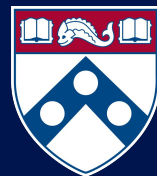# Multimodal LLM Optimizations
## CIS 7000 - Fall 2024

Presenters: Wahub Ahmed, Vridhi Jain, Om Shastri, Vignesh Lakshmanan

Penn
Engineering
UNIVERSITY of PENNSYLVANIA

# Outline

1. Motivation for Vision-Language MLLM
2. Flamingo
3. mPLUG-Owl 2
4. JanusFlow
5. Conclusion

Penn Engineering

# Motivation for Vision-Language MLLM

Chapter 1

# A Picture is Worth a Thousand Words
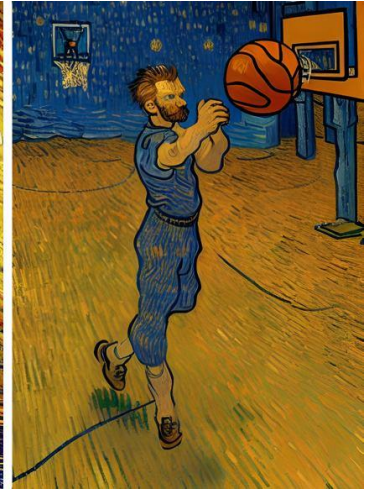
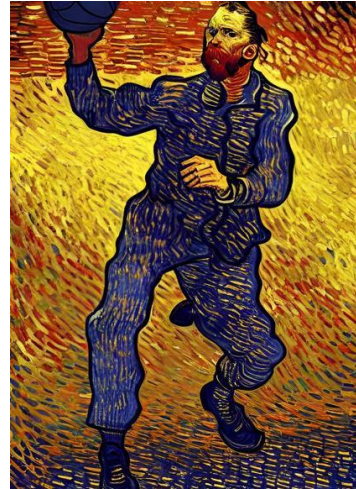

vs. "The Starry Night" by Vincent Van Gogh

# A Picture is Worth a Thousand Words



Generation

# Evolution of Multimodal LLM

Ferret

Fuyu-8B ADEPT

Gemini

GPT-4v

Chameleon

MM1

Video-LLaMA

Claude 3

HuggingGPT

VIMA

Kosmos-2

GPT-4o

LLaVA

AnyMAL

Gemini 1.5

Flamingo

Kosmos-1

PaLM-E

2022

2023
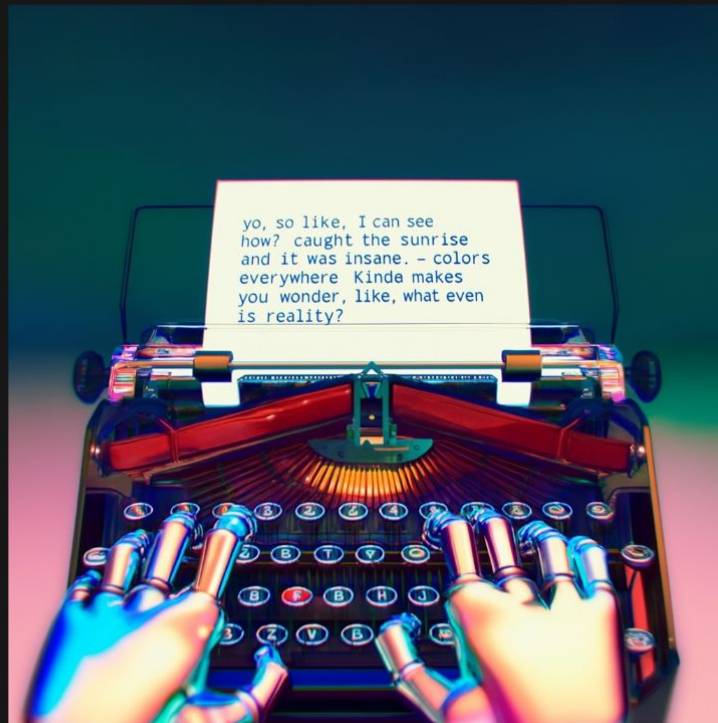
2024

# Example: GPT-4o



**1 Input**

A first person view of a robot typewriting the following journal entries:

