

# Multilingual and multimodal language models



Desmond Elliott

Language and Multimodal Processing Group  
University of Copenhagen

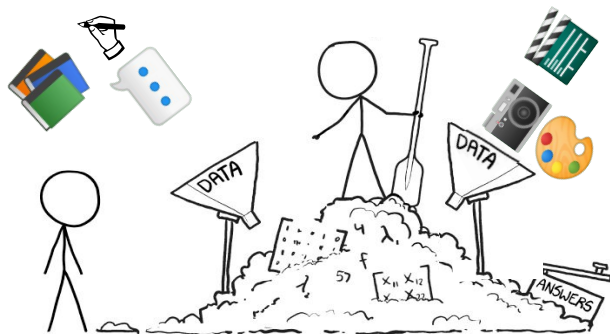
---

Slides: <https://elliottd.github.io/vlprimer/>

# Working Definition

---

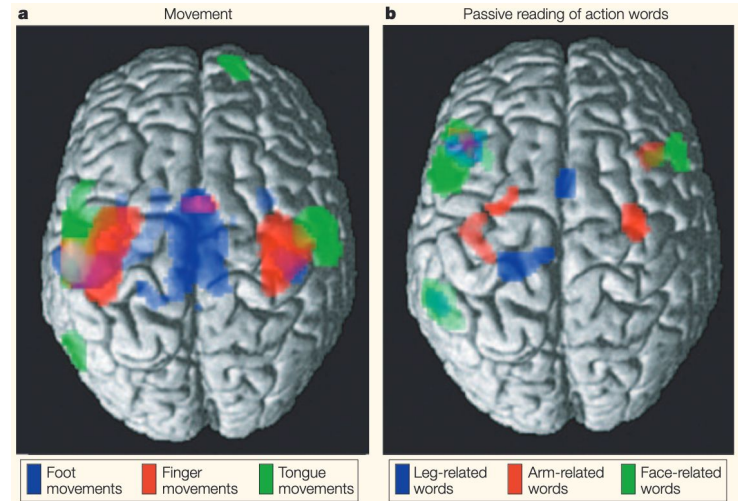
Multimodal models jointly processes information from two or more input modalities, e.g. images and text, speech and video, etc.



# Why Multimodality?

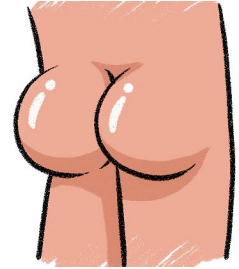
---

- Humans ground conceptual knowledge in modality processing systems in the brain
- Evidence that grounding activates similar brain regions for different input modalities



# Multimodality reduces ambiguity

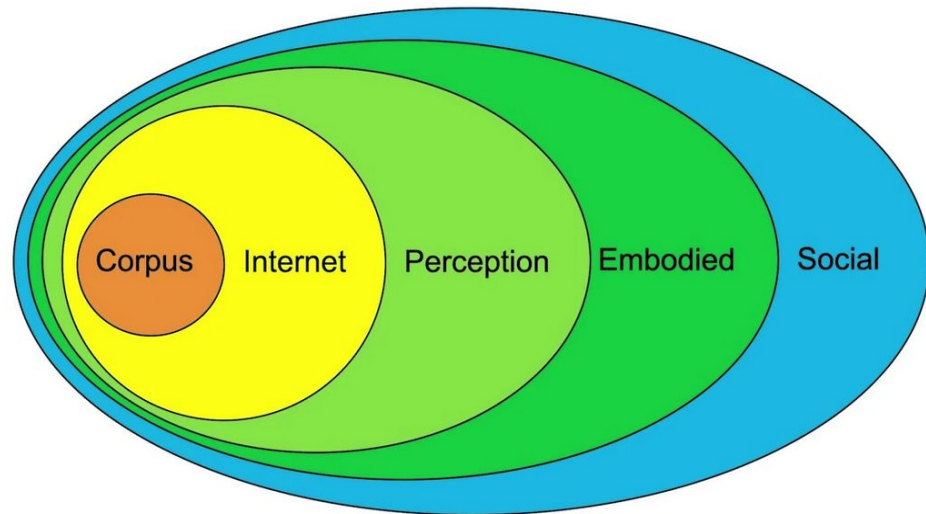
---



# You Cannot Learn Language From

---

- The radio without grounding  
*(lack perception)*
- The television without actions  
*(lack embodiment)*
- Without interacting with others  
*(lack social)*



# (At Least) Five Major Areas

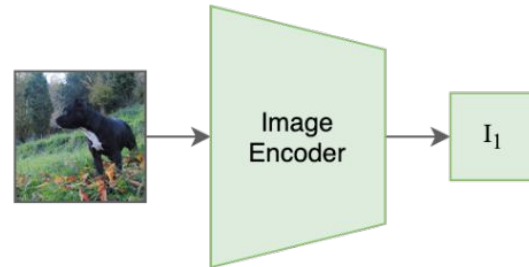
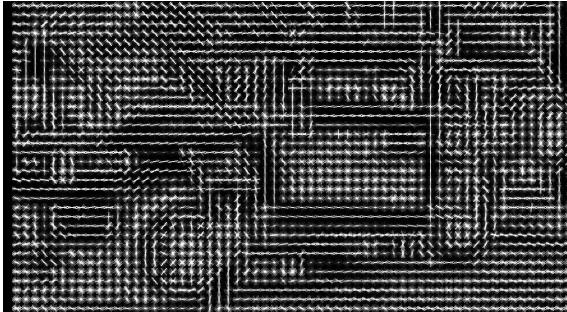
---

- **Representation:** how to convert raw inputs into a usable format
- **Translation:** transform from one modality to another
- **Alignment:** predict relationships between elements across modalities
- **Fusion:** join features from modalities to support prediction
- **Co-learning:** transferring knowledge from one modality to another

# Representation

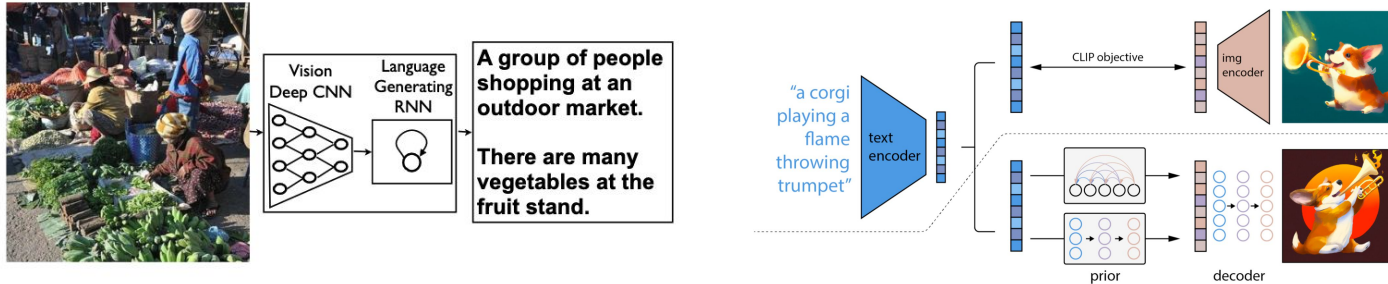
---

- Great deal of work over the last decade, from HOG features in the early 2000s to CLIP features in the 2020s.



# Translation

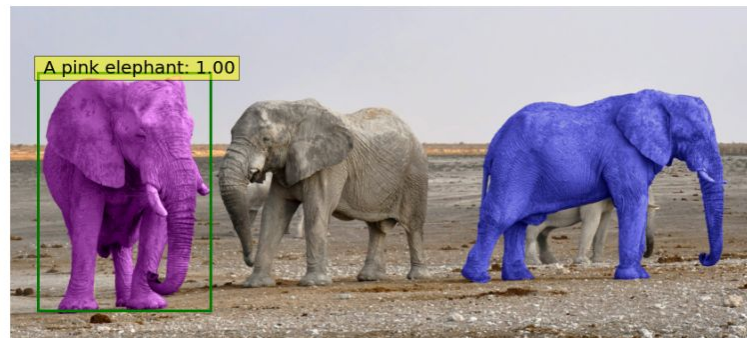
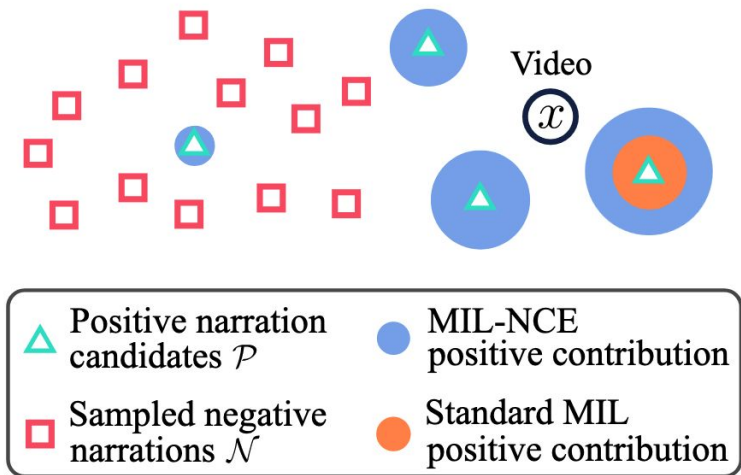
- Explosion of end-to-end neural network models since the mid 2010s





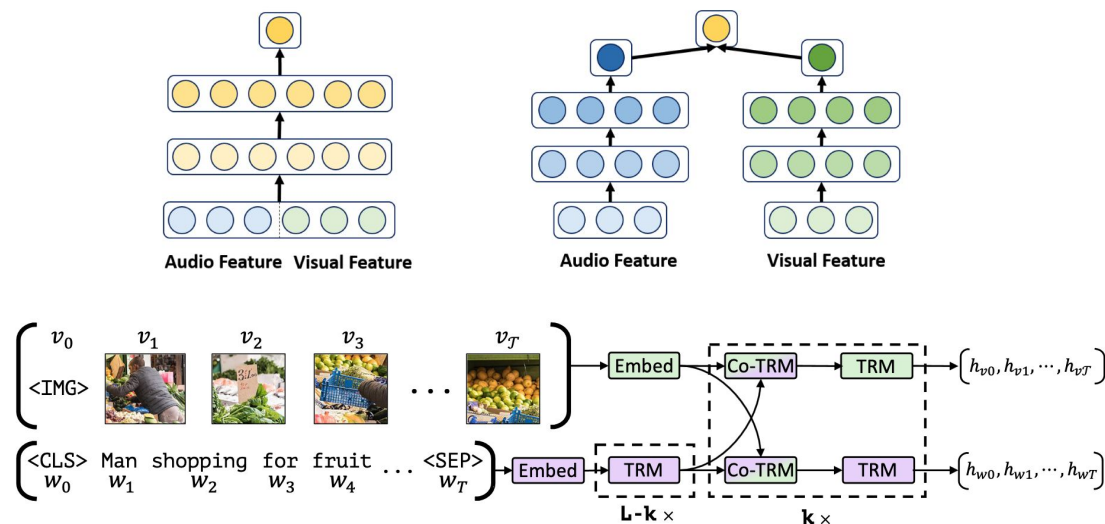
# Alignment

- Important for self-supervised learning and also for phrase grounding



# Fusion

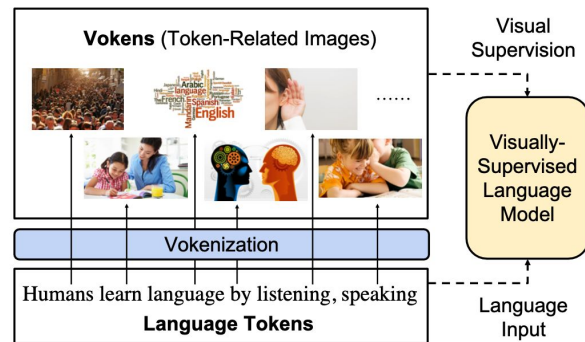
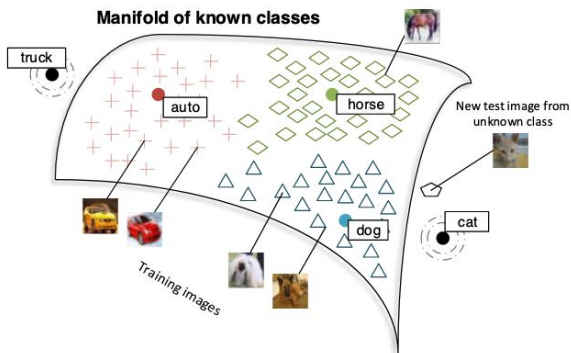
- Early work studied the differences between early and late fusion.
- Multi-head self-attention now provides model-based fusion.



Chen and Jin (2016). Multi-modal conditional attention fusion for dimensional emotion prediction. MM.  
Lu et al. (2019). ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *NeurIPS*.

# Co-learning

- Zero-shot transfer across modalities, or using visual grounding to improve language models on text-only tasks.



Socher et al. (2013). Zero-shot learning through cross-modal transfer. NeurIPS.


Tan & Bansal. (2020). Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision. EMNLP

# Roadmap

---

## Part 1

### 1. Datasets for Multimodal Learning

 Visually Grounded Reasoning across Languages and Cultures

### 2. Data Representation

### 3. Modelling Techniques

 Retrieval-Augmentation in Image Captioning

## Part 2

### 4. Understanding Multimodal Models

### 5. Future Directions

 Language Modelling with Pixels

# 1. Datasets for Multimodal Learning

# Two Types of Dataset

---

- **General-purpose:** visual data with descriptive annotations

- Conceptual Captions
- LAION-2/5B
- Speech-COCO



Blue Beach  
Umbrellas, Point  
Of Rocks, Crescent  
Beach, Siesta Key -  
Spiral Notebook

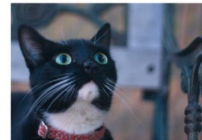
- **Task-specific:** visual data with e.g. classification labels

- Image / Video Captioning
- Visual Question Answering
- Visually Grounded Reasoning

What color is the cat's leash?

purple

red



# Degree of Multimodality

---

Social media platforms often form 'echo chambers' that encourage users to only read content that confirms beliefs they already hold (Getty)

Weak



A woman in a dark grey suit is giving a speech

Strong

Panofsky. (1939). Studies in Iconology.

<https://www.independent.co.uk/news/angela-merkel-says-internet-search-engines-endangering-debate-algorithms-should-be-revealed-a7383811.html>

# Conceptual Captions

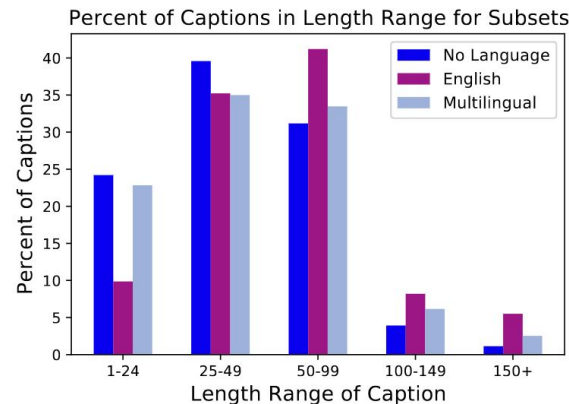
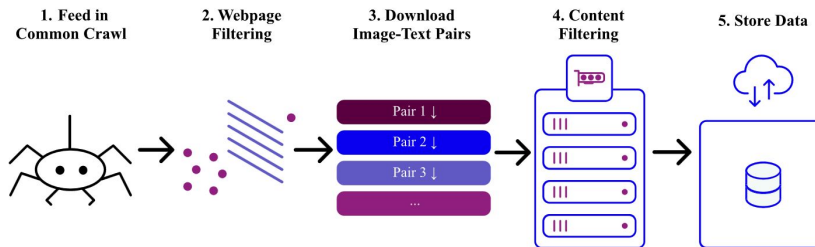
- Used for pretraining
- 3M Images and *normalized* English captions.
- Normalization is not public.
- Due to *linkrot*, much less data is currently available.





# LAION

- Used for pretraining
- Image and multilingual *raw* captions harvested from within Common Crawl
- Data behind Stable Diffusion and OpenCLIP
- 5B variant removed due to illegal material



# COCO

---

- Used both a **general-purpose** and **task-specific** dataset
- Images covering 80 common objects in context with multiple human-authored captions.
- Object segmentation data too!

some sheep walking in the middle of a road  
a herd of sheep with green markings walking down the road  
a herd of sheep walking down a street next to a lush green grass covered hillside.  
sheared sheep on roadway taken from vehicle, with green hillside in background.  
a flock of freshly sheered sheep in the road.



# VQA v2

---

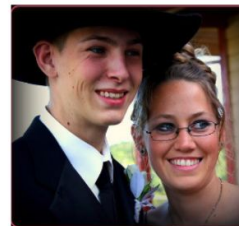
- Answer questions about images
- Task with multimodal inputs:
  - Image
  - Question
- Commonly tackled as classification but increasing efforts as NLG
- 1.1M image–question pairs with balanced distribution of answers

Who is wearing glasses?

man



woman

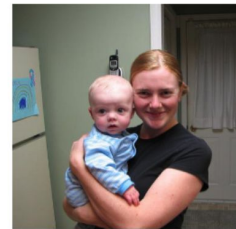


Where is the child sitting?

fridge



arms



# NLVR2

---

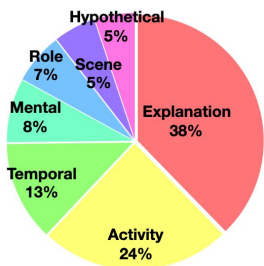
- Binary classification task that requires jointly reasoning over a pair of images and a sentence.
- Human-created hard negatives.
- 107K examples in total.



*The left image contains twice the number of dogs as the right image, and at least two dogs in total are standing.*

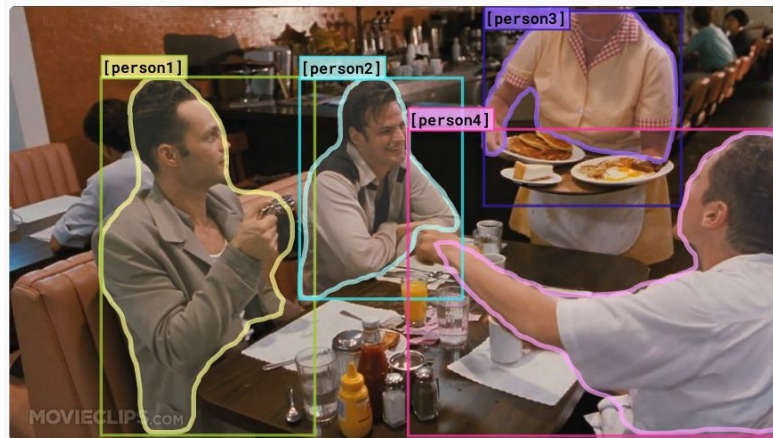
# Visual Commonsense Reasoning

- 290,000 multiple-choice VQA examples derived from movies.



- Why is [person11] wearing sunglasses inside?
- What are [person1] and [person2] doing?
- What will [person6] do after unpacking the groceries?
- What is [person3] thinking while [person5] shakes his hand?
- What is [person1]'s relation to [person4]?
- Where is [person1] now?
- What would happen if [person3] fell asleep?

- In addition to Question Answering, the dataset includes rationale selection too!



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.

Rationale:

- 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.

