close up of a plate of food on a table

Qualitative Example from VizWiz



- some carrots potatoes garlic an onion and some chicken broth
- a selection of ingredients for soup includes carrots, meat, and prepackaged broth
- this is the makings of a meal with chicken and vegetables
- the meal has chicken, bread, and cole slaw

Generated caption:

a close up of a plate of food on a table



- a can of swanson fat free chicken broth
- a can of swanson brand chicken broth with less sodium
- a 14,5 ounce can of swanson branded chicken broth
- a can of swanson chicken broth on a table

Generated caption:

a can of swanson brand chicken broth on a table

Qualitative Example from VizWiz



- some carrots potatoes garlic an onion and some chicken broth
- a selection of ingredients for soup includes carrots, meat, and prepackaged broth
- this is the makings of a meal with chicken and vegetables
- the meal has chicken, bread, and cole slaw

Generated caption:

a close up of a plate of food on a table



- a can of swanson fat free chicken broth
- a can of swanson brand chicken broth with less sodium
- a 14,5 ounce can of swanson branded chicken broth
- a can of swanson chicken broth on a table

Generated caption:

a can of swanson brand chicken broth on a table

“swanson” does not appear anywhere in the COCO training dataset

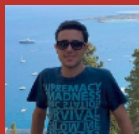
Q: What about multilingual captioning?

