



# MoLe-VLA: Dynamic Layer-skipping Vision Language Action Model via Mixture-of-Layers for Efficient Robot Manipulation

Rongyu Zhang<sup>1,2,3,4\*</sup>, Menghang Dong<sup>3\*</sup>, Yuan Zhang<sup>3</sup>, Liang Heng<sup>3</sup>, Xiaowei Chi<sup>5</sup>, Gaole Dai<sup>3,4</sup>, Li Du<sup>1</sup>, Yuan Du<sup>1</sup> Shanghang Zhang<sup>3</sup>

<sup>1</sup>Nanjing University; <sup>2</sup>The Hong Kong Polytechnic University;

<sup>3</sup>State Key Laboratory of Multimedia Information Processing, School of Computer Science, Peking University; <sup>4</sup>Beijing Academy of Artificial Intelligence;

<sup>5</sup>The Hong Kong University of Science and Technology

\* Equal contribution, Corresponding author, Project web page: [MoLe-VLA-Web](#)

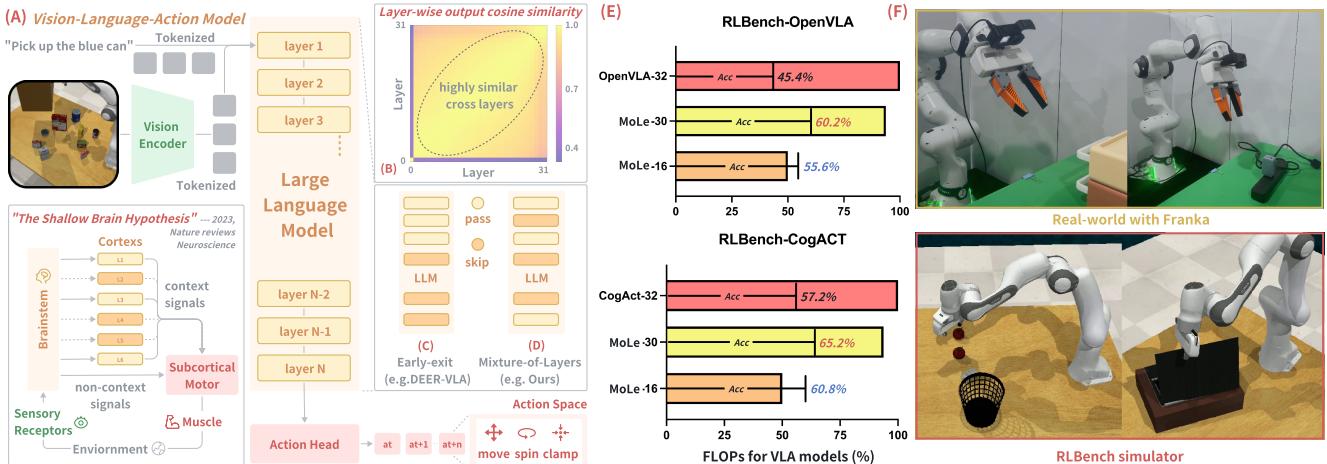


Figure 1. Overview of our proposed **MoLe-VLA**: Our proposed framework integrates dynamic layer activation, a novel Spatial-Temporal Aware Router (STAR), and self-knowledge distillation (CogKD) to achieve efficient and adaptive performance in robotic applications. MoLe reduces computational costs while enhancing model performance, enabling resource-constrained platforms to benefit from MLLMs.

## Abstract

Multimodal Large Language Models (MLLMs) excel in understanding complex language and visual data, enabling generalist robotic systems to interpret instructions and perform embodied tasks. Nevertheless, their real-world deployment is hindered by substantial computational and storage demands. Recent insights into the homogeneous patterns in the LLM layer have inspired sparsification techniques to address these challenges, such as early exit and token pruning. However, these methods often neglect the critical role of the final layers that encode the semantic information most relevant to downstream robotic tasks. Aligning with the recent breakthrough of the Shallow Brain Hypothesis (SBH) in neuroscience and the mixture of experts in model sparsification, we conceptualize each LLM layer as an expert and pro-

pose a **Mixture-of-Layers Vision-Language-Action model (MoLe-VLA or simply MoLe)** architecture for dynamic LLM layer activation. We introduce a Spatial-Temporal Aware Router (STAR) for MoLe to selectively activate only parts of the layers based on the robot's current state, mimicking the brain's distinct signal pathways specialized for cognition and causal reasoning. Additionally, to compensate for the cognition ability of LLM lost in MoLe, we devise a cognition self-knowledge distillation (CogKD) to enhance the understanding of task demands and generate task-relevant action sequences by leveraging cognition features. Extensive experiments in both RL Bench simulation and real-world environments demonstrate the superiority of MoLe-VLA in both efficiency and performance, achieving performance improvement of 8% mean success rate across ten tasks while reducing at most  $\times 5.6$  computational costs in LLM.

