



---

# **Vision-Language Pretraining: Current Trends and the Future**

Aishwarya Agrawal & Damien Teney & Aida Nematzadeh

---

ACL  
22 May 2022



**Please check our website for our final slide deck:**

**<https://vlp-tutorial-acl2022.github.io/>**

# Why Multimodal Pretraining?

The ability to ground language to vision—multimodal pretraining—is a fundamental aspect of both language & vision.



Q: What are the people waiting for?

A: bus

# Why Multimodal Pretraining?

**Train once, use multiple times.** Multimodal features are useful across a range of multimodal tasks and applications.



Q: What are the people waiting for?

A: bus



**Goal:** Give an overview of ingredients needed for working on multimodal problems (particularly vision and language). Also, discuss some of the open problems.

### **Plan for today:**

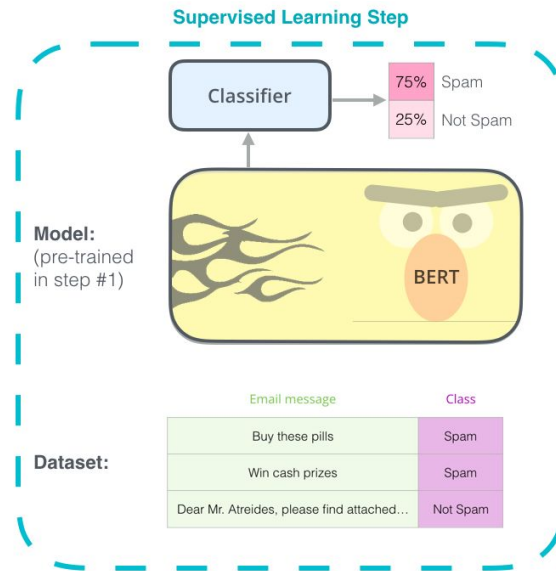
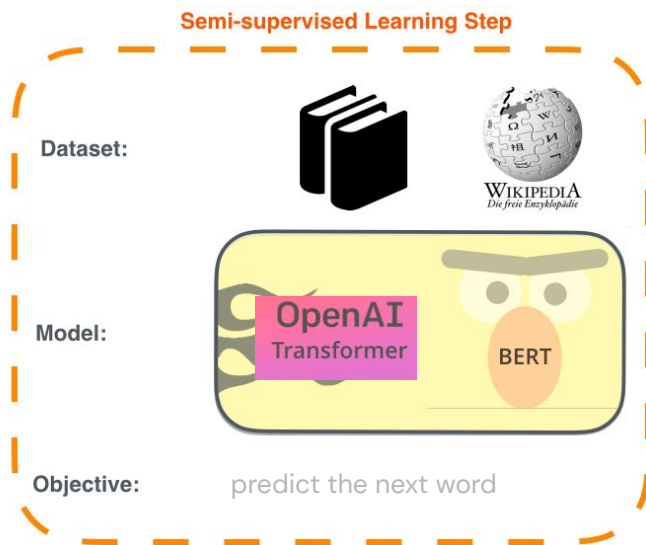
- Before the Pretraining Era (Aishwarya, 55 min)
- Vision-Language Pretraining (Aida, 60 min)
- Beyond Statistical Learning (Damien, 55 min)



# Vision-Language Pretraining

Aida Nematzadeh

# Success of Pretraining in NLP



Performance gain is due to architecture innovations & larger data. [Peters et al., 2018; Howard & Ruder, 2018; Devlin et al., 2018; Radford et al., 2018; Raffel et al., 2019; Rae et al. 2022]

# Similar Models for Multimodal Pretraining?

**Dataset:**

"The scenic route through mountain ranges includes these unbelievably coloured mountains."



**Model:**



**Objective:**

predict the next word

**Other objectives?**

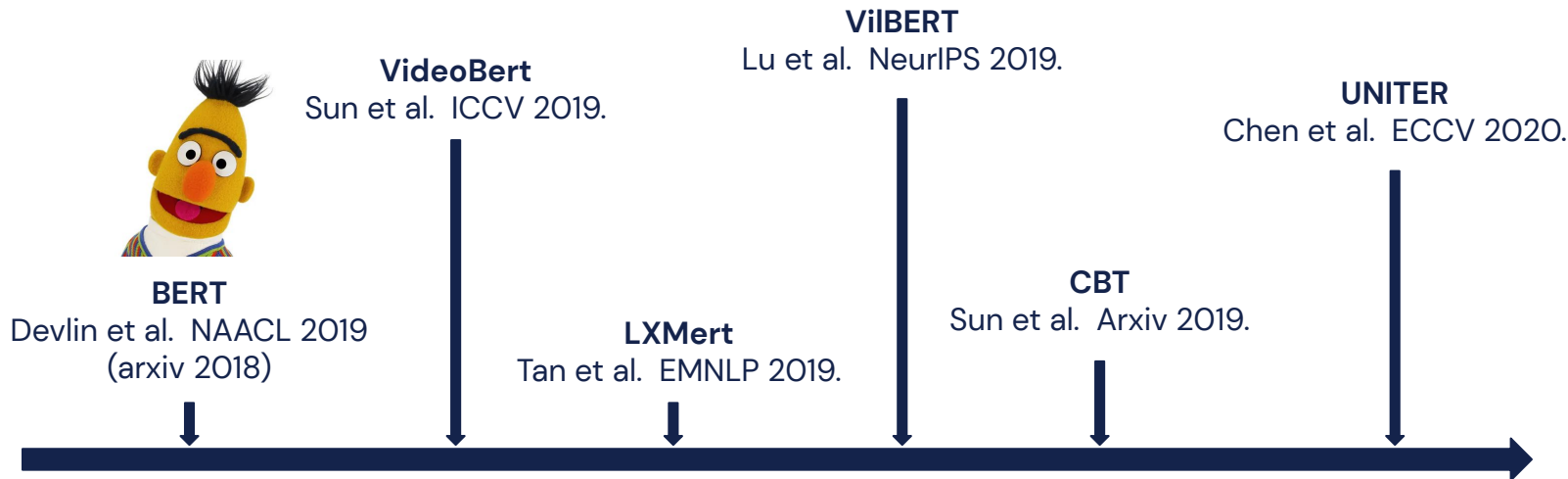
Dataset: image-text pairs where a given text describes its image.

Model: attention mechanisms over both image and text; preprocessing images to "visual tokens".

Objective: loss functions specific to the image modality and image-text pairs.



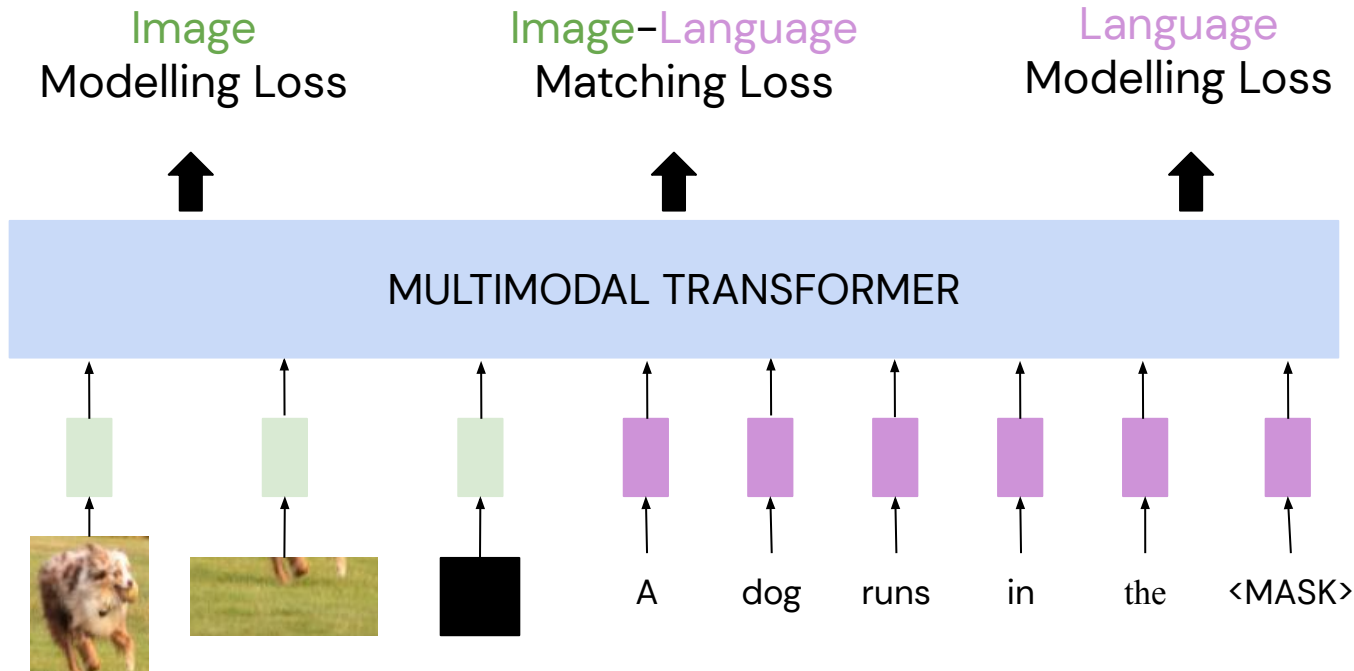
# Multimodal Pretraining: How it Started