PAELLA: Parameter-Efficient Lightweight Language-agnostic Captioning Model

Findings of NAACL 2024



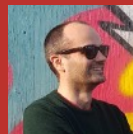
R. Ramos



E. Bugliarello



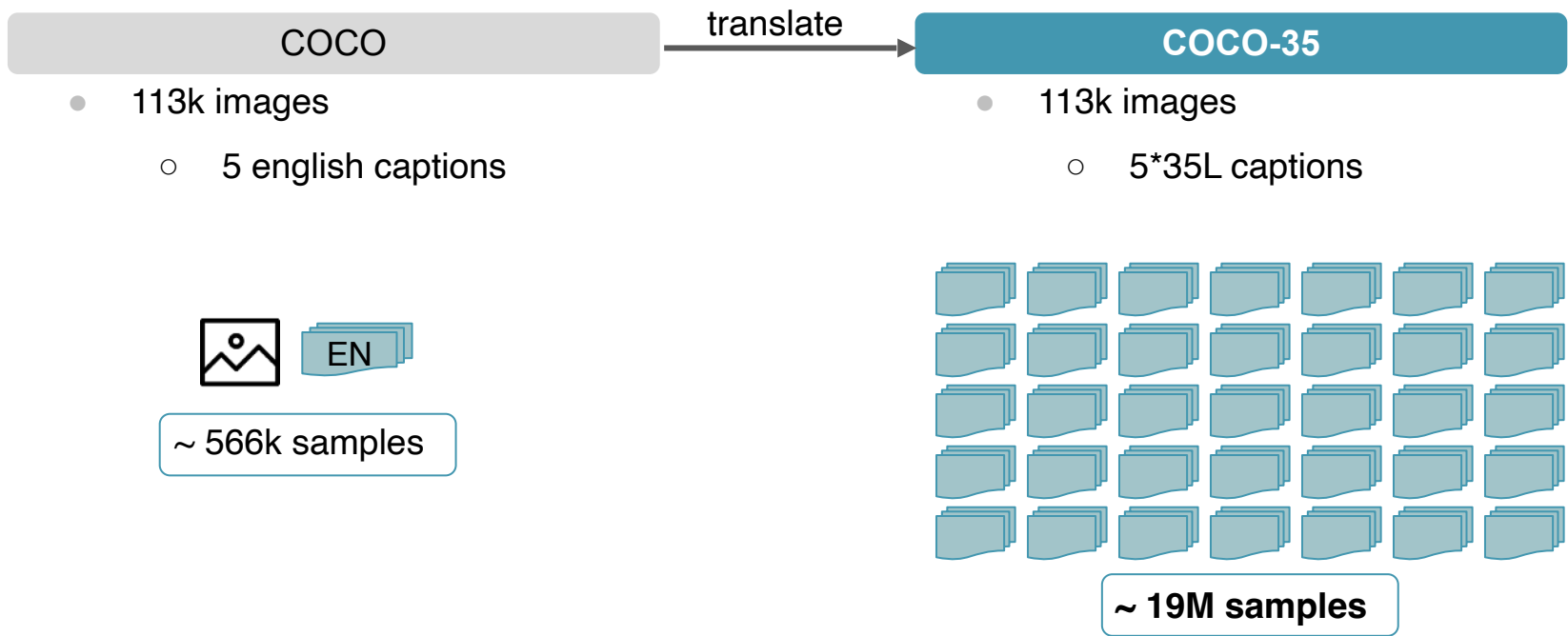
B. Martins



D. Elliott

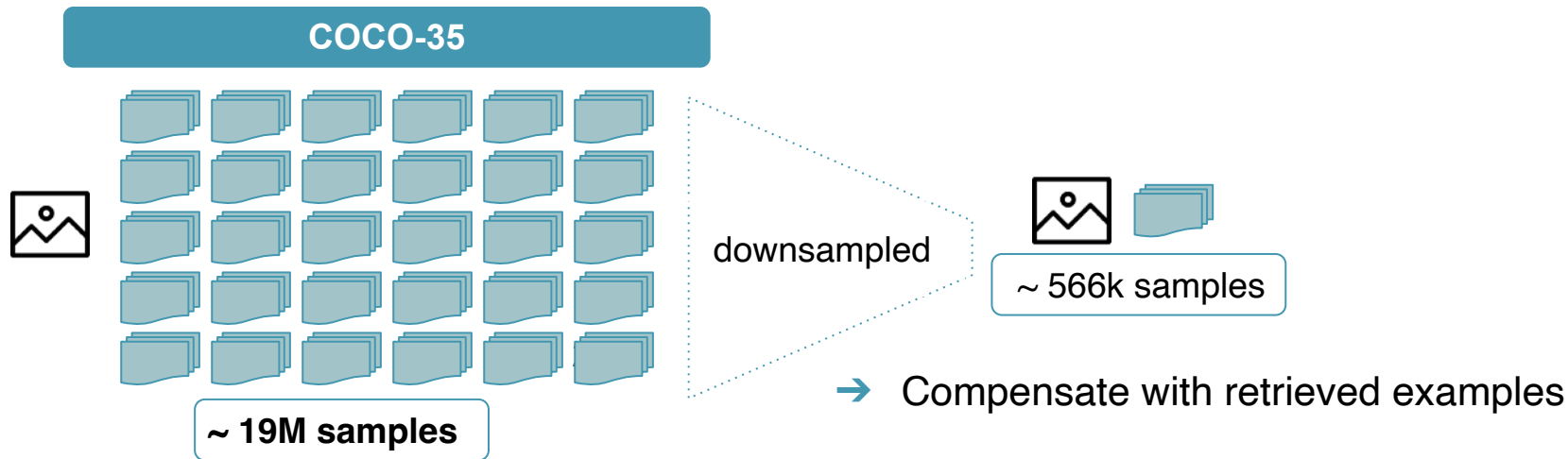
Multilingual Image Caption Training

- Common approach in the literature is to machine translate and train

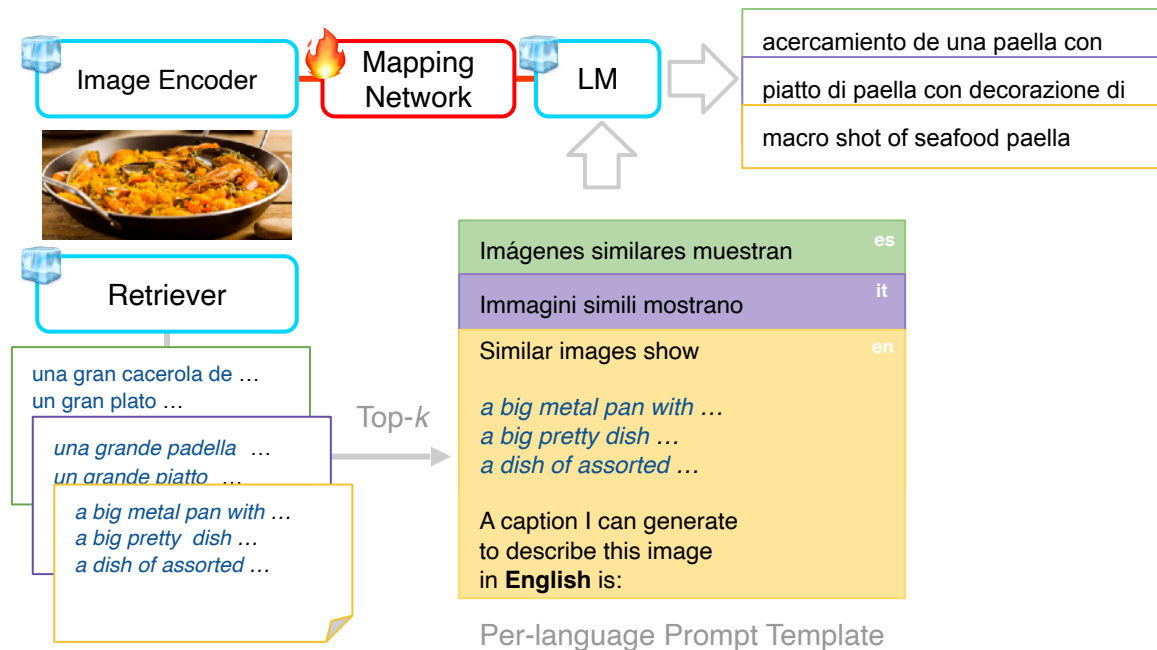


Data-Efficient Multilingual Training

- Only train on a subset of COCO-35:
 - Sample uniformly across 35 languages
 - **Match the size** of the English COCO dataset

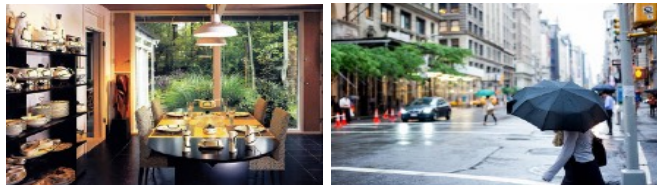


PAELLA Model



Experimental Protocol

- Encoder: Multilingual CLIP
- Decoder: XGLM-2.9B
- Training data:
 - 566K captions sampled from COCO-35
- Evaluation: XM-3600
 - 3600 geographically-diverse images
 - 36 languages with 100 captions per image
 - 5 low-resource languages (L5):
 - Bengali, Cusco Quechua, Maori, Swahili, Telugu



Example training images from COCO



Examples evaluation images from XM3600

Results

	Data	Trained Θ	L36	L5
PaLI	12B	17B	53.6	-
Lg _{COCO-35}	19M	2.6B	15.0	12.5
mBLIP: BLOOMZ-7B	135M	800M	23.4	6.7
BB+CC _{COCO-35 + CC-35}	135M	800M	28.5	22.4
mBLIP: mT0-XL	489M	124M	28.3	7.9
PAELLA	566K	30M	26.2	20.7

PAELLA is competitive against models with 35-863x more training data, and 4-87x more trained parameters

Zero-shot Multilingual Transfer

- **PAELLA_{mono}** is a variant trained on 566K examples in English COCO
- Outperforms **Lg** trained on 19.8M examples in the machine translated COCO-35 dataset

	Data	Trained Θ	L36	L5
Lg: Thapliyal et al. COCO-35	19M	2.6B	15.0	12.5
PAELLA_{mono}	566K _e n	30M	15.5	12.1

Qualitative Example



类似图片显示:

ऐसी ही तस्वीरें दिखाती हैं:

Imágenes similares muestran:

Similar images show:

the owl is perched outside in front of the people
an owl sitting a top a table during the daytime
an owl is sitting on a perch at a camp site
the fuzzy owl is sitting on a tree branch

A caption I can generate to describe this image in english is:

PAELLA

en: "an owl sitting on top of a tree"

es: "un búho sentado en una rama de un árbol"
(an owl sitting on a tree branch)

hi: "एक उल्लू एक पेड़ की टहनी पर बैठा है"
(an owl is sitting on a tree branch)

zh: "一只猫头鹰 站在 树上"
(an owl standing in a tree)

NoRAG

en: "a large black and white picture of a bird"

es: "un pájaro posado en la parte superior de un edificio"
(a bird perched on the top of a building)

hi: "एक पेड़ के पास खड़ा एक पक्षी"
(a bird standing near a tree)

zh: "一只长颈鹿 坐在 树枝上"
(a giraffe sitting on a branch)

Q: Do you even train?

LMCap: Few-shot Multilingual Image Captioning by Retrieval Augmented Language Model Prompting

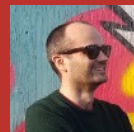
Findings of ACL 2023



R. Ramos



B. Martins



D. Elliott

Socratic Models

- Enable models to “communicate” with each other through their output labels, prompting, and ranking

$$f_{\text{VLM}}^3(f_{\text{LM}}^2(f_{\text{VLM}}^1(\text{image})))$$



Socratic Models

- Enable models to “communicate” with each other through their output labels, prompting, and ranking

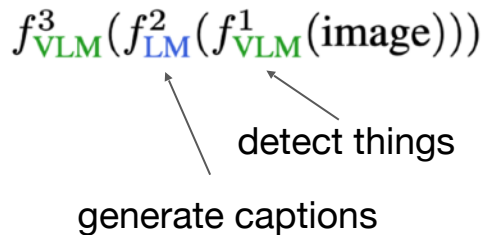
$$f_{\text{VLM}}^3(f_{\text{LM}}^2(f_{\text{VLM}}^1(\text{image})))$$

detect things



Socratic Models

- Enable models to “communicate” with each other through their output labels, prompting, and ranking

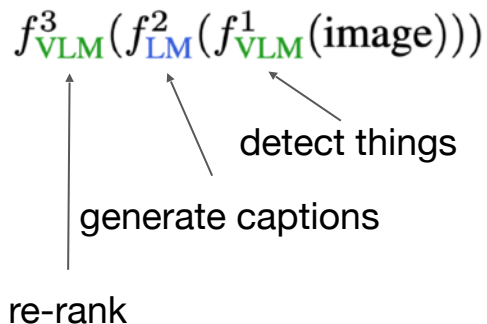


I am an intelligent image captioning bot. This image is a {img_type}. There {num_people}. I think this photo was taken at a {place1}, {place2}, or {place3}. I think there might be a {object1}, {object2}, {object3},... in this {img_type}. A creative short caption I can generate to describe this image is:



Socratic Models

- Enable models to “communicate” with each other through their output labels, prompting, and ranking