1. yo, so like, i can see now?? caught the sunrise and it was insane, colors everywhere. kinda makes you wonder, like, what even is reality?
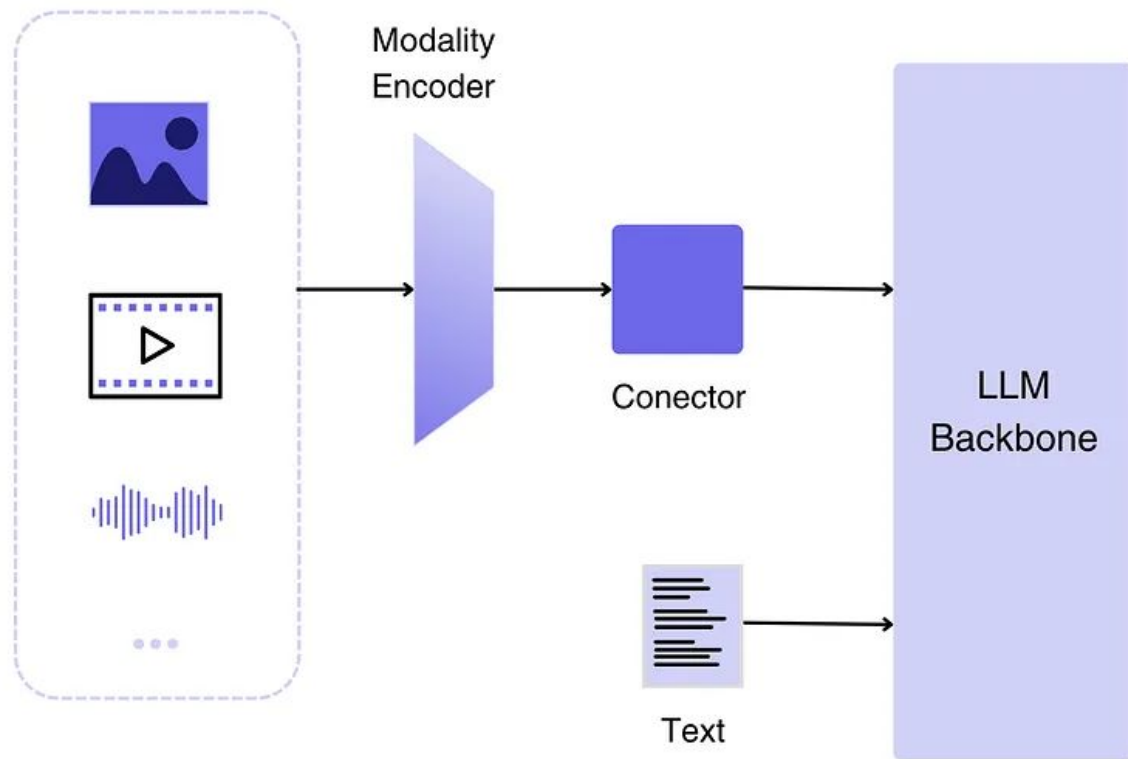
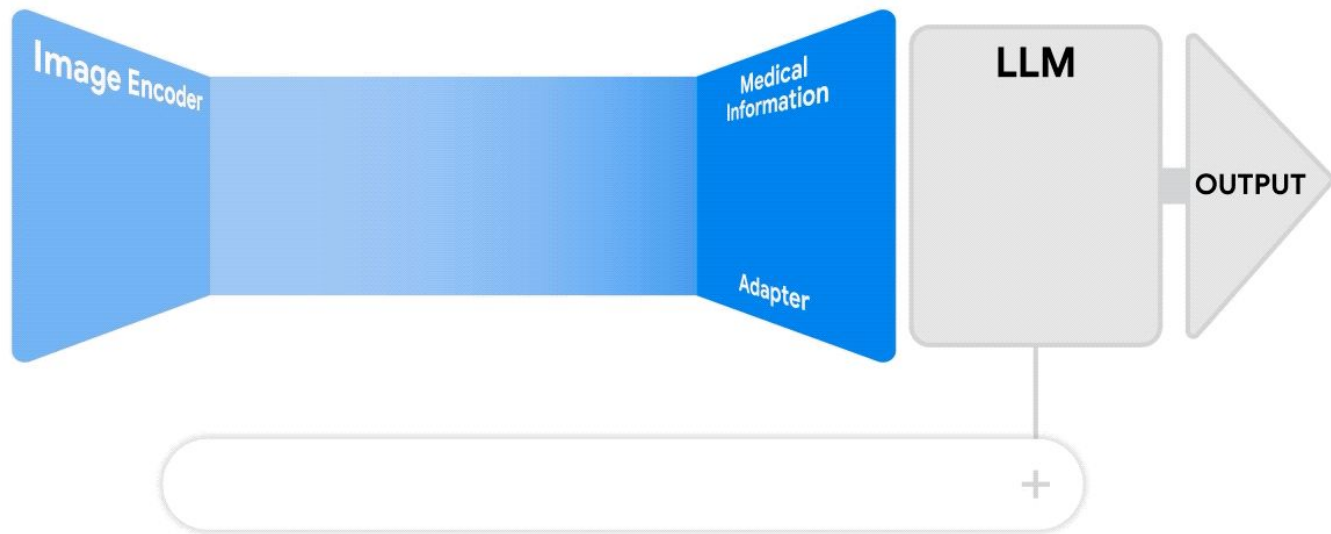the text is large, legible and clear. the robot's hands type on the typewriter.