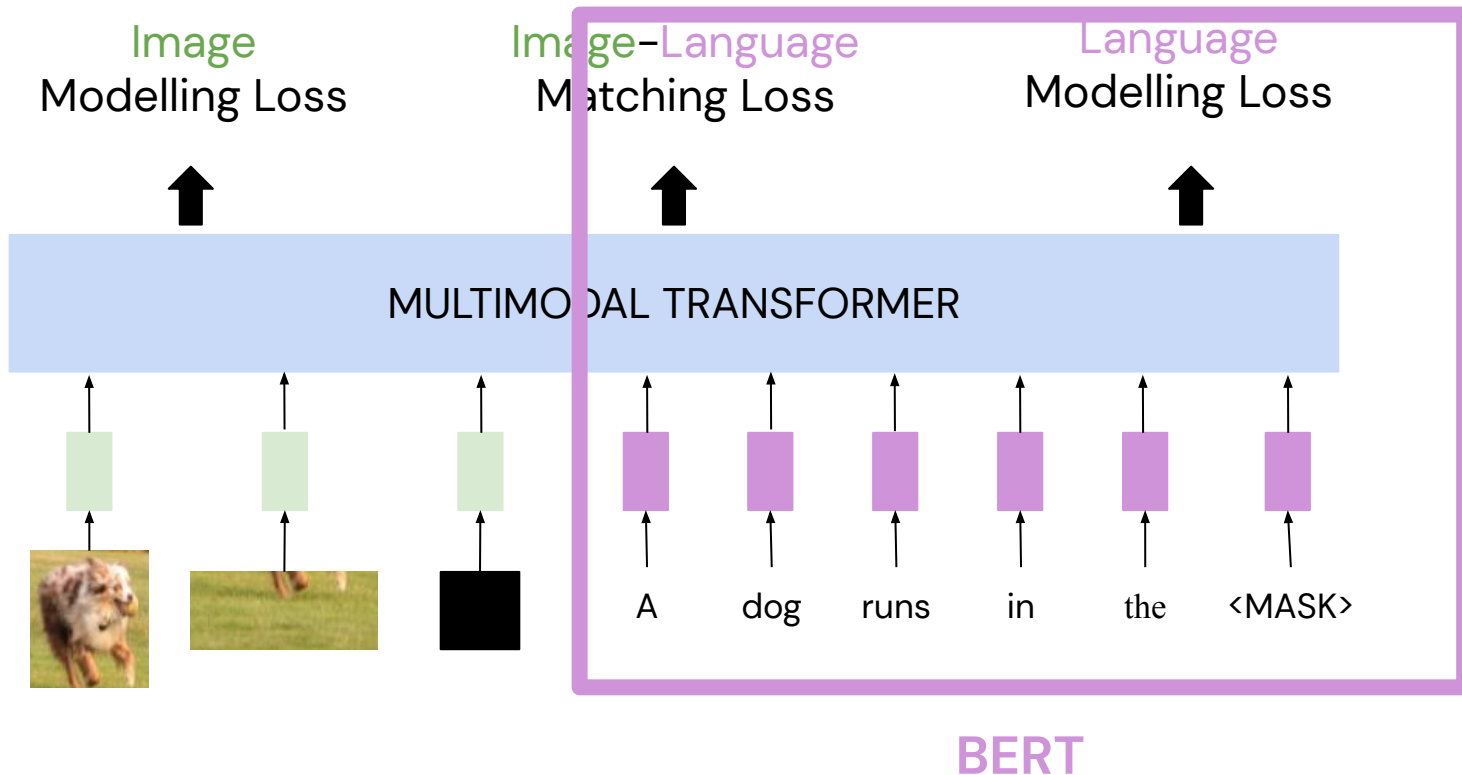
Models share the same “backbone” with slight differences in loss design, preprocessing, etc.

They achieve the state-of-the-art results in a range of tasks.

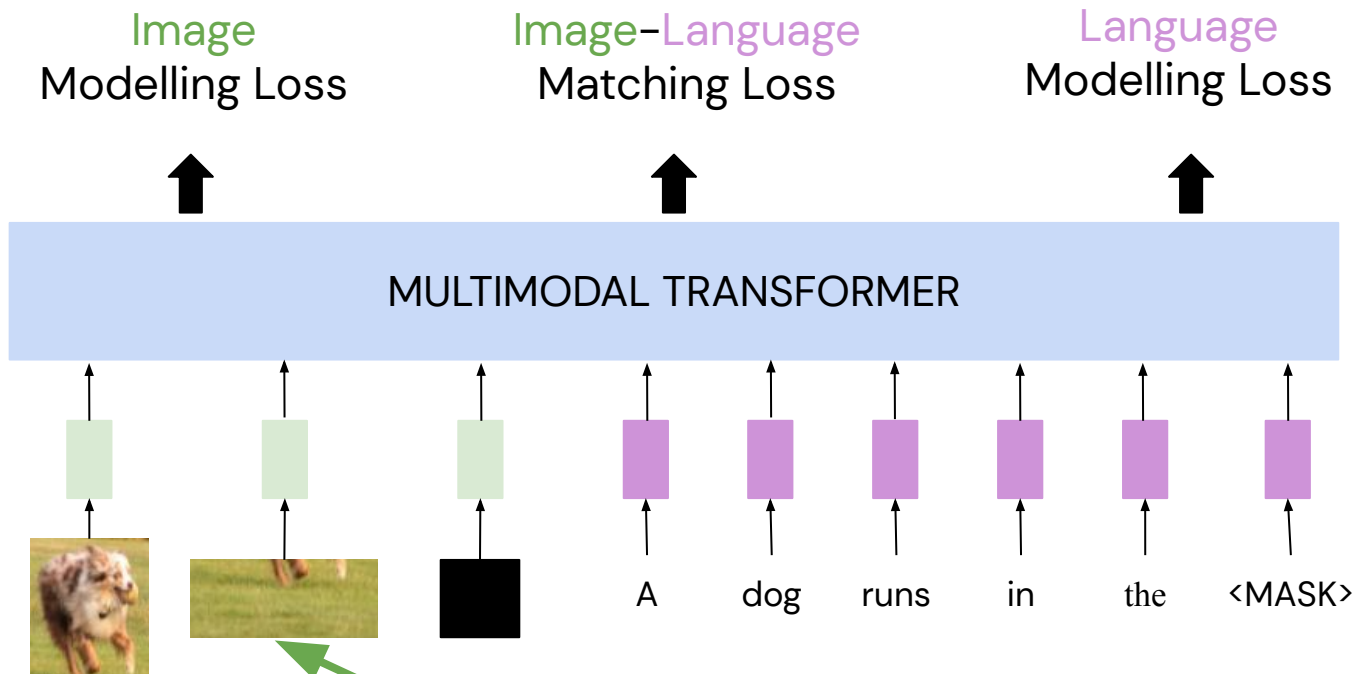
# Multimodal Transformers (Joint Encoders)



# Multimodal Transformers (Joint Encoders)



# Multimodal Transformers (Joint Encoders)



visual "word":  
bounding box

VOLTA: a general framework (with 5  
English and 4 multilingual MMT). 12



# What Contributes to these Models' Success?

Are results due to advances in the architecture or large pretraining datasets?

Are the “adopted” losses from language models good enough?

Is the cross-talk between modalities (via attention) important?

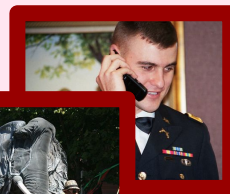
What makes a good pretraining dataset?

# Evaluation: Zero-Shot Image Retrieval

**Zero-shot** image retrieval directly evaluates the goodness of **pretrained** representations.

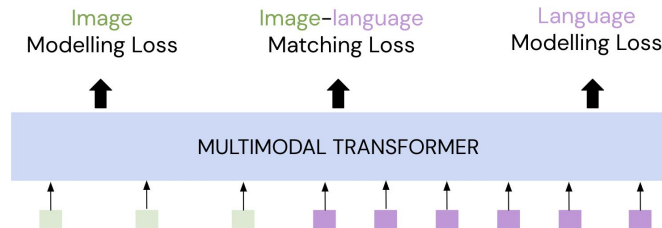
## Image Retrieval (IR)

*"Grey haired man in black  
and yellow tie."*





# Typical Loss Functions



Language/image modeling: masked language/region modeling

Image–language matching: binary classification or contrastive formulation

$$L_{i2t} = -\frac{1}{N} \sum_i \log \frac{\exp(x_i^\top y_i / \sigma)}{\sum_{j=1}^N \exp(x_i^\top y_j / \sigma)}$$

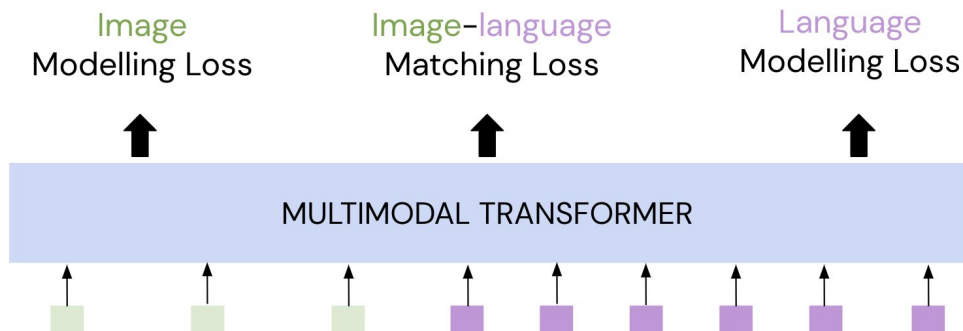
$$L_{t2i} = -\frac{1}{N} \sum_i \log \frac{\exp(y_i^\top x_i / \sigma)}{\sum_{j=1}^N \exp(y_i^\top x_j / \sigma)}$$

[Taken from [ALIGN](#)]

# Are All Losses Needed? [Hendricks et al. TACL 2021]

R@1	No <b>Language</b> Modeling Loss	No <b>Image</b> Modeling Loss	All Loses
Zeroshot Flickr	15.0	<b>41.1</b>	<b>40.7</b>

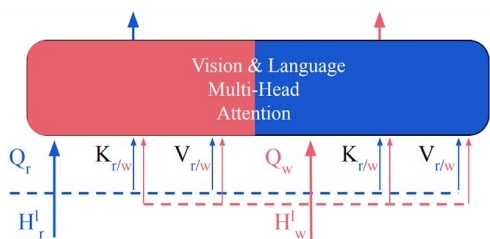
With the right hyper-parameters,  
**image modeling loss** is not needed.



Vision-and-Language  
or  
**Vision-for-Language?**  
[Frank et al, 2021]

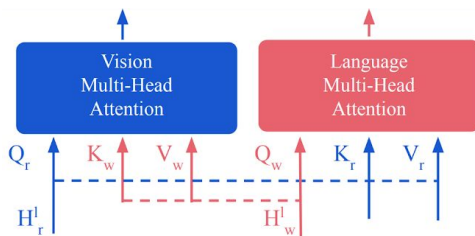


# Different Attention Mechanisms [Hendricks et al. TACL 2021]



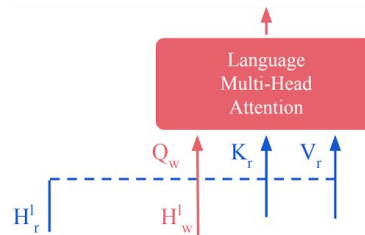
*Merged attention*

Each modality attends to **both** modalities.



*Coattention*

Each modality *attends only* to the other modality (two asymmetric attentions).