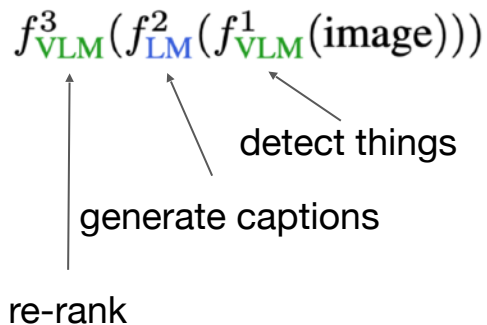


I am an intelligent image captioning bot. This image is a {img_type}. There {num_people}. I think this photo was taken at a {place1}, {place2}, or {place3}. I think there might be a {object1}, {object2}, {object3},... in this {img_type}. A creative short caption I can generate to describe this image is:



Socratic Models

- Enable models to “communicate” with each other through their output labels, prompting, and ranking

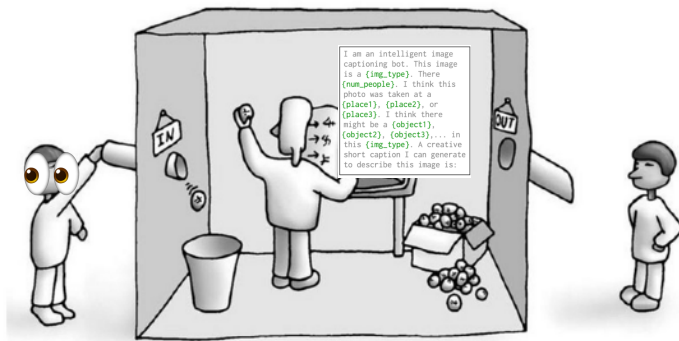


I am an intelligent image captioning bot. This image is a {img_type}. There {num_people}. I think this photo was taken at a {place1}, {place2}, or {place3}. I think there might be a {object1}, {object2}, {object3},... in this {img_type}. A creative short caption I can generate to describe this image is:



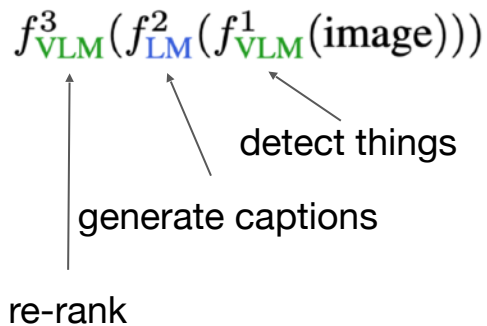
SM (ours): This image shows an inviting dining space with plenty of natural light.

ClipCap: A wooden table sitting in front of a window.

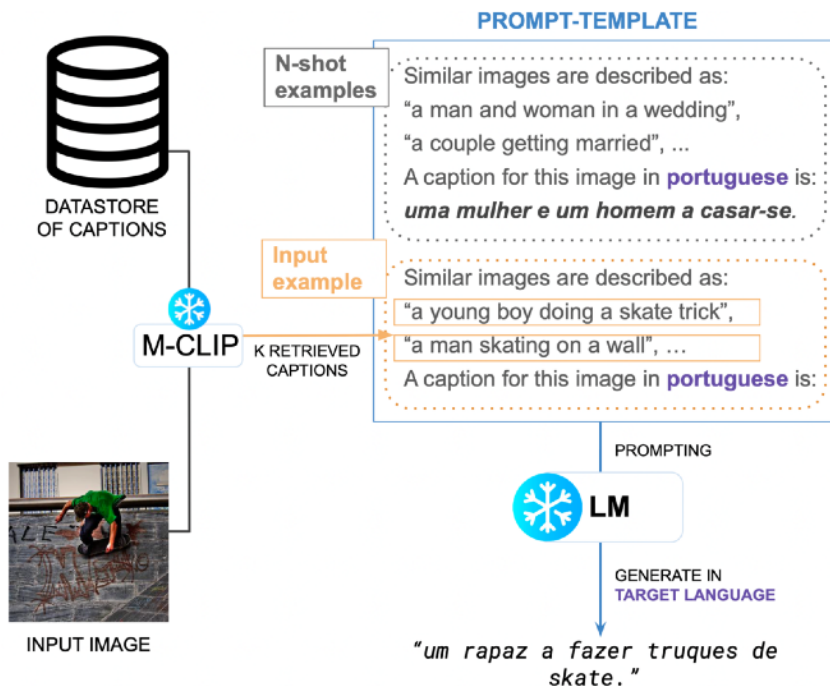


What does it mean to only understand symbols as defined by other symbols?

Multilingual Captioning with Retrieval Augmentation



- Prompt with N-shot examples of transforming captions into the target language



Multilingual Captioning with Retrieval Augmentation

$$f_{\text{VLM}}^3(f_{\text{LM}}^2(f_{\text{VLM}}^1(\text{image})))$$

re-rank
generate captions
detect things



DATASTORE
OF CAPTIONS



M-CLIP



INPUT IMAGE

N-shot
examples

PROMPT-TEMPLATE

Similar images are described as:
"a man and woman in a wedding",
"a couple getting married", ...
A caption for this image in **portuguese** is:
uma mulher e um homem a casar-se.

Input
example

Similar images are described as:
"a young boy doing a skate trick",
"a man skating on a wall", ...
A caption for this image in **portuguese** is:

K RETRIEVED
CAPTIONS

PROMPTING

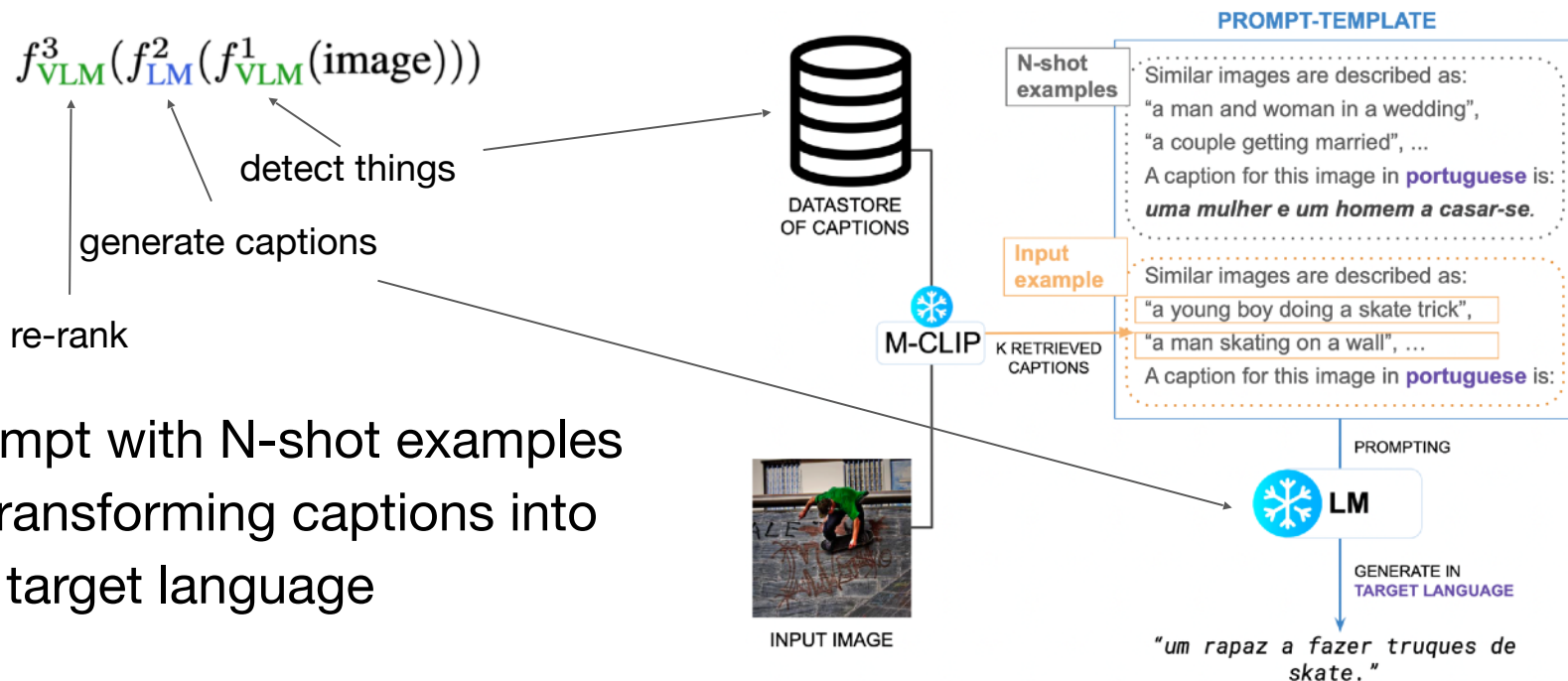


GENERATE IN
TARGET LANGUAGE

"um rapaz a fazer truques de skate."

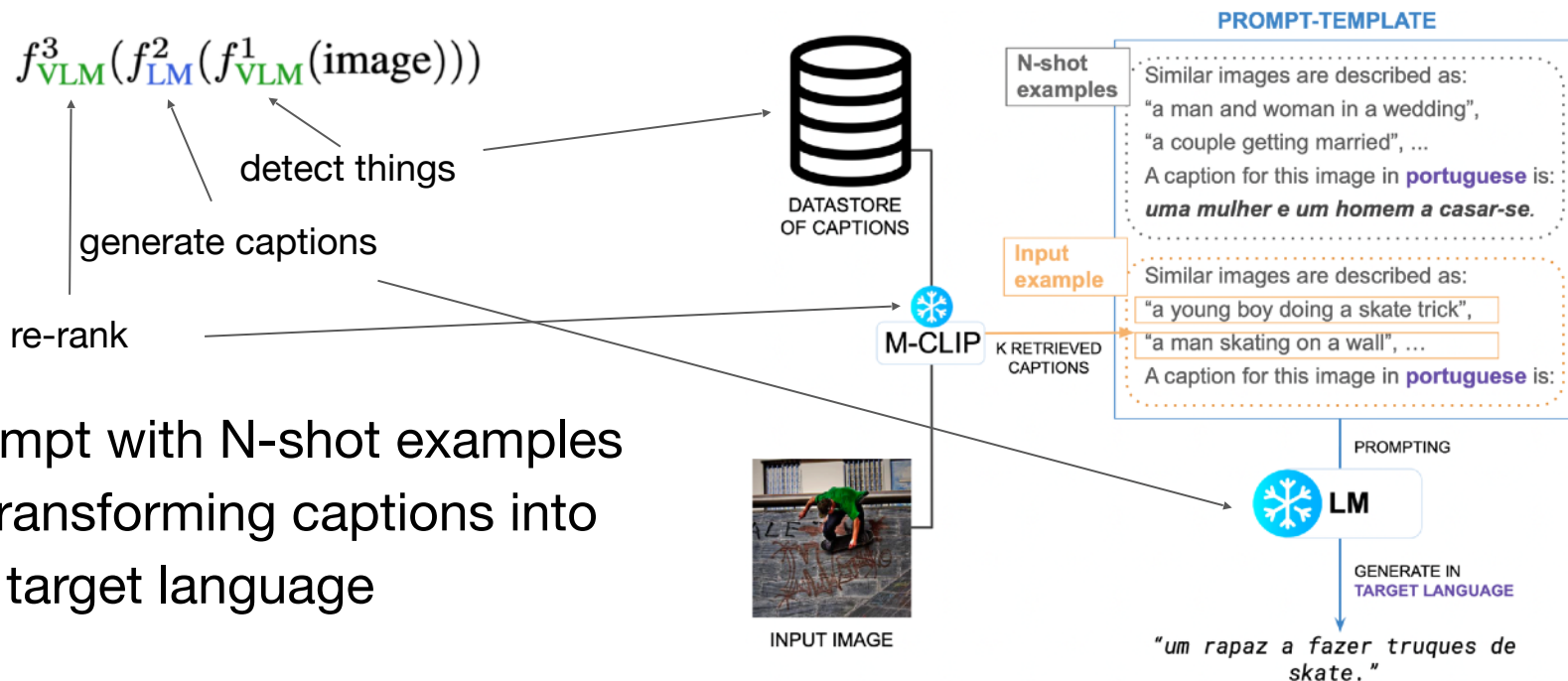
- Prompt with N-shot examples of transforming captions into the target language

Multilingual Captioning with Retrieval Augmentation



- Prompt with N-shot examples of transforming captions into the target language

Multilingual Captioning with Retrieval Augmentation

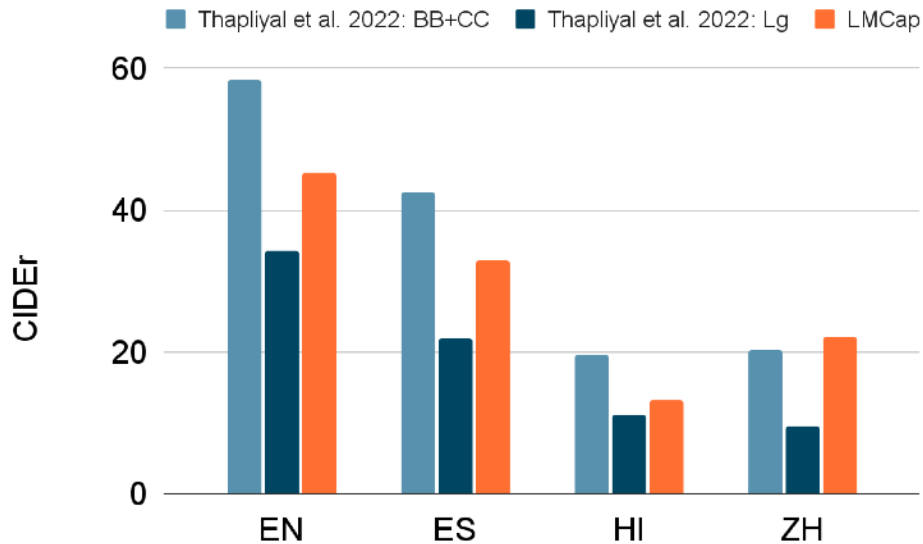


- Prompt with N-shot examples of transforming captions into the target language

Experimental Setup

- XGLM Language Model 564M - 7.6B params
- Multilingual CLIP (LAION)
- Experiments on XM3600 (Thapliyal et al. 2022)
 - 100 images in 36 languages
- **No training or fine-tuning on any captioning data.**

Results



Competitive performance
compared to supervised models

Params	RAM	en	es	hi	zh
564M	6G	0.411	0.094	0.030	0.146
1.7B	12G	0.637	0.143	0.066	0.272
2.9B	16G	0.767	0.454	0.334	0.584
7.5B	22G	0.787	0.489	0.365	0.644

Need at least 2.9B parameter
decoder for multilingual generation

Qualitative Example

Retrieved Examples



two people and a kid skiing along a trail

an adult and two children are cross country skiing

two men and a little boy are skiing on a snowy spot

two adults on skis with a child on skis between them

Qualitative Example

Retrieved Examples



two people and a kid skiing along a trail

an adult and two children are cross country skiing

two men and a little boy are skiing on a snowy spot

two adults on skis with a child on skis between them

Generated Captions

ENG: two people and a kid skiing along a trail

ESP: dos hombres y un niño esquiando en una pista de nieve

ZHO: 两个大人和一个小男孩在雪地上滑雪

Q: How does all of this work?