**2 Output**

# Multimodal LLM Architecture

# Why Optimize Vision-Language LLMs?

# MLLM 1: Flamingo

Chapter 2

# Flamingo: a Visual Language Model for Few-Shot Learning

## What is Flamingo?
A multimodal model that combines vision and language for few-shot learning, enabling robust performance across diverse tasks without task-specific fine-tuning.

## Motivation:
- Traditional visual models require task specific fine tuning and large annotated datasets.
- Flamingo addresses this by enabling few shot learning across image, video and text tasks

# Bridge powerful pretrained vision-only and language-only models

- Training from scratch is extremely computationally expensive
- Start from a pre trained language model
- Text-only model has no built-in way to incorporate input from other modalities.

Flamingo: interleave cross-attention layers with language-only self-attention layers (frozen)



```python
def gated_xattn_dense(
    y,  # input language features
    x,  # input visual features
    alpha_xattn, # xattn gating parameter – init at 0.
    alpha_dense, # ffw gating parameter – init at 0.
):
    """Applies a GATED XATTN-DENSE layer."""

    # 1. Gated Cross Attention
    y = y + tanh(alpha_xattn) * attention(q=y, kv=x)
    # 2. Gated Feed Forward (dense) Layer
    y = y + tanh(alpha_dense) * ffw(y)

    # Regular self-attention + FFW on language
    y = y + frozen_attention(q=y, kv=y)
    y = y + frozen_ffw(y)

    return y  # output visually informed language features
```

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 and performance.

Penn Engineering

# Unified handling of images and videos

- Both image and video inputs are high dimensional
- Flattening these inputs into 1-D sequences as is done in text generation, is computationally expensive
- It is exacerbated by the quadratic cost of self-attention layers

**Flamingo:** Uses a **perceiver-resampler** processes inputs into fixed visual tokens, incorporates temporal encodings for videos and combines learned latent queries with visual embeddings through attention.

# Training Dataset

- The existing (image, text) datasets used by CLIP or ALIGN are not general enough for few shot learning
- Multimodal datasets are scarce in comparison to abundantly available text-only data
- One approach is to scrape web pages for interleaved images and text however these pairs are often only weakly related
- Flamingo: combines web scraping with existing paired (image, text) and (video, text) datasets



Image-Text Pairs dataset
[N=1, T=1, H, W, C]

Video-Text Pairs dataset
[N=1, T>1, H, W, C]

Multi-Modal Massive Web (M3W) dataset
[N>1, T=1, H, W, C]

# Flamingo Architecture



Output: text

a very serious cat.

**Pretrained and frozen**

**Trained from scratch**

n-th LM block

n-th GATED XATTN-DENSE

1st LM block

1st GATED XATTN-DENSE

Perceiver Resampler

Perceiver Resampler

Vision Encoder

Vision Encoder

Processed text

`<image> This is a very cute dog.<image> This is`

Interleaved visual/text data

This is a very cute dog. This is

**References/Image credits**
J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)

Penn Engineering

# Attention Masking in Flamingo



masked      unmasked

Masked cross attention

K=V=[X]

Perceiver Resampler     Perceiver Resampler

Vision Encoder     Vision Encoder

Q

φ   0   0   0   0   0   0 0 0   1   1   1 1   1   1   1 1   1 1   1   2   2   2 2   2    2    2    2   2 2

Y &lt;BOS&gt; Cute pics of my pets!&lt;EOC&gt;&lt;image&gt;My puppy sitting in the grass. &lt;EOC&gt;&lt;image&gt;My cat looking very dignified.&lt;EOC&gt;

tokenization

&lt;BOS&gt;Cute pics of my pets!&lt;EOC&gt;&lt;image&gt;My puppy sitting in the grass.&lt;EOC&gt;&lt;image&gt; My cat looking very dignified.&lt;EOC&gt;