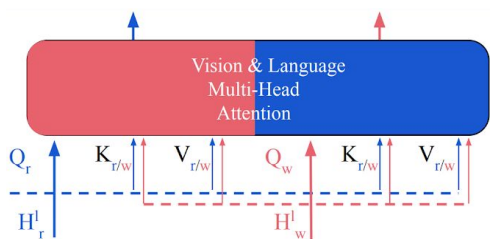


*Asymmetric attention*

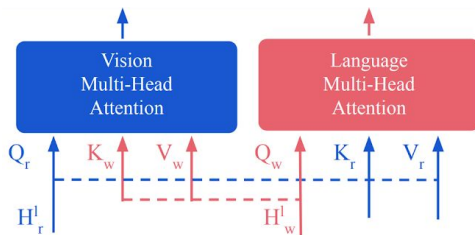
Only one modality (e.g., language) attends to the **other** modality (e.g., image).

multimodal attention

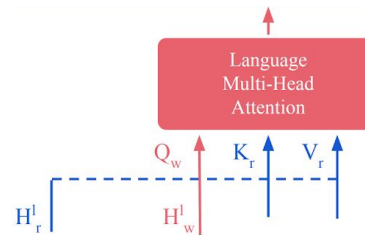
# Different Attention Mechanisms [Hendricks et al. TACL 2021]



*Merged attention*



*Coattention*

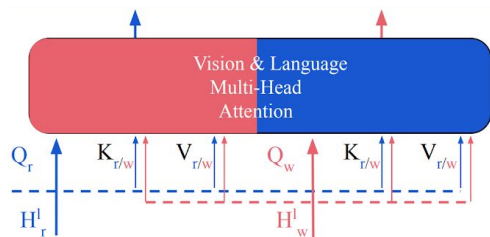


*Asymmetric attention*

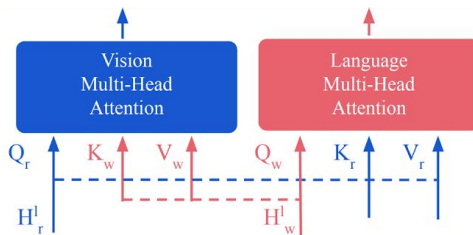
Similar  
performance

R@1	merged	coattention
Zeroshot Flickr	40.0	<b>41.9</b>

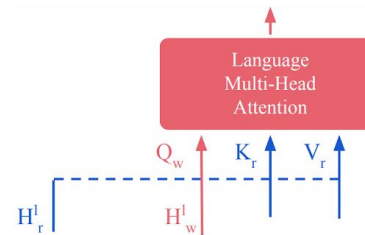
# Different Attention Mechanisms [Hendricks et al. TACL 2021]



*Merged attention*



*Coattention*



*Asymmetric attention*

R@1	merged	coattention	asymmetric (language)	asymmetric (image)
Zeroshot Flickr	40.0	<b>41.9</b>	33.6	31.6



# Multimodal Attention > Depth / Size [Hendricks et al. TACL 2021]

← 6 multimodal layers & 12 attention heads →

R@1	coattention	asymmetric (language)	asymmetric (image)	coattention with 1 multimodal layer	coattention with 6 attention heads
Zeroshot Flickr	<b>41.9</b>	33.6	31.6	37.2	39.9

Depth and number of parameters alone are not enough.



# What Contributes to these Models' Success?

Are results due to advances in the architecture or large pretraining datasets?

Are the “adopted” losses from language models good enough? No, we need better image modeling losses.

Is the cross-talk between modalities (via attention) important? Yes, multimodal attention is important.

What makes a good pretraining dataset?

# Pretraining Datasets

## MSCOCO



"The two people are walking down the beach."

## MSCOCO/OI Narratives



"In this image we can see a bridge and sea. In the background, we can see trees and the sky. We can see so many people on the bridge. At the bottom of the image, we can see two people. We can see stairs in the right bottom of the image ..."

## Visual Genome



small round yellow frisbee, man has cast on his arm, concrete trail path in the park, man wearing black sunglasses

manually annotated

## Conceptual Captions



"The *scenic route* through mountain ranges includes these unbelievably coloured mountains."

## SBU Captions



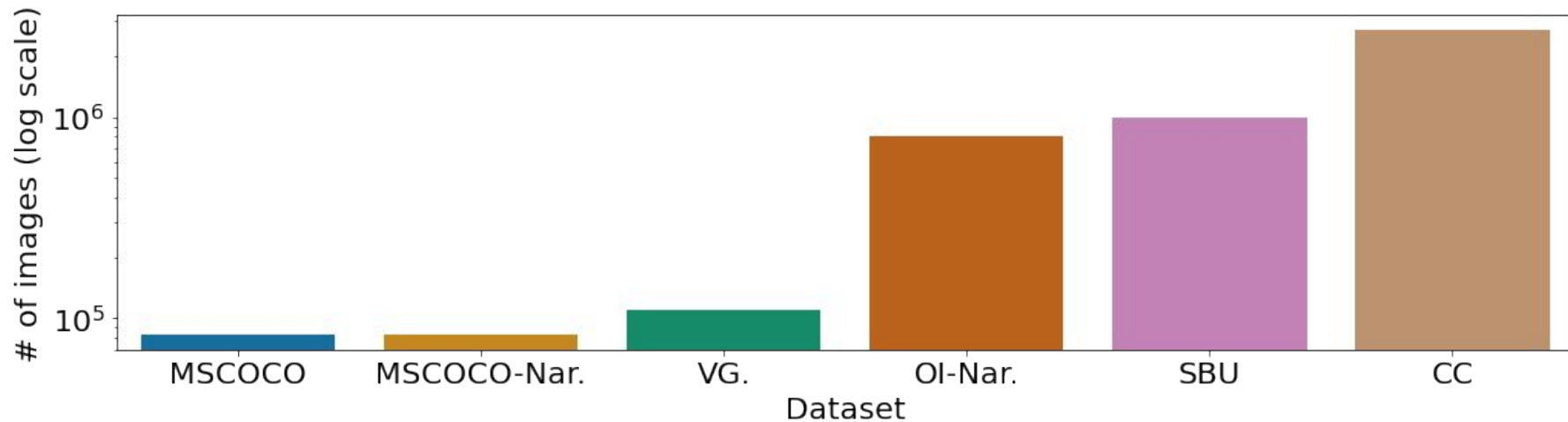
"*King Arthur's* beheading rock - right on the sidewalk in the middle of *town*".

from "the wild"

Noisier image-text correspondence but larger



# Dataset Considerations: Size



# Dataset Considerations: Language



**COCO style caption:** "Single black dog sitting on the grass"

**Narratives style caption:** "The dog is black and brown. The collar is red. ... The dog is on the grass. ..."

**Genome style caption:** "Black dog"



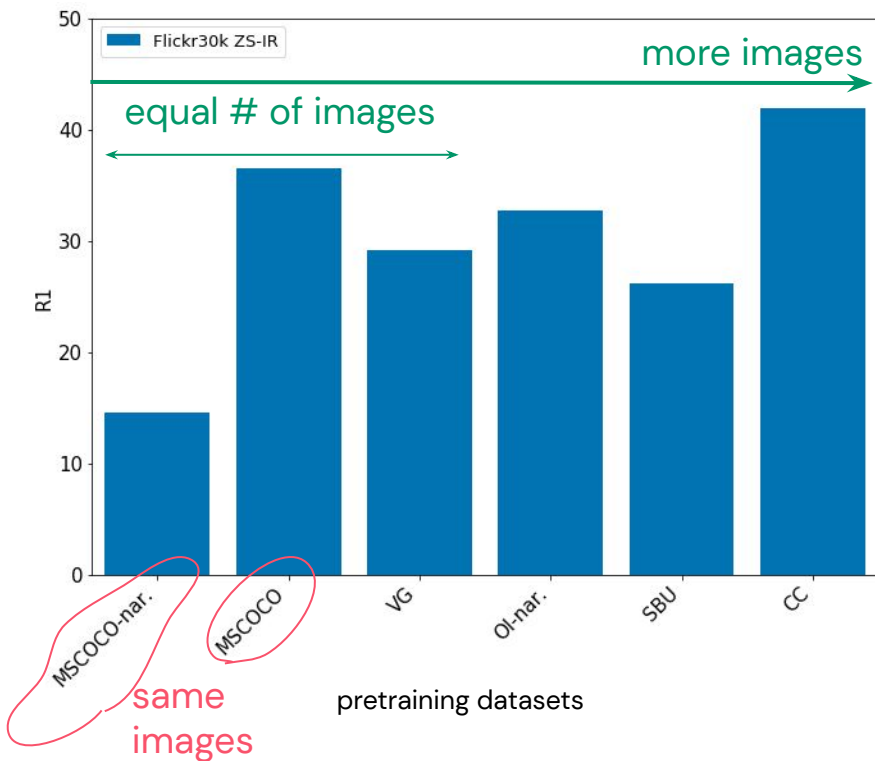
# Dataset Considerations: Noise



"Single *black dog sitting* on the *grass*"

"A *person* takes a *dog* on a *walk* near the *river*."

# Image Retrieval: Language Similarity Matters

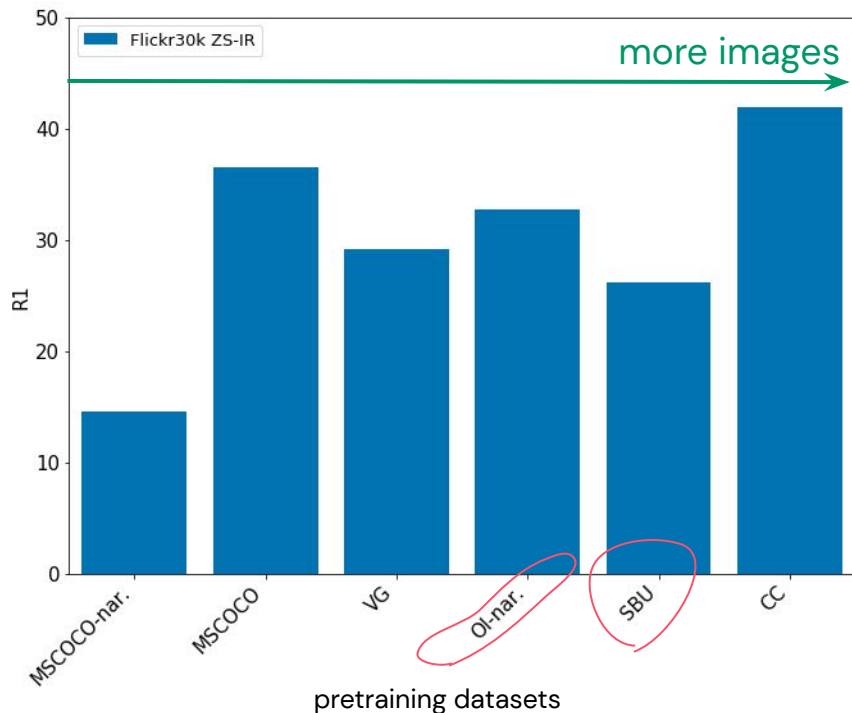


Performance is **not** directly correlated with the number of images.

Language similarity in pretraining & test (measured by perplexity) explains the difference in the results.

Hendricks et al. "Data, Architecture, or Losses: What Contributes Most to Multimodal Transformer Success?" TACL 2021.