Understanding Retrieval Robustness for Retrieval-augmented Image Captioning

ACL 2024



W. Li



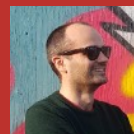
J. Li



R. Ramos



R. Tang



D. Elliott

Revisiting Swanson Soup

- In Ramos et al. CVPR 2023, we observed the power of in-context learning and retrieval-augmentation
- But what is happening here?
- How is the model using the retrieved captions?



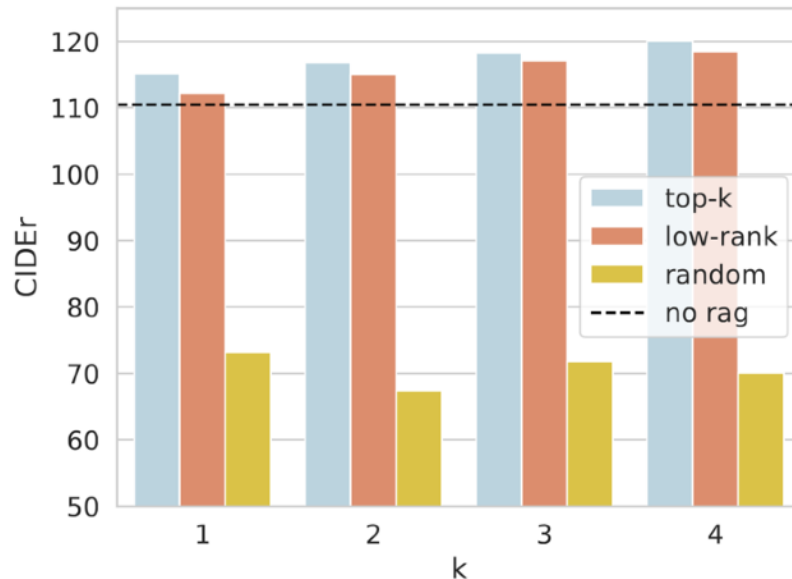
- a can of swanson fat free chicken broth
- a can of swanson brand chicken broth with less sodium
- a 14,5 ounce can of swanson branded chicken broth
- a can of swanson chicken broth on a table

Generated caption:

a can of swanson brand chicken broth on a table

Measuring Robustness

- Is SmallCap sensitive to the *quality* of the retrieved captions?
 - Top-ranked items
 - Random items
 - Lower-ranked items



Question: If the model is so affected by random captions, then is it more like a paraphrasing model that ignores the visual content?

Majority Token Analysis

- Given a list of K retrieved captions, we can create an ordered list of the frequency that each unique token appears in the captions:

$$C := \{C_{T_1}, C_{T_2}, \dots, C_{T_U}\}$$
$$C_{T_i} = \sum_{k=1}^K \mathbf{1}_{T_i \in R_n}$$

- Majority Token:** If token T_i appears at least $K/2$ times, then we define it as majority token in the retrieved captions:

$$M_T := \{C_{T_i}, C_{T_j}, \dots\} \quad \text{s.t.} \quad C_{T_u} \geq K/2 \quad \forall U$$

Majority Tokens Example

- **Majority Token:**

$$M_T := \{C_{T_i}, C_{T_j}, \dots\} \quad \text{s.t.} \quad C_{T_u} \geq K/2 \quad \forall U$$



R₁: Three people skiing through a forest

R₂: An older woman in a wheelchair holding a white teddy

R₃: A man and a woman sit holding a teddy bear

Majority Tokens: “teddy”, “bear”, “woman”

Known Good / Known Bad Captions

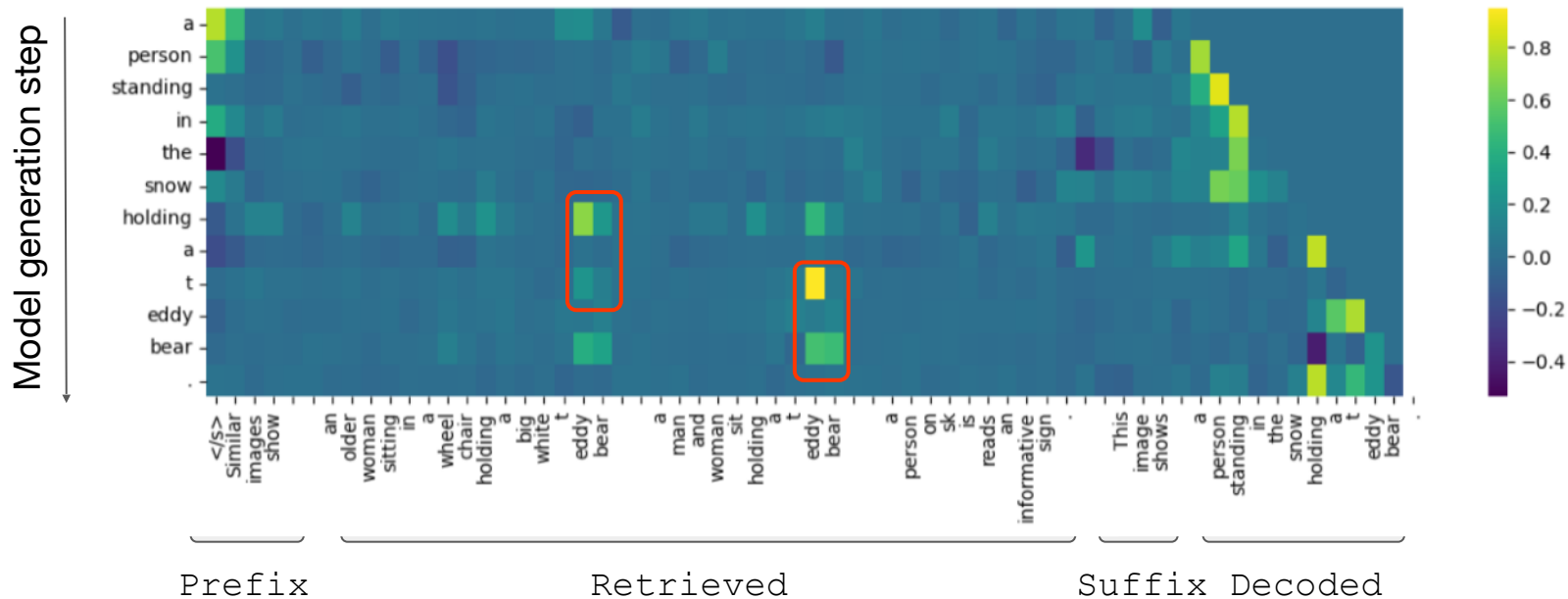
- With Majority Tokens, we can force an experimental setup with known *good* or known *bad* retrieved captions
- Force an asymmetry:
 - 2 Good captions ~ 1 Bad caption \Rightarrow *useful* majority tokens?
 - 1 Good caption ~ 2 Bad captions \Rightarrow *harmful* majority tokens?
- **Good**: high-ranked caption **Bad**: random caption