## 1. Introduction

The rapid advancements in multimodal large language models (MLLMs) [2, 3, 17, 27, 32] have demonstrated their ability to integrate complex language and visual representations, inspiring the development of generalist robots and embodied agents capable of vision-language comprehension, human interaction, and flexible problem-solving in manipulation tasks. Preliminary vision language action (VLA) models [16, 19, 21, 25], such as RT-2 [6] and OpenVLA [16], have shown the feasibility of using MLLMs for end-to-end robotic control, enabling robust policies and emergent abilities, including generalization to unseen objects and understanding novel commands. However, deploying MLLMs in real-world robotic systems faces significant challenges due to their high computational demands, including substantial memory usage, power consumption, and time delays, which conflict with robotic platform resource-constrained and real-time requirements. For example, a 7B VLA model running on a commercial-grade GPU like the RTX 4090 generally achieves an inference frequency of approximately 5 – 12 Hz, which falls significantly short of the 50 – 1000 Hz control frequency required by the Franka robotic arm.

Recent studies [34, 43] have uncovered significant redundancy in LLM layer, particularly in robotic tasks, where homogeneous patterns across layers lead to high computational costs with limited performance gains. For instance, DeeR [43] demonstrated that using all 24 layers of the Flamingo [21] model improves task success rates by only 3.2% compared to using six layers, while computational costs increase 4x on the Calvin LH-MTLC [28]. Similarly, our analysis of OpenVLA [16] with RLBench [13] in Fig. 1 (A) reveals that cosine similarity between consecutive layer outputs exceeds 90%, while features from the first and last layers differ significantly. This suggests the potential for skipping adjacent layers to reduce computation but also highlights the limitations of early-exit strategies [10, 43], as shown in Fig. 1 (B), where discarding deeper layers risks losing critical semantic information. Inspired by the Shallow Brain Hypothesis (SBH) [37], which suggests that the brain balances deep hierarchical structures with shallow, parallel cortico-subcortical loops for cognition and causal reasoning, we propose a selective layer activation strategy in VLA models. As shown in Fig. 1 (C), our approach mirrors the brain’s dynamic depth-parallelism balance, activating only task-relevant layers to enhance efficiency and adaptability, embodying principles of SBH in VLA model design.

In this paper, we introduce a Mixture-of-Layers Vision-Language-Action model (**MoLe-VLA**) incorporating a novel layer-selection router at the input stage of LLMs for its sparsity. Our design emulates the brain’s decision-making process described in the SBH by dynamically selecting optimal forward pathways with varying layer combinations. Inspired by the routing mechanism in mixture-of-experts

(MoE) [23, 46, 47], which enables horizontal expert-wise activation within a single LLM layer, we extend this concept vertically to achieve layer-wise activation. Specifically, we treat each LLM layer as an independent expert and utilize a biologically inspired router to manage layer skipping, mimicking the brain’s selective activation of cortico-subcortical loops. Unlike Mixture-of-Depth (MoD) [34], which assigns input tokens to different experts and risks token-wise inconsistencies due to varying perception levels across layers, our proposed MoLe dynamically selects the most relevant layers while processing input features holistically.

Traditional MoE or MoD routers, which rely on simple linear layers, often fail to capture critical spatial-temporal information necessary for reasoning in dynamic, embodied intelligence tasks. To address this limitation, we propose the *Spatial-Temporal Aware Router (STAR)*, which independently processes spatial features from visual inputs and temporal dependencies from textual inputs. By combining these essential properties into a unified representation, STAR aligns the selection of LLM layers with the demands of the current environment. STAR dynamically activates the most relevant layers by generating softmax probabilities for each layer and selecting the top- $k$  layers with the highest probabilities. By fully leveraging spatial-temporal information, STAR ensures accurate and efficient adaptation to the dynamic nature of embodied intelligence tasks, achieving optimal performance with reduced computational overhead.

Nonetheless, skipping certain layers inevitably reduces the cognitive expressiveness of the model. To address this, we propose *Cognitive self-Knowledge Distillation (CogKD)*, a novel approach to preserve grasping ability while mitigating cognitive collapse. In CogKD, the original full-layer model serves as the teacher, while the MoLe layer-skipping model acts as the student. Inspired by [19], we introduce a learnable *cognition token*, which efficiently integrates visual tokens and language guidance to enhance comprehension of task demands and produce task-relevant action sequences. By analyzing the similarity between cognition tokens and student tokens, we identify tokens of interest (ToIs) that represent task-critical information the student needs to learn. These ToIs provide precise guidance for adaptively re-weighting the distillation process, ensuring the student model focuses on key cognitive features while consistently benefiting from the layer-skipping efficiency.