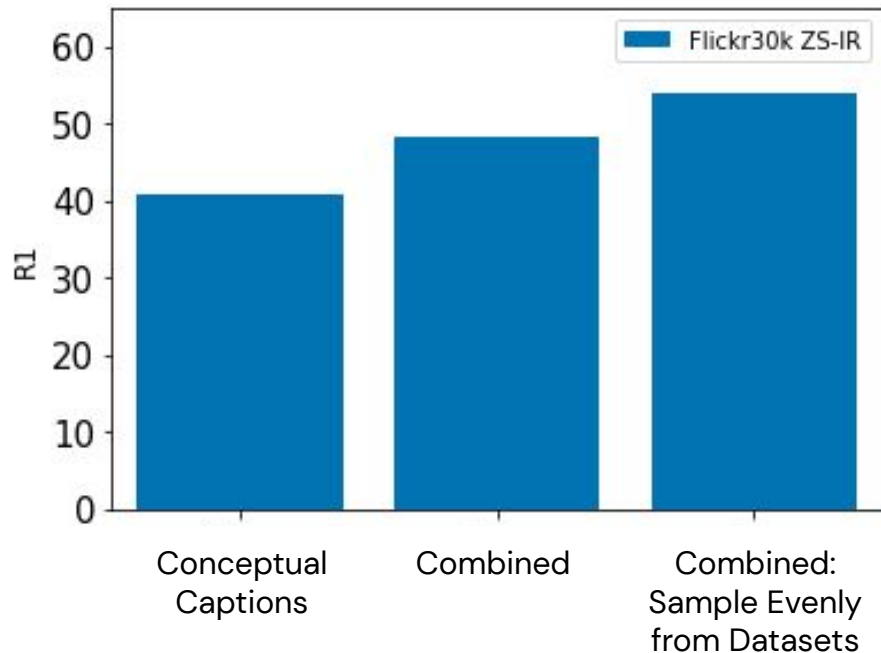
# Image Retrieval: Noise Matters



SBU is larger than Open Images and has lower perplexity, but is still worse. However, SBU has more noise, meaning the language does not always describe the image content.

Hendricks et al. "Data, Architecture, or Losses: What Contributes Most to Multimodal Transformer Success?" TACL 2021.

# Image Retrieval: Dataset Sampling Matters

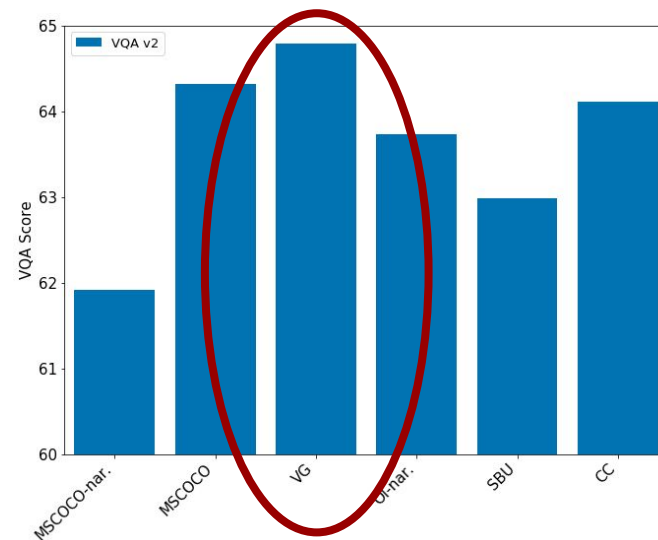
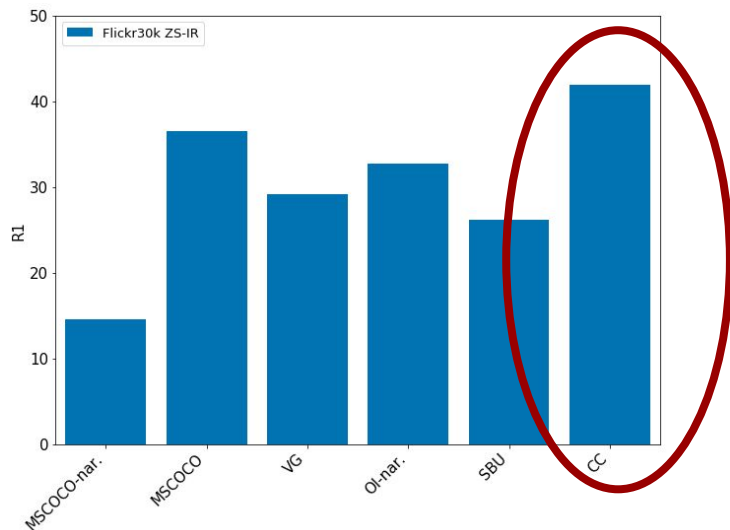


Combining datasets does lead to better results, but how we sample from combined datasets matter.

MSCOCO is a good dataset for pretraining; sampling method which weights MSCOCO images higher does better.

Hendricks et al. "Data, Architecture, or Losses: What Contributes Most to Multimodal Transformer Success?" TACL 2021.

# “Best” Dataset is Task Dependent



Best datasets are different for IR (Conceptual Captions is best) and VQA (VG is best)



# What Contributes to these Models' Success?

Are results due to advances in the architecture or large pretraining datasets?

Are the “adopted” losses from language models good enough? No, we need better image modeling losses.

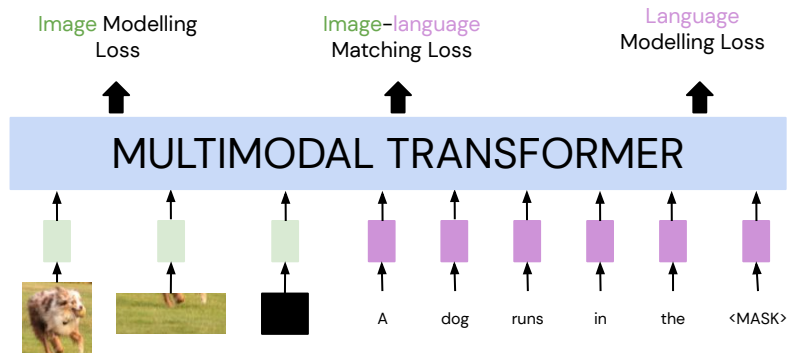
Is the cross-talk between modalities (via attention) important? Yes, multimodal attention is important.

What makes a good pretraining dataset? The level of noise and language matter.

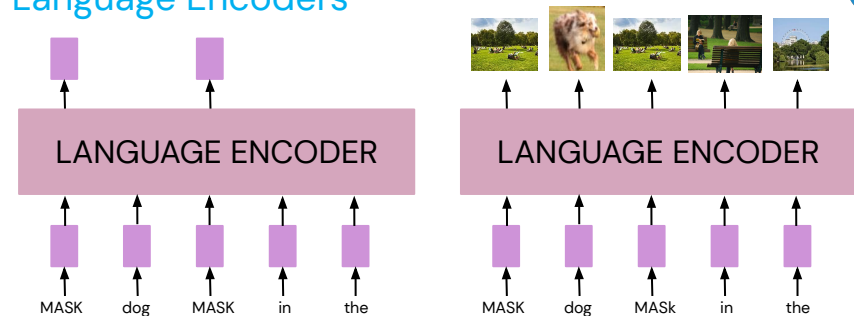
We have released our pretrained models!



