

# Multimodal LLMs

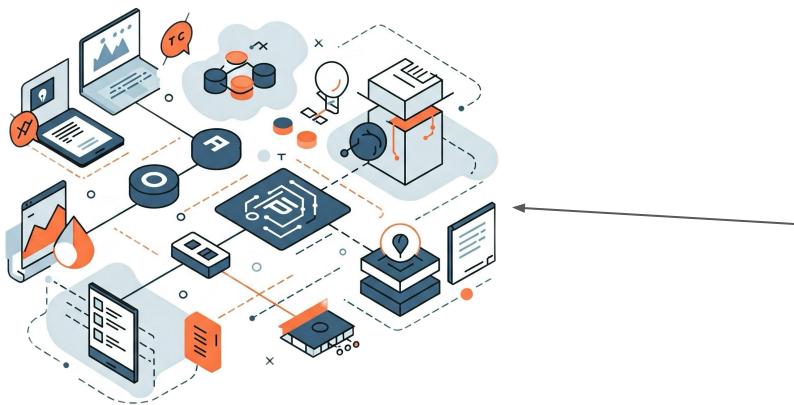


Chulalongkorn University (Online)  
April 21, 2025

Soravit “Beer” Changpinyo  
Google DeepMind

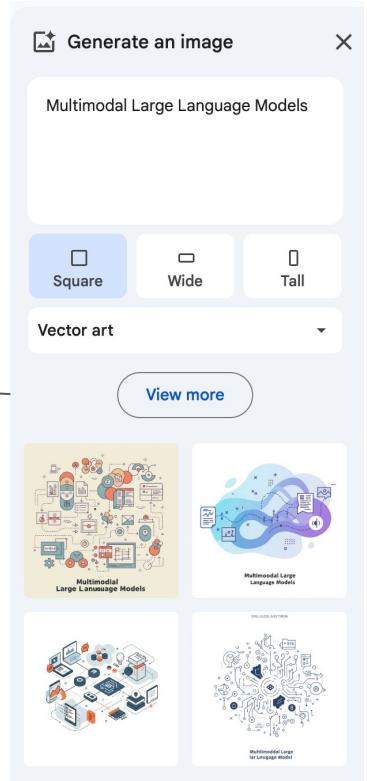
Opinions are my own  
Compiled from multiple sources with credit given in slides

# Multimodal LLMs



Chulalongkorn University (Online)  
April 21, 2025

Soravit “Beer” Changpinyo  
Google DeepMind



# **LLMs**

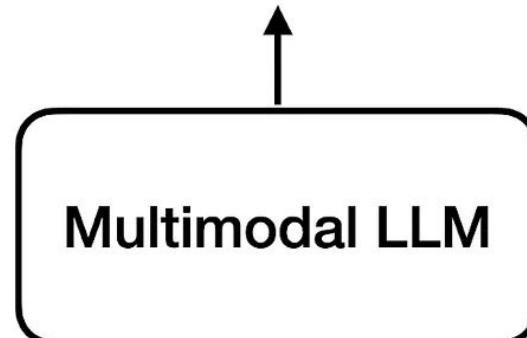
# Multimodal LLMs

# Multimodal LLMs possess

the ability to **process** and **generate**  
information across different data modalities  
(text, images, audio, video, etc.)

**Output:**

*The image depicts a rotary phone, which...*



**Input:**

Audio

*What does  
this...*



Text

Images



Videos

# OCR

Translate ▾

Google Docs

Download ▾

Copy

อยดุงคือชื่อของเทือกเขาสูงทางตอน

เหนือของจังหวัดเชียงราย ตั้นแต่น

สูงสุดยอดในสยาม

ณ เทือกเขาสูงชันแห่งนี้เป็นอาณา  
จักรของการผสมกลมกลืนระหว่างภูมิ  
อากาศที่หนาวเหน็บ และความงามตามของ  
ทิวเขียวสูงชันที่ลดหลั่นໄ่เรียงกันไปบน  
พื้นที่กว้าง闊มีไร้เขตอำเภอแม่จัน

อยดุงคือชื่อของเทือกเขาสูงทางตอน  
เหนือของจังหวัดเชียงราย ตั้นแต่น  
สูงสุดยอดในสยาม  
ณ เทือกเขาสูงชันแห่งนี้เป็นอาณา  
จักรของการผสมกลมกลืนระหว่างภูมิ  
อากาศที่หนาวเหน็บ และความงามตามของ  
ทิวเขียวสูงชันที่ลดหลั่นໄ่เรียงกันไปบน  
พื้นที่กว้าง闊มีไร้เขตอำเภอแม่จัน

# ASR



# Image Captioning & Visual QA

10-09-19

## Chrome's new AI feature solves one of the web's eternal problems

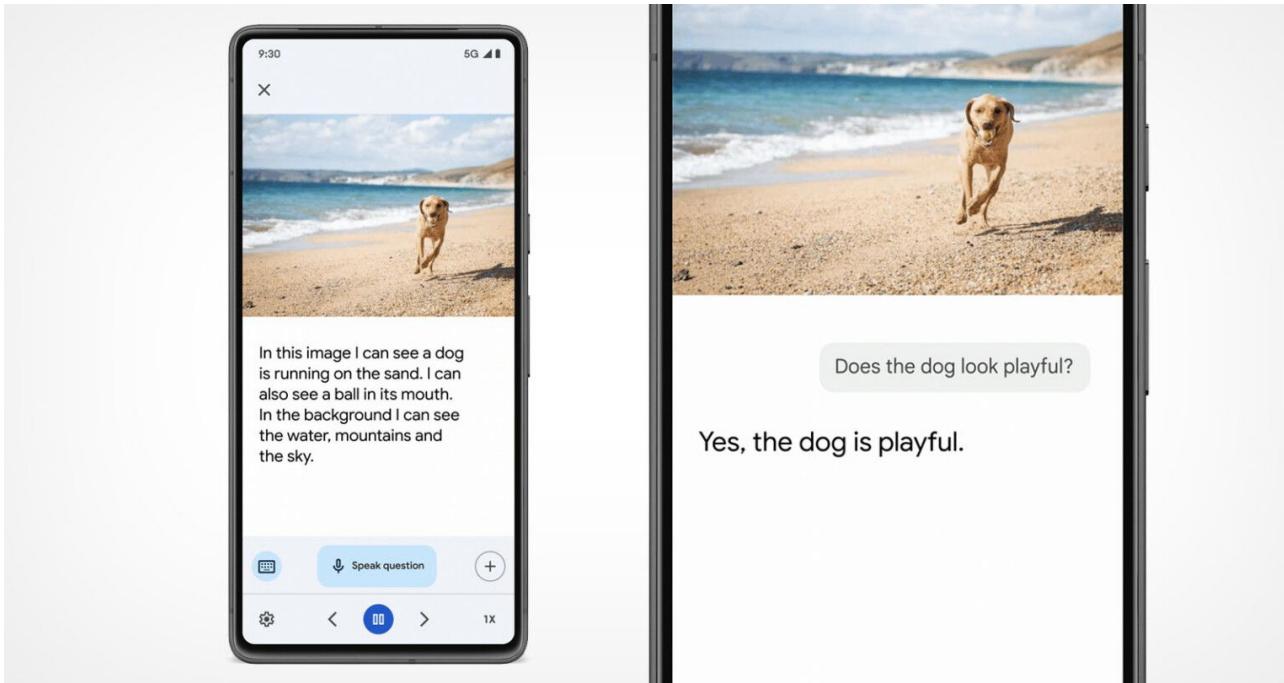
To help blind and low-vision users, Google is using machine learning to generate descriptions for millions of images.



# Google Can Now Describe and Answer Questions About Images

MAY 18, 2023

JARON SCHNEIDER



# ASL



<https://www.lifeprint.com/asl101/>

# ASL



<https://www.lifeprint.com/asl101/>

# Multimodal LLMs possess

the ability to **process** and **generate**  
information across different data modalities  
(text, images, audio, video, etc.)

**Generate an image of a dessert that looks like hair**



Generated by Gemini (Imagen 3)

# Image Generation

<text>



# Image Generation

Generate an image of  
the Songkran festival  
with people throwing  
mud instead of water



# Image Generation

<text>



# Image Generation

Make this happen in  
Paris, France instead



## A Problem (input)

You are given two strings  $s$  and  $t$ , both consisting of lowercase English letters. You are going to type the string  $s$  character by character, from the first character to the last one.

When typing a character, instead of pressing the button corresponding to it, you can press the 'Backspace' button. It deletes the last character you typed among those that aren't deleted yet (or does nothing if there are no characters in the current string). For example, if  $s$  is "abcbcd" and you press Backspace instead of typing the first and the fourth characters, you will get the string "bd" (the first press of Backspace deletes no character, and the second press deletes the character 'c'). Another example, if  $s$  is "abcaa" and you press Backspace instead of the last two letters, then the resulting text is "a".

Your task is to determine whether you can obtain the string  $t$ , if you type the string  $s$  and press 'Backspace' instead of typing several (maybe zero) characters of  $t$ .

### Input

The first line contains a single integer  $q$  ( $1 \leq q \leq 10^5$ ) — the number of test cases.

The first line of each test case contains the string  $s$  ( $1 \leq |s| \leq 10^5$ ). Each character of  $s$  is a lowercase English letter.

The second line of each test case contains the string  $t$  ( $1 \leq |t| \leq 10^5$ ). Each character of  $t$  is a lowercase English letter.

It is guaranteed that the total number of characters in the strings over all test cases does not exceed  $2 \cdot 10^5$ .

### Output

For each test case, print "YES" if you can obtain the string  $t$  by typing the string  $s$  and replacing some characters with presses of "Backspace" button, or "NO" if you cannot.

You may print each letter in any case (YES, yes, Yes will all be recognized as positive answer, NO, no and nO will all be recognized as negative answer).

### Example

Input	Output
4	
ababa	YES
ba	NO
ababa	NO
bb	YES
aaa	
aaaa	
aababa	
ababa	

### Note

Consider the example test from the statement.

In order to obtain "ba" from "ababa", you may press Backspace instead of typing the first and the fourth characters.

There's no way to obtain "bb" while typing "ababa".

There's no way to obtain "aaaa" while typing "aaa".

In order to obtain "aababa" while typing "aababa", you have to press Backspace instead of typing the first character, then type all the remaining characters.

## AlphaCode



## B Solution (output)

```
t=int(input())
for i in range(t):
    s=input()
    t=input()
    a=[]
    b=[]
    for j in s:
        a.append(j)
    for j in t:
        b.append(j)
    a.reverse()
    b.reverse()
    c=[]
    while len(b)!=0 and len(a)!=0:
        if a[0]==b[0]:
            c.append(b.pop(0))
            a.pop(0)
        elif a[0]!=b[0] and len(a)!=1:
            a.pop(0)
            a.pop(0)
        elif a[0]==b[0] and len(a)==1:
            a.pop(0)
            if len(b)==0:
                print("YES")
            else:
                print("NO")
```

First the solution reads the two phrases.

If the letters at the end of both phrases don't match, the last letter must be deleted. If they do match we can move onto the second last letter and repeat.

If we've matched every letter, it's possible and we output that.

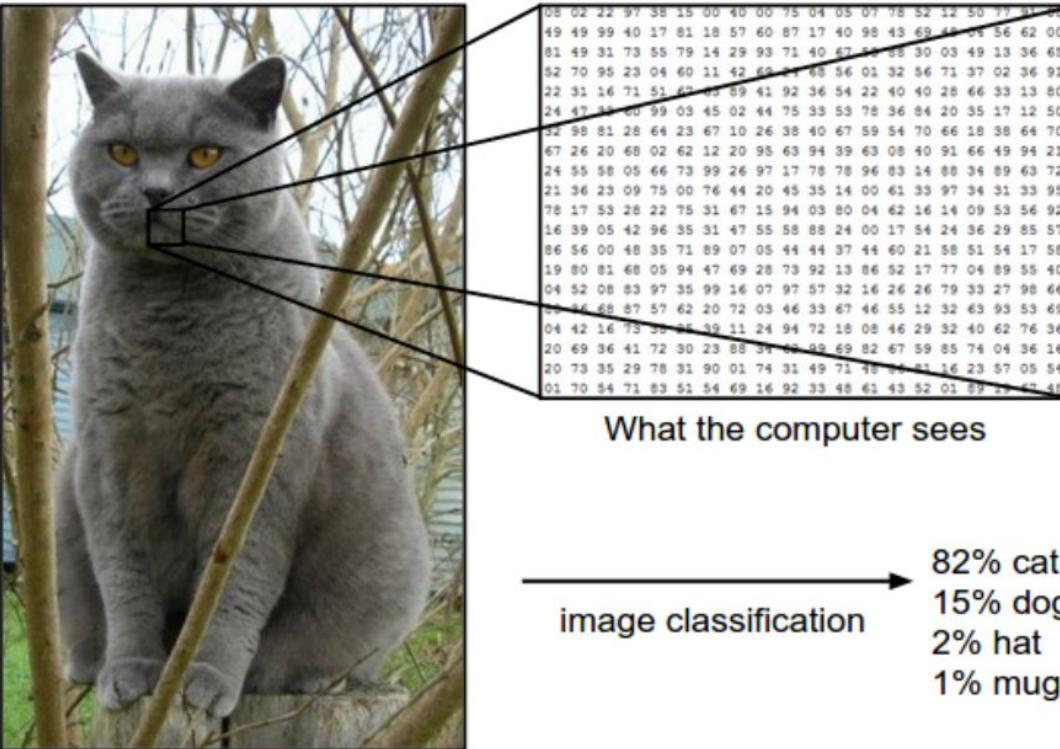
Backspace deletes two letters. The letter you press backspace instead of, and the letter before it.



**This Lecture:**  
**image understanding focus**  
**(where much attention is)**  
**(what I am more aware of)**

# Outline

- Fundamentals & Evolution
- Multimodal LLMs: Modeling, Data, Evaluation
- Example: Gemini



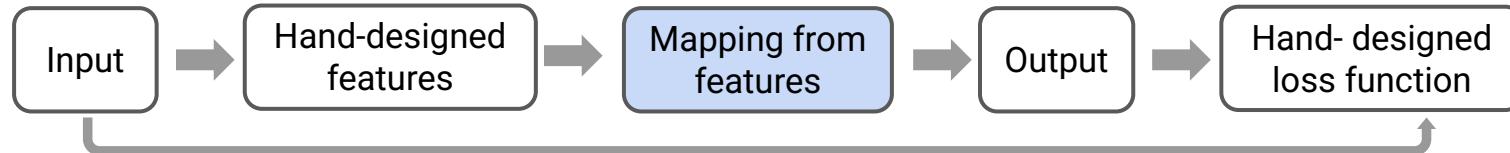
[Image credit: Andrej Karpathy]

## Rule-based systems



Learnable part of  
the system

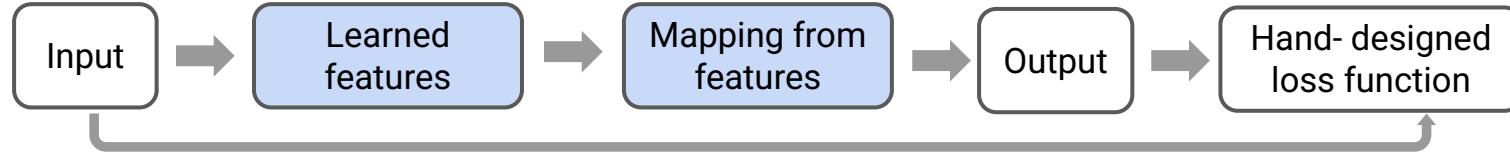
## Classical machine learning



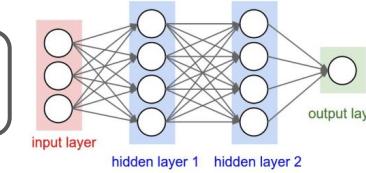
logistic regression



## Deep learning: (self-)supervised learning



Feedforward neural net



[PDF] Learning image attributes using the Indian Buffet Process

[PDF] brown.edu

S Changpinyo, E Sudderth

2012 • static.cs.brown.edu

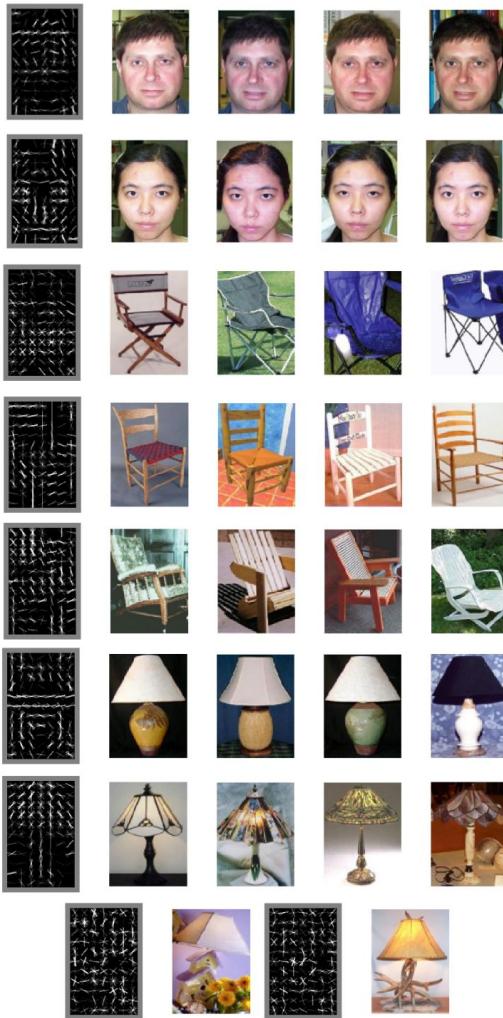
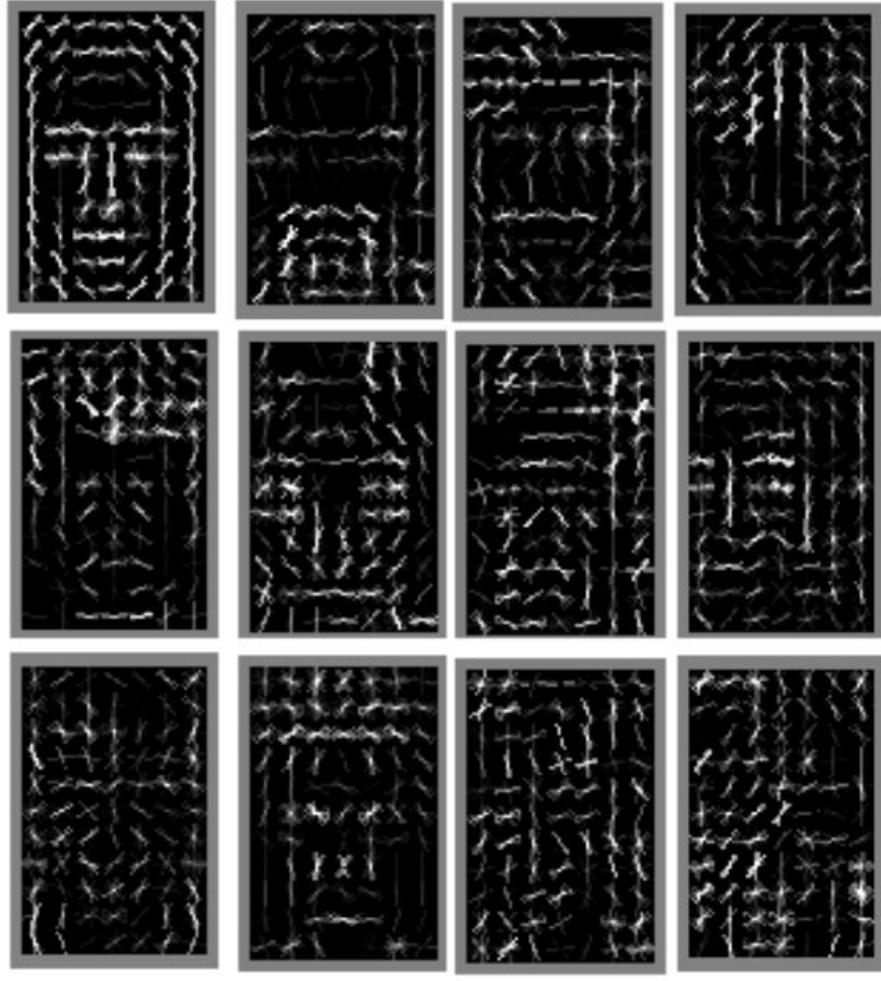
## Abstract

In the domain of object recognition and image classification, a recent trend is to use image properties or attributes to represent the images. Most of the proposed models in the past require that the number of attributes and attribute semantics be specified in advance. In this paper, we propose a generative model for image attributes that combine attribute-based vision models and feature-based nonparametric models. We learn the model using Gibbs sampling. Qualitatively, we demonstrate the learned attributes of images in three categories. Quantitatively, we show that our model outperforms simple baseline methods in image retrieval and transfer learning tasks.

static.cs.brown.edu

SHOW LESS ^

☆ Save ⚡ Cite Cited by 10 Related articles All 6 versions ☰



**2012**

**0.28**



**35**

**2010**



**2011**



**29**

**2012**



**2013**

**123**

**2014**

**157**

**2015**

**172**

**0.03**

**2016**



**Number of  
Entries**



**Classification  
Errors (top-5)**

## **SUN, 131K**

[Xiao et al. '10]

## **LabelMe, 37K**

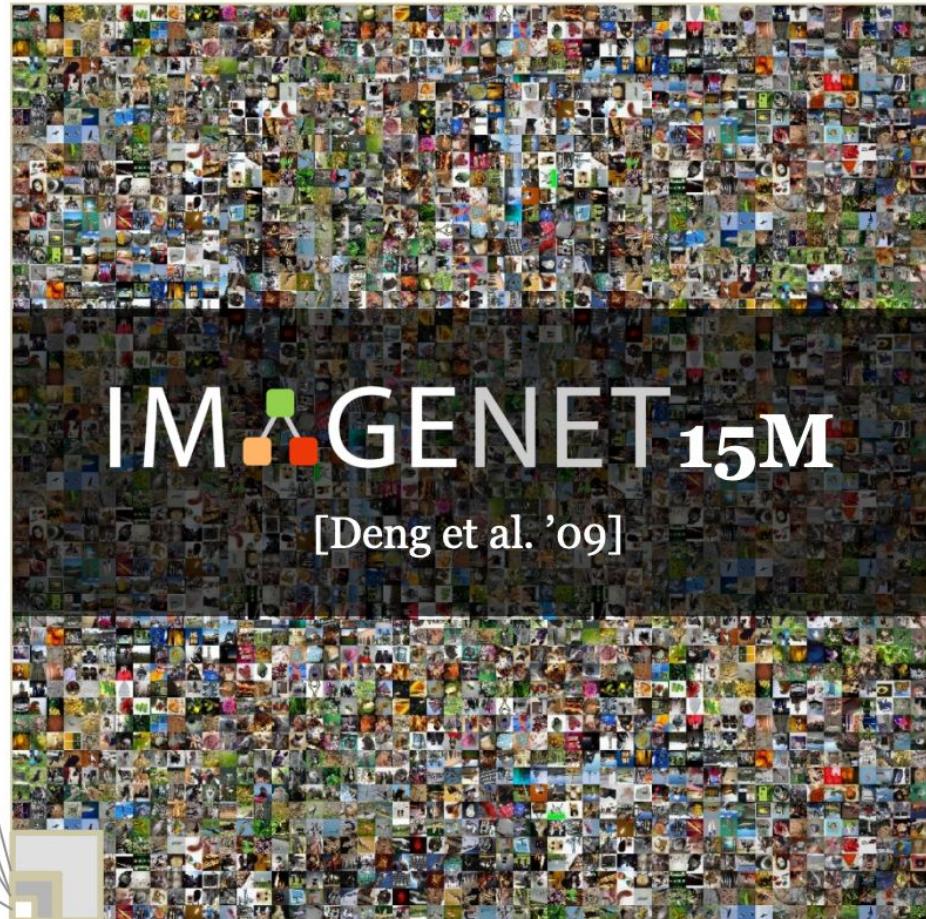
[Russell et al. '07]

