

Multi Modal LLMs

By Madhava Prasath

**Additional Reads on @himanshu's Blog where he
covered VAE really well**



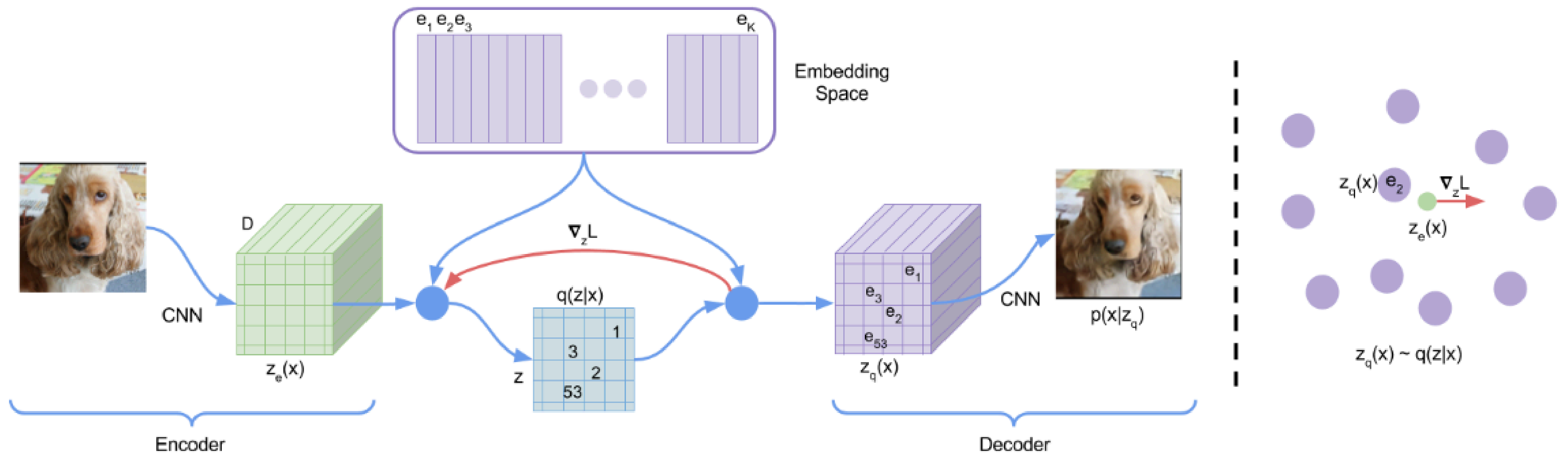
[Link here](#)

Table of Contents in Today's lecture

Prerequisites :

- **VQ-VAE**
- **Decoder mask on LLMs**
- **Video Poet**
- **Flamingo**
- **Llava 3D**

$$L = \log p(x|z_q(x)) + \|\text{sg}[z_e(x)] - e\|_2^2 + \beta \|z_e(x) - \text{sg}[e]\|_2^2,$$



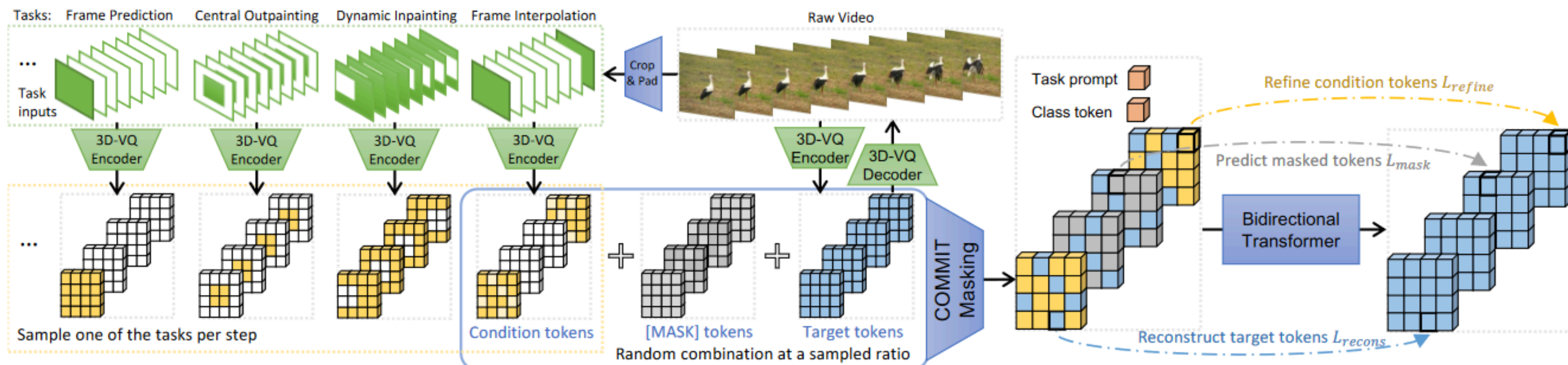
$K = 512$

Why VQ-GAN Matters?

- **Stable Diffusion uses it**
- **A whole new paradigm - Latent Diffusion Models**

Decoder mask on LLMs

Mag-ViT



At a sampled mark ratio, we randomly replace target tokens \mathbf{z}_i , with either 1) the condition token $\tilde{\mathbf{z}}_i$, if the corresponding supervoxel of \mathbf{z}_i contains condition pixels; or 2) the special [MASK] token, otherwise. Formally, we compute the *multivariate* conditional mask $\mathbf{m}(\cdot \mid \tilde{\mathbf{z}})$ as

$$\mathbf{m}(\mathbf{z}_i \mid \tilde{\mathbf{z}}_i) = \begin{cases} \tilde{\mathbf{z}}_i & \text{if } s_i \leq s^* \wedge \neg \text{ispad}(\tilde{\mathbf{z}}_i) \\ [\text{MASK}] & \text{if } s_i \leq s^* \wedge \text{ispad}(\tilde{\mathbf{z}}_i) \\ \mathbf{z}_i & \text{if } s_i > s^* \end{cases} \quad (2)$$

$$\mathcal{L}(\mathbf{V}; \theta) = \mathbb{E}_{\rho, \hat{\mathbf{V}}} \mathbb{E}_{\mathbf{m} \sim p_{\mathcal{M}}} \left[\sum_i -\log p_{\theta}(z_i \mid [\rho, \mathbf{c}, \bar{\mathbf{z}}]) \right] \quad (3)$$

We can decompose the loss in Eq. (3) into three parts according to Eq. (2): $\mathcal{L}_{\text{refine}}$ refines the task-specific condition tokens, $\mathcal{L}_{\text{mask}}$ predicts masked tokens, and $\mathcal{L}_{\text{recons}}$ reconstructs target tokens. Let $\bar{\mathbf{c}} = [\rho, \mathbf{c}, \bar{\mathbf{z}}]$ for simplicity,

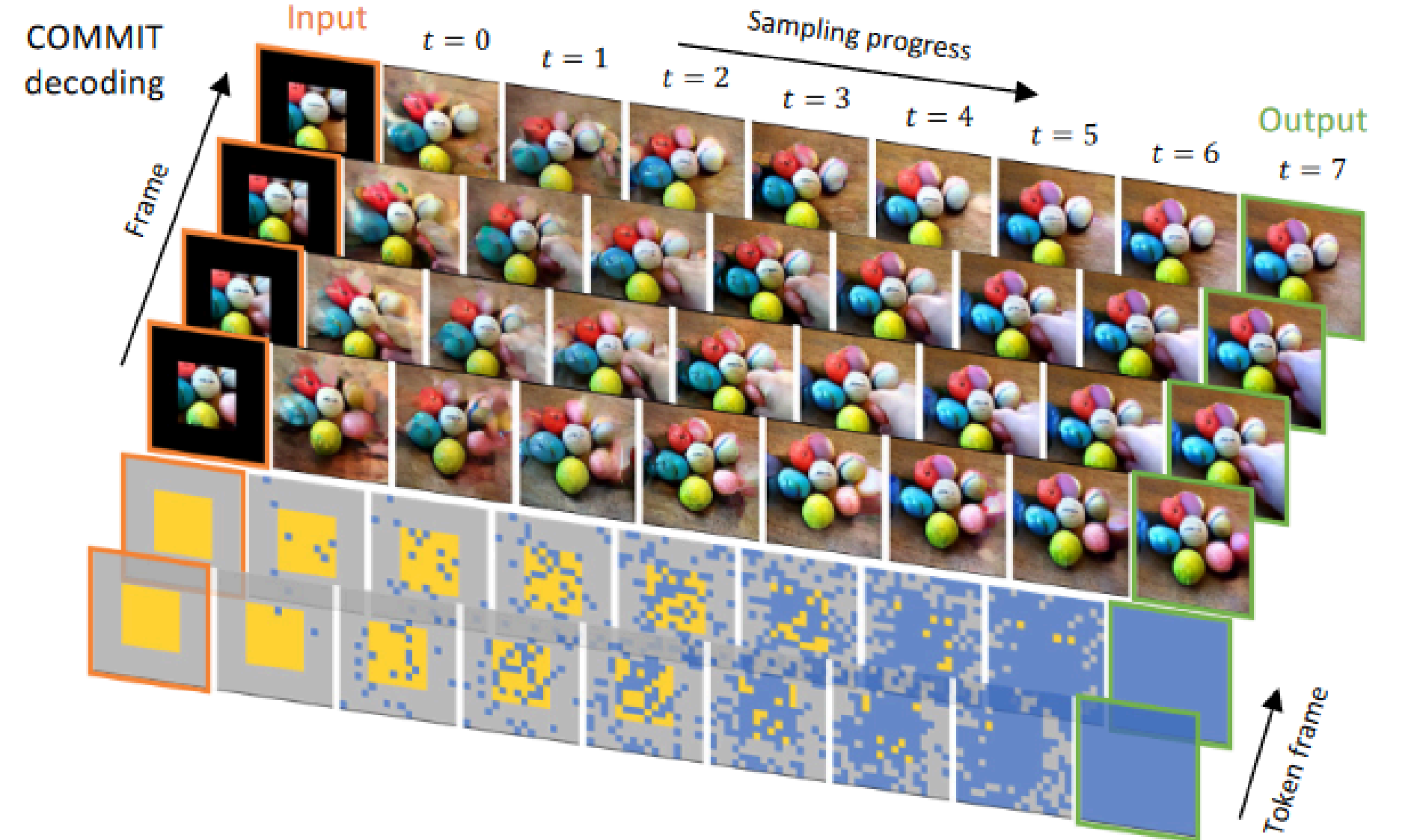
$$\begin{aligned} \sum_{i=1}^N -\log p_{\theta}(z_i \mid [\rho, \mathbf{c}, \bar{\mathbf{z}}]) &= \underbrace{\sum_{\bar{\mathbf{z}}_i = \tilde{\mathbf{z}}_i} -\log p_{\theta}(z_i \mid \bar{\mathbf{c}})}_{\text{Refine condition tokens } \mathcal{L}_{\text{refine}}} \\ &+ \underbrace{\sum_{\bar{\mathbf{z}}_i = [\text{MASK}]} -\log p_{\theta}(z_i \mid \bar{\mathbf{c}})}_{\text{Predict masked tokens } \mathcal{L}_{\text{mask}}} + \underbrace{\sum_{\bar{\mathbf{z}}_i = \mathbf{z}_i} -\log p_{\theta}(z_i \mid \bar{\mathbf{c}})}_{\text{Reconstruct target tokens } \mathcal{L}_{\text{recons}}} \end{aligned} \quad (4)$$

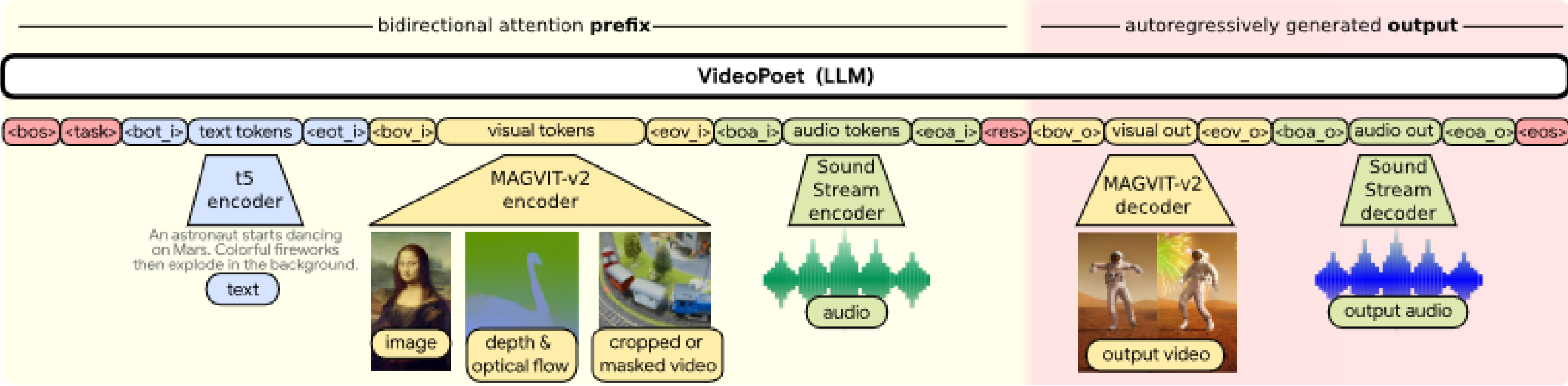
Algorithm 1 Non-autoregressive Decoding by COMMIT

Input: prefix ρ and \mathbf{c} , condition $\tilde{\mathbf{z}}$, steps K , temperature T

Output: predicted visual tokens $\hat{\mathbf{z}}$

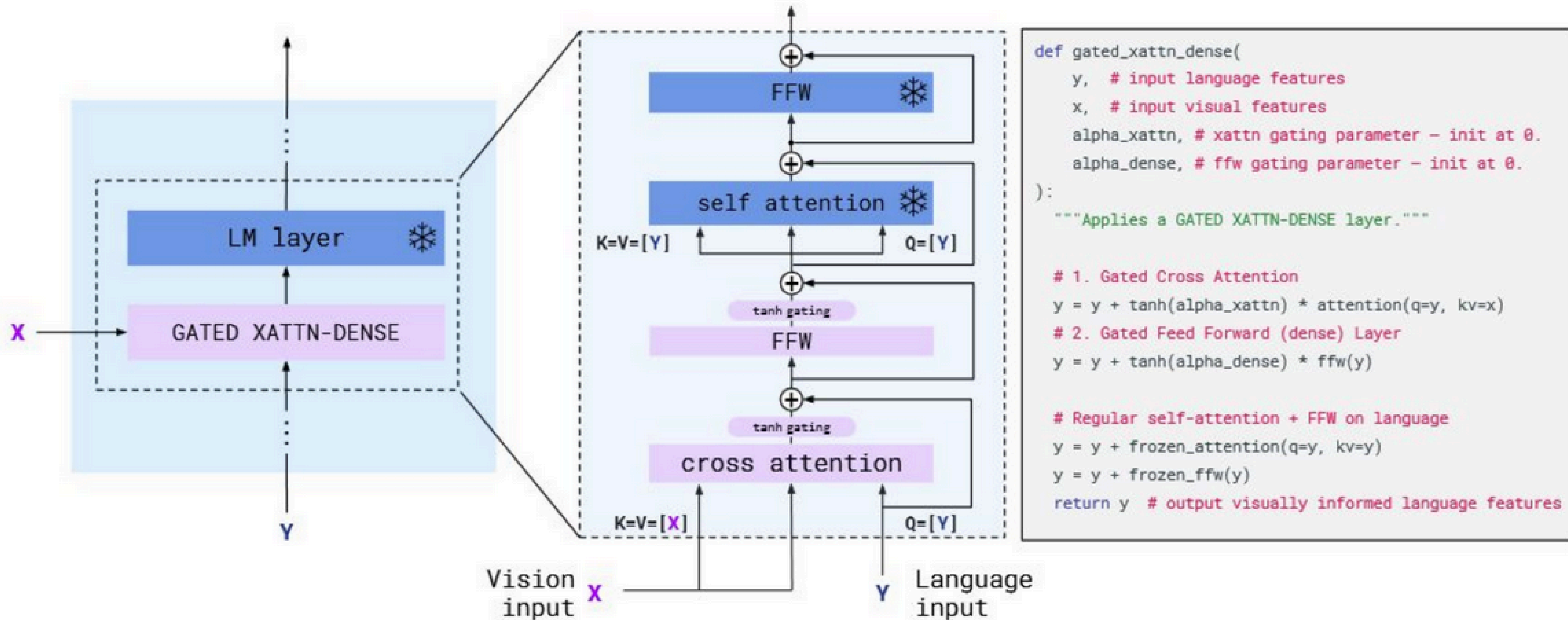
- 1: $\mathbf{s} = \mathbf{0}, s^* = 1, \hat{\mathbf{z}} = \mathbf{0}^N$
- 2: **for** $t \leftarrow 0, 1, \dots, K - 1$ **do**
- 3: $\bar{\mathbf{z}} \leftarrow \mathbf{m}(\hat{\mathbf{z}} \mid \tilde{\mathbf{z}}; \mathbf{s}, s^*)$
- 4: $\hat{z}_i \sim p_{\theta}(z_i \mid [\rho, \mathbf{c}, \bar{\mathbf{z}}]), \forall i$ where $s_i \leq s^*$
- 5: $s_i \leftarrow p_{\theta}(\hat{z}_i \mid [\rho, \mathbf{c}, \bar{\mathbf{z}}]), \forall i$ where $s_i \leq s^*$
- 6: $s_i \leftarrow s_i + T(1 - \frac{t+1}{K}) \text{Gumbel}(0, 1), \forall i$ where $s_i < 1$
- 7: $s^* \leftarrow$ The $\lceil \gamma(\frac{t+1}{K})N \rceil$ -th smallest value of \mathbf{s}
- 8: $s_i \leftarrow 1, \forall i$ where $s_i > s^*$
- 9: **end for**



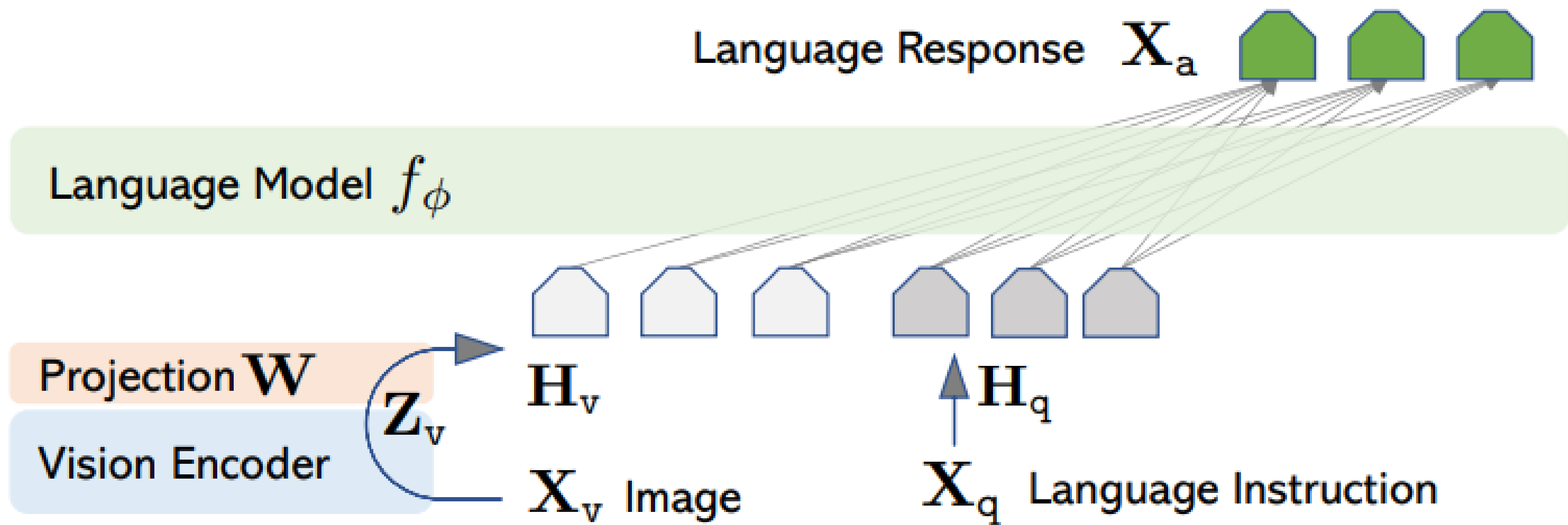


**Lets say you are not rich as Google and You have Vision Encoder and LLM
you have to align it**

Flamingo

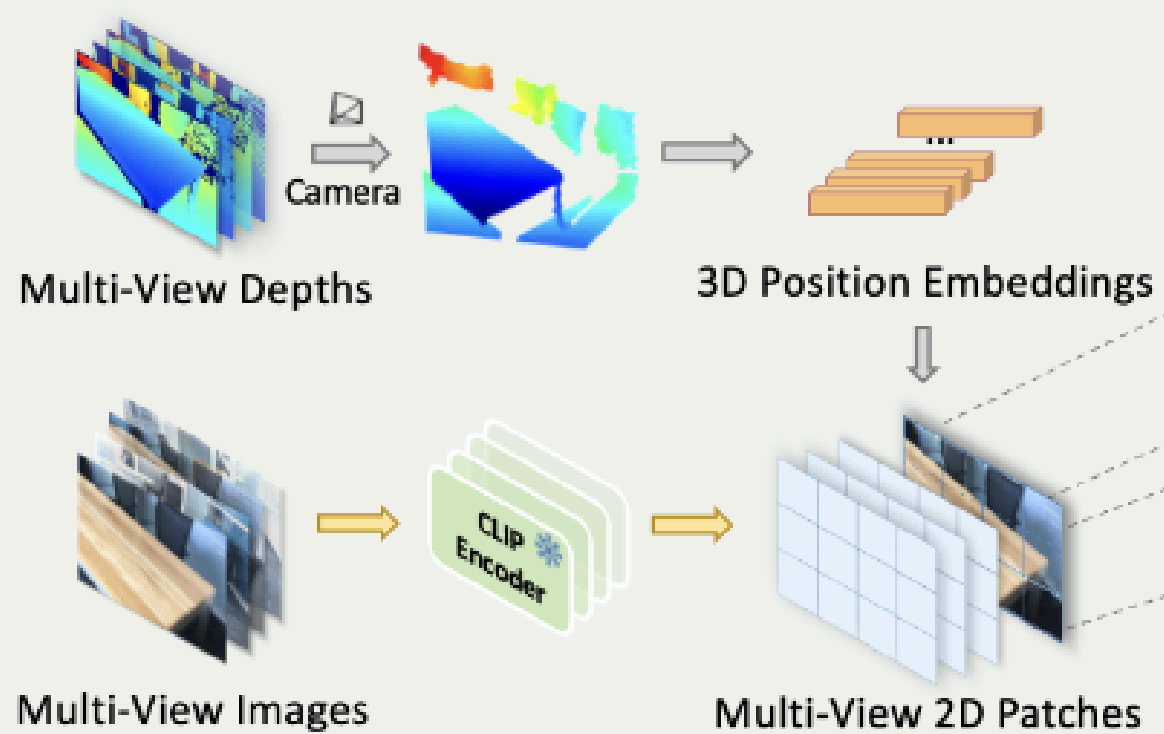


Llava Generally

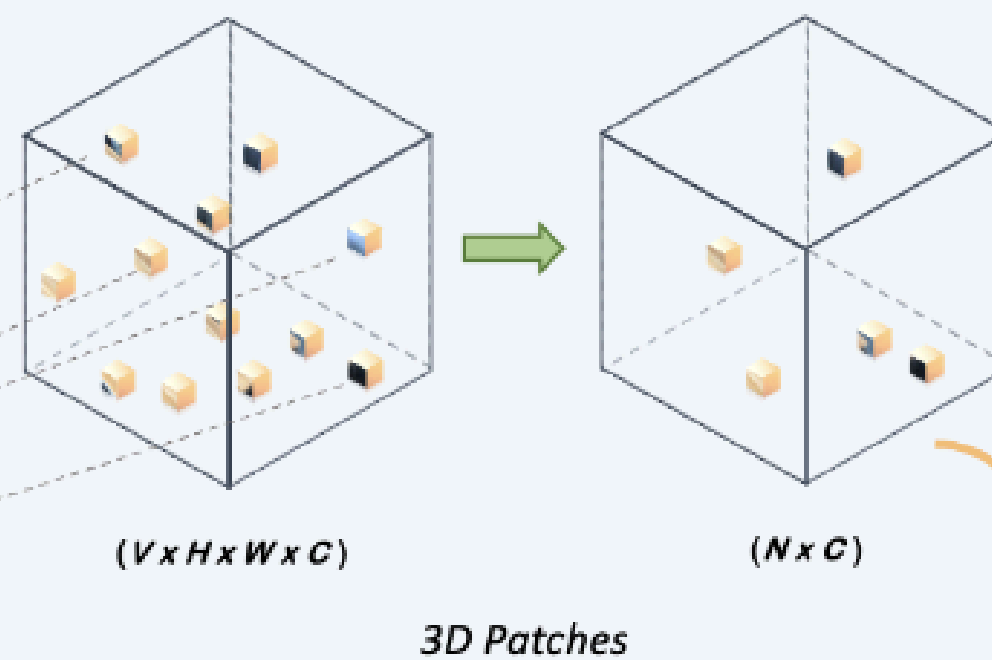


Llava - 3D

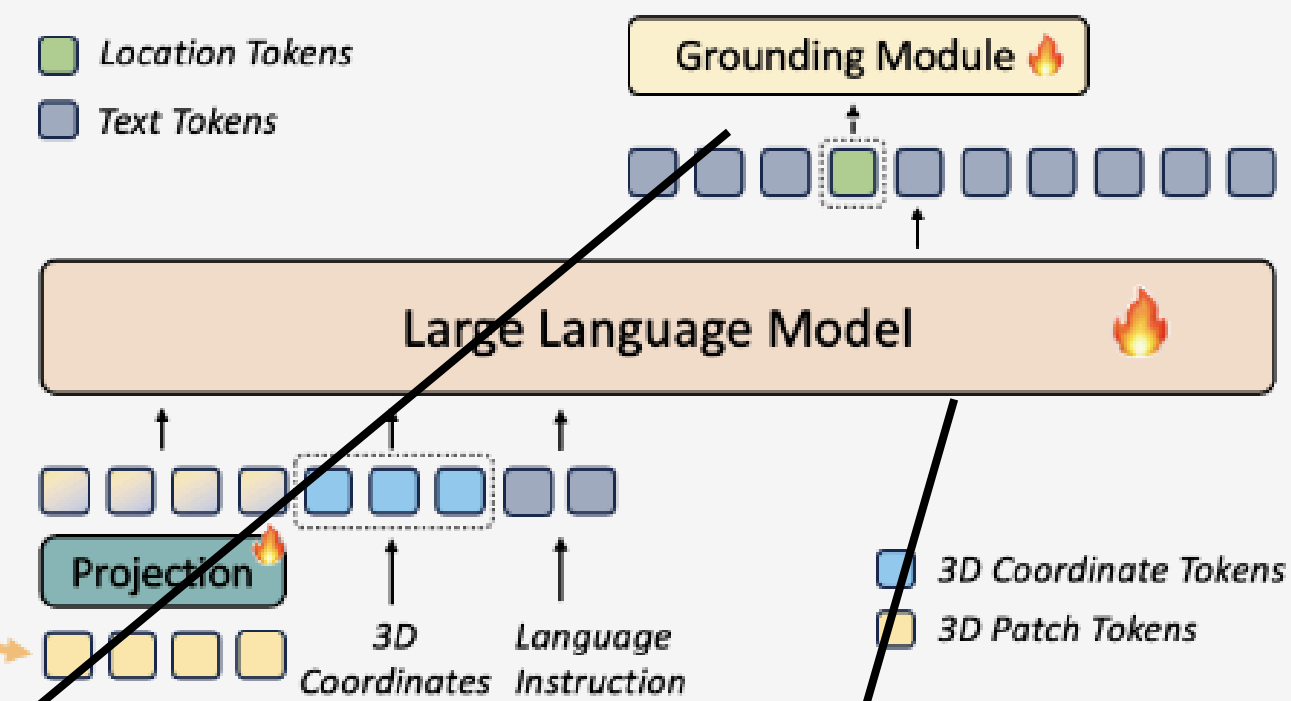
3D Patch Construction



3D Patch Pooling



3D-aware Position Encoding & Decoding



1

2

3

3.2. 3D Patch

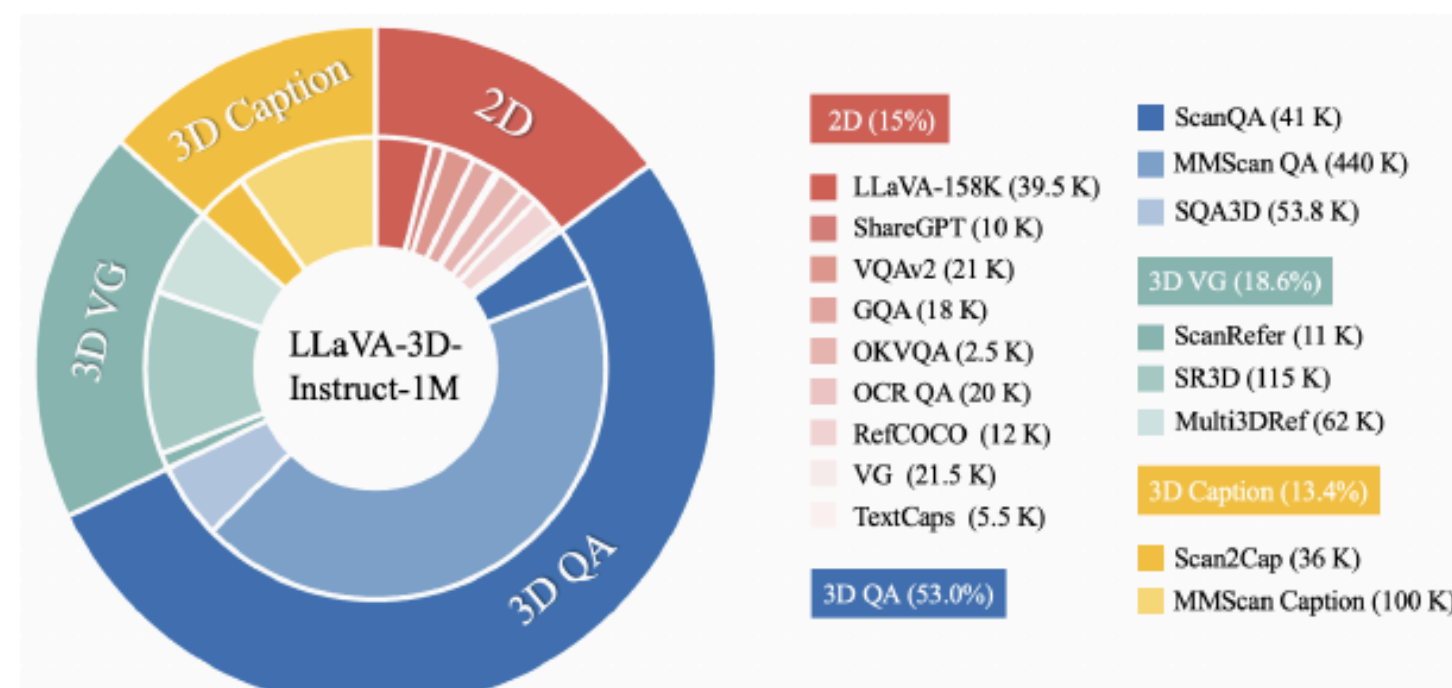
Our 3D Patch representations are built upon the 2D patch features X'_v extracted from multi-view images with CLIP visual encoder to leverage the strong visual-semantic alignment. To construct the 3D Patches, we inject the 3D position information into the aforementioned 2d patches so that the 3D Patches can explicitly model 3D spatial information while preserving the semantic information from 2D patches. As illustrated in left block of Fig. 2, given the multi-view 2D patch features $X'_v \in \mathbb{R}^{V \times c \times w \times h}$, we obtain their 3D positions $P \in \mathbb{R}^{V \times 3 \times w \times h}$ in the 3D world, using nearest neighbor depth and known camera intrinsic and extrinsic parameters. The 3D positions P are then encoded into 3D position embeddings $P' \in \mathbb{R}^{V \times w \times h \times d}$ through a learnable two-layer MLP, which are subsequently added to the 2D patch visual tokens, resulting in the 3D patches $X'_{3D} \in \mathbb{R}^{V \times w \times h \times d}$.

$$X'_{3D} = X'_v + \text{MLP}(P') \quad (1)$$

Stage 1: 3D Patch Language Alignment. During the first training stage, we use the region-level and scene-level caption data that describe spatial relationships among 3D objects to align the 3D patches with the LLM for enhanced 3D spatial comprehension. At this stage, the input multi-view images used are selected from sequences that correspond

to specific regions or entire scenes. We freeze the vision encoder and LLM parameters, and only train the projection layer and 3D position embedding layer, encouraging efficient alignment between 3D patch features and text space. Since 3D patches are derived from CLIP features augmented with 3D positional information, the alignment between 3D Patch and LLM converges rapidly.

Stage 2: Task Instruction Tuning. During the instruction-tuning stage, LLaVA-3D is optimized to respond to complex 3D V&L tasks and maintain its inherent 2D image reasoning and instruction-following capabilities. To facilitate this capability, we collect the **LLaVA-3D-Instruct-1M** dataset, a hybrid collection of 2D and 3D data specifically tailored for instruction tuning. The overall distribution of the dataset collection is shown in Fig 3, more details about the instructional tuning datasets are listed in the appendix. The 2D Data of LLaVA-3D-Instruct-1M is derived from existing LLaVA-1.5 instruction tuning data, ensuring the preservation of 2D image comprehension and vision-language conversation abilities. When tuning with 2D data, we keep the vision encoder frozen and jointly train the projection layer and LLM. The 3D Data of LLaVA-3D-Instruct-1M, on the other hand, comprises data from diverse 3D QA, 3D dense captioning, and 3D grounding tasks. During the 3D data instruction tuning, the 3D position embedding layer will be added to jointly train with the other modules. Additionally, for tasks where the instruction contains 3D coordinate information or requires 3D bounding box outputs, the corresponding encoding and decoding modules will be trained together. During instruc-



tion tuning, the 3D data pathway includes the 3D position embeddings and 3D patches, while the 2D data pathway is the original LLaVA. All modules except for the 3D position embeddings to construct 3D patches are shared across 2D and 3D data. This training setup ensures that LLaVA-3D is capable of processing both 2D and 3D visual tokens effectively, and is adaptive to various task formulations and

5.1. Implementation Details

LLaVA-3D is built upon the LLaVA-1.5-7B, utilizing their pre-trained weights from the HuggingFace library. For 3D tasks, we add the 3D position embeddings to the 2D patch visual tokens, and utilize the voxelization pooling strategy to reduce 3D patch number before passing the input visual tokens to the projection layer and LLM. The number of views V is set to be 20 and voxel size is set to 0.2m. Due to the LLM context length limitation, the maximum number of 3D patch tokens after 3D-aware pooling is set to 3096. For 2D tasks, LLaVA-3D functions the same as LLaVA. All experiments are conducted on $8 \times 80\text{G}$ A100 GPUs. We train our model for 1 epoch with a learning rate of $1\text{e-}3$ and a batch size of 32 in stage 1, and fine-tune on the collected LLaVA-3D-Instruct-1M dataset, with a learning rate of $2\text{e-}5$ and a batch size of 16 in stage 2.

Cheap self promotion - :)

A Comprehensive Survey of Mamba Architectures for Medical Image Analysis: Classification, Segmentation, Restoration and Beyond

SHUBHI BANSAL*, Indian Institute of Technology Indore, India

SREEHARISH A*, R.M.D. Engineering College, Kavaraipettai, India

MADHAVA PRASATH J*, R.M.D. Engineering College, Kavaraipettai, India

MANIKANDAN S*, R.M.D. Engineering College, Kavaraipettai, India

SREEKANTH MADISETTY*, Jio Platforms Limited, India

MOHAMMAD ZIA UR REHMAN*, Indian Institute of Technology Indore, India

CHANDRAVARDHAN SINGH RAGHAW*, Indian Institute of Technology Indore, India

GAURAV DUGGAL, Birla Institute of Technology & Science Pilani, India

¹ NAGENDRA KUMAR†, Indian Institute of Technology Indore, India

[Paper Link](#)

[Github link here](#)

Awesome Mamba Papers on Medical Domain

This repository includes all the referenced papers, along with the sections covered in the survey titled **A Comprehensive Survey of Mamba Architectures for Medical Image Analysis: Classification, Segmentation, Restoration, and Beyond**. If you encounter any issues with the links or notice any errors, please feel free to reach out via email at 21204028@rmd.ac.in

While some papers may be repeated, our objective is to highlight the papers available in each specific sections. For the best experience, we recommend keeping this repository open alongside the topics you are currently reading. If you wish to explore further, this resource can serve as a useful tool.

Any doubts ?

Feel free to ping me on discord @maddy

or

Linkedin - Madhava Prasath