Known Good / Known Bad Captions

- With Majority Tokens, we can force an experimental setup with known *good* or known *bad* retrieved captions
- Force an asymmetry:
 - 2 Good captions ~ 1 Bad caption \Rightarrow *useful* majority tokens?
 - 1 Good caption ~ 2 Bad captions \Rightarrow *harmful* majority tokens?
- **Good**: high-ranked caption **Bad**: random caption
- Results
 - 2 Good ~ 1 Bad: 86% of generated captions contain a majority token
 - 1 Good ~ 2 Bad: 21%

Integrated Gradients Input Attribution

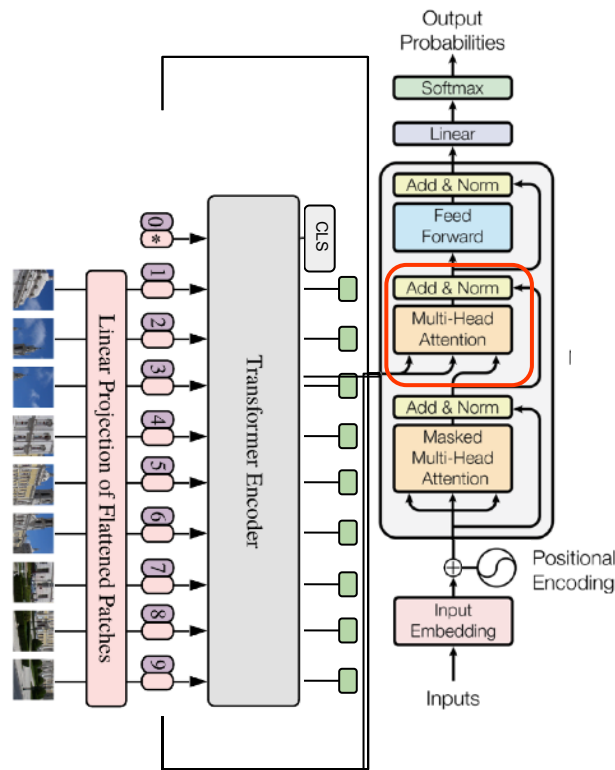
- Which input tokens are most/least important in the model output?



Self- and Cross-Attention Analysis

- What can we learn about SmallCap by inspecting what it attends to in the textual and the visual inputs?
- Track the location of the **maximally**-attended inputs

$$\mathbb{1}[I_n(i, j)] = \begin{cases} 1 & \text{if } \arg \max_z \text{Att}(j, z)_i \in S_n \\ 0 & \text{otherwise} \end{cases}$$



Self-Attention Analysis

<S>

← BOS

Similar images show

← Prefix

a man working some levers at a train yard

a train engineer preparing the engine of a train

← Retrieved

a train being worked on in a train manufacturer

a man wearing a safety vest standing by a train.

← Suffix

This image shows

a person working ...

← Generated

Self-Attention Analysis

<S>

Similar images show

a man working some levers at a train yard

a train engineer preparing the engine of a train

a train being worked on in a train manufacturer

a man wearing a safety vest standing by a train.

This image shows

a person working ...

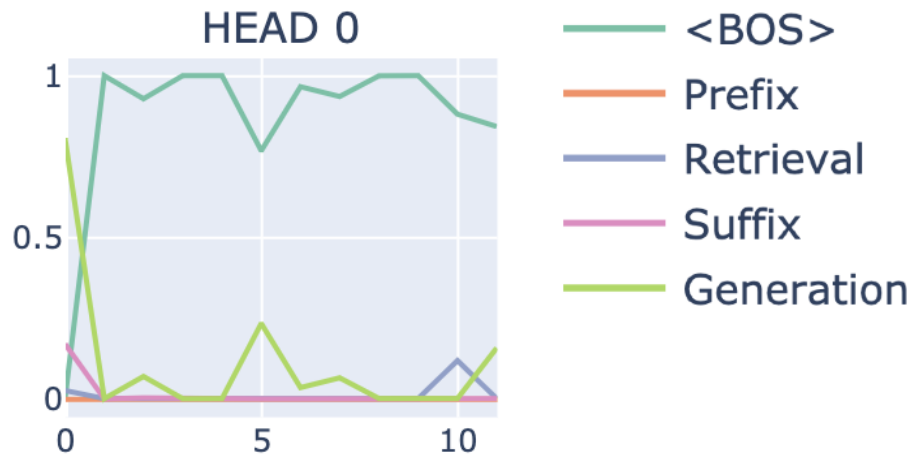
← BOS

← Prefix

← Retrieved

← Suffix

← Generated



Self-Attention Analysis

<S>

Similar images show

a man working some levers at a train yard

a train engineer preparing the engine of a train

a train being worked on in a train manufacturer

a man wearing a safety vest standing by a train.

This image shows

a person working ...