## **PASCAL VOC, 30K**

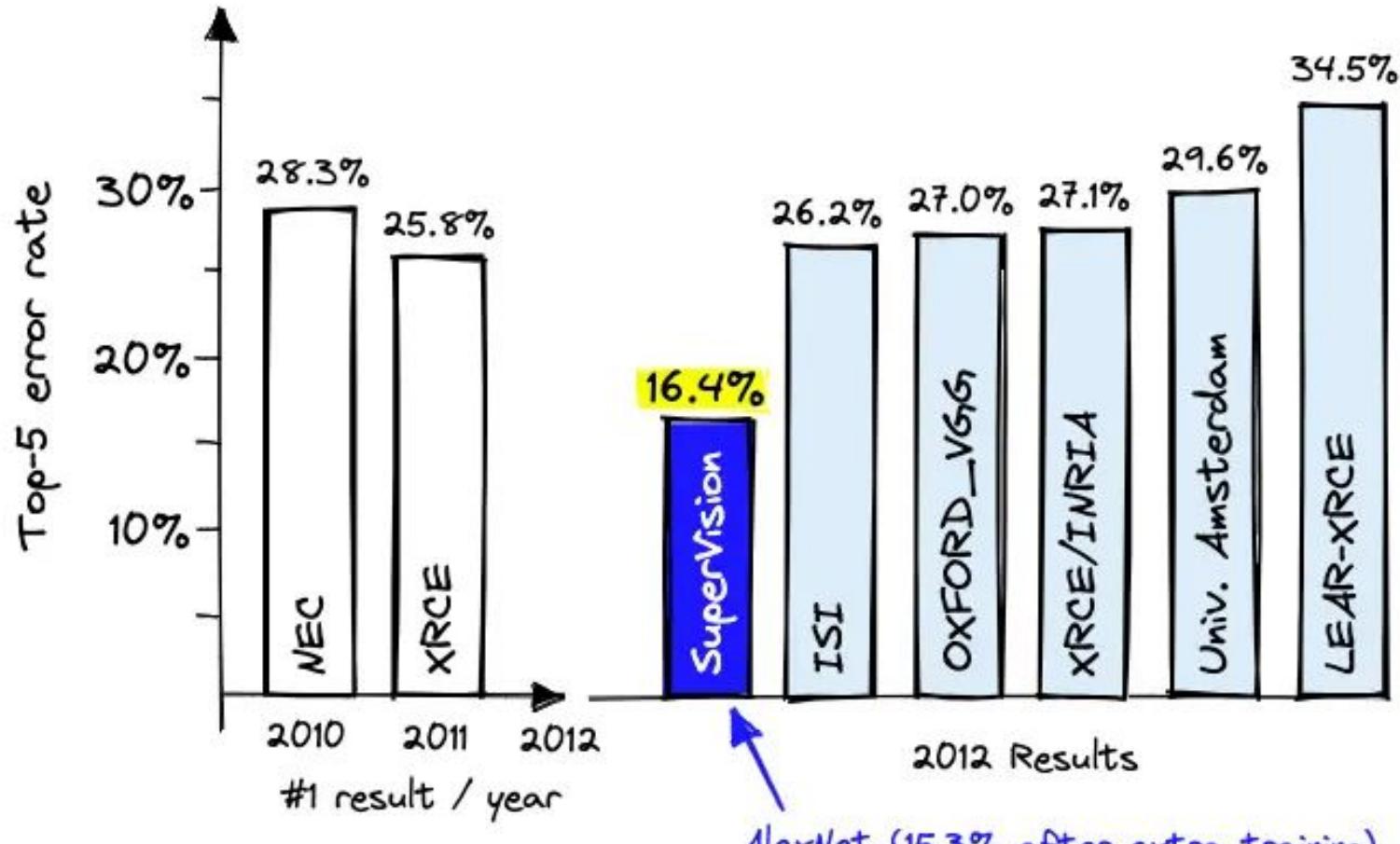
[Everingham et al. '06-'12]

## **Caltech101, 9K**

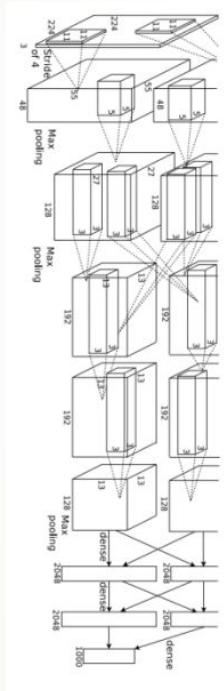
[Fei-Fei, Fergus, Perona, '03]



[https://image-net.org/static\\_files/files/imagenet\\_ilsvrc2017\\_v1.0.pdf](https://image-net.org/static_files/files/imagenet_ilsvrc2017_v1.0.pdf)

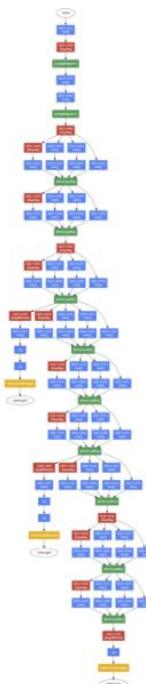


# “AlexNet”



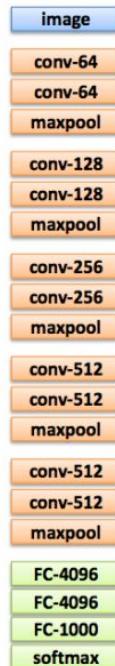
[Krizhevsky et al. NIPS 2012]

# “GoogLeNet”



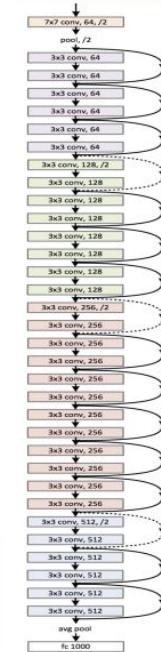
[Szegedy et al. CVPR 2015]

# “VGG Net”



[Simonyan & Zisserman,  
ICLR 2015]

# “ResNet”



[He et al. CVPR 2016]

## Imagenet classification with deep convolutional neural networks

A Krizhevsky, I Sutskever... - Advances in neural ..., 2012 - proceedings.neurips.cc

... We trained a large, **deep convolutional** neural **network** to **classify** the 1.2 million high-resolution images in the **ImageNet** LSVRC-2010 contest into the 1000 different classes. On the test ...

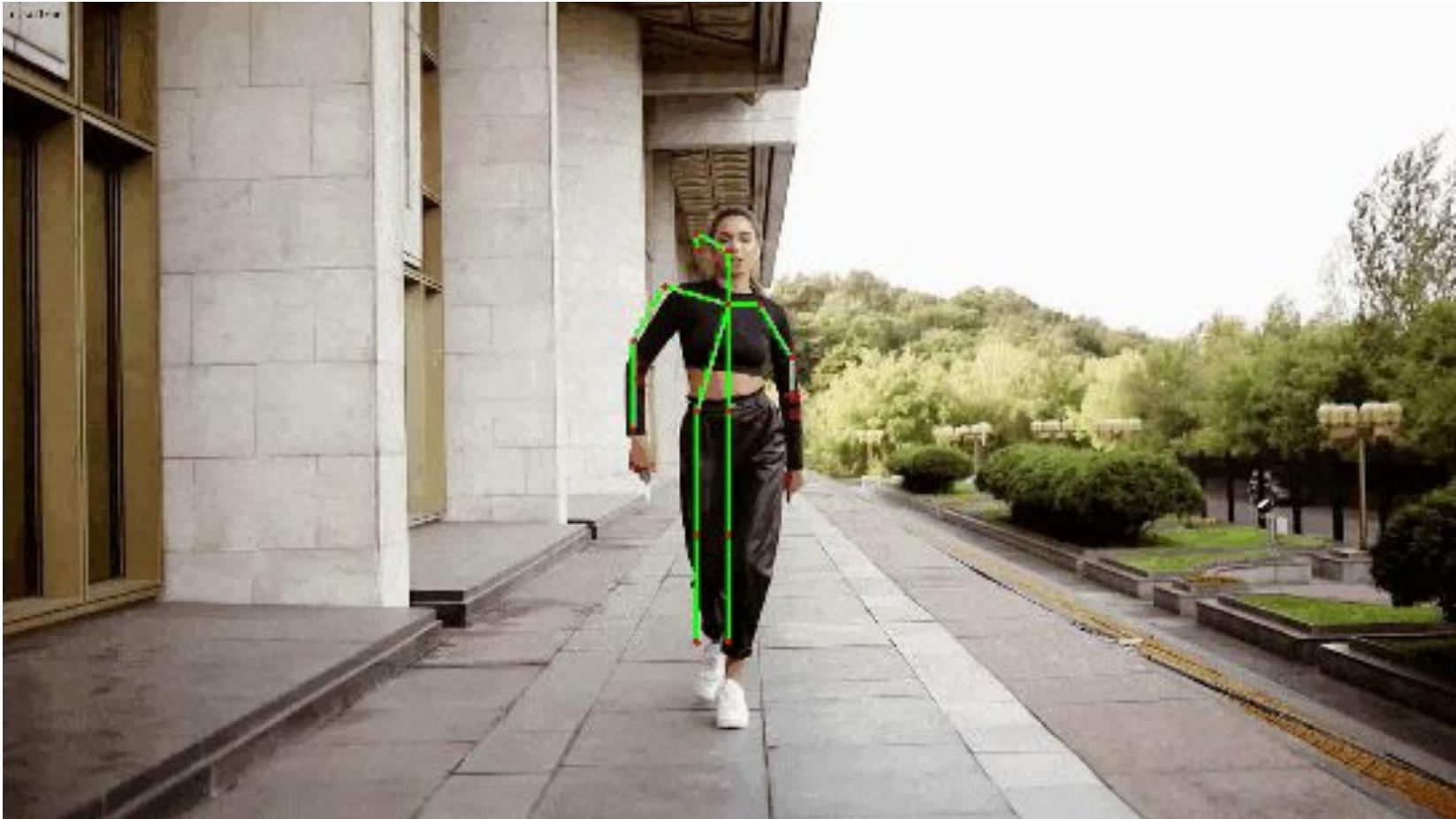
☆ Save ⚡ Cite Cited by 142407 Related articles All 89 versions ☰

## Deep residual learning for image recognition

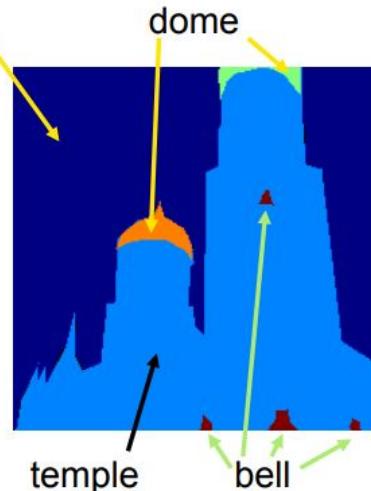
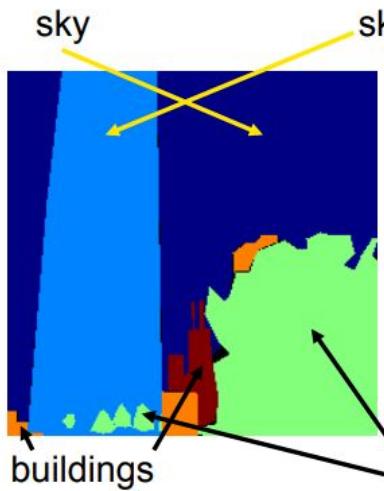
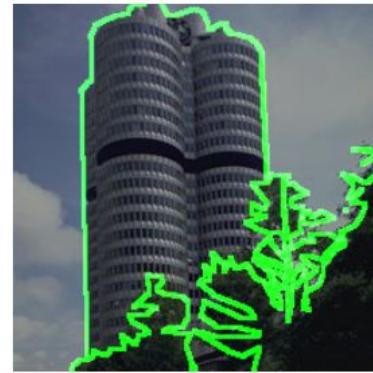
K He, X Zhang, S Ren, J Sun - ... and pattern recognition, 2016 - openaccess.thecvf.com

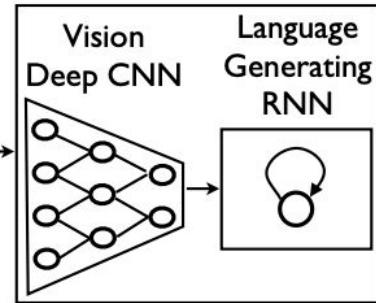
... **Deeper** neural **networks** are more difficult to train. We present a **residual learning** framework to ease the training of **networks** that are substantially **deeper** than those used previously. ...

☆ Save ⚡ Cite Cited by 265023 Related articles All 53 versions ☰



From [https://cs.brown.edu/courses/cs195-5/spring2011/lectures/2011-01-27\\_overview.pdf](https://cs.brown.edu/courses/cs195-5/spring2011/lectures/2011-01-27_overview.pdf)





A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

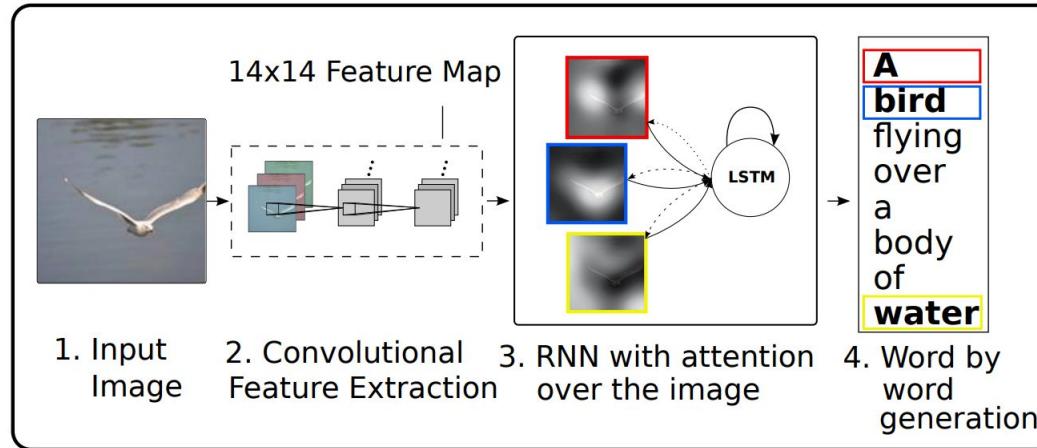
## Show and tell: A neural image caption generator

[O Vinyals, A Toshev, S Bengio... - Proceedings of the IEEE ..., 2015 - cv-foundation.org](#)

Automatically describing the content of an image is a fundamental problem in artificial intelligence that connects computer vision and natural language processing. In this paper, we ...

[☆ Save](#) [⤙ Cite](#) [Cited by 8225](#) [Related articles](#) [All 26 versions](#) [⤚⤚](#)

*Figure 1.* Our model learns a words/image alignment. The visualized attentional maps (3) are explained in Sections 3.1 & 5.4



## Show, attend and tell: Neural image caption generation with visual attention

[K Xu, J Ba, R Kiros, K Cho, A Courville...](#) - International ..., 2015 - proceedings.mlr.press

... to the task at hand, **and** we **show** how learning to **attend** at different locations in order to ...  
attention mechanism **and** a “soft” deterministic attention mechanism. We also **show** how one ...

☆ Save ⚡ Cite Cited by 13429 Related articles All 24 versions ➞

# Aside: My paper in 2013

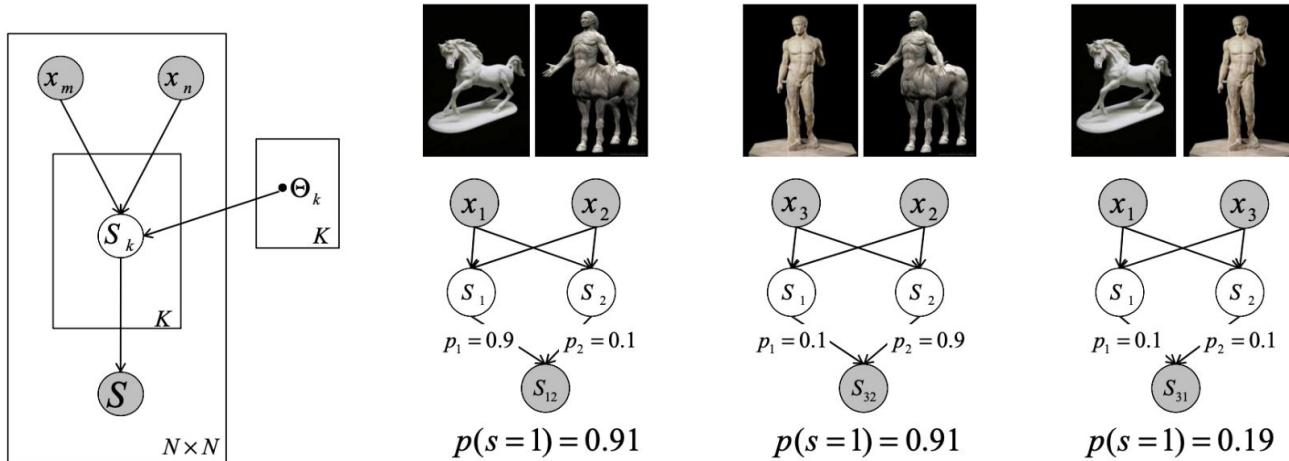


Figure 1: Similarity Component Analysis and its application to the example of CENTAUR, MAN and HORSE. SCA has  $K$  latent components which give rise to local similarity values  $s_k$  conditioned on a pair of data  $x_m$  and  $x_n$ . The model's output  $s$  is a combination of all local values through an OR model (straightforward to extend to a noisy-OR model).  $\Theta_k$  is the parameter vector for  $p(s_k|x_m, x_n)$ . See texts for details.

# Aside: My paper in 2016

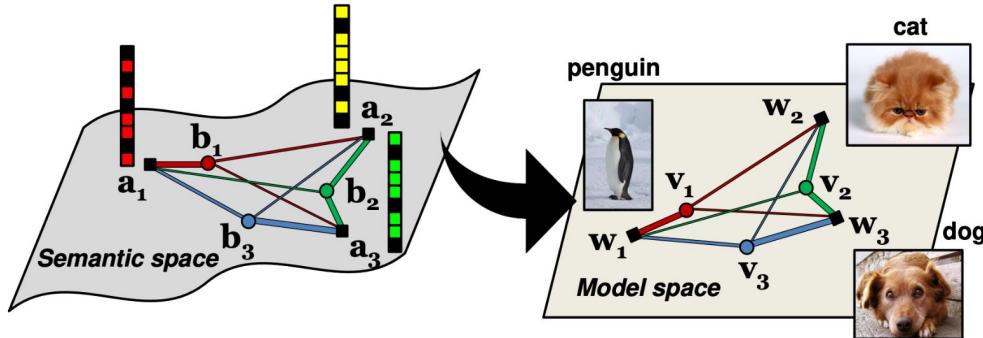


Figure 1: Illustration of our method for zero-shot learning. Object classes live in two spaces. They are characterized in the semantic space with semantic embeddings ( $as$ ) such as attributes and word vectors of their names. They are also represented as models for visual recognition ( $ws$ ) in the model space. In both spaces, those classes form weighted graphs. The main idea behind our approach is that these two spaces should be aligned. In particular, the coordinates in the model space should be the projection of the graph vertices from the semantic space to the model space — preserving class relatedness encoded in phantom classes ( $b$  and  $v$ ) to connect seen and unseen classes — classifiers for the phantom classes for real classes. In particular, the synthesis takes the form of convex combination.

**Advantage of deep features** It is also clear from Table 4 that, across all methods, deep features significantly boost the performance based on shallow features. We use GoogLeNet and AlexNet (numbers in parentheses) and GoogLeNet generally outperforms AlexNet. It is worthwhile to point out that the reported results under

# Aside: My paper in 2017

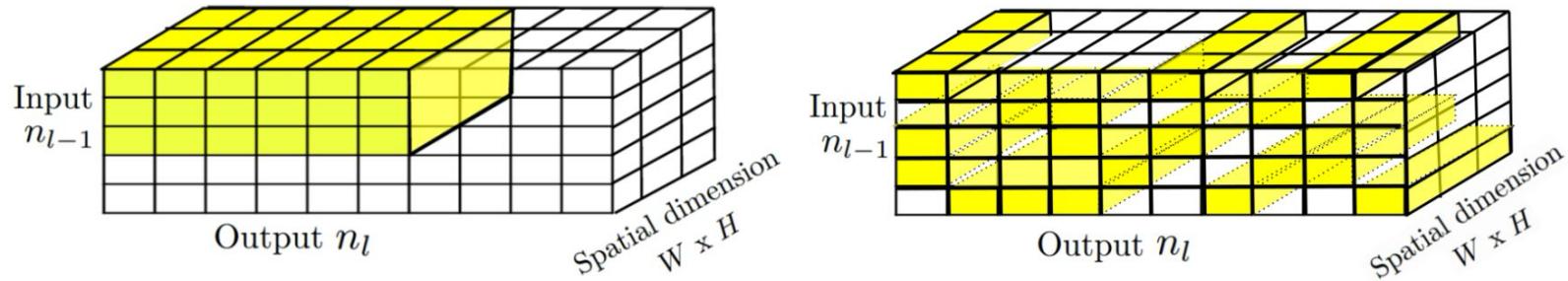
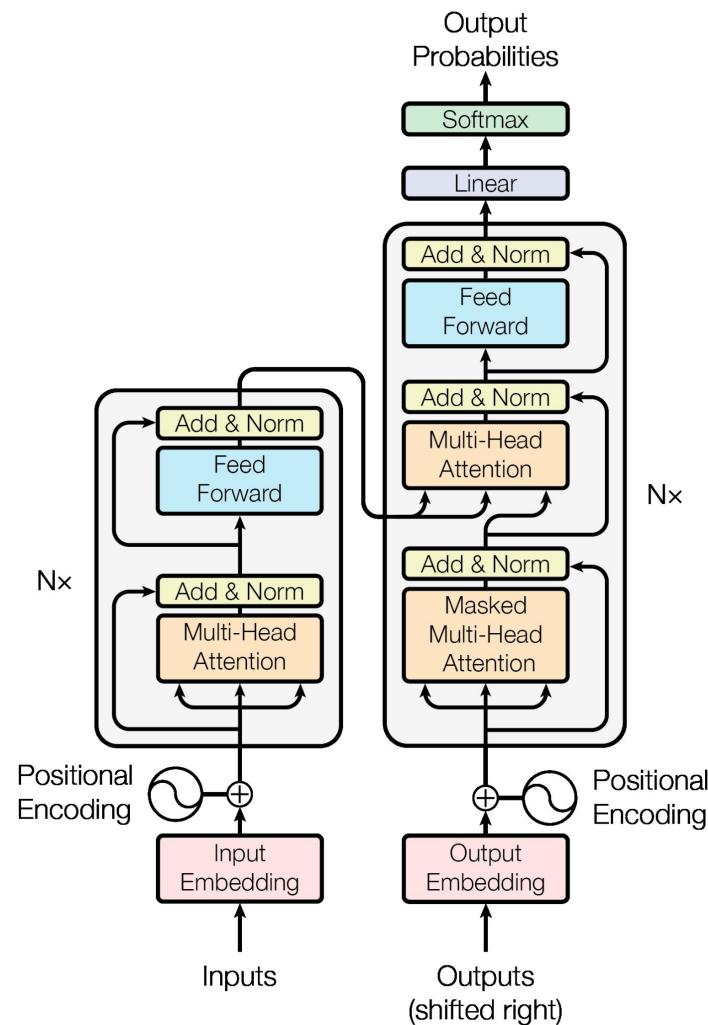


Figure 1: Connection tensors of depth multiplier (left) and sparse random (right) approaches for  $n_{l-1} = 5$  and  $n_l = 10$ . Yellow denotes active connections. For both approaches, the connection pattern is the same across spatial dimension and fixed before training. However, in the sparse random approach, each output channel is connected to a (possibly) different subset of input channels, and vice versa.

**2017**

**What's the most significant event in  
the AI community in 2017?**



# Attention is all you need

[A Vaswani, N Shazeer, N Parmar... - Advances in neural ...](#), 2017 - proceedings.neurips.cc

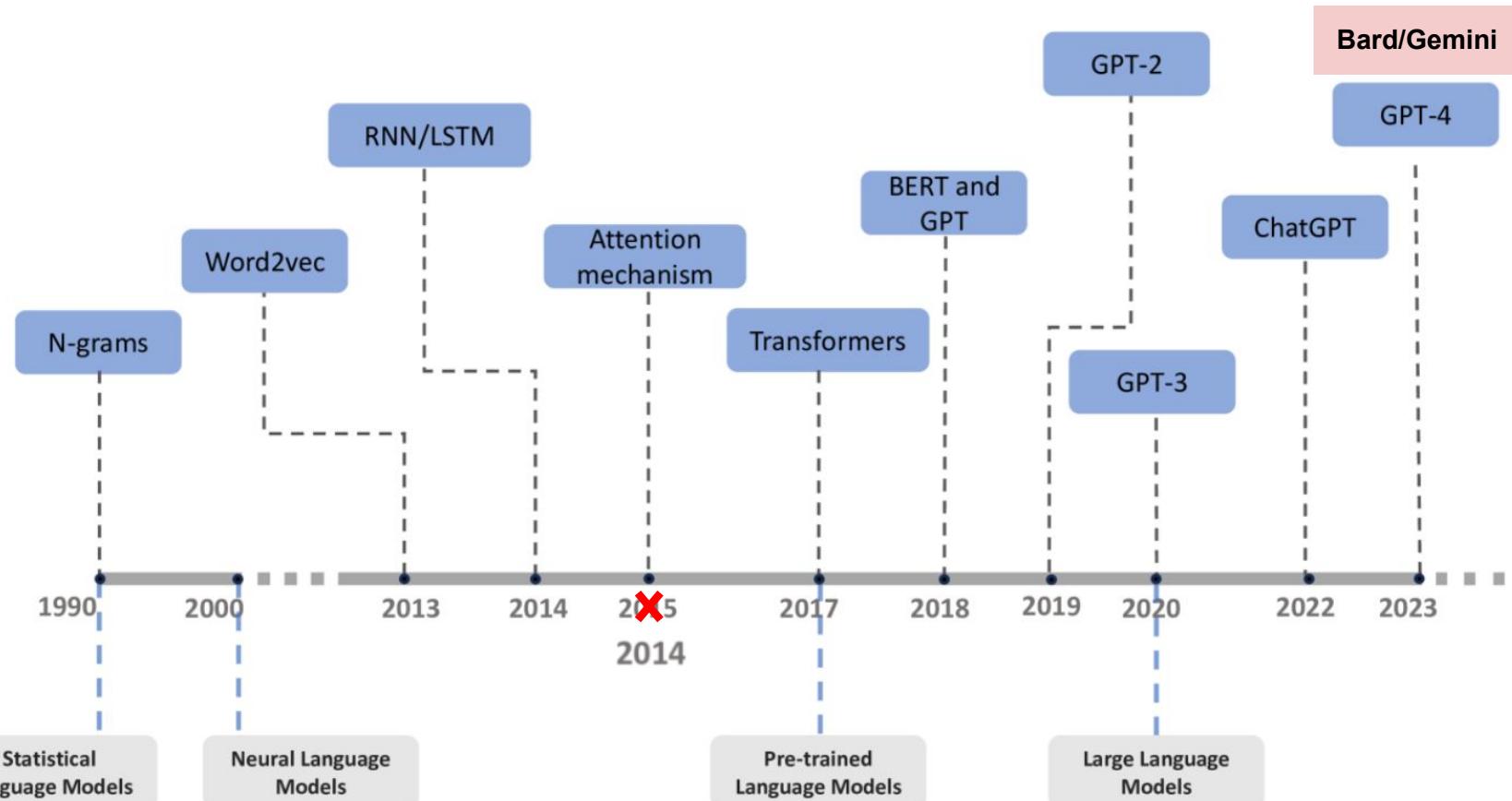
... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent

... **We** implement this inside of scaled dot-product **attention** by masking out (setting to  $-\infty$ ) ...

 Save  Cite Cited by 176713 Related articles All 73 versions 

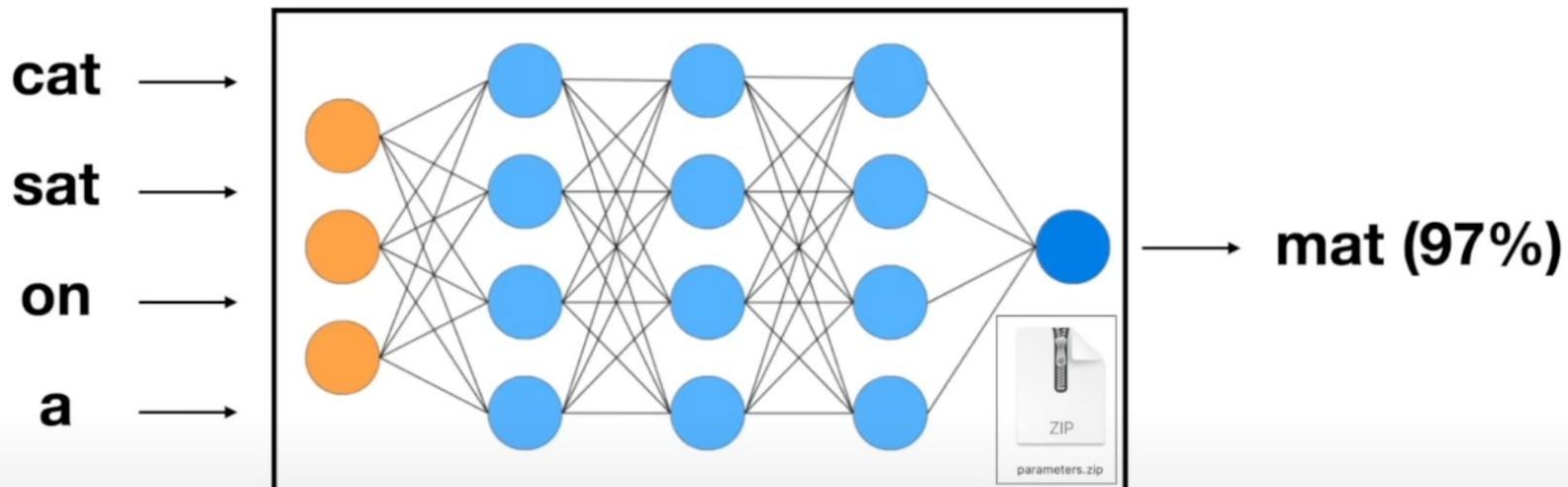
## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.



# Neural Network

Predicts the next word in the sequence.



e.g. context of 4 words

predict next word

Slide Credit: [Andrej Karpathy](#)

# Next word prediction forces the neural network to learn a lot about the world:

Ruth Marianna Handler (née Mosko; November 4, 1916 – April 27, 2002) was an American businesswoman and inventor. She is best known for inventing the Barbie doll in 1959,<sup>[2]</sup> and being co-founder of toy manufacturer Mattel with her husband Elliot, as well as serving as the company's first president from 1945 to 1975.<sup>[3]</sup>

The Handlers were forced to resign from Mattel in 1975 after the Securities and Exchange Commission investigated the company for falsifying financial documents.<sup>[3][4]</sup>

## Early life [ edit ]

Ruth Marianna Mosko<sup>[5][2][3]</sup> was born on November 4, 1916, in Denver, Colorado, to Polish-Jewish immigrants Jacob Moskowicz, a blacksmith, and Ida Moskowicz, née Rubenstein.<sup>[6]</sup>

She married her high school boyfriend, Elliot Handler, and moved to Los Angeles in 1938, where she found work at Paramount.<sup>[7]</sup>

Ruth Handler	
	
Born	Ruth Marianna Mosko November 4, 1916 Denver, Colorado, U.S.
Died	April 27, 2002 (aged 85) <sup>[1]</sup> Los Angeles, California, U.S.

## Process

"Many words don't map to one token: indivisible."

## Shape

[ ]

## Process

"Many words don't map to one token: indivisible."

↓ Tokenization

Unicode characters like emojis may be split.

[7085, 2456, 836, 470, 3975, 284, 530, 11241, 25, 773, 452, 12843, 13]

## Shape

[ ]

↓

[length]

## Process

"Many words don't map to one token: indivisible."

↓ Tokenization

Unicode characters like emojis may be split.

[7085, 2456, 836, 470, 3975, 284, 530, 11241, 25, 773, 452, 12843, 13]

↓ Embedding

2.3	-3.2	8.3	5.4	2.1	3.9	-8.9	3.8	3.9	3.3
4.5	5.9	4.5	7.1	1.0	5.3	5.0	3.1	0.7	5.0
...	...	...	...	...	...	...	...	...	...
3.8	1.2	3.8	9.0	9.3	3.1	4.2	0.8	9.2	5.8

## Shape

[ ]

↓

[length]

↓

[d\_model, length]

## Process

"Many words don't map to one token: indivisible."



Unicode characters like emojis may be split.

[7085, 2456, 836, 470, 3975, 284, 530, 11241, 25, 773, 452, 12843, 13]



2.3	-3.2	8.3	5.4	2.1	3.9	-8.9	3.8	3.9	3.3
4.5	5.9	4.5	7.1	1.0	5.3	5.0	3.1	0.7	5.0
...	...	...	...	...	...	...	...	...	...
3.8	1.2	3.8	9.0	9.3	3.1	4.2	0.8	9.2	5.8

## Shape

[ ]



[length]



3.2	-2.3	3.8	4.5	1.2	9.3	-9.8	8.3	9.3	3.3
5.4	9.5	5.4	1.7	0.1	3.5	0.5	1.3	7.0	0.5
...	...	...	...	...	...	...	...	...	...
8.3	2.1	8.3	0.9	3.9	1.3	2.4	8.0	2.9	8.5

[d\_model, length]



[d\_model, length]

## Process

"Many words don't map to one token: indivisible."

↓ Tokenization

Unicode characters like emojis may be split.

[7085, 2456, 836, 470, 3975, 284, 530, 11241, 25, 773, 452, 12843, 13]

↓ Embedding

2.3	-3.2	8.3	5.4	2.1	3.9	-8.9	3.8	3.9	3.3
4.5	5.9	4.5	7.1	1.0	5.3	5.0	3.1	0.7	5.0
...	...	...	...	...	...	...	...	...	...
3.8	1.2	3.8	9.0	9.3	3.1	4.2	0.8	9.2	5.8

## Shape

[ ]

↓

[length]

↓ N Transformer layers

3.2	-2.3	3.8	4.5	1.2	9.3	-9.8	8.3	9.3	3.3
5.4	9.5	5.4	1.7	0.1	3.5	0.5	1.3	7.0	0.5
...	...	...	...	...	...	...	...	...	...
8.3	2.1	8.3	0.9	3.9	1.3	2.4	8.0	2.9	8.5

[d\_model, length]

↓

[d\_model, length]

↓ Loss function (predict next token given previous)

↓

## Batched Process

Many words don't map to one token: indivisible.

↓ Tokenization

Many words don't map to one token: indivisible.  
Unicode characters like emojis may be split.

```
[7085, 2456, 836, 470, 3975, 284, 530, 11241, 25, 773, 452, 12843, 13]  
[3118, 291, 1098, 3435, 588, 795, 13210, 271, 743, 307, 6626]]
```

↓ Embedding

8.2	2.0	6.9	9.1	8.1	3.1	5.1	4.4	3.7	0.1	
2.3	-3.2	8.3	5.4	2.1	3.9	-8.9	3.8	3.9	3.3	0
4.5	5.9	4.5	7.1	1.0	5.3	5.0	3.1	0.7	5.0	
...	...	...	...	...	...	...	...	...	...	8
3.8	1.2	3.8	9.0	9.3	3.1	4.2	0.8	9.2	5.8	

↓ N Transformer layers

8.2	2.0	6.9	9.1	8.1	3.1	5.1	4.4	3.7	0.1	
3.2	-2.3	3.8	4.5	1.2	9.3	-9.8	8.3	9.3	3.3	0
5.4	9.5	5.4	1.7	0.1	3.5	0.5	1.3	7.0	0.5	
...	...	...	...	...	...	...	...	...	...	8
8.3	2.1	8.3	0.9	3.9	1.3	2.4	8.0	2.9	8.5	

↓ Loss function (predict next token given previous)

## Batched Shape

[batch]

↓

[batch, length]

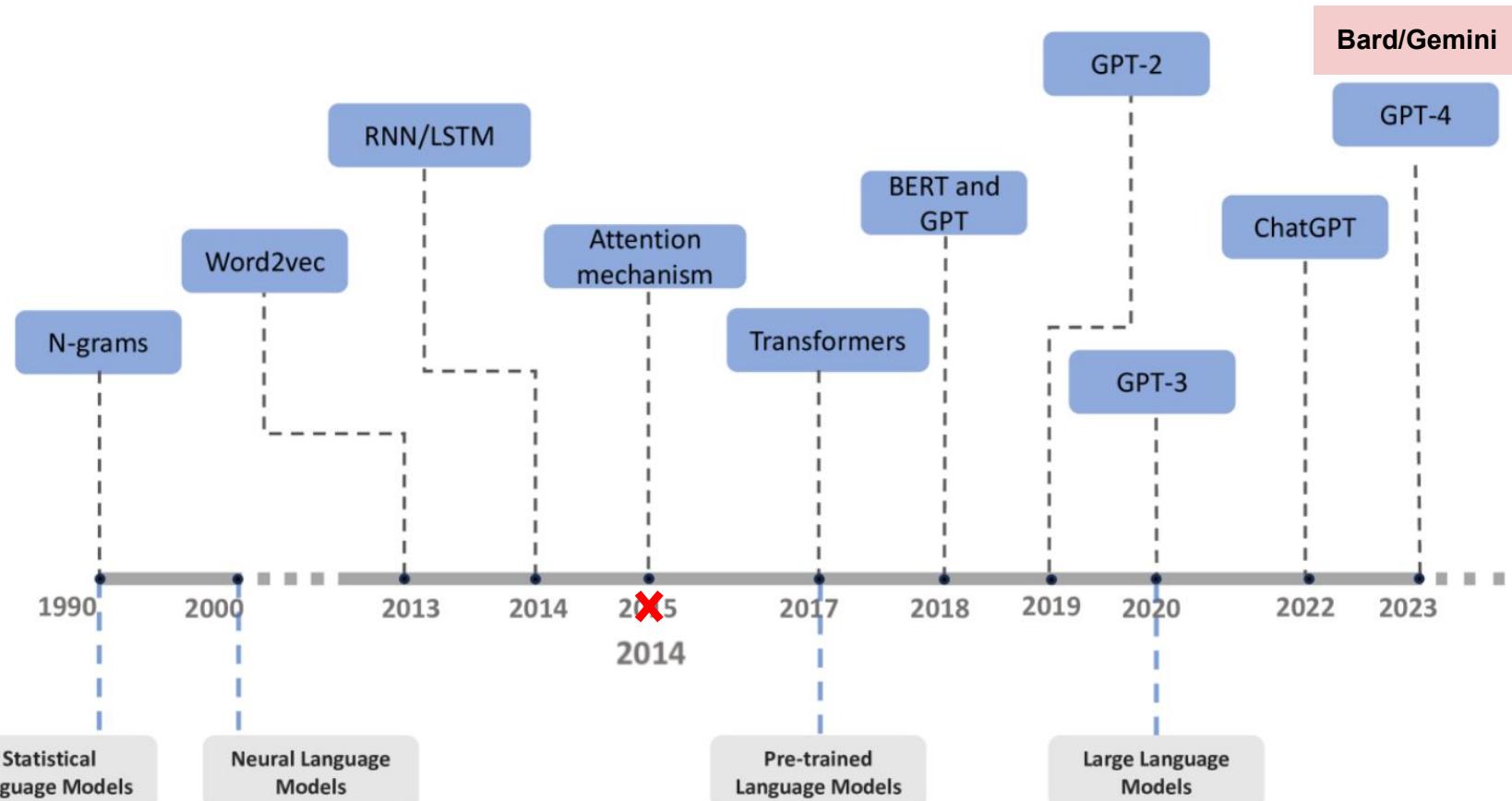
↓

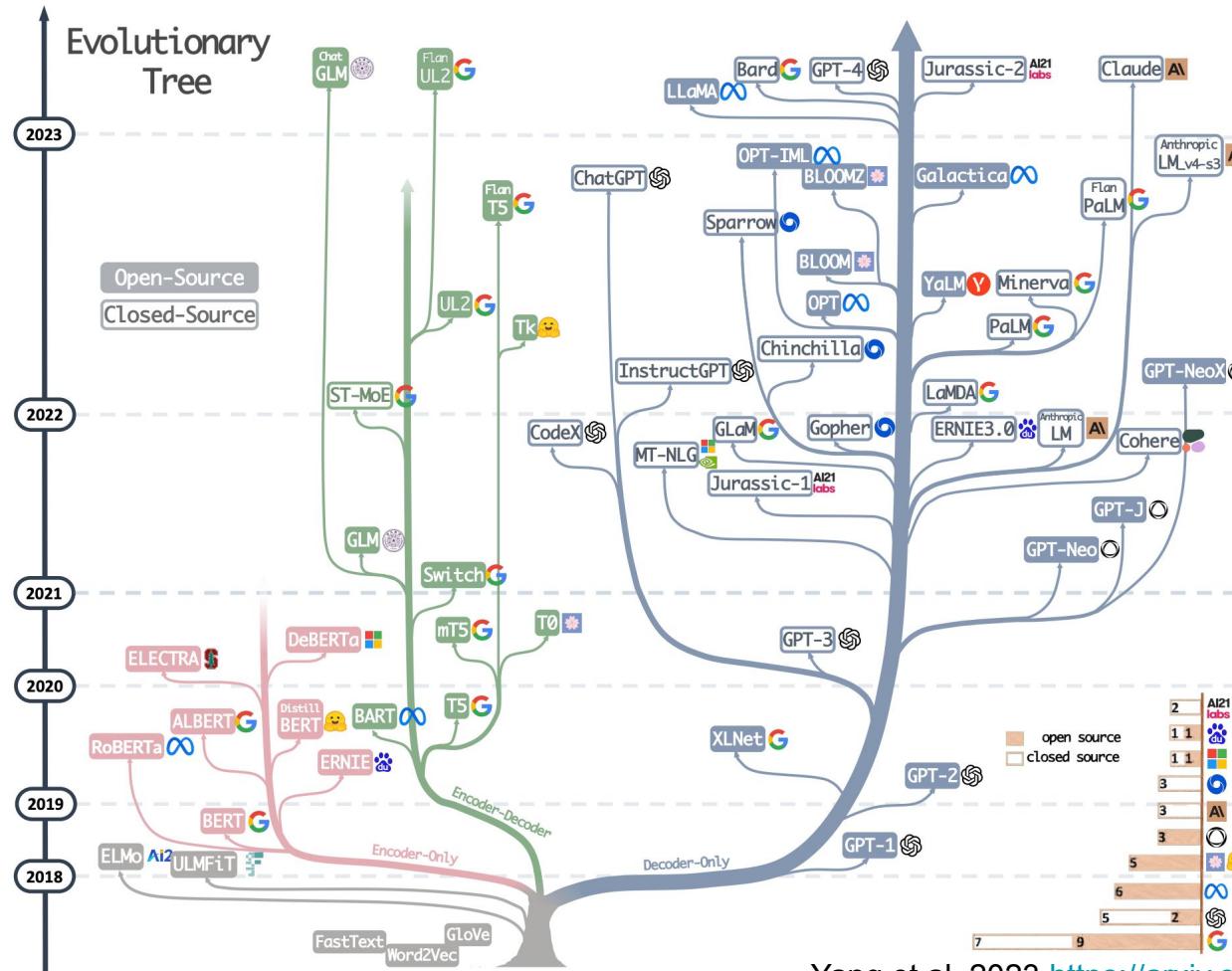
[batch, d\_model, length]

↓

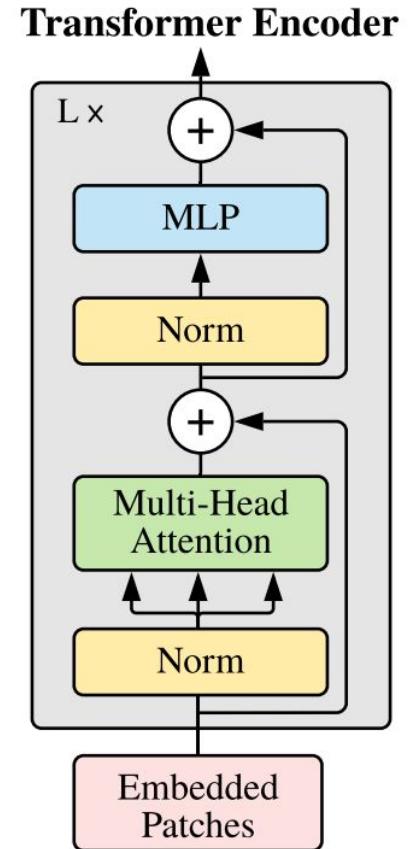
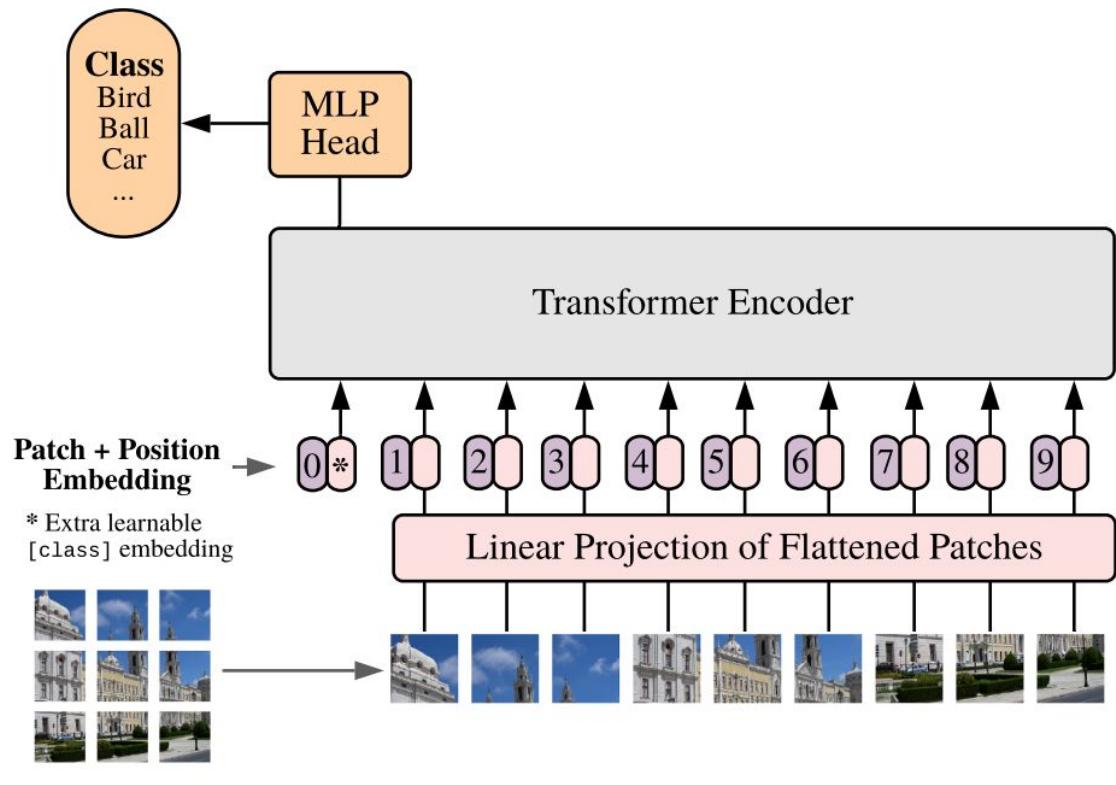
[batch, d\_model, length]

↓

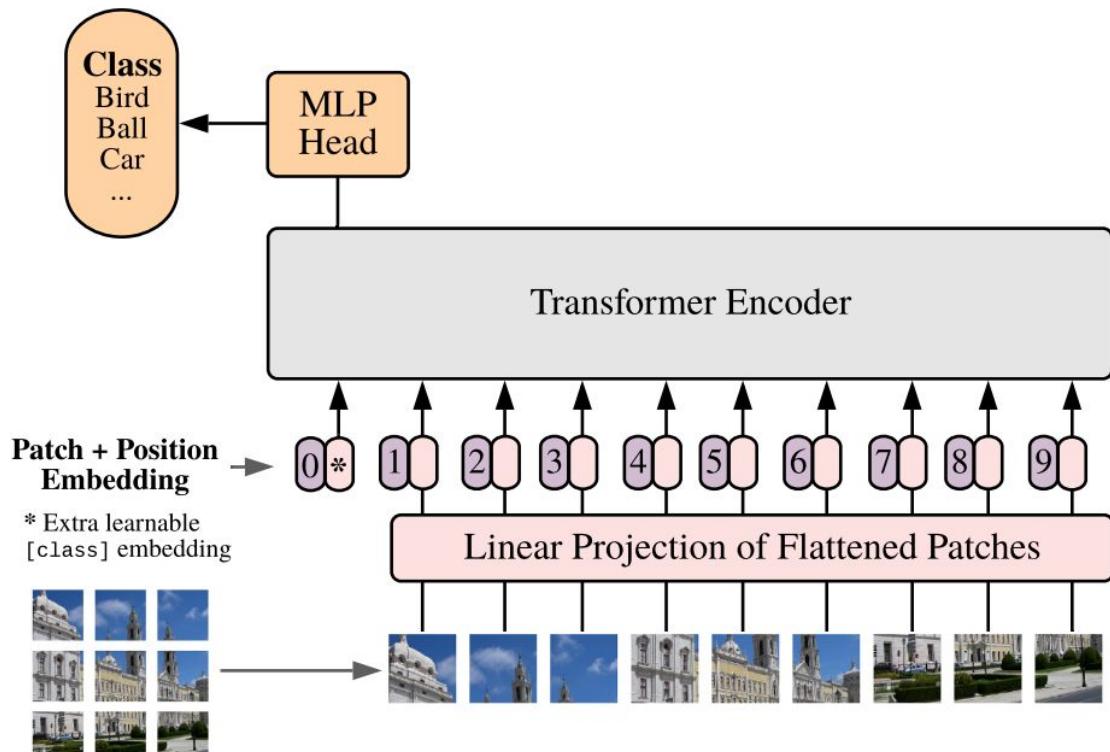




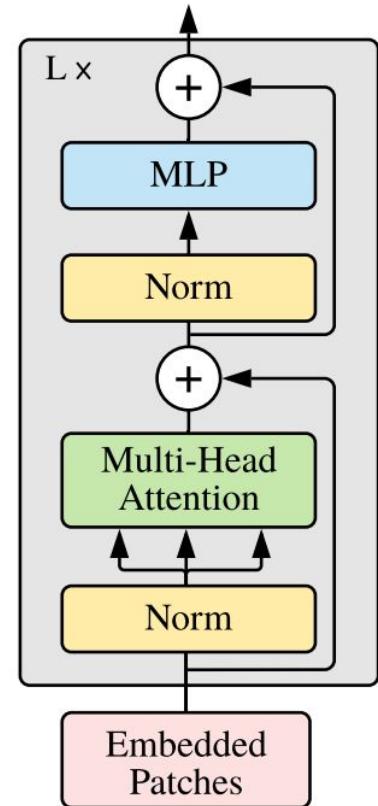
**2020**



## Vision Transformer (ViT)



## Transformer Encoder



# AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>

<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising

Google Research, Brain Team

{adosovitskiy, neilhoulsby}@google.com

## ABSTRACT

While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision remain limited. In vision, attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. We show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained on large amounts of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-100, VTAB, etc.), Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.<sup>1</sup>

[PDF] An image is worth 16x16 words: **Transformers** for image recognition at scale

[A Dosovitskiy, L Beyer, A Kolesnikov...](#) - arXiv preprint arXiv ..., 2020 - arxiv.org

... To begin to understand how the **Vision Transformer** processes image data, we analyze its internal representations. The first layer of the **Vision Transformer** linearly projects the flattened ...

☆ Save 99 Cite Cited by 60237 Related articles All 21 versions »

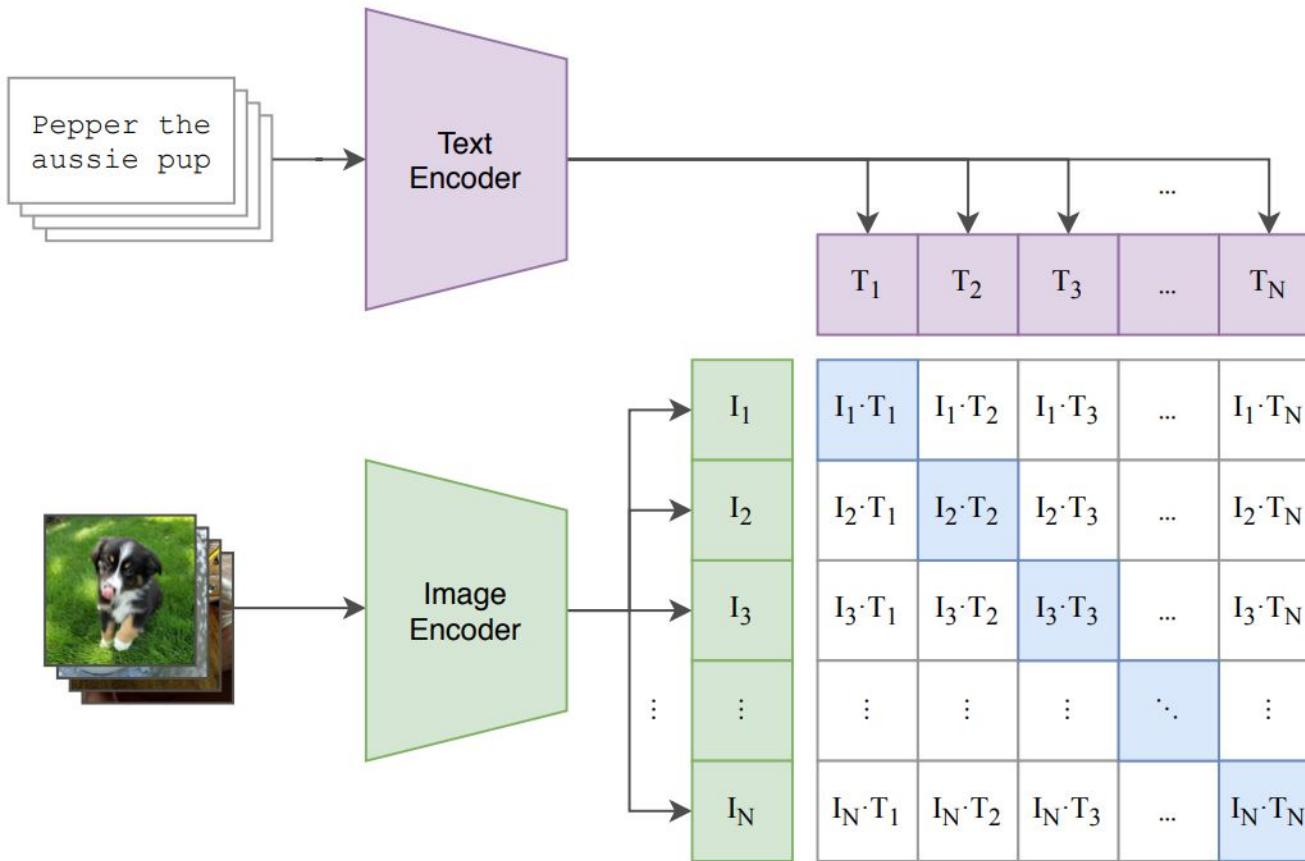
**Learning transferable visual models from natural language supervision**

[A Radford, JW Kim, C Hallacy...](#) - ... machine learning, 2021 - proceedings.mlr.press

... We speculate this is due to **natural language** providing wider **supervision** for **visual** concepts involving verbs, compared to the noun-centric object **supervision** in ImageNet. ...

☆ Save 99 Cite Cited by 33015 Related articles All 22 versions »

## (1) Contrastive pre-training



# Summary

- **Vision:** CNNs significantly advanced state-of-the-art **image recognition** in **2012**, then going on to improve other computer vision related tasks
- **Vision+Language:** Attention mechanisms developed in **2014-2015**
- **Language:** Transformers started to replace **RNNs** for language tasks since **2017**, starting with **machine translation**
- **Vision:** Transformers as image encoders, **VITs**, popularized in **2020**
- **Vision+Language** took whatever worked best **at the moment**

# **Multimodal LLMs**

## **Modeling & Data & Evaluation**

# The Revolution of Multimodal Large Language Models: A Survey

**Davide Caffagni<sup>1\*</sup>, Federico Cocchi<sup>1,2\*</sup>, Luca Barsellotti<sup>1\*</sup>, Nicholas Moratelli<sup>1\*</sup>,  
Sara Sarto<sup>1\*</sup>, Lorenzo Baraldi<sup>2\*</sup>, Lorenzo Baraldi<sup>1</sup>, Marcella Cornia<sup>1</sup>, and Rita Cucchiara<sup>1,3</sup>**

<sup>1</sup>University of Modena and Reggio Emilia, Italy

<sup>2</sup>University of Pisa, Italy

<sup>3</sup>IIT-CNR, Italy

<sup>1</sup>{name.surname}@unimore.it    <sup>2</sup>{name.surname}@phd.unipi.it

Model	LLM	Visual Encoder	V2L Adapter	VIstr. Tuning	Main Tasks & Capabilities
BLIP-2 (Li et al., 2023g)	FlanT5-XXL-11B★	EVA ViT-g	Q-Former	✗	Visual Dialogue, VQA, Captioning, Retrieval
FROMAGe (Koh et al., 2023b)	OPT-6.7B★	CLIP ViT-L	Linear	✗	Visual Dialogue, Captioning, Retrieval
Kosmos-1 (Huang et al., 2023b)	Magneto-1.3B◊	CLIP ViT-L	Q-Former*	✗	Visual Dialogue, VQA, Captioning
LLaMA-Adapter V2 (Gao et al., 2023)	LLaMA-7B▲	CLIP ViT-L	Linear	✗	VQA, Captioning
OpenFlamingo (Awadalla et al., 2023)	MPT-7B★	CLIP ViT-L	XAttn LLM	✗	VQA, Captioning
Flamingo (Alayrac et al., 2022)	Chinchilla-70B★	NFNet-F6	XAttn LLM	✗	Visual Dialogue, VQA, Captioning
PaLI (Chen et al., 2023j)	mT5-XXL-13B♦	ViT-e	XAttn LLM	✗	Multilingual, VQA, Captioning, Retrieval
PaLI-X (Chen et al., 2023h)	UL2-32B♦	ViT-22B	XAttn LLM	✗	Multilingual, VQA, Captioning
LLaVA (Liu et al., 2023e)	Vicuna-13B♦	CLIP ViT-L	Linear	✓	Visual Dialogue, VQA, Captioning
MiniGPT-4 (Zhu et al., 2023a)	Vicuna-13B★	EVA ViT-g	Linear	✓	VQA, Captioning
mPLUG-Owl (Ye et al., 2023c)	LLaMA-7B▲	CLIP ViT-L	Q-Former*	✓	Visual Dialogue, VQA
InstructBLIP (Dai et al., 2023)	Vicuna-13B★	EVA ViT-g	Q-Former	✓	Visual Dialogue, VQA, Captioning
MultiModal-GPT (Gong et al., 2023)	LLaMA-7B▲	CLIP ViT-L	XAttn LLM	✓	Visual Dialogue, VQA, Captioning
LaVIN (Luo et al., 2023)	LLaMA-13B▲	CLIP ViT-L	MLP	✓	Visual Dialogue, VQA, Captioning
Otter (Li et al., 2023b)	LLaMA-7B★	CLIP ViT-L	XAttn LLM	✓	VQA, Captioning
Kosmos-2 (Peng et al., 2023)	Magneto-1.3B◊	CLIP ViT-L	Q-Former*	✓	Visual Dialogue, VQA, Captioning, Referring, REC
Shikra (Chen et al., 2023f)	Vicuna-13B♦	CLIP ViT-L	Linear	✓	Visual Dialogue, VQA, Captioning, Referring, REC, GroundCap
Clever Flamingo (Chen et al., 2023b)	LLaMA-7B▲	CLIP ViT-L	XAttn LLM	✓	Visual Dialogue, VQA, Captioning
SVIT (Zhao et al., 2023a)	Vicuna-13B♦	CLIP ViT-L	MLP	✓	Visual Dialogue, VQA, Captioning
BLIVA (Hu et al., 2024)	Vicuna-7B★	EVA ViT-g	Q-Former+Linear	✓	Visual Dialogue, VQA, Captioning
IDEFICS (Laurençon et al., 2024)	LLaMA-65B★	OpenCLIP ViT-H	XAttn LLM	✓	Visual Dialogue, VQA, Captioning
Owen-VL (Bai et al., 2023b)	Owen-7B♦	OpenCLIP ViT-bigG	Q-Former*	✓	Visual Dialogue, Multilingual, VQA, Captioning, REC
StableLLaVA (Li et al., 2023i)	Vicuna-13B♦	CLIP ViT-L	Linear	✓	Visual Dialogue, VQA, Captioning
Ferret (You et al., 2023)	Vicuna-13B♦	CLIP ViT-L	Linear	✓	Visual Dialogue, Captioning, Referring, REC, GroundCap
LLaVA-1.5 (Liu et al., 2023d)	Vicuna-13B♦	CLIP ViT-L	MLP	✓	Visual Dialogue, VQA, Captioning
MiniGPT-v2 (Chen et al., 2023e)	LLaMA-2-7B▲	EVA ViT-g	Linear	✓	Visual Dialogue, VQA, Captioning, Referring, REC, GroundCap
Pink (Xuan et al., 2023)	Vicuna-7B▲	CLIP ViT-L	Linear	✓	Visual Dialogue, VQA, Captioning, Referring, REC, GroundCap
CogVLM (Wang et al., 2023c)	Vicuna-7B♦	EVA ViT-E	MLP	✓	Visual Dialogue, VQA, Captioning, REC
DRESS (Chen et al., 2023l)	Vicuna-13B▲	EVA ViT-g	Linear	✓	Visual Dialogue, VQA, Captioning
LION (Chen et al., 2023d)	FlanT5-XXL-11B★	EVA ViT-g	Q-Former+MLP	✓	Visual Dialogue, VQA, Captioning, REC
mPLUG-Owl2 (Ye et al., 2023d)	LLaMA-2-7B♦	CLIP ViT-L	Q-Former*	✓	Visual Dialogue, VQA, Captioning
SPHINX (Lin et al., 2023b)	LLaMA-2-13B♦	Mixture	Linear	✓	Visual Dialogue, VQA, Captioning, Referring, REC, GroundCap
Honeybee (Cha et al., 2023)	Vicuna-13B♦	CLIP ViT-L	ResNet blocks	✓	Visual Dialogue, VQA, Captioning
VILA (Lin et al., 2023a)	LLaMA-2-13B♦	CLIP ViT-L	Linear	✓	Visual Dialogue, VQA, Captioning
SPHINX-X (Gao et al., 2024)	Mixtral-8×7B♦	Mixture	Linear	✓	Visual Dialogue, Multilingual, VQA, Captioning, Referring, REC

Table 1: Summary of generalist MLLMs for vision-to-language tasks. For each model, we indicate the LLM used in its best configuration as shown in the original paper (◊: LLM training from scratch; ♦: LLM fine-tuning; ▲: LLM fine-tuning with PEFT techniques; ★: frozen LLM). The \* marker indicates variants to the reported vision-to-language adapter, while gray color indicates models not publicly available.

<b>Model</b>	<b>LLM</b>	<b>Visual Encoder</b>	<b>V2L Adapter</b>	<b>VInstr. Tuning</b>
BLIP-2 (Li et al., 2023g)	FlanT5-XXL-11B★	EVA ViT-g	Q-Former	X
FROMAGe (Koh et al., 2023b)	OPT-6.7B★	CLIP ViT-L	Linear	X
Kosmos-1 (Huang et al., 2023b)	Magneto-1.3B◊	CLIP ViT-L	Q-Former*	X
LLaMA-Adapter V2 (Gao et al., 2023)	LLaMA-7B▲	CLIP ViT-L	Linear	X
OpenFlamingo (Awadalla et al., 2023)	MPT-7B★	CLIP ViT-L	XAttn LLM	X
Flamingo (Alayrac et al., 2022)	Chinchilla-70B★	NFNet-F6	XAttn LLM	X
PaLI (Chen et al., 2023j)	mT5-XXL-13B♦	ViT-e	XAttn LLM	X
PaLI-X (Chen et al., 2023h)	UL2-32B♦	ViT-22B	XAttn LLM	X

Exploring the Limits of Transfer Learning with a Unified  
Text-to-Text Transformer



**mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer**

Mihir Kale Linting Xue\* Noah Constant\* Adam Roberts\*  
Rami Al-Rfou Aditya Siddhant Aditya Barua Colin Raffel  
*Google Research*

AN IMAGE IS WORTH 10,000 WORDS:  
TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy\*<sup>†</sup>, Lucas Beyer\*, Alexander Kolesnikov\*, Dirk Weissenborn\*,  
Xiaohua Zhai\*, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby\*<sup>†</sup>

\*equal technical contribution, <sup>†</sup>equal advising  
Google Research, Brain Team

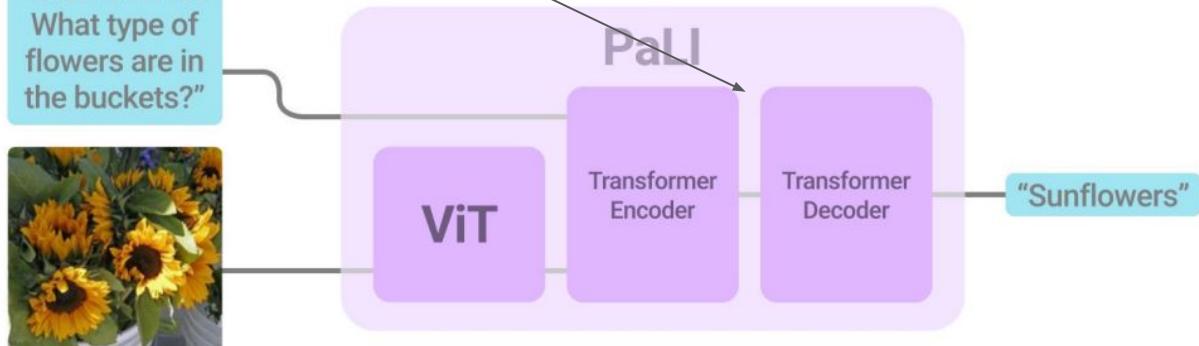
{adosovitskiy, neilhoulsby}@google.com

"Answer in EN:  
What type of  
flowers are in  
the buckets?"



# PALI: A JOINTLY-SCALED MULTILINGUAL LANGUAGE-IMAGE MODEL

Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski  
Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer  
Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari  
Gaurav Mishra, Linting Xue, Ashish Thapliyal, James Bradbury, Weicheng Kuo  
Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme  
Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, Radu Soricut  
Google Research\*



# UL2

## UL2: Unifying Language Learning Paradigms

Yi Tay\* Mostafa Dehghani\*

Vinh Q. Tran<sup>#</sup> Xavier Garcia<sup>#</sup> Jason Wei<sup>#</sup> Xuezhi Wang<sup>#</sup> Hyung Won Chung<sup>#</sup>

Siamak Shakeri<sup>#</sup> Dara Bahri<sup>b</sup> Tal Schuster<sup>b</sup> Huaixiu Steven Zheng<sup>△</sup>

Denny Zhou<sup>△</sup> Neil Houlsby<sup>△</sup> Donald Metzler<sup>△</sup>

Google Brain

## Scaling Vision Tr... ViT-22B

## lion Parameters

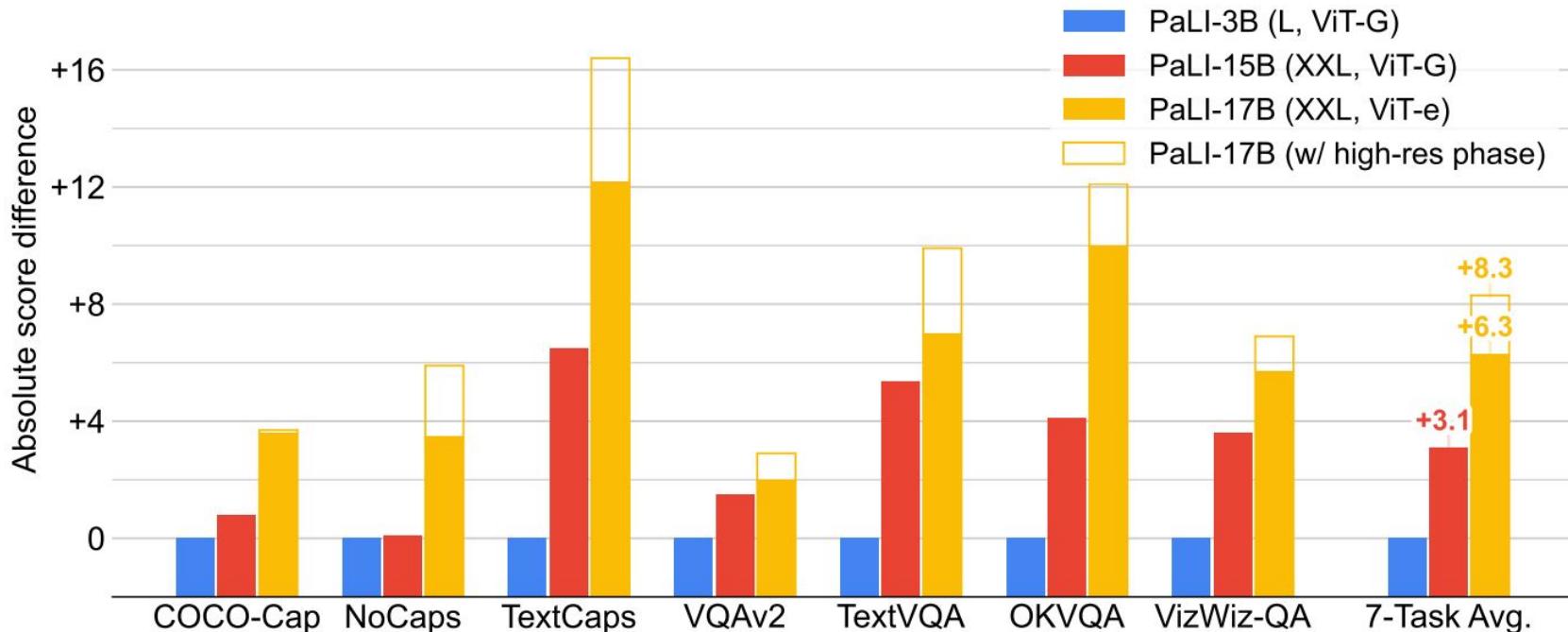
Mostafa Dehghani\* Josip Djolonga\* Basil Mustafa\* Piotr Padlewski\* Jonathan Heek\*  
Justin Gilmer Andreas Steiner Mathilde Caron Robert Geirhos Ibrahim Alabdulmohsin  
Rodolphe Jenatton Lucas Beyer Michael Tschannen Anurag Arnab Xiao Wang  
Carlos Riquelme Matthias Minderer Joan Puigcerver Utku Evcı Manoj Kumar  
Sjoerd van Steenkiste Gamaleldin F. Elsayed Aravindh Mahendran Fisher Yu  
Avital Oliver Fantine Huot Jasmijn Bastings Mark Patrick Collier Alexey A. Gritsenko  
Vighnesh Birodkar Cristina Vasconcelos Yi Tay Thomas Mensink Alexander Kolesnikov  
Filip Pavetić Dustin Tran Thomas Kipf Mario Lučić Xiaohua Zhai Daniel Keysers  
Jeremiah Harmsen Neil Houlsby\*  
Google Research

## PaLI-X: On Scaling up a Multilingual Vision and Language Model

Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu,  
Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, Siamak Shakeri,  
Mostafa Dehghani, Daniel Salz, Mario Lucic, Michael Tschannen, Arsha Nagrani,  
Hexiang Hu, Mandar Joshi, Bo Pang, Ceslee Montgomery, Paulina Pietrzyk,  
Marvin Ritter, AJ Piergiovanni, Matthias Minderer, Filip Pavetic, Austin Waters,  
Gang Li, Ibrahim Alabdulmohsin, Lucas Beyer, Julien Amelot, Kenton Lee,  
Andreas Peter Steiner, Yang Li, Daniel Keysers, Anurag Arnab, Yuzhong Xu,  
Keran Rong, Alexander Kolesnikov, Mojtaba Seyedhosseini,  
Anelia Angelova, Xiaohua Zhai, Neil Houlsby, Radu Soricut

Google Research  
pali-communications@google.com

# Example: PaLI



# Example: PaLI-X

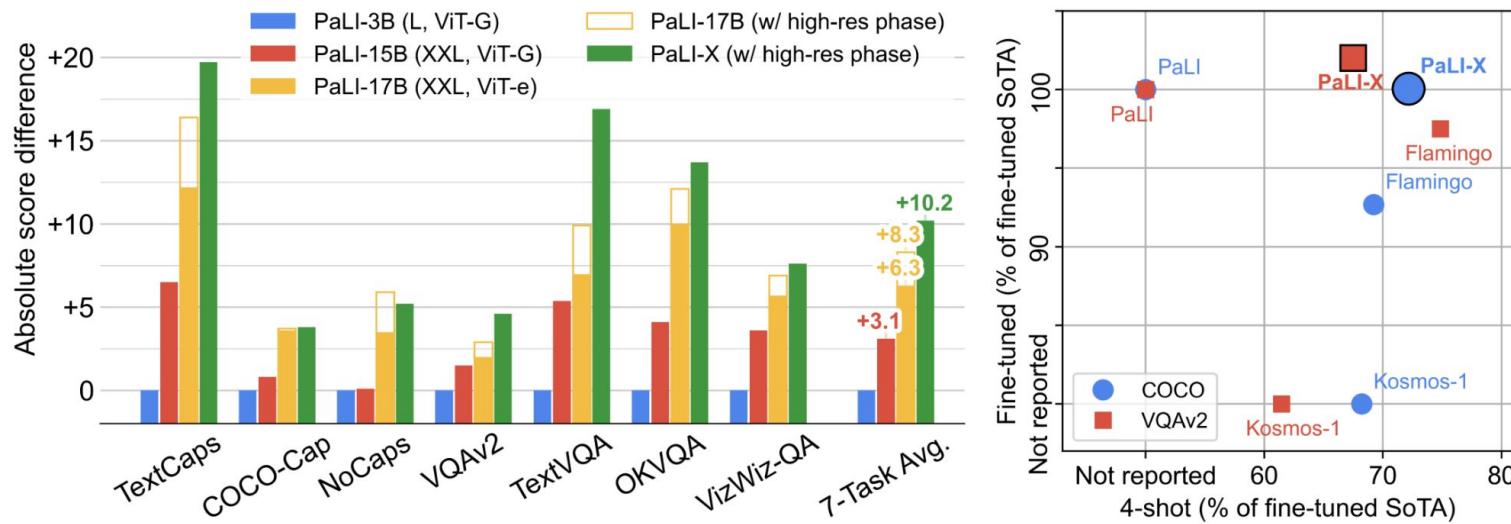


Figure 1: [Left] Comparing PaLI-X against PaLI on image-captioning and VQA benchmarks. [Right] The Pareto frontier between few-shot and fine-tuned performance, comparing PaLI-X with PaLI [5], Flamingo [10], and Kosmos-1 [11].

# Pre-Training For Scene-Text Understanding

## Main idea

**Scene-text modeling by transferring from off-the-shelf OCR**

*SplitOCR*

## PreSTU: Pre-Training for Scene-Text Understanding

Jihyung Kil<sup>1\*</sup> Soravit Changpinyo<sup>2</sup>  
Xi Chen<sup>2</sup> Hexiang Hu<sup>2</sup> Sebastian Goodman<sup>2</sup> Wei-Lun Chao<sup>1</sup> Radu Soricut<sup>2</sup>  
<sup>1</sup>The Ohio State University <sup>2</sup>Google Research

ICCV 2023

# PaLI (ICLR 2023) → PaLI-X (CVPR 2024)

## 3.2 Pretraining Data and Mixture

**Training mixture** To accommodate diverse tasks in the image-language space, we train PaLI using a mixture of eight pre-training tasks. This mixture is designed to span a range of general capabilities useful for downstream tasks. **Span corruption on text-only data** uses the same technique described by Xue et al. (2021) the pre-training obj into two parts,  $\langle \text{cap}, \langle \text{lang} \rangle \rangle$ , with the alt-text string in language  $\langle \text{lang} \rangle$  as the target, based on the Conceptual Captions (Sharma et al., 2018) training data and machine translated alt-texts. **OCR on WebLI OCR-text data** uses the concatenation of the annotated OCR texts in language  $\langle \text{lang} \rangle$  (Kil et al., 2022) produced by publicly available automatic service for the input image. **English and Cross-Lingual VQA** is VQ<sup>2</sup>A-CC3M (Changpinyo et al., 2022a), translated in the same way as CC3M-35L.

answers in all instances here, as the English-native answers for VQA are often errors to perform out-of-context automatic translation. **English and Cross-generation (VQG)** is also based on native and translated VQ<sup>2</sup>A-CC3M-35I we use only English answers here. **English-only Object-Aware (OA) VQA** is based on VQA triplets derived from automatically-produced, non-exhaustive object labels, inspired by Piergiovanni et al. (2022a). The QA pairs include listing all the objects in the image and whether a subset of objects are in the image. To create these examples, we require object-level annotations, for which we use Open Images (Kuznetsova et al., 2020). **Object detection** is a generative object-detection task inspired by Chen et al. (2021; 2022).

## OCR → SplitOCR

The main pretraining data for our model is based on WebLI [5], consisting of roughly one billion images with alt-texts from the web and OCR annotations (using the GCP Vision API), covering over 100 languages. In addition to WebLI  $\langle \text{image}, \text{text} \rangle$  pairs, we introduce here *Episodic WebLI* data, where each episode corresponds to a set of such pairs. We aim to have each episode contain loosely related images (i.e., they are clustered according to their URL field), so as to encourage attention and 400M

### (iv) split-ocr [24] on WebLI OCR annotations;

The pretraining mixture consists of the following data and objectives: (i) span corruption on text-only data (15% of tokens); (ii) split captioning on WebLI alt-text data [21, 5]; (iii) captioning on CC3M [22] on native and translated alt-text data (over the same 35 languages covered by XM3600 [23]); (iv) split-ocr [24] on WebLI OCR annotations; (v) visual-question-answering objective over  $\langle \text{image}, \text{question}, \text{answer} \rangle$  pairs generated using the VQ<sup>2</sup>A method [25] over the CC3M training split, over native and translated text (same 35 language pairs); (vi) visual-question-generation me pairs as above; (vii) visual-question-answering objective over  $\langle \text{image}, \text{question} \rangle$  pairs using the Object-Aware method [26] (English only); (viii) captioning on images (target alt-text predicted from the remaining alt-text and images); (ix) captioning on 4-pair examples (resembling Episodic WebLI and using VQ<sup>2</sup>A-CC3M pairs), with the answer target conditioned on the other pairs of  $\langle \text{image}, \text{question}, \text{answer} \rangle$  data. (x) pix2struct objective, introduced in [27], targeting page layout and structure using screenshot images paired with DOM-tree representations of html pages. (xi) Captioning on short video data, using the VTP data [10] (using four frames per video). (xii) object-detection objective on WebLI data, whereby an OWL-ViT model [28] (L/14) is used to annotate WebLI images, resulting in hundreds of pseudo object labels and bounding boxes per image. (xiii) image-token prediction objective, whereby we tokenize WebLI images (256×256 resolution) using a ViT-VQGAN [29] model with patch size 16×16 (256 tokens per image); this objective is framed as a 2D masked-token task (i.e., fill-in the missing grid pieces, with the corresponding image pixels also masked). Note that the image-token prediction objective is added mainly as a condition to check whether it adversarially impacts the performance on language-output tasks; our ablation experiments show that is does not.

<b>Model</b>	<b>LLM</b>	<b>Visual Encoder</b>	<b>V2L Adapter</b>	<b>VInstr. Tuning</b>
BLIP-2 (Li et al., 2023g)	FlanT5-XXL-11B★	EVA ViT-g	Q-Former	X
FROMAGe (Koh et al., 2023b)	OPT-6.7B★	CLIP ViT-L	Linear	X
Kosmos-1 (Huang et al., 2023b)	Magneto-1.3B◊	CLIP ViT-L	Q-Former*	X
LLaMA-Adapter V2 (Gao et al., 2023)	LLaMA-7B▲	CLIP ViT-L	Linear	X
OpenFlamingo (Awadalla et al., 2023)	MPT-7B★	CLIP ViT-L	XAttn LLM	X
Flamingo (Alayrac et al., 2022)	Chinchilla-70B★	NFNet-F6	XAttn LLM	X
PaLI (Chen et al., 2023j)	mT5-XXL-13B♦	ViT-e	XAttn LLM	X
PaLI-X (Chen et al., 2023h)	UL2-32B♦	ViT-22B	XAttn LLM	X

# Adaptors

- Linear projection (Linear) & Multi-layer perceptrons (MLP)
- Querying Transformers (Q-Former)
- Cross Attention Mechanisms (XAttn)

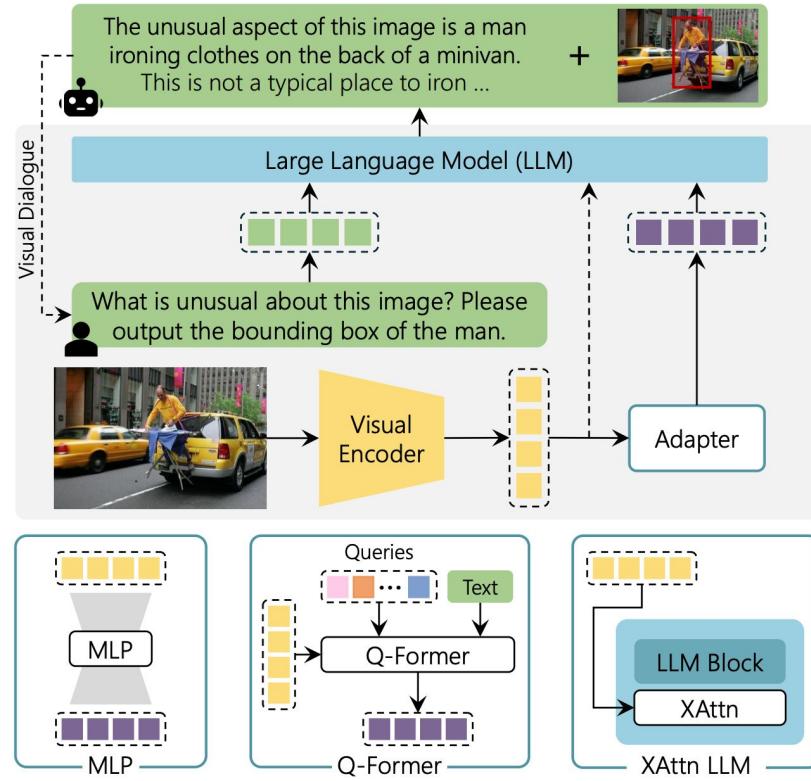
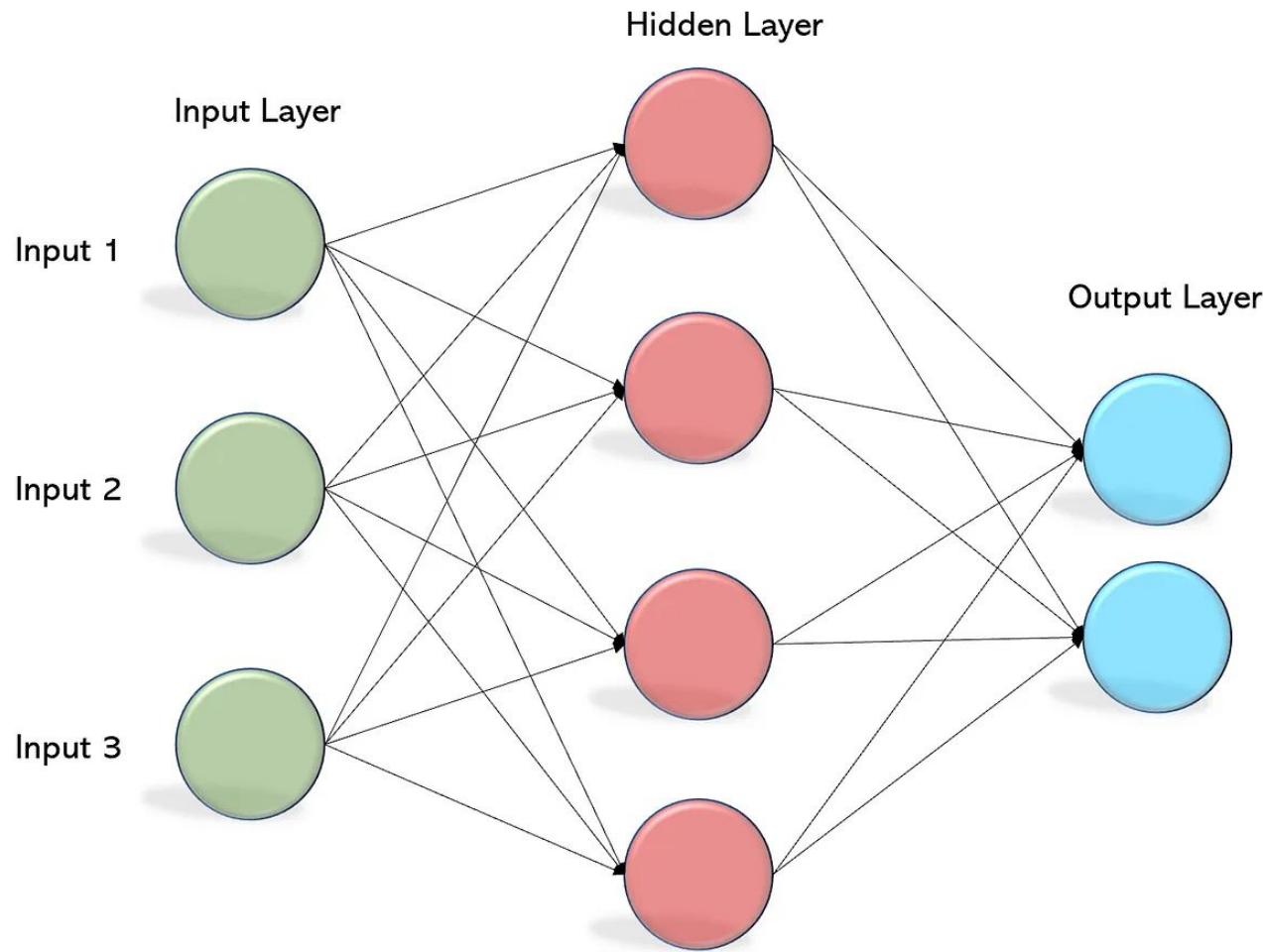
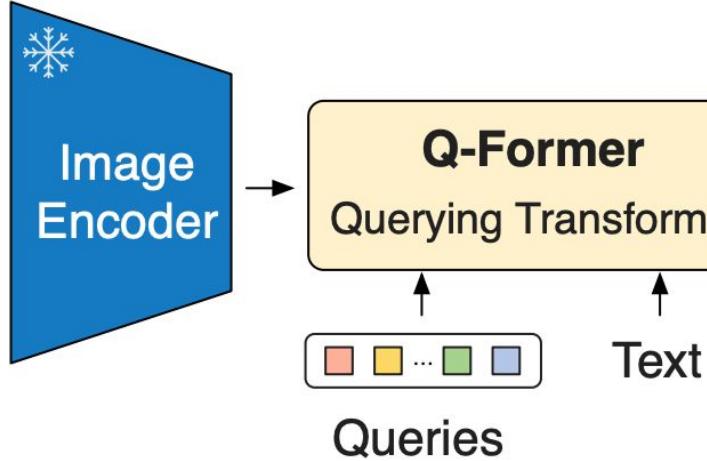


Figure 1: General architecture of Multimodal Large Language Models (MLLMs), composed of a visual encoder, a language model, and an adapter module that connects visual inputs to the textual space.

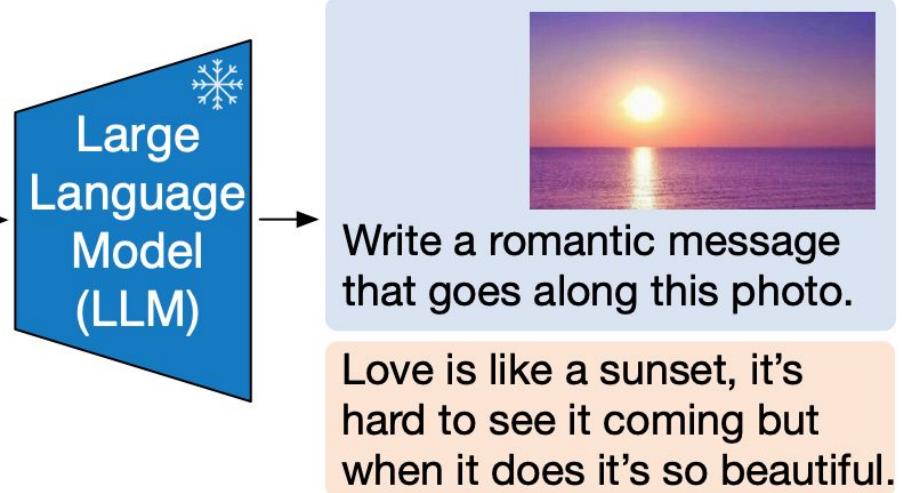


## Vision-and-Language Representation Learning

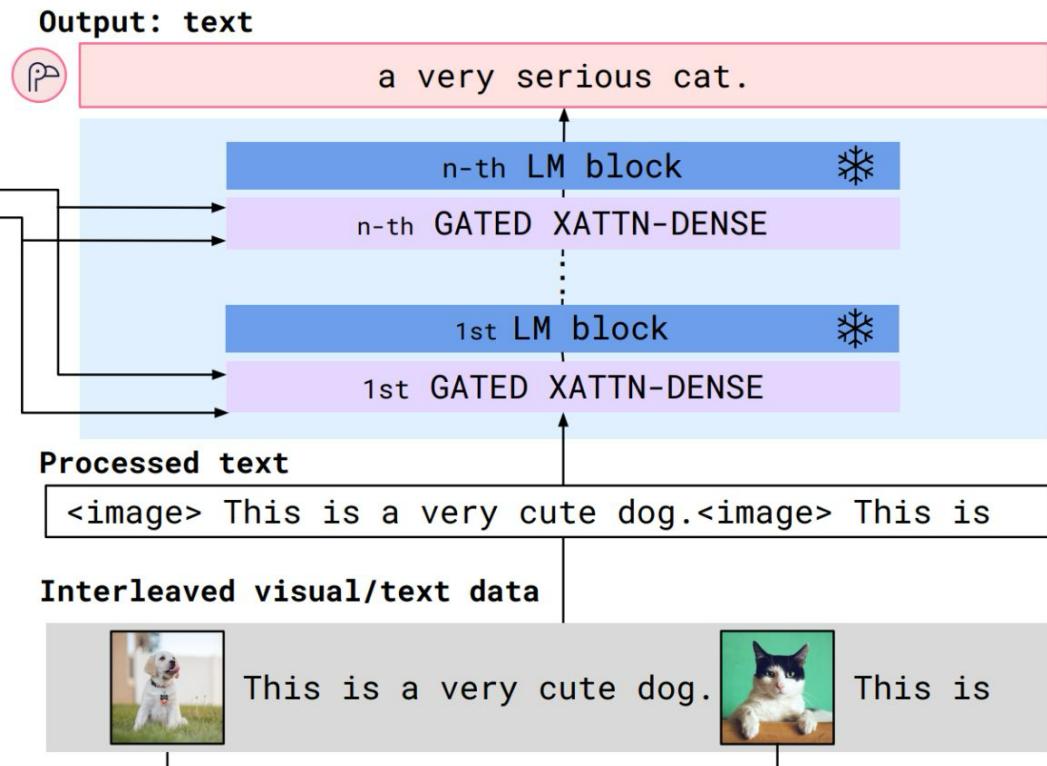
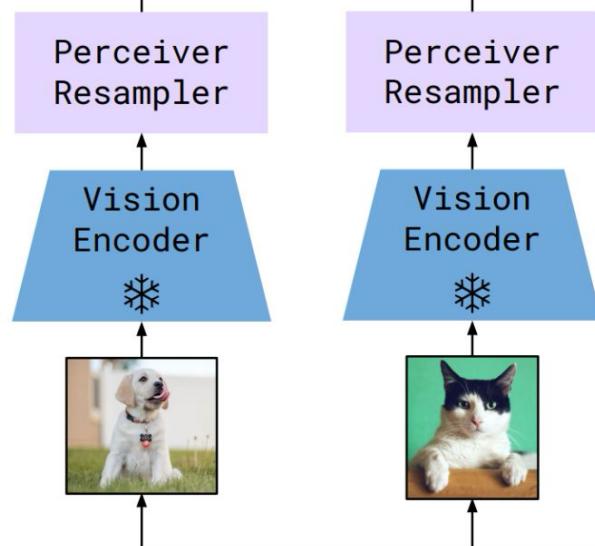
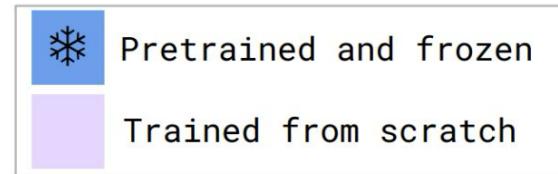


**Bootstrapping Pre-trained Image Models**

## Vision-to-Language Generative Learning



**Bootstrapping Pre-trained Large Language Models (LLMs)**



**Figure 3: Flamingo architecture overview.** Flamingo is a family of visual language models (VLMs) that take as input visual data interleaved with text and produce free-form text as output.

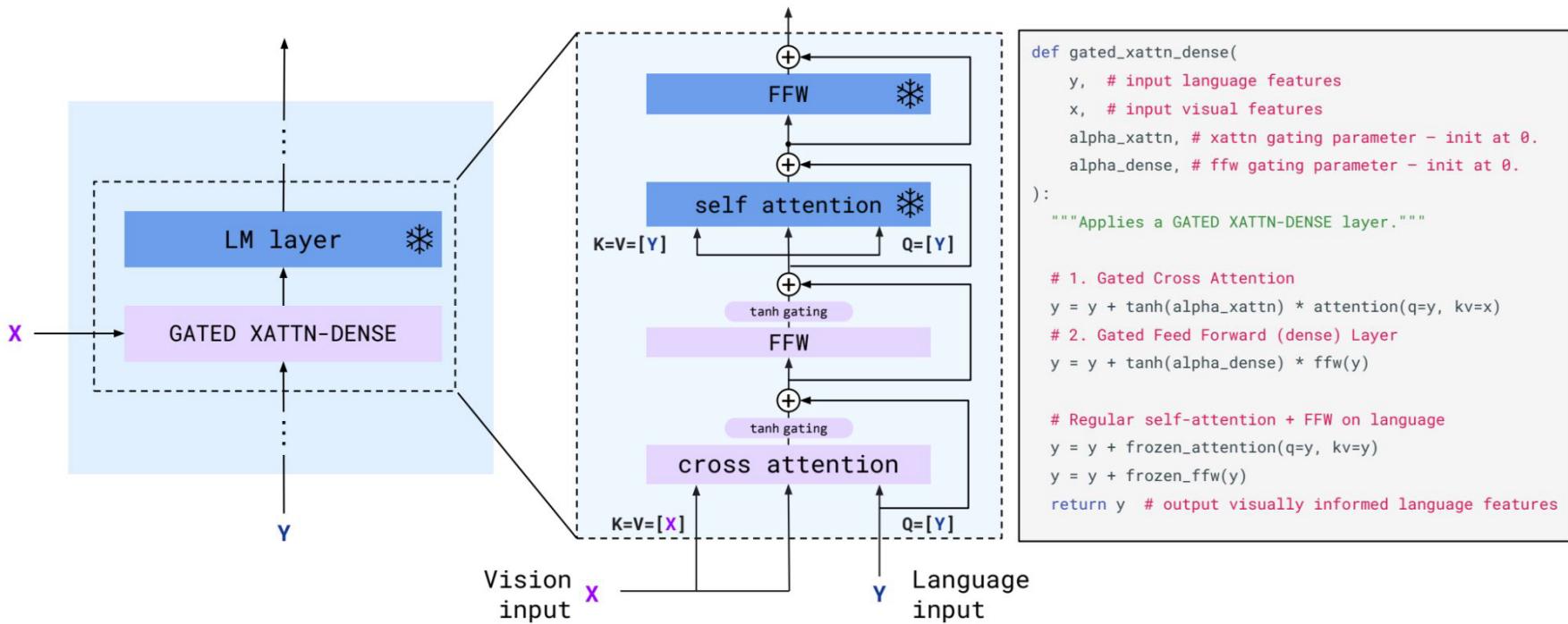


Figure 4: **GATED XATTN-DENSE layers.** To condition the LM on visual inputs, we insert new cross-attention layers between existing pretrained and frozen LM layers. The keys and values in these layers are obtained from the vision features while the queries are derived from the language inputs. They are followed by dense feed-forward layers. These layers are *gated* so that the LM is kept intact at initialization for improved stability

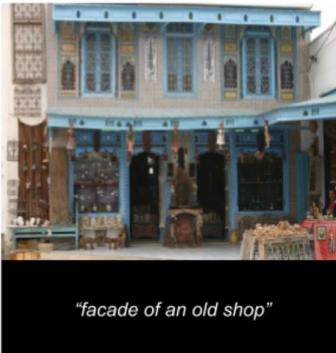
# Data

- Pre-Training: Web (Later: Books, Video, Synthetic, etc.)
- Post-Training: Mostly human-annotated

# **Pre-Training**

# Large-Scale Image+(Alt)text data

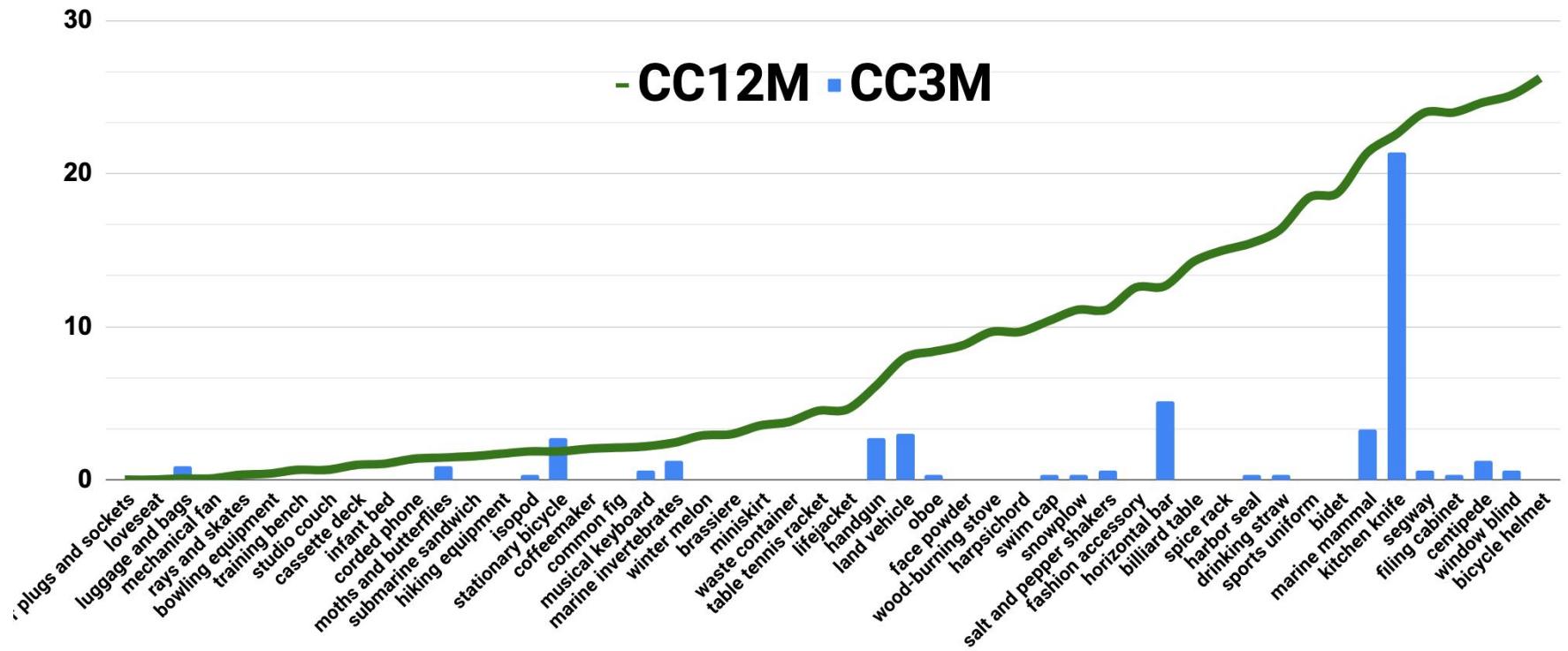
Conceptual Captions  
[Sharma et al. 18]



Conceptual 12M  
[Changpinyo et al. 21]



# Long-Tail Concepts (*Out-of-Domain* Classes from nocaps)



## Nocaps Example



No Pre-Training

CC3M Pre-Training

CC12M Pre-Training

A man in a military uniform holding a **microphone**.

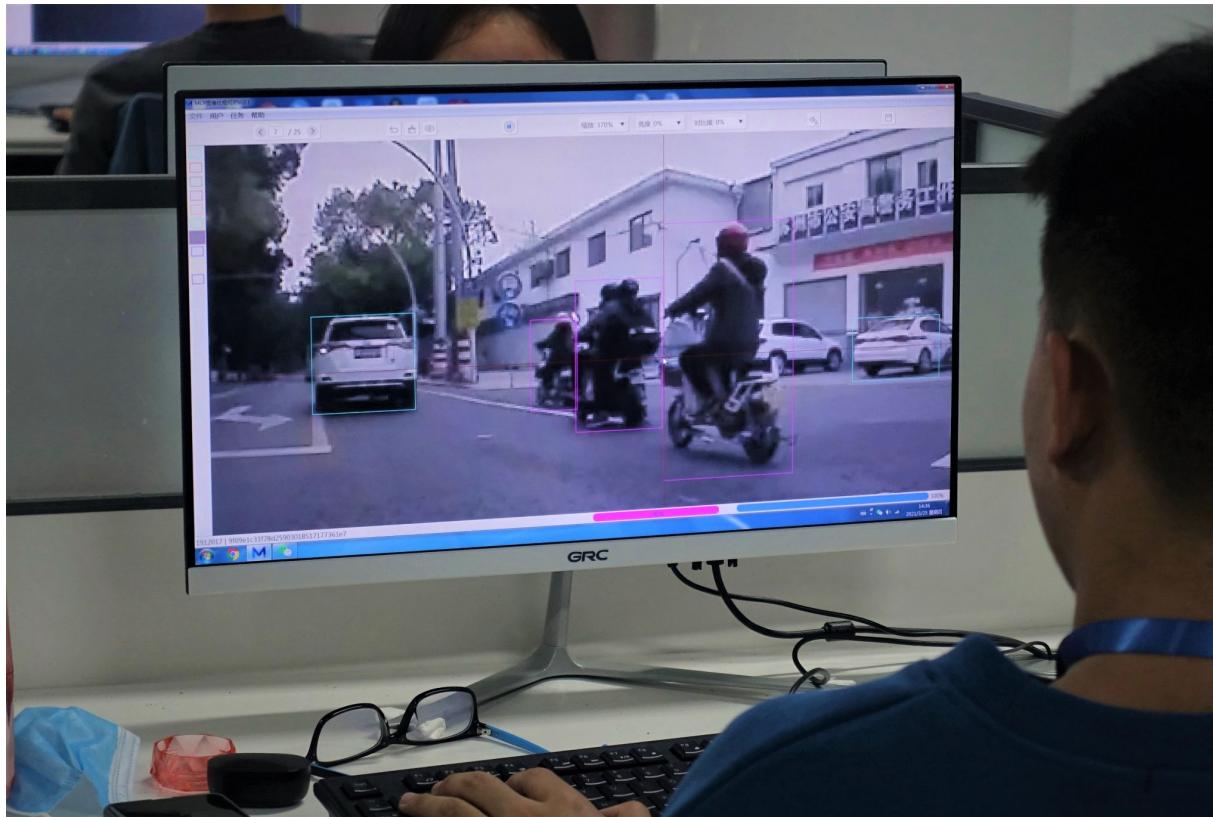
A man in a military uniform holds a **sword**.

A man in a **kilt** playing the **bagpipes**.



# **Post-Training**

## Appen, Sama, iMerit, Scale AI, Labelbox, and Amazon SageMaker Ground Truth...



# Evaluation

- **Early days:** Academic benchmarks focusing on captioning and question answering, with automatic metrics
- **These days:** Also real-world prompts, using humans and LLMs to evaluate long(er) responses. Other considerations: safety, fairness, localization, multi-turn, etc.

# Image Captioning

## MSCOCO Captions (2015)



A large bus sitting next to a very tall building.



Bunk bed with a narrow shelf sitting underneath it.

## nocaps (2019)



1. A **man** sitting in the saddle on a **camel**.
2. A **person** is sitting on a **camel** with another **camel** behind him.
3. A **man** with long hair and blue jeans sitting on a **camel**.
4. **Man** sitting on a **camel** with a standing **camel** behind them.
5. Long haired **man** wearing sitting on blanket draped **camel**
6. A **camel** stands behind a sitting **camel** with a **man** on its back.
7. The standing **camel** is near a sitting one with a **man** on its back.
8. Someone is sitting on a **camel** and is in front of another **camel**.
9. Two **camels** in the dessert and a **man** sitting on the sitting one.
10. Two **camels** are featured in the sand with a **man** sitting on one of the seated **camels**.

## Localized Narratives (2020)



In the front portion of the picture we can see a dried grass area with dried twigs. There is a woman standing wearing a light blue jeans and ash colour long sleeve length shirt. This woman is holding a black jacket in her hand. On the other hand she is holding a balloon which is peach in colour. on the top of the picture we can see a clear blue sky with clouds. The hair colour of the woman is brownish.

# VQA

## VQAv2 (2017)

Who is wearing glasses?  
man  
woman



Where is the child sitting?  
fridge  
arms



Is the umbrella upside down?  
yes  
no



How many children are in the bed?  
2  
1



Q: Does this foundation have any sunscreen?  
A: yes



Q: What is this?  
A: 10 euros



Q: What color is this?  
A: green



Q: Please can you tell me what this item is?  
A: butternut squash red pepper soup



Q: Is it sunny outside?  
A: yes



Q: Is this air conditioner on fan, dehumidifier, or air conditioning?  
A: air conditioning

## GQA (2019)



- A1. Is the **tray** on top of the **table** black or light brown? light brown  
 A2. Are the **napkin** and the **cup** the same color? yes  
 A3. Is the small **table** both oval and wooden? yes  
 A4. Is there any **fruit** to the left of the **tray** the **cup** is on top of? yes  
 A5. Are there any **cups** to the left of the **tray** on top of the **table**? no  
 B1. What is the brown **animal** sitting inside of? **box**  
 B2. What is the large **container** made of? cardboard  
 B3. What **animal** is in the **box**? **bear**  
 B4. Is there a **bag** to the right of the green **door**? no  
 B5. Is there a **box** inside the plastic **bag**? no

### Vehicles and Transportation



Q: What sort of vehicle uses this item?  
A: firetruck

### Brands, Companies and Products



Q: When was the soft drink company shown first created?  
A: 1898

### Objects, Material and Clothing



Q: What is the material used to make the vessels in this picture?  
A: copper

### Sports and Recreation



Q: What is the sports position of the man in the orange shirt?  
A: goalie

### Cooking and Food



Q: What is the name of the object used to eat this food?  
A: chopsticks

### Geography, History, Language and Culture



Q: What days might I most commonly go to this building?  
A: Sunday

### People and Everyday Life



Q: Is this photo from the 50's or the 90's?  
A: 50's

### Plants and Animals



Q: What phylum does this animal belong to?  
A: chordate, chordata

### Science and Technology



Q: How many chromosomes do these creatures have?  
A: 23

### Weather and Climate



Q: What is the warmest outdoor temperature at which this kind of weather can happen?  
A: 32 degrees

## OK-VQA (2019)



Q: Combien de nems sont servis dans une assiette rectangulaire blanche avec de la sauce soja?

A: 5, 5 nems, cinq, cinq nems

(Q: *How many spring rolls are served in a white rectangular plate with soy sauce?*

A: 5, 5 spring rolls, five, five spring rolls)



Q: Ce fel de cafea este în ceașca albă?

A: cu crema, cafea cu lapte, cafea cu spumă, cafea cu spumă de lapte, cafea neagră cu spumă

(Q: *What kind of coffee is in the white cup?*

A: with cream, coffee with milk, coffee with foam, coffee with milk foam, black coffee with foam)



Q: काली कार किस सतह पर दिखाई दे रही हैं?

A: सड़क पर, काली सड़क पर, सड़क पे, काली सड़क पे

(Q: *On which surface is the black car visible?*

A: on the road, on the black road, on the road, on the black road)



Q: กระโปรงบนตุ๊กตาสีอะไร?

A: สีขาว, ขาว

(Q: *What color is the skirt on the doll?*

A: white, white)



Q: מה השם של העוף זהה?

A: קזואר

(Q: *What is the name of this bird?*

A: Cassowary)



Q: 除了可回收物, 垃圾桶上还印有什么?

A: 其他垃圾

(Q: *Besides recyclables, what else is printed on the trash can?*

A: other waste)

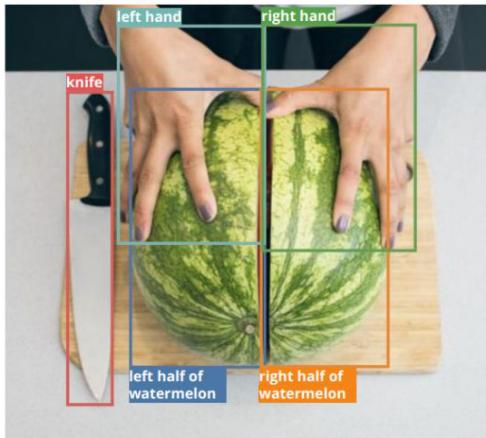
Model	Text Caps	VizWiz Cap	Text VQA	VizWiz VQA	ST VQA	OCR VQA	Info VQA	Doc VQA	AI2D QA	Chart Words	Screen2 Cap	Widget Cap	OVEN	Info Seek
<i>with OCR pipeline input</i>														
SoTA	160.4 [5]	124.7 [5]	73.67 [52]	73.3 [5]	79.9 [5]	67.5 [53]	47.4 [54]	84.7 [54]	38.5 [45]	45.5 [46]	-	-	-	-
PaLI-X	<b>163.7</b>	<b>125.7</b>	<b>80.78</b>	<b>74.6</b>	<b>84.5</b>	<b>77.3</b>	<b>54.8</b>	<b>86.8</b>	<b>81.4</b>	<b>72.3</b>	-	-	-	-
<i>without OCR pipeline input</i>														
SoTA	145.0 [9]	120.8 [9]	67.27 [9]	70.7 [5]	75.8 [9]	71.3 [27]	40.0 [27]	76.6 [27]	42.1 [27]	70.5 [8]	109.4 [27]	141.8 [20]	20.0 [47]	17.7 [48]
PaLI-X	<b>147.0</b>	<b>122.7</b>	<b>71.44</b>	<b>70.9</b>	<b>79.9</b>	<b>75.0</b>	<b>49.2</b>	<b>80.0</b>	<b>81.2</b>	<b>70.9</b>	<b>127.9</b>	<b>153.0</b>	<b>23.1</b>	<b>21.8</b>

Table 2: Results on benchmarks more focused on text understanding capabilities. For OVEN [47] & InfoSeek [48], we follow the proposed  $224 \times 224$  resolution settings for fair comparison.

	MSR-VTT		Activity-Net		VATEX	SMIT	NExT-QA
Method	Cap. [55]	QA [60]	Cap. [57]	QA [61]	Cap. [56]	Cap. [58]	QA [59]
Prior SOTA	75.9	<b>47.4</b>	52.5	44.7	94.0 <sup>†</sup>	28.1 <sup>‡</sup>	33.5 <sup>§</sup>
	GIT2 [9]	Flamingo [10]	PDVC [62]	VINLU [63]	GIT2 [9]	MV-GPT [64]	Flamingo 32shot [10]
PaLI-X (8fr)	74.6	46.9	49.0	48.4	66.0	42.5	37.0
PaLI-X (16fr)	<b>76.8</b>	<b>47.1</b>	<b>54.9</b>	<b>49.4</b>	69.3	<b>43.5</b>	<b>38.3</b>

Table 5: Results for Video Captioning and Video-QA using 8 frames (8fr) or 16 frames (16fr). <sup>†</sup>GIT2 uses Self-Critical Sequence Training to directly optimize the CIDEr metric for VATEX. <sup>‡</sup>SMIT has not been used for video captioning before, we apply MV-GPT [64] and report results on the test set. <sup>§</sup>Numbers were obtained using 32-shot; since Flamingo 32-shot outperforms fine-tuning SOTA on this open-ended QA task, they did not conduct further fine-tuning experiments for this task.

**prompt:** detect knife and left half of watermelon and right half of watermelon and right hand and left hand  
**output:** 298 284 722 493 left half of watermelon and 298 483 721 681 right half of watermelon and 302 188 779 260 knife and 202 265 535 493 left hand and 201 485 546 723 right hand



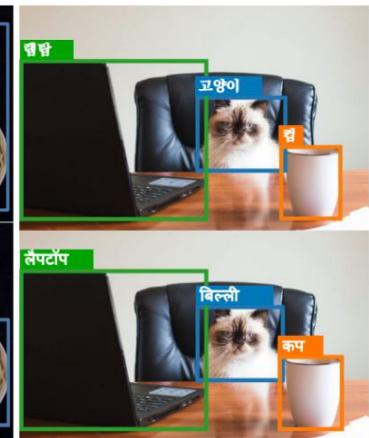
**prompt:** detect window and cat and bike and door and tree and this is paradise  
**output:** 259 487 423 550 window and 259 839 424 904 window and 706 625 787 778 bike and 550 839 759 942 door and 256 363 424 430 window and 258 241 422 305 window and 254 715 422 774 window and 254 115 422 177 window and 590 120 653 811 this is paradise and 549 131 588 161 vent



**prompt 1:** detect bowls  
**output:** 193 56 827 980 bowls  
**prompt 2:** detect bowl  
**output:** 213 55 594 431 bowl and 192 429 536 772 bowl and 511 322 822 630 bowl and 467 630 826 979 bowl



**prompt (KO):** detect 랩탑 and 컵 and 고양이  
**output:** 489 447 779 679 고양이 and 676 659 962 818 컵 and 349 0 964 478 랩탑  
**prompt (HI):** detect लैपटॉप and कप and बिल्ली  
**output:** 489 447 779 672 बिल्ली and 679 656 965 822 कप and 352 0 966 479 लैपटॉप



*Credits: Watermelon/Cat; Sarah Pflug (burst), Bowls; ariesandrea (flickr), Wall; Matthew Henry (burst)*

Figure 2: Examples demonstrating multilingual, OCR and other capabilities transferred to detection.

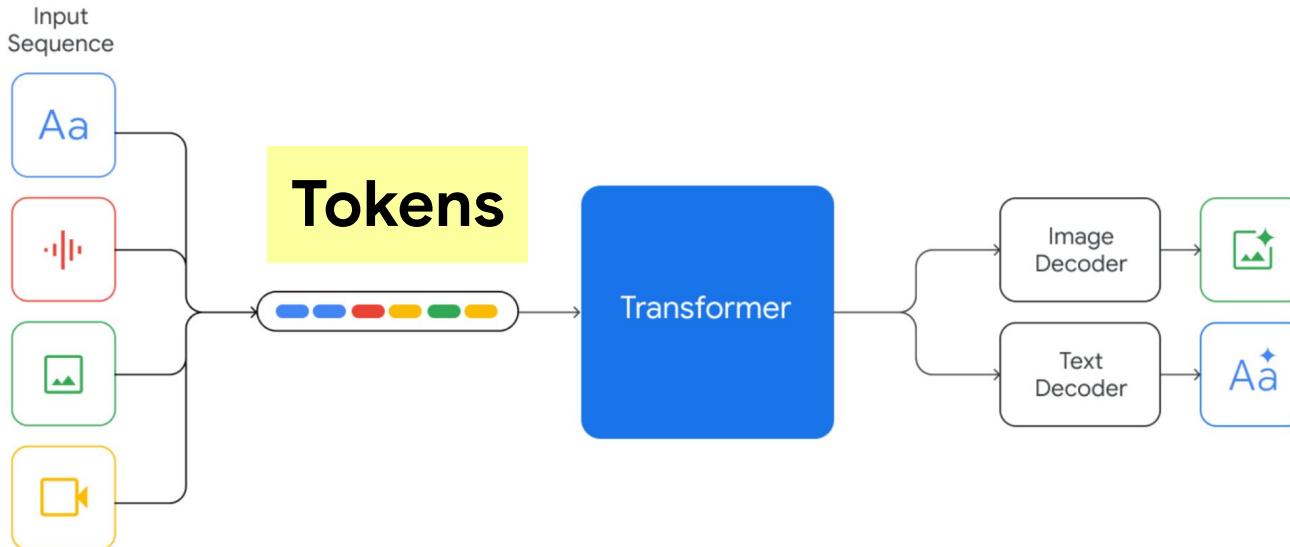
# Summary

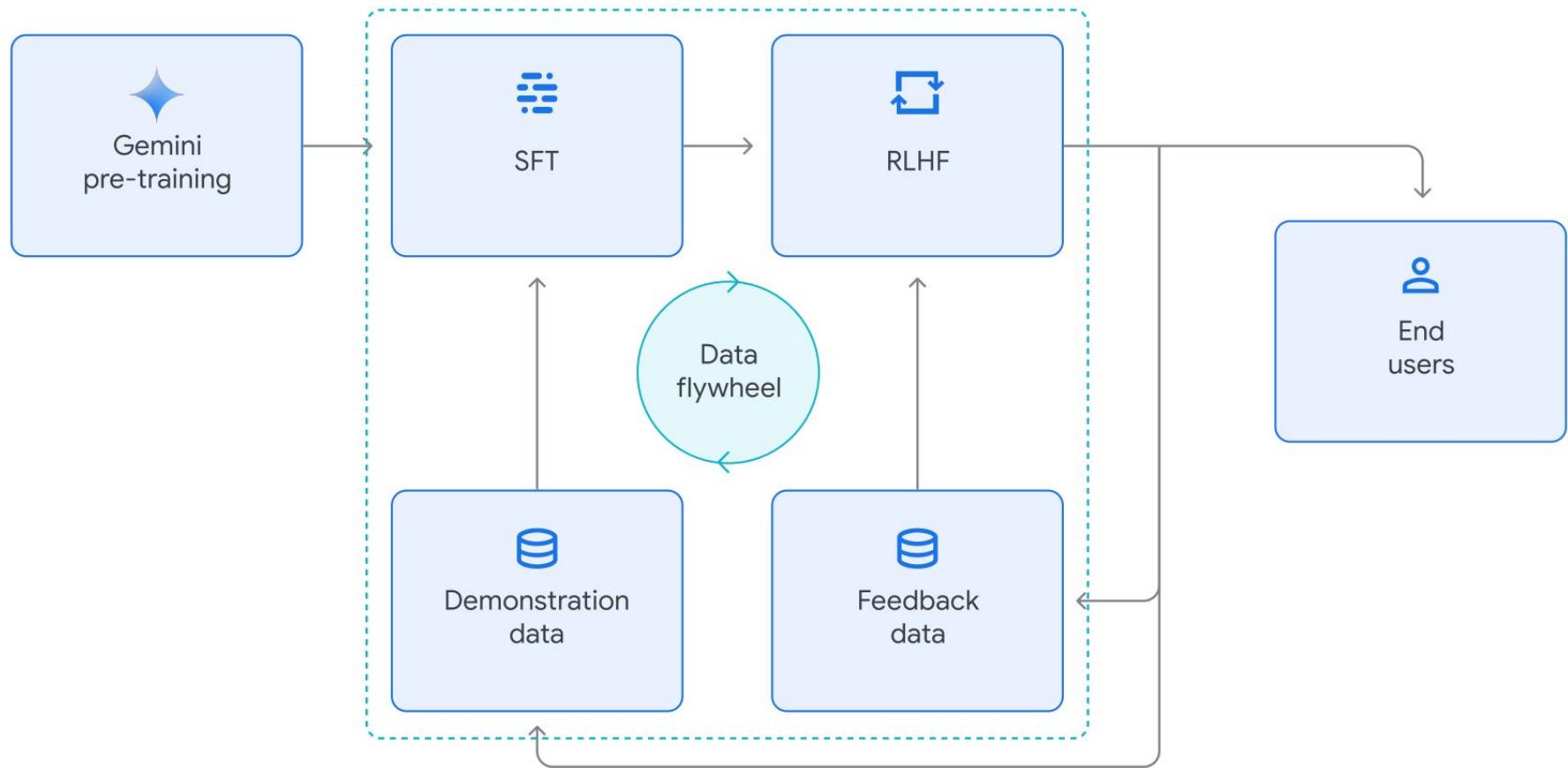
- **Vision+Language** took whatever worked best **at the moment**
- **Extensive engineering** goes into modeling (bridging the gaps between modalities), data, and evaluation.
- **Data scale and diversity** are important
- **Metrics for “hill-climbing”** are important

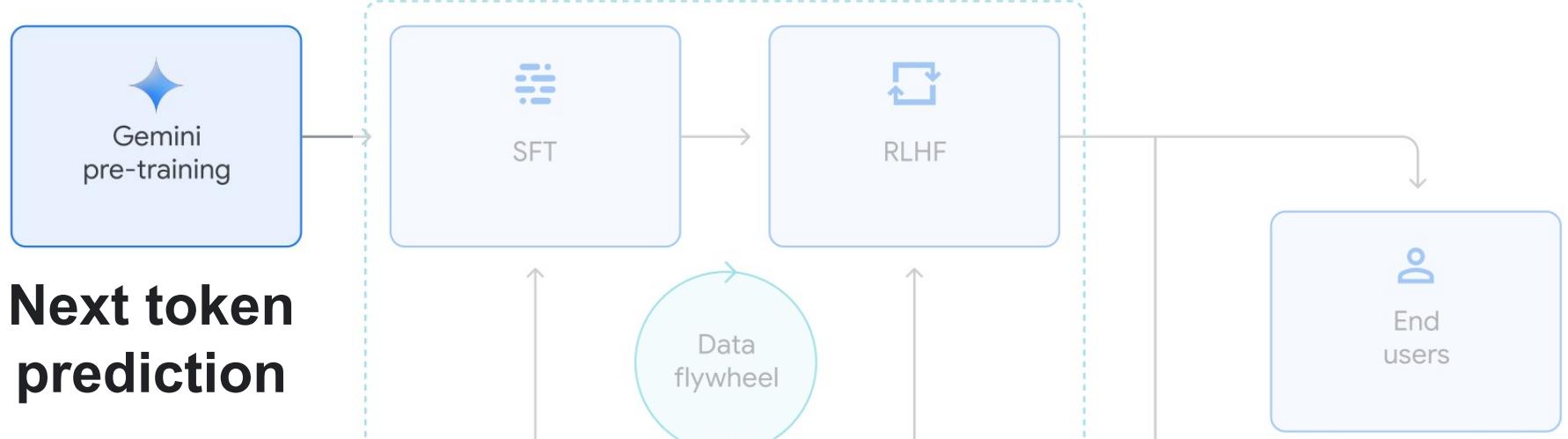
# **Gemini**

# Gemini: A Family of Highly Capable Multimodal Models

Gemini Team, Google<sup>1</sup>







## Next token prediction

### 4. Pre-Training Dataset

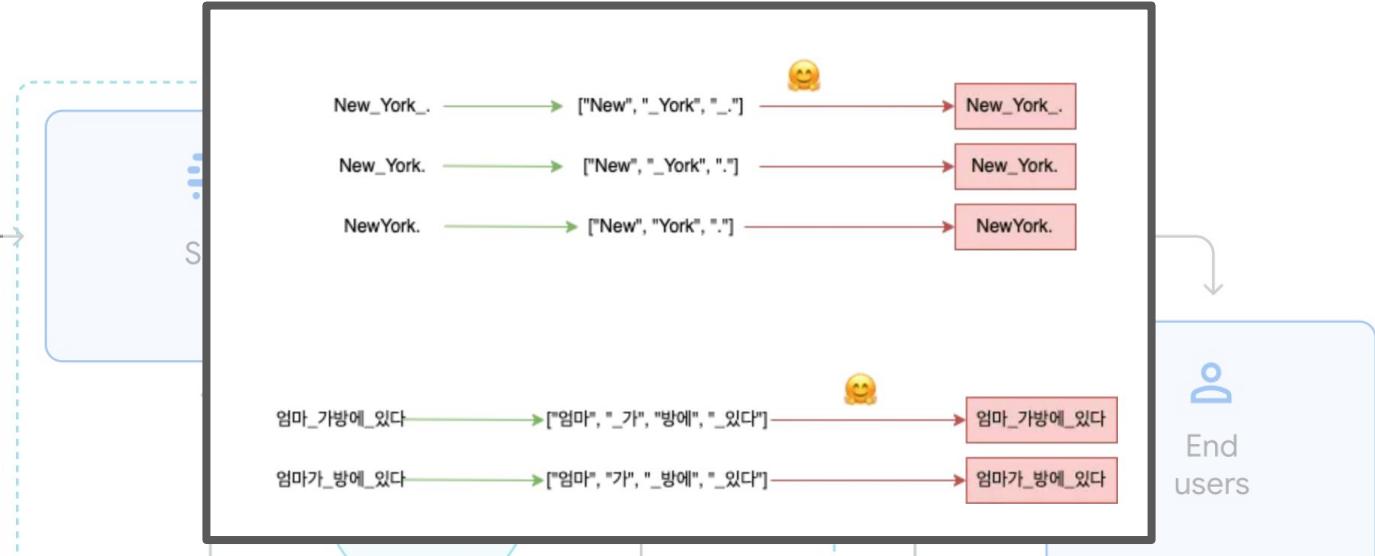
Gemini models are trained on a dataset that is both multimodal and multilingual. Our pre-training dataset uses data from web documents, books, and code, and includes image, audio, and video data.



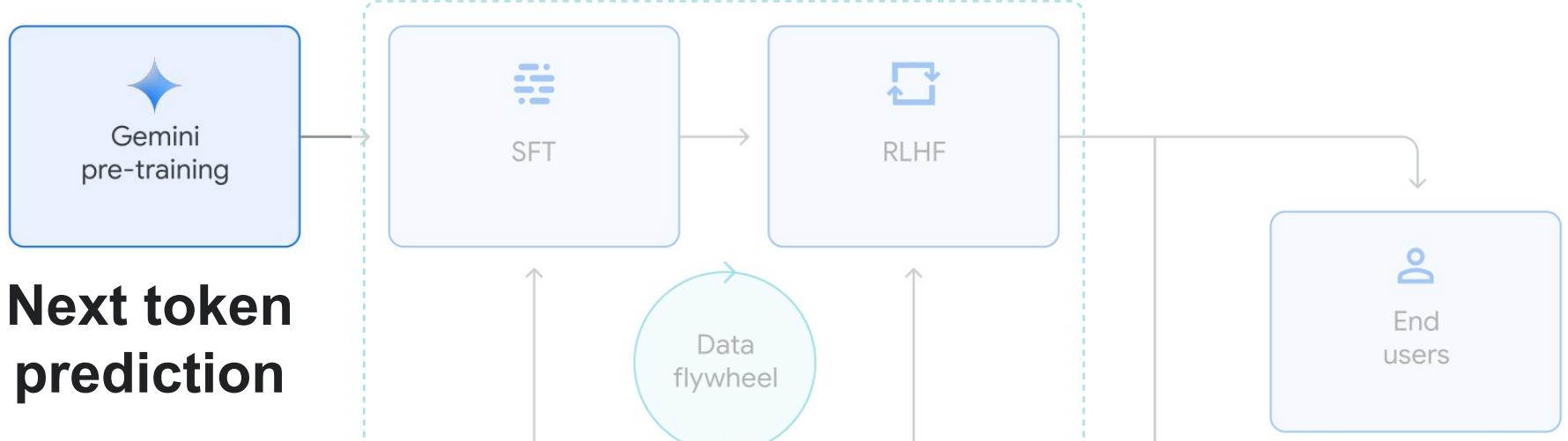
Gemini  
pre-training

## Next token prediction

We use the SentencePiece tokenizer ([Kudo and Richardson, 2018](#)) and find that training the tokenizer on a large sample of the entire training corpus improves the inferred vocabulary and subsequently improves model performance. For example, we find Gemini models can efficiently tokenize non-Latin scripts which can, in turn, benefit model quality as well as training and inference speed.



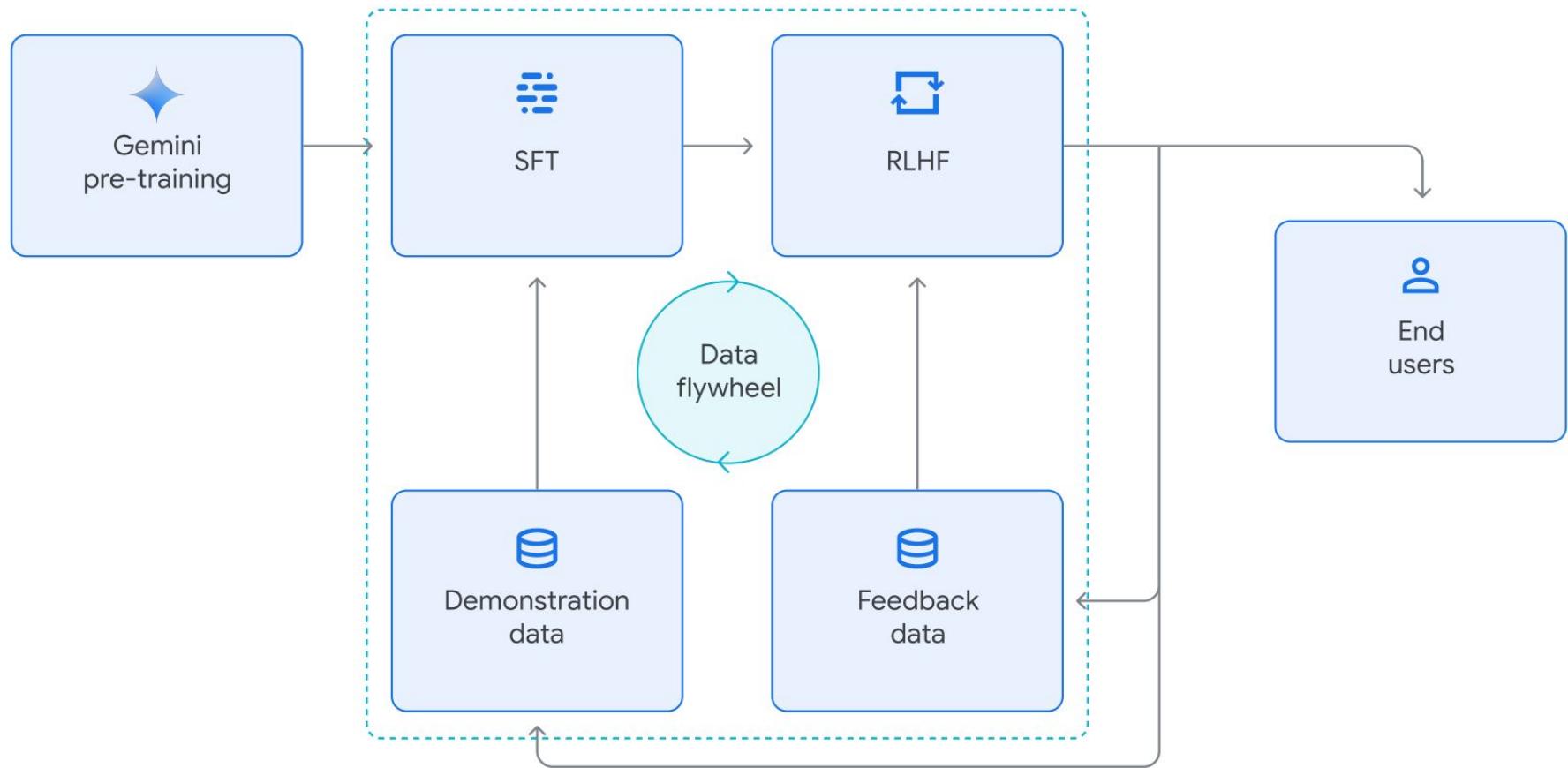
End users

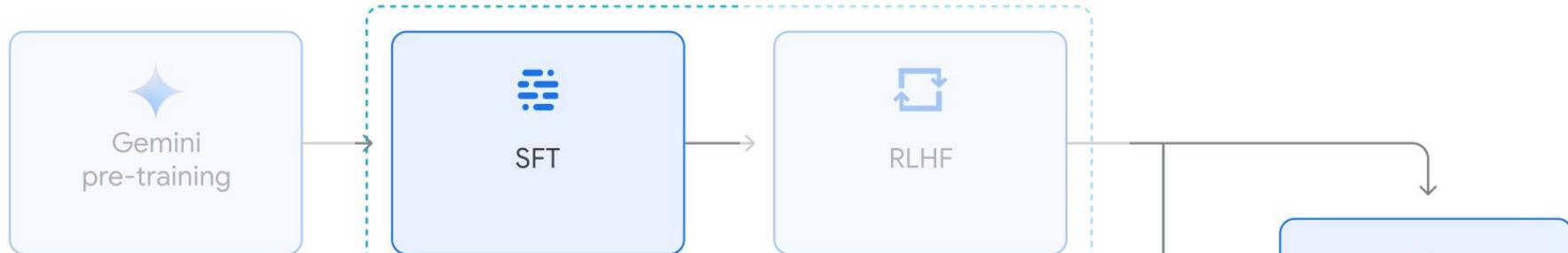


## Next token prediction

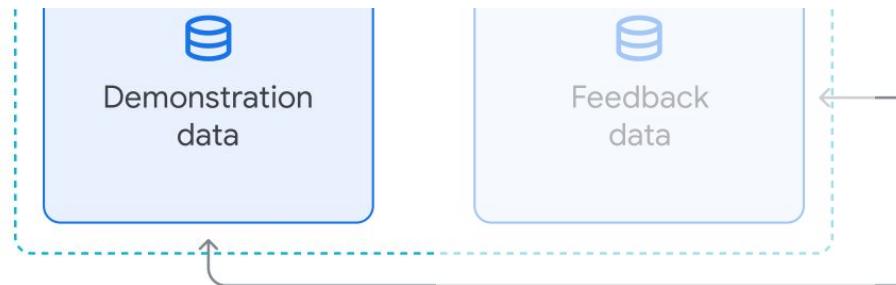
The number of tokens used to train the largest models were determined following the approach in [Hoffmann et al. \(2022\)](#). The smaller models are trained for significantly more tokens to improve performance for a given inference budget, similar to the approach advocated in [Touvron et al. \(2023a\)](#).







**(2) SFT on Demonstration Data:** SFT trains the model to output a desired target response given a prompt. Our Demonstration Data target responses can be directly written by a human expert, or generated by a model and in some cases revised or reviewed by a human. Additionally, we use data analysis tools and heuristics to ensure high data diversity across capabilities, use cases, and semantic clusters.



# Instruction Tuning

Input: text

Is the following sentence  
acceptable?  
“The course is jumping well.”

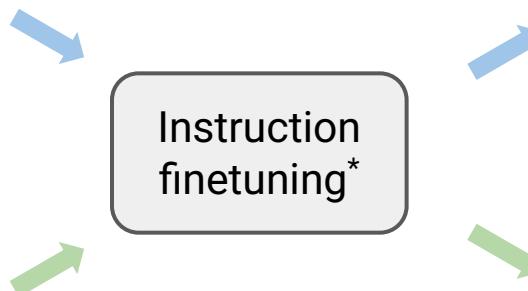
On the scale of 1 to 5, how similar  
are the following two sentences?  
  
1. The rhino grazed on the grass.  
2. A rhino is grazing in a field.

Output: text

“It is not  
acceptable”

“3.8”

Instruction  
finetuning\*



Wei et al., 2021  
Sanh et al. 2021

Slide Credit: [Hyung Won Chung](#)

# Finetuning tasks

## TO-SF

Commonsense reasoning  
Question generation  
Closed-book QA  
Adversarial QA  
Extractive QA  
Title/context generation  
Topic classification  
Struct-to-text  
...

**55 Datasets, 14 Categories, 193 Tasks**

## Muffin

Natural language inference  
Code instruction gen.  
Program synthesis  
Dialog context generation  
Closed-book QA  
Conversational QA  
Code repair  
...

**69 Datasets, 27 Categories, 80 Tasks**

## CoT (Reasoning)

Arithmetic reasoning	Explanation generation
Commonsense Reasoning	Sentence composition
Implicit reasoning	...

**9 Datasets, 1 Category, 9 Tasks**

## Natural Instructions v2

Cause effect classification  
Commonsense reasoning  
Named entity recognition  
Toxic language detection  
Question answering  
Question generation  
Program execution  
Text categorization  
...

**372 Datasets, 108 Categories, 1554 Tasks**

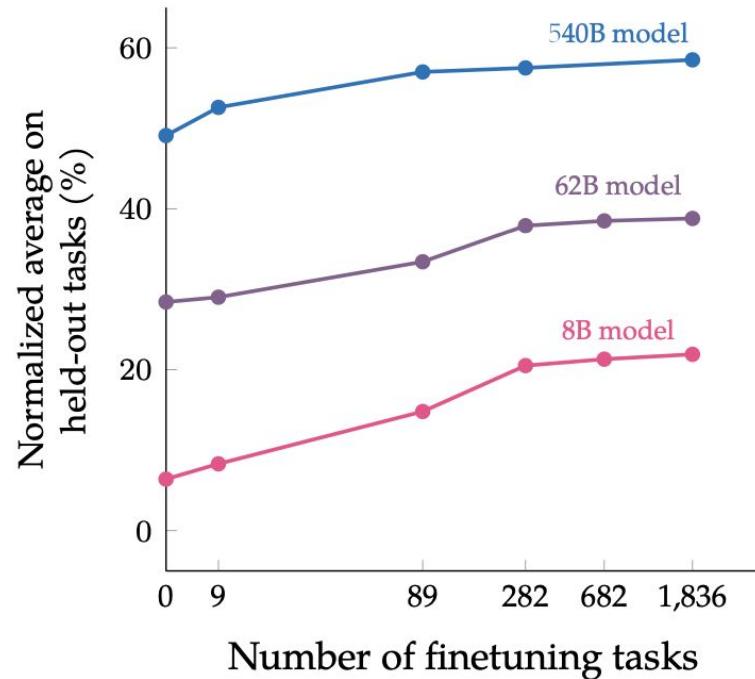
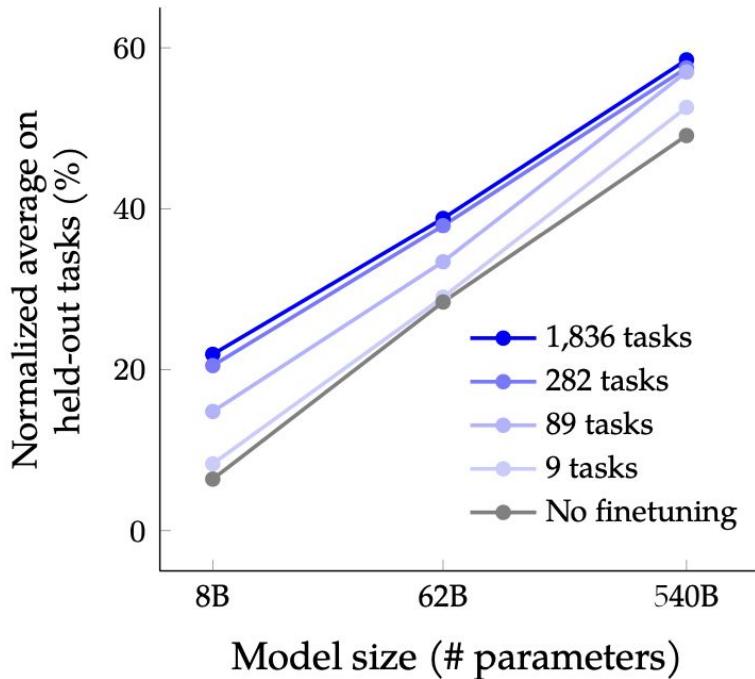
- ❖ A Dataset is an original data source (e.g. SQuAD).
- ❖ A Task Category is unique task setup (e.g. the SQuAD dataset is configurable for multiple task categories such as extractive question answering, query generation, and context generation).
- ❖ A Task is a unique <dataset, task category> pair, with any number of templates which preserve the task category (e.g. query generation on the SQuAD dataset.)

## Instruction finetuning on 1836 (!! academic tasks

Chung et al., 2022

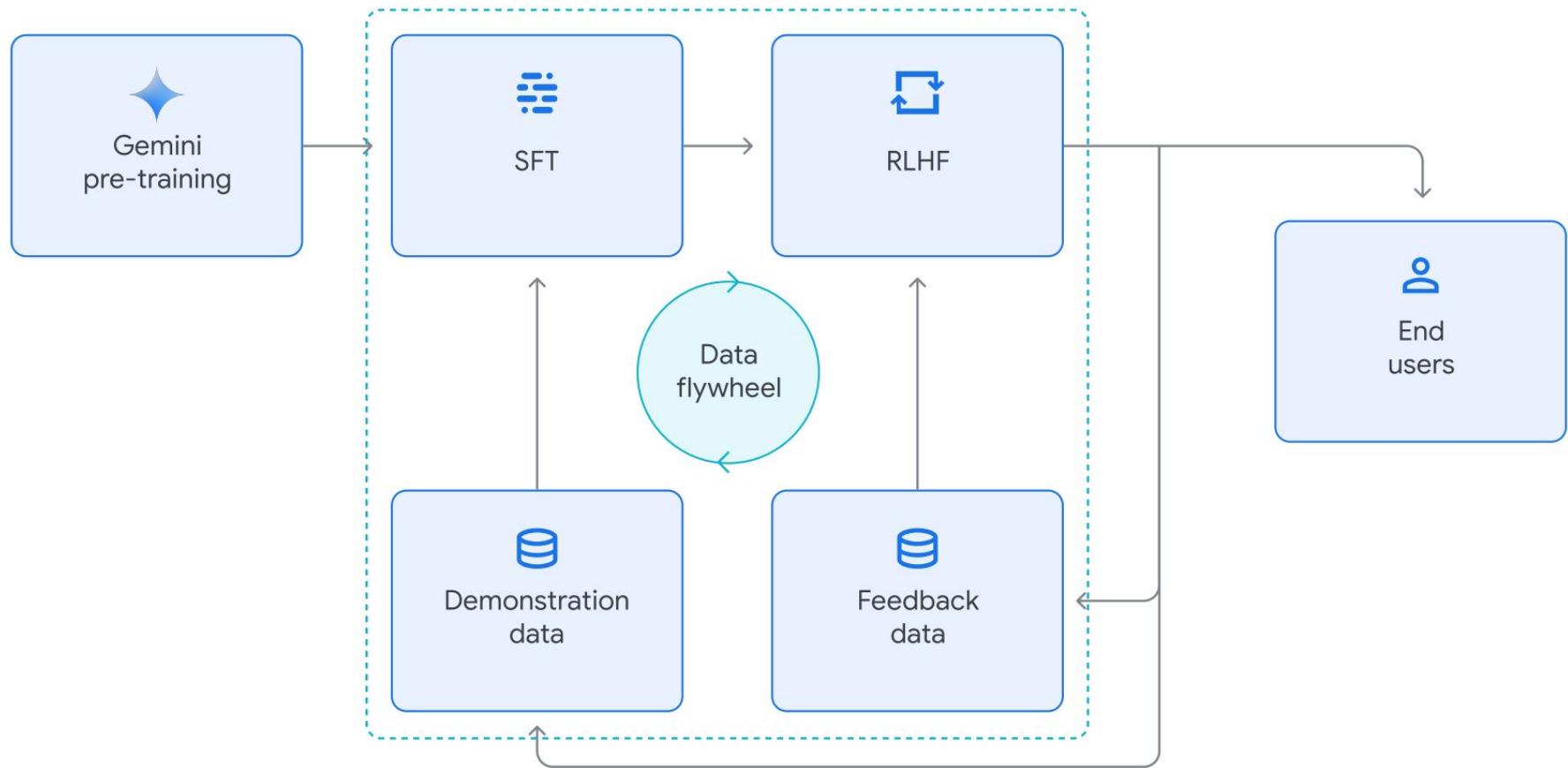
Slide Credit: [Hyung Won Chung](#)

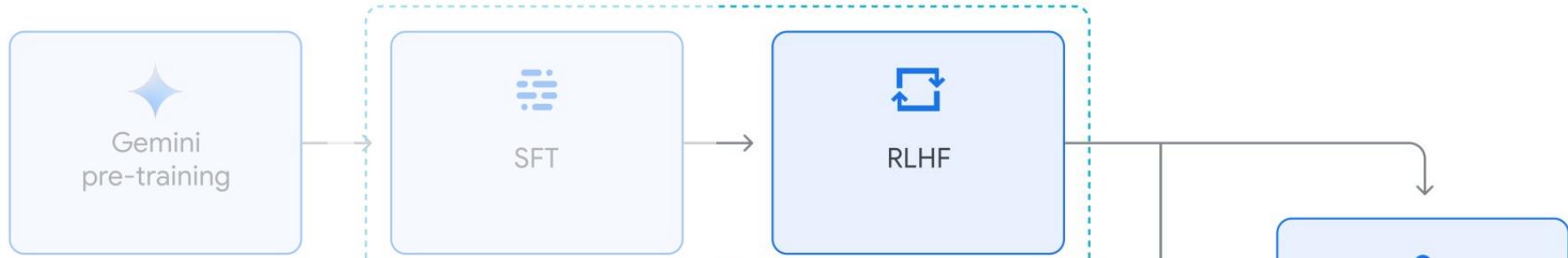
# Scaling the number of tasks and model size improves the



Chung et al., 2022

Slide Credit: [Hyung Won Chung](#)





**(4) RLHF:** Applying reinforcement learning from human feedback (RLHF) to our models provides further gains over SFT alone. Our approach creates an iterative process in which RL continually pushes the boundaries of the RM, while the RM is continuously improved through evaluation and data collection, leading to progressive improvements in both.



# Reward Model (RM)

## (Relative Preference)

**Input**

Explain the moon landing to a 6 year old in a few sentences

**Completion 1**

The Moon is a natural satellite of the Earth. It is the fifth largest moon in the Solar System and the largest relative to the size of its host planet.



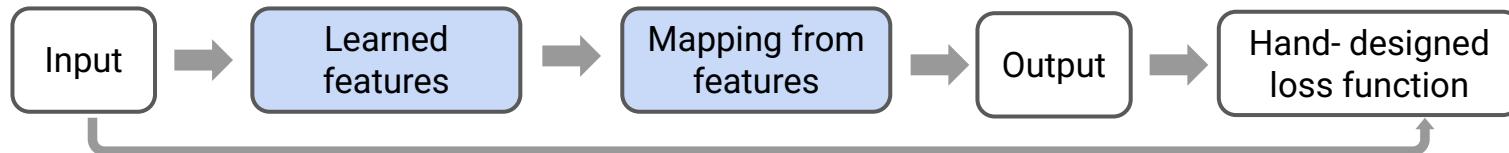
**Completion 2**

People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

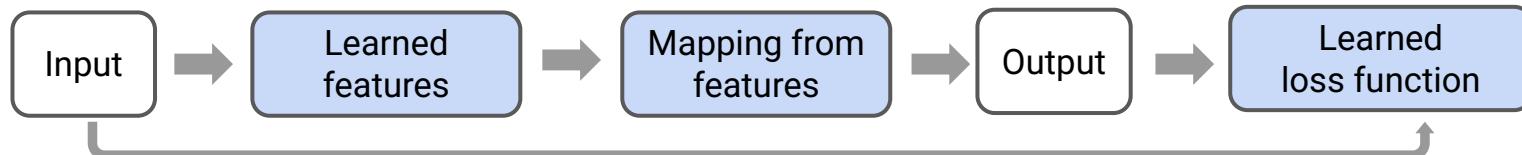
# Why RLHF?

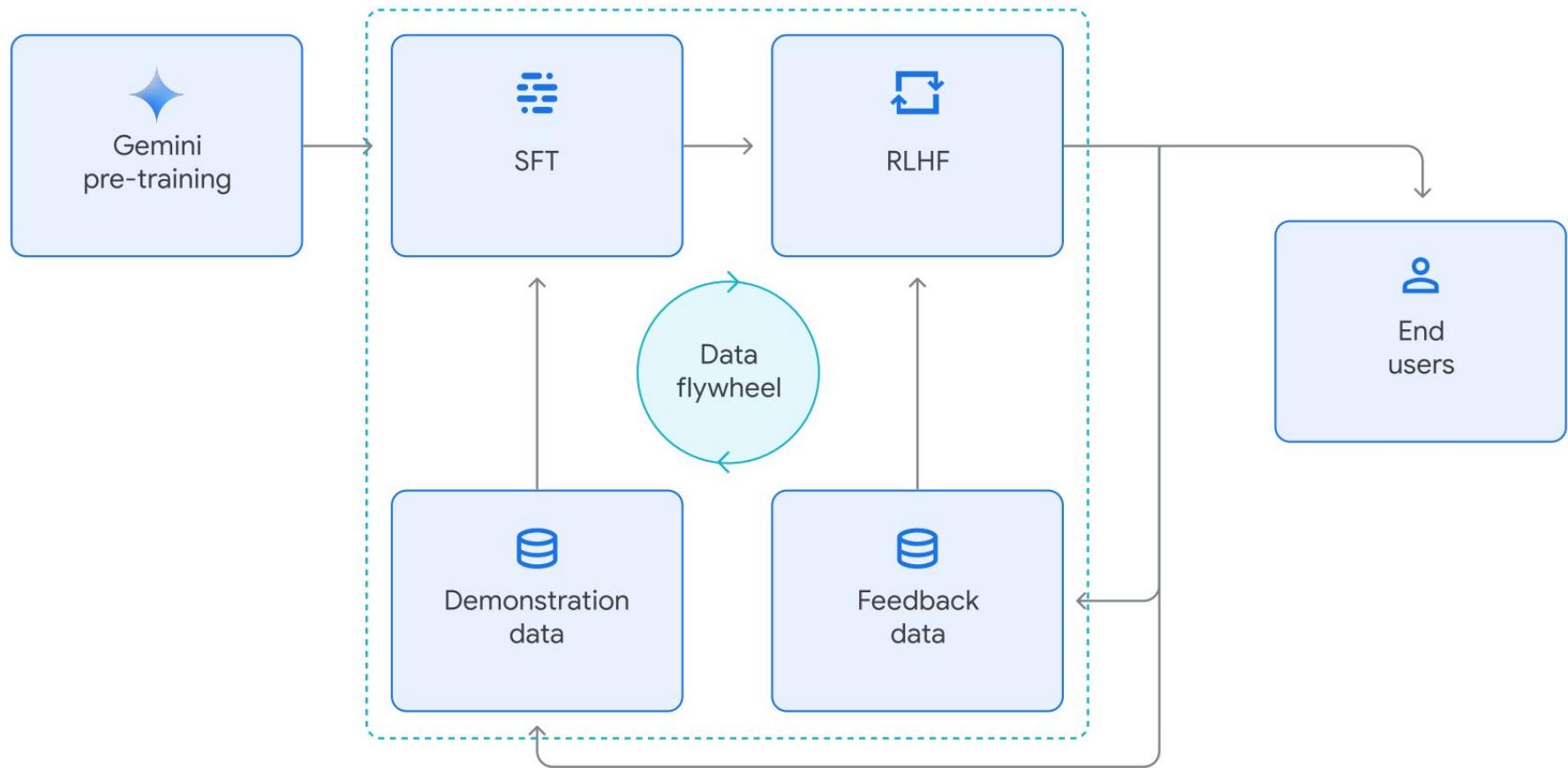
Learnable part of  
the system

Deep learning: (self-)supervised learning



Deep learning: other RL formulations





# Gemini Team Organization

## Gemini Leads

Rohan Anil, *Co-Lead, Text*  
Sebastian Borgeaud, *Co-Lead, Text*  
Jean-Baptiste Alayrac, *Co-Lead, MM Vision*  
Jiahui Yu, *Co-Lead, MM Vision*  
Radu Soricuț, *Co-Lead, MM Vision*  
Johan Schalkwyk, *Lead, MM Audio*  
Andrew M. Dai, *Co-Lead, Data*  
Anja Hauth, *Co-Lead, Data*  
Katie Millican, *Co-Lead, Data*  
David Silver, *Co-Lead, Fine-Tuning*  
Melvin Johnson, *Lead, Instruction Tuning*  
Ioannis Antonoglou, *Co-Lead, RL Techniques*  
Julian Schrittwieser, *Co-Lead, RL Techniques*  
Amelia Glaese, *Lead, Human Data*  
Jilin Chen, *Lead, Safety*  
Emily Pitler, *Co-Lead, Tool Use*  
Timothy Lillicrap, *Co-Lead, Tool Use*  
Angeliki Lazaridou, *Co-Lead, Eval*  
Orhan Firat, *Co-Lead, Eval*  
James Molloy, *Co-Lead, Infra*  
Michael Isard, *Co-Lead, Infra*  
Paul R. Barham, *Co-Lead, Infra*  
Tom Henigan, *Co-Lead, Infra*  
Benjamin Lee, *Co-Lead, Codebase & Parallelism*  
Fabio Viola, *Co-Lead, Codebase & Parallelism*  
Malcolm Reynolds, *Co-Lead, Codebase & Parallelism*  
Yuanzhong Xu, *Co-Lead, Codebase & Parallelism*  
Ryan Doherty, *Lead, Ecosystem*  
Eli Collins, *Lead, Product*  
Clemens Meyer, *Co-Lead, Operations*  
Eliza Rutherford, *Co-Lead, Operations*  
Erica Moreira, *Co-Lead, Operations*  
Kareem Ayoub, *Co-Lead, Operations*  
Megha Goel, *Co-Lead, Operations*

- Text
- Multimodal Vision
- Multimodal Audio
- Data
- Fine-Tuning
- Instruction-Tuning
- RL Techniques
- Human Data
- Safety
- Tool Use
- Eval
- Infra
- Codebase & Parallelism
- Ecosystem
- Product
- Operations

# Introducing Gemini

*By Demis Hassabis, CEO and Co-Founder of Google DeepMind, on behalf of the Gemini team*

AI has been the focus of my life's work, as for many of my research colleagues. Ever since programming AI for computer games as a teenager, and throughout my years as a neuroscience researcher trying to understand the workings of the brain, I've always believed that if we could build smarter machines, we could harness them to benefit humanity in incredible ways.

This promise of a world responsibly empowered by AI continues to drive our work at Google DeepMind. For a long time, we've wanted to build a new generation of AI models, inspired by the way people understand and interact with the world. AI that feels less like a smart piece of software and more like something useful and intuitive — an expert helper or assistant.

Today, we're a step closer to this vision as [we introduce Gemini](#), the most capable and general model we've ever built.

Gemini is the result of large-scale collaborative efforts by teams across Google, including our colleagues at Google Research. It was built from the ground up to be multimodal, which means it can generalize and seamlessly understand, operate across and combine different types of information including text, code, audio, image and video.

<https://blog.google/technology/ai/google-gemini-ai/#introducing-gemini>

# Introducing Gemini

*By Demis Hassabis, CEO and Co-Founder of Google DeepMind, on behalf of the Gemini team*

AI has been the focus of my life's work, as for many of my research colleagues. Ever since programming AI for computer games as a teenager, and throughout my years as a neuroscience researcher trying to understand the workings of the brain, I've always believed that if we could build smarter machines, we could harness them to benefit humanity in incredible ways.

This promise of a world responsibly empowered by AI continues to drive our work at Google DeepMind. For a long time, we've wanted to build a new generation of AI models, inspired by the way people understand and interact with the world. AI that feels less like a smart piece of software and more like something useful and intuitive — an expert helper or assistant.

Today, we're a step closer to this vision as [we introduce Gemini](#), the most capable and general model

~~Gemini is the result of large scale collaborative efforts by teams across Google, including our~~

colleagues at Google Research. It was built from the ground up to be multimodal, which means it can generalize and seamlessly understand, operate across and combine different types of information including text, code, audio, image and video.

# Multimodal

- *Data at all stages*



[BradyFU / Awesome-Multimodal-Large-Language-Models](#) (Public)

[Notifications](#)

[Fork 779](#)

[Star 12.2k](#)

- [Awesome Datasets](#)
  - [Datasets of Pre-Training for Alignment](#)
  - [Datasets of Multimodal Instruction Tuning](#)
  - [Datasets of In-Context Learning](#)
  - [Datasets of Multimodal Chain-of-Thought](#)
  - [Datasets of Multimodal RLHF](#)
  - [Benchmarks for Evaluation](#)
  - [Others](#)

# Multimodal

- *Data at all stages*  
*Example: Llama 3*

## 7.1.1 Image Data

Our image encoder and adapter are trained on image-text pairs. We construct this dataset via a complex data processing pipeline that consists of four main stages: **(1)** quality filtering, **(2)** perceptual de-duplication, **(3)** resampling, and **(4)** optical character recognition. We also apply a series of safety mitigations.

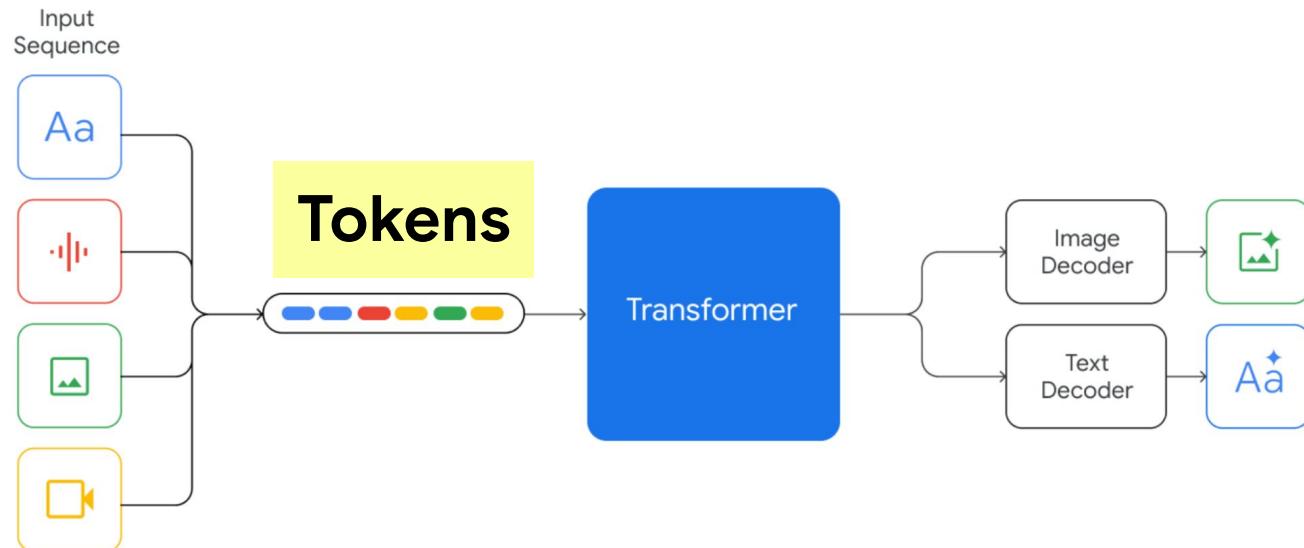
**Transcribing documents.** To improve the performance of our models on document understanding tasks, we render pages from documents as images and paired the images with their respective text. The document text is obtained either directly from the source or via a document parsing pipeline.

**Annealing data.** We create an annealing dataset by resampling the image-caption pairs to a smaller volume of ~350M examples using n-grams. Since the n-grams resampling favor richer text descriptions, this selects a higher-quality data subset. We augment the resulting data with ~150M examples from five additional sources:

- **Visual grounding.** We link noun phrases in the text to bounding boxes or masks in the image. The grounding information (bounding boxes and masks) are specified in the image-text pair in two ways. (1) We overlay boxes or masks with marks on the image and use marks in the text as reference, akin to set-of-marks (Yang et al., 2023a). (2) We insert normalized  $(x_{\min}, y_{\min}, x_{\max}, y_{\max})$  coordinates directly into the text, demarcated by special tokens.
- **Screenshot parsing.** We render screenshots from HTML code and task the model with predicting the code that produced a specific element in the screenshot, akin to Lee et al. (2023). The element of interest is indicated in the screenshot via a bounding box.
- **Question-answer pairs.** We include question-answer pairs, enabling us to use volumes of question-answering data that are too large to be used in model finetuning.
- **Synthetic captions.** We include images with synthetic captions that were generated by an early version of the model. Compared to original captions, we find that synthetic captions provide a more comprehensive description of images than the original captions.
- **Synthetically-generated structured images.** We also include synthetically generated images for a variety of domains such as charts, tables, flowcharts, math equations and textual data. These images are accompanied by a structured representation such as the corresponding markdown or LaTeX notation. Besides improving recognition capabilities of the model for these domains, we find this data useful to generate question-answer pairs via the text model for finetuning.

# Multimodal

- *How to process multimodal input?*



# Multimodal

- *How to process multimodal input?*

Here's how tokens are calculated for images:

- **Gemini 1.0 Pro Vision:** Each image accounts for 258 tokens.
- **Gemini 1.5 Flash and Gemini 1.5 Pro:**
  - If both dimensions of an image are less than or equal to 384 pixels, then 258 tokens are used.
  - If one dimension of an image is greater than 384 pixels, then the image is cropped into tiles. Each tile size defaults to the smallest dimension (width or height) divided by 1.5. If necessary, each tile is adjusted so that it's not smaller than 256 pixels and not greater than 768 pixels. Each tile is then resized to 768x768 and uses 258 tokens.

<https://cloud.google.com/vertex-ai/generative-ai/docs/multimodal/image-understanding#image-requirements>

# Multimodal

- *Image generation?*
- *Video? Audio? Code?*
- *Interleaved multimodal inputs?*
- *Interleaved multimodal outputs?*

# Efficiency

- *Model architecture*
- *Parallelization and Hardware*
- *Precision*

# Efficiency

- *Model architecture*

## 3. Model Architecture

### 3.1. Gemini 1.5 Pro

Gemini 1.5 Pro is a sparse mixture-of-expert (MoE) Transformer-based model that builds on Gemini 1.0's (Gemini-Team et al., 2023) research advances and multimodal capabilities. Gemini 1.5 Pro also builds on a much longer history of MoE research at Google (Clark et al., 2022; Du et al., 2022; Fedus et al., 2021; Lepikhin et al., 2020; Riquelme et al., 2021; Shazeer et al., 2017; Zoph et al., 2022) and language model research in the broader literature (Anil et al., 2023b; Anthropic, 2023a; Brown et al., 2020; Chowdhery et al., 2023b; Hoffmann et al., 2022; Jiang et al., 2024; Kim et al., 2021; OpenAI, 2023a; Rae et al., 2021; Raffel et al., 2020; Roller et al., 2021; Thoppilan et al., 2022; Touvron et al., 2023a,b; Vaswani et al., 2017). MoE models use a learned routing function to direct inputs to a subset of the model's parameters for processing. This form of conditional computation (Bengio et al., 2013; Davis and Arel, 2014; Jacobs et al., 1991) allows models to grow their total parameter count while keeping the number of parameters that are activated for any given input constant.

Gemini Team, Google<sup>1</sup>

**Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context**

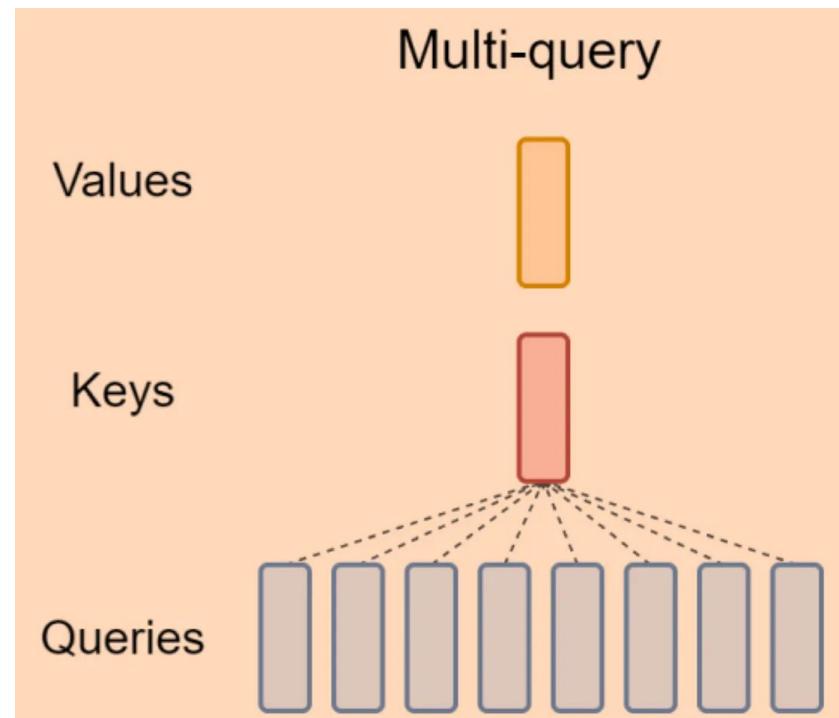
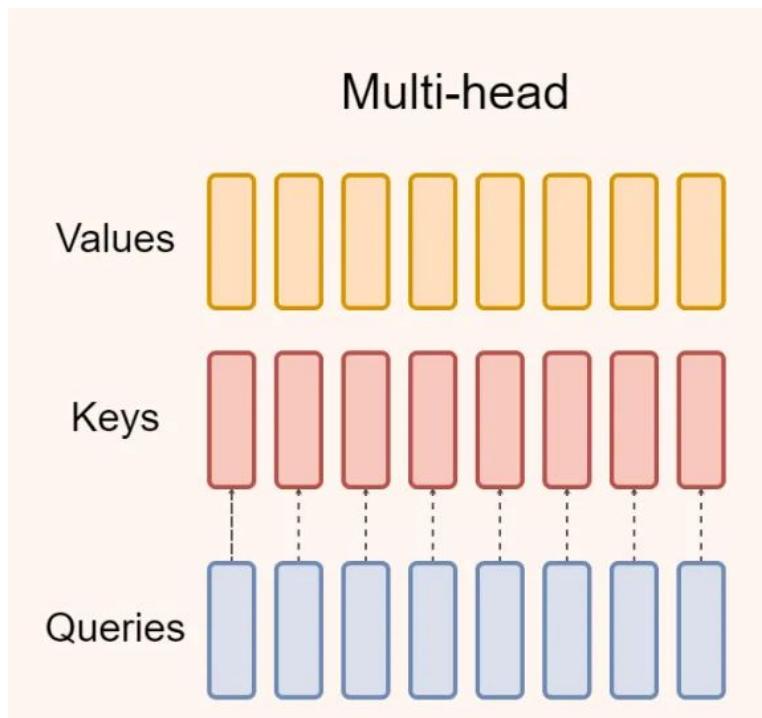
# Efficiency

- *Attention mechanisms*

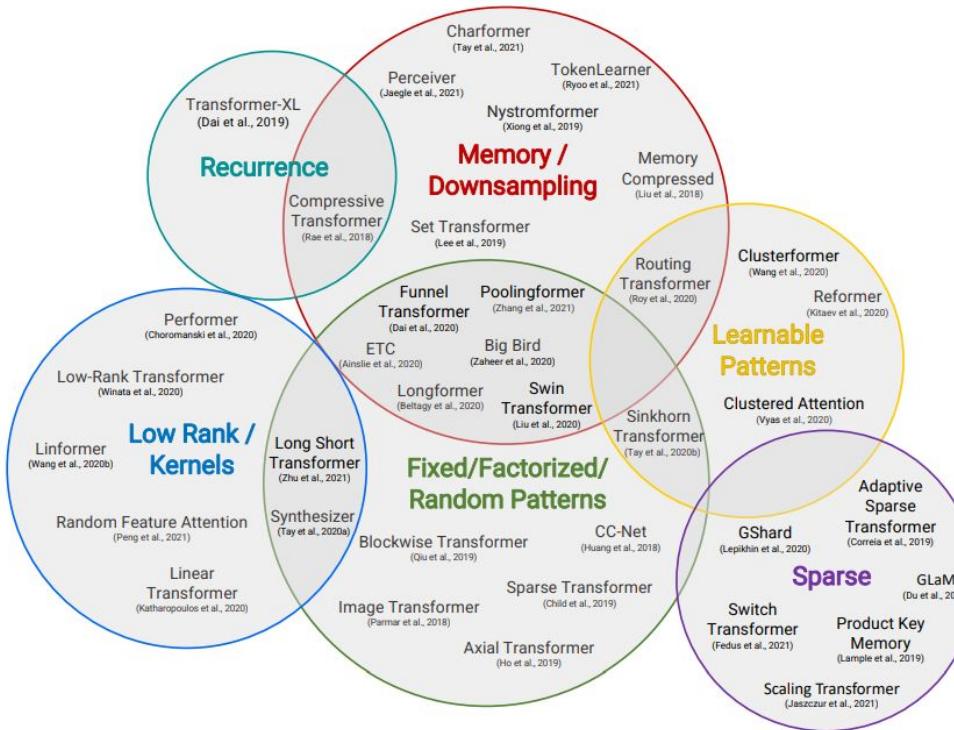
## 2. Model Architecture

Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google’s Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)).

# Multi-Query Attention



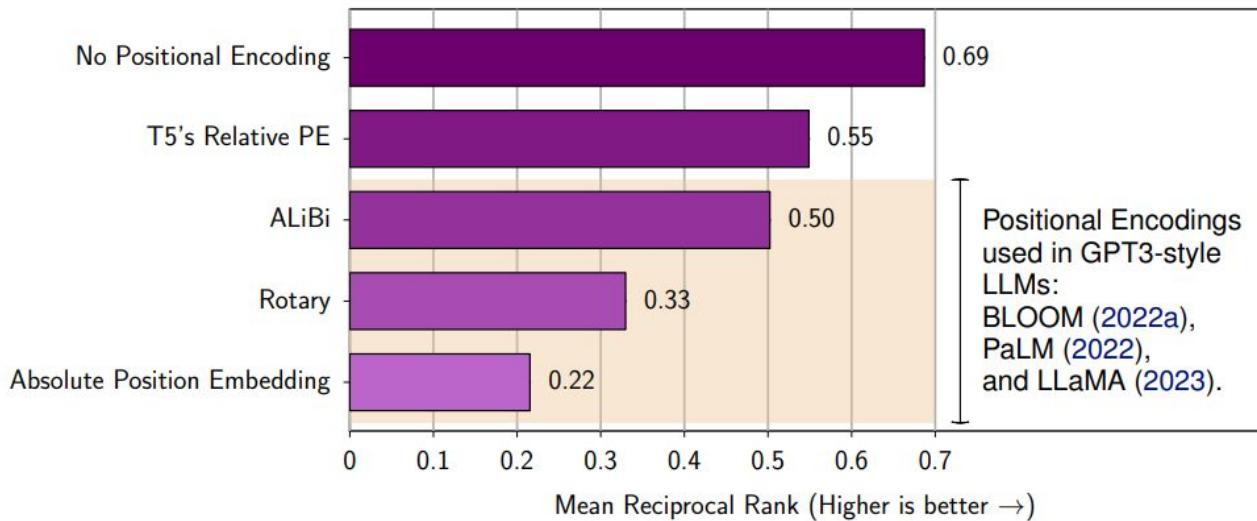
# Efficiency



Tay et al. 2022

# Efficiency

- *Positional Encoding & Length Generalization*



# Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context

Gemini Team, Google<sup>1</sup>

In practice, 1 million tokens would look like:

- 50,000 lines of code (with the standard 80 characters per line)
- All the text messages you have sent in the last 5 years
- 8 average length English novels
- Transcripts of over 200 average length podcast episodes

# Summary

- **Devil is in the detail:** infrastructure, architecture, data, evaluation. **Massive engineering!**
- **Multimodal brings further challenges:** data sources, how to bridge the gap between different modalities (e.g. tokenization, encoders, adaptors), how to process long inputs, etc.
- **R&D** are on-going on multimodal + X, where X = multilinguality, reasoning, coding, agent, etc.
- Will **enrich/automate existing pipelines** and **enable new applications**



Nando de Freitas ✅

@NandoDF

🔗 ...

RL is not all you need, nor attention nor Bayesianism nor free energy minimisation, nor an age of first person experience. Such statements are propaganda.

You need thousands of people working hard on data pipelines, scaling infrastructure, HPC, apps with feedback to drive benchmarks and data, tons of research and engineering on generative models, data mixtures, ablations, RL/selftraining, etc etc and we will probably need lots of people working hard to figure out safety, causal world models, awareness, models that create abstractions comparable to infinity and zero and use these to predict the existence of things like black holes and suggest experiments to verify such hypothesis, or come up with novel engineering designs to generate energy more efficiently, robotics, etc etc.

# **Applications (Brief)**

# Generative AI is a hammer and no one knows what is and isn't a nail



Colin Fraser · [Follow](#)

35 min read · Feb 22, 2024

---

1.3K

21



<https://medium.com/@colin.fraser/generative-ai-is-a-hammer-and-no-one-knows-what-is-and-isnt-a-nail-4c7f3f0911aa>

## Use cases

### Applications

Code assistant

Flutter code generator

#### Content search

Data exploration agent

Writing assistant

Slides reviewer

# Build an AI content search with Docs Agent

[Send feedback](#)

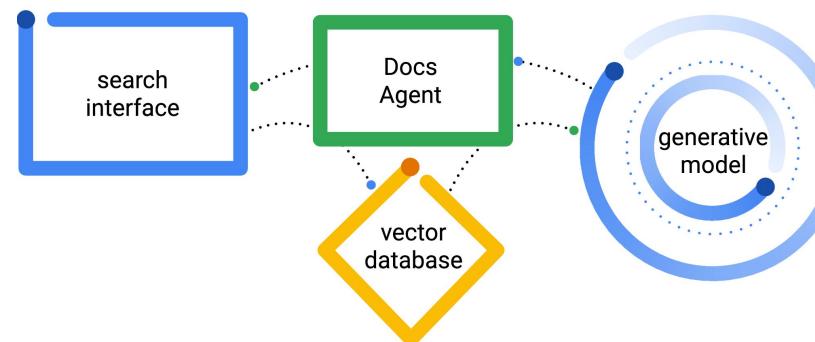
Searching for information is one of the most common uses of artificial intelligence (AI) generative models. Building a conversational search interface for your content using AI allows your users to ask specific questions and get direct answers.



**Note:** This example is updated for use with Gemini API.

This tutorial shows you how to build an AI-powered, conversational search interface for your content. It's based on [Docs Agent](#), an open source project that uses Google Gemini API to create a conversational search interface, *without training a new AI model or doing model tuning with Gemini models*. That means you can get this search capability built quickly and use it for small and large content sets.

For a video overview of the project and how to extend it, including insights from the folks who build it, check out: [AI Content Search | Build with Google AI](#). Otherwise you can get started extending the project following the instructions below.



<https://ai.google.dev/gemini-api/tutorials/docs-agent>



<https://www.youtube.com/watch?v=LTJb76UHuJg> Based on this video, what affects the performance the most?

Regenerate draft ⚡ 🔍

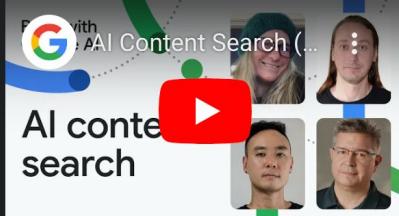


YouTube ▾

According to the video, splitting the text into meaningful chunks, such as headings and related paragraphs, is the most important thing for good performance.

YouTube

what affects the performance the most



AI Content Search

(RAG) with Docs ...

Google for Developers

Learn how to build an AI-powered conversational search interface for your content using the Googl...



Can you add some flowers to the table?

Type something



Run ⌘ ↵





# What will you build?

Push Gemini to the limits of what AI can do using the Gemini API



## Image Editing

Edit an image of croissants.



## Visual Story

Generate a story with images.



## Birthday Card

Design a custom birthday card.

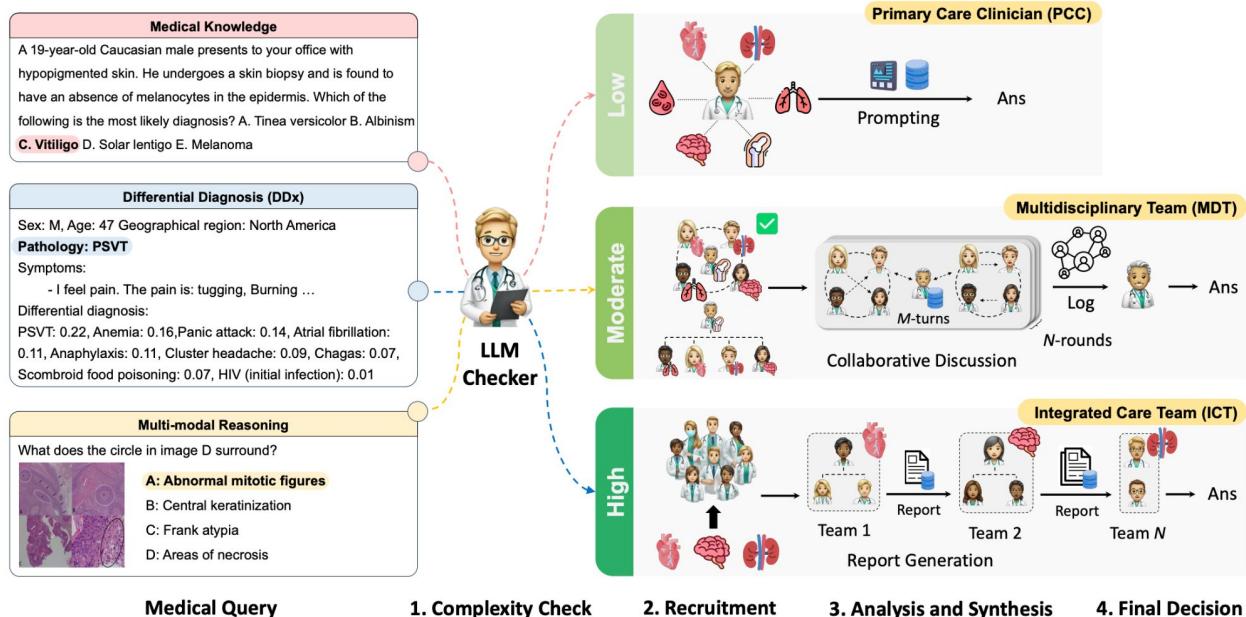
Give me a recipe for a chocolate chip cookie.



Run ↘↔



# Agents

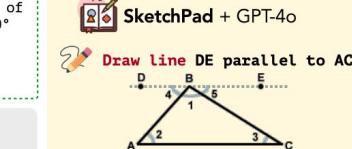
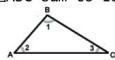


**Figure 1: Medical Decision-Making Agents (MDAgents) framework.** Given a medical query from different medical datasets, the framework performs 1) medical complexity check, 2) recruitment, 3) analysis and synthesis, and 4) decision-making steps.

# Agents

## Geometry

Prove the angles of  $\triangle ABC$  sum to  $180^\circ$



GPT-4o  
Start by assuming, for contradiction, that the sum is not  $180^\circ$ . We'll introduce a new angle,  $\angle 4$  ... X

(a)

## Math Function

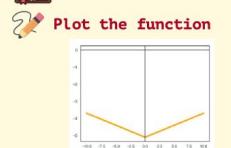
Is  $f(x)$  an convex function?

$$f(x) = 0.14 |x| - 5.09$$



It is convex for  $x > 0$  but concave for  $x < 0$  X

SketchPad + GPT-4o  
Plot the function



convex because line between any two points lie above the graph. ✓

(b)

## Visual Search

According to the board, what is the type of the shop?



The sign in the image reads "ceremonial matcha.". it's a cafe. X

SketchPad + GPT-4o  
Use sliding window to find the board



I see a board in patch 2 and 3 and will zoom in

The zoomed-in image shows the text "corepower YOGA". So it's a yoga studio ✓

(c)

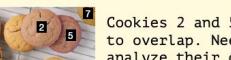
## Spatial Reasoning

Are there cookies stacked on other cookies?



No

SketchPad + GPT-4o  
Use segmentation & marking



Cookies 2 and 5 seem to overlap. Need to analyze their depth

Use depth estimation

Warmer colors indicate closer proximity, suggesting they are stacked on cooler-colored ones. ✓

(d)

# Conclusion

- Fundamentals & Evolution
- Multimodal LLMs: Modeling, Data, Evaluation
- Example: Gemini

# Conclusion

- Fundamentals & Evolution
- Multimodal LLMs: Modeling, Data, Evaluation
- Example: Gemini

**But this is just the beginning!**

- Learn, explore, research, experiment (perhaps with LLMs) to advance the sciences and engineering
- Tailor (multimodal) LLMs to your specific tasks or domains