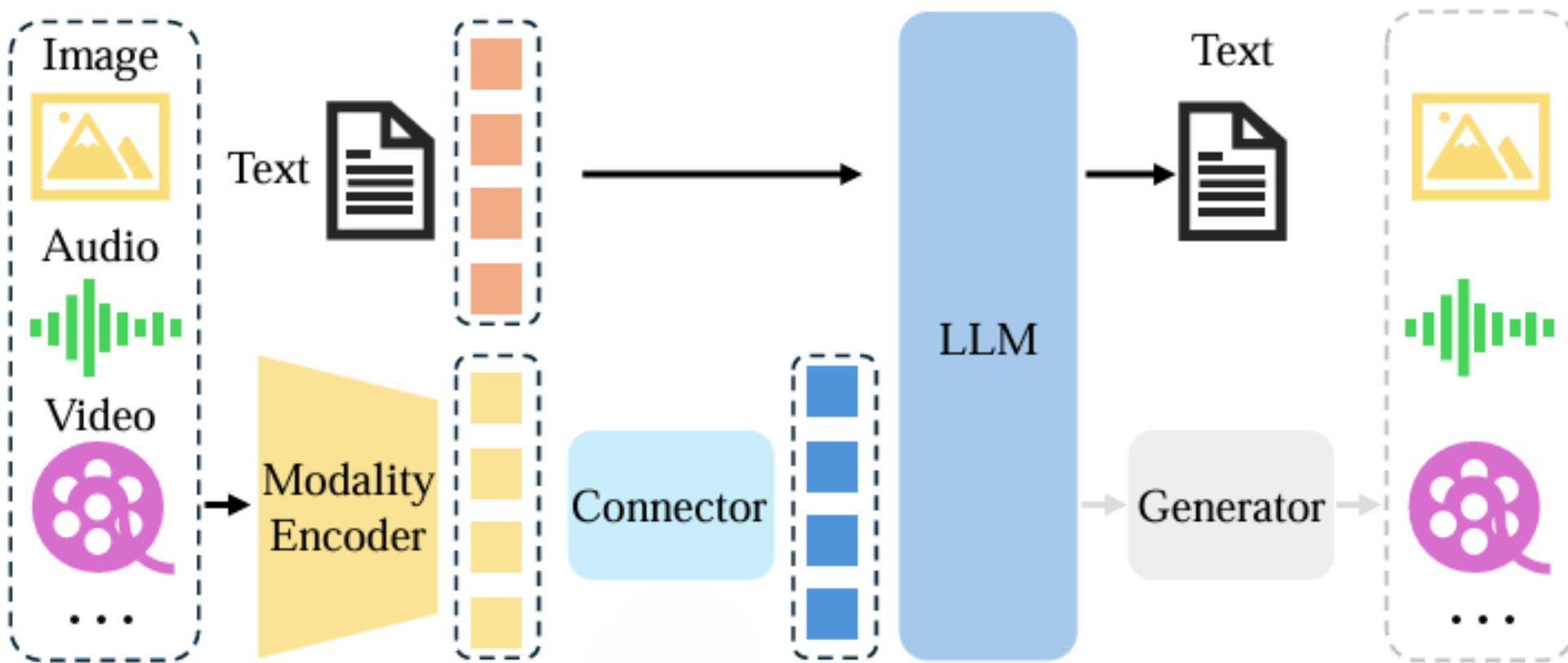


M-LLMs

On the way to AGI



Static Benchmarks: GQA



Pattern: What/Which <type> [do you think] <is> <dobject>, <attr> or <decoy>?

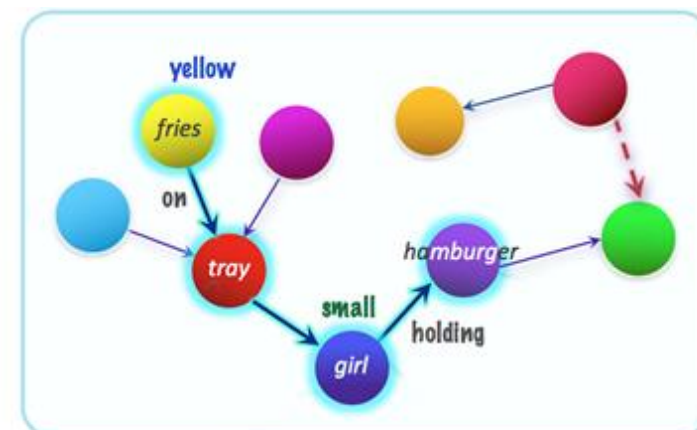
Program: Select: <dobject> → Choose <type>: <attr>|<decoy>

Reference: The food on the red object left of the small girl that is holding a hamburger

Decoy: brown

What color is the food on the red object left of the small girl that is holding a hamburger, yellow or brown?

Select: hamburger → Relate: girl, holding → Filter size: small → Relate: object, left → Filter color: red → Relate: food, on → Choose color: yellow | brown



Graph Normalization

- Ontology construction
- Edge Pruning
- Object Augmentation
- Global Properties

Question Generation

- Patterns Collection
- Compositional References
- Decoys Selection
- Probabilistic Generation

Sampling and Balancing

- Distribution Balancing
- Type-Based Sampling
- Deduplication

Entailments Relations

- Functional Programs
- Entailment Relations
- Recursive Reachability

New Metrics

- Consistency
- Validity & Plausibility
- Distribution
- Grounding



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

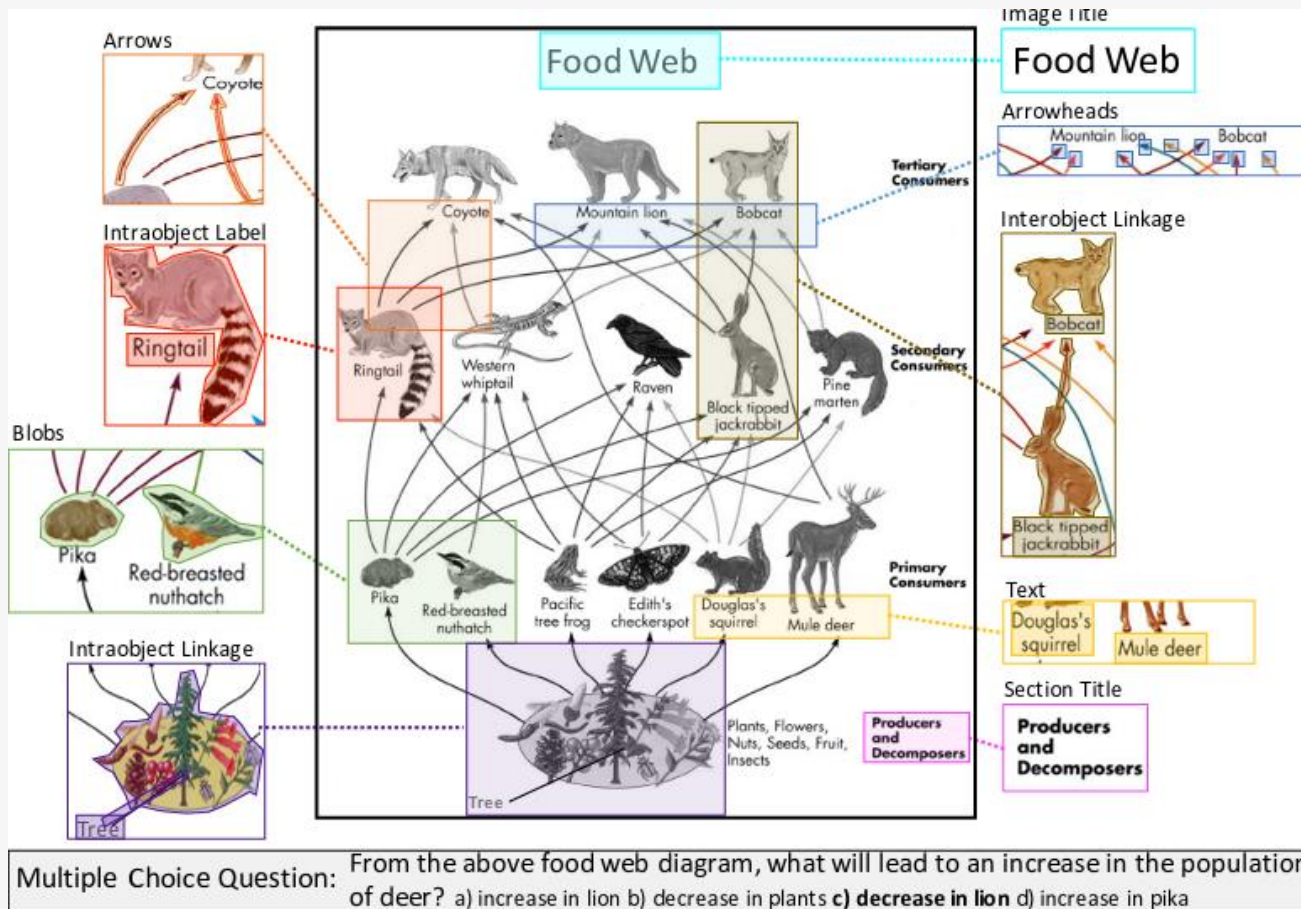
B3. What **animal** is in the **box**? **bear**

B4. Is there a **bag** to the right of the green **door**? no

B5. Is there a **box** inside the plastic **bag**? no

Static Benchmarks: GQA

- Questions are generated using a graph of images
- 22.6M questions for 113K images
- Evaluation metrics: accuracy and 5 more


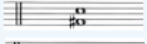
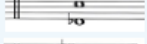
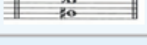

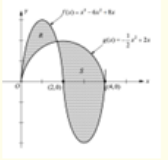

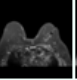
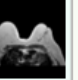

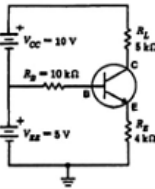
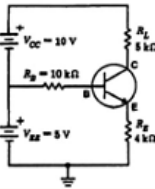


AI2D

- 15K multiple choice questions for 5K school level diagrams
- Parsing graph is available
- Evaluation metric: accuracy

MMMU

- 11.5K questions from 6 university disciplines
- Answers were extracted using regexps
- Evaluation metric: accuracy

Art & Design	Business	Science
<p>Question: Among the following harmonic intervals, which one is constructed incorrectly?</p> <p>Options:</p> <p>(A) Major third </p> <p>(B) Diminished fifth </p> <p>(C) Minor seventh </p> <p>(D) Diminished sixth </p> <p>Subject: Music; Subfield: Music; Image Type: Sheet Music; Difficulty: Medium</p>	<p>Question: ...The graph shown is compiled from data collected by Gallup . Find the probability that the selected Emotional Health Index Score is between 80.5 and 82?</p> <p>Options:</p> <p>(A) 0 (B) 0.2142 (C) 0.3571 (D) 0.5</p> <p>Subject: Marketing; Subfield: Market Research; Image Type: Plots and Charts; Difficulty: Medium</p>	<p>Question:  The region bounded by the graph as shown above. Choose an integral expression that can be used to find the area of R.</p> <p>Options:</p> <p>(A) $\int_0^{1.5} [f(x) - g(x)] dx$ (B) $\int_0^{1.5} [g(x) - f(x)] dx$ (C) $\int_0^2 [f(x) - g(x)] dx$ (D) $\int_0^2 [g(x) - x(x)] dx$</p> <p>Subject: Math; Subfield: Calculus; Image Type: Mathematical Notations; Difficulty: Easy</p>
Health & Medicine	Humanities & Social Science	Tech & Engineering
<p>Question: You are shown subtraction , T2 weighted  and T1 weighted axial  from a screening breast MRI. What is the etiology of the finding in the left breast?</p> <p>Options:</p> <p>(A) Susceptibility artifact (B) Hematoma (C) Fat necrosis (D) Silicone granuloma</p> <p>Subject: Clinical Medicine; Subfield: Clinical Radiology; Image Type: Body Scans: MRI, CT.; Difficulty: Hard</p>	<p>Question: In the political cartoon, the United States is seen as fulfilling which of the following roles? </p> <p>Option:</p> <p>(A) Oppressor (B) Imperialist (C) Savior (D) Isolationist</p> <p>Subject: History; Subfield: Modern History; Image Type: Comics and Cartoons; Difficulty: Easy</p>	<p>Question: Find the VCE for the circuit shown in   Neglect VBE</p> <p>Answer: 3.75</p> <p>Explanation: ...IE = [(VEE) / (RE)] = [(5 V) / (4 k-ohm)] = 1.25 mA; VCE = VCC - IERL = 10 V - (1.25 mA) 5 k-ohm; VCE = 10 V - 6.25 V = 3.75 V</p> <p>Subject: Electronics; Subfield: Analog electronics; Image Type: Diagrams; Difficulty: Hard</p>

MMMU: A Massive Multi-discipline Multimodal Understanding and Reasoning Benchmark for Expert AGI (2023)

<https://arxiv.org/abs/2311.16502>

TextVQA



What does it say near the star on the tail of the plane?

Ground Truth

jet

Prediction

nothing

(a)



What is the time on bottom middle phone?

Ground Truth

15:20

Prediction

12:00

(b)



What is the top oz?

Ground Truth

16

Prediction

red

(c)



What is the largest denomination on table?

Ground Truth

500

Prediction

unknown

(d)

- 45K questions for 26K images
- 10 answers per question
- Evaluation metric: CQA accuracy (100% in case 3 people provided the same answer)

Towards VQA Models That Can Read (2019)

<https://arxiv.org/abs/1904.08920>

<u>INDICATIONS FOR IMPLANTATION</u>				
	WOVEN	KNITTED	VELOURS	TOTAL
Aneurism	1 *	7 *	6 *	14

Q: What is the underlined heading just above the table?

GT: *Indications for implantation*

M4C best: *indications for implantation*

BERT best: *total aneurism*

Human: *indications for implantation*

BUSINESS EXPENSE VOUCHER	
Employee Name Charles A. Blixt	Account Number 71614
Mailing Address (If applicable) Sr. VP/GC 11803 Executive	Extension Number (910) 741-0673

Q: What is the Extension Number as per the voucher?

GT: *(910) 741-0673*

M4C best: *963.12*

BERT best: *(910) 741-0673*

Human: *(910) 741-0673*



Q: How many boxed illustrations are there ?

GT: *9*

M4C best: *4*

BERT best: *4*

Human: *9*

DocVQA

- 50K questions for 12K images
- Documents from 1960-2000
- Evaluation metrics: Average Normalized Levenshtein Similarity, Accuracy

Image



Model A

Model B



Given my horse's location on this map, what is the quickest route to reach it?

Given my horse's location on this map, what is the quickest route to reach it?

Based on the map, your current location is marked by the pin/stamp icon at Lookout Landing. The quickest route to reach your horse, which is marked with a horse icon at Rito Village, would be to head northwest along the roads/paths shown on the map. It looks like you'll need to traverse the area labeled "Tabantha Frontier" to get from your current spot to Rito Village where your horse is located. The exact path is a bit unclear from this zoomed out view, but generally heading northwest from Lookout Landing should get you to your horse's position at Rito Village by the most direct route possible given the geography shown.

I'm sorry, but I don't have access to your specific map or location information. Could you please provide more context or details so I can assist you better?

Rules

- Upload any image, and send any question to two anonymous models and vote for the better one!
- Or you could also click "Sample Input" to get a random example from public benchmarks such as VisIT-Bench.
- Single image multi-round chat is allowed, you can continue to send question until you identify a winner.

Reason for your vote

Please input your reason for response preference here before clicking the model choice button.

A is better

B is better

Tie

Both are bad

WildVision: Evaluating Vision-Language Models in the Wild with Human Preferences (2024)
<https://arxiv.org/abs/2406.11069>

High-Level Classification

1. Deep Fusion

```
graph TD; A([1. Deep Fusion]) --> B[1.1. Standard]; A --> C[1.2. Custom];
```

1.1. Standard

Cross-Attention
- Flamingo

1.2. Custom

Layers
- LLaVA
- miniGPT-4
- Gemini

2. Early Fusion

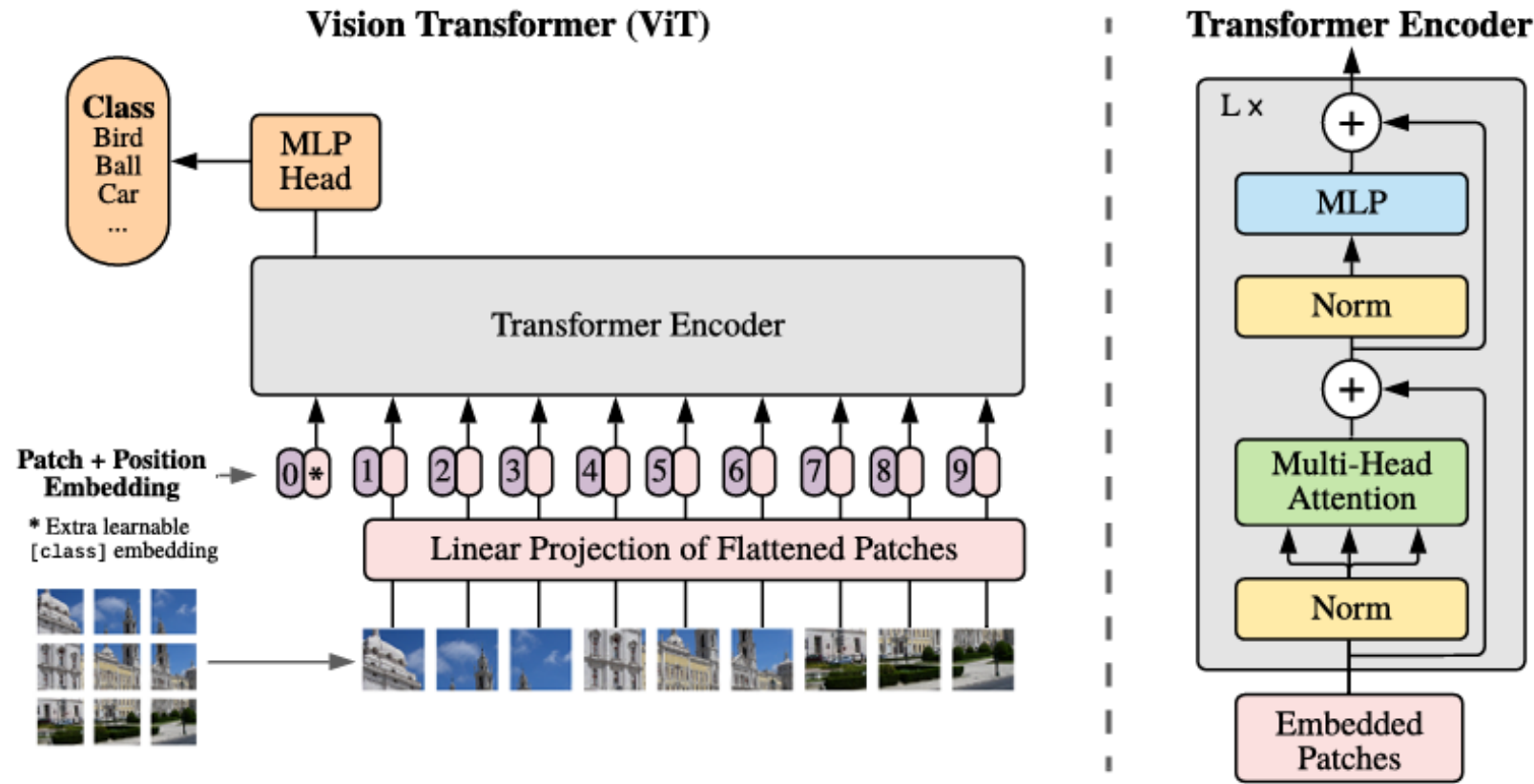
```
graph TD; A([2. Early Fusion]) --> B[2.1. Non-tokenized]; A --> C[2.2. Tokenized];
```