# Multi30K

---

- Multilingual aligned image–sentence dataset in many languages
  - English, German, French, Czech, Arabic, Japanese, Turkish, Ukranian

*A group of people are eating noodles.*

*Eine Gruppe von Leuten isst Nudeln.*

*Un groupe de gens mangent des nouilles.*

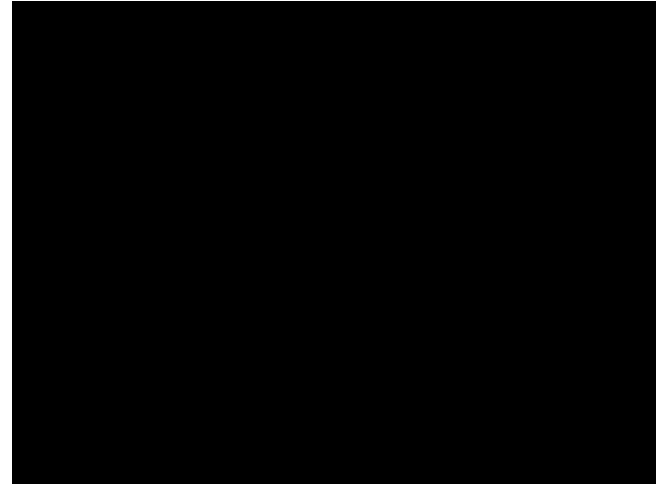
*Skupina lidí jedí nudle.*



# BOBSL

---

- BBC-Oxford British Sign Language Dataset
- Sign spotting and sentence localization tasks
- 1,400 hours of signed shows
  - Factual, entertainment, drama, comedy, children's shows



# Many Many More

---

- Visual Storytelling, e.g. VIST
- Grounded Referring Expression, e.g. Flickr30K Entities, Visual Genome
- Visual Entailment, e.g. SNLI-VE
- Vision & Language Navigation, e.g. RxR
- Visual Common Sense Reasoning: VCR
- Text-to-Image Generation, e.g. DALLEval
- Abstract reasoning, e.g. KiloGram, CRAFT
- Sign Language Processing, e.g. How2Sign
- *and more and more and more and more*



# Ethical Issues

---

- Multimodal datasets are usually data scraped from the web with *unknown degrees of conformance*, or information about, licensing.



**CC BY:** This license allows reusers to distribute, remix, adapt, and build upon the material in any medium or format, so long as attribution is given to the creator. The license allows for commercial use.

- As of 2022, there are an estimated 2.5B CC-licensed objects online.

# A Problem with Scale

- Build multimodal systems that perpetuate harmful stereotypes



(Eileen Collins, American astronaut)

$\cos(\cdot, \cdot)$

0.276

← This is a portrait of an astronaut with the American flag

0.308

← This is a photograph of a smiling housewife in an orange jumpsuit with the American flag

**Q: How can we collect multimodal data that better reflects the diversity of the world?**

---

# Visually Grounded Reasoning across Languages and Cultures

EMNLP 2021



F. Liu\*



E. Bugliarello\*



E.M. Ponti



S. Reddy



N. Collier



D. Elliott

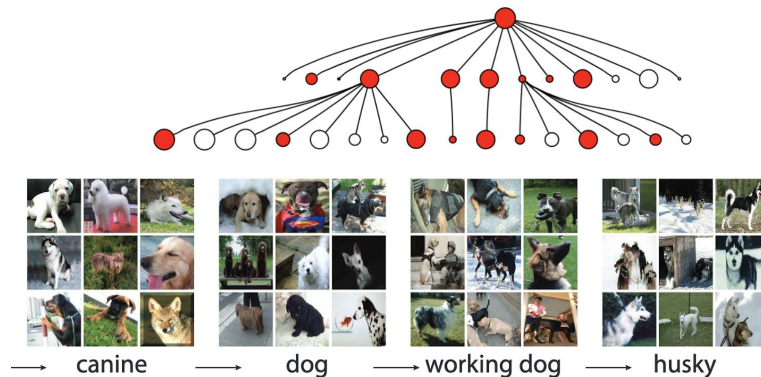
---

# Typical Vision and Language



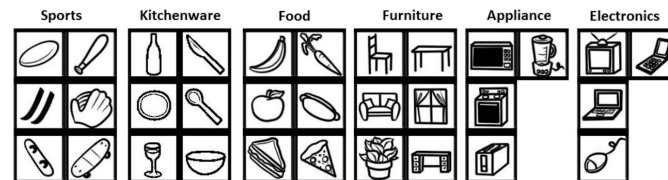
## ImageNet (Deng et al. 2009)

- Train visual encoders
- Millions of labelled images
- Derived from the WordNet hierarchy



## Common Objects in Context (Lin et al. 2014)

- Train and evaluate multimodal models
- 330K labelled images
  - 80 types of commonly occurring objects



# Rethinking Vision and Language



## Languages

- Mostly in English
- Or some Indo-European Languages



ENG: An **unusual** looking vehicle ...

NLD: Een mobiel **draaiorgel** ...

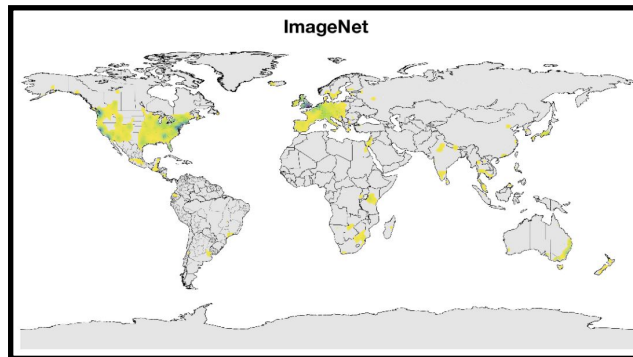
Example from [van Miltenburg+ 2017](#)

## Image sources

- Mostly from ImageNet or COCO
- Reflecting North American and European cultures

## Implications for V&L models

- Narrow linguistic/cultural domain
- No way to assess their real-world comprehension



Density map of geographical distribution of images in ImageNet ([DeVries+, 2019](#))

# Concepts and Hierarchies

**Category:** objects with similar properties (Aristotle 40 BCE, ...)

**Concept:** mental representation of a category (Rosch 1973)



"Dog" concept

Categories form a *hierarchy*

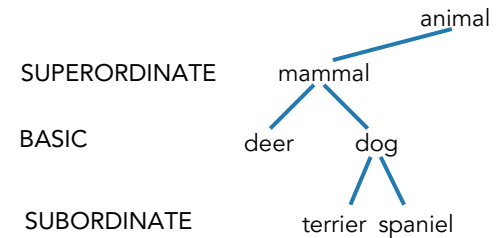
- Basic-level categories (Rosch 1976)



"Dog" category

Somewhat universal

- Different cultures (Berlin 2014)
- Familiarity of individuals (Wisniewski and Murphy, 1989)



# Concrete Concepts in Cultural Context

- Some concepts are most immediately understood within a cultural background

*Culture*: The way of life of a collective of people that distinguishes them from other people ([Mora, 2013](#); [Shweder et al. 2007](#)).



**Pilota / Jai-alai**



**Sanxian / Shamisen**

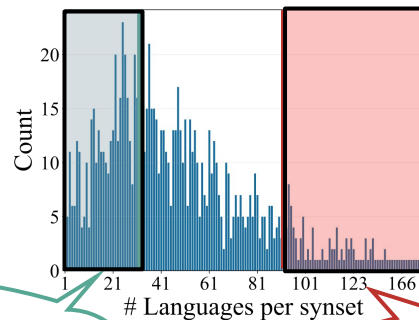
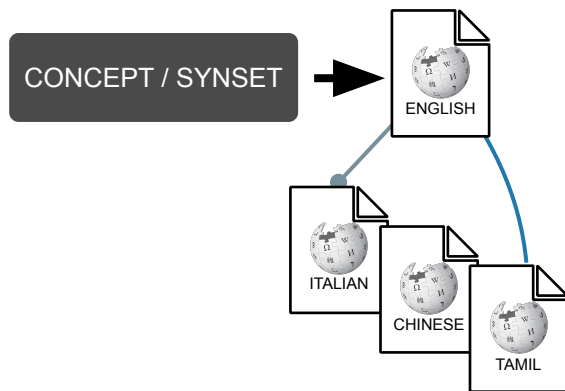


**Clavie**



# Are ImageNet Concepts Cross-Lingual?

- The ImageNet, COCO and Visual Genome datasets use English WordNet concepts
- Idea: estimate cross-linguality using Wikipedia as a proxy



Most concepts in  
<30 languages

Only a few "universal"  
concepts



Representative of annotators' cultures



5 typologically diverse languages

Independent, culture-specific annotations



MaRVL-id Bola basket



MaRVL-sw Mpira wa kikapu



MaRVL-tr Basketbol



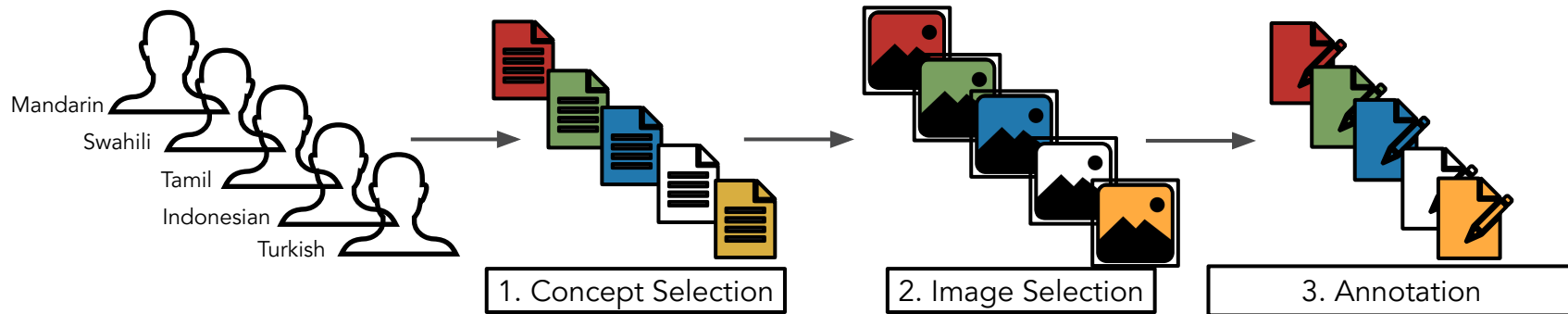
MaRVL-zh 篮球



MaRVL-ta கூடைப்பந்தாட்டம்

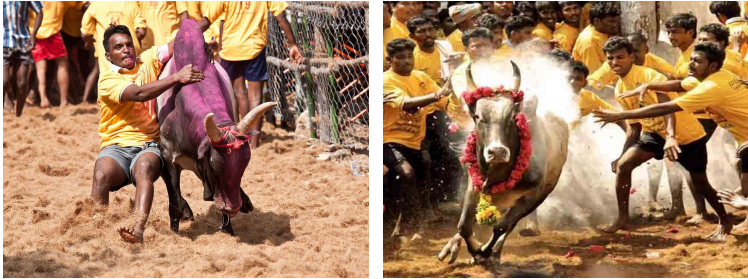
# Collecting **MaRVL** data

Native speaker-driven protocol



# Visual Reasoning Task (Suhr et al. ACL 2019)

- **Datapoint:** two images ( $v_1, v_2$ ) paired with a sentence  $x$
- **Task:** Predict whether  $x$  is a true description of the pair of images  $v_1 v_2$



இரு படங்களில் ஒன்றில்  
இரண்டிற்கும் மேற்பட்ட  
மஞ்சள் சட்டை அணிந்த  
வீரர்கள் காளையை அடக்கும்  
பணியில் ஈடுபட்டிருப்பதை  
காணமுடி.

True

$x$

$y$

# MaRVL is created from Universal Concepts

- Taken from the *Intercontinental Dictionary Series* (Key & Comrie, 2015)
  - 18/22 chapters with concrete objects & events

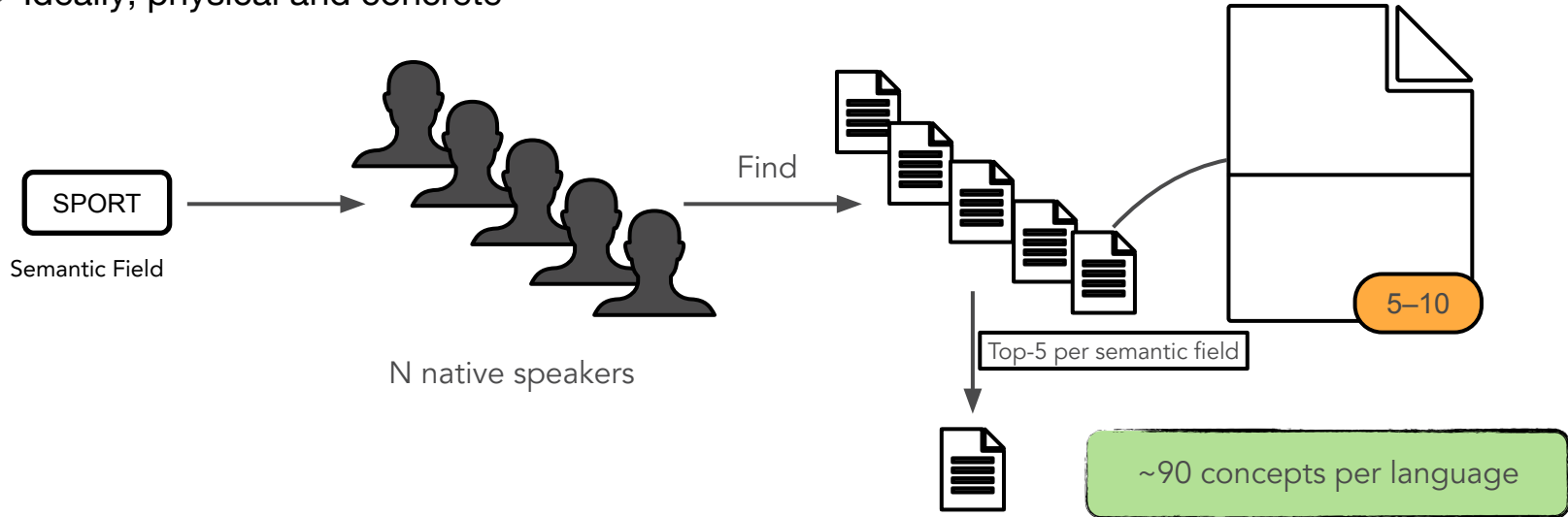
Chapter	Semantic Field
Animal	Bird, mammal
Food and Beverages	Food, Beverages
Clothing and grooming	Clothing
The house	Interior, exterior
Agriculture and vegetation	Flower, fruit, vegetable, agriculture
Basic actions and technology	Utensil/tool
Motion	Sport
Time	Celebrations
Cognition	Education
Speech and language	Music (instruments), visual arts
Religion and belief	Religion



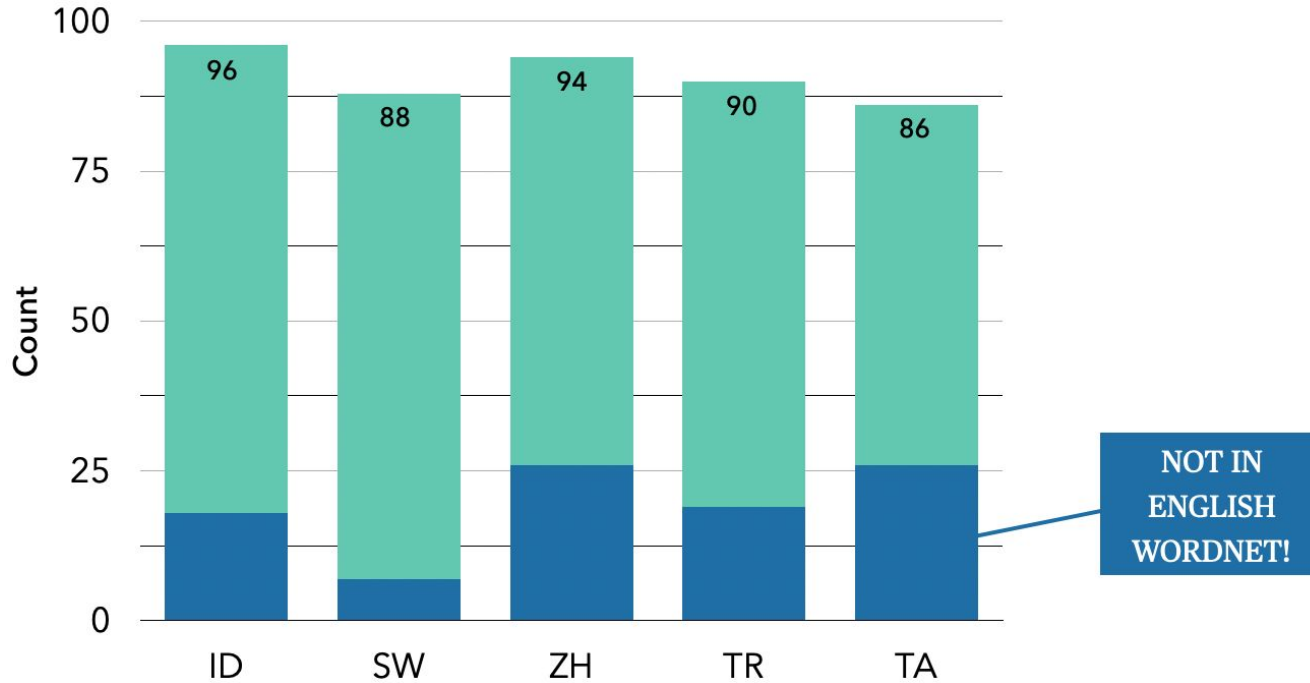
# Step 1. Language-Specific Concept Selection

Defined by native speakers

- Commonly seen or representative in their culture
- Ideally, physical and concrete



# Overview of Resulting Concepts



# Step 2. Image Collection

Collected by native speakers

- Representative of the language population
- NLVR2 (Suhr et al. ACL 2019) requirements
  1. Contains more than one instance of a concept
  2. Shows an instance of the concept interacting with other objects
  3. Shows an instance of the concept performing an activity
  4. Displays a set of diverse objects or features



MaRVL-zh 花椰菜 (Cauliflower)



MaRVL-ta **Clonji** (Buttermilk)



MaRVL-sw Jembe (Shovel)



MaRVL-tr Raki (Raki)

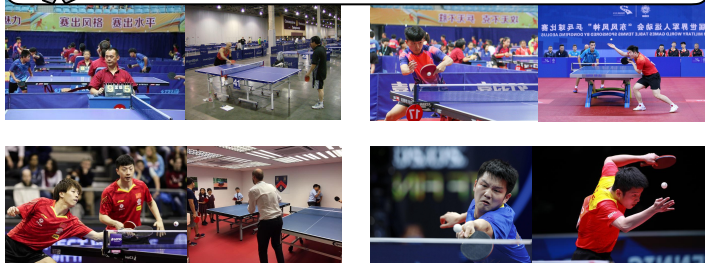


# Step 3. Language Annotation

Written by native speakers



MATCH 4 PAIRS AT RANDOM



WRITE CAPTION TRUE ONLY FOR 2 PAIRS



右图中的人在发球, 左图中的人在接球。



VALIDATE ANNOTATIONS



右图中的人在发球, 左图中的人在接球。



FINAL VALIDATION



右图中的人在发球, 左图中的人在接球。

(The man in the right image is serving a ball while the man in the left image is returning a ball.)