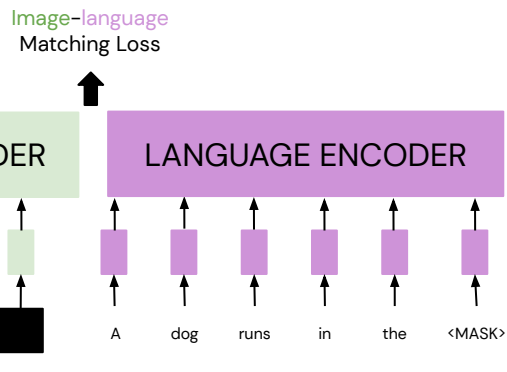
## Joint Encoders



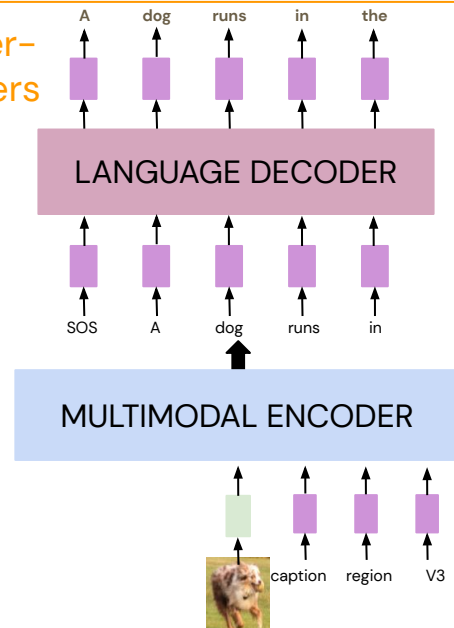
## Language Encoders



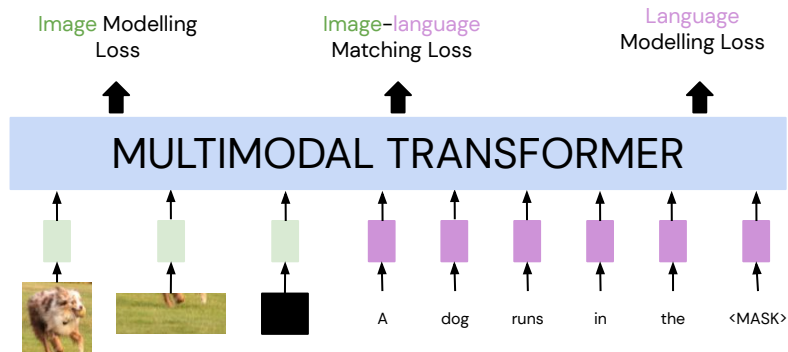
## Dual Encoders



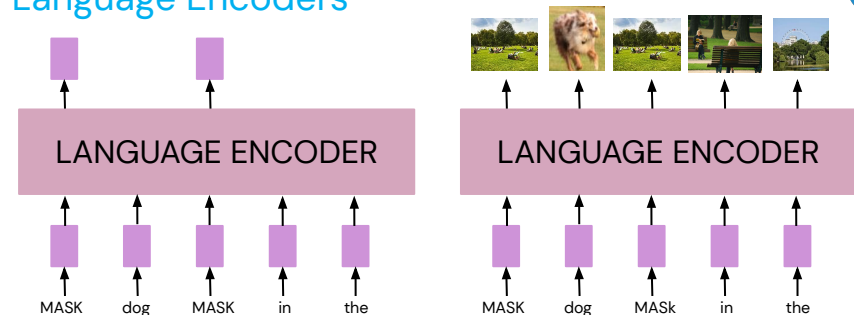
## Encoder-Decoders



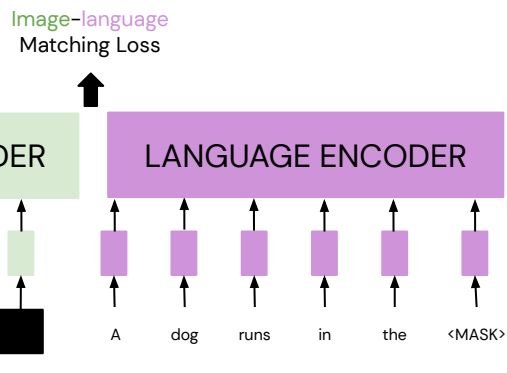
## Joint Encoders



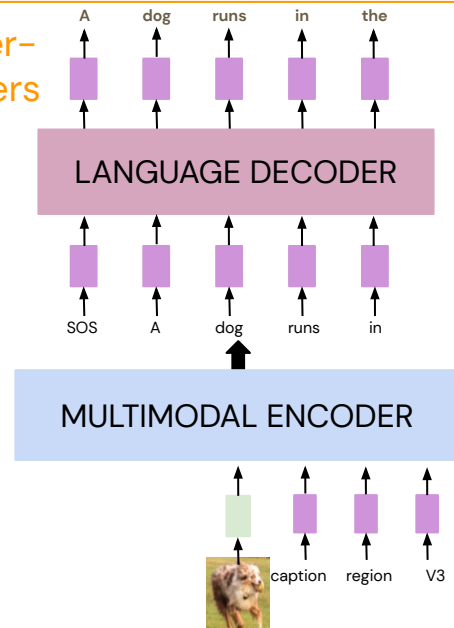
## Language Encoders



## Dual Encoders



## Encoder-Decoders

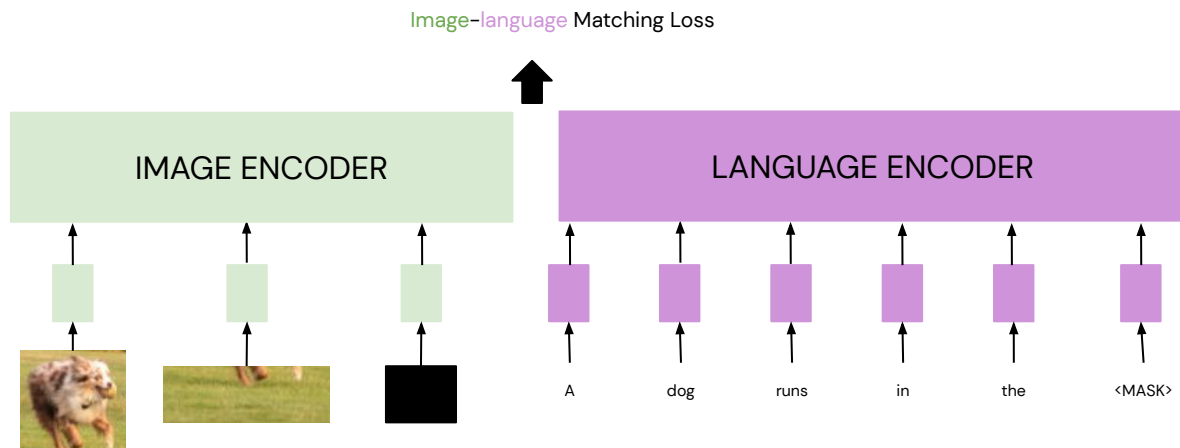




# Dual Encoders

Two separate encoders for image and language modalities; no cross-talk between the two. [Weston et al., 2011; Frome et al., 2013; Kiros et al., 2014]

Very successful for retrieval tasks [Chowdhury et al., 2018; Miech, Alayrac, et al.2020]





# Recent Large-Scale Dual Encoders [Radford et al, 2021; Jia et al, 2021]

CLIP [Radford et al, 2021] and ALIGN [Jia et al, 2021]: Larger **models** & **datasets**

How to collect large-scale datasets?

# Pretraining Datasets: Refresher

## MSCOCO



"The two people are walking down the beach."

## MSCOCO/OI Narratives



"In this image we can see a bridge and sea. In the background, we can see trees and the sky. We can see so many people on the bridge. At the bottom of the image, we can see two people. We can see stairs in the right bottom of the image ..."

## Visual Genome



small round yellow frisbee, man has cast on his arm, concrete trail path in the park, man wearing black sunglasses

manually annotated

## Conceptual Captions



"The *scenic route* through mountain ranges includes these unbelievably coloured mountains."

## SBU Captions



"*King Arthur's* beheading rock - right on the sidewalk in the middle of *town*".

from "the wild"

Noisier image-text correspondence but larger



# Recent Large-Scale Dual Encoders [Radford et al, 2021; Jia et al, 2021]

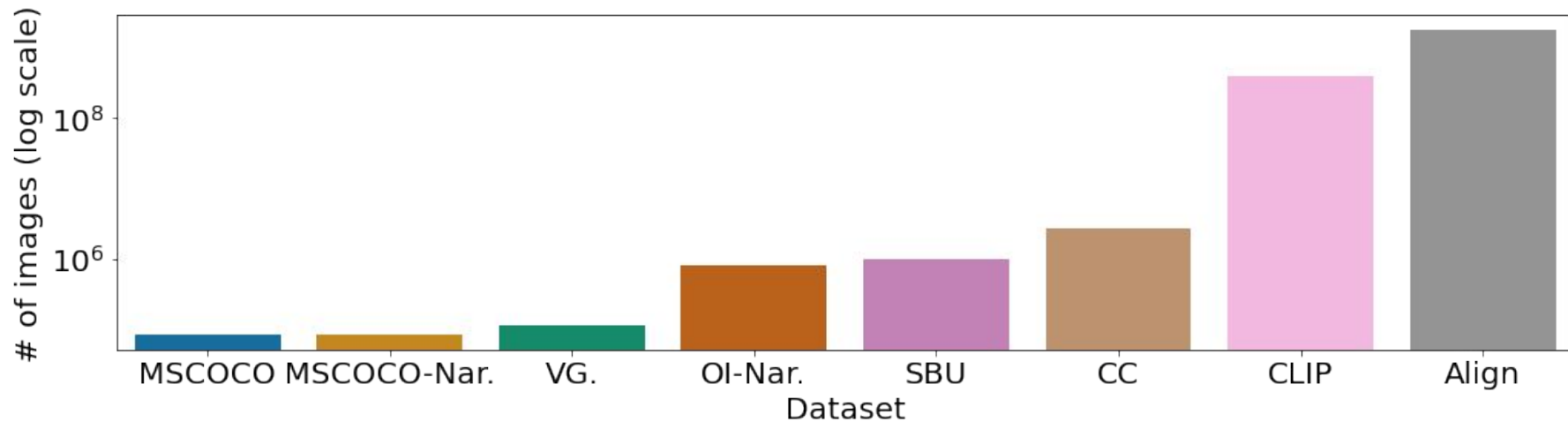
CLIP [Radford et al, 2021] and ALIGN [Jia et al, 2021]: Larger **models** & **datasets**

How to collect large-scale datasets?

- ALIGN removes any filtering to increase the size (1.8B) → noisier.
  - The same pipeline as Conceptual Captions (CC).
- CLIP uses heuristics to clean the data (400M).

Tradeoff between data size & noise: CC (3M) > ALIGN (3M/6M) on MSCOCO retrieval. [Jia et al, 2021]

# Dataset Considerations: Size





# Recent Large-Scale Dual Encoders [Radford et al, 2021; Jia et al, 2021]

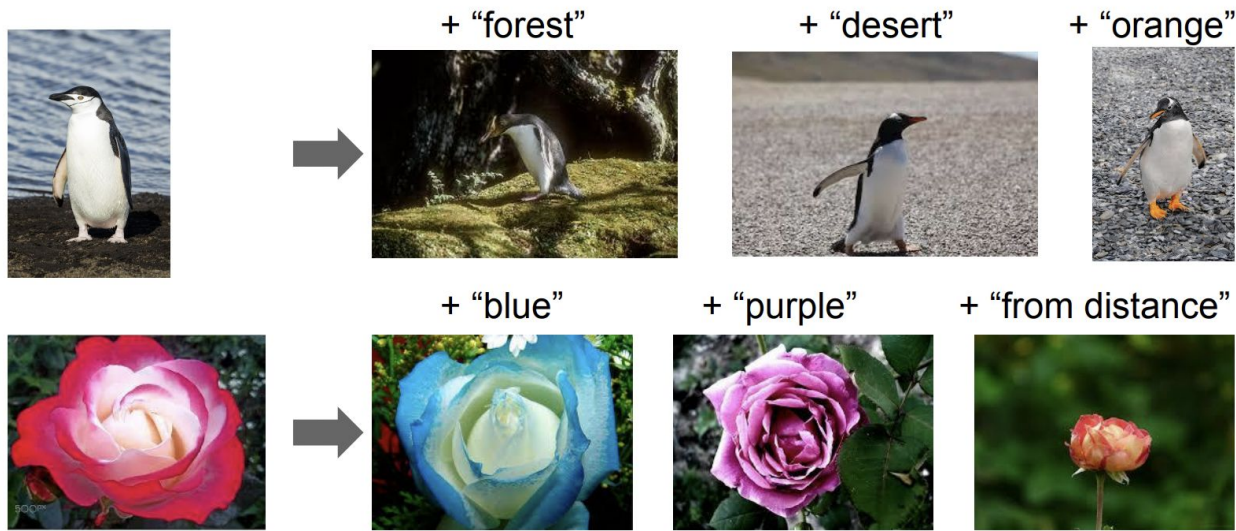
CLIP [Radford et al, 2021] and ALIGN [Jia et al, 2021]: Larger **models** & **datasets**

Use similar contrastive losses; ALIGN uses label smoothing that can be helpful with dataset noise.

Perform zero-shot image classification as a image-text retrieval task.

# Qualitative Examples from ALIGN

Image retrieval with image +/- text queries



# Qualitative Examples from ALIGN

Image retrieval with fine-grained queries.

“Lombard street ...”

“view from bottom”



“view from top”



“bird’s eye view”



“in heavy rain”

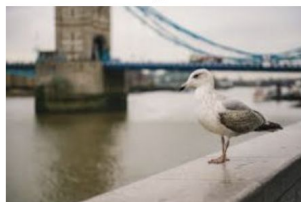


“seagull in front of ...”

“Golden Gate Bridge”



“London Tower Bridge”



“Sydney Harbour Bridge”



“Rialto Bridge”

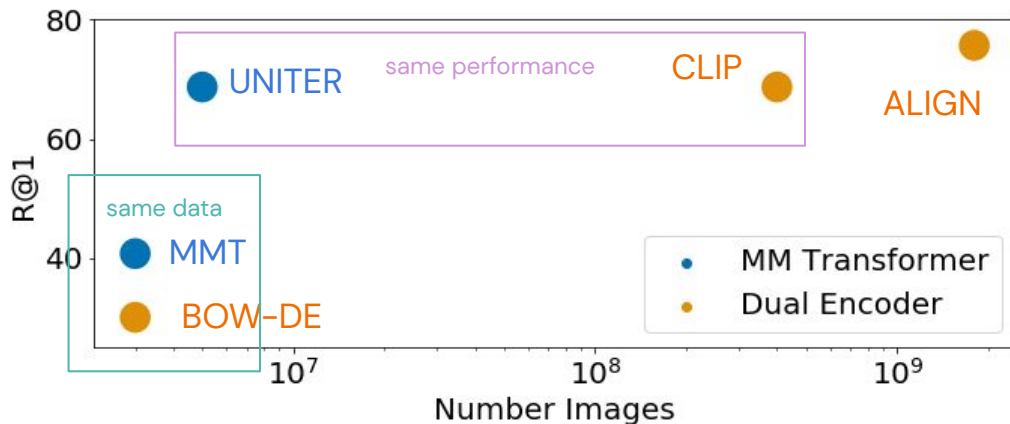




# Dual Encoders Vs. Multimodal Transformers

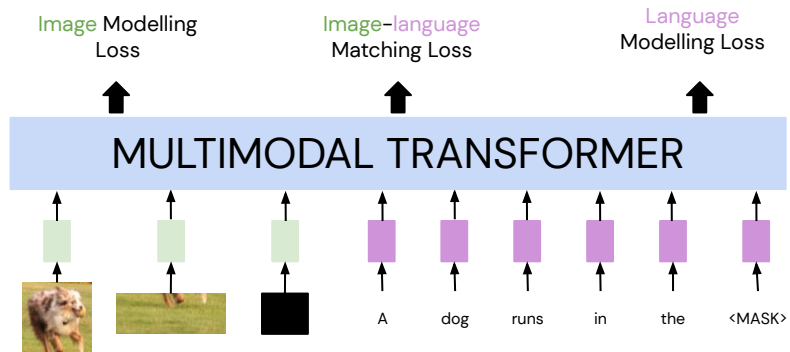
Dual encoders are easier to scale since they can reuse image/language features across pairs.