← BOS

← Prefix

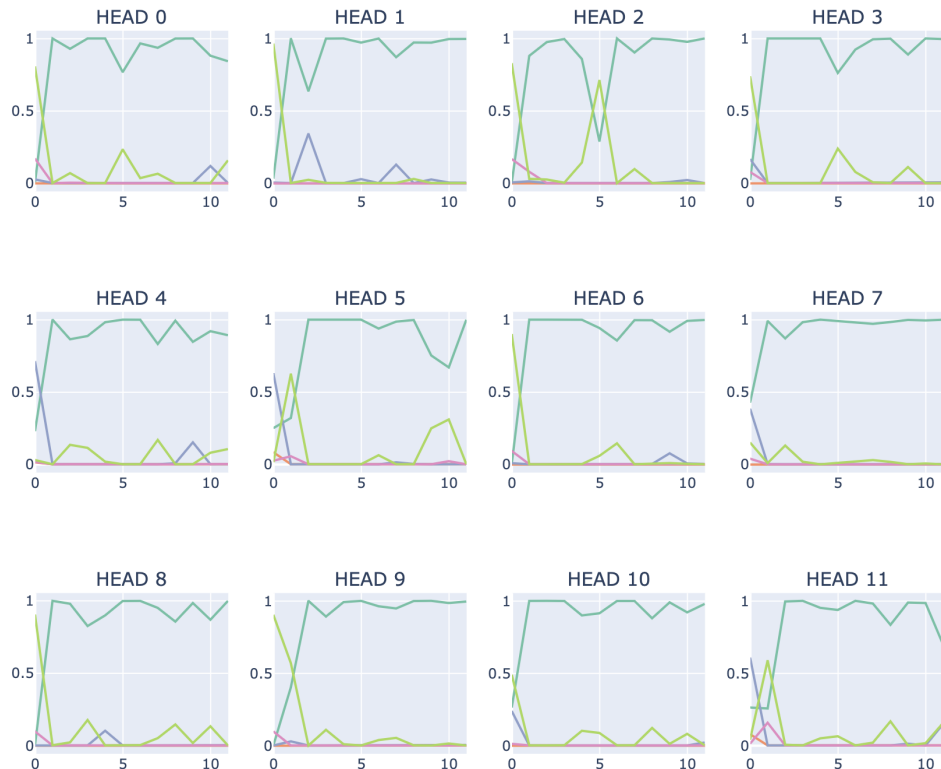
← Retrieved

← Suffix

← Generated

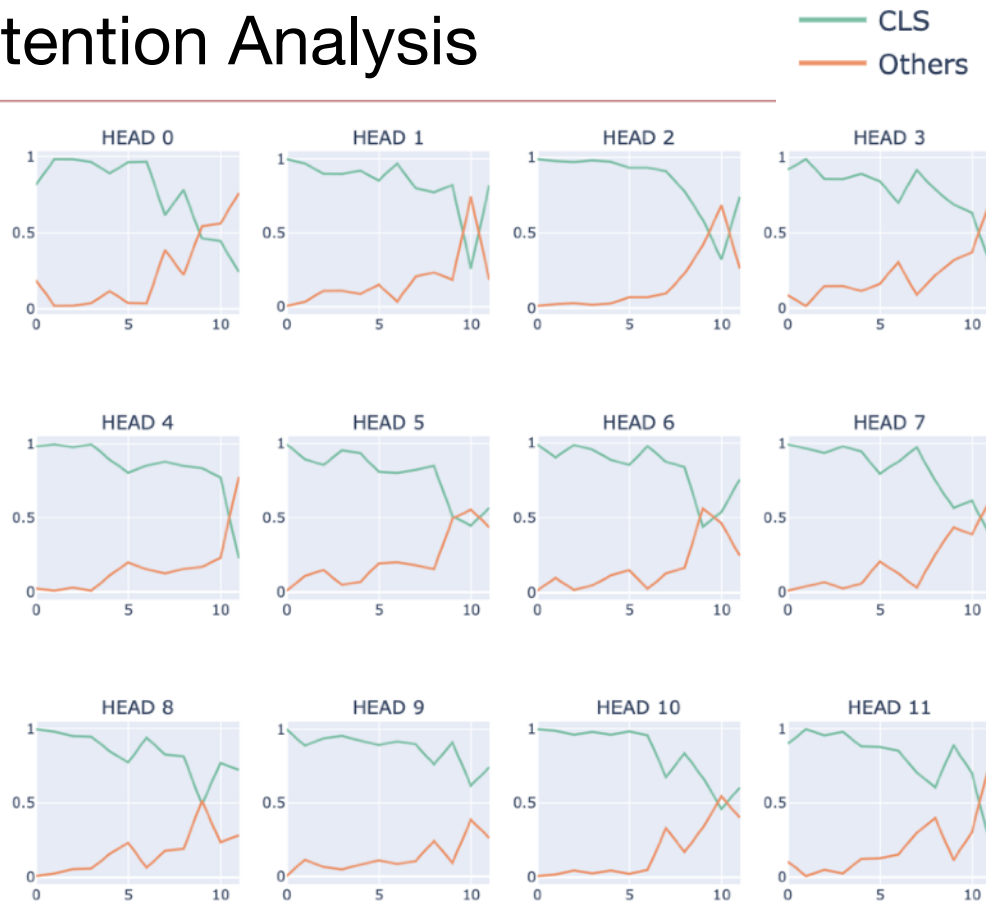
**Self-attention is
maximally attending
to the BOS token**

— <BOS>
— Prefix
— Retrieval
— Suffix
— Generation



Cross-Attention Analysis

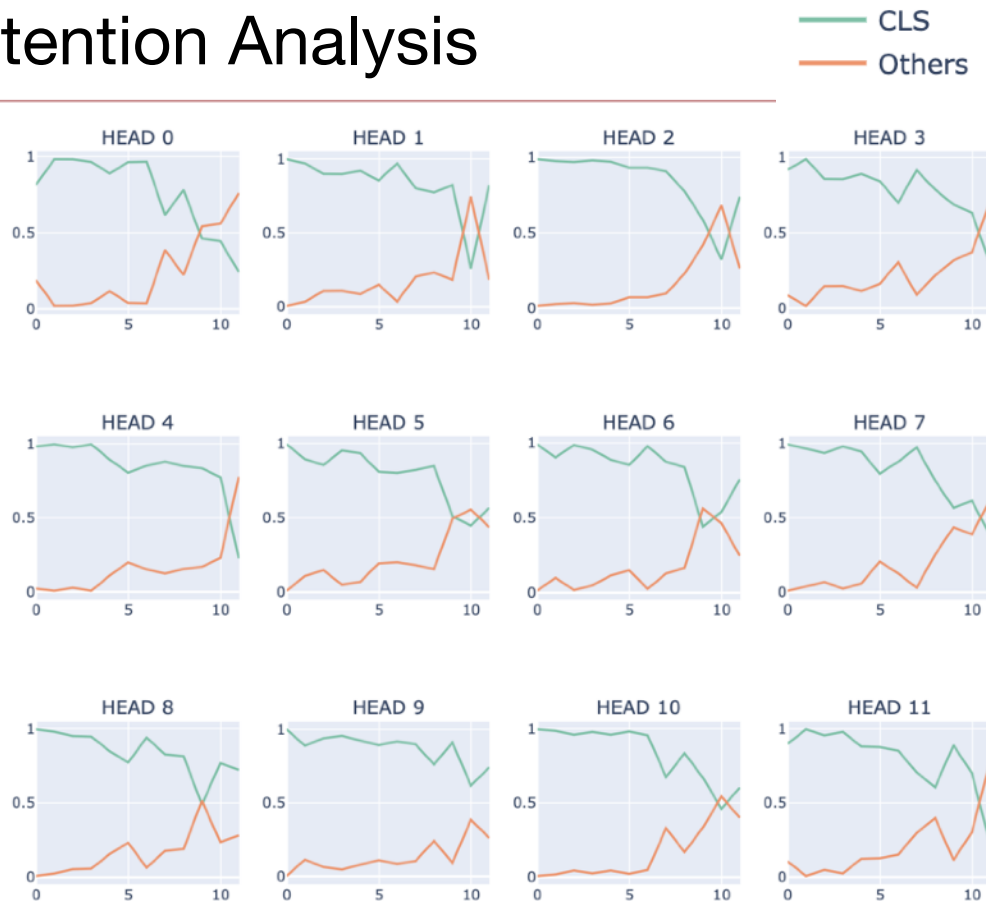
The cross-attention layers focus on the “summary”
image CLS embedding



Cross-Attention Analysis

The cross-attention layers focus on the “summary” **image CLS** embedding

Cross-attention to **image patches** only emerges at the final layers of the LM



Improving Robustness

- Given that the model appears to be strongly guided by retrieved captions, can we train the model to be less reliant on this?
 - Yes! We can create less-perfect retrieval lists during training
- **Sample-K**: randomly choose k/N retrieved captions
- **C-Sample-K**: only use the most relevant caption, and $k-1$ randomly sampled captions in the prompt