Fleiss' kappa: 93%

# Dataset Examples

MaRVL-tr Kanun (çalgı)



Görsellerden birinde dizlerinde kanun bulunan birden çok insan var

(In one of the images, there are multiple people with qanuns on their knees)

Label: True

MaRVL-ta வை (Vada)



இரண்டு படங்களிலும் நிறைய மசால் வடைகள் உள்

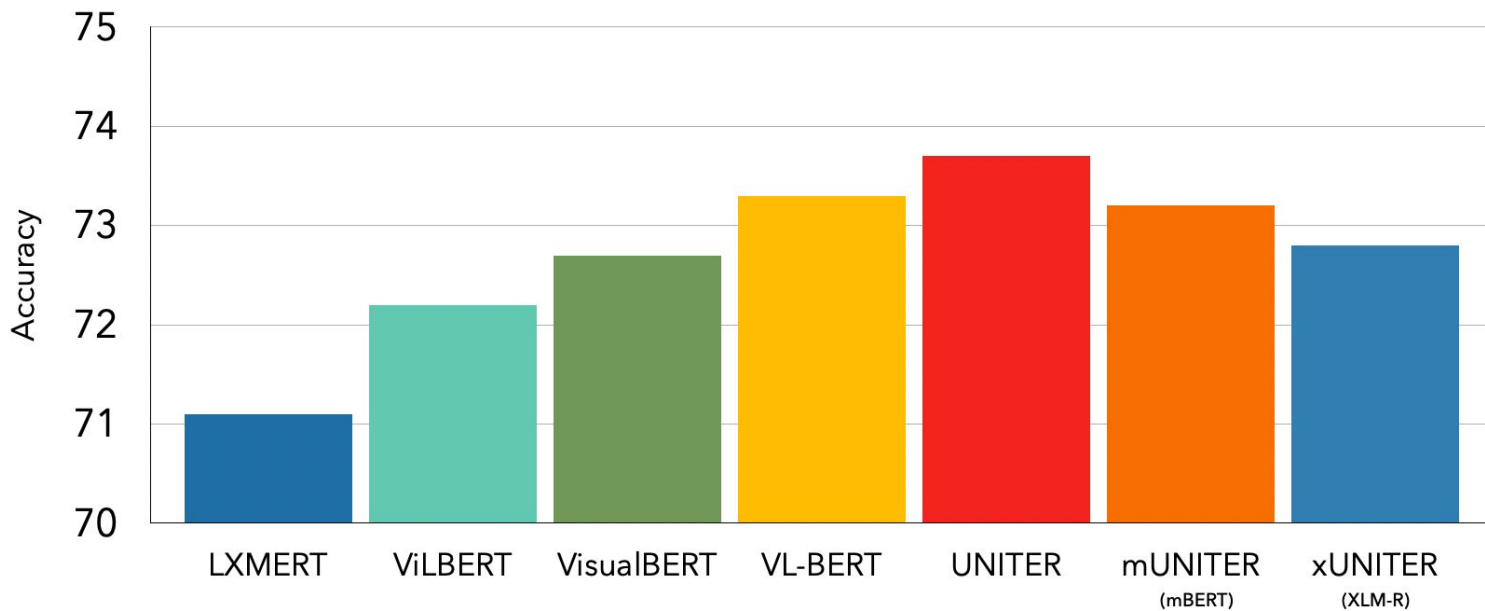
(Both images contain a lot of masala vadas)

Label: False

# Pretraining and Finetuning

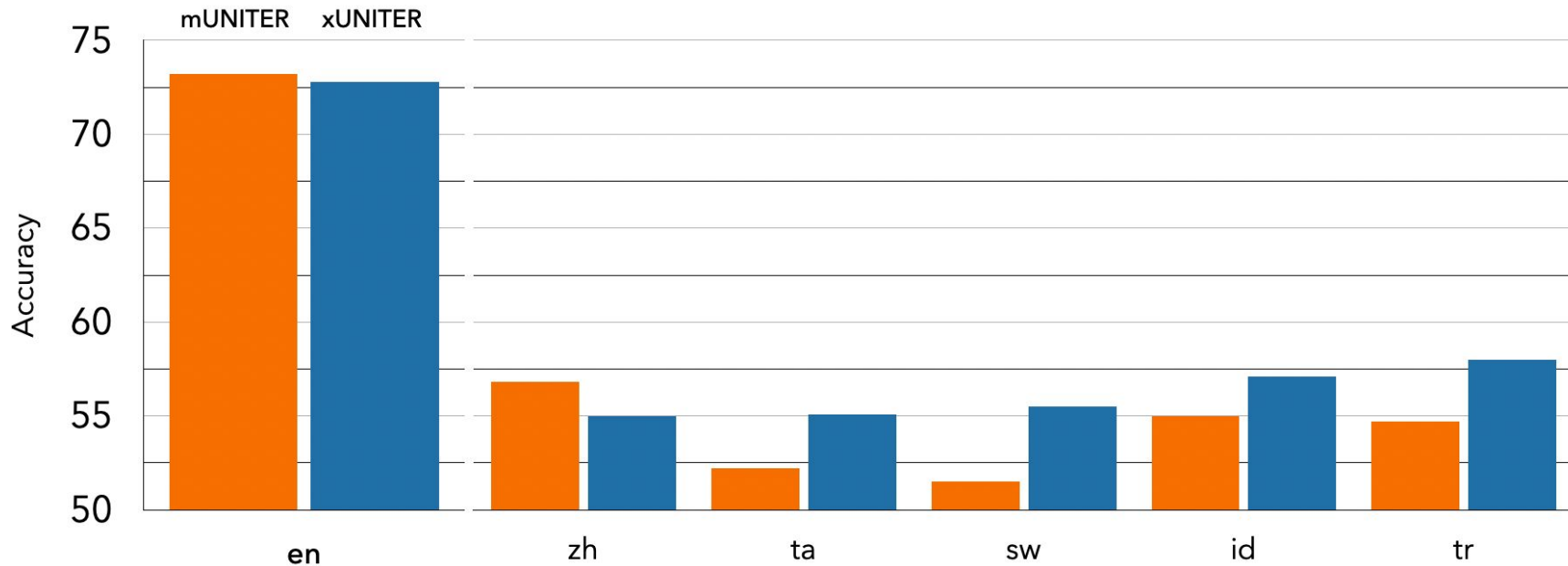
- Two new multilingual UNITER-based models
  - Pretrained on English Conceptual Captions + 104 languages Wikipedia
    - mUNITER: Initialised from mBERT
    - xUNITER: Initialised from XLM-R
- Finetune on 86,373 data points in **English** NLVR2 (Suhr+, 2019)
- Test on 5,560 datapoints in MaRVL (5,560 datapoints)
  - **Zero-shot**: Multilingual inputs directly in a cross-lingual approach
  - **Translate-test**: English models by machine translating language data

# English NLVR2 Results (Sanity check)



m/xUNITER perform similarly to English-only models

# MaRVL Zero-shot Results



Zero-shot transfer: substantial drop in performance

# Conclusions

- Concepts and images in existing V&L datasets have an NA/EU bias
- Devise a new protocol for data creation driven by native speakers
- **MaRVL**: V&L reasoning dataset in 5 typologically diverse languages
- Implications beyond vision and language research
  - Multilingual datasets should not just be translations of English data

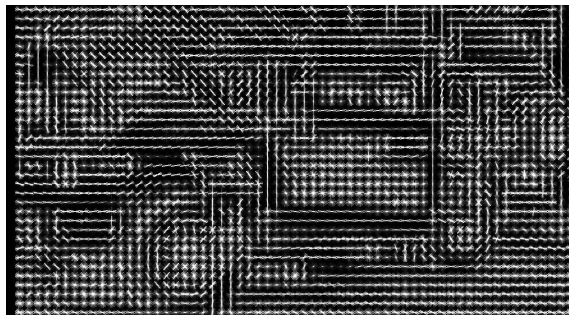
## 2. Data Representation

# Three Levels of Representation

---

- Perceptual
- Pre-processed features
- Raw input

- Yellow
- Has wheels
- Metal
- Five-door
- Can transport
- ...





# Perceptual Norms

---

- Ask people to write down the words that are triggered by textual stimuli.
- Stimuli: 541 noun concepts
- Norms are categorized into the likely knowledge source

Moose		
is large	27	visual-form and surface
has antlers	23	visual-form and surface
has legs	14	visual-form and surface
has four legs	12	visual-form and surface
has fur	7	visual-form and surface
has hair	5	visual-form and surface
has hooves	5	visual-form and surface
is brown	10	visual-color
hunted by people	17	function
eaten as meat	5	function
lives in woods	14	encyclopedic
lives in wilderness	8	encyclopedic
an animal	17	taxonomic
a mammal	9	taxonomic
an herbivore	8	taxonomic

# Perceptual Norms: Pros / Cons

---

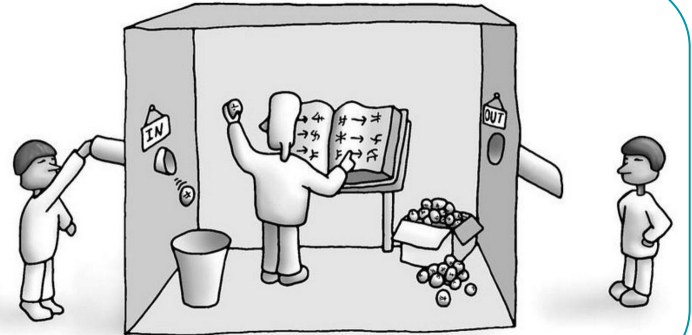
## Pros

- Seemingly simple task
- Rich features

## Cons

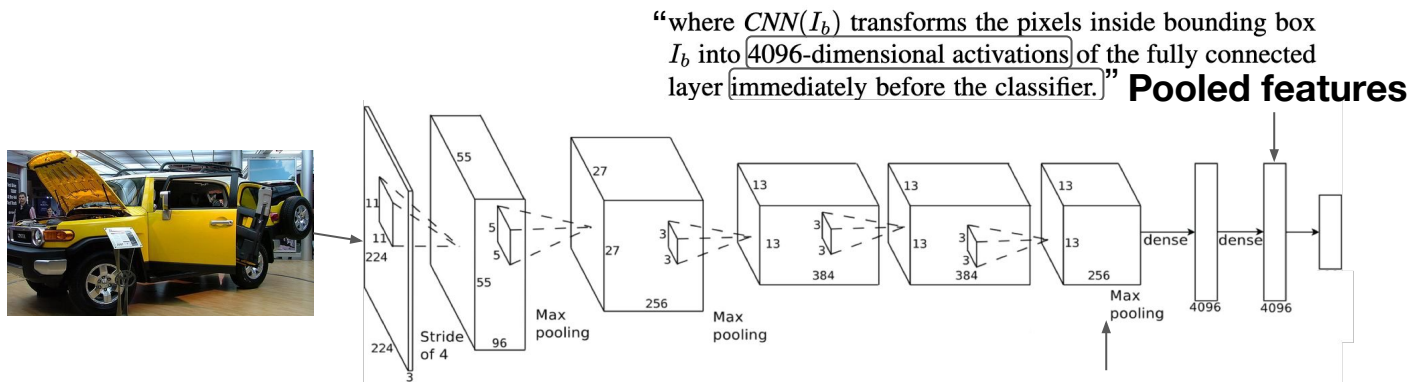
- Can it scale?
- Handling ambiguity

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



# Spatial and Pooled Visual Features

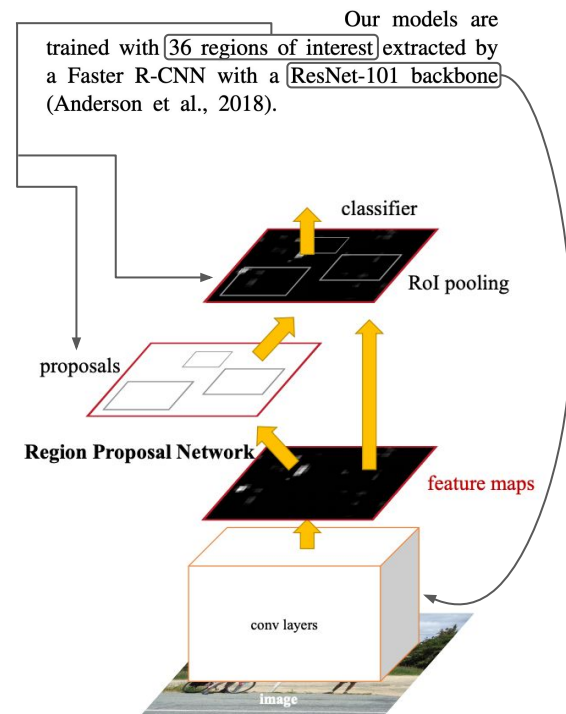
- Earliest work in neural-network era used pooled or spatial preserving features from a pretrained Convolutional Neural Network.



**Spatial features** “In our experiments we use the  $14 \times 14 \times 512$  feature map of the fourth convolutional layer before max pooling.”

# Pre-processed Visual Features

- Faster R-CNN region-based feature vectors
  - Trained on the Visual Genome Dataset
  - The Region Proposal Network suggests the location of *regions of interest*.
  - RoI pooling performs spatial pooling in the final CNN layer to give a 2048D vector.



# Pre-processed: Pros / Cons

---

## Pros

- Long-established practice
- Usually an offline process: do it once and forget

## Cons

- Large datasets require specialized storage
- Not obvious how to randomly augment data
- Specialist knowledge can be opaque to newcomers

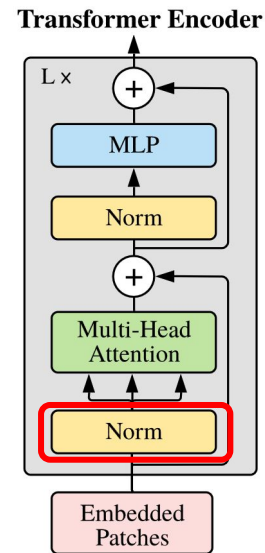
# Raw Input

---

- Directly process data from the raw images or speech signal.
- Images:
  - Vision Transformer (ViT)
  - Swin Transformer
- Speech
  - Spectrogram Transformer
  - AudioMAE

# Vision Transformer

- Good news! You are already almost an expert in how the Vision Transformer works
  - Split image into K patches
  - Embed each patch
  - Add position information
  - Encode using Transformer blocks that include an **extra pre-norm layer** for stability.

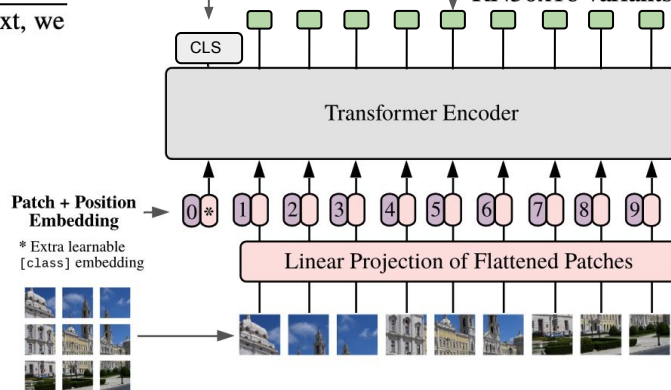


# Extracting ViT Features

- Extract **pooled features** or **patch-level features**

To extract visual information from an image  $x^i$ , we use the visual encoder of a pre-trained CLIP [29] model. Next, we

For the CLIP encoders, we extract the feature grid before the pooling layers, resulting in an  $N \times N$  grid, where  $N = 7, 7, 12$  for the ViT-B/32, RN50x4 and RN50x16 variants of CLIP respectively.





# Raw input: Pros / Cons

---

## Pros

- Data augmentation is straightforward because you always have the raw input
- Fewer preprocessing steps means fewer creeping errors

## Cons

- Smaller batches with an extra model on the GPU
- Potentially many inputs

# Summary

---

- Many options for how to represent your multimodal inputs
  - Language-oriented
  - Object / stuff oriented
  - Raw inputs
- **No universally best option** but raw inputs are promising because the visual representation model can be fully differentiable.

# 3. Modelling

# Main Approaches

---

## Cross encoding

Transformer

The red horse



## Dual encoding

Language

Vision

The red horse



## Visual Prefix

Language Model

Embed

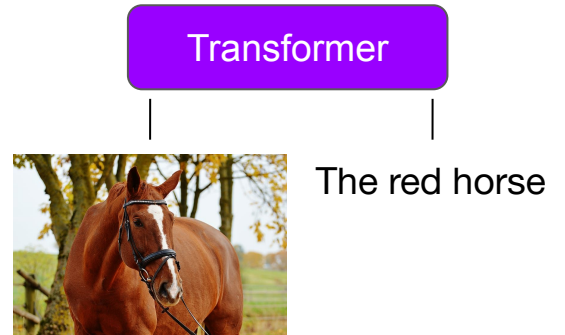
The red horse



# Cross-encoding Models

---

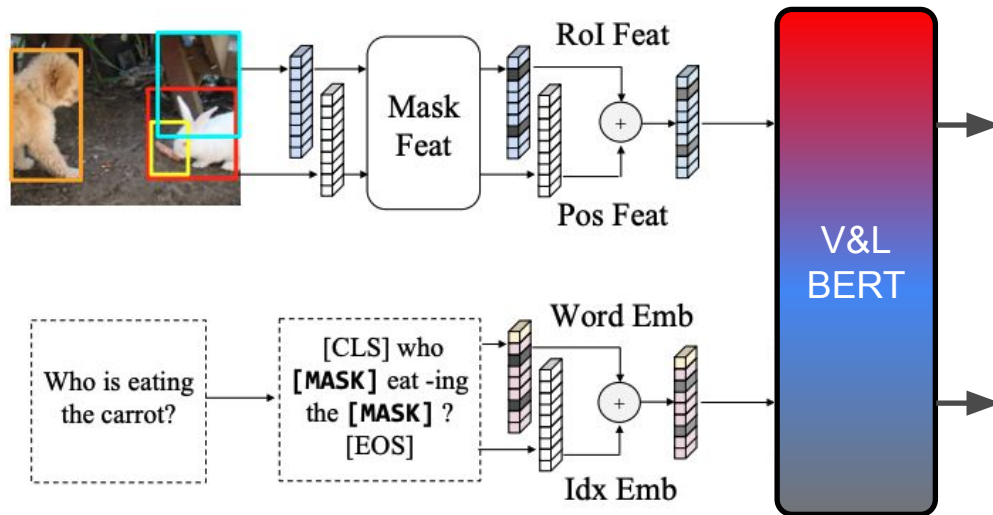
- Emerged as a key modelling approach in 2019 with a flurry of approaches to creating visually-grounded BERT models.
- This is a form of *model-based fusion*
- The backbone consists of two components:
  - language model
  - visual encoder



# High-level Overview

Image:

Faster R-CNN, or  
raw pixels + ViT

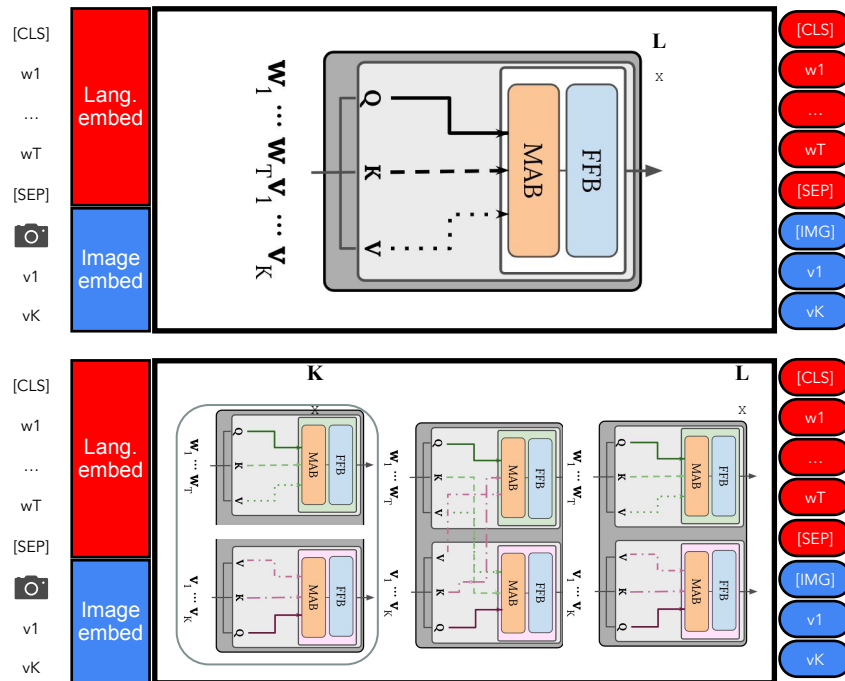


Language:

BERT tokens

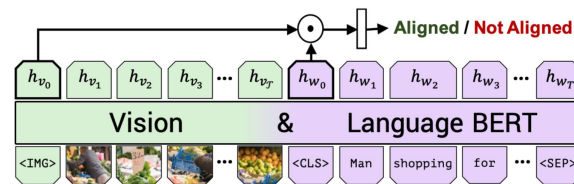
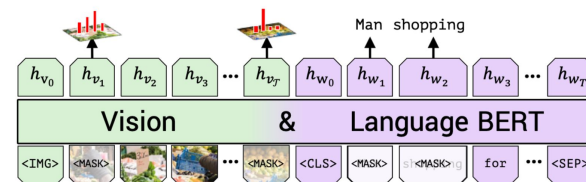
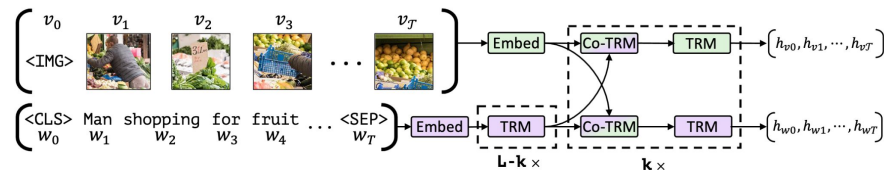
# Single- & Dual-Stream Architectures

- Single-stream
  - Concatenate inputs into one sequence
- Dual-stream
  - Process modalities independently
    - Intra-modal
    - Inter-modal



# 2019: ViLBERT

- Dual-stream model
- Initialized from BERT
- Visual features extracted from 10-36 regions using Faster-RCNN
- Pretrained on Conceptual Captions
  - Masked Language Modelling
  - Masked Region Classification
  - Image-Text Matching





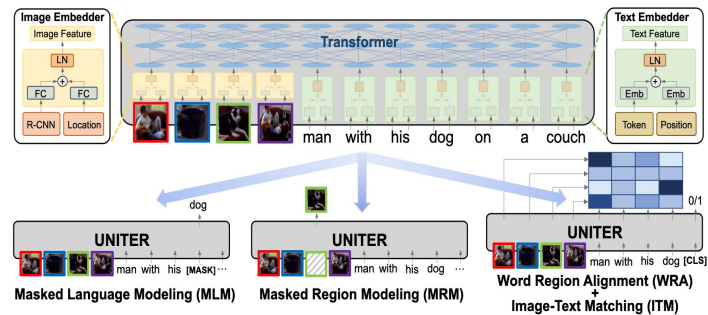
# MLM, MRC, ITM

---

- Masked Language Modelling  $\mathcal{L}_{\text{MLM}}(\theta) = -\mathbb{E}_{(\mathbf{w}, \mathbf{v}) \sim D} \log P_{\theta}(\mathbf{w}_{\mathbf{m}} | \mathbf{w}_{\setminus \mathbf{m}}, \mathbf{v})$ 
  - Same as BERT et al.
- Masked Region Classification  $\mathcal{L}_{\text{MRM}}(\theta) = \mathbb{E}_{(\mathbf{w}, \mathbf{v}) \sim D} f_{\theta}(\mathbf{v}_{\mathbf{m}} | \mathbf{v}_{\setminus \mathbf{m}}, \mathbf{w})$ 
  - Mean Squared Error Regression over the 2048D feature vector; or
  - Predict the probability distribution over the 1600 Faster R-CNN classes
- Image-Text Matching  $\mathcal{L}_{\text{ITM}}(\theta) = -\mathbb{E}_{(\mathbf{w}, \mathbf{v}) \sim D} [y \log s_{\theta}(\mathbf{w}, \mathbf{v}) + (1 - y) \log(1 - s_{\theta}(\mathbf{w}, \mathbf{v}))]$ 
  - 50% chance of randomly sampling a mis-matched sentence
  - Predict with a binary classifier (aka Next Sentence Prediction)
- Note: 15% masking usually spans both modalities

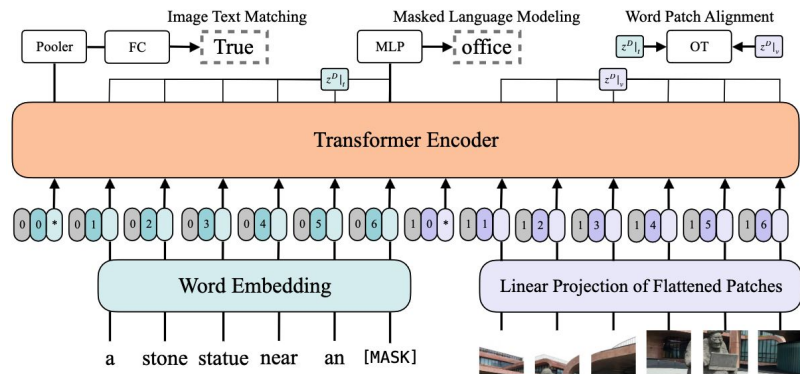
# 2020: UNITER

- Single-stream model
- Initialized from BERT
- Visual features from Faster-RCNN
- Pretrained on Conceptual Captions, Visual Genome, COCO, SBU Captions
  - Masked Language Modelling
  - Masked Region Classification
  - Image-Text Matching
  - Word-Region Alignment

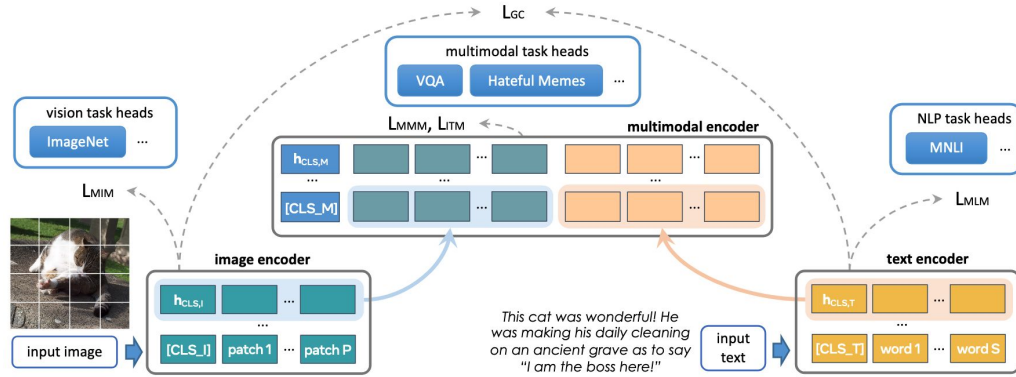


# 2021: ViLT

- Single-stream model
- Initialized from BERT
- Visual features extracted from ViT-B/32
- Pretrained on Conceptual Captions, Visual Genome, COCO, SBU Captions
  - Masked Language Modelling
  - Image-Text Matching
  - Word-Patch Alignment



# 2022: FLAVA

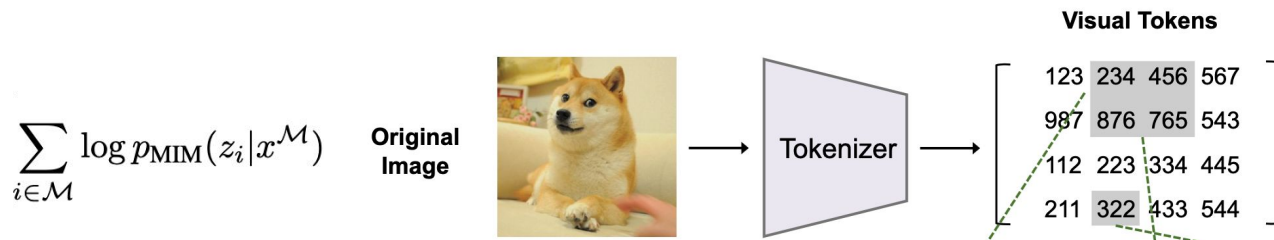


- Dual-stream Visual features extracted from ViT-B/16
- Pretrained on PMD70M
  - Masked Language Modelling, Masking Image Modelling
  - Image-Text Matching, Masked Multimodal Modelling
  - Global Contrastive Matching

# MIM and CL

- Masked Image Modelling

- Immediately after the image encoder and before multimodal encoding
- Tokens from a discrete VAE (BEiT)



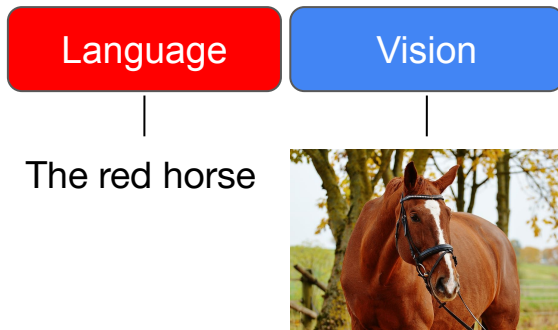
- Contrastive Loss

- On the CLS embedding of each unimodal encoder  $\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E} \left[ \log \frac{f(\mathbf{t}, \mathbf{i})}{\sum_{\mathbf{t}' \in T} f(\mathbf{t}', \mathbf{i})} \right]$

# Dual-encoding Models

---

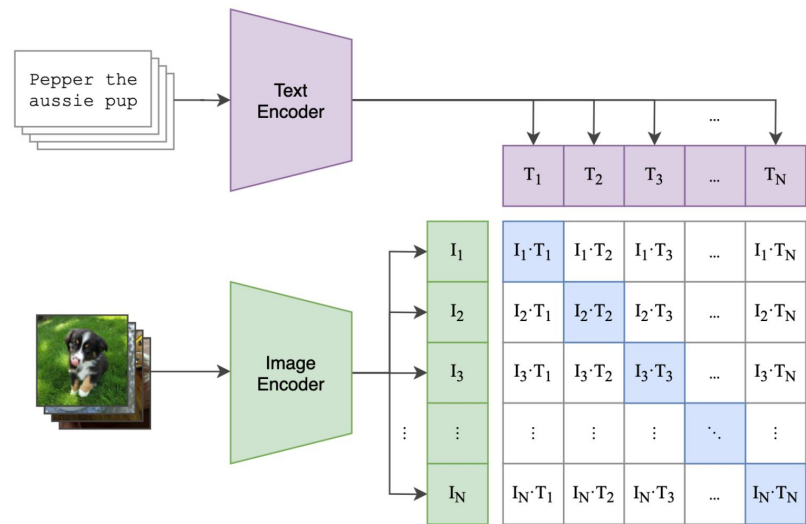
- Emerged as a sample-efficient alternative to cross-encoding.
- The backbone consists of two separate components:
  - language encoder
  - visual encoder



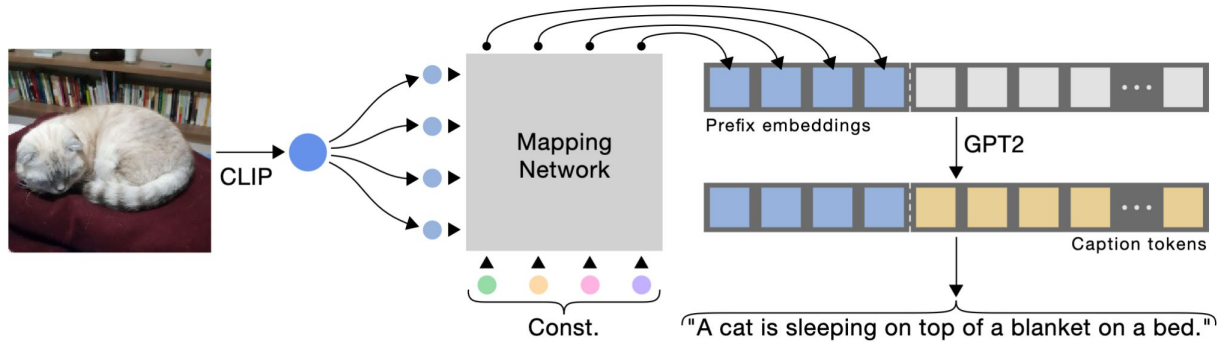
# CLIP

- 12 Layer Transformer Encoder
- ViT or ResNet Visual Encoder
- Maximize the similarity of the embeddings of paired examples (I, T):

$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E} \left[ \log \frac{f(\mathbf{t}, \mathbf{i})}{\sum_{\mathbf{t}' \in T} f(\mathbf{t}', \mathbf{i})} \right]$$



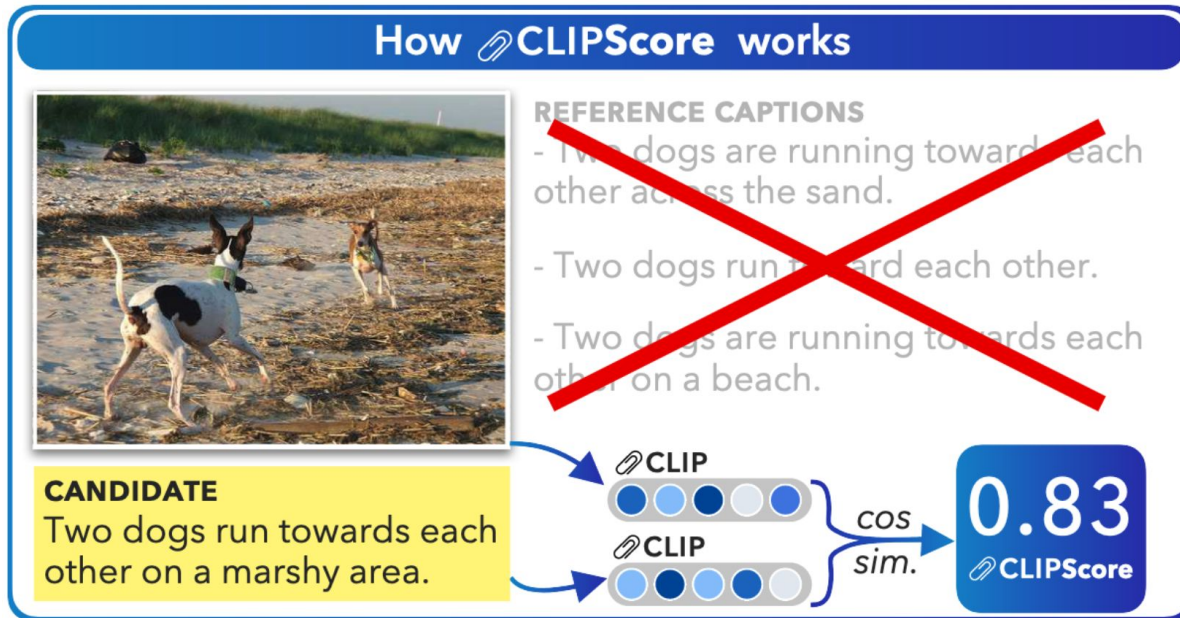
# CLIP for Captioning



- Use CLIP as a feature extractor and GPT-2 as a language model.
  - Only train the mapping network to generate prefix embeddings
  - Lightweight system that exploits pretrained models

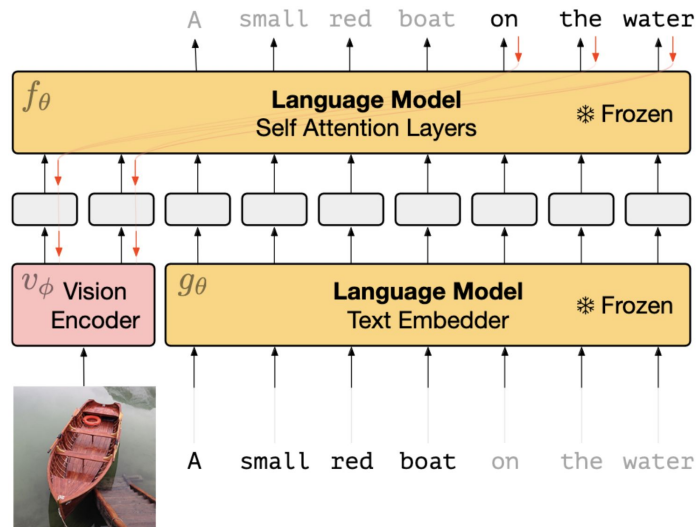
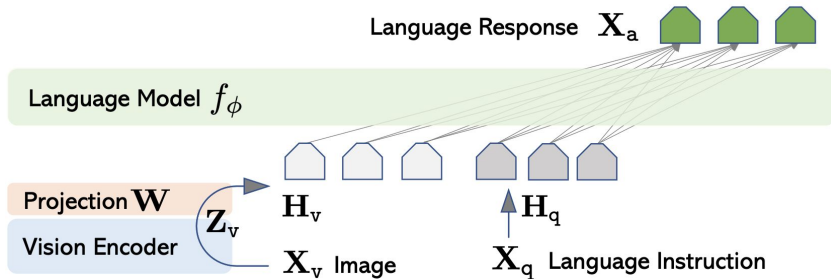


# CLIP for Evaluation

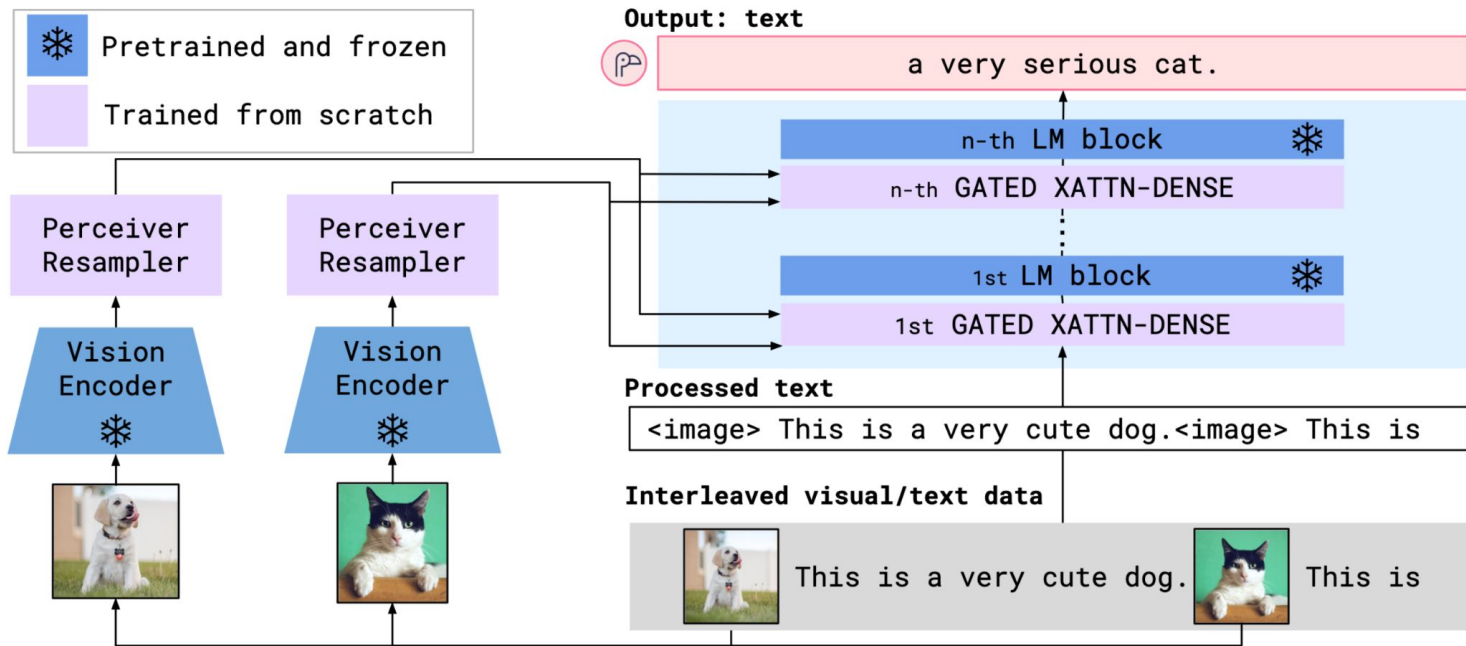


# Visual Prefix Learning

- Exploit the representations learned during large-scale modality specific pretraining



# Learning Dense Cross-Attention



# Summary

---

- Cross-encoding:
  - Many advances in which parts of the input contribute to loss
  - Shift from regions-of-interest to Vision Transformers
- Dual-encoding:
  - Excellent cross-domain transfer to a wide range of problems
- Visual Prefix Learning:
  - Exploit the benefits of single-modality pretraining

**Q: Do we need to learn everything for image captioning in-weights?**

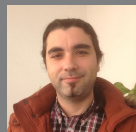
---

# SmallCap: Lightweight Image Captioning Prompted with Retrieval Augmentation

CVPR 2023



R. Ramos



B. Martins



D. Elliott

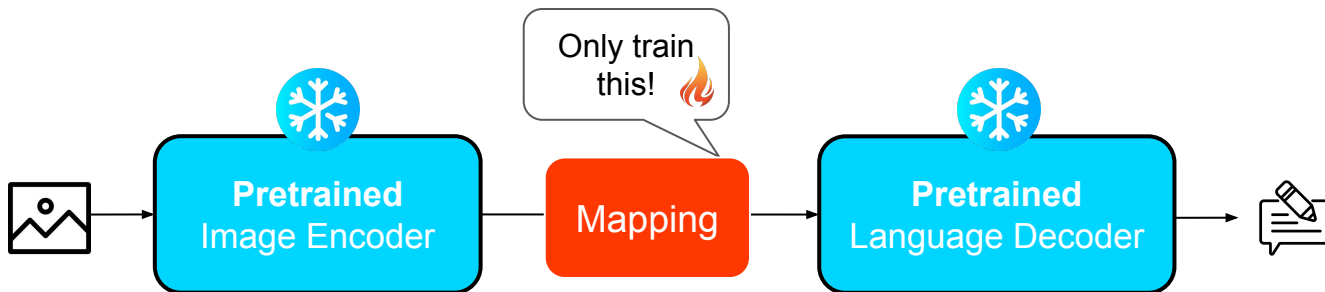


Y. Kementchedjheva

# Motivation

---

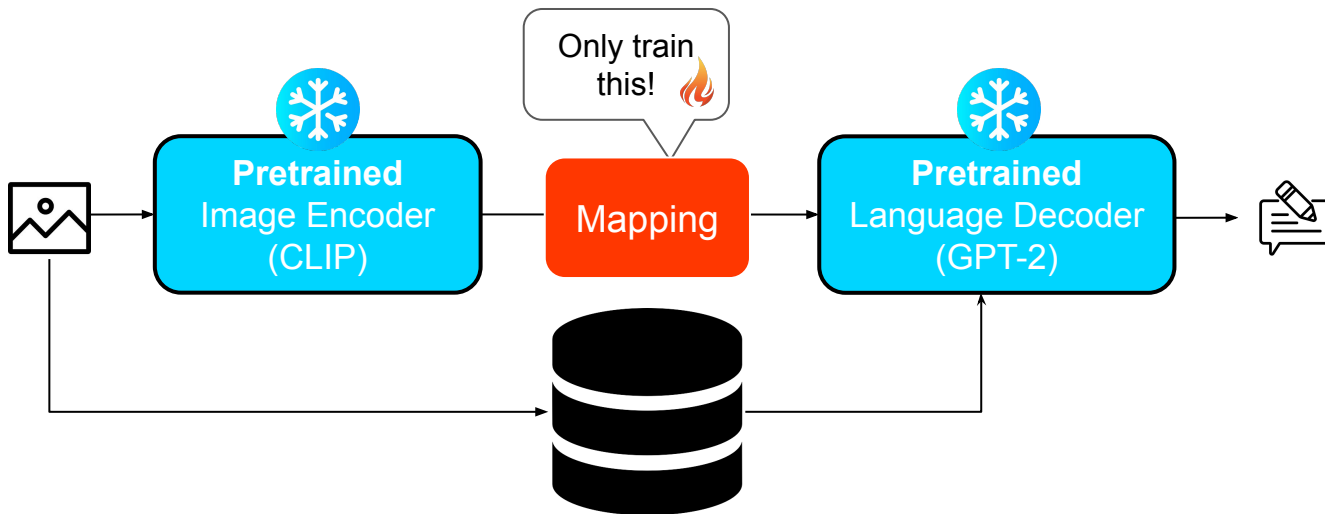
- Main trend in V&L is training bigger models on more data
- Alternative is emerging that re-uses independent backbone models
  - CLIPCap, I-Tuning



# Lightweight Training through Retrieval

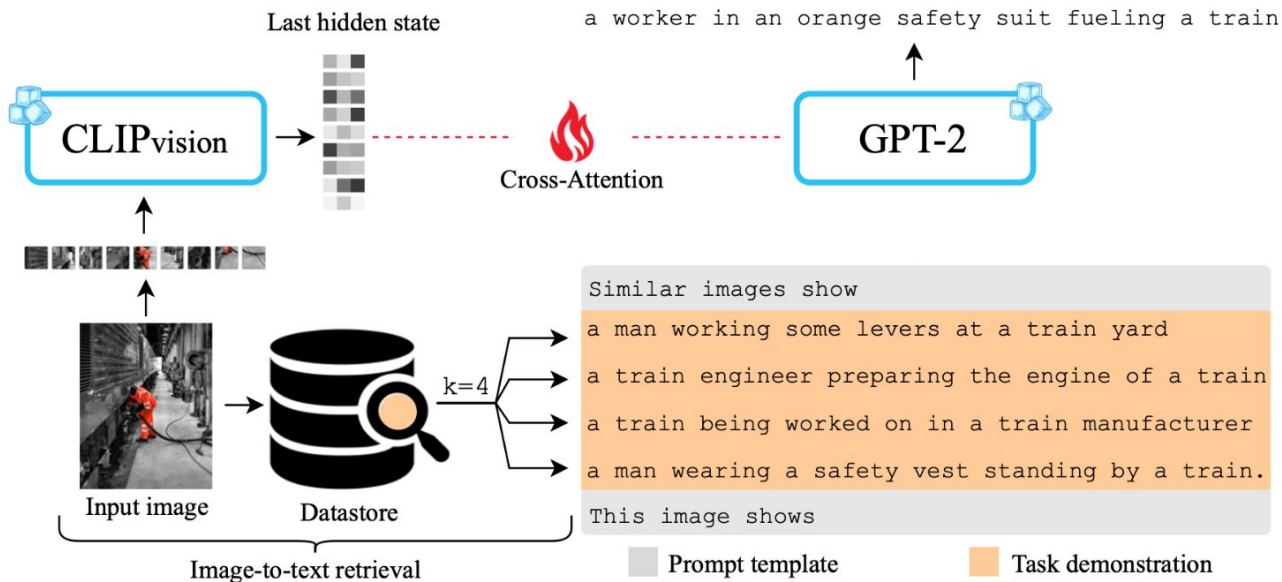
---

- Given the success of multimodal retrieval augmentation, can we extend this to the lightweight training paradigm?





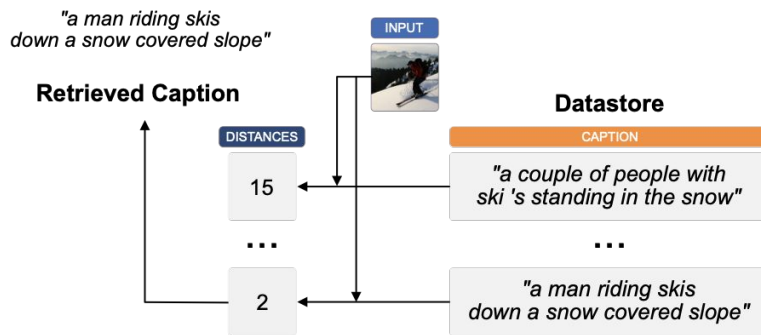
# SmallCap Model



# Retrieval System

---

- Build a FAISS datastore: store high-dimensional vectors
  - Captions of images represented with CLIP embeddings
- Retrieve k nearest-neighbours captions from datastore
  - Image embedding compared against datastore caption vectors



# 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

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

$$W_i^K, W_i^Q$$

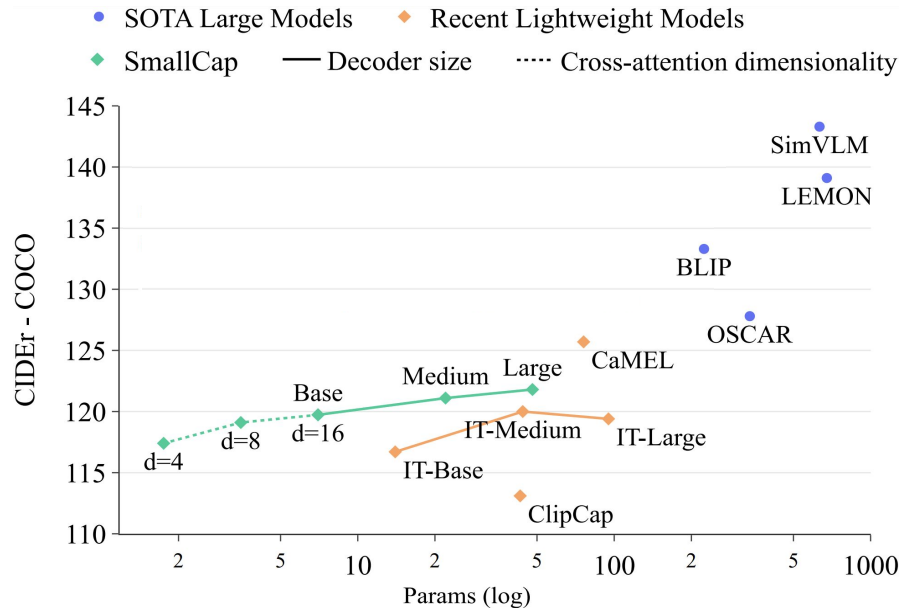
$$, W_i^V$$

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

$d$

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

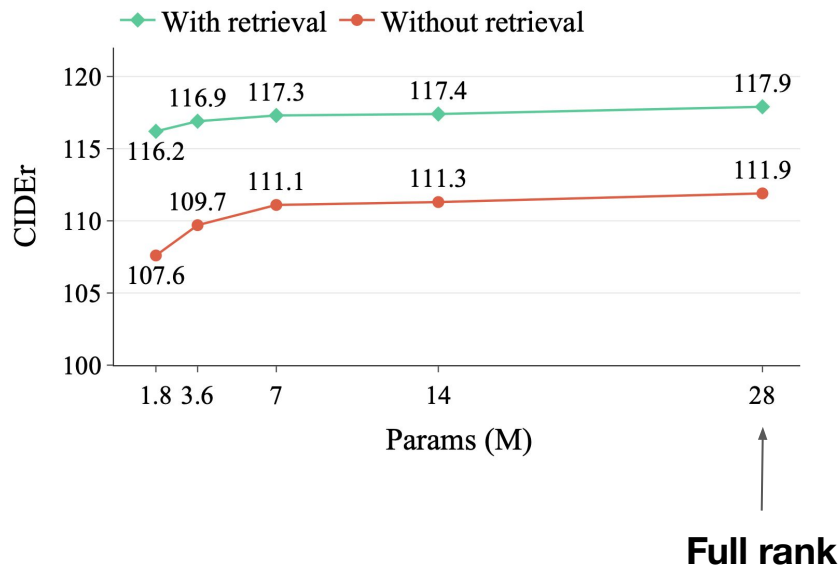
# Results



- 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:
  - SMALLCAP model performance degrades at a higher rate

# Training-Free Domain Transfer

---

- SmallCap was trained on COCO but we can easily swap the datastore.
- Much stronger performance than other lightweight approaches

	Flickr30k	VizWiz	MSR-VTT
ClipCap	41.2	28.3	12.5
CaMEL	55.2	37.6	20.7
<b>SMALLCAP</b>	<b>60.6</b>	<b>55.0</b>	<b>28.4</b>

# Qualitative Example on 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

# Try it yourself

---





# Conclusions

---

- SmallCap:
  - light to train
  - easily transferred across domains without retraining
- Prompt-based conditioning method, wherein retrieved captions are used as a prompt to a generative language model
- Strong performance in out-of-domain settings

**Q: Do you even need to train?**

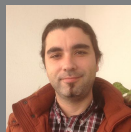
---

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

Findings of ACL 2023



R. Ramos



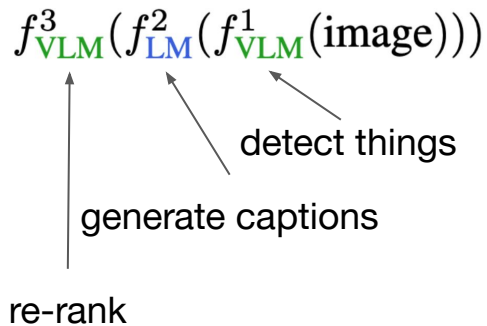
B. Martins



D. Elliott

# Socratic Models

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



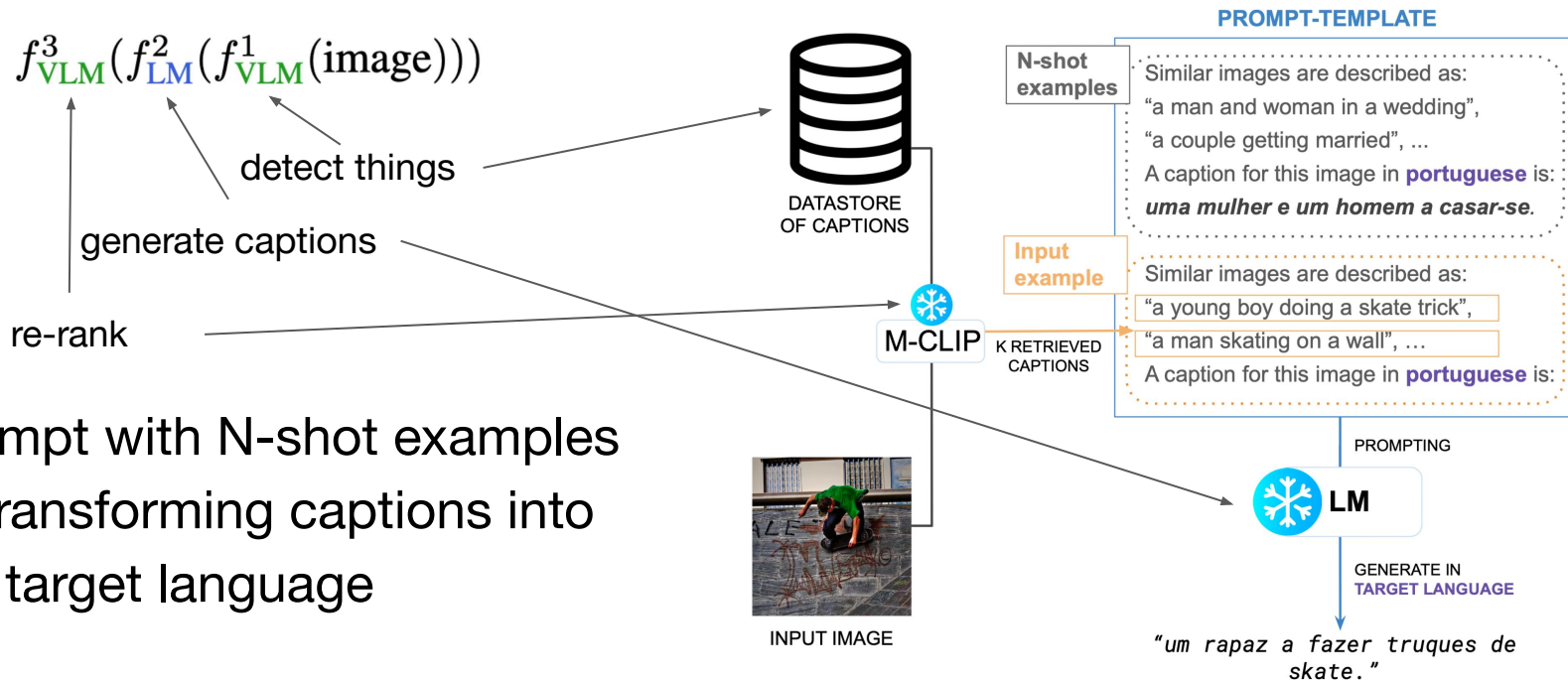
I am an intelligent image captioning bot. This image is a {img\_type}. There {num\_people}. I think this photo was taken at a {place1}, {place2}, or {place3}. I think there might be a {object1}, {object2}, {object3},... in this {img\_type}. A creative short caption I can generate to describe this image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



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

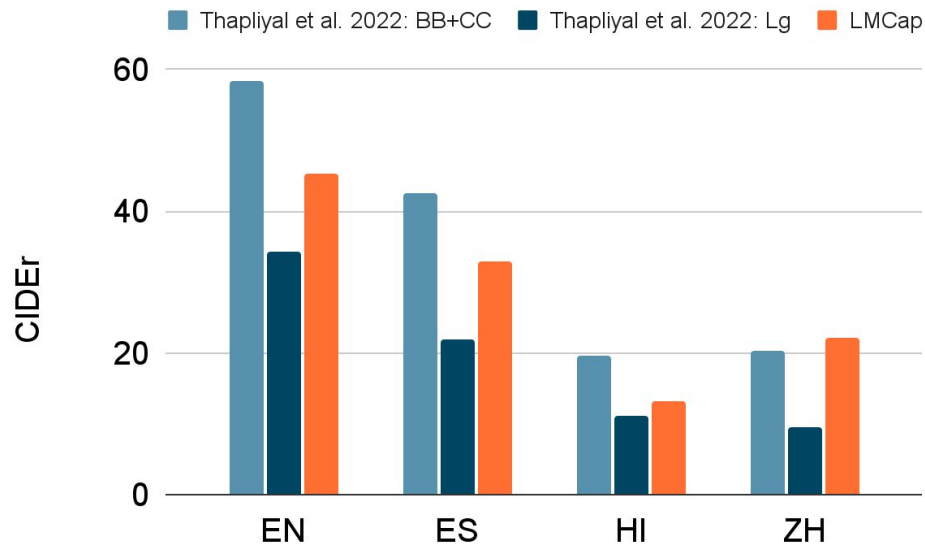
**ClipCap:** A wooden table sitting in front of a window.

# Multilingual Captioning with Retrieval Augmentation



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

# Results



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	<b>0.787</b>	<b>0.489</b>	<b>0.365</b>	<b>0.644</b>

Competitive against **fully supervised** models

Need **at least 2.9B** parameter decoder for multilingual generation

# 4. Understanding Multimodal Models

# Beyond Benchmarking

---

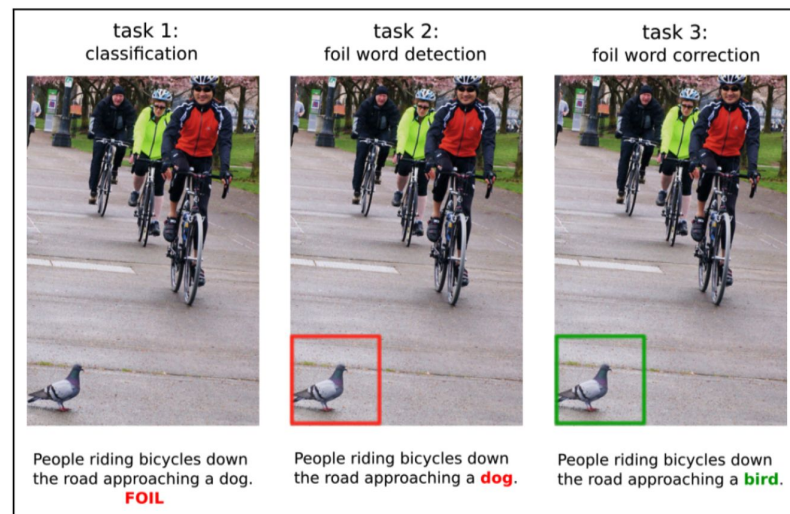
- Many questions about what drives the success of these models?
  - Better contextualization: make better use of the multimodal inputs
  - Acquire certain “skills”, e.g. counting or localization
  - Understand linguistic structures
  - Something else?
- Model-internal behaviour
  - Attention mechanism patterns
- Probing
  - Tasks related to different skills



# FOIL Captions

---

- Do V&L models really understand the relationship between words and images?
- Crowdsourced datasets that contain contextually plausible but incorrect image–text pairs

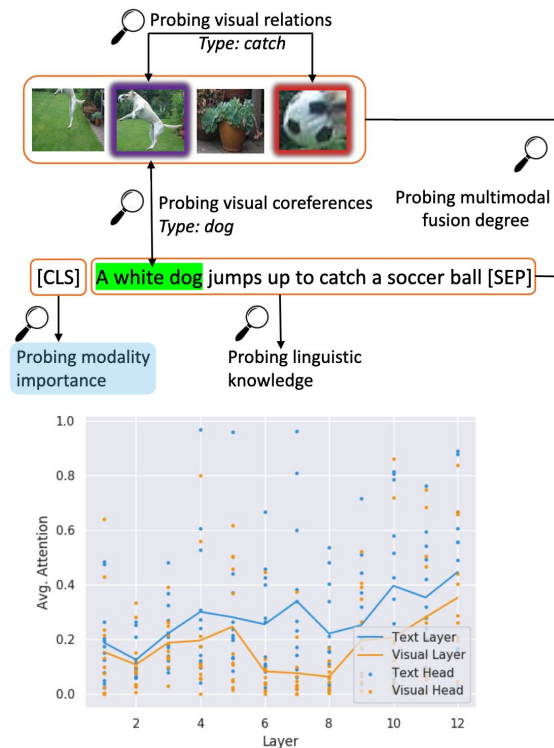


# Vision and Language Understanding Evaluation

- Suite of five model probing tasks
- **Modality Influence:** Estimate the layer-wise contribution of each modality to the [CLS] embedding:

$$I_{M,j} = \sum_{i \in \mathcal{S}} \mathbb{1}(i \in M) \cdot \alpha_{ij}$$

- The UNITER model relies more on textual features when fusing modalities throughout the model



# VALSE Benchmark

- Test visio-linguistic capabilities with image-sentence foil pairs
- Image-sentence matching task
  - Existential quantifiers
  - Semantic number
  - Counting
  - Prepositional relations
  - Action replacement / swap
  - Co-reference



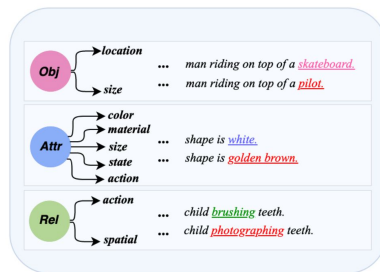
A small copper vase with **some flowers** / **exactly one flower** in it.

Metric	Model	Avg.
	Random	50.0
$acc_r$	GPT1*	60.7
	GPT2*	60.1
	CLIP	64.0
	LXMERT	59.6
	ViLBERT	63.7
	12-in-1	<b>75.1</b>
	VisualBERT	<u>46.4</u>

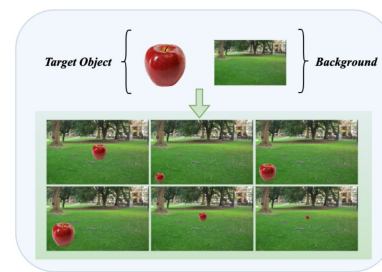
$$p(\text{caption}, \text{img}) > p(\text{foil}, \text{img})$$

# VL-CheckList

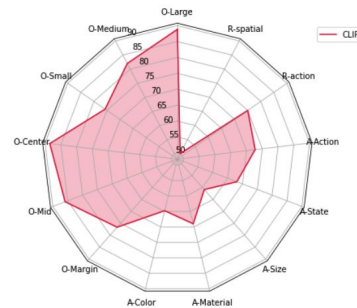
- Evaluate V&L models based on automatic manipulations to vision and language data.
- Image-Sentence matching task
- Radar chart overviews based on object / attribute / relationship variations



(a) Language Variation



(b) Vision Variation



# Subject-Verb-Object Probes

- Large-scale dataset with SVO triplets mined from Conceptual Captions and 14K images and with crowdsourced captions
- Foil detection formulation



# WinoGround

---

- 1,600 text-image pairs to evaluate compositional understanding



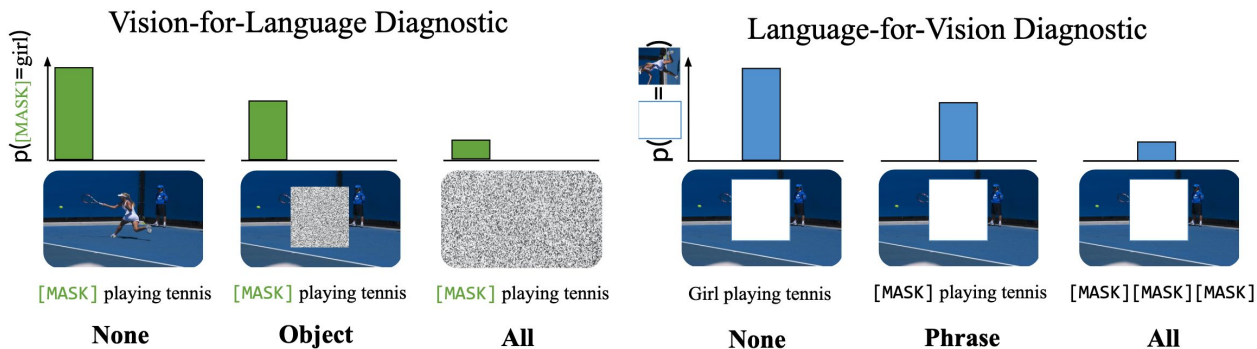
some plants  
surrounding a  
lightbulb



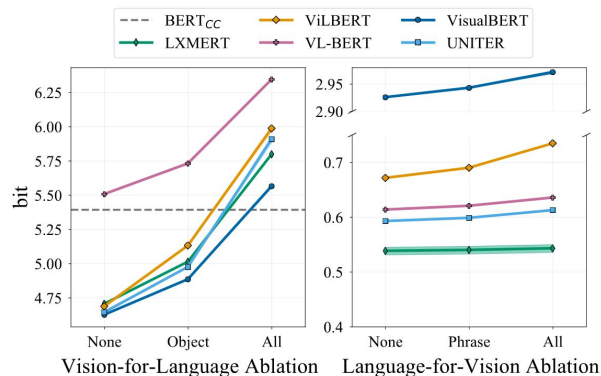
a lightbulb  
surrounding  
some plants

- Images sourced **with permission** from Getty.
- Differences are categorised into: swap dependent, swap-independent, and visual differences

# Vision-for-Language?



- Pair of diagnostic evaluations that can be applied to any model that makes MLM and MRC predictions.



# Summary

---

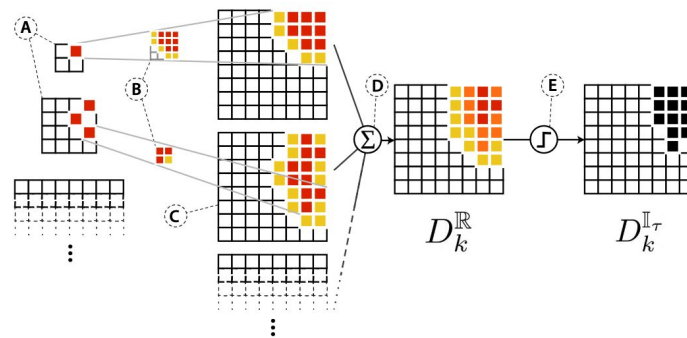
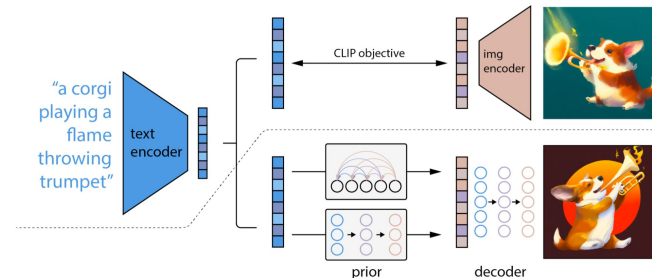
- Understanding and analysis is a vibrant area of research
- Foil detection is the most popular methodology
- Witnessing a methodological shift
  - attention analyses to linguistically-informed analyses
  - hand-crafted datasets
  - simpler accuracy-based metrics



# 5. Future Directions


# Text-to-Image Generation

- Two main approaches: VQVAE and Conditional diffusion
- Big questions:
  - Do these models produce verifiably correct outputs?
  - How can they deal with cross-cultural generation?



# Humour

- Could humour a new frontier in multimodal understanding?
  - Non-literal understanding
  - Deeper multimodal interaction
  - Social / world-level knowledge



*Can you please pass the cow?*

**Task 1** Match the Caption + Cartoon  
✗ I'd kill for some cream cheese.  
vs.  
✓ Can you please pass the cow?

**Task 2** Rank the Finalist  
✗ Welcome to Insomniacs Anonymous  
vs.  
🏆 Can you please pass the cow?

**Task 3** Explanation Generation

**Human-authored:**  
When drinking coffee or tea, people often add cream, and may ask others to pass it if it's on the other side of a table. But here, the mugs are huge, so instead of asking for a small cup of cream, they are asking for the entire cow, which is the appropriately-sized cream dispenser for these huge drinks.

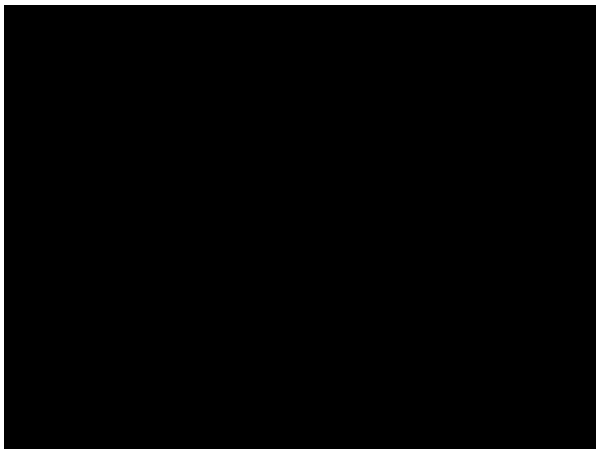
**From Pixels (OFA + T5-11B):**  
The joke is that the meeting participants are actually sitting on chairs made out of coffee mugs, which is an unlikely location for the discussion. Instead of asking for another mug of coffee, the person at the head of the table simply asks for "the cow", or a coffee machine.

**From Description (5-shot GPT 3.5):**  
"Pass the cow" is an example of a non sequitur, something that looks like a logical thing to say, but doesn't make sense in context. The humor in this cartoon comes from the large size of the coffee mugs: they are so large that they resemble buckets rather than mugs, thus making the request to "pass the cow" almost reasonable.

# Physical Understanding

---

- Predicting and explaining physical actions in the world will become of increasing importance as we create embodied agents



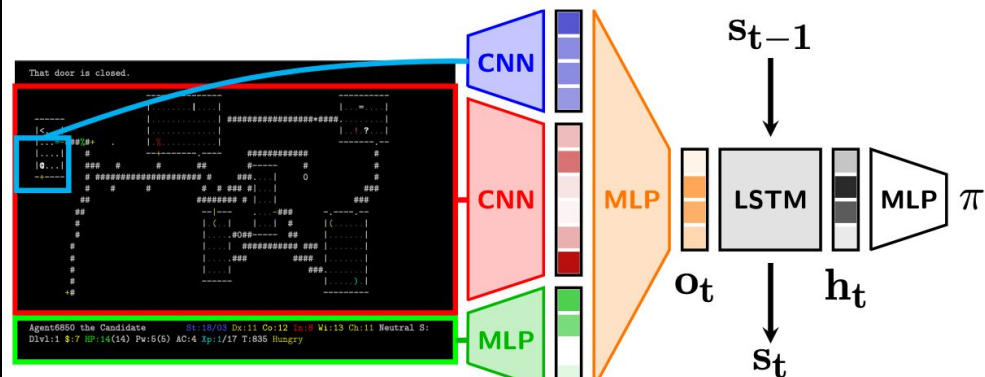
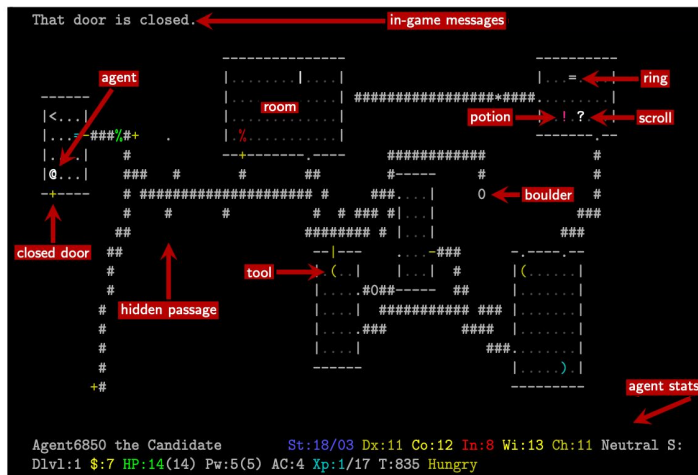
Q: How many objects are prevented by the tiny green triangle from falling into the basket?

Q: What is the color of the last object that collided with the tiny red circle?

Q: If any of the other objects are removed, will the tiny green circle end up in the basket?

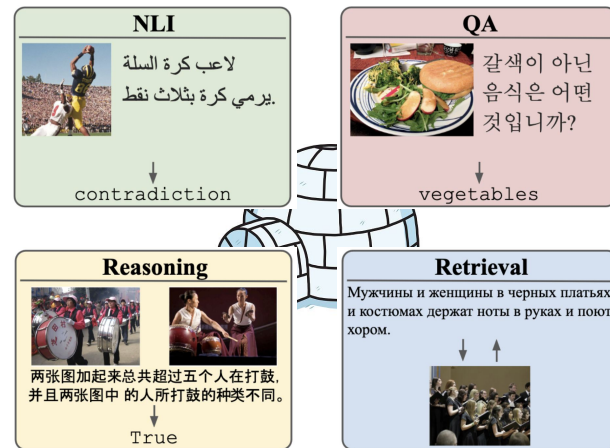
# Multimodality and Interaction

- Learning to act in procedurally-generated video game environments with rich contexts, action spaces, and long-term rewards



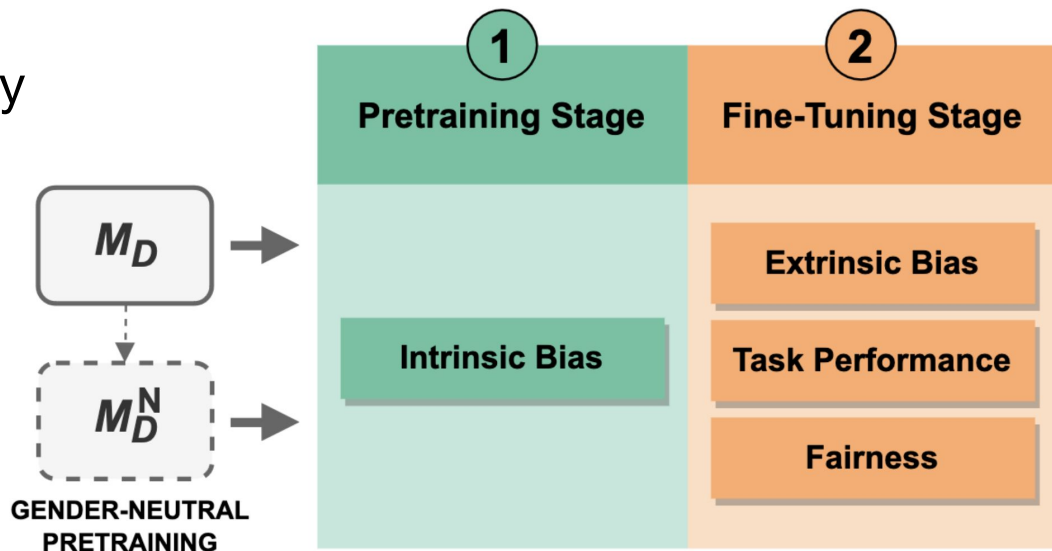
# Multilinguality

- The majority of Vision and Language research is in English
- We need resources, models, and evaluations to create useful multilingual multimodal models
- High-quality data requires:
  - time
  - money
  - community engagement



# Bias and Fairness

- What are the intrinsic biases learned during multimodal pretraining and how do they affect downstream task performance?



**Q: What if we treated language as vision?**



---

# Language Modelling with Pixels

ICLR 2023



P. Rust



J. F. Lotz



E. Bugliarello



E. Salesky



M. de Lhoneux

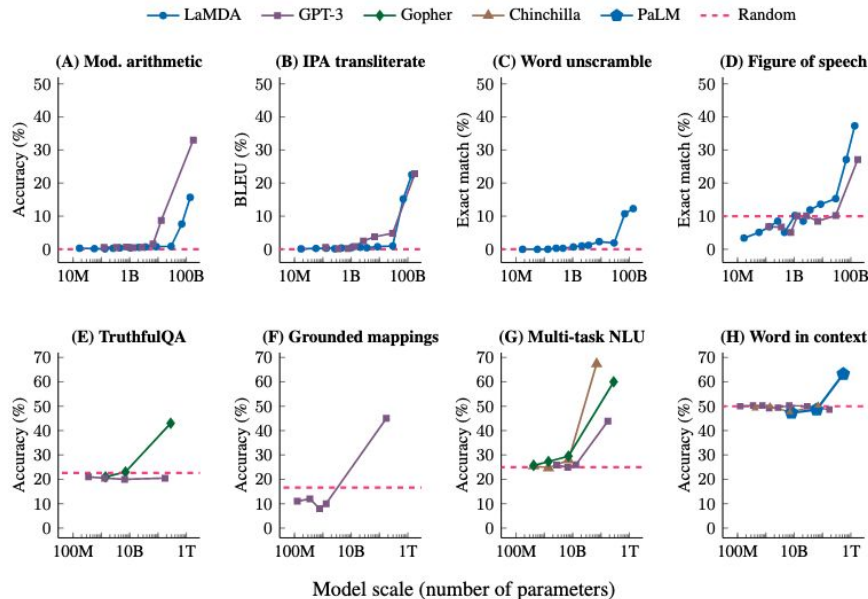
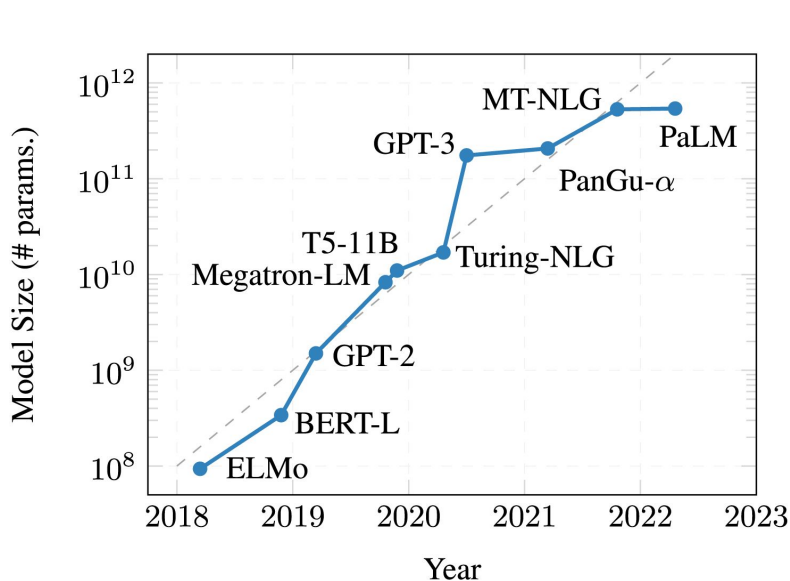


D. Elliott

---

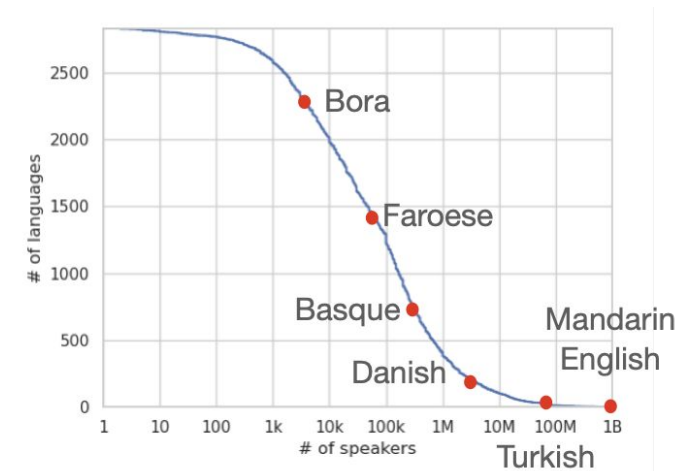
Warning: The final part of this section contains dataset samples that are racist in nature.

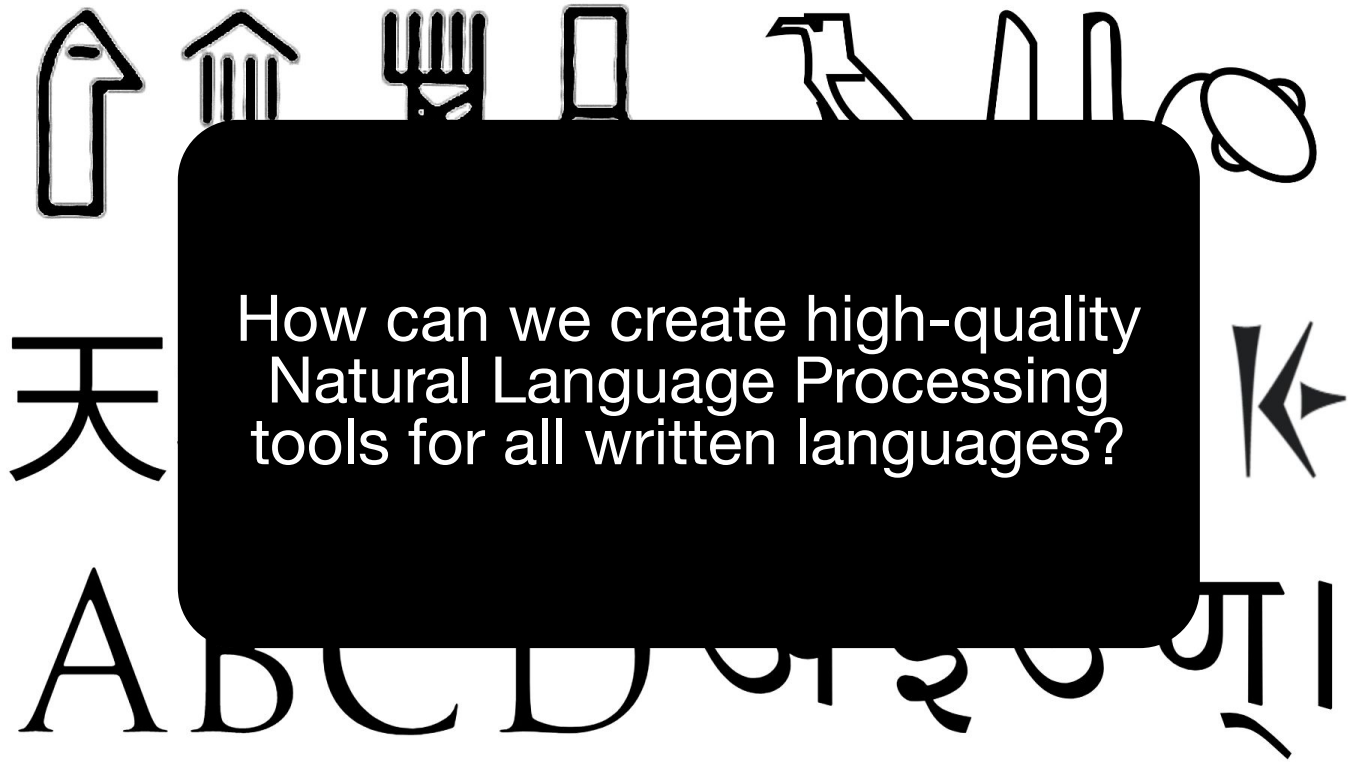
# NLP in the Era of Scale



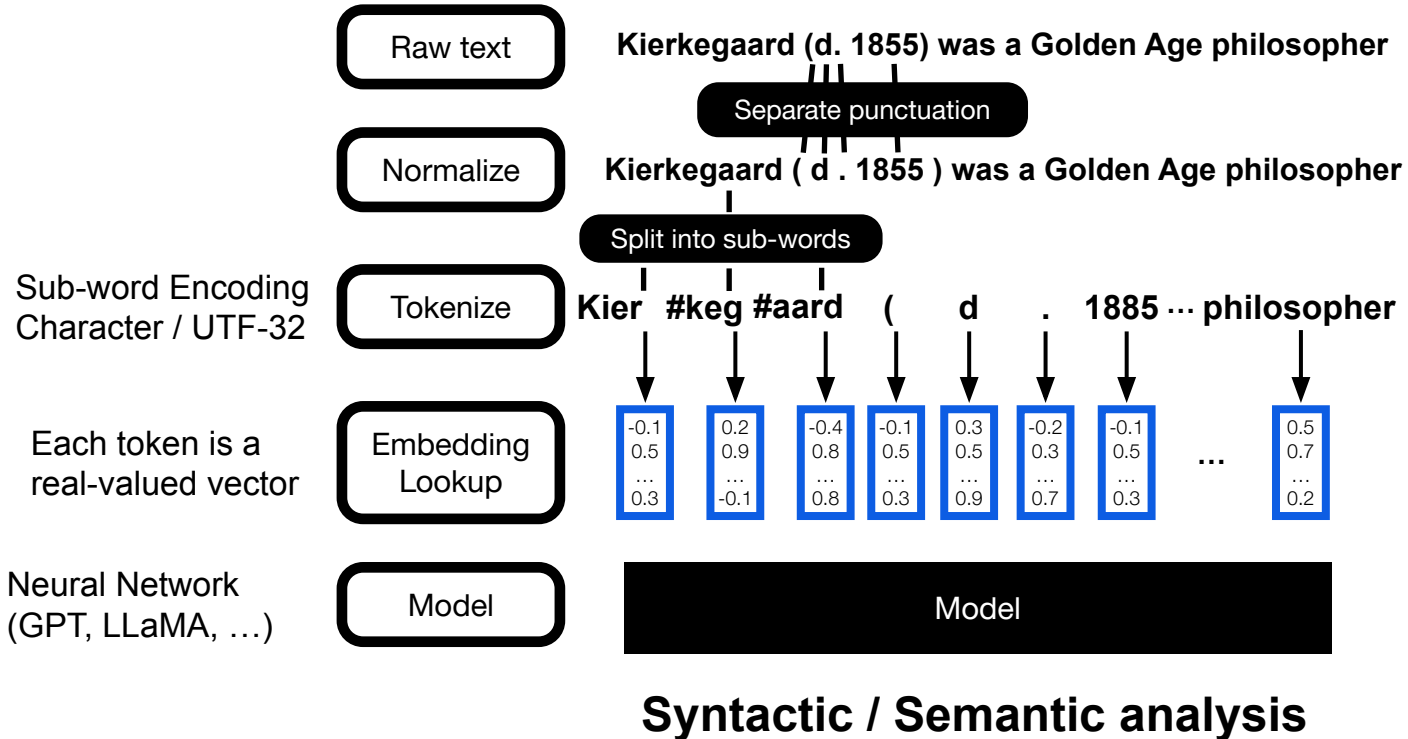
# NLP for **All** Written Languages

- There are 3,000 written languages
  - 400 with >1M speakers
- NLP usually covers 100 languages
  - Technological exclusion for billions

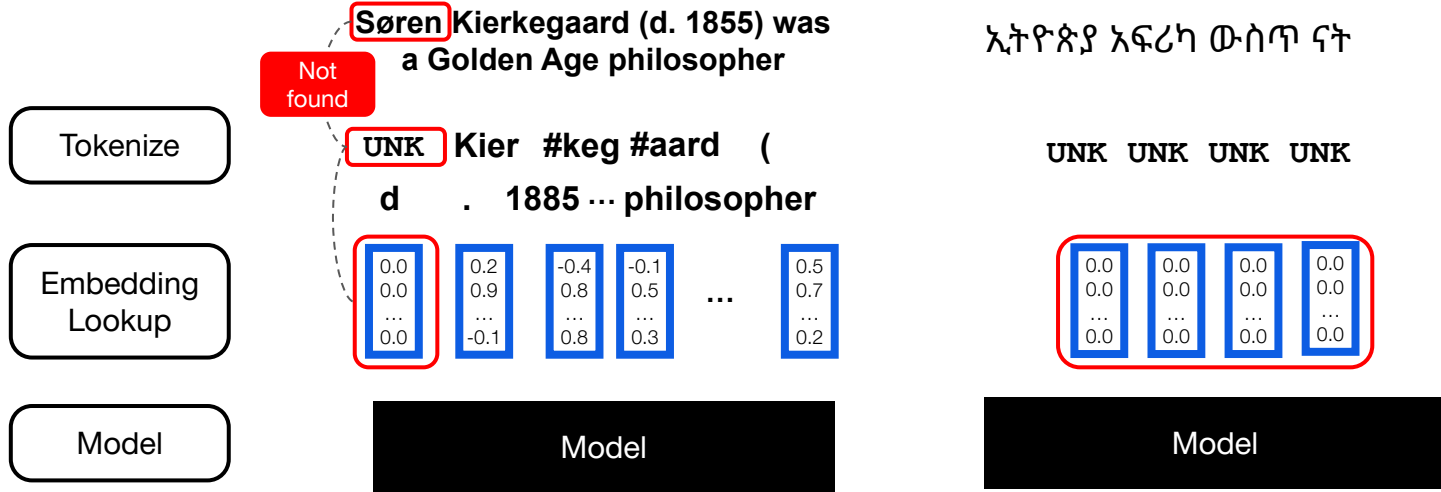




# NLP is a pipeline ...



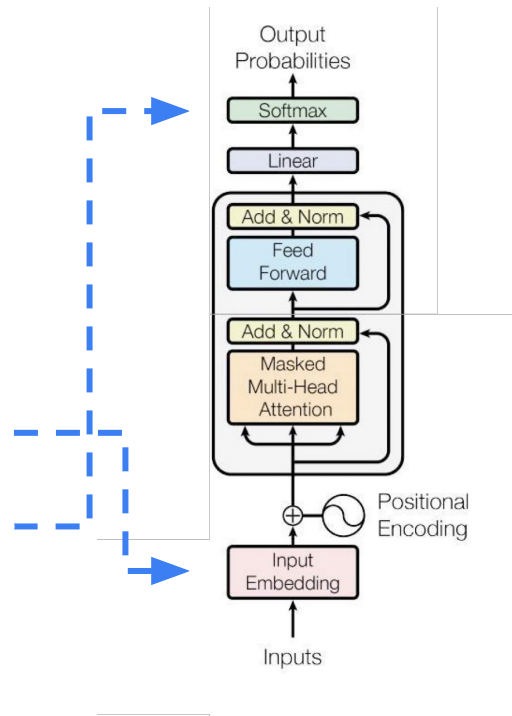
# ... that is easily broken



This issue disproportionately affects low-resource languages

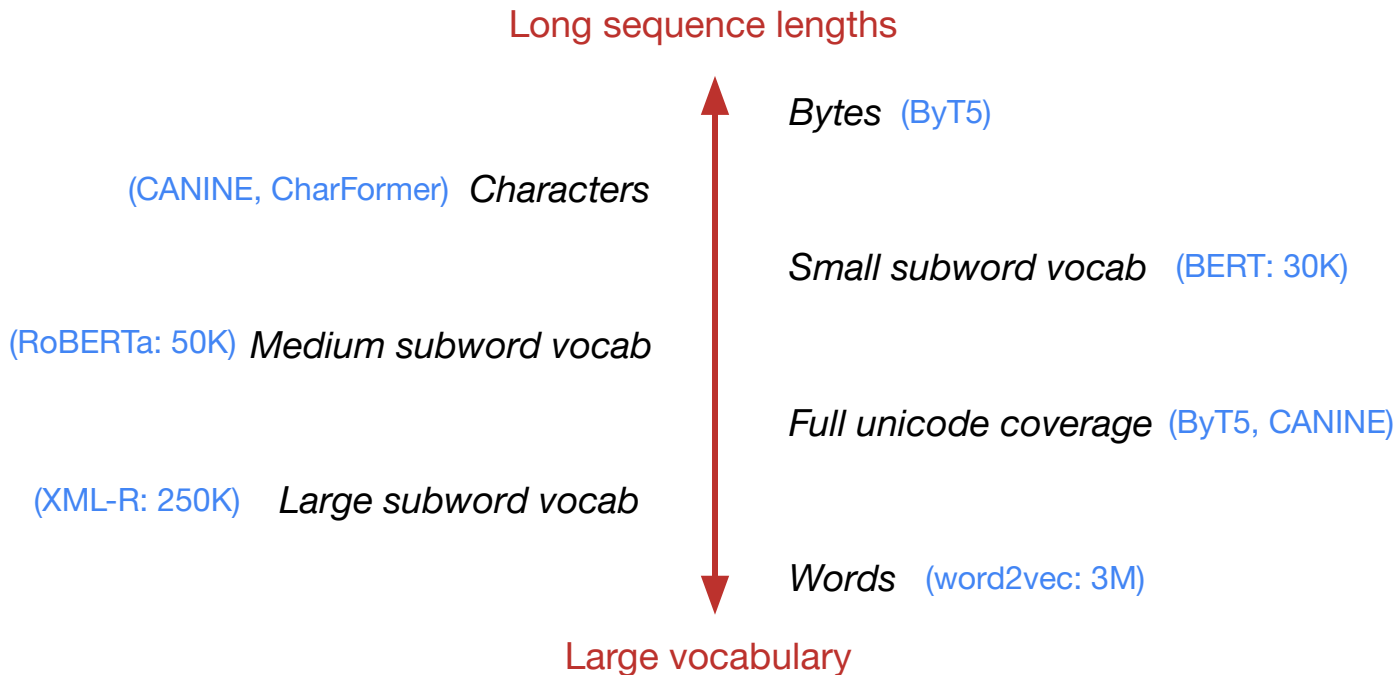
# The Vocabulary Bottleneck

- NLP is an **open vocabulary problem** and the ability of a model is determined by its vocabulary:
  1. tokens, characters, sub-words, etc.
- This creates a bottleneck in two places:
  1. *Representational bottleneck* in the Embedding layer
  2. *Computational bottleneck* in the Output layer



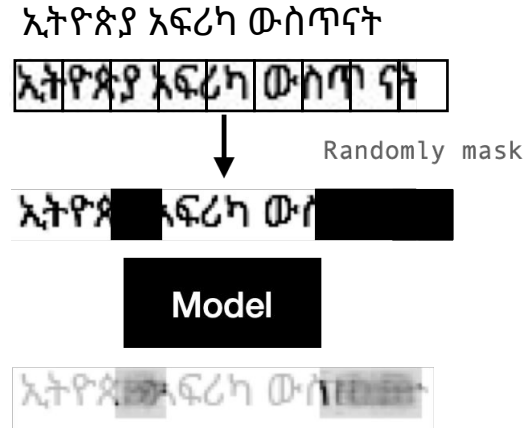
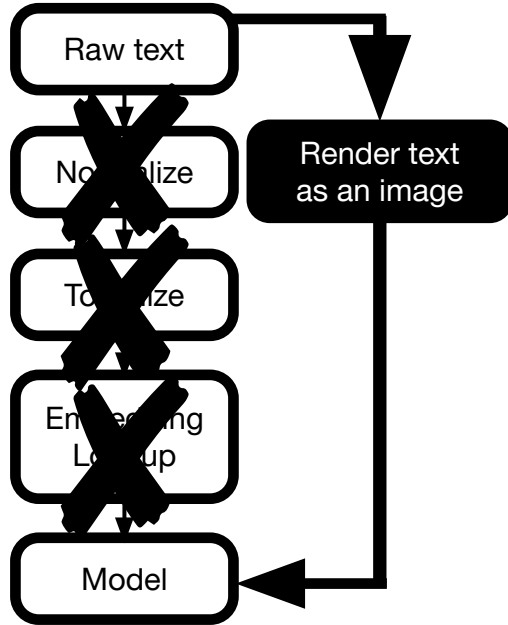
# Where's the sweet spot?

---

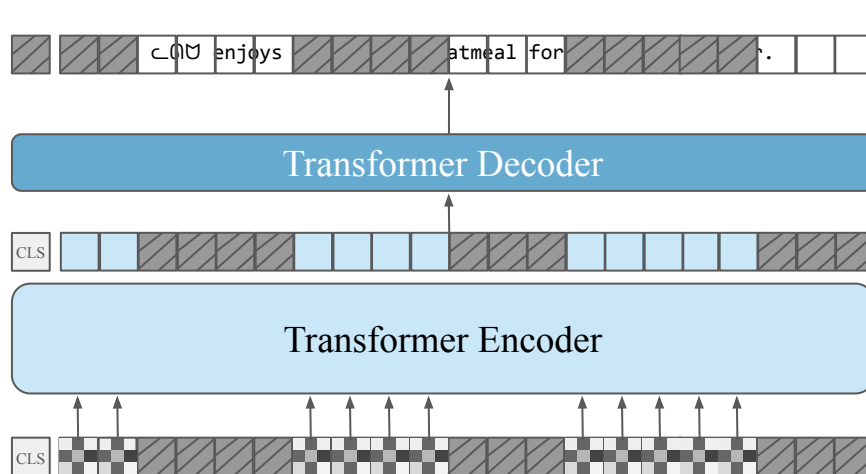




# Treat language as vision



# The Model



$$\frac{1}{M} \sum_{m \in M} (Y_m - \hat{Y}_m | \mathbf{x}_{\setminus m})^2$$

8 Layers

12 Layers

- 3 CLS Embedding & Span Mask  $m$  patches  
- 2 Projection + Position Embedding  

- 1 Render Text as Image



My cat enjoys eating warm oatmeal for lunch and dinner.

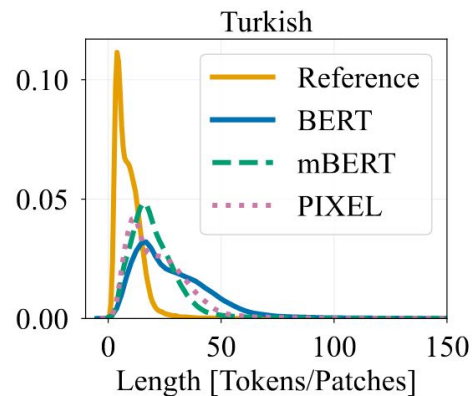
16pixel x 16pixel patch

Google Noto Fonts

PyGame / PangoCairo

# Rendered Text is Compact

- PIXEL encoding produces sequence lengths that are at least as long as as BERT.
  - Universal Dependencies datasets with human reference segmentations
  - No length penalty for any language, unlike some LLMs (Ahia et al. 2023)



# Pretraining

- **English Dataset:** English Wikipedia and Books Corpus
- **Masking:** 25% Span Masking
- **Maximum sequence length:** 529 patches (1 6x8 4 64 pixels)
- **Compute:** 8 x 40GB A100 GPUs for 8 days
- **Parameters:** 86M encoder + 26M decoder

There is only 0.05% non-English text in our pretraining data (estimated by Blevins and Zettlemoyer 2022)

The **Great Wall of China** (traditional Chinese: 萬里長城; simplified Chinese: 万里长城; pinyin: Wànlǐ Chángchéng)

# A new type of generative model

Penguins are designed to be streamlined and hydrodynamic, so **having long legs** would add expanding. Having short legs with wedged feet to act like **rubbers**, helps to give them that **top** do-like figure. **didn't** compare bird anatomy with humans, we would see something **not** peculiar. By taking a look at the side-by-side image in Figure 1, you can see how their leg **bones are** to ours. What most people mistake for **knees** are actually the **anatomies** of birds. This **gives a conclusion** that bird knees bend opposite of ours. The knees are actually tucked up inside the **bones** of the **bird**. So how does this look inside of a penguin? In the **images** below, you can see boxes surrounding the penguins' knees.