But they are not sample efficient.

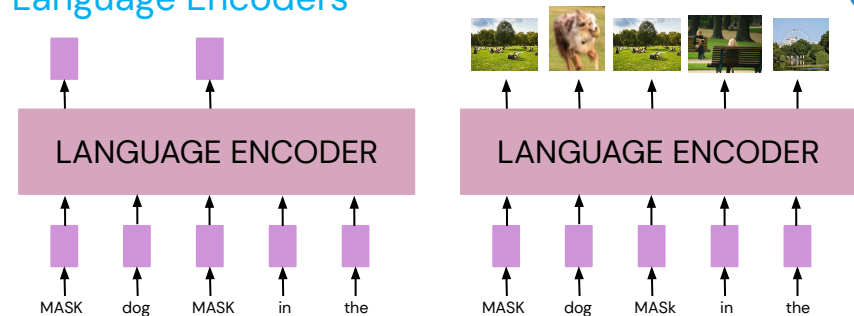


BOW-DE: [Miech & Alayrac et al. CVPR 2021]  
 MMT: [Hendricks et al. TACL 2021]  
 UNITER: [Chen et al. ECCV 2020]  
 CLIP: [Radford et al. Arxiv 2021]  
 ALIGN: [Jia et al. Arxiv 2021]

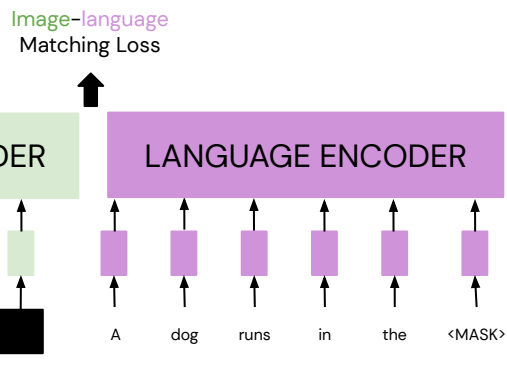
## Joint Encoders



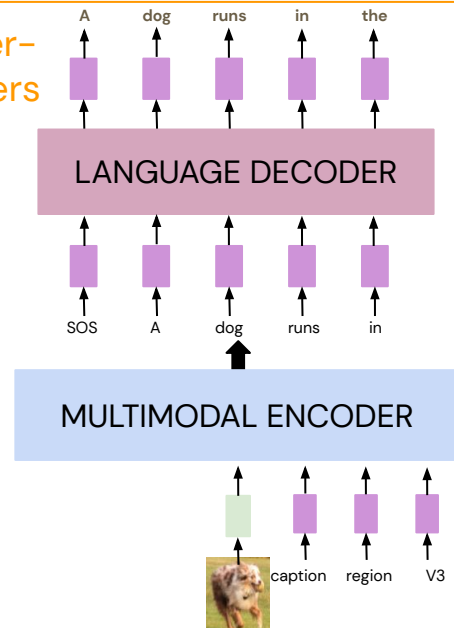
## Language Encoders



## Dual Encoders



## Encoder-Decoders



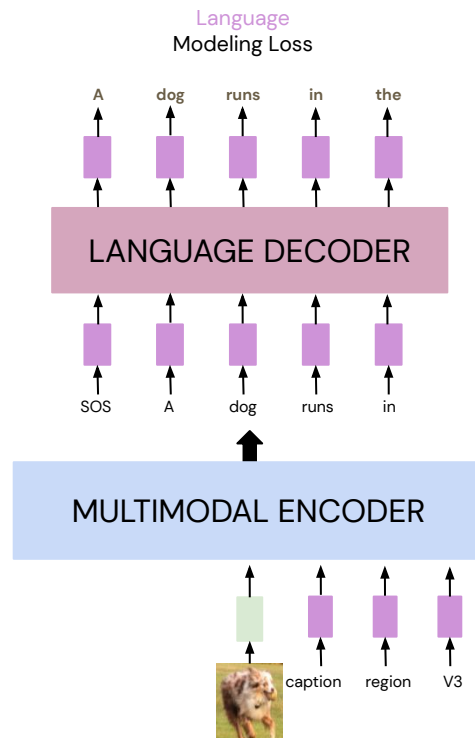
# Encoder-Decoders

An image captioning setup:

- Replace the image encoder with a multimodal one
- Virtex, VL-BART(T5), SimVLM [Desai & Johnson, 2020; Cho *et al*, 2021; Wang *et al*, 2022]

Uses **language** as supervision for **vision** or multimodal pretraining.

Requires less images than the image-classification setting.

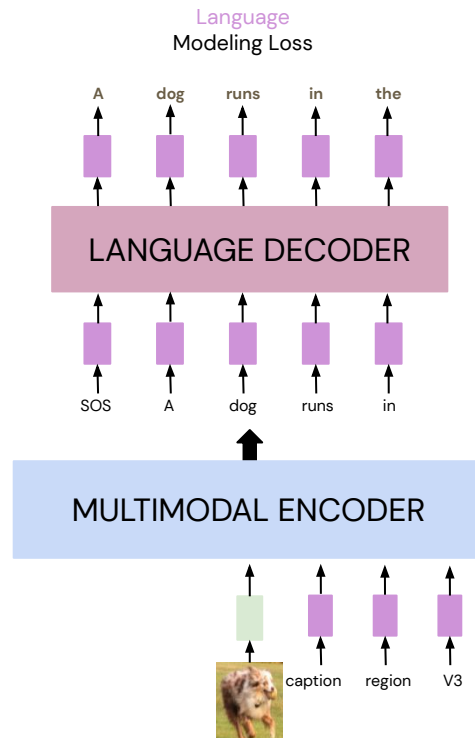


# Encoder-Decoders

An image captioning setup:

- Replace the image encoder with a multimodal one
- Virtex, VL-BART(T5), SimVLM [Desai & Johnson, 2020; Cho *et al*, 2021; Wang *et al*, 2022]

Uses **language** as supervision for **vision** or multimodal pretraining.



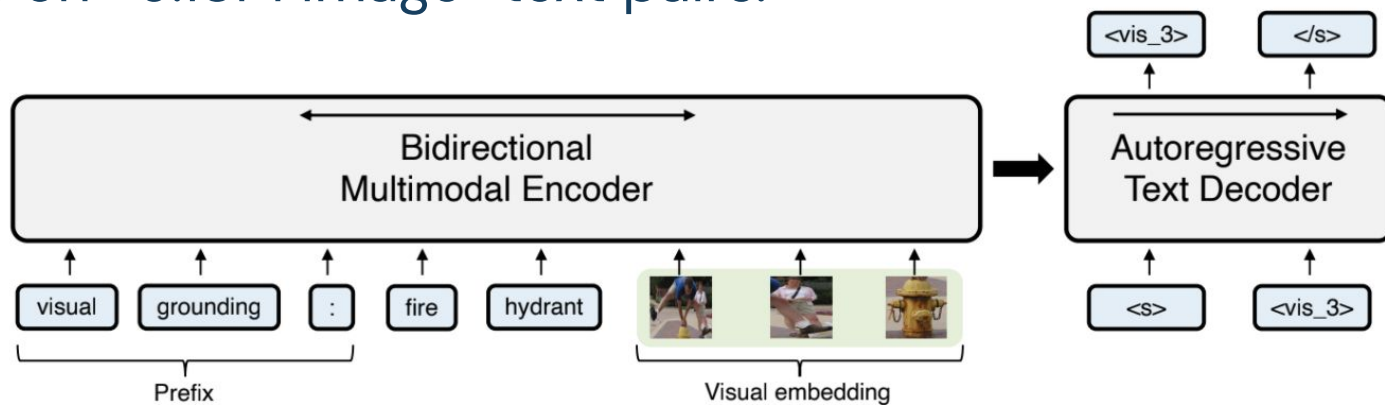
# VL-BART (VL-T5)<sup>[Cho et al, 2021]</sup>

Unifies tasks as text generation (w/ task-specific prefixes).

- Parameters for each task are separately optimized.

Builds on a pretrained language model (BART or T5).

Trains on ~9.18M image–text pairs.

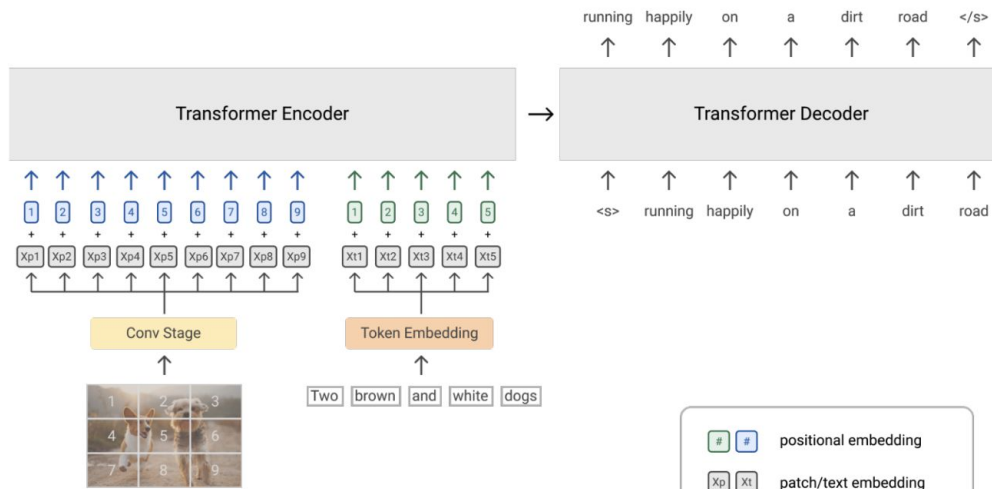


# SimVLM [Wang et al, 2022]

Unifies tasks as text generation.

Removes object detection supervision.

Trains on large-scale noisy image-text data (ALIGN).





# Combining Frozen (Pretrained) Models

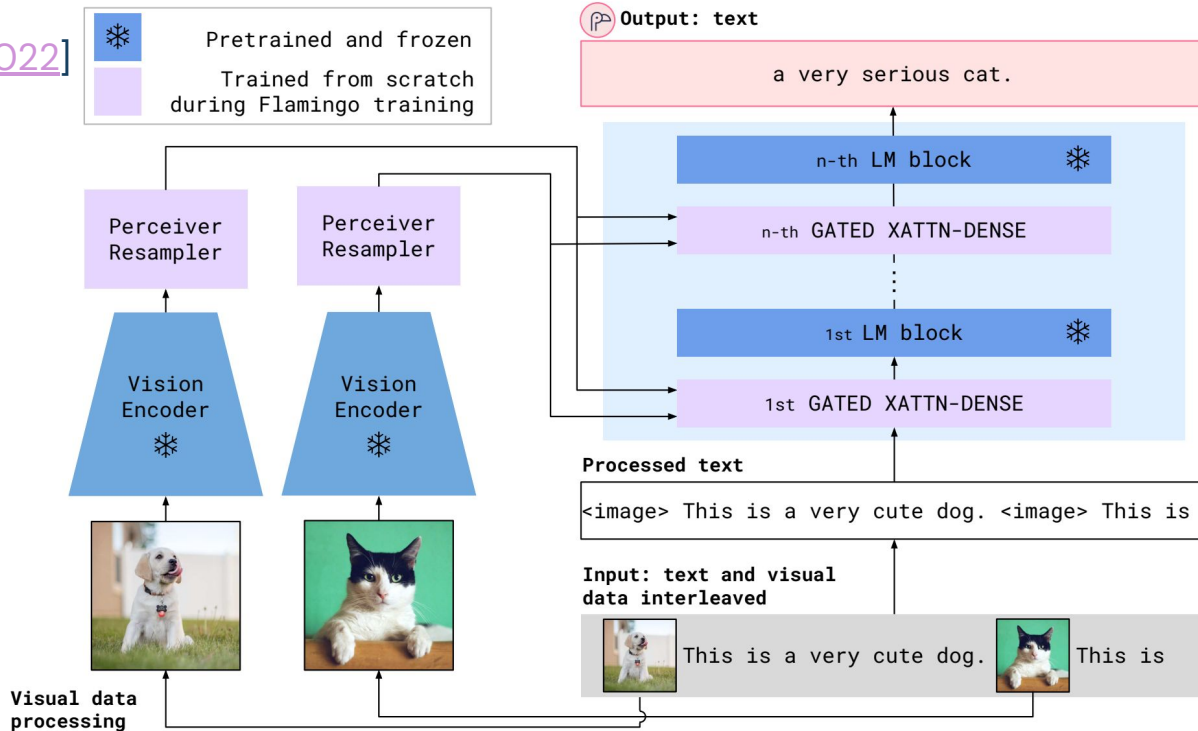
Given the cost of pretraining large models, can we reuse and combine existing vision and/or language models?

- Frozen [[Tsimpoukelli et al, 2021](#)]
- MAGMA [[Eichenberg et al, 2021](#)]
- Flamingo [[Alayrac et al, 2022](#)]

But, need to map different feature spaces (e.g., image features to the language ones) → Train adapter layers

# Combining Frozen (Pretrained) Models

Flamingo [Alayrac et al, 2022]







# Moving Towards Generative Models

## Pros.




- Unified modeling of vision–language tasks
- Better out-of-distribution generalization in VQA [Agrawal et al, 2022]

## Cons.

- How to accurately evaluate the generated text? [Agrawal et al, 2022]
- Inheriting existing biases of pretrained language (and vision) models.

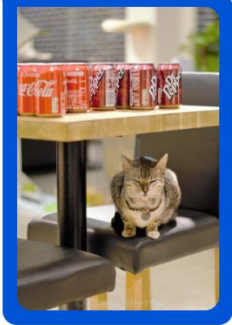
# Qualitative Examples from Flamingo

Language model bias can result in hallucinations.

Input Prompt	 Question: What is on the phone screen? Answer:	 Question: What can you see out the window? Answer:	 Question: Whom is the person texting? Answer:
Output	A text message from a friend.	A parking lot.	The driver.

# Qualitative Examples from Flamingo

Mistakes in spatial understanding but correct language use.




What is below the table?

It looks like a cat.

What is on top of the table?

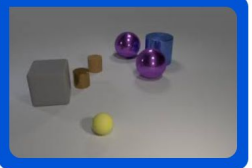
It looks like a bunch of soda cans.



This is a picture of three objects. The first one is a red ball, the second one is a blue cube, and the third one is a yellow cube.

well done! describe their arrangement with respect to one another

The red ball is on top of the blue cube, and the yellow cube is on top of the blue cube.



What shapes are present in the image?

This is a picture of a cube, a sphere, a cylinder and a cone.

What is behind the first purple sphere?

It's a cube.

What color is it?

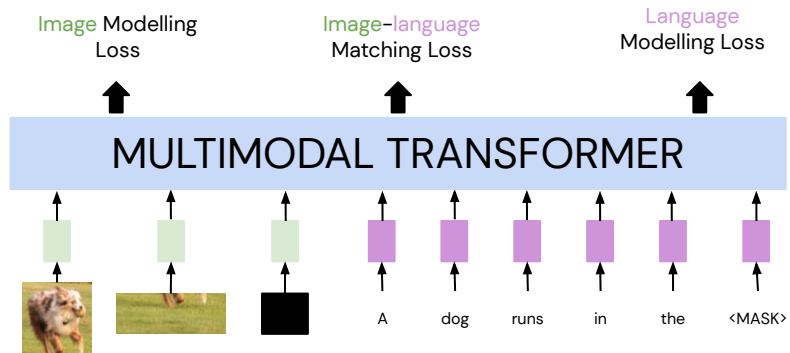
It's a light blue color.

What is the closest object?

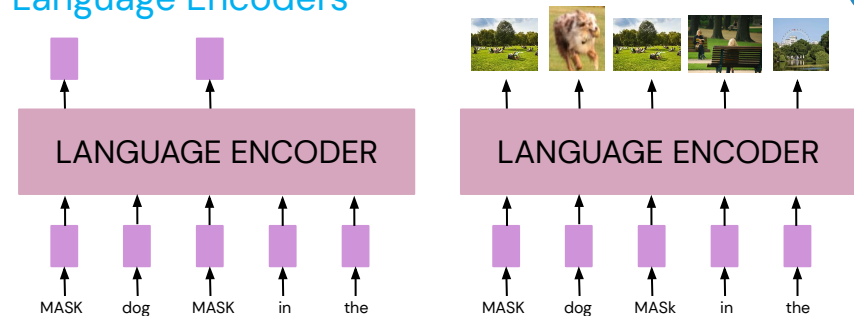
The closest object is the first purple sphere.

[Examples from JB Alayrac]

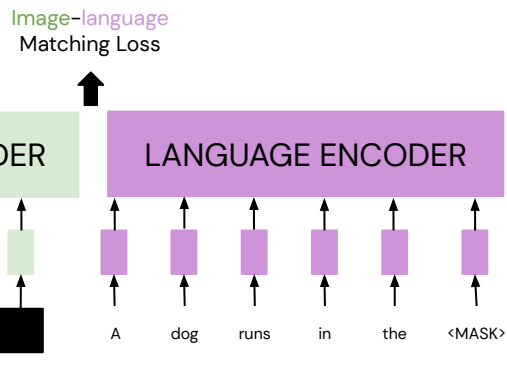
## Joint Encoders



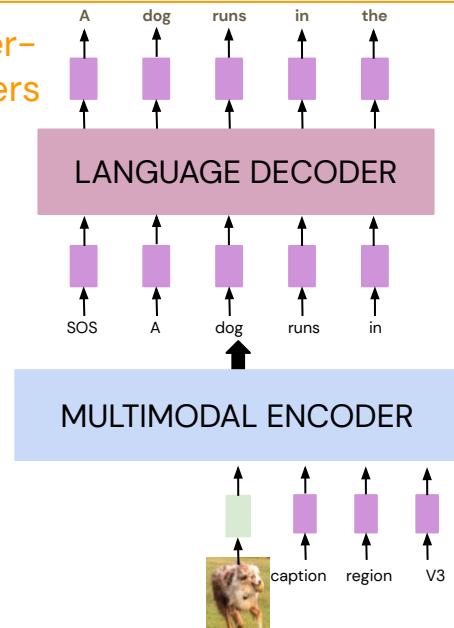
## Language Encoders



## Dual Encoders



## Encoder-Decoders

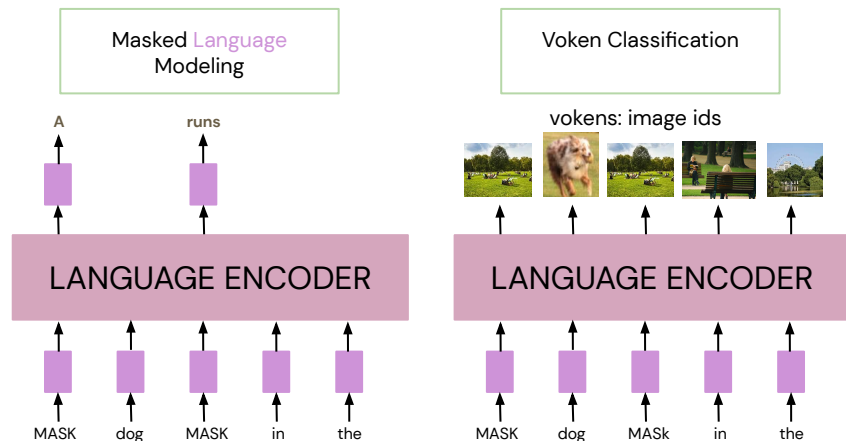
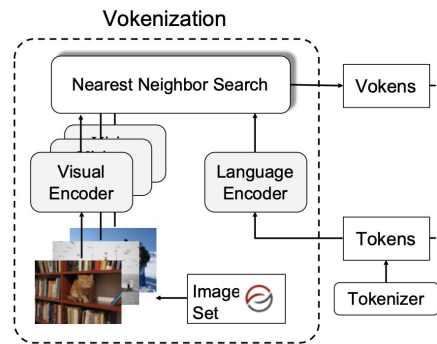


# Language Encoders

A language modeling setup:

- Vokenization: map each language token to a visual token (voken) [Tan & Bansal, 2020]

Uses **vision** as supervision for **language** pretraining.





# Summary of Different Approaches

How to evaluate pretrained models?

- Use task-specific heads for each downstream task (e.g., ViLBERT, LXMERT, UNITER, OSCAR, VinVL).
- Treat all downstream tasks as language generation with no task-specific head (e.g., VL-T5, VL-BART, SimVLM).



# Summary of Different Approaches

How are the features used (other than vision-language tasks)?

- In vision tasks (e.g., [VirTex](#), [CLIP](#), [ALIGN](#))
- In language tasks, including multilingual data (e.g., Vokenization, M3P, VL-T5, SimVLM)



# Towards a Better Evaluation of Pretrained Models

Performance after fine-tuning depends on the the size of fine-tuning data and other experimental set-up [Yogatama et al., 2019].

Recent work has shifted focus to few- and zero-shot evaluation.

Other approaches

- Evaluate for out-of-distribution generalization (transfer)
- Probe for certain capabilities (e.g., verb understanding)

See [my talk on evaluation](#) if you are interested!