The effectiveness of MoLe in both performance and efficiency enhancement is demonstrated in real-world and RL-Bench simulation environments based on various VLA models against state-of-the-art baselines. Extensive robotic experiments show that MoLe reduces the computational costs by  $\times 5.6$  while improving model performance by up to 8%. The key contributions of this work are summarized as:

- We draw inspiration from the Shallow Brain Hypothesis to develop a MoLe framework, which mimics the signal flow

in the human brain and enables dynamic layer activation via a router to improve model efficiency.

- We propose a novel layer-decision router, STAR, which fully leverages the spatial-temporal information from robotic inputs to make more accurate activation decisions.
- We introduce a self-knowledge distillation paradigm, CogKD, to recover cognitive information lost due to layer-skipping in sparse LLMs, enhancing overall performance.

## 2. Related works

### 2.1. Vision language action model

The remarkable success of LLMs [26, 29, 33, 39] and VLMs [1, 15, 18, 24, 31] has driven the rapid development of VLA models [5, 16, 19], which extend VLMs by incorporating action generation. VLA models aim to bridge the gap between perception and action, enabling machines to not only interpret and understand visual and textual inputs but also generate and execute actions based on that understanding [4, 22]. By integrating visual and linguistic information, these models produce more complex, context-aware outputs tailored to real-world environments, advancing their applicability in dynamic and embodied intelligence tasks.

### 2.2. Efficient multimodal large language models

With the advancement of VLA models, improving inference efficiency has become a critical area of research. Existing efforts can be categorized into three main strategies: efficient architectural design, model compression, and dynamic networks. Liu et al.[25] leverage the Mamba model[11] to enable efficient fine-tuning and inference, achieving pose prediction speeds 7× faster than existing robotic MLLMs in both simulation and real-world experiments. Wang et al.[40] utilizes a lightweight model with only 93M parameters while retaining 98.4% of its performance and delivering a 2.2× speedup. Yue et al.[43] propose a dynamic inference framework with multi-exit architectures, allowing early computation termination based on task-specific requirements. However, existing early-exit methods often overlook the significance of the final layers, which carry greater semantic relevance to downstream tasks. Building on dynamic networks, our work integrates knowledge distillation to achieve a layer-skipping mechanism, optimizing model performance while reducing redundant computations.

### 2.3. Sparse mixture-of-experts

While activation sparsity has been widely explored [20, 45], sparse MoE model architecture has shown significant advantages in LLMs. [35] demonstrated their ability to efficiently utilize vast numbers of parameters by activating only a small portion of the computation graph during inference. In the LLMs and VLMs era, MoE has become a widely adopted and effective architecture [9, 46, 47]. For example, [23]

achieves performance comparable to LLaVA-1.5-7B on various visual understanding benchmarks and even surpasses LLaVA-1.5-13B on the object hallucination benchmark, using only 3B sparsely activated parameters. Additionally, [34] employs a router to dynamically choose between computational paths, such as a standard block’s computation or a residual connection. While our model shares similarities with [34], we differ by employing a router to select all standard block computations, enabling a more comprehensive approach to layer activation.

## 3. Methods

### 3.1. Preliminary: Mixture-of-Experts

The MoE paradigm enhances model capacity while maintaining computational efficiency via conditional computation. For an input  $\mathbf{x} \in \mathbb{R}^d$ , a standard MoE layer is defined as:

$$\text{MoE}(\mathbf{x}) = \sum_{i=1}^{N_e} G_i(\mathbf{x}) \cdot E(\mathbf{x}), \quad (1)$$

where  $N_e$  is the number of experts,  $E : \mathbb{R}^d \rightarrow \mathbb{R}^d$  represents the  $i$ -th expert network, and  $G(\mathbf{x}) = \{G_1(\mathbf{x}), \dots, G_{N_e}(\mathbf{x})\}$  is the gating function satisfying  $\sum_{i=1}^{N_e} G_i(\mathbf{x}) = 1$ . The gating weights are computed as:

$$G(\mathbf{x}) = \text{Softmax}(\mathbf{W}_g \cdot \mathbf{x} + \mathbf{b}_g), \quad (2)$$

where  $\mathbf{W}_g \in \mathbb{R}^{N_e \times d}$  and  $\mathbf{b}_g \in \mathbb{R}^{N_e}$  are learnable parameters. To improve efficiency, sparse gating with top- $k$  selection is often applied. To address load imbalance, where too many inputs are routed to a few experts, a load balance loss  $\mathcal{L}_{lb}$  is introduced:

$$\mathcal{L}_{lb} = \frac{1}{N_e} \sum_{i=1}^{N_e} \left( \frac{\sum_{n=1}^N v_i(\mathbf{x}_n)}{\sum_{n=1}^N v_i(\mathbf{x}_n) + \epsilon} \right)^2, \quad (3)$$

where  $v_i(\mathbf{x}_n) = 1$  if the  $i$ -th expert is selected for input  $\mathbf{x}_n$  by the top- $k$  gating mechanism, and  $v_i(\mathbf{x}_n) = 0$  otherwise. This loss encourages balanced expert utilization and improves computational efficiency.

### 3.2. Mixture-of-Layers: MoLe-VLA

**Vision language action model.** Tasked with a language instruction  $\mathbf{l}$  with a length  $L$ , a robot receives an observation  $\mathbf{o}_t$  from sensors (e.g., RGB image from the camera) at timestep  $t$  to predict the action space of a gripper with 7 degrees of freedom (DoF) to execute:

$$\mathbf{a}_t^* = [\Delta x, \Delta y, \Delta z, \Delta \phi, \Delta \theta, \Delta \psi, g], \quad (4)$$

where  $\Delta x, \Delta y$ , and  $\Delta z$  are the relative translation offsets of the end effector,  $\Delta \phi, \Delta \theta, \Delta \psi$  denote the rotation changes, and  $g \in \{0, 1\}$  indicates the gripper open/close state.

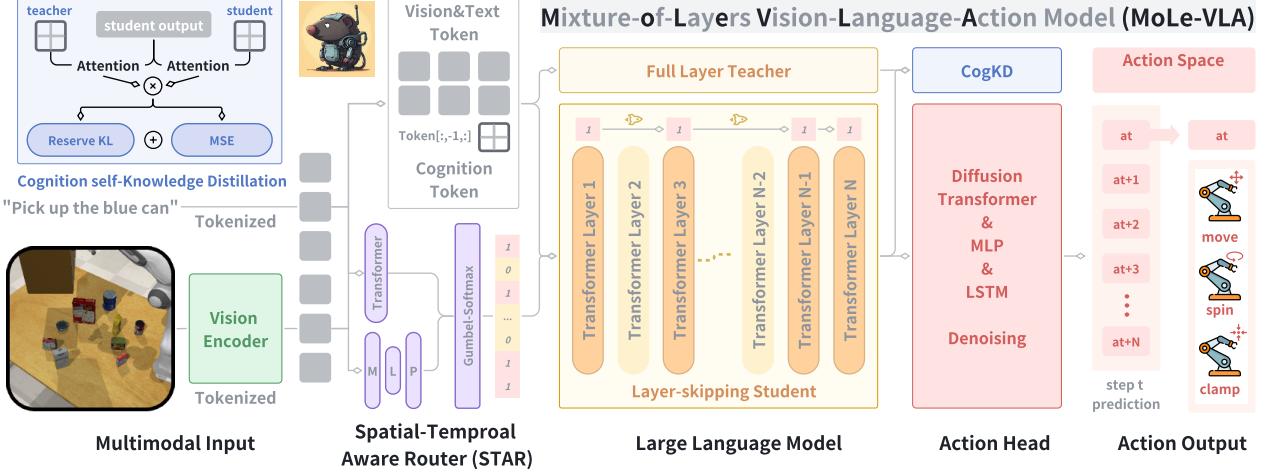


Figure 2. **The overall framework of MoLe-VLA.** Our proposed *Mixture of Layers (MoLe)* architecture consists of a *Spatial-Temporal Aware Router (STAR)* and a devised *Cognition self-Knowledge Distillation (CogKD)* for vision language action models.

Our basic VLA model mainly consists of a vision encoder  $\mathcal{E}$ , an MLLM  $\pi$ , and an action module  $\mathcal{A}$ . The vision encoder  $\mathcal{E}$  comprises DINO-v2 [30] and Siglip [44], which encodes an input image  $o_t$  into a sequence of informative tokens  $v_t$ . For multimodal fusion, an MLLM is established on top of the visual representations generated by the vision encoder  $\mathcal{E}$ , which functions as an effective multimodal feature extractor  $\pi$ , formalized as follows:

$$\mathbf{f}_t = \pi(\mathbf{l}, \mathcal{E}(o_t)), \quad (5)$$

where the output  $\mathbf{f}_t$  represents the hidden state sequence from the last layer of our MLLM at timestep  $t$ , corresponding to the cognition token. This serves as a condition for the subsequent action module to interpret and derive the desired actions. Following CogAct [19], our action module  $\mathcal{A}$  takes the cognition feature  $e_t^c$  extracted from the output feature  $\mathbf{f}_t$  as input and predicts the final actions  $a_t^*$ .

Our vision, language, and action modules are trained end-to-end by minimizing the mean squared error between the predicted noises from the action module and the ground truth noises. Taking the diffusion head as an example, the loss function is defined as:

$$\mathcal{L}_{task} = \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1), i} \|\hat{\epsilon}^i - \epsilon\| \quad (6)$$

where  $\hat{\epsilon}^i$  is the predicted noise for the noisy action  $a_t^*$  at the  $i$ 's denoising step, and  $\epