# Module 10 - Extra

# Multimodal
# Large Language Models
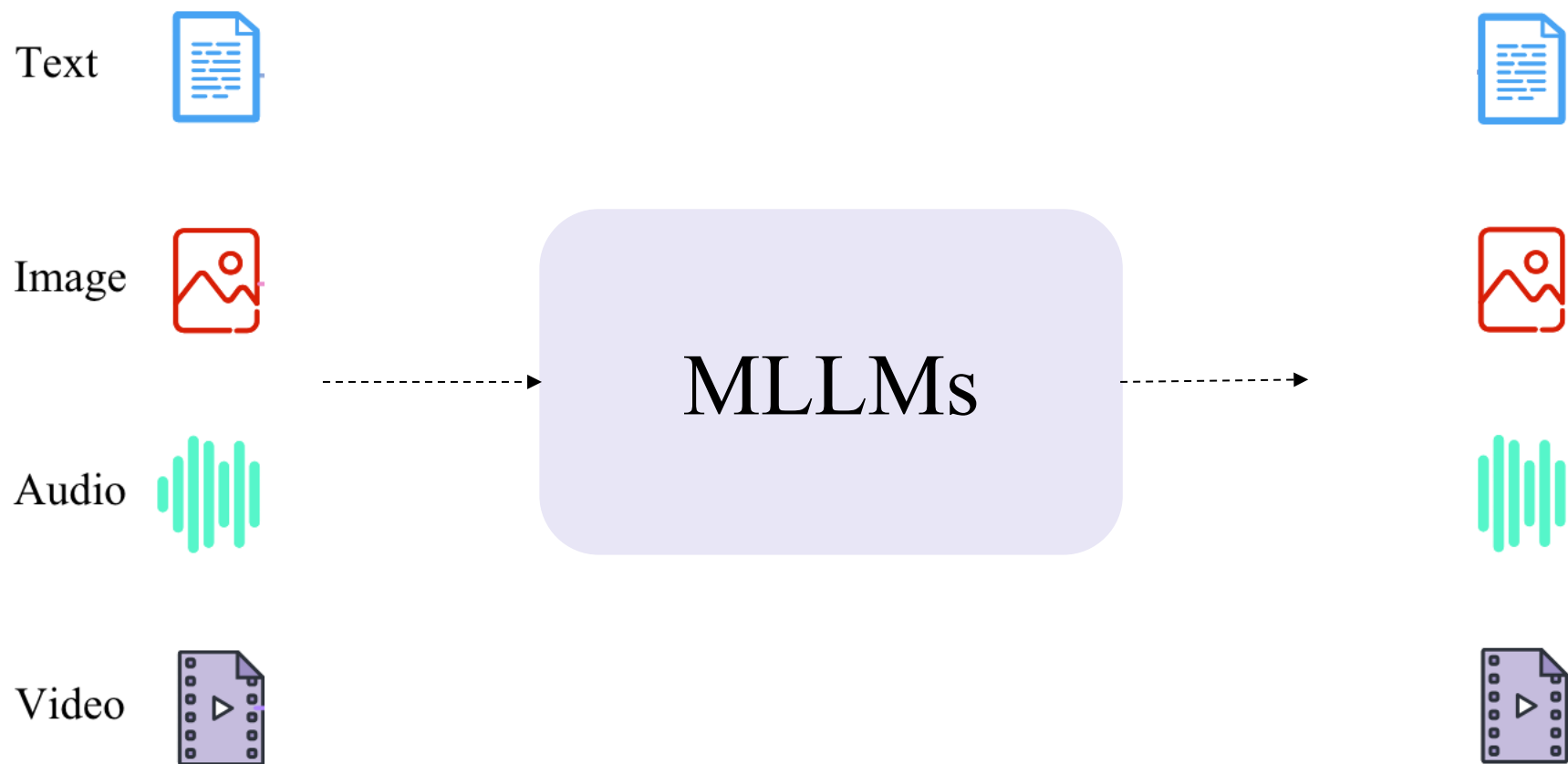
**AI VIET NAM**

**Nguyen Quoc Thai**

**AI VIET NAM**
@aivietnam.edu.vn

**!** **Multimodal Large Language Models**
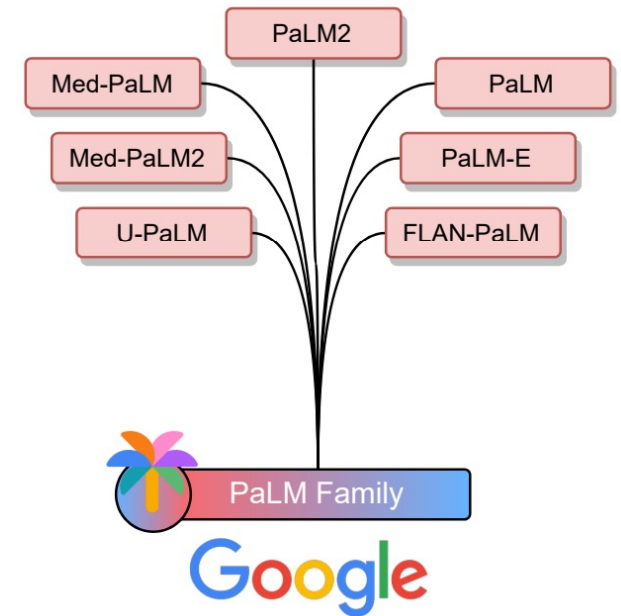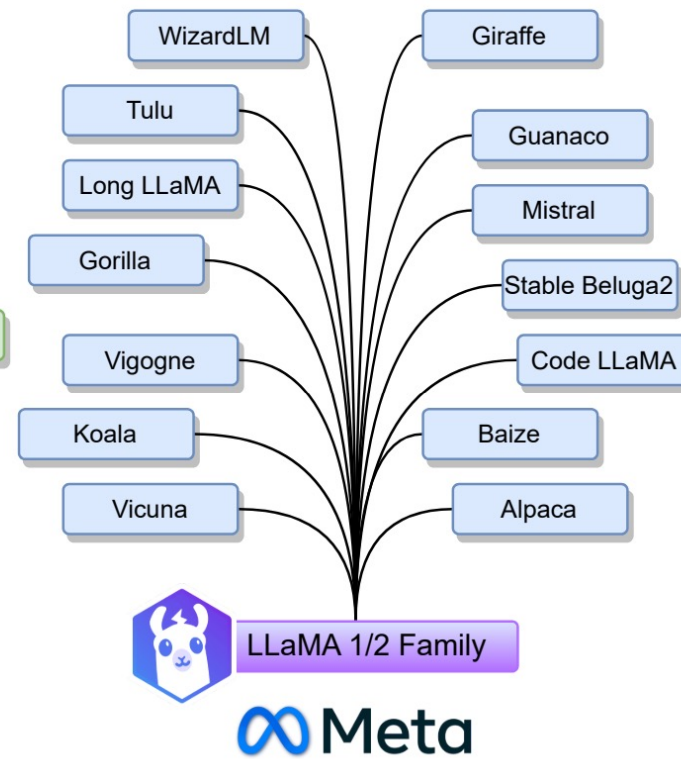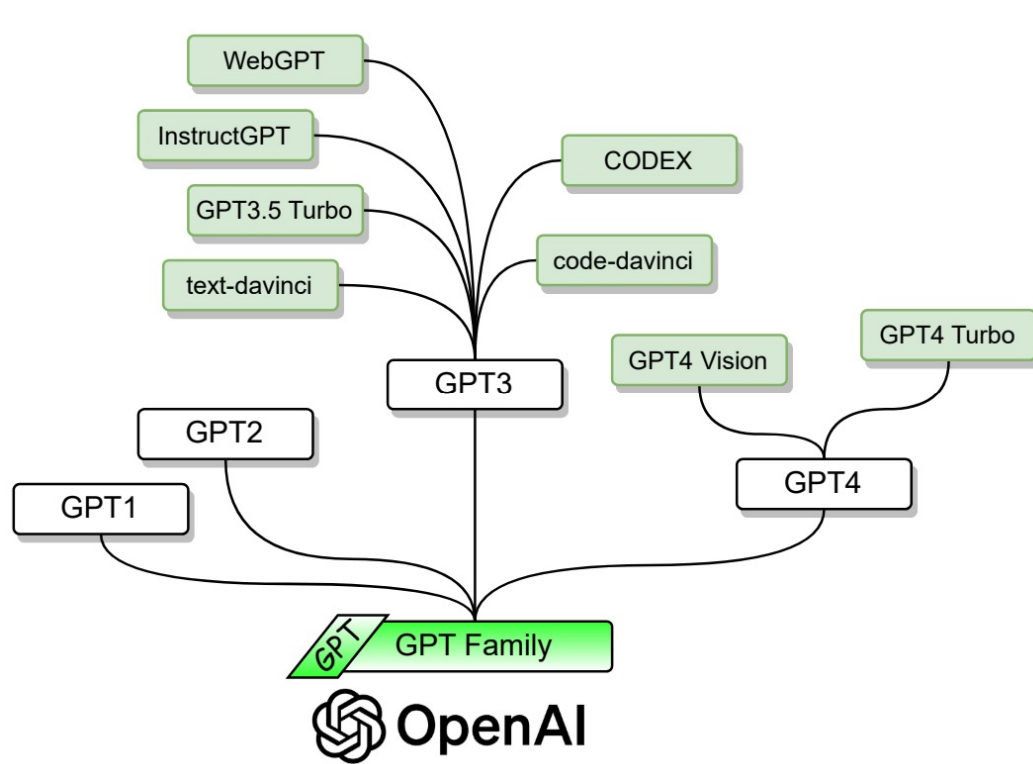
Text

Image

Audio

Video

MLLMs

# Outline

- ➢ **Introduction**
- ➢ **Multimodal Large Language Models**
- ➢ **BLIP-2**
- ➢ **NExT-GPT: Any-to-Any MLLM**

# Introduction

**Large Language Models**

# Introduction

**!** **Large Language Models**

➢ "Very" large LMs: models of 100+ billion parameters

GPT3 (175B), BLOOM (176B), PaLM (540B), GLaM (1200B)…

➢ Data scale: usually in the order of trillions of tokens
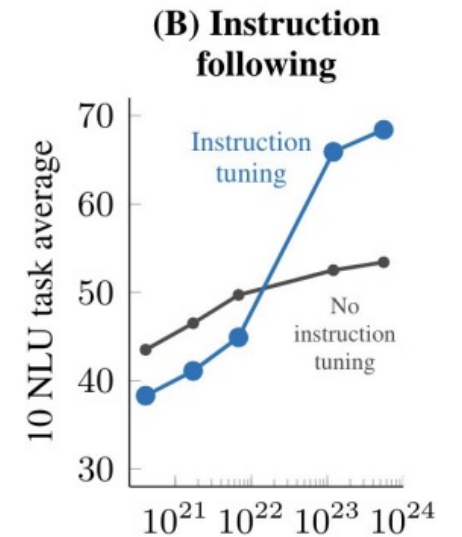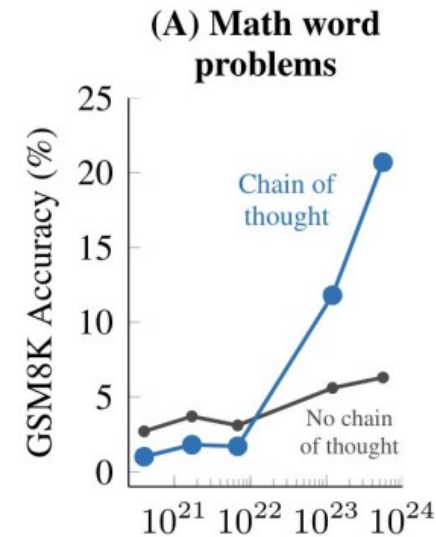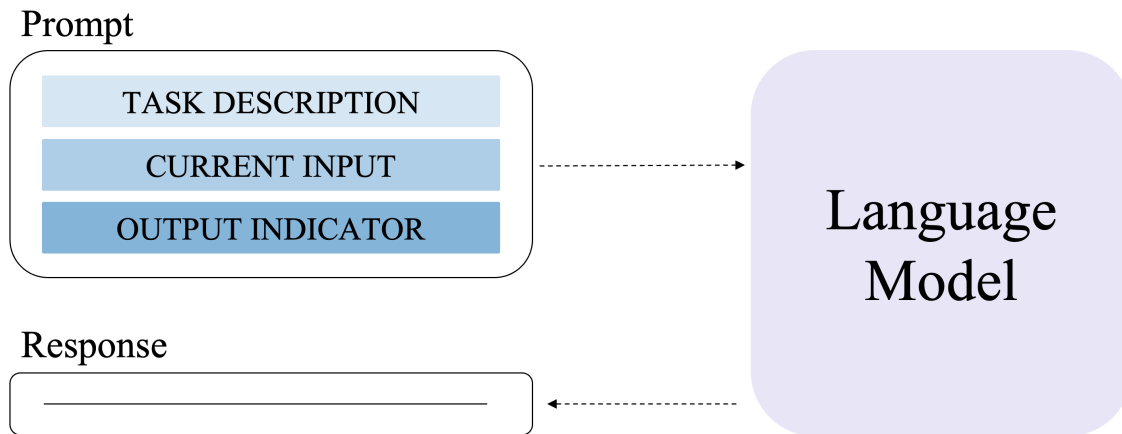
GPT3 (0.5 trillion tokens), LLaMa (1.4 trillion tokens)