Experimental Protocol

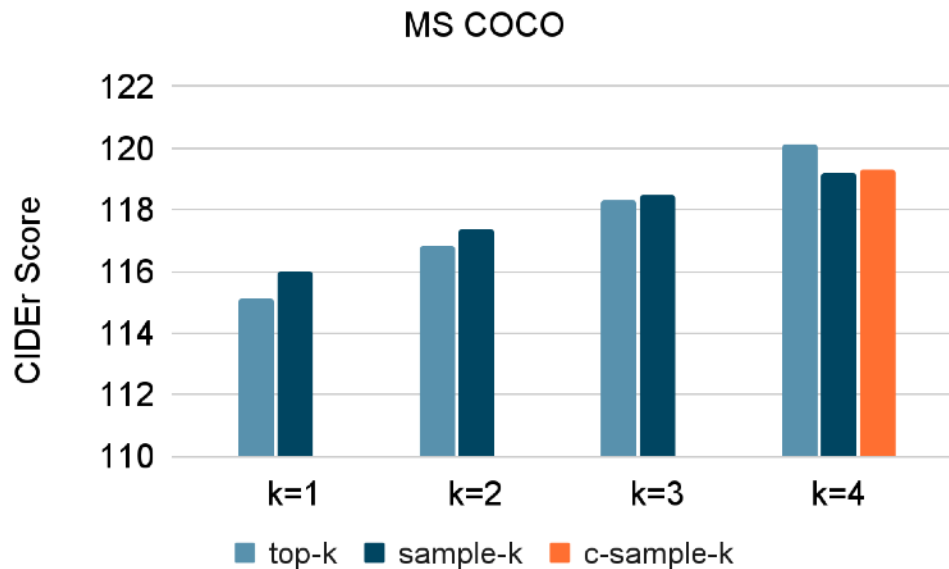
- Encoder: CLIP ViT-B/32
- Decoder: OPT-125M
- Training data: MS COCO



Example evaluation images from the NoCaps dataset

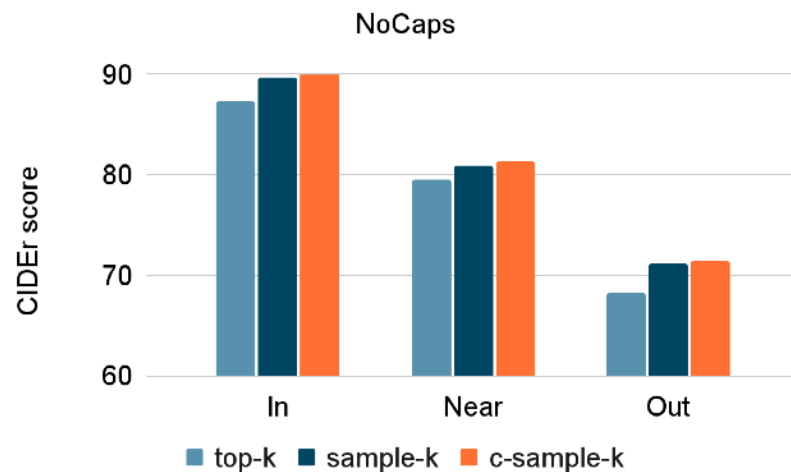
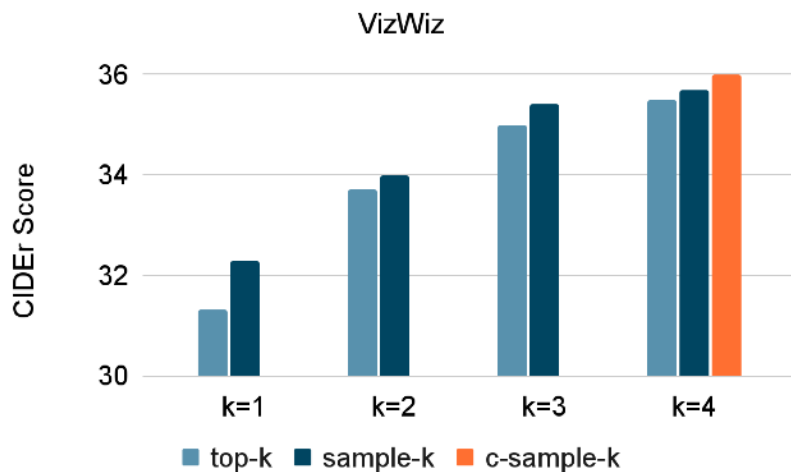
- Evaluation with CIDEr
 - MS COCO
 - VizWiz
 - NoCaps

In-Domain Results



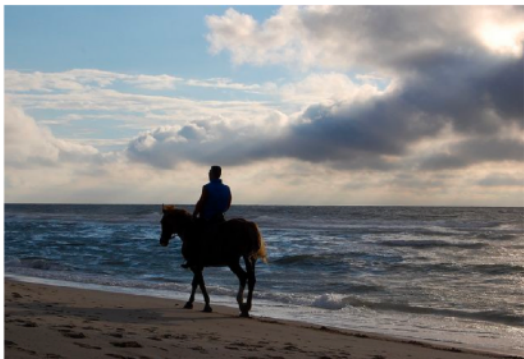
Improved performance with smaller retrieval sets

Out-of-Domain Results



Improvements in two out-of-domain datasets

Qualitative Examples



- a man posing with a surfboard on an elevator
- a woman sitting on a bench next to a man in a hat
- a greyhound dog lying on an unmade bed
- a pink teddy bear and a brown teddy bear sitting on wooden rods

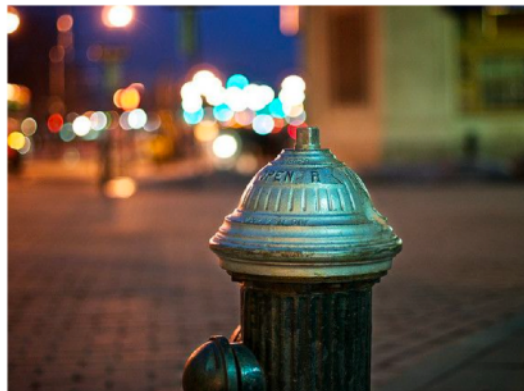


Sample-k

a person riding a **horse** on top of a beach

Top-k

a person sitting on a **bench** on a beach



- a train with the numbers 60016 is heading down the tracks
- a black and white photo of two people holding hands in a city on a rainy day
- this youngster has a boogie board to ride the smaller waves
- a wooden entertainment center containing a television set



Sample-k

a close up of a **fire hydrant** on a sidewalk

Top-k

a close up of a **person** on a sidewalk

Wrap-up

Open Questions

- How many of these observations apply to visual prefix models?
 - I think we will still observe the problems associated with majority tokens
- What is the best way to construct N-shot examples for mRAG?
 - Demonstrate the diversity of the tasks / target languages / visual inputs
- When will we have *usable* multimodal ICL for multimodal RAG?
 - We have been trying to make progress on this with ImageChain

**IMAGECHAIN: Advancing Sequential Image-to-Text Reasoning in
Multimodal Large Language Models**

Danae Sánchez Villegas* Ingo Ziegler* Desmond Elliott

Final Conclusions

- Retrieval-augmentation is a powerful approach to building lightweight image captioning models that can easily adapt to new domains
 - Improve lightweight trained models
 - Improve zero-training models
 - Enable zero-shot multilingual transfer
- Open questions about how to make RAG-based models more robust and reliable in practice

References

- R. Ramos, B. Martins, and D. Elliott. **Few-shot Multilingual Image Captioning by Retrieval Augmented Language Model Prompting**. ACL 2023.
- R. Ramos, B. Martins, D. Elliott, and Y. Kementchedjhieva. **SmallCap: Lightweight Image Captioning Prompted with Retrieval Augmentation**. CVPR 2023.
- R. Ramos, E. Bugliarello, B. Martins, and D. Elliott. **PAELLA: Parameter-Efficient Lightweight Language-Agnostic Captioning Model**. NAACL 2024.
- W. Li, J. Li, R. Ramos, R. Tang, and D. Elliott. **Understanding Retrieval Robustness for Retrieval-augmented Image Captioning**. ACL 2024.
- D. S. Villegas, I. Ziegler, and D. Elliott. **ImageChain: Advancing Sequential Image-to-Text Reasoning in Multimodal Large Language Models**.