2.1. Non-tokenized

- BLIP

2.2. Tokenized

- Kosmos-1
- Fuyu



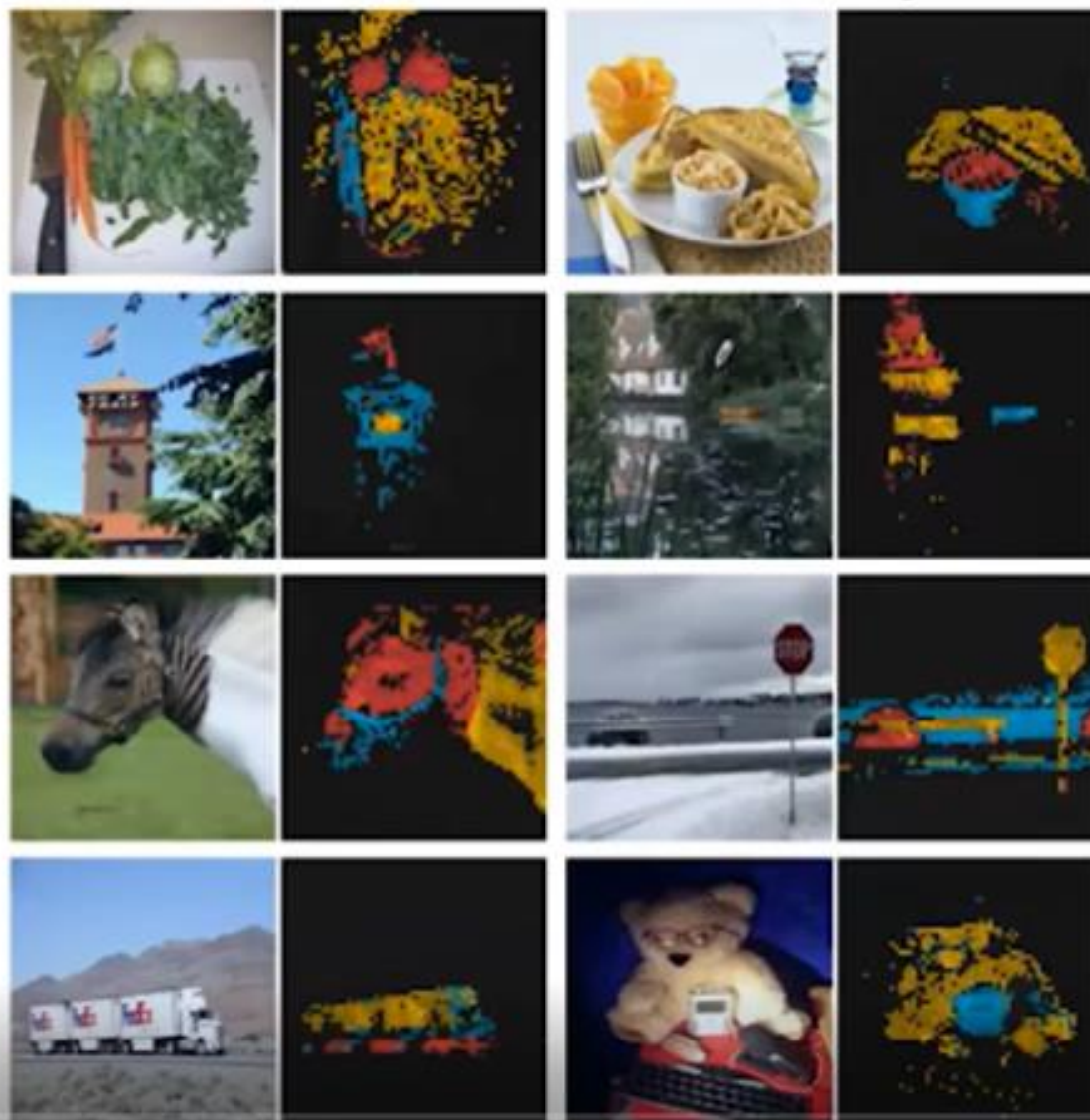
An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale (2021) <https://arxiv.org/abs/2010.11929>

Input

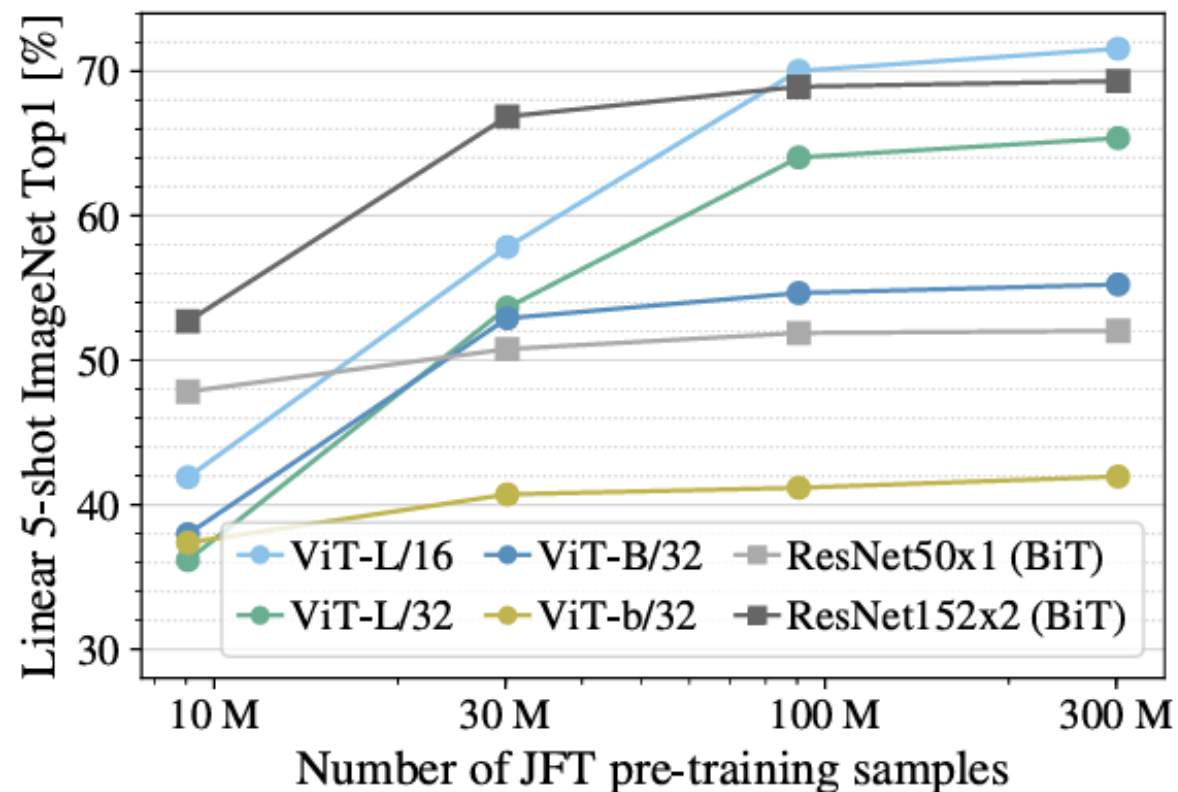
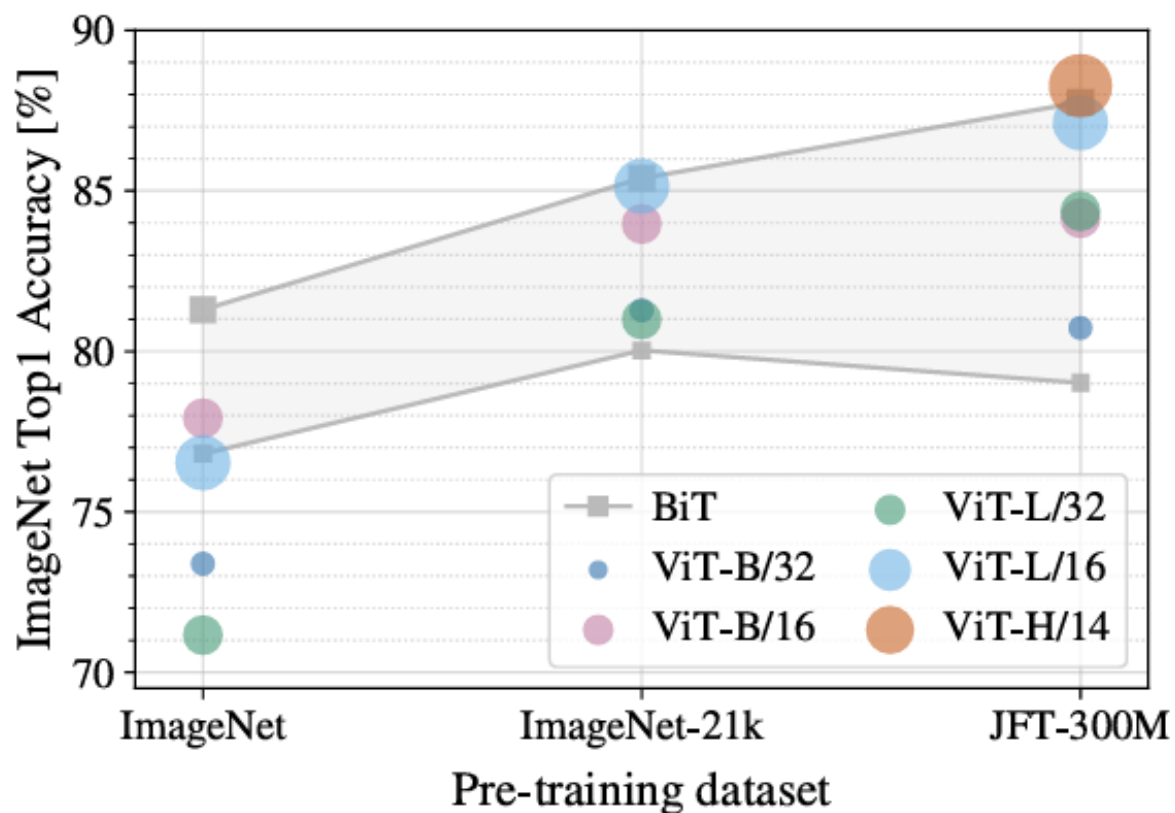
Attention



MSA visualization, 8×8 patches

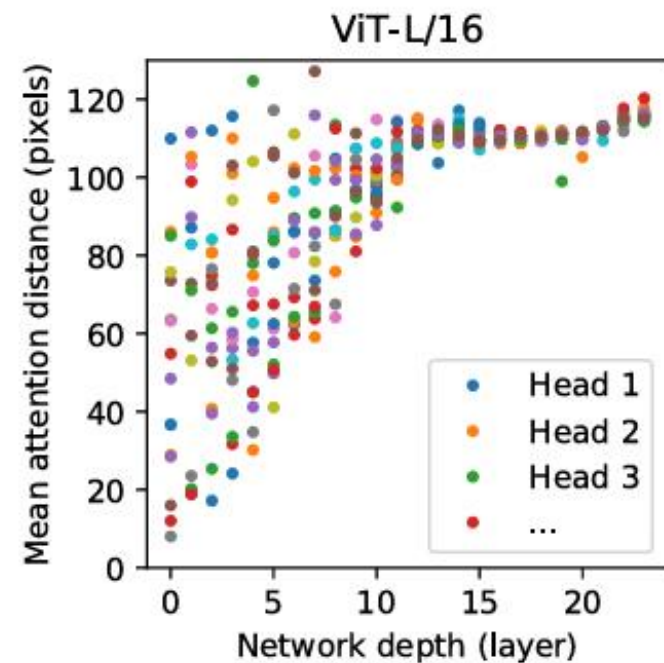
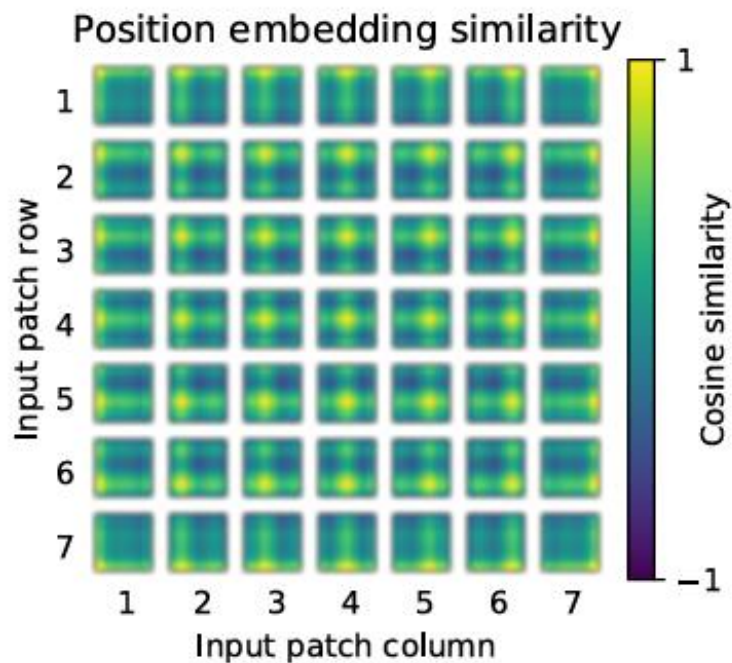
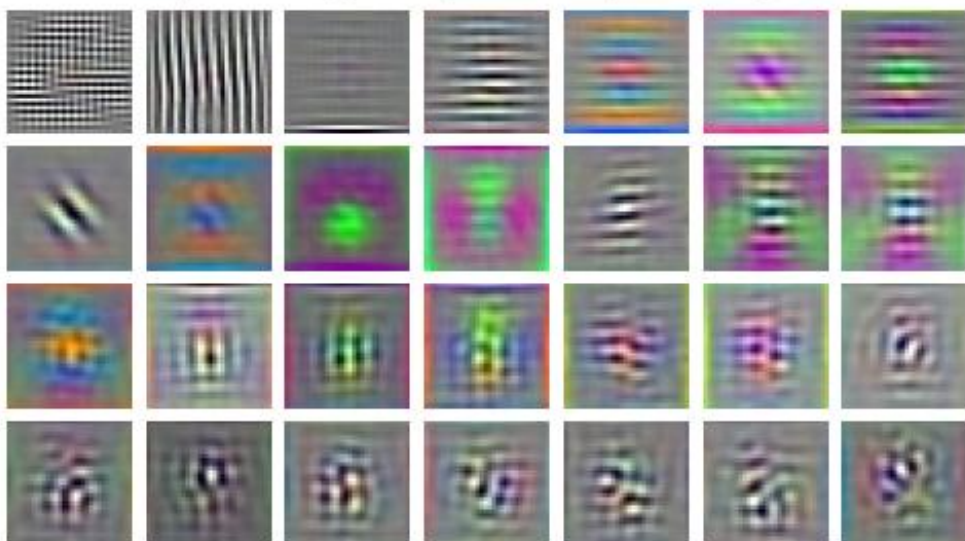


Vision Transformer

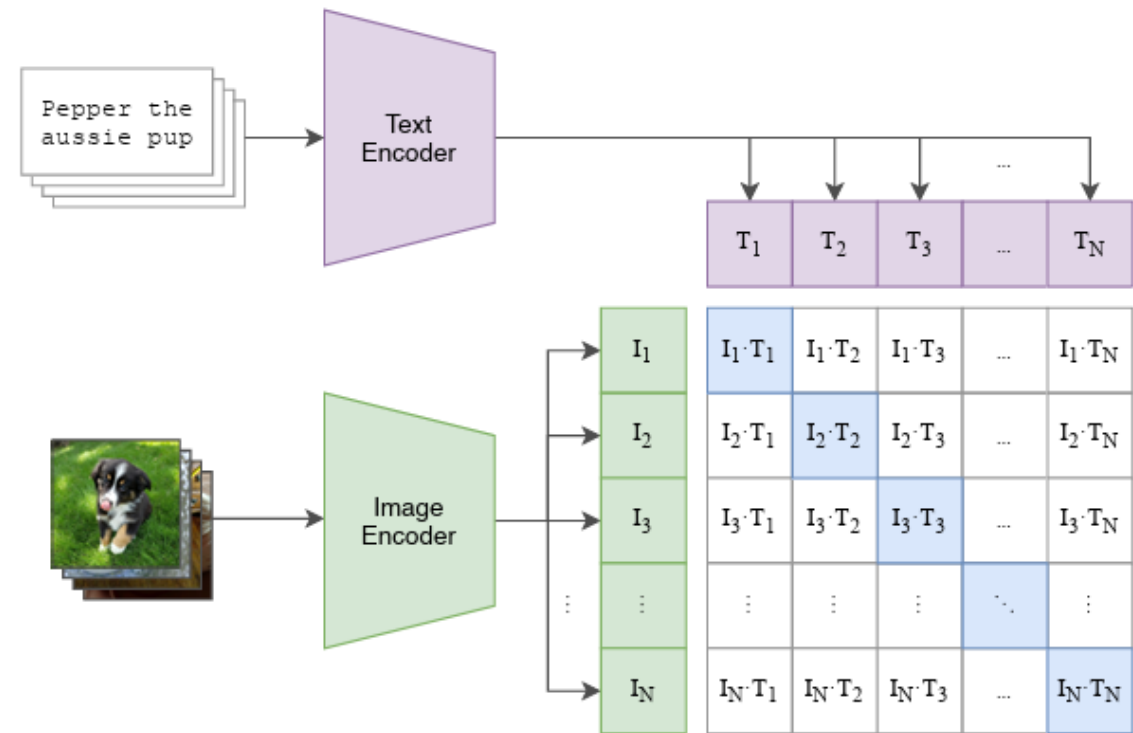
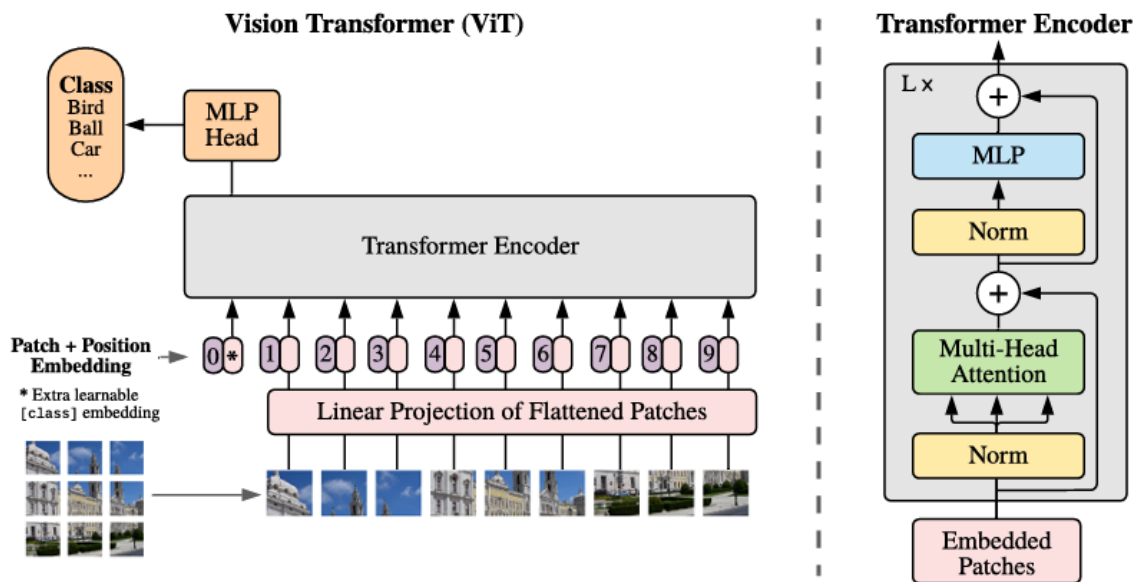


Vision Transformer

RGB embedding filters
(first 28 principal components)



CLIP

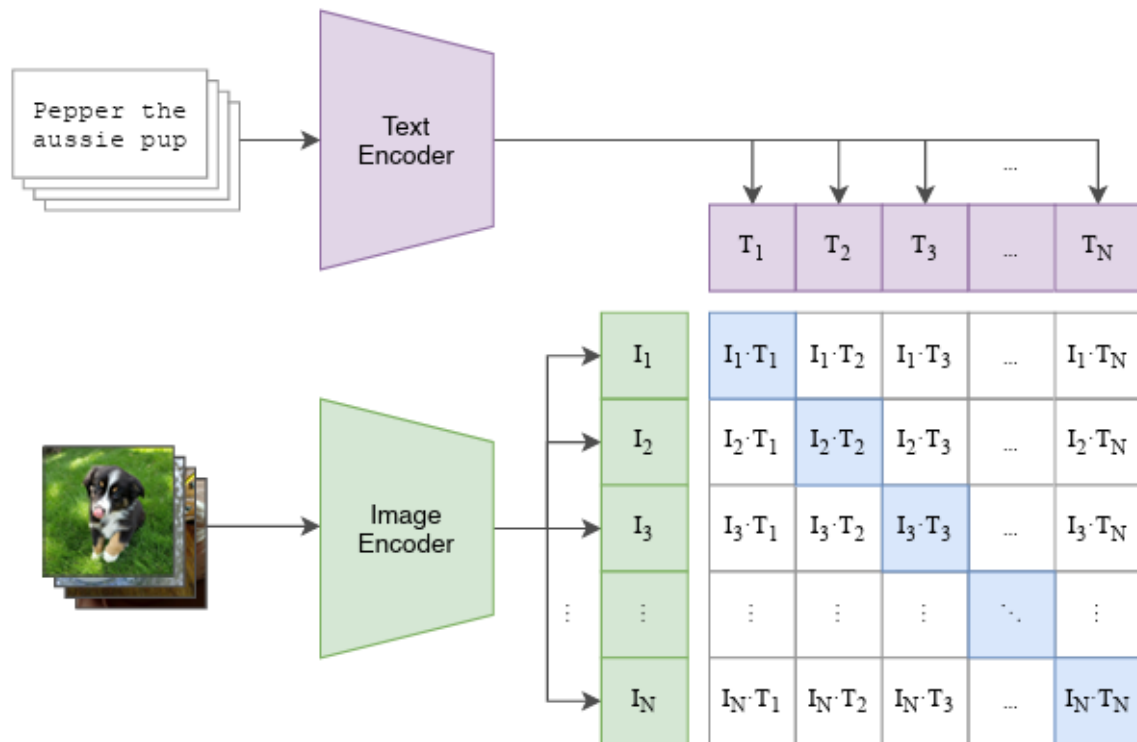


$$-\frac{1}{2|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \left(\overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_i \cdot \mathbf{y}_j}}}^{\text{image} \rightarrow \text{text softmax}} + \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_j \cdot \mathbf{y}_i}}}^{\text{text} \rightarrow \text{image softmax}} \right)$$

400M (image, text), 500 V100 GPU,

Learning Transferable Visual Models From Natural Language Supervision (2021) <https://arxiv.org/abs/2103.00020>

SigLIP



$$-\frac{1}{2|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \left(\overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_i \cdot \mathbf{y}_j}}}^{\text{image} \rightarrow \text{text softmax}} + \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_j \cdot \mathbf{y}_i}}}^{\text{text} \rightarrow \text{image softmax}} \right)$$

$$-\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \underbrace{\log \frac{1}{1 + e^{z_{ij}(-t\mathbf{x}_i \cdot \mathbf{y}_j + b)}}}_{\mathcal{L}_{ij}}$$