# **Introduction**

! **Large Language Models**

➢ The promise: one single model to solve many NLP tasks
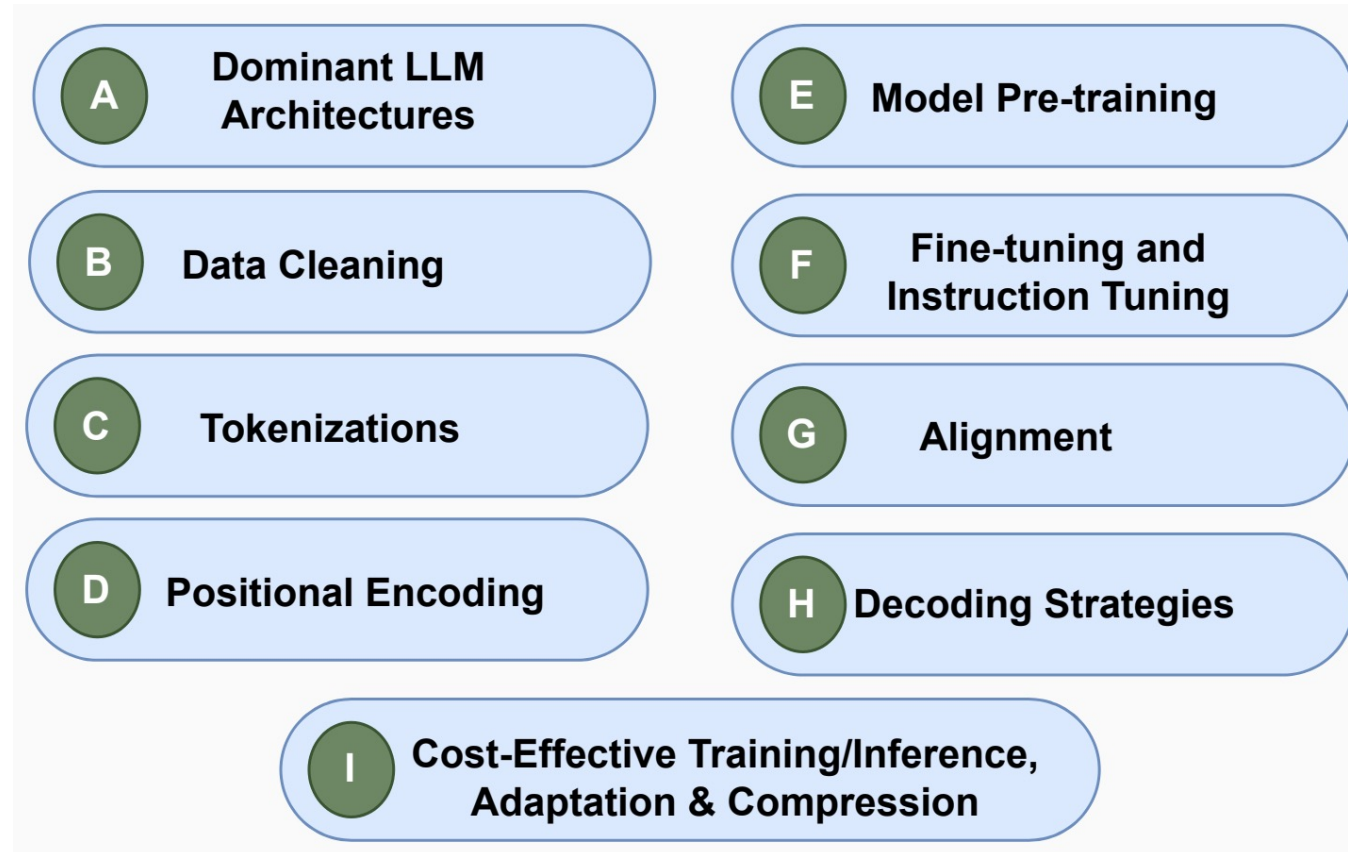
➢ Emergent properties in LLMs

Prompt

| TASK DESCRIPTION |
| CURRENT INPUT |
| OUTPUT INDICATOR |

Language Model

Response

_____



(A) Math word problems

Chain of thought

No chain of thought

GSM8K Accuracy (%)

