模式分析与机器智能
工业和信息化部重点实验室
MIIT Key Laboratory of
Pattern Analysis & Machine Intelligence

ParNeC | 模式识别与神经计算研究组
PAttern Recognition and NEural Computing

# Multimodal Models

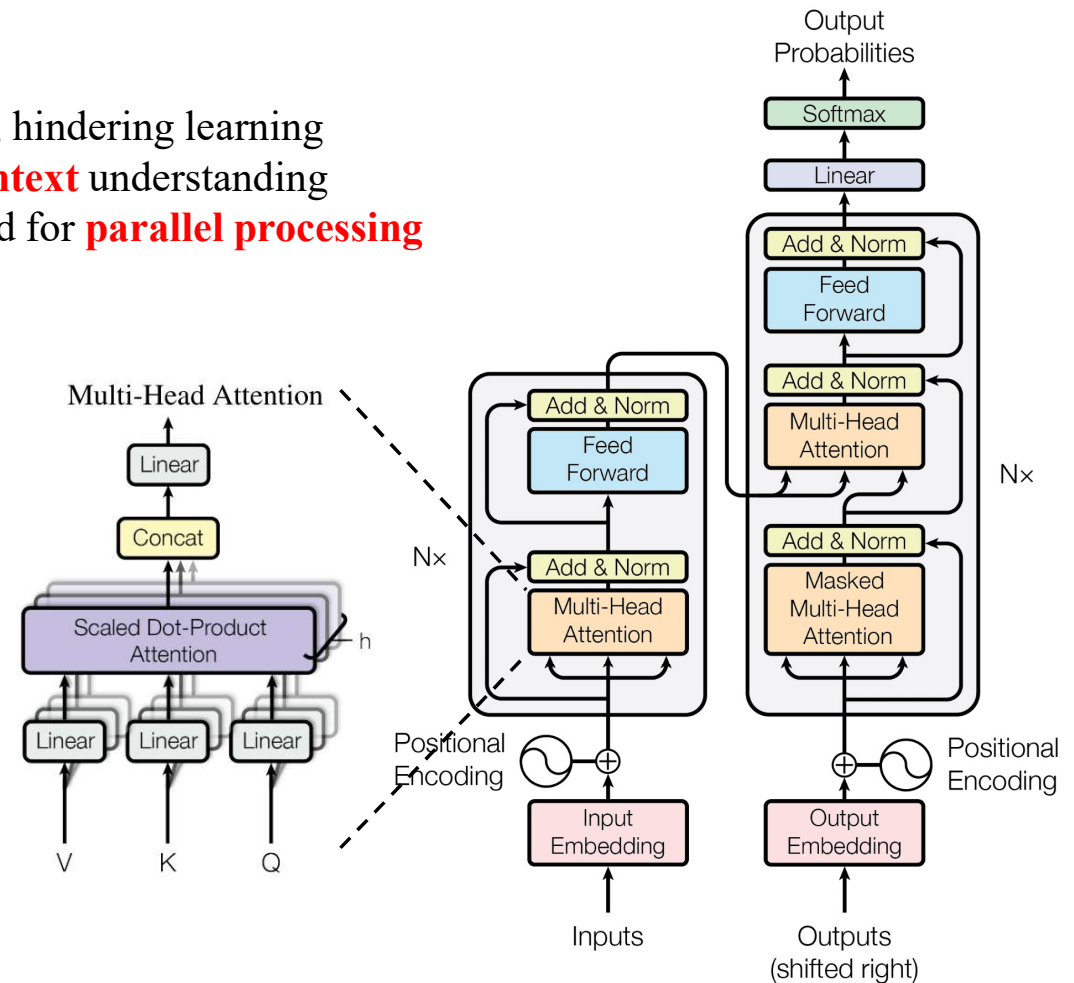## 3. Multi-modal Large Language Model

# Background

## Transformer

Drawbacks of RNN and LSTM:

➤ Prone to vanishing and exploding gradient problems, hindering learning
➤ More biased toward recent data instead of **global context** understanding
➤ Not able to take advantage of modern GPUs designed for **parallel processing**

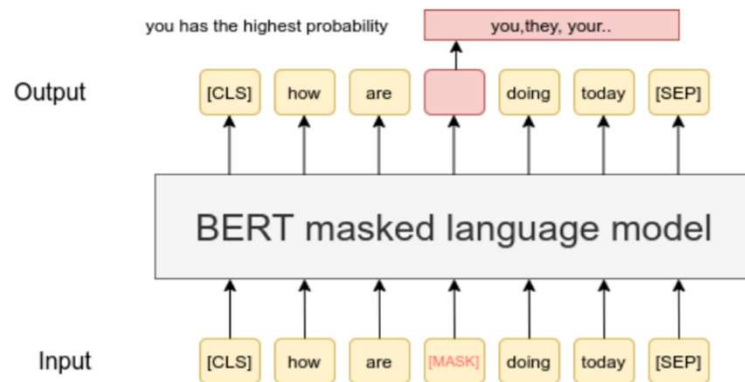Main contributions of the Transformer paper:

➤ Encoder-decoder structure
   Self attention and cross attention
➤ Multi-head attention
   Each head could attend to different feature
➤ Causal modeling
   Causal mask in the decoder prevents positions
   from attending to future positions
➤ Positional embedding
   Adding perturbation to the sequence to
   differentiate the order of tokens

Attention is All You Need (June 2017 Google)

ParNeC 模式识别与神经计算研究组
PAttern Recognition and NEural Computing

## Different Architecture for Different NLP Tasks

BERT – Masked Language Modeling (Encoder only)

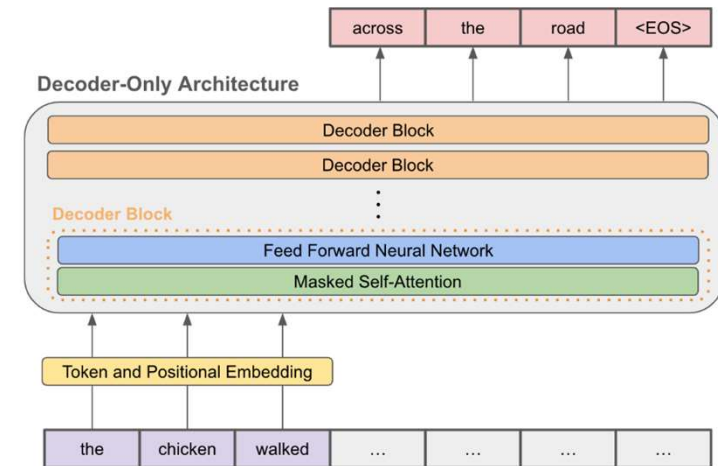

GPT – Causal Language Modeling (Decoder only)



**Pros:**
Better understanding for global information. Suitable for almost every NLP task like sentiment analysis, token classification etc.
**Cons:**
Fixed context window size.
Not suitable for generative tasks.
High annotation cost.

**Pros:**
Best choice for generative tasks.
Extrapolable context window size.
**Cons:**
Limited understanding for global context.

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (October 2018)
GPT: Improving Language Understanding by Generative Pre-Training (June 2018)

**ParNeC** 模式识别与神经计算研究组
PAttern Recognition and NEural Computing

## **Training Steps for LLMs**

➢ Pretraining
Use causal loss to perform next-token prediction
training on large corpus of data

➢ Instruction tuning
Supervised training on instruction tuning dataset
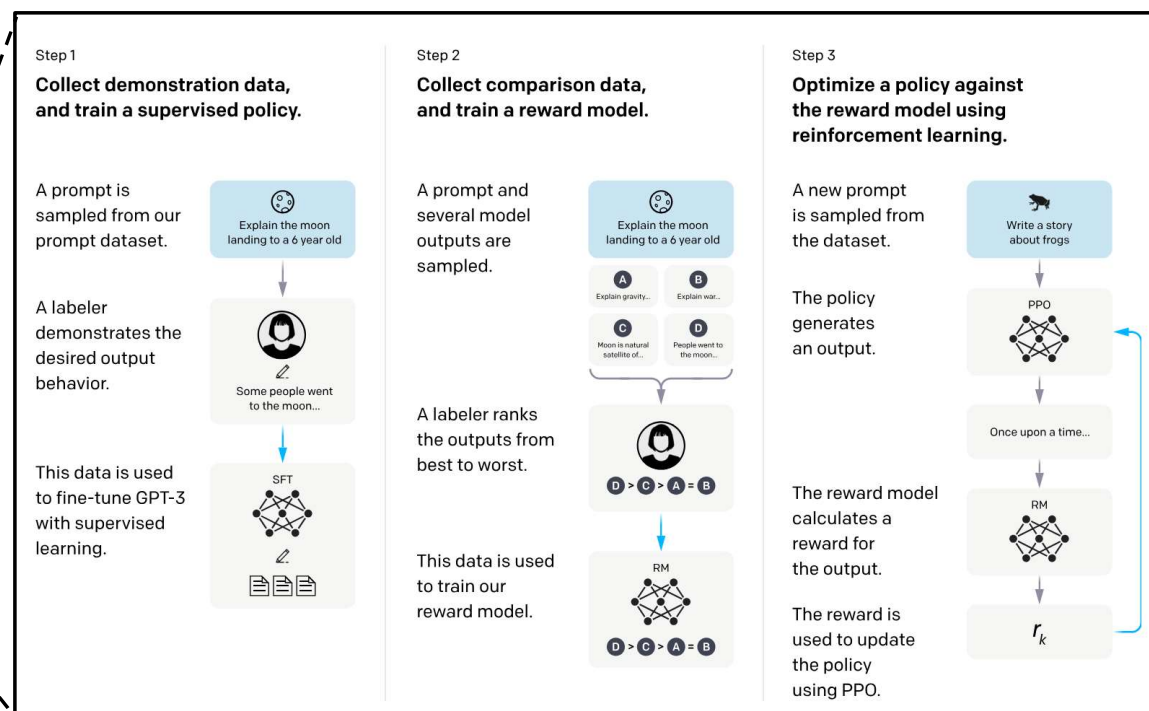to make LLM better follow users' instructions

$$L_1 = -\sum_t \log P_\theta(x_{t+1} \mid x_{t:1}),$$

➢ Alignment. LLM should generate contents
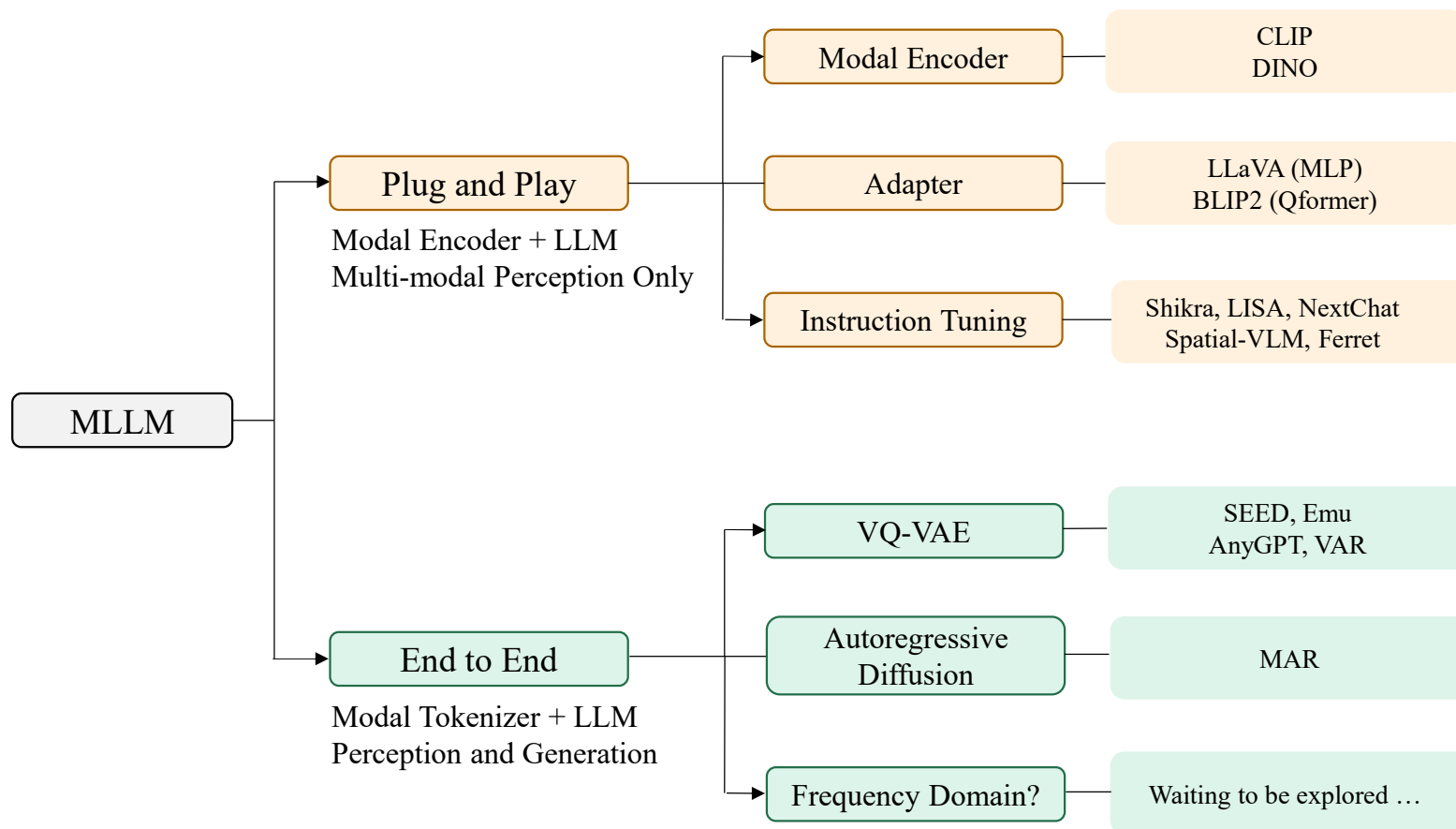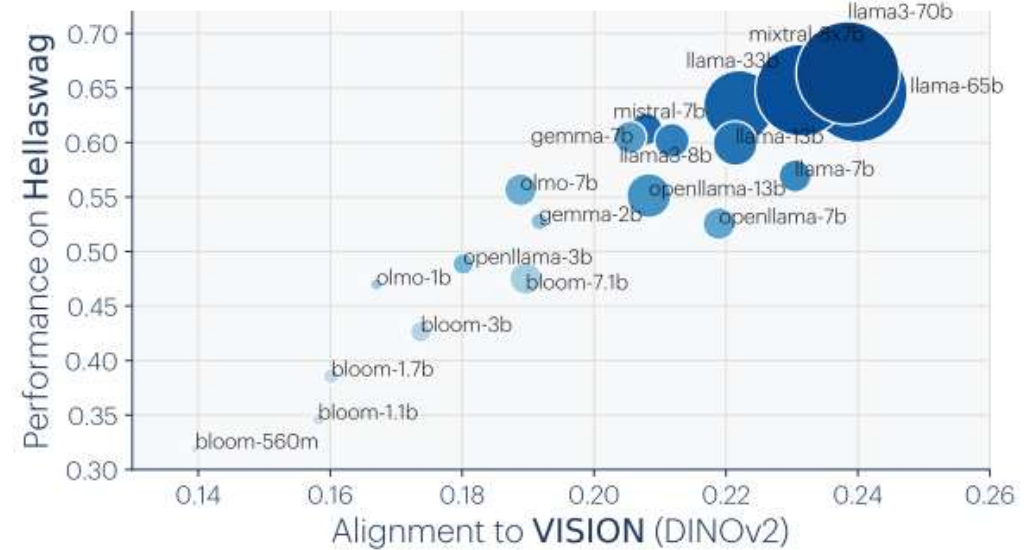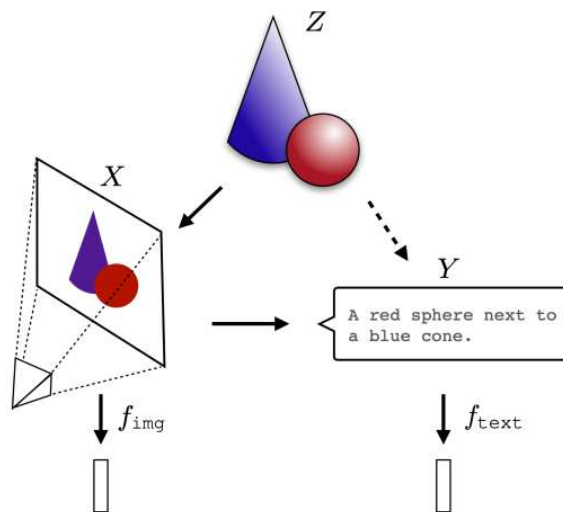that align with human's preference (safety…)
RLHF (PPO)
DPO
…



**Step 1**
**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon…

This data is used to fine-tune GPT-3 with supervised learning.

SFT

**Step 2**
**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A: Explain gravity… B: Explain war…
C: Moon is natural satellite of… D: People went to the moon…

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

**Step 3**
**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time…

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

Training language models to follow instructions with human feedback (Mar 2022)

# Introduction

## Content

**ParNeC** | 模式识别与神经计算研究组
PAttern Recognition and NEural Computing

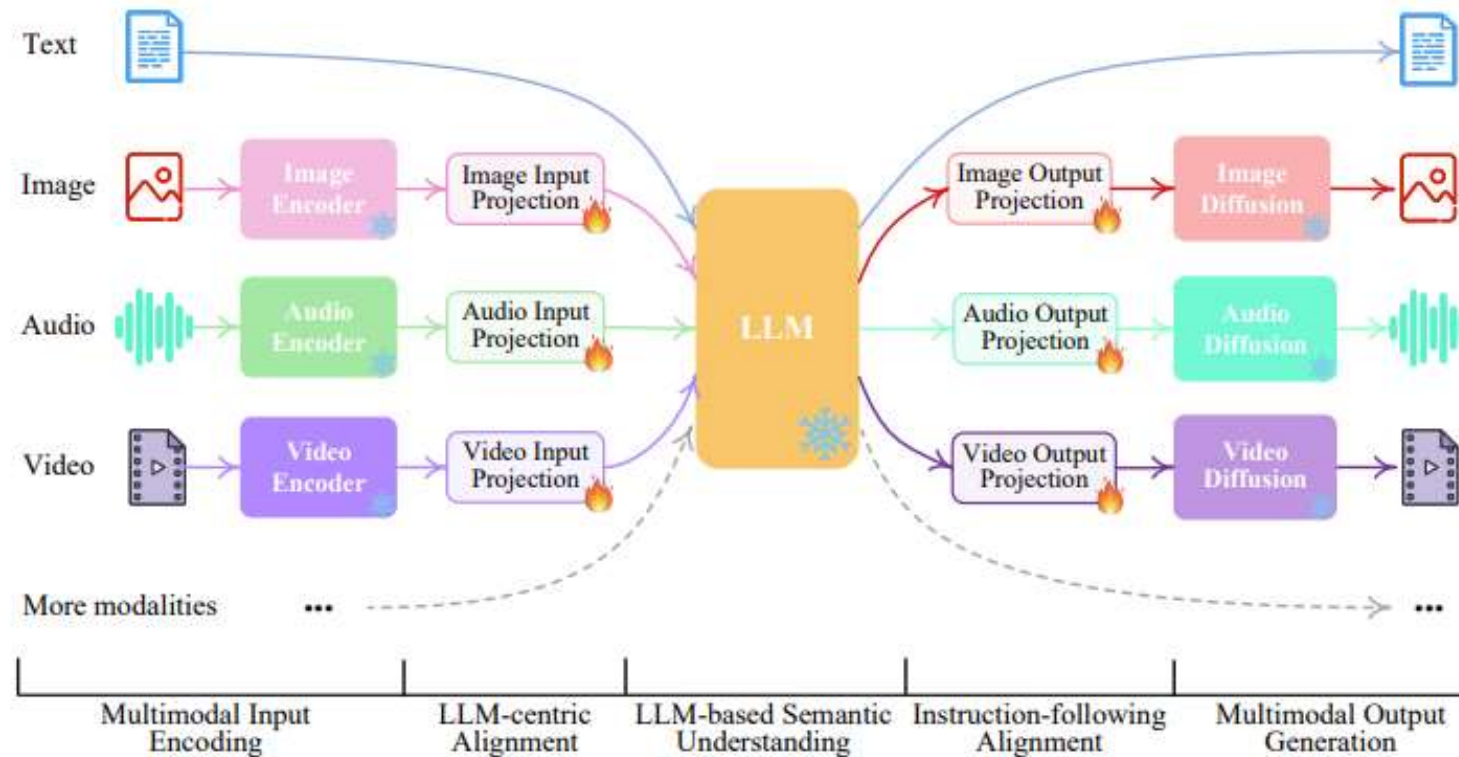## The Meta Question: Why Could Multi-modal Reasoning Exist?



The representations across multiple modalities will converge on a shared representation of Z, and scaling model size, as well as data and task diversity, drives this convergence.

$\Longrightarrow$

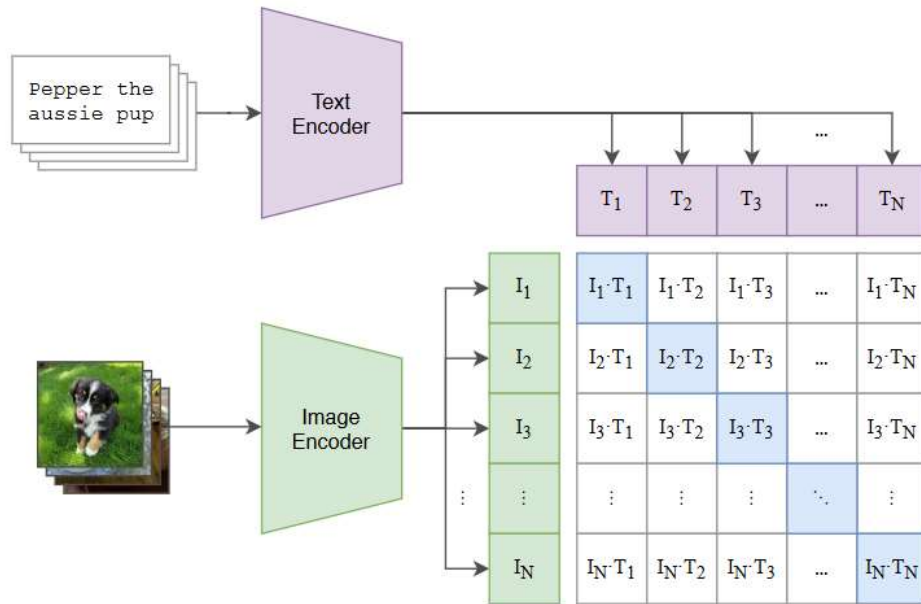We can build MLLM using existing **visual encoder** and **LLM** with suitable adapter to deal with modal alignment.

The Platonic Representation Hypothesis (May 2024)

Connect an LLM with multimodal adaptors and different diffusion decoders, enabling it to perceive inputs and generate outputs in arbitrary combinations of **text**, **image**, **video**, and **audio**

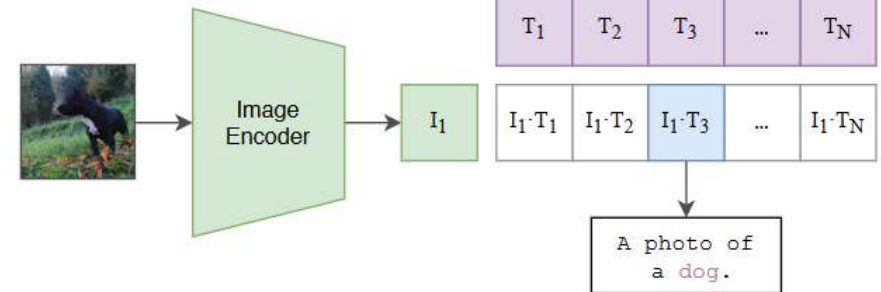NExT-GPT: Any-to-Any Multimodal LLM (Sep 2023)

**CLIP**



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

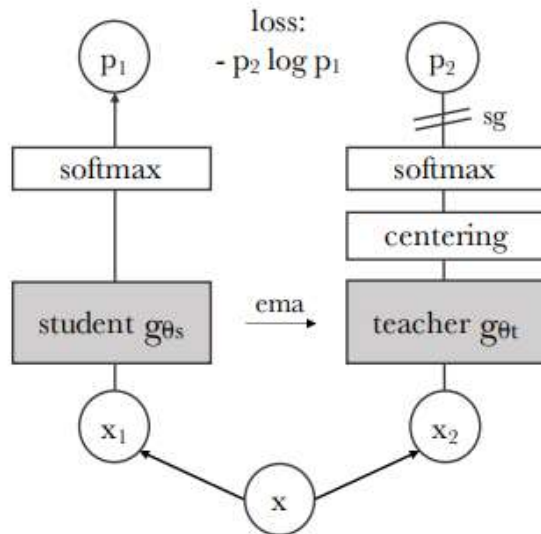- ➢ Pretrained on 400 million image-text pairs just using **contrastive** loss
- ➢ Powerful **zero-shot** capability
- ➢ The most popular visual encoder for MMLM

Learning Transferable Visual Models From Natural Language Supervision (Feb 2021 OpenAI)

## DINO

**Motivation**: The success of Transformers in NLP was the use of self-supervised pretraining (BERT, GPT)



Figure 1: **Self-attention from a Vision Transformer with** $8 \times 8$ **patches trained with no supervision.** We look at the self-attention of the [CLS] token on the heads of the last layer. This token is not attached to any label nor supervision. These maps show that the model automatically learns class-specific features leading to unsupervised object segmentations.

From a given image, generate a set $V$ of different views. This set contains two global views, $x_1^g$ and $x_2^g$ and several local views of smaller resolution. All **crops** are passed through the **student** while only the **global views** are passed through the **teacher.** Align the predictions of two networks on feature dimension $K$.

$$\text{loss\_v1} = \sum_{x \in \{x_1^g, x_2^g\}} \sum_{\substack{x' \in V \\ x' \neq x}} H(P_t(x), \boxed{P_s(x')}) \longrightarrow P_s(x)^{(i)} = \frac{\exp(g_{\theta_s}(x)^{(i)}/\tau_s)}{\sum_{k=1}^{K} \exp(g_{\theta_s}(x)^{(k)}/\tau_s)}$$

$$\text{loss v2} = \text{loss\_v1} + \sum_{i=1}^{N} m_i \cdot P_{\boldsymbol{\theta'}}^{\text{patch}}(\boldsymbol{u}_i)^{\text{T}} \log P_{\boldsymbol{\theta}}^{\text{patch}}(\hat{\boldsymbol{u}}_i)$$

Image-level        Patch-level

DINO: Emerging Properties in Self-Supervised Vision Transformers (May 2021 Meta)
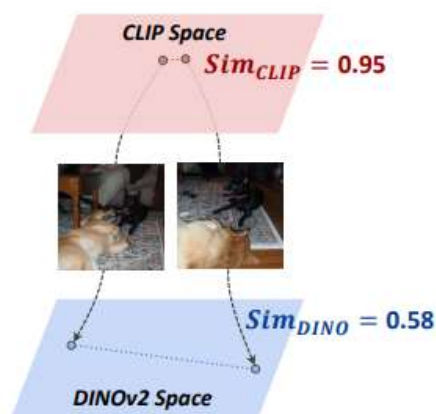DINOv2: Learning Robust Visual Features without Supervision (Feb 2024 Meta)

## CLIP v.s. DINO   How to Choose Between These Two?

**Motivation:** *CLIP-blind pairs* – images that CLIP perceives as similar despite their clear visual differences, is the main source of incorrect answers and hallucinated explanations.



Step 1

**Finding CLIP-blind ⌀ pairs.**

Discover image pairs that are proximate in CLIP feature space but distant in DINOv2 feature space.

**CLIP Space** $Sim_{CLIP} = 0.95$

$Sim_{DINO} = 0.58$

**DINOv2 Space**

Step 2

**Spotting the difference between two images.**

For a CLIP-blind pair, a human annotator attempts to spot the visual differences and formulates questions.

"The dog's head in the left image is resting on the carpet, while the dog's head in the right image is lying on the floor."

Formulating questions and options for both images.

Where is the yellow animal's head lying in this image?
(a) Floor  (b) Carpet

Step 3

**Benchmarking multimodal LLMs.**

Evaluate multimodal LLMs using a CLIP-blind image pair and its associated question.

Where is the yellow animal's head lying in this image?
(a) Floor  (b) Carpet

(b) Carpet ✓     (b) Carpet ✗

✗ (no score for this pair)

The model receives a score only when **both** predictions for the CLIP-blind pair are correct.

DINO is more suitable for fine-grained level perception, object detection, semantic segmentation
Merging visual representations from CLIP and DINO t leads to improved performance in visual grounding tasks.

Eyes Wide Shut? Exploring the Visual Shortcomings of Multimodal LLMs (Apr 2024)

# Adapter

## LLaVA (MLP)

**Motivation**: Use a two-layer MLP to project visual feature into language space.



Figure 1: LLaVA network architecture.



Stage 1: Pre-training for Feature Alignment. Optimize MLP only on short VQA pairs.

Stage 2: Fine-tuning End-to-End. Use 150K GPT-generated multimodal instruction-following data, plus around 515K VQA data from academic-oriented tasks, to teach the model to follow multimodal instructions.

Visual Instruction Tuning (Dec 2023)

## BLIP2 (QFormer)

**Motivation**: Optimize learnable queries (Query2Label, Perceiver, finetuning CLIP) to merge and down-sample visual feature
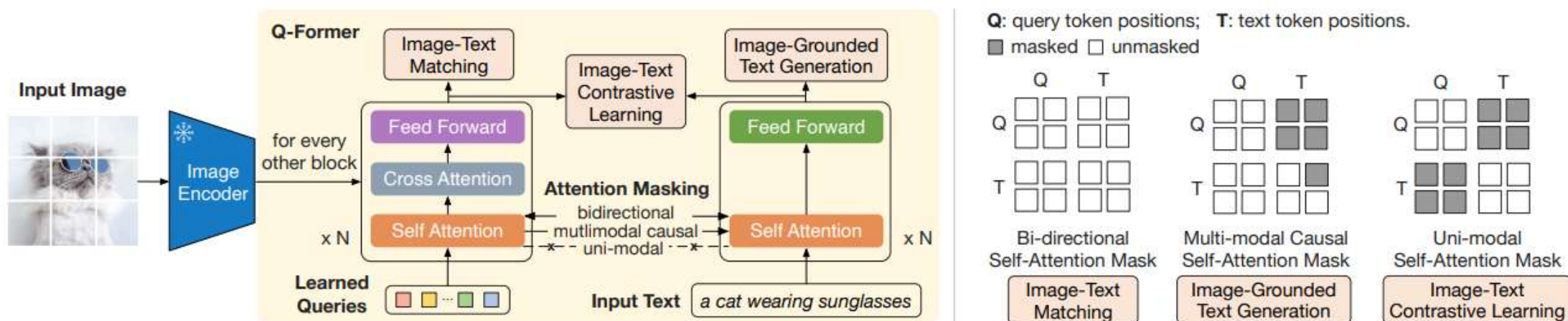


**Image-Text Contrastive Learning (ITC)** learns to align image representation and text representation by contrasting the image-text similarity of a positive pair against those of negative pairs in a batch.
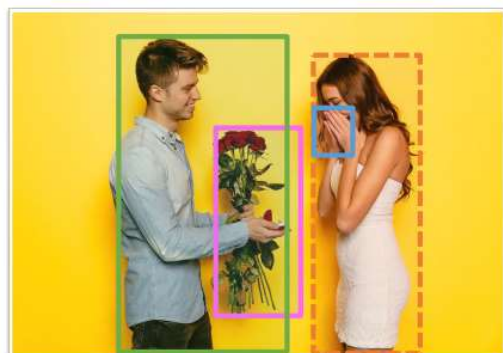
**Image-grounded Text Generation (ITG)** loss trains the Q-Former to generate texts, given input images as the condition.

**Image-Text Matching (ITM)** aims to learn fine-grained alignment between image and text representation. It is a binary classification task where the model is asked to predict whether an image-text pair is positive (matched) or negative (unmatched).

## Shikra

**Motivation**: Endow LLM with fine-grained perception capability like object detection, semantic segmentation
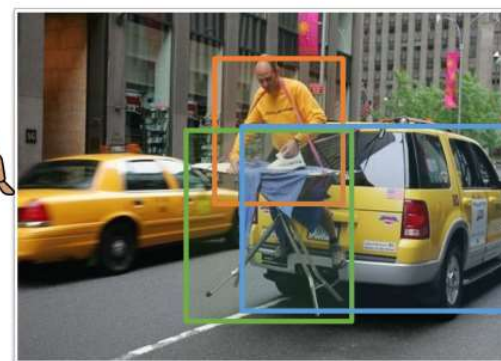


> How is the person in the picture feeling[0.610,0.251,0.825,0.829]? How did you find out?

> The person in the picture [0.610,0.251,0.825,0.829] is feeling surprised and happy. We can infer this from the fact that she is covering her face with her hands [0.612,0.352,0.690,0.444] while looking at the man [0.216,0.223,0.558,0.832] offering her a bunch of red roses [0.412,0.390,0.588,0.756].

> What is unusual about this image? give coordinates [xmin,ymin,xmax,ymax] for the items you reference.

> The unusual aspect of this image is the man [0.392,0.254,0.652,0.530] standing on top of the ironing board [0.338,0.392,0.668,0.756] in the back of the yellow SUV [0.452,0.384,0.998,0.738]. The scene is unusual and unexpected, as one would typically not expect to see someone ironing clothes while standing on top of the car.
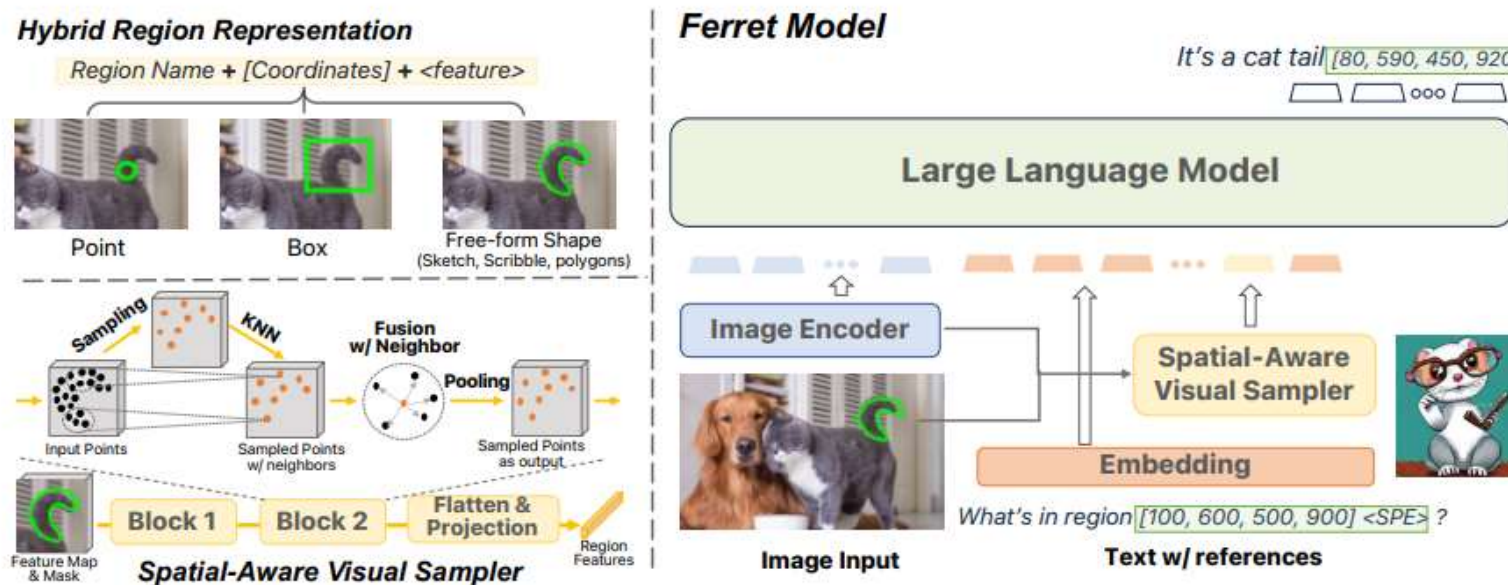
Add more object detection datasets in pretraining stage
Use GPT4 to generate QA pairs with 2d bounding box to perform instruction tuning

Shikra: Unleashing Multimodal LLM's Referential Dialogue Magic (Jul 2023)

## Ferret

**Motivation**: Endow LLM with fine-grained perception capability like object detection, semantic segmentation



**Spatial-aware visual sampler.** The shape of the referred regions can be quite varied, not limited to just points or rectangle boxes. Given extracted image feature map Z and the binary region mask M, we first randomly sample N positive points inside M. For each point, its feature is obtained by bilinear interpolation.
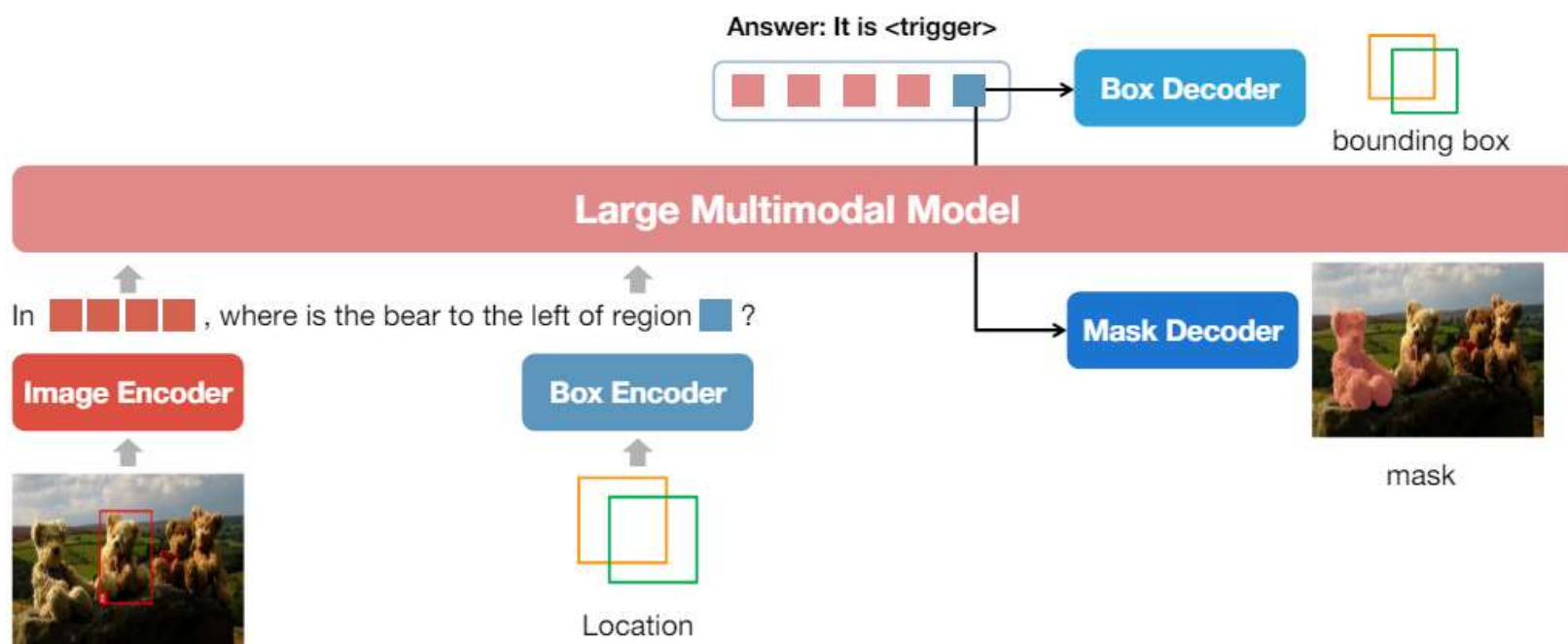
2D coordinates of point x

$$h_{ik} = \sigma([\theta([\mathbf{Z}(x_{ik}) - \mathbf{Z}(x_i); C(x_{ik}) - C(x_i)]); \mathbf{Z}(x_i); C(x_i)]),$$

Ferret: Refer And Ground Anything Anywhere at Any Granularity (Oct 2023)
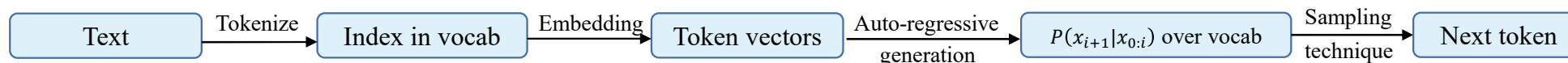
## NeXT-Chat

**Motivation**: Endow LLM with fine-grained perception capability like object detection, semantic segmentation



Introduce a special token, denoted as `<trigger>`, which serves to trigger the localization (detection/segmentation)
Incorporate box encoder and decoder to extract fine-grained localization information
Connect with SAM to perform semantic segmentation

NExT-Chat: An LMM for Chat, Detection and Segmentation (Dec 2023)

## One Unified MLLM Architecture

| Text | $\xrightarrow{\text{Tokenize}}$ | Index in vocab | $\xrightarrow{\text{Embedding}}$ | Token vectors | $\xrightarrow[\text{generation}]{\text{Auto-regressive}}$ | $P(x_{i+1}|x_{0:i})$ over vocab | $\xrightarrow[\text{technique}]{\text{Sampling}}$ | Next token |

How about using light-weight **modal tokenizer** to discrete feature in other modalities as well?

➢ Throw away modal adapter. LLM will process all modal data in one feature space.

➢ Generate image, text, audio… all in auto-regressive manner

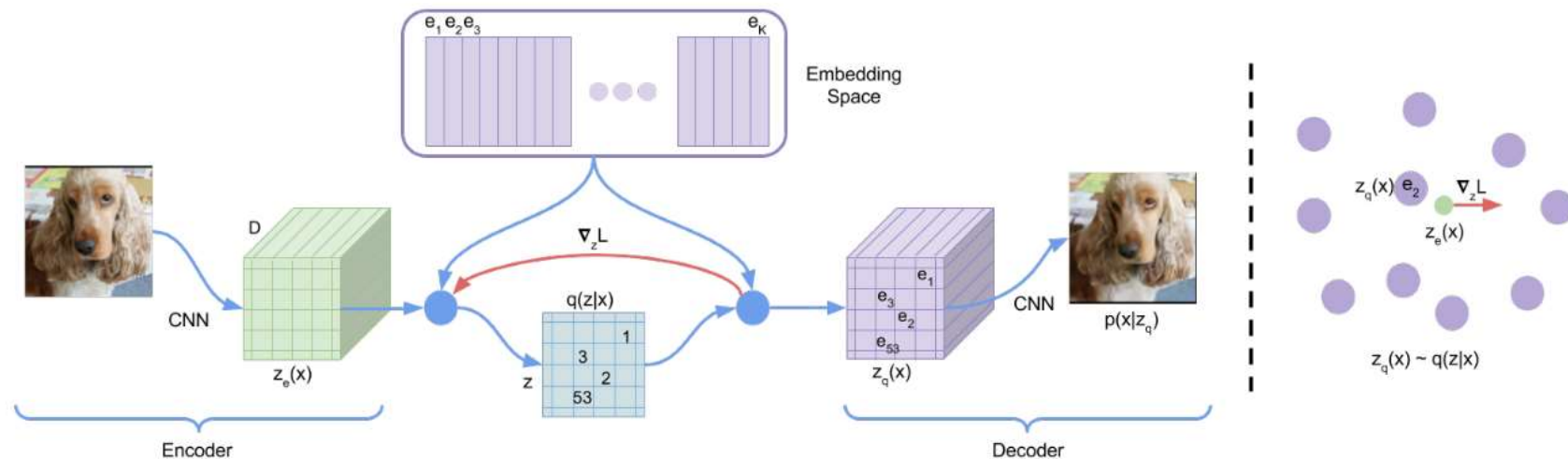Tone perception, real-time conversation, change any voice type you like, image generation…



May 13 2024 GPT-4o Release

**How to Tokenize Image?  Pixel by Pixel? Slow and high cost!**

**Motivation:** Vector Quantized-Variational-AutoEncoder (VQ-VAE), a simple yet powerful **generative** model that learns visual discrete representations, differs from VAEs in two key ways: the encoder network outputs discrete feature



Define a latent embedding space $e \in R^{K \times D}$ where $K$ is the size of the discrete latent space (i.e., a $K$-way categorical), and $D$ is the dimensionality of each latent embedding vector $e_i$. Take an image input $x$ to calculate $z_e(x)$ from encoder, and then find the nearest $e_i$ in codebook

$$z_q(x) = e_k, \quad \text{where} \quad k = \boxed{\text{argmin}_j \|z_e(x) - e_j\|_2} \xrightarrow{\text{No gradient}} \|x - decoder(z + sg[z_q - z])\|_2^2 + \beta\|sg[z] - z_q\|_2^2 + \gamma\|z - sg[z_q]\|_2^2$$

Neural Discrete Representation Learning (May 2018)

<span style="color:red">Can we use VQ-VAE as image tokenizer? Yes and No</span>

**What Stands in The Way of A Unified MLLM Architecture?**

Causal modeling for image

**What's causality?**

*The system will only laugh if you tickle it…*

Text/audio generation: the future token is dependent on past tokens. That's why we can generate text/audio sequence in a auto-regressive manner. (1D causal system)
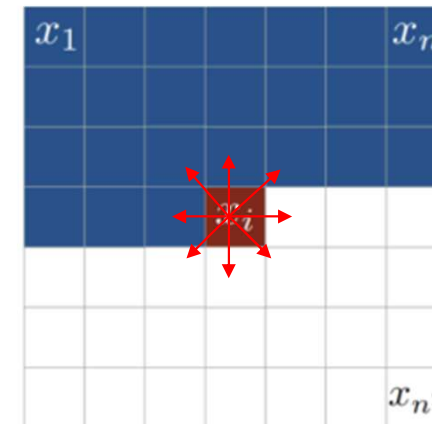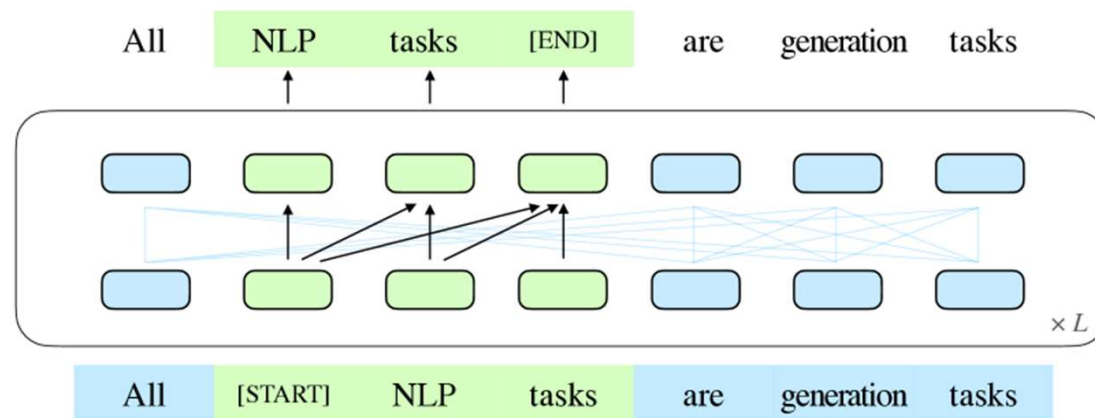


Image generation: A pixel in image space has relation with other 8 pixels in its neighborhood. So 1D causal modeling is a wrong representation for image. (2D semi-causal system)

**How to define causality for image?**
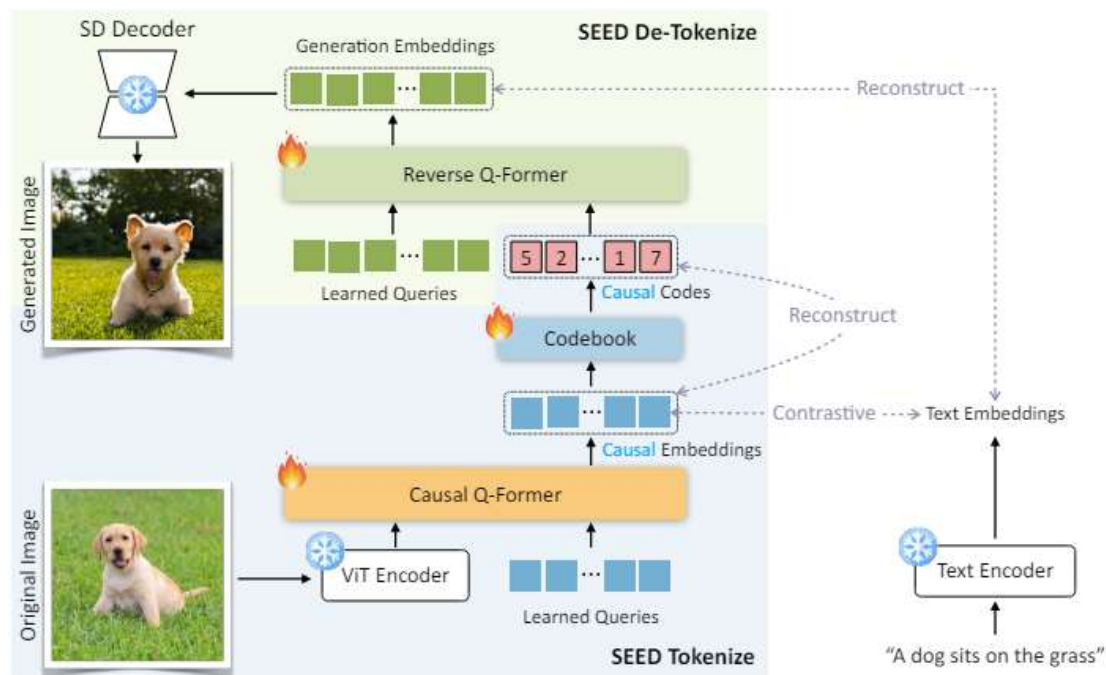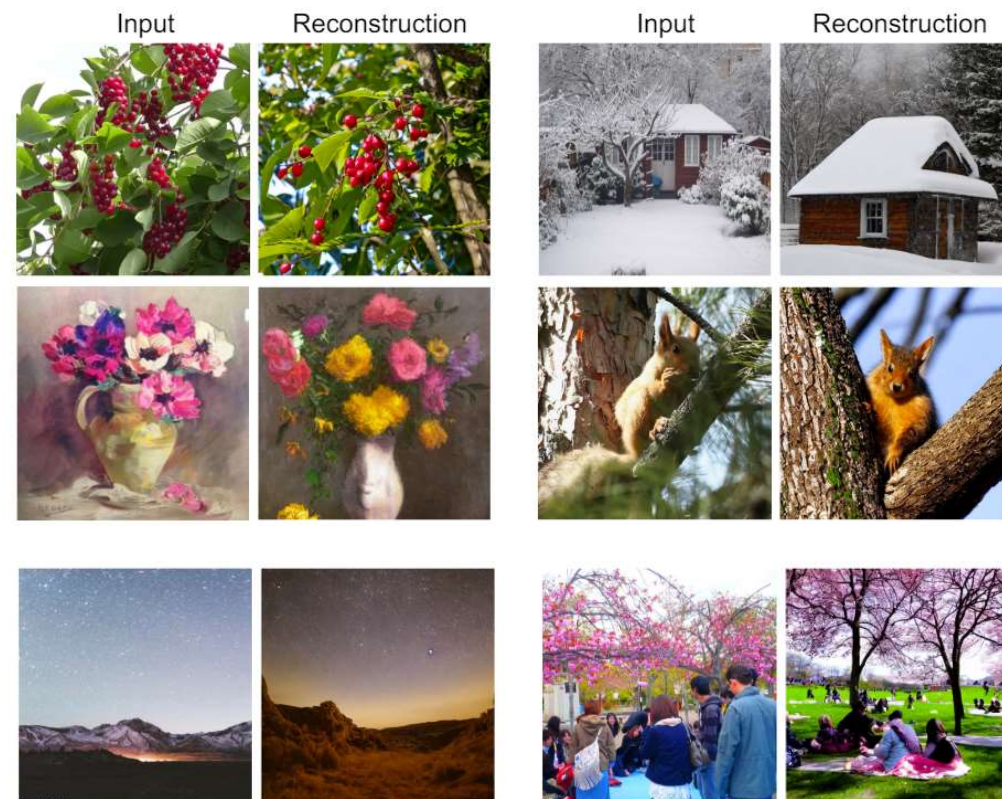
## Causal Modeling for Image



Figure 2: Overview of our **SEED** tokenizer, which produces discrete visual codes with causal dependency and high-level semantics.

Training Stage I: Causal Q-Former. Use causal mask in Q-Former to optimize learnable queries by contrastive loss (final token only).

Training Stage II: Visual Quantization and De-tokenization. Train a VQ codebook to discretize the causal embeddings and then employ a Reverse Q-Former to reconstruct the textual features of a frozen stable diffusion model from discrete codes.

Planting a SEED of Vision in Large Language Model (Aug 2023)

**Encode image feature into 1D sequence, losing 2D spatial information!**

# Causal Modeling for Image

**Motivation:** autoregressive image generation can be viewed as coarse-to-fine "next-scale prediction" or "next-resolution prediction", diverging from the standard raster-scan "next-token prediction".
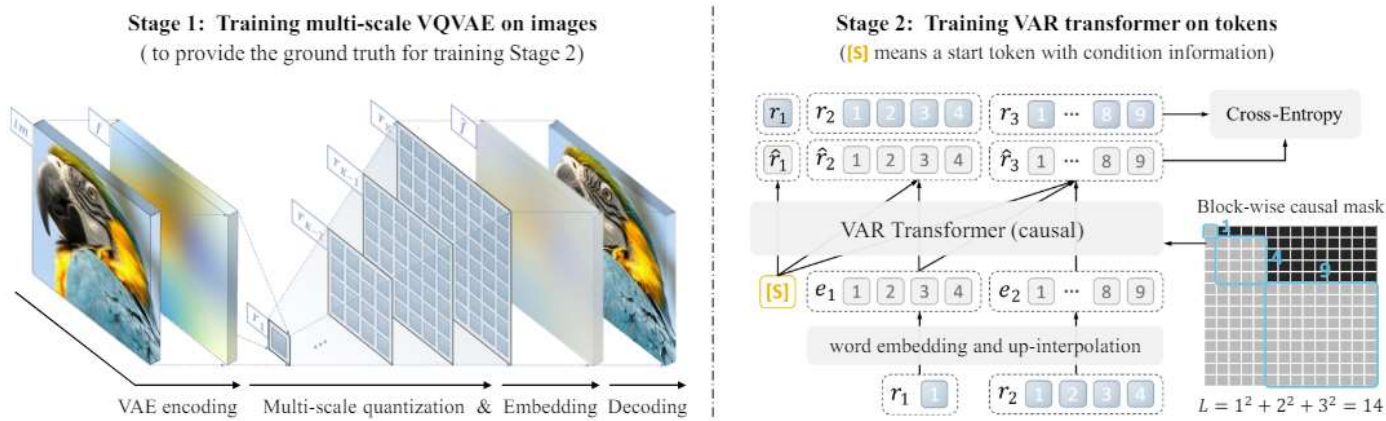


Figure 4: **VAR involves two separated training stages. Stage 1:** a multi-scale VQ autoencoder encodes an image into $K$ token maps $R = (r_1, r_2, \ldots, r_K)$ and is trained by a compound loss (5). For details on "Multi-scale quantization" and "Embedding", check Algorithm 1 and 2. **Stage 2:** a VAR transformer is trained via next-scale prediction (6): it takes $([s], r_1, r_2, \ldots, r_{K-1})$ as input to predict $(r_1, r_2, r_3, \ldots, r_K)$. The attention mask is used in training to ensure each $r_k$ can only attend to $r_{\leq k}$. Standard cross-entropy loss is used.

The autoregressive unit should be an entire token map, rather than a single token. Quantize a feature map $f \in \mathbb{R}^{h \times w \times C}$ into $K$ multi-scale token maps $(r_1, r_2, \ldots, r_K)$, each at a increasingly higher resolution $h_k \times w_k$, culminating in $r_K$ matches the original feature map's resolution $h \times w$.

$$p(r_1, r_2, \ldots, r_K) = \prod_{k=1}^{K} p(r_k \mid r_1, r_2, \ldots, r_{k-1}),$$

Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction (Jun 2024)

## Does Image Modeling Have to Be Discrete?

**Motivation:** model the per-token probability distribution using a diffusion procedure, which allows us to apply autoregressive models in a continuous-valued space using *Diffusion Loss*.
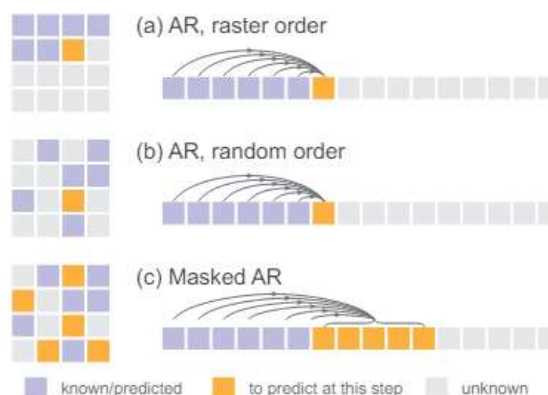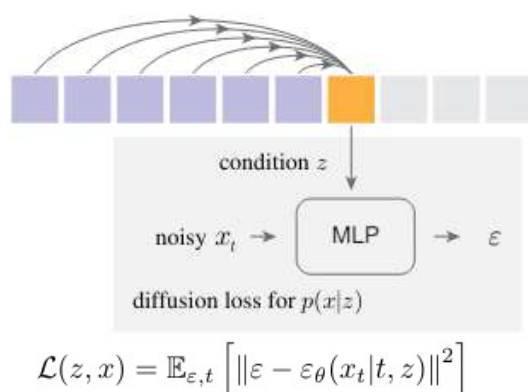


Figure 3: **Generalized Autoregressive Models.** (a) A standard, raster-order autoregressive model predicts one next token based on the previous tokens. (b) A random-order autoregressive model predicts the next token given a random order. It behaves like randomly masking out tokens and then predicting one. (c) A Masked Autoregressive (MAR) model predicts multiple tokens simultaneously given a random order, which is conceptually analogous to masked generative models [4, 29]. In all cases, the prediction of one step can be done by causal or bidirectional attention (Figure 2).

$$\mathcal{L}(z, x) = \mathbb{E}_{\varepsilon, t}\left[\|\varepsilon - \varepsilon_\theta(x_t|t, z)\|^2\right]$$

Instead of using diffusion models for representing the joint distribution of all pixels or all tokens, in our case, the diffusion model is for representing the distribution for **each token**. Similar to MAE, predict multiple tokens based on previous tokens.

$$p(x^1, ..., x^n) = p(X^1, ..., X^K) = \prod_{k}^{K} p(X^k \mid X^1, ..., X^{k-1}).$$
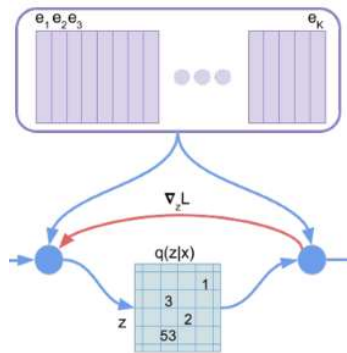
$X^k = \{x^i, x^{i+1}..., x^j\}$ is a set of tokens to be predicted at the k-th step

**MAE + Diffusion**
**Use diffusion to replace MAE decoder**

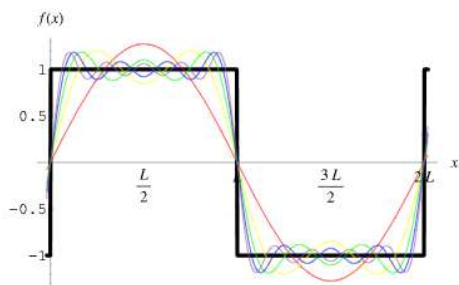Autoregressive Image Generation without Vector Quantization (Jun 2024)
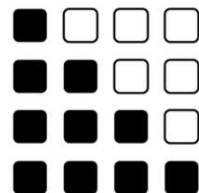
# End to End

## How About Frequency Domain?



Hard to ensure the independence among vectors (bases) in codebook
Not an efficient information encoding technique, leading to detail loss in image
Still unclear how to define causality in image space

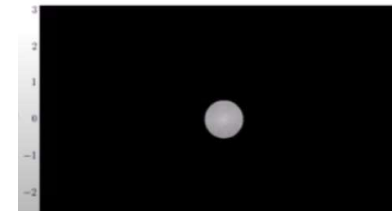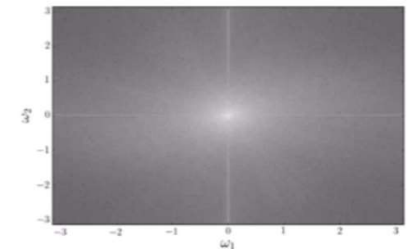Use FFT to transform images into frequency domain.



- Sinusoidal bases are linear independent to each other
- In 2D image space, there is a **causal** relationship between low frequency basis (blur) and high frequency (clear) basis.

$$f_q = e_1 f_1 + e_2 f_2 + e_3 f_3 + \cdots + e_n f_n$$

**Codebook only consists of coefficients**

$$F(k, l) = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} f(i, j) \, e^{-\iota 2\pi(\frac{ki}{N} + \frac{lj}{N})}$$
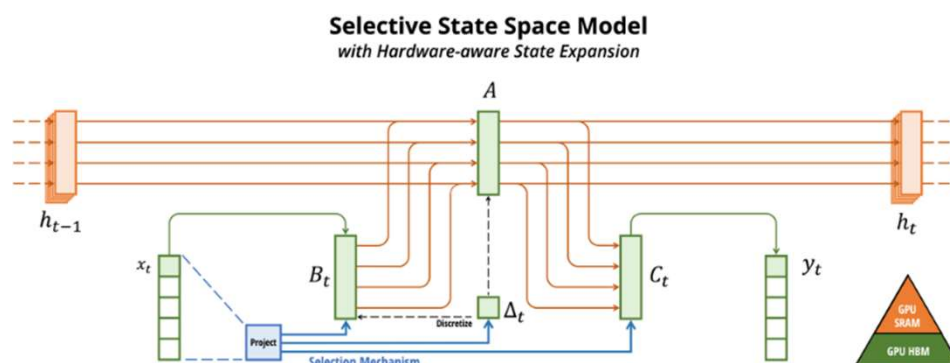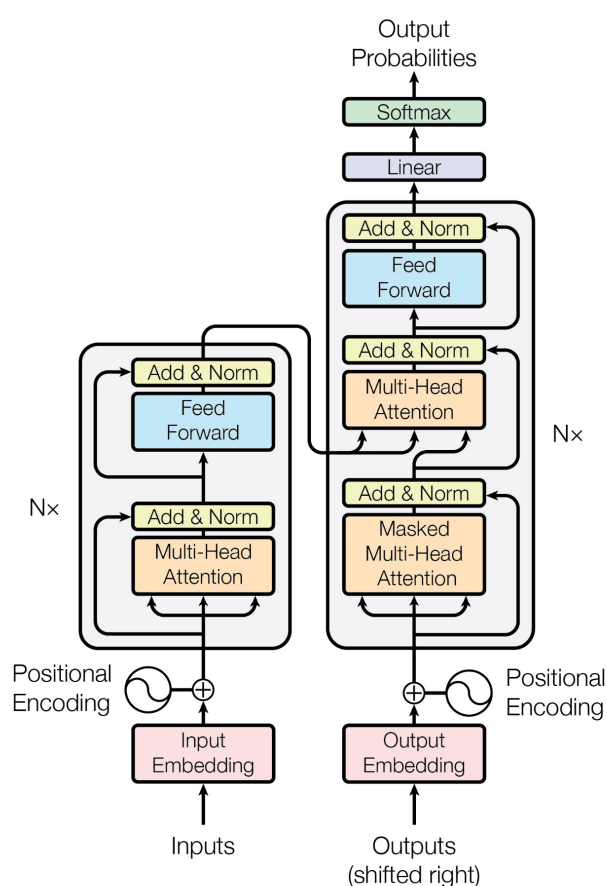
# Alternative

## How Long Will Transformer Dominate?

➢ Simple but efficient. The basic operation is matrix multiplication, making parallel computing possible. (ZeRO, Megatron…)

➢ Suitable for almost any modal data.

$O(L^2)$ inference cost. Struggle to achieve long-context window understanding. (KV-Cache, GQA, RoPE, NTK-aware interpolation … )
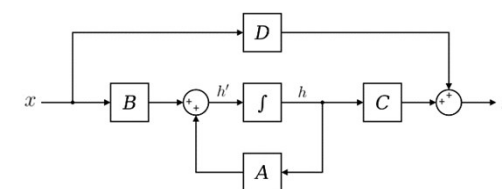
Any challenger?



**Selective State Space Model**
*with Hardware-aware State Expansion*

Technical structure of Mamba

$$h'(t) = Ah(t) + Bx(t) \quad (1a)$$
$$y(t) = Ch(t) \quad (1b)$$

$$h_t = \overline{A}h_{t-1} + \overline{B}x_t \quad (2a)$$
$$y_t = Ch_t \quad (2b)$$

$$\overline{K} = (C\overline{B}, C\overline{AB}, \ldots, C\overline{A}^k\overline{B}, \ldots) \quad (3a)$$
$$y = x * \overline{K} \quad (3b)$$

Mamba: Linear-Time Sequence Modeling with Selective State Spaces (May 2024)

Linear-Time Sequence Modeling

## Next-Token Prediction: The Path to AGI?

In computer vision, there has been a similar pattern. Early methods conceived of vision as searching for edges, or generalized cylinders, or in terms of SIFT features. But today all this is discarded. Modern deep-learning neural networks use only the notions of **convolution** and certain kinds of invariances, and perform much better.

➢ AI researchers have often tried to build knowledge into their agents
➢ this always helps in the short term, and is personally satisfying to the researcher,
➢ in the long run it plateaus and even inhibits further progress,
➢ breakthrough progress eventually arrives by an opposing approach based on scaling computation by search and learning.

We should build in only the meta-methods that can find and capture this arbitrary complexity. Essential to these methods is that they can find good approximations, but the search for them should be by our methods, not by us. We want AI agents that can discover like we can, not which contain what we have discovered.

# Thanks