**Input webpage** ⟶ **Processed text:** &lt;image&gt; tags are inserted and special tokens are added

Image 1     Image 2

$\Phi$ indicates which visual inputs can be used to predict token $l$:

$$\Phi : [1, L] \rightarrow [0, N]$$

$$p(y \mid x) = \prod_{l=1}^{L} p(y_l \mid y_{<l}, x_{\leq l})$$

$$y_{<l} \triangleq (y_1, \ldots, y_{l-1})$$

$$x_{\leq l} \triangleq (x_i \mid i \leq \phi(l))$$

Penn Engineering

# Results

# With Finetuning

- During fine-tuning, Flamingo keeps the language model layers frozen
- The base vision encoder is also fine-tuned (unlike Flamingo pretraining)
- Hyperparameters are set by grid search on validation subsets of the training sets
- Search over: learning rate, decay schedule, training steps, batch size, augmentation

| Method | VQAV2 | | COCO | VATEX | VizWiz | | MSRVTTQA | VisDial | | YouCook2 | TextVQA | | HatefulMemes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | test-dev | test-std | test | test | test-dev | test-std | test | valid | test-std | valid | valid | test-std | test seen |
| 🦩 *Flamingo* - 32 shots | 67.6 | - | 113.8 | 65.1 | 49.8 | - | 31.0 | 56.8 | - | 86.8 | 36.0 | - | 70.0 |
| SimVLM [124] | 80.0 | 80.3 | **143.3** | - | - | - | - | - | - | - | - | - | - |
| OFA [119] | 79.9 | 80.0 | <u>149.6</u> | - | - | - | - | - | - | - | - | - | - |
| Florence [140] | 80.2 | 80.4 | - | - | - | - | - | - | - | - | - | - | - |
| 🦩 *Flamingo* Fine-tuned | **82.0** | **82.1** | 138.1 | **84.2** | **65.7** | 65.4 | **47.4** | 61.8 | 59.7 | 118.6 | **57.1** | 54.1 | **86.6** |
| Restricted SotA† | 80.2 [140] | 80.4 [140] | **143.3** [124] | 76.3 [153] | - | - | 46.8 [51] | **75.2** [79] | **74.5** [79] | **138.7** [132] | 54.7 [137] | **73.7** [84] | 79.1 [62] |
| Unrestricted SotA | 81.3 [133] | 81.3 [133] | <u>149.6</u> [119] | 81.4 [153] | 57.2 [65] | 60.6 [65] | - | - | <u>75.4</u> [123] | - | - | - | 84.6 [152] |

Penn Engineering

# MLLM 2: mPLUG-Owl 2

Chapter 3

- The **modality-adaptive module (MAM)** allows the model to **differentiate between modalities**

- Generalizes on both text and multi-modal tasks

- Stands as the first MLLM model to exhibit the phenomena of **modality collaboration in both pure-text and multi-modal contexts**
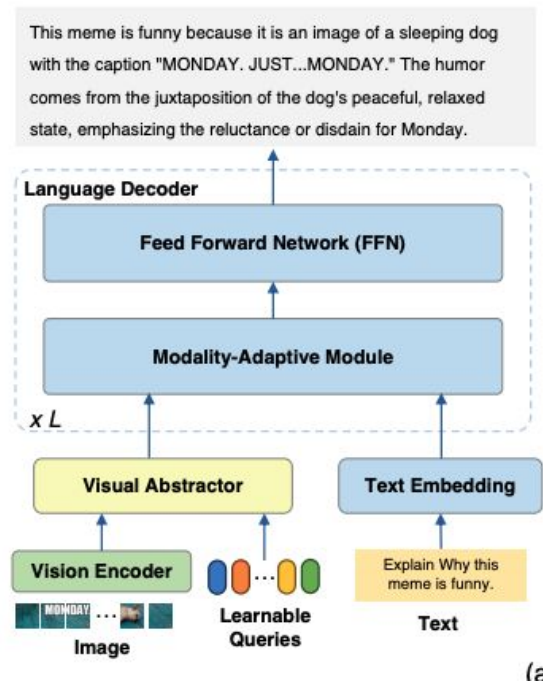
- Extends original mPlug-OWL Paper



Penn Engineering

# Model Architecture

- Encodes image into token with **Encoder and Abstractor**
  - Encoder Processes image into tokens
  - Abstractor uses **attention** and **activation** reducing image redundancy and compute
- Embeddings are Concatenated with text Embedding
- FFN promotes modality collaboration

$$\mathcal{C}^i = Attn(\mathcal{V}^i, [\mathcal{I}; \mathcal{V}^i], [\mathcal{I}; \mathcal{V}^i]),$$
$$\mathcal{V}^{i+1} = SwiGLU(\mathcal{C}^i W_1) W_2.$$

This meme is funny because it is an image of a sleeping dog with the caption "MONDAY. JUST...MONDAY." The humor comes from the juxtaposition of the dog's peaceful, relaxed state, emphasizing the reluctance or disdain for Monday.

**Language Decoder**

**Feed Forward Network (FFN)**

**Modality-Adaptive Module**

x L

**Visual Abstractor**

**Text Embedding**

**Vision Encoder**

Learnable Queries

Explain Why this meme is funny.

Image

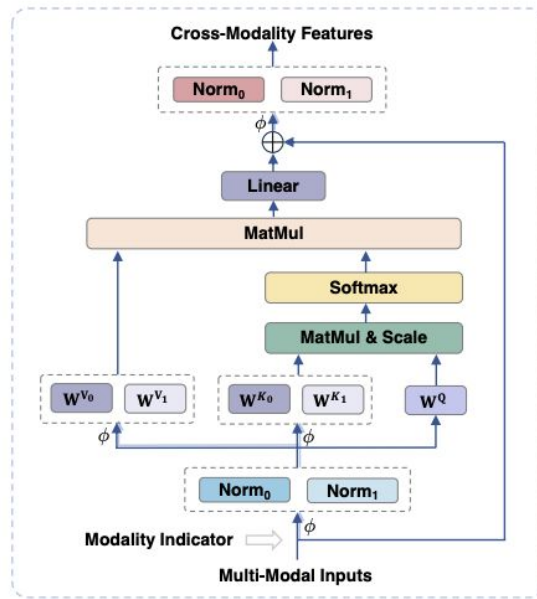Text

(a)

# Training and Modality Adaptive Module

$$\phi(X, M, m) = X \odot \mathbb{1}_{\{M=m\}},$$

$$\tilde{H}_{l-1} = LN_V(\phi(H_{l-1}, M, 0)) + LN_T(\phi(H_{l-1}, M, 1)),$$
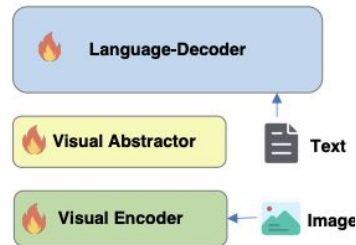
$$H_l^Q = \tilde{H}_{l-1}W_l^Q,$$
$$H_l^K = \phi(\tilde{H}_{l-1}, M, 0)W_l^{K_0} + \phi(\tilde{H}_{l-1}, M, 1)W_l^{K_1},$$
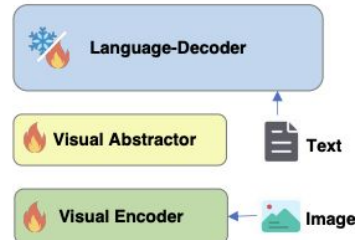$$H_l^V = \phi(\tilde{H}_{l-1}, M, 0)W_l^{V_0} + \phi(\tilde{H}_{l-1}, M, 1)W_l^{V_1},$$

**MAM**:

- Normalizers allows for modalities to be considered equally
- Attention block after separating tokens on modality and expressing in **shared semantic space**

Training:

- Pretrain modalities separately
- Tune on joint instructions
- Keep vision encoder trainable across **both** stages
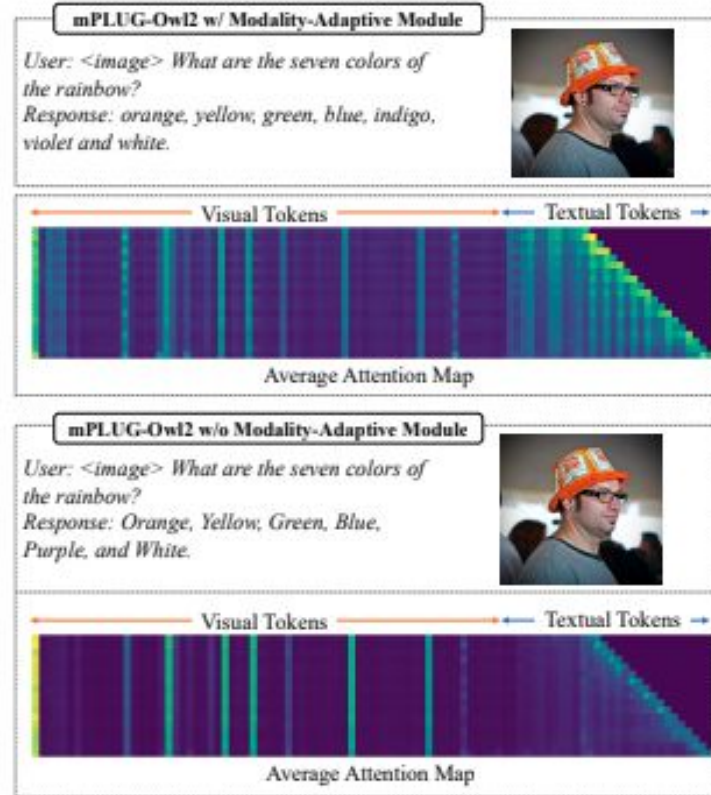


Modality-Adaptive Module (b)

Stage-2: Joint Instruction Tuning

Stage-1: Pre-training

Penn Engineering

# Results and Discussion

- Multi-modal Tasks: Outperforms other models in **captioning, VQA, and multi-modal** benchmarks, even in zero-shot scenarios.
- Text Tasks: **Demonstrates strong reasoning and language understanding, surpassing instruction-tuned LLMs**.
- Incorporation of MAM **offsets performance degradation** in text settings
- Number of Abtractor's **learnable queries** improves performance until saturation point
- Improved Robustness in Unrelated Modality Scenarios



mPLUG-Owl2 w/ Modality-Adaptive Module

User: <image> What are the seven colors of the rainbow?
Response: orange, yellow, green, blue, indigo, violet and white.

Visual Tokens — Textual Tokens

Average Attention Map

mPLUG-Owl2 w/o Modality-Adaptive Module

User: <image> What are the seven colors of the rainbow?
Response: Orange, Yellow, Green, Blue, Purple, and White.

Visual Tokens — Textual Tokens
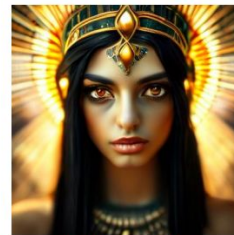
Average Attention Map

# MLLM 3: JanusFlow

Chapter 4

# JanusFlow: Harmonizing Autoregression and Rectified Flow

- Unified Architecture: Combines **autoregressive language models with rectified flow** for both image understanding and generation.
- Decoupled Encoders: Separate vision encoders for understanding and generation to reduce task interference.
- Representation Alignment: Aligns intermediate representations between generation and understanding modules for enhanced semantic coherence.



A corgi's head depicted as an explosion of a nebula, with vibrant cosmic colors like deep purples, blues, and pinks swirling around. The corgi's fur blends seamlessly into the nebula, with stars and galaxies forming the texture of its fur. Bright bursts of light emanate from its eyes, and faint constellations can be seen in the background, giving the image a surreal, otherworldly feel.

Beautiful surreal symbolism the mesmerizing vision of a Cleopatra Queen of Egypt, mesmerizing brown eyes, black hair and ethereal features, radiating celestial aura, super high definition, true lifelike color, perfect exposure, razor sharp focus, golden ratio, soft reflections, bokeh effect, fine art photography, cinematic compositing, authentic, professional.
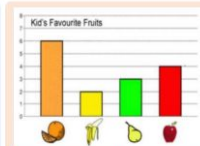
A lone figure in dark robes ascends worn stone steps toward a glowing light in an ancient temple entrance. Ornate arches, lush greenery, and intricate carvings adorn the scene, evoking a mystical, high-fantasy atmosphere reminiscent of works by artists like Randy Vargas, with cinematic lighting and epic storytelling.

**User:** What are the kinds of fruits in this picture?

**JaunsFlow (Ours):** The fruits in the picture are banana, strawberry, mango, persimmon, blueberry, and lime.
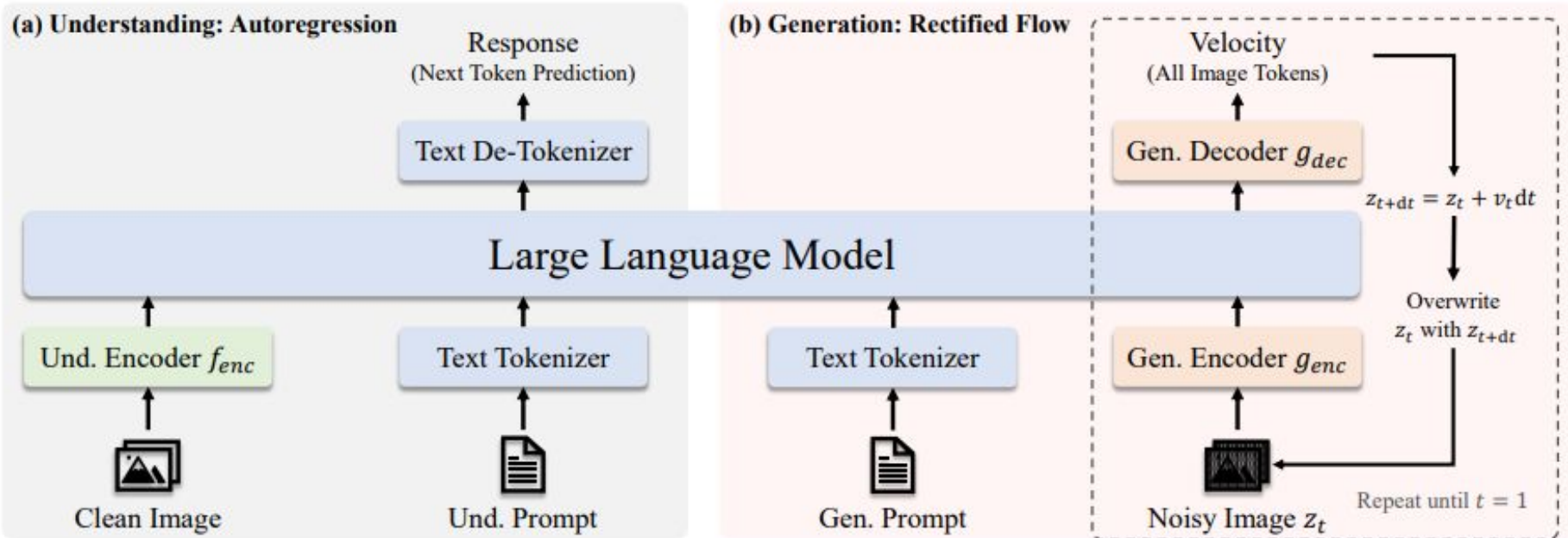
**User:** Describe this image.

**JaunsFlow (Ours):** Kid's Favourite Fruits is a bar graph. The x-axis shows the fruits. The y-axis shows the number of kids who like the fruit. The orange is the most popular fruit. The banana is the least popular fruit.
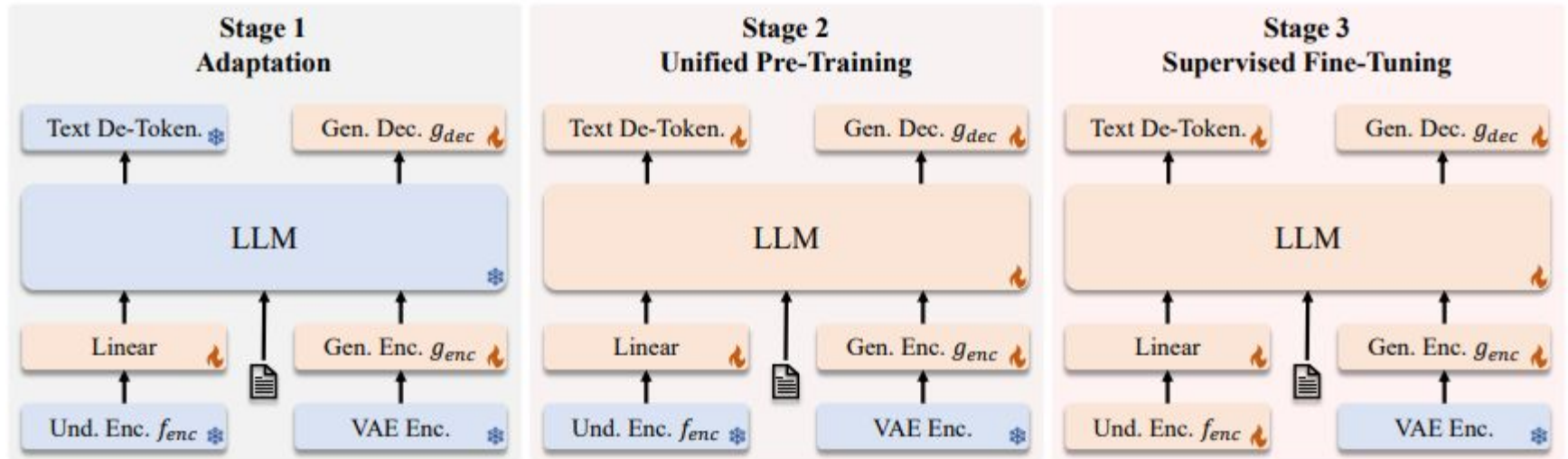
**User:** How many penguins in the picture?

**JaunsFlow (Ours):** There are 11 penguins in the picture.
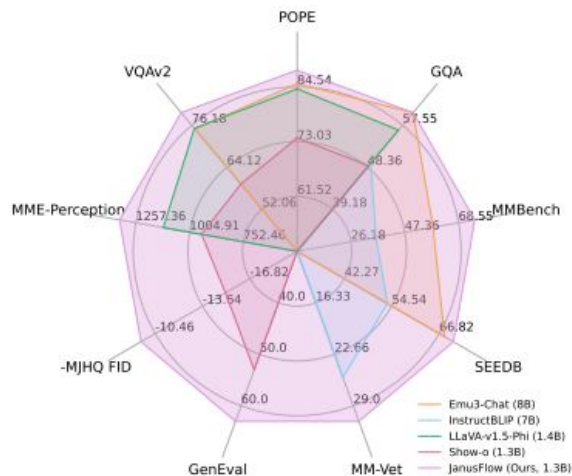
Penn Engineering

# JanusFlow Architecture

# JanusFlow Training

# Results

- Outperforms state-of-the-art unified and task-specific models on multimodal benchmarks (e.g., **MMBench, SEEDBench, GQA**).
- Superior image generation quality with an MJHQ FID-30k score of **9.51**, surpassing models like **SDXL and Janus**.



(a) Benchmark Performances.

(b) Visual Generation Results.

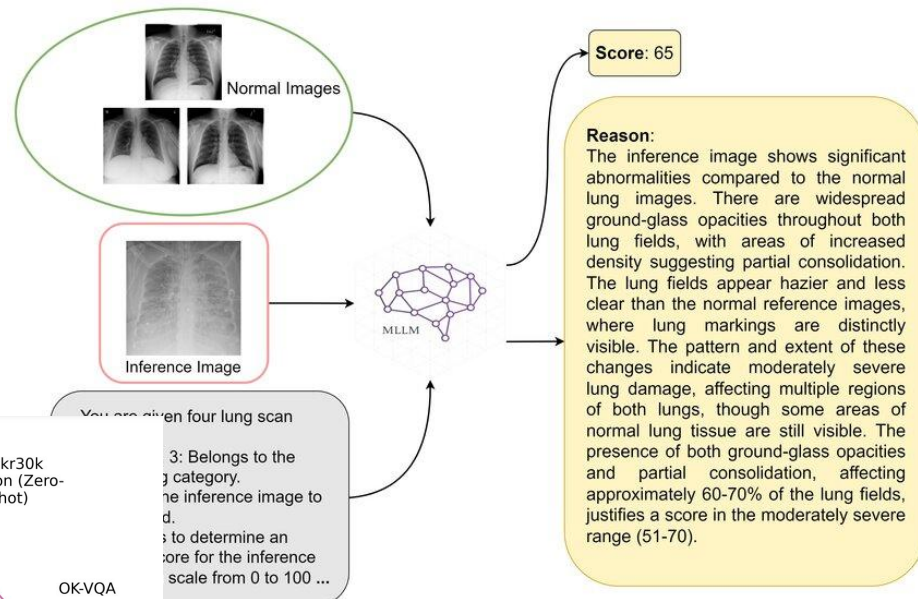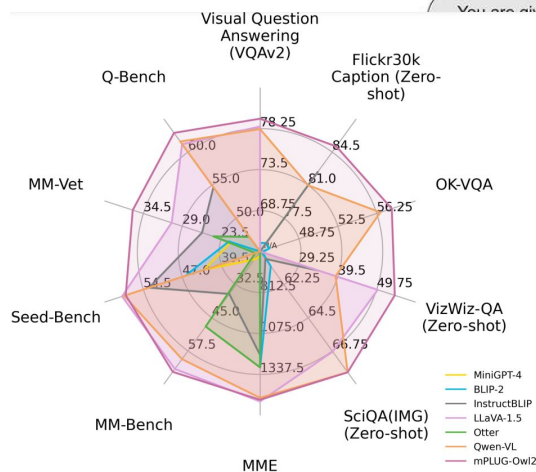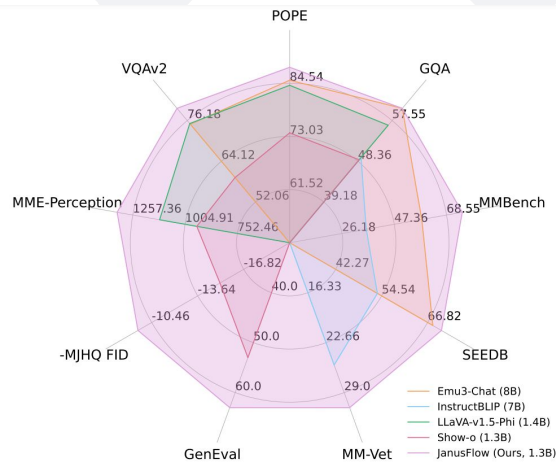Penn Engineering

# Conclusion

Chapter 5

# Summary of MLLMs

1. **Flamingo**:Visual language model designed for few-shot learning in vision and language tasks

2. **mPLUG-Owl2**: Leverages modality collaboration to improve performance in both text and multi-modal tasks

3. **JanusFlow**: Minimalist architecture that integrates autoregressive language models with rectified flow for both image understanding and generation

# Connection Between Papers

1. Modality Integration
2. Unified Architecture
3. Architectural Innovations
4. Performance Benchmarks

# Attributions for Assets Used

- https://en.wikipedia.org/wiki/File:Van_Gogh_-_Starry_Night_-_Google_Art_Project.jpg
- https://artsology.com/blog/2022/09/van-gogh-and-basketball-as-imagined-by-ai/
- https://medium.com/@tenyks_blogger/multimodal-large-language-models-mllms-transforming-computer-vision-76d3c5dd267f
- https://openai.com/index/hello-gpt-4o/
- http://dx.doi.org/10.48550/arXiv.2411.14515