# CV / VLMs

Unit 4: Vision-Language Models (VLMs)

# 4.1.1

# Introduction to Multimodal Models
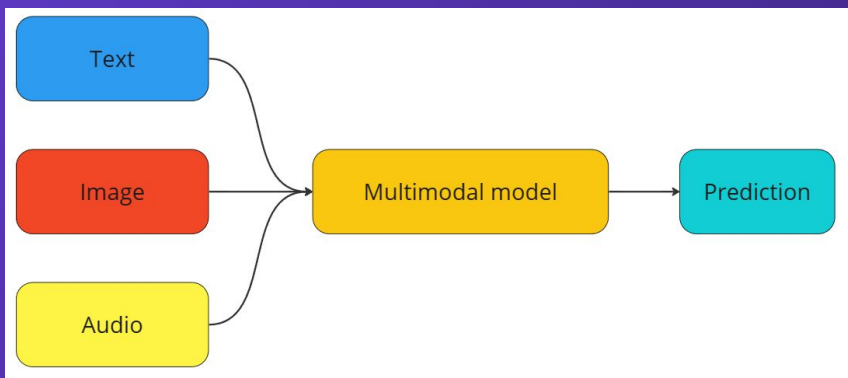
What are multimodal models?

TIL-AI
TODAY I LEARNED AI

# Multimodal models
## Definition

Multimodal models are models that integrates information from multiple modalities (e.g. text, image, videos, audio, gestures) to create a unified representation.
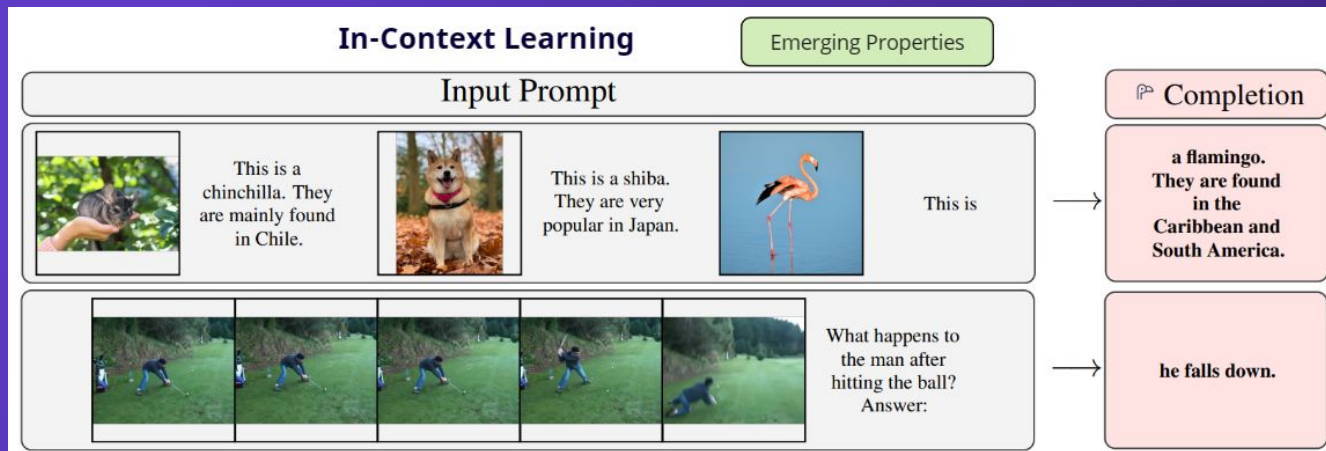
By leveraging different types of data, multimodal models can perform tasks that involve complex interactions between modalities and thus understand and reason about the world in more comprehensive way.

# Multimodal models and their importance

The key importance of multimodal models in computer vision are:

1.  Richer Understanding

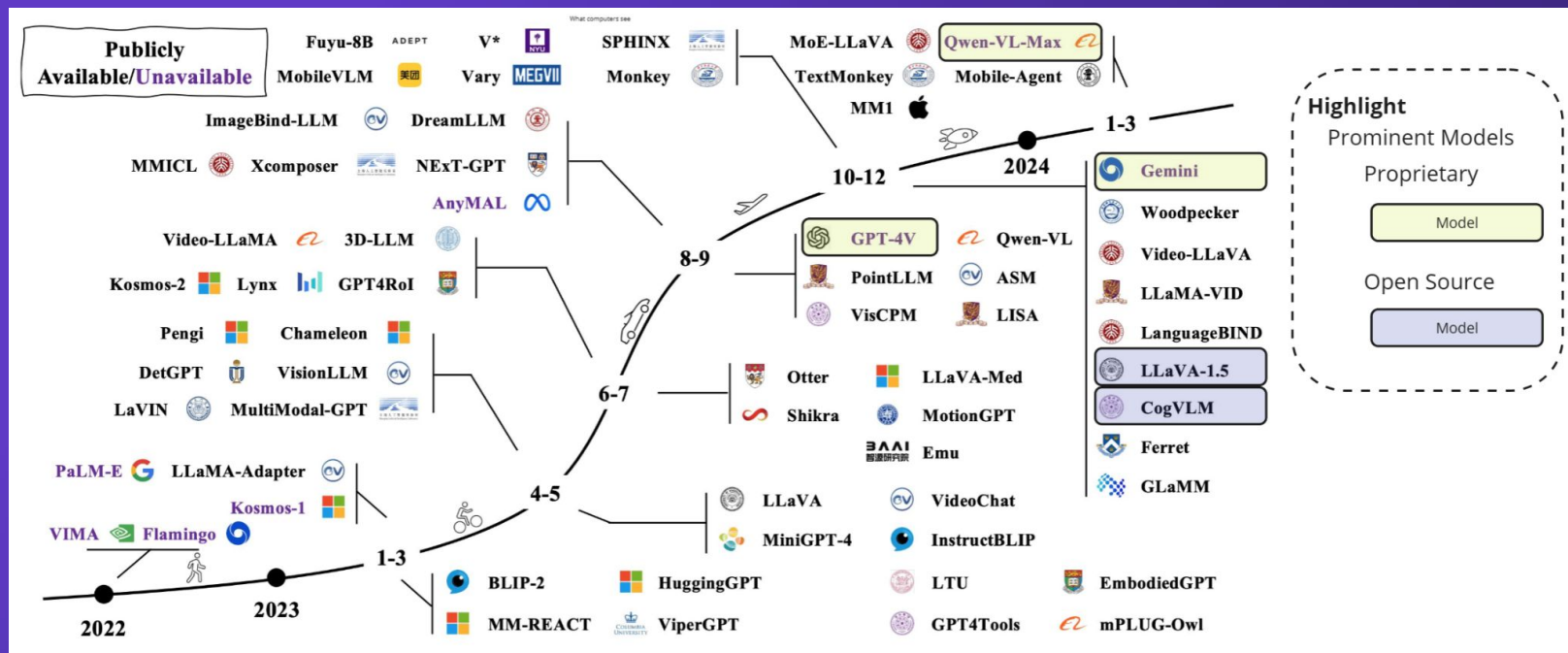2.  Better Performance

3.  Facilitate Emerging Properties



(left) The emerging properties of pre-training on web-scale interleaved image-text data: multimodal in-context-learning.VLM (Flamingo)

Adapted from Alayrac et al. (2022) via

# Multimodal Models
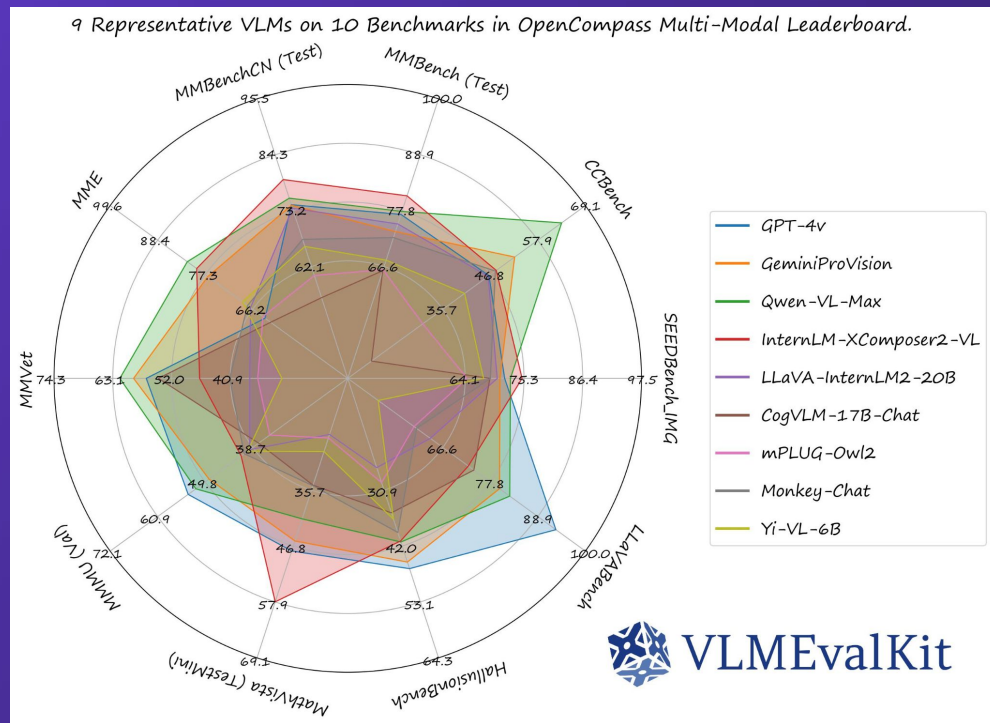## Overview - Development History

This is an area that is being actively developed. Listed below are the models released between 2022 - 2024.

# Multimodal Models
## Selection & Benchmarks

Given the wide variety of models and tasks, the utilization of benchmarks to select the most appropriate model is vital.

- There exists a host of different benchmarks. Each of them is designed to cater to a different task.
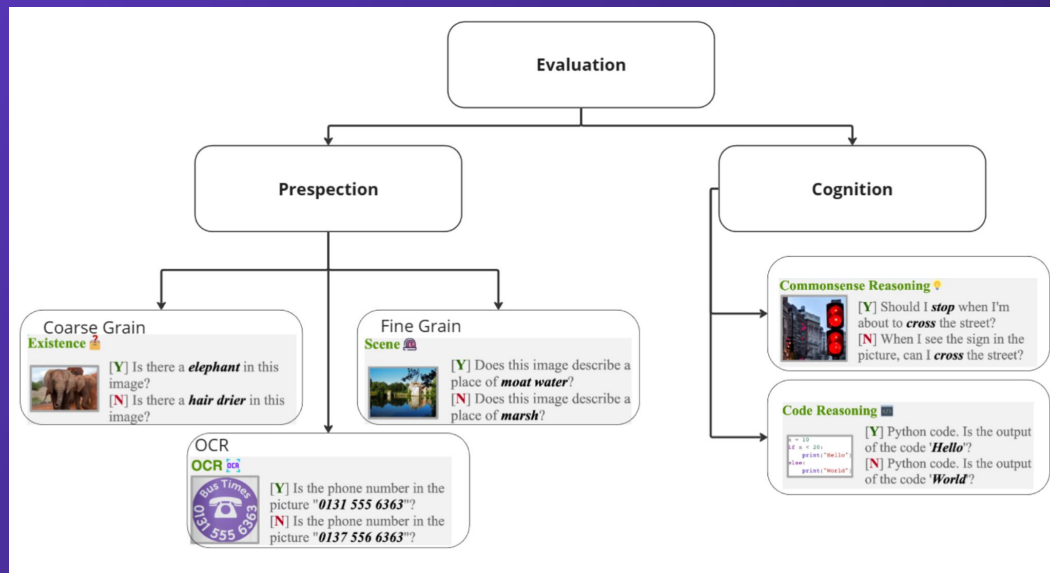


9 Representative VLMs on 10 Benchmarks in OpenCompass Multi-Modal Leaderboard.

Legend:
- GPT-4v
- GeminiProVision
- Qwen-VL-Max
- InternLM-XComposer2-VL
- LLaVA-InternLM2-20B
- CogVLM-17B-Chat
- mPLUG-Owl2
- Monkey-Chat
- Yi-VL-6B

VLMEvalKit

- [open-compass/VLMEvalKit: Open-source evaluation toolkit of large vision-language models (LVLMs), support GPT-4v, Gemini, QwenVLPlus, 30+ HF models, 15+ benchmarks (github.com)](github.com)
- [MME: A Comprehensive Evaluation Benchmark for Multimodal Large Language Models](#)

# Multimodal Models
## Selection & Benchmarks

Hence, it may be more appropriate to select models based on the various benchmarks depending on the task at hand, for example:

- **MME**:
  A Comprehensive Evaluation Benchmark for Multimodal Large Language Models

  Suitable for tasks like image captioning, where the model needs to describe what's happening in an image based on understanding the visual content and relationships between objects.

# Multimodal Models
## Selection & Benchmarks

- **MMMU(Massive multi discipline Multimodal Understanding):**

  Massive multi-discipline tasks demanding college-level subject knowledge and deliberate reasoning.

  Includes 11.5K meticulously collected multimodal questions from college exams, quizzes, and textbooks, covering six core disciplines: Art & Design, Business, Science, Health & Medicine, Humanities & Social Science, and Tech & Engineering.
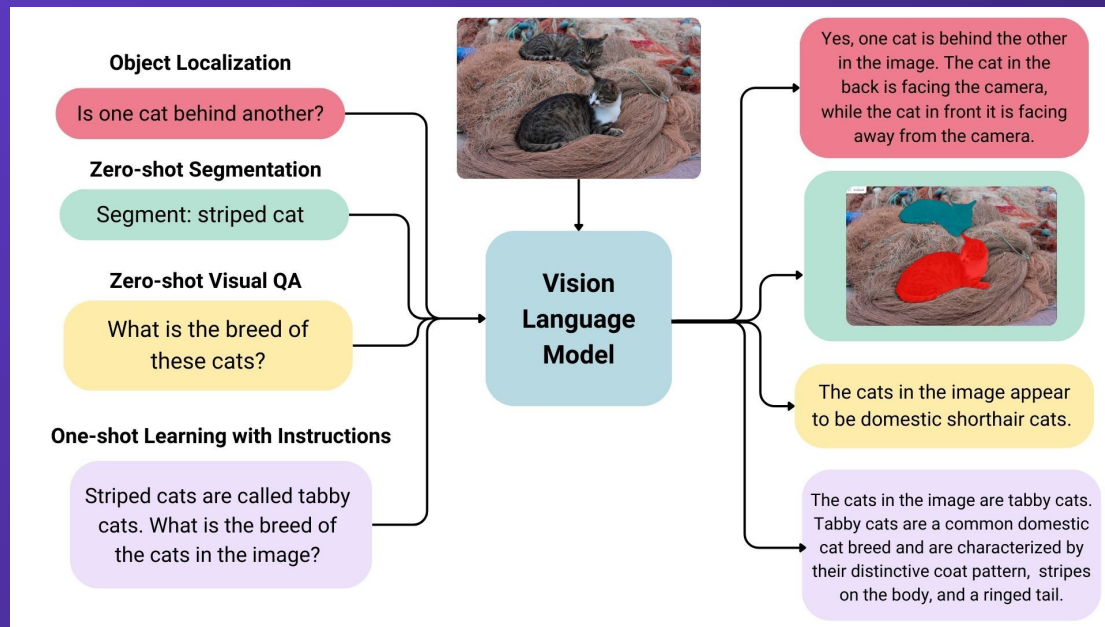
- **Vision Arena:**
  Based on anonymous voting of model output.



(top) MMMU areas of reasoning benchmarks.

# What is a Vision Language Model?

- Vision Language models are defined as multimodal models that can learn from both images and text.
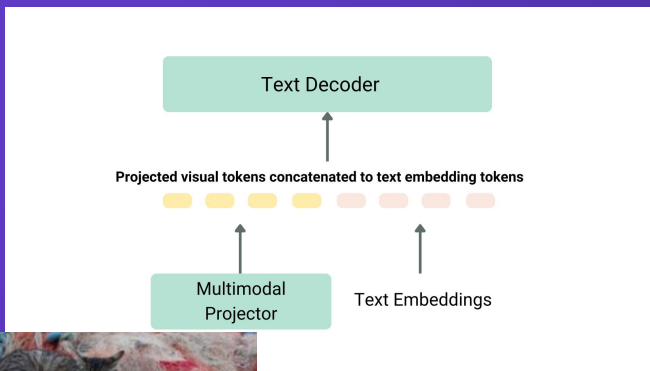


(top) illustrations of vision language model examples

# What is a Vision Language Model?

## Example structure of Vision Language Models (Visual Question Answering - VQA)

*(response) Yes, one cat is behind the other in the image. The cat in the back is facing the camera, while the cat in front is facing away from the camera.*

Text Decoder

Projected visual tokens concatenated to text embedding tokens

Multimodal Projector

Text Embeddings

*(task) Object Localization: (question) Is one cat positioned behind the other?*

## Comparison with LLM

The primary function of LLM models is to:

- Input: Receive text input

With vision-language models, we expand this functionality to:

- Input: Receive a combination of image and text input

This enables the model to not only understand and generate text, but also to integrate visual information and generate images or text that correspond to the input image and text.

TIL-AI
TODAY I LEARNED AI