$10^{21}$ $10^{22}$ $10^{23}$ $10^{24}$

(B) Instruction following

Instruction tuning

No instruction tuning

10 NLU task average

$10^{21}$ $10^{22}$ $10^{23}$ $10^{24}$

6

AI VIET NAM
@aivietnam.edu.vn

**Large Language Models**

# **Introduction**

! **Large Language Models**



A — Dominant LLM Architectures

B — Data Cleaning

C — Tokenizations

D — Positional Encoding

E — Model Pre-training

F — Fine-tuning and Instruction Tuning

G — Alignment

H — Decoding Strategies

I — Cost-Effective Training/Inference, Adaptation & Compression

8

**AI VIET NAM**
@aivietnam.edu.vn

**!** **Large Language Models**

➢ Solve many NLP tasks

Prompt

| Machine Translation |
|:---:|
| … |
| Reasioning |

LLMs

Response

9

# **Introduction**

! **Multimodal Large Language Models**

Text

Image

MLLMs

Audio

Video

# Outline

# Multimodal LLMs

@aivietnam.edu.vn

**The milestones of Multimodal LLMs**

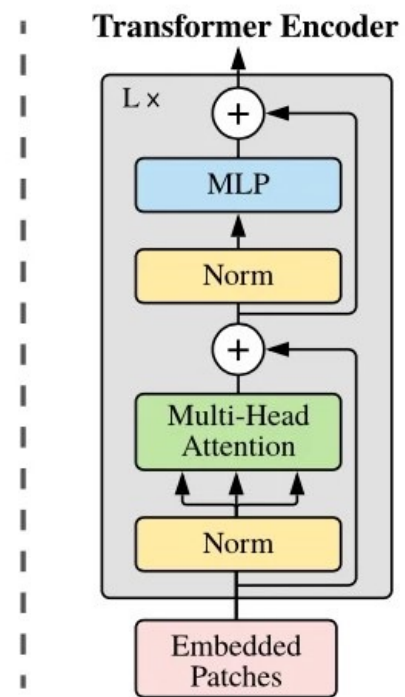12

# Multimodal LLMs

**!** **Architecture**

# Multimodal LLMs

**!** **Architecture – Modility Encoder**

➢ Encode inputs from diverse modalities to obtain corresponding features
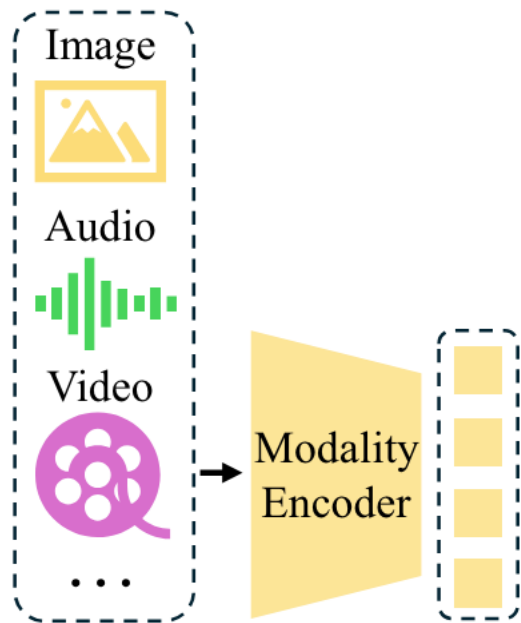
# Multimodal LLMs

**Architecture – Modility Encoder**

➢ Image/Video Encoder: ViT

# Multimodal LLMs

**Architecture – Modility Encoder**

➢ Image/Video Encoder: ViT/ CLIP ViT/ Eva-CLIP ViT

# Multimodal LLMs

**Architecture – Modility Encoder**

➢ Audio Encoder: C-Former / HuBERT / BEATs / Whisper / **CLAP**

**Architecture – Modility Encoder**

➢ Audio Encoder: C-Former / HuBERT / BEATs / Whisper / **CLAP**

# Multimodal LLMs

**!** **Architecture – Modility Encoder**

➢ IMAGEBLIND: One Embedding Space To Bind Them All

➢ Join embedding space enables novel multimodal capabilities

# Multimodal LLMs

**! Architecture – Modility Encoder**

➢ IMAGEBLIND: One Embedding Space To Bind Them All

➢ Join embedding space enables novel multimodal capabilities

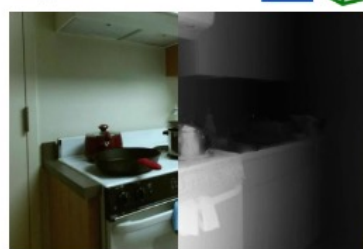**Architecture – Connecter (Input Projector)**

➢ Align the encoded features of other modalities with the text feature

➢ Linear Projector / MLP / Cross-attention / Q-Former / P-Former

# Multimodal LLMs

**Architecture – LLMs**

➢ LLMs: PaLM, LLaMA, Vicuna, Qwen,…

**AI VIET NAM**
@aivietnam.edu.vn

! **Architecture – Output Projector**

➢ Output Projector: maps the signal token representation from LLM into features



➢ MLP / Tiny Transformer

23

## ! Architecture – Modality Generator

➤ Product outputs in distinct modalities

➤ **Stable Diffusion Model**

# Multimodal LLMs

**Architecture – Modality Generator**

➢ **Stable Diffusion Model for Image**



25

# Multimodal LLMs

**Architecture – Modality Generator**

➤ **Stable Diffusion Model for Audio (AudioLDM)**



(a) Training and sampling process of AudioLDM

(b) Audio inpainting with AudioLDM

(c) Audio style transfer with AudioLDM

# Multimodal LLMs

! **Architecture**

# Multimodal LLMs

**Training Strategy – Pre-Training**

➢ Align different modalities and learn multimodal world knowledge

➢ Entails large-scale text-paired data

Input: <image>
Response: {caption}

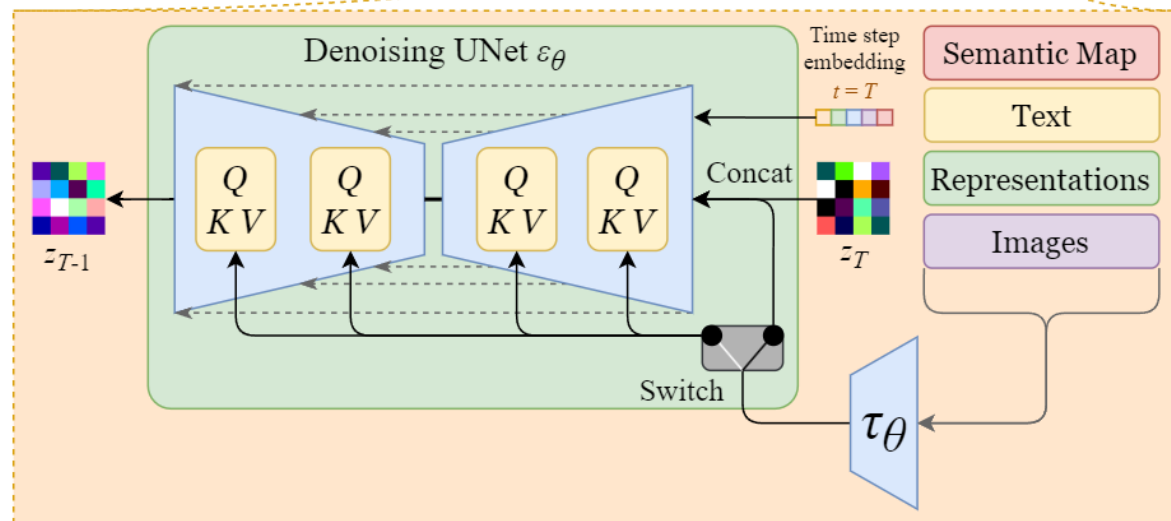| Dataset | Samples | Date |
|---|---|---|
| **Coarse-grained Image-Text** | | |
| CC-3M [84] | 3.3M | 2018 |
| CC-12M [85] | 12.4M | 2020 |
| SBU Captions [86] | 1M | 2011 |
| LAION-5B [87] | 5.9B | Mar-2022 |
| LAION-2B [87] | 2.3B | Mar-2022 |
| LAION-COCO [88] | 600M | Sep-2022 |
| COYO-700M [90] | 747M | Aug-2022 |
| **Fine-grained Image-Text** | | |
| ShareGPT4V-PT [83] | 1.2M | Nov-2023 |
| LVIS-Instruct4V [91] | 111K | Nov-2023 |
| ALLaVA [92] | 709K | Feb-2024 |
| **Video-Text** | | |
| MSR-VTT [93] | 200K | 2016 |
| **Audio-Text** | | |
| WavCaps [94] | 24K | Mar-2023 |

! **Training Strategy – Instruction-Tuning**

Below is an instruction that describes a task. Write a response that appropriately completes the request

Instruction: <instruction>
Input: {<image>, <text>}
Response: <output>

**!** **Training Strategy – Instruction-Tuning**

- <Image> {Question}
- <Image> Question: {Question}
- <Image> {Question} A short answer to the question is
- <Image> Q: {Question} A:
- <Image> Question: {Question} Short answer:
- <Image> Given the image, answer the following question with no more than three words. {Question}
- <Image> Based on the image, respond to this question with a short answer: {Question}. Answer:
- <Image> Use the provided image to answer the question: {Question} Provide your answer as short as possible:
- <Image> What is the answer to the following question? "{Question}"
- <Image> The question "{Question}" can be answered using the image. A short answer is

# Multimodal LLMs

! **Training Strategy – Instruction-Tuning**

| Dataset | Sample | Modality | Source | Composition |
|---------|--------|----------|--------|-------------|
| LLaVA-Instruct | 158K | I + T → T | MS-COCO | 23K caption + 58K M-T QA + 77K reasoning |
| LVIS-Instruct | 220K | I + T → T | LVIS | 110K caption + 110K M-T QA |
| ALLaVA | 1.4M | I + T → T | VFlan, LAION | 709K caption + 709K S-T QA |
| Video-ChatGPT | 100K | V + T → T | ActivityNet | 7K description + 4K M-T QA |
| VideoChat | 11K | V+T → T | WebVid | description + summarization + creation |
| Clotho-Detail | 3.9K | A + T → T | Clotho | caption |

**AI VIET NAM**
@aivietnam.edu.vn

! **Training Strategy – Alignment Tuning**

# Multimodal LLMs

## ! SOTA MLLMs

| Model | I/O | Modality Encoder | Input Projector | LLM | Output Projector | Modality Generator |
|-------|-----|------------------|-----------------|-----|------------------|--------------------|
| BLIP-2 | IT => T | CLIP ViT | Q-Former Linear | Flan-T5 OPT | - | - |
| LLaVA | IT => T | CLIP ViT | Linear | Vicuna | - | - |
| miniGPT-4 | IT => T | Eva-CLIP ViT | Q-Former Linear | Vicuna | - | - |
| InstructBLIP | IVT => T | ViT | Q-Former Linear | Flan-T5 Vicuna | - | - |
| Next-GPT | IVAT => IVAT | ImageBlind | Linear | Vicuna | Tiny Transformer | Stable Diffusion Model |
| ModaVerse | IVAT => IVAT | ImageBlind | Linear | LLaMA2 | MLP | Stable Diffusion Model |

# Outline

- ➤ **Introduction**
- ➤ **Multimodal Large Language Models**
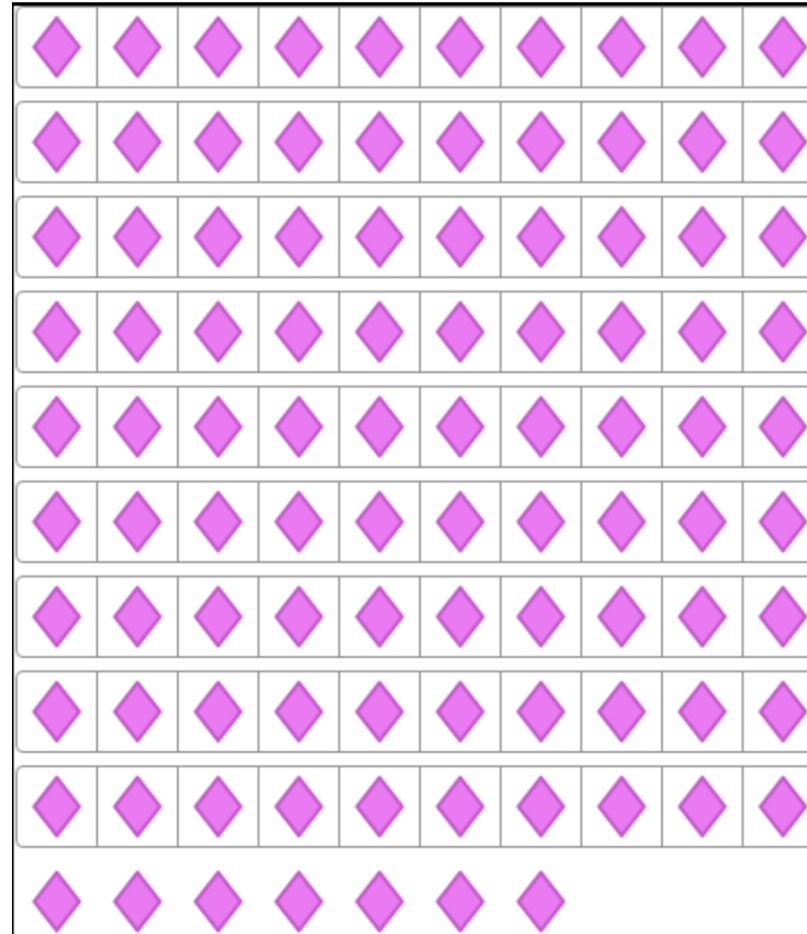- ➤ **BLIP-2 for Visual Question Answering**
- ➤ **NExT-GPT: Any-to-Any MLLM**

# BLIP-2 for VQA

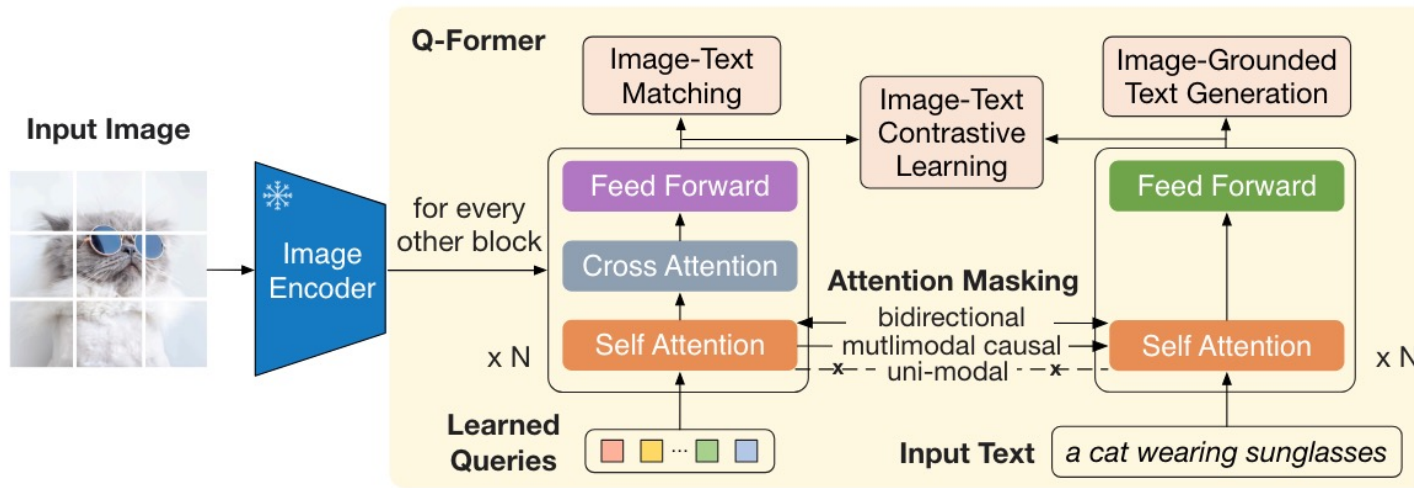**VQA Dataset**
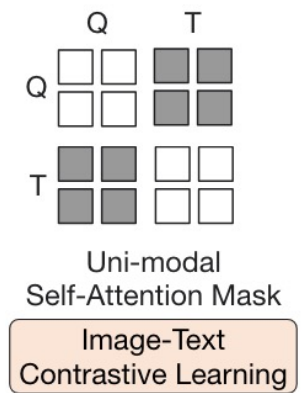
Question:
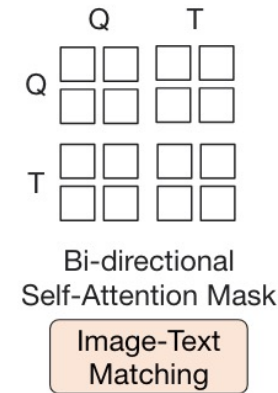How many diamonds are there?