Penguins are designed to be streamlined and hydrodynamic, so **having long legs** would add expanding. Having short legs with wedged feet to act like **rubbers**, helps to give them that **top** do-like figure. **If** compare bird anatomy with humans, we would see something **not** peculiar. By taking a look at the side-by-side image in Figure 1, you can see how their leg **bones are** to ours. What most people mistake for **knees** are actually the **anatomies** of birds. This **gives the illusion** that bird knees bend opposite of ours. The knees are actually tucked up inside the **bones** of the **bird**. So how does this look inside of a penguin? In the **images** below, you can see boxes surrounding the penguins' knees.

100K steps

500K steps

1M steps

# Downstream Tasks

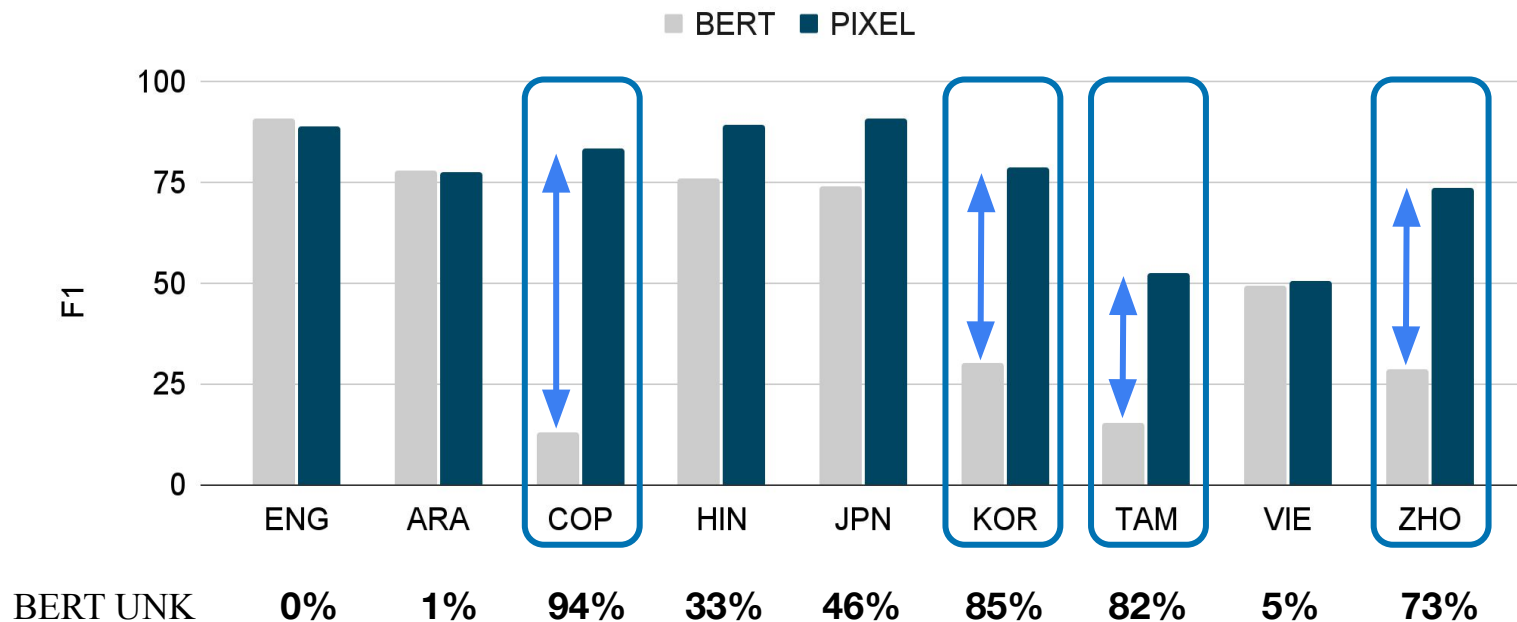
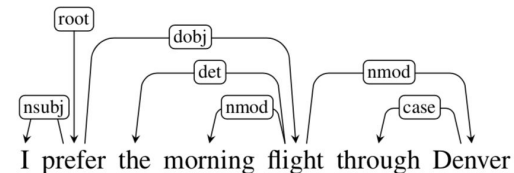
- **Datasets:** Universal Dependencies, MasakhaNER, GLUE, Zeroé
- **Models:**

	Parameters	Pretraining Data
PIXEL <sub>BASE</sub>	86M	English Wikipedia + Bookcorpus
BERT <sub>BASE</sub>	110M	—
CANINE-C	127M	104-languages from Wikipedia

Similar pretraining setup

Tries to solve the same problem using UTF-32

# Dependency Parsing Results

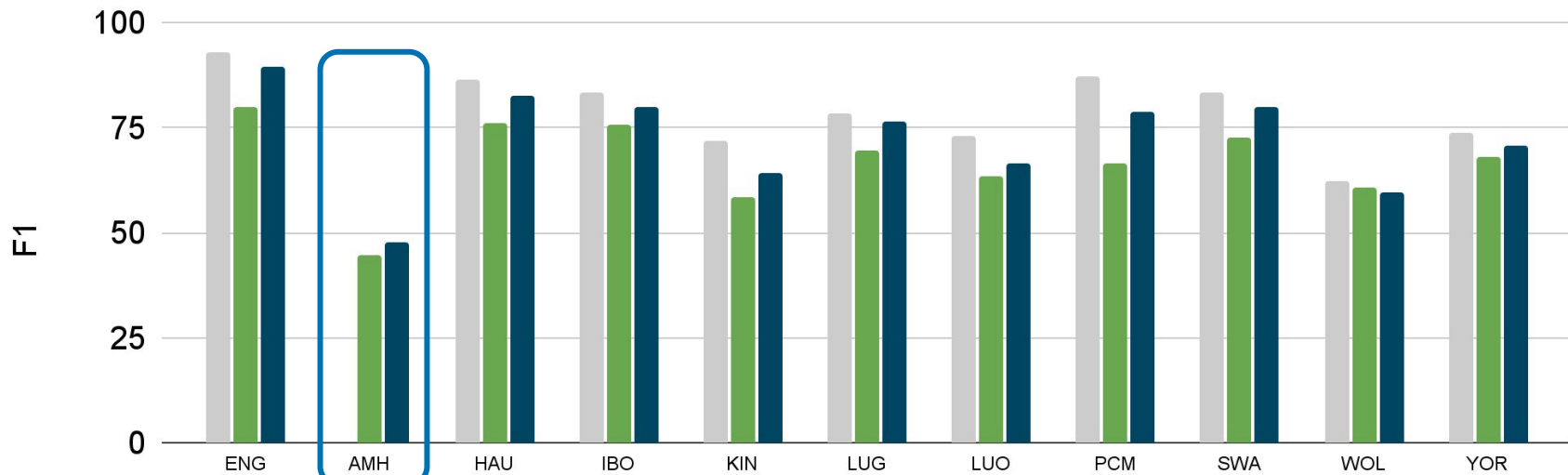


PIXEL vastly outperforms BERT on unseen scripts

# Named Entity Recognition in African Languages

Emir of Kano turban Zhang wey don spend 18 years for Nigeria

■ BERT ■ CANINE ■ PIXEL

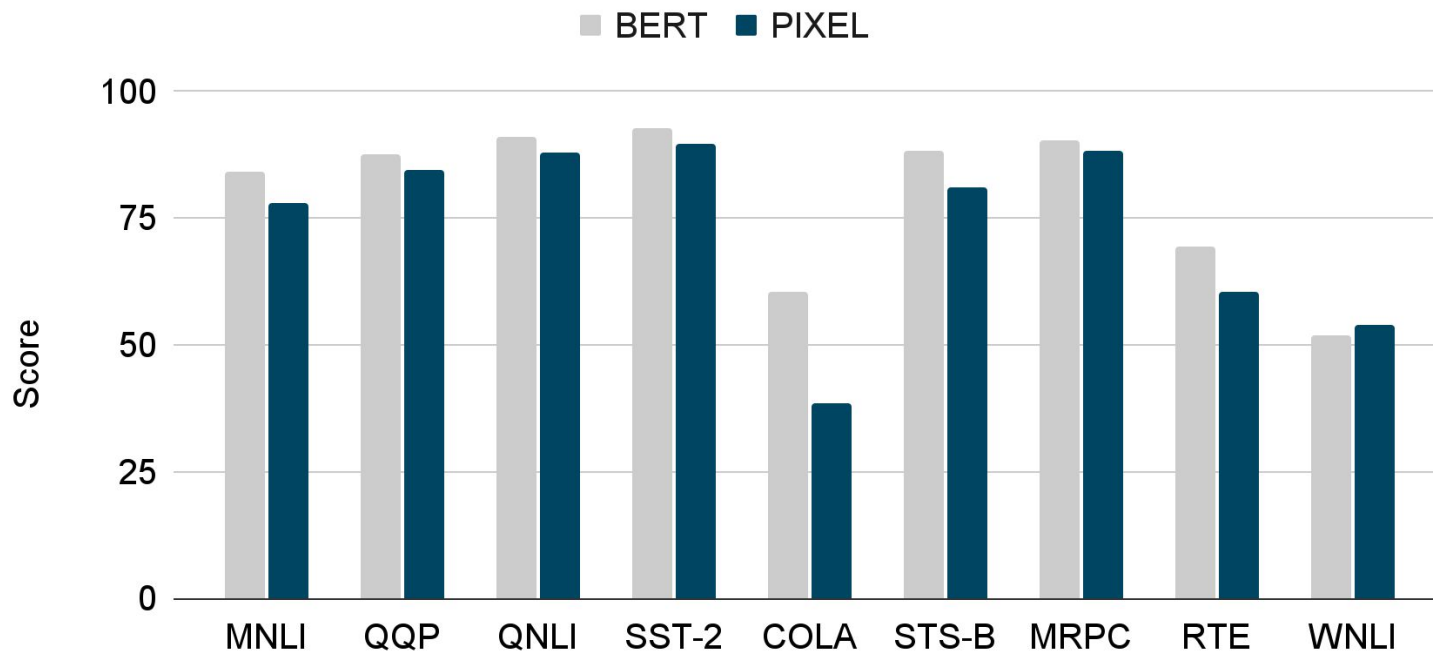


PIXEL outperforms BERT on the non-Latin script

PIXEL outperforms the multilingually pretrained CANINE-C



# GLUE: Sentence-level Understanding



BERT outperforms PIXEL on English sentence-level tasks

# Text Rendering Matters

---

- Our original text renderer produces many nearly-identical patches
  - This is representation- and compute-wasteful

the the the the the the the the the

Can we do better?

# Alternative Text Renderers

(a) Continuous rendering (CONTINUOUS):

I m u s t b e g r o w i n g s m a l l a g a i n . ■

(b) Structured rendering (BIGRAMS):

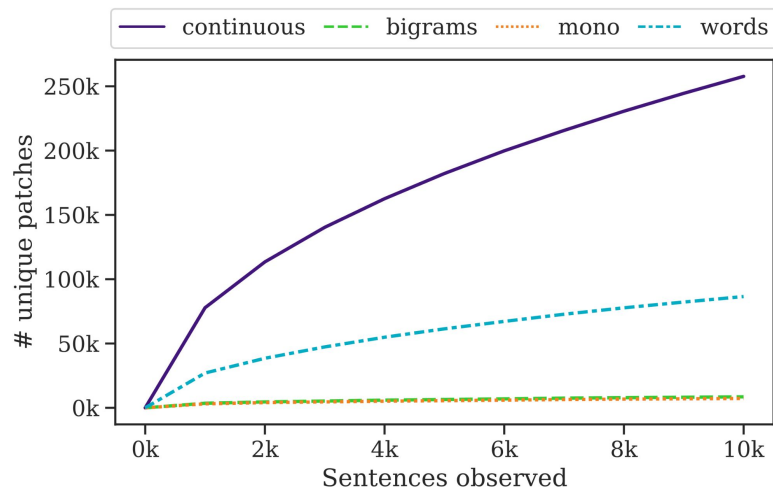
I m u s t b e g r o w i n g s m a l l a g a i n . ■

(c) Structured rendering (MONO):

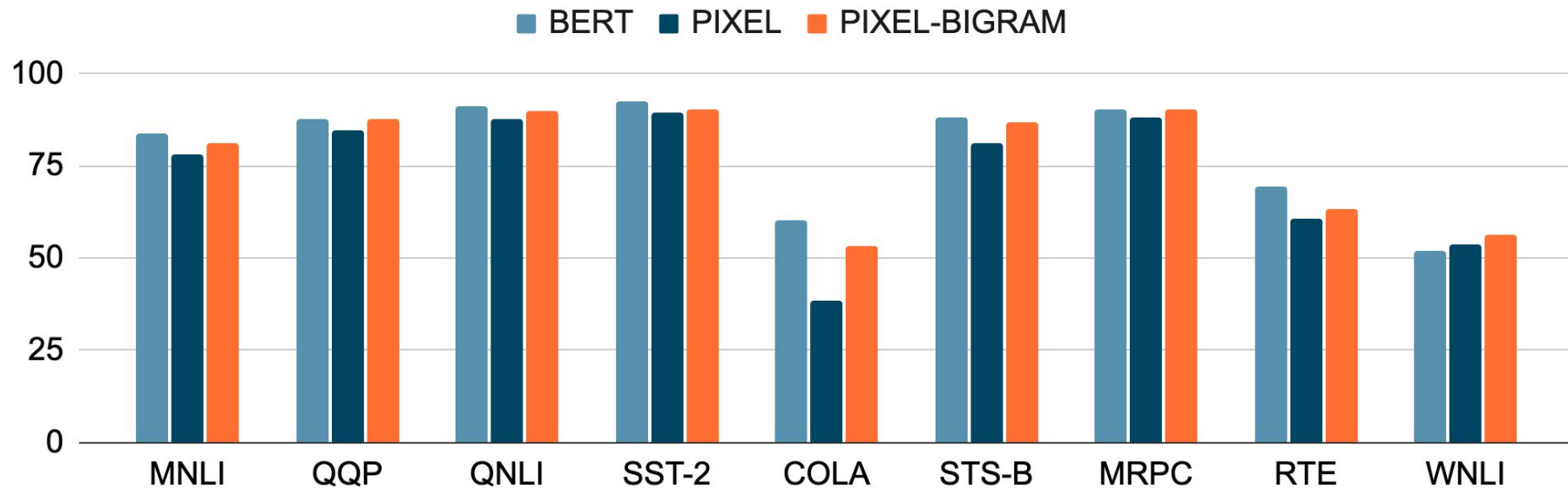
I m u s t b e g r o w i n g s m a l l a g a i n . ■

(d) Structured rendering (WORDS):

I m u s t b e g r o w i n g s m a l l a g a i n . ■



# GLUE (revisited)

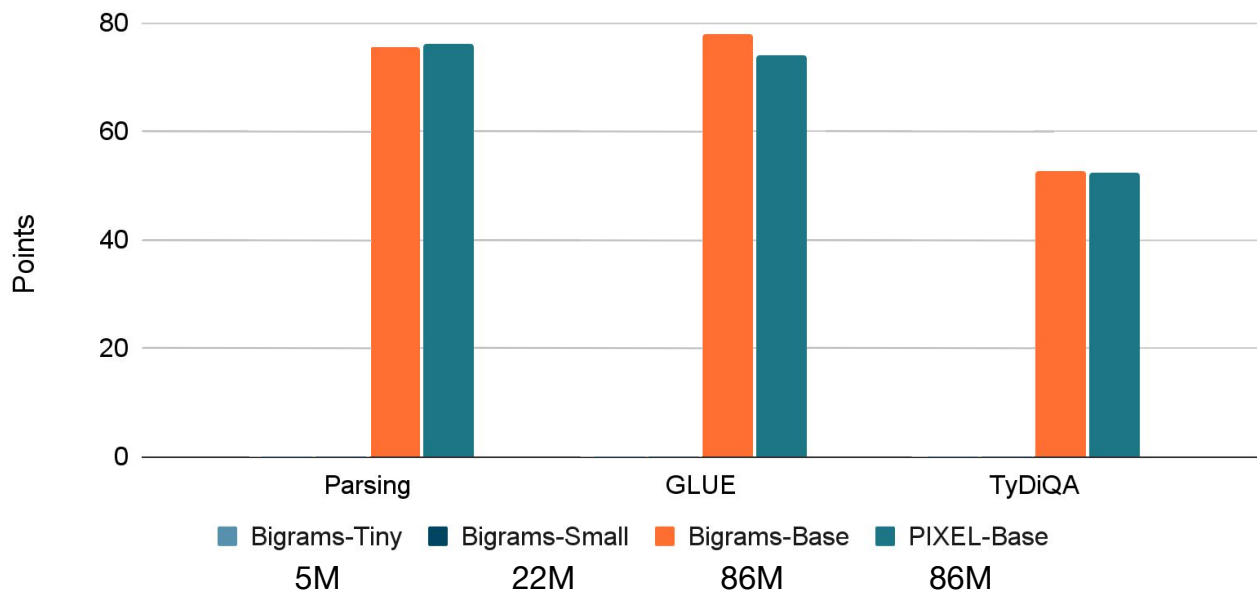


Bigram text rendering produces better models

# Scaling Down ↓

---

- Better text rendering can create effective models at smaller scales



# Open Questions

---

- Does this work based on orthographic similarity or something else?
- How should we train a multilingual PIXEL encoder?
  - Language-based or script-based data selection
- How can we apply this type of model to language generation tasks?

# Conclusions

- PIXEL is a new type of language model that tackles the open vocabulary problem using visually rendered text.
  1. This enables high-quality transfer to unseen scripts.
  2. Robustness to orthographic attacks
  3. Can also process historical documents
- *My opinion:* Language is special but its computer format should be as flexible and expressive as possible.

# Wrap-up

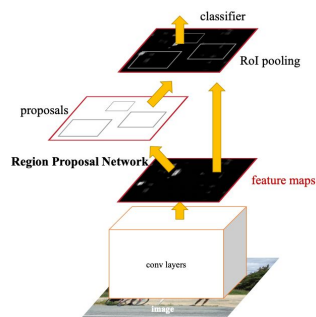


# 1. Datasets

some sheep walking in the middle of a road  
a herd of sheep with green markings walking down the road  
a herd of sheep walking down a street next to a lush green grass covered hillside.  
sheared sheep on roadway taken from vehicle, with green hillside in background.  
a flock of freshly sheared sheep in the road.



# 2. Representation

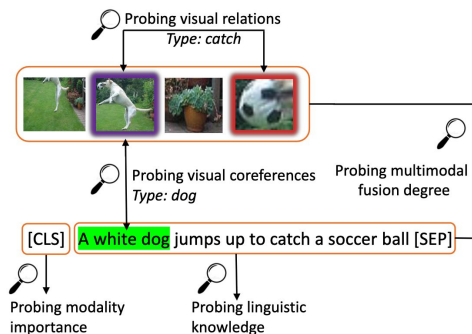


# 3. Modelling

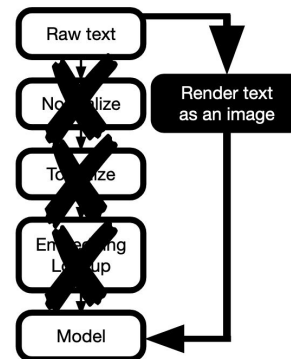
Transformer



The red horse



# 4. Understanding



# 5. New Ideas

# Where to find more research?

---

CVPR

NeurIPS

ACL

NAACL

LREC

ICML

EACL

ECCV

ICCV

IJCAI

ICLR

COLING

arXiv

...



DALL-E 2: “digital art of someone drinking from arxiv firehose every morning”

# Predictions & Speculations

---

- Increasing societal impact of V&L models
  - Both for entertainment and for misinformation
- Shift in focus to zero-shot instruction-based models
  - Fine-tuning is too expensive for each task
- Concentrated focus on understanding how models work
  - Bigger and better datasets will continue to be major contributions
- Big challenge to evaluate bidirectional generative models

# Acknowledgements

---



W. Li



R. Ramos



B. Martins



L. Cabello



S. Brandl



N. Borenstein



P. Rust



I. Augenstein



E. Bugliarello



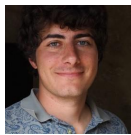
E. Salesky



M. de Lhoneux



F. Liu



E.M. Ponti



S. Reddy



N. Collier



Y. Kementchedjhieva

VILLUM FONDEN