# Answering Questions from Blind People



Q: What are the people waiting for?

A: bus



Q: What is this?

A: 10 euros.

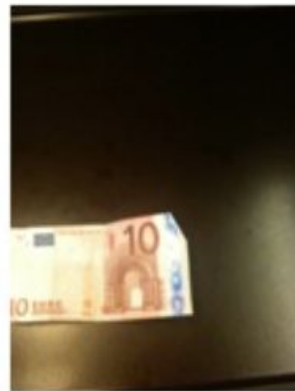
VizWiz is a benchmark curated from visually-impaired users.

# Evaluate in a Transfer Setting [Agrawal et al, 2022]

Fine-tune a multimodal transformer on one dataset (VQAv2),  
test on another one (VizWiz): we observe **~26** drop in accuracy.



Q: What are the people waiting for?  
A: bus



Q: What is this?  
A: 10 euros.

# Winoground: Visio-Linguistic Compositionality



(a) there is [a mug] in [some grass]



(c) a person [sits] and a dog [stands]



(e) it's a [truck] [fire]



(b) there is [some grass] in [a mug]



(d) a person [stands] and a dog [sits]



(f) it's a [fire] [truck]

[Thrush et al, 2022]

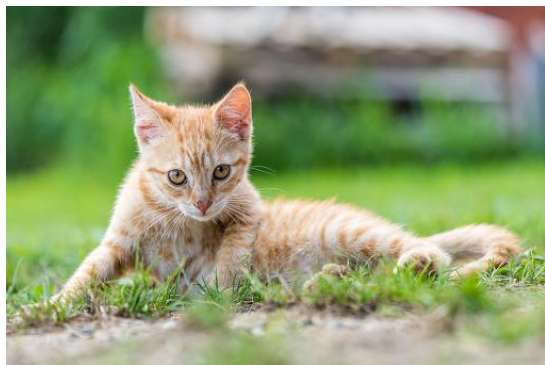
*Object*

*Relation*

*Both*

# What Image Retrieval Tests

Order images with respect to their match to a sentence.



A person is riding a horse.

Subject

Verb

Object

Does not require fine-grained multimodal understanding.



# What SVO-Probes Tests [Hendricks et al., Findings of ACL 2021]

A person is **riding** a horse

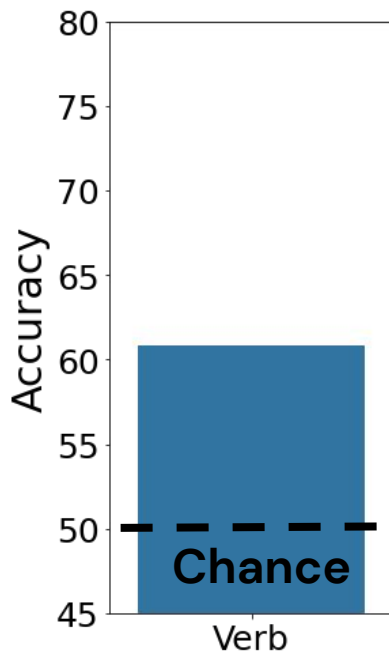


Correctly classify both the **positive** & **negative** examples.

We have released our dataset! 🎉🎉

# Do MMTs Have Fine-grained Verb Understanding?

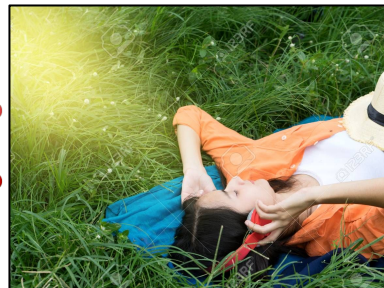
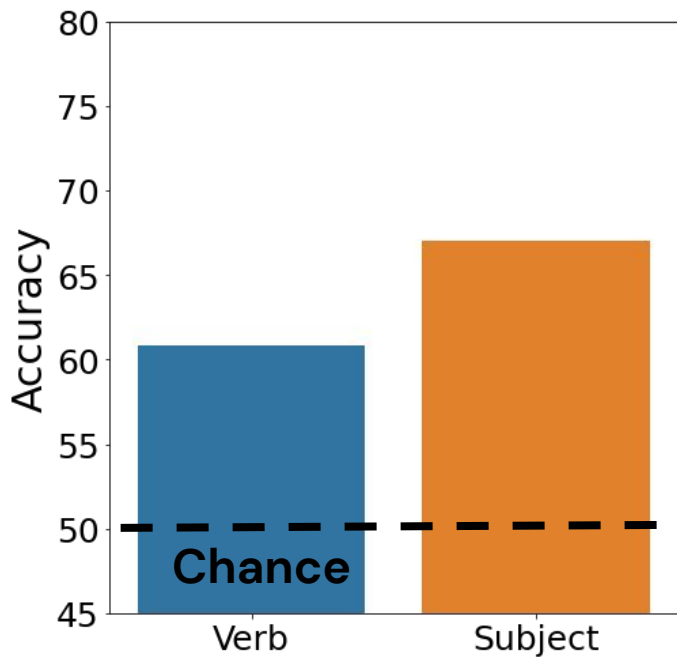
A woman **lying** with a dog





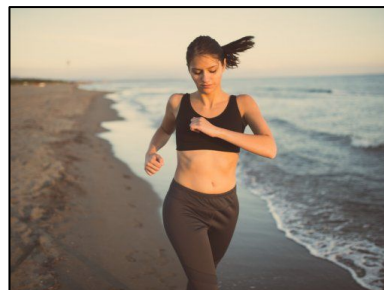
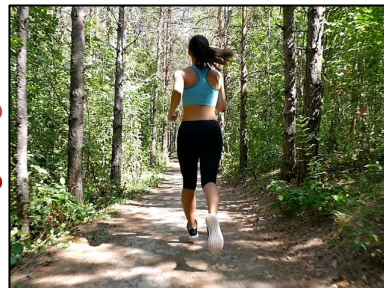
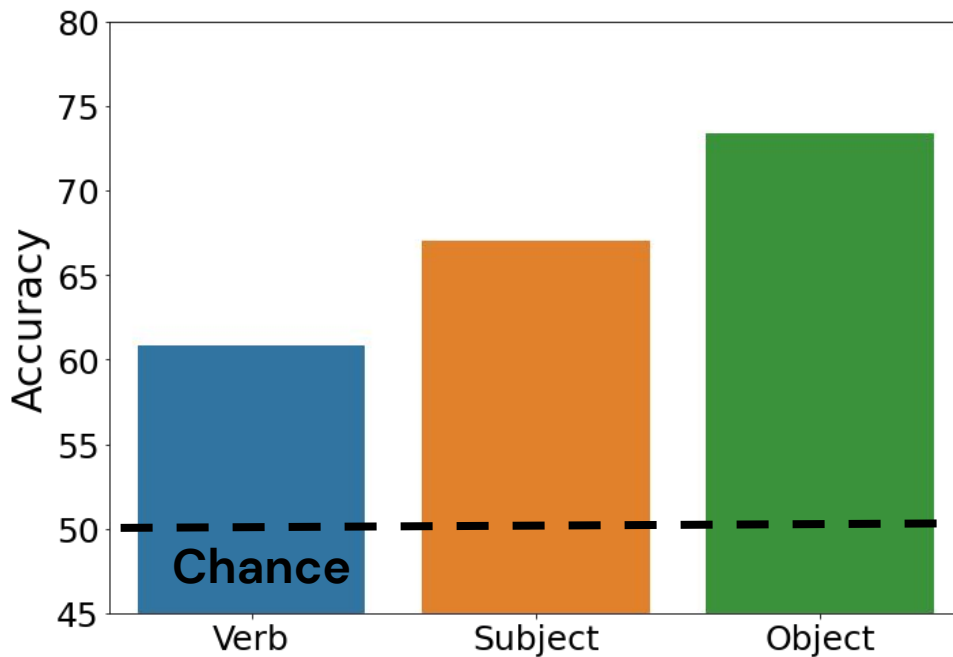
# Do MMTs Have Fine-grained Verb Understanding?

A **animal** lays in the grass



# Do MMTs Have Fine-grained Verb Understanding?

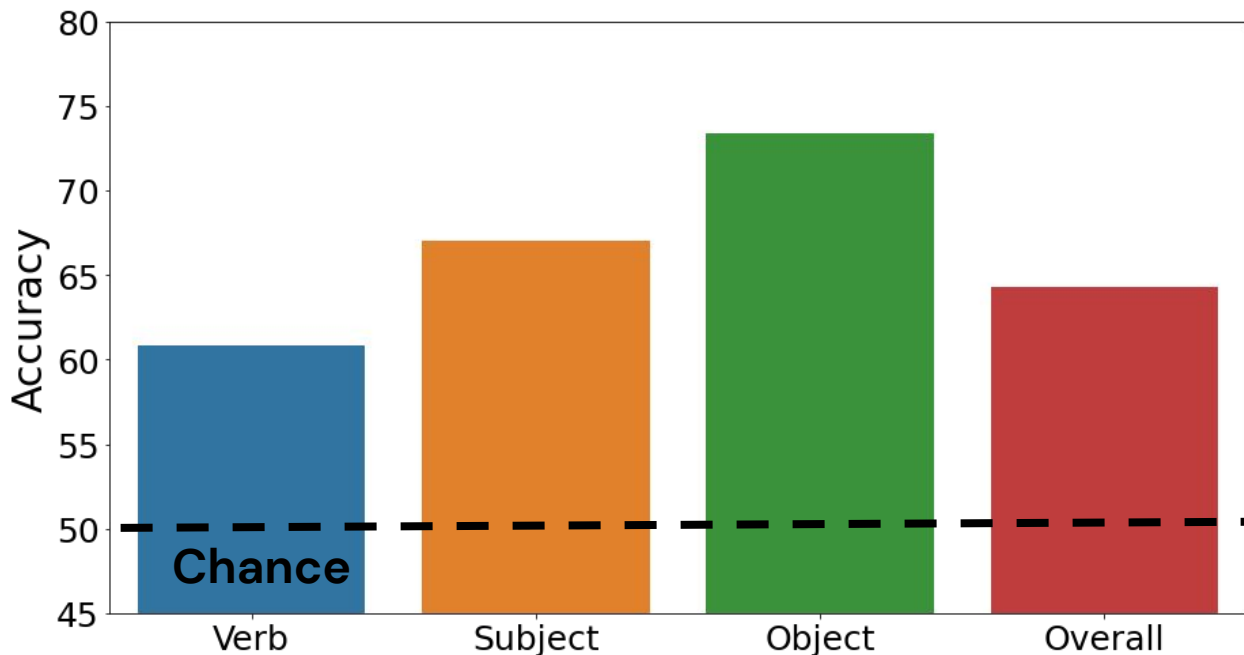
A woman jogs on the **beach**







# Do MMTs Have Fine-grained Verb Understanding?



Overall MMT  
performance 64.3 --  
lots of room for  
improvement!



**To build stronger models, we need to  
better evaluate them first.**

Thanks for listening!

**Questions?**