Response:
97

**!** **BLIP-2 - Training**

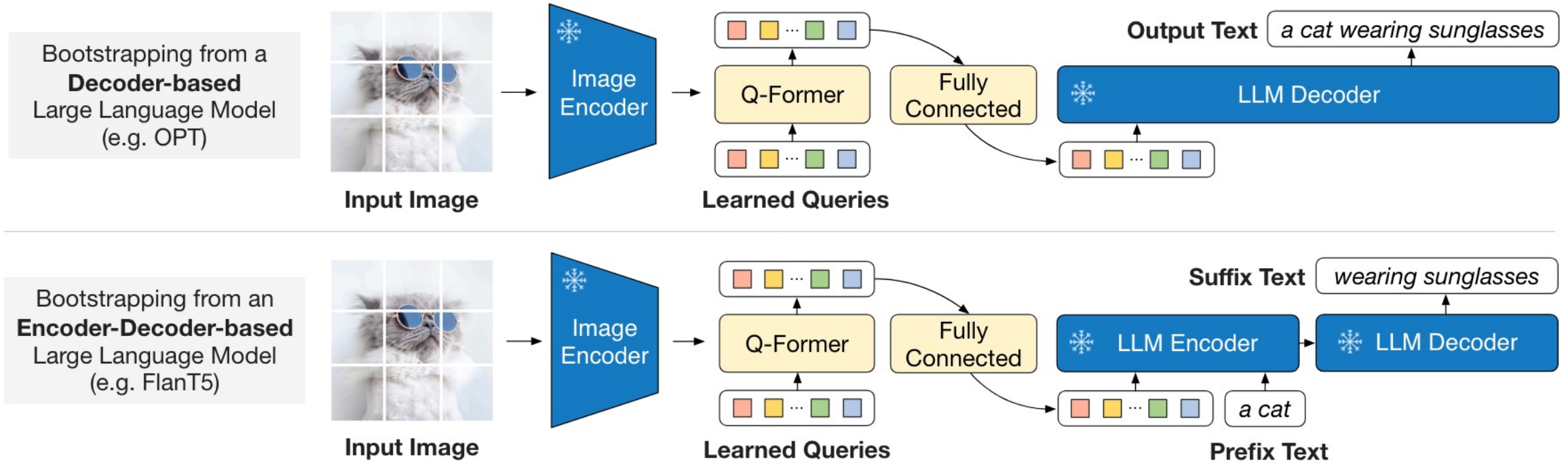| Model | I/O | Modality Encoder | Input Projector | LLM | Output Projector | Modality Generator |
|-------|-----|------------------|-----------------|-----|------------------|--------------------|
| BLIP-2 | IT => T | CLIP ViT | Q-Former Linear | Flan-T5 OPT | - | - |

**AI VIET NAM**
@aivietnam.edu.vn
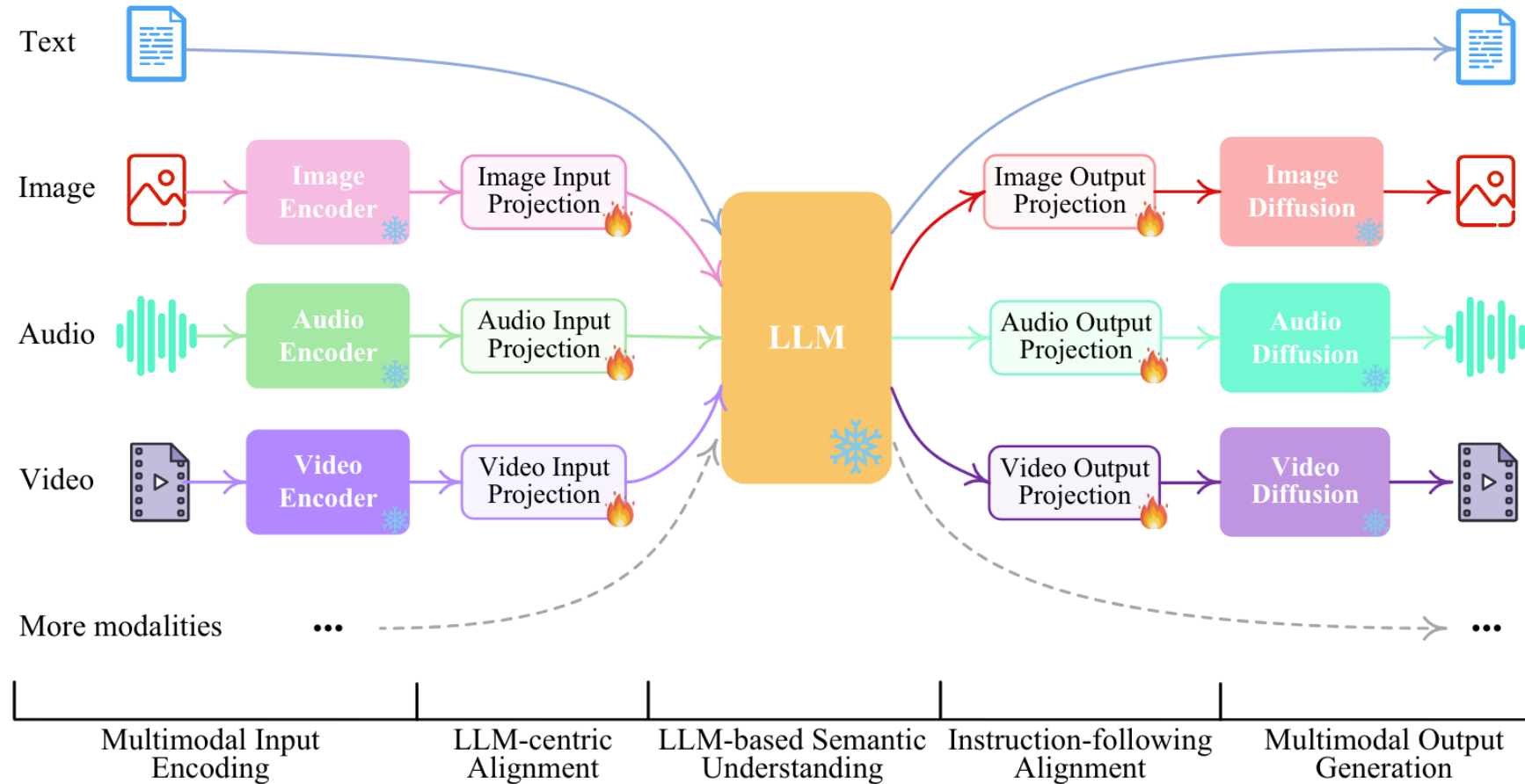
## ! BLIP-2 - Inference

# Multimodal LLMs

! BLIP-2 - Demo

# Outline

- ➤ **Introduction**
- ➤ **Multimodal Large Language Models**
- ➤ **BLIP-2 for Visual Question Answering**
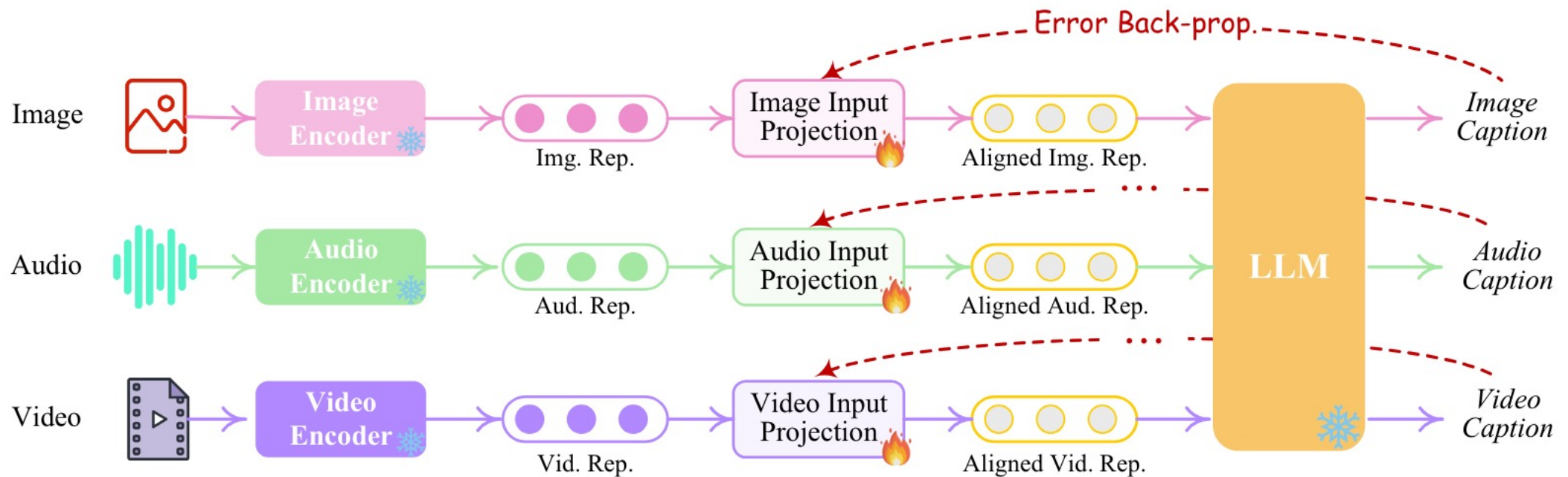- ➤ **NExT-GPT: Any-to-Any MLLM**

# NExT-GPT

**NExT-GPT**

# NExT-GPT

**!** **NExT-GPT**

| Model | I/O | Modality Encoder | Input Projector | LLM | Output Projector | Modality Generator |
|---|---|---|---|---|---|---|
| Next-GPT | IVAT => IVAT | ImageBlind | Linear | Vicuna | Tiny Transformer | Stable Diffusion Model |

| | Encoder | | Input Projection | | LLM | | Output Projection | | Diffusion | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Name | Param | Name | Param | Name | Param | Name | Param | Name | Param |
| **Text** | — | — | — | — | | | — | — | — | — |
| **Image** | | | | | Vicuna [12] | 7B❄️ | Transformer | 31M🔥 | SD [68] | 1.3B❄️ |
| **Audio** | ImageBind [25] | 1.2B❄️ | Linear | 4M🔥 | (LoRA | 33M🔥) | Transformer | 31M🔥 | AudioLDM [51] | 975M❄️ |
| **Video** | | | | | | | Transformer | 32M🔥 | Zeroscope [8] | 1.8B❄️ |

41

## NExT-GPT: Lightweight Multimodal Alignment Learning

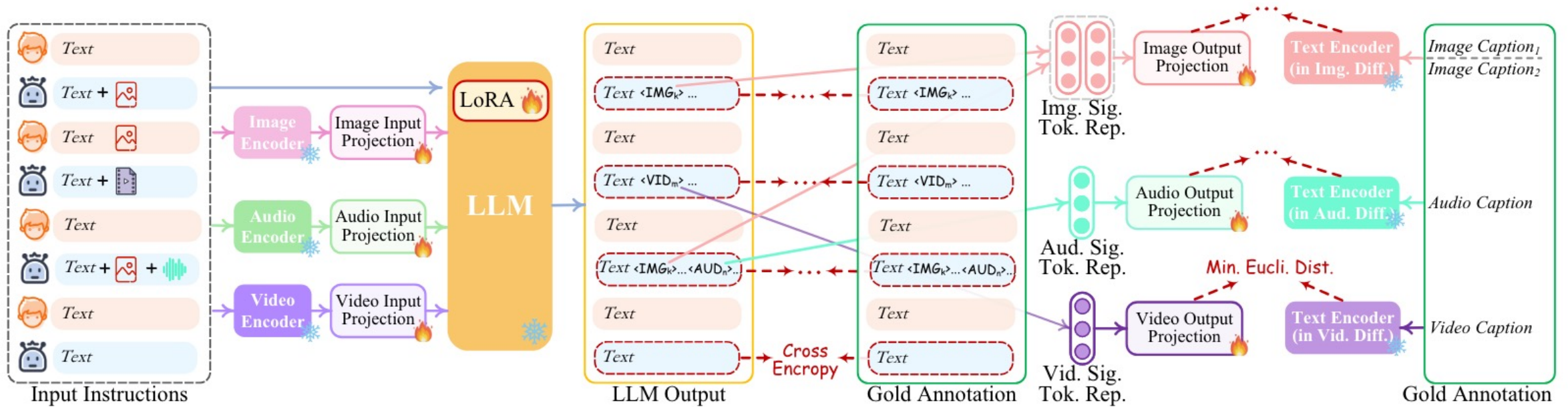➢ Encoding-side LLM-centric Multimodal Alignment

**AI VIET NAM**
@aivietnam.edu.vn

**!** **NExT-GPT: Lightweight Multimodal Alignment Learning**

➢ Decoding-side Instruction-following Alignment



43

# NExT-GPT

**!** **NExT-GPT: Lightweight Multimodal Alignment Learning**

➤ Modality-switching Instruction Tuning

**!** **NExT-GPT - Demo**

45

# NExT-GPT - Demo

# NExT-GPT

**!** **NExT-GPT - Demo**

## ! NExT-GPT - Demo

# NExT-GPT

## ! NExT-GPT - Demo

# Thanks!

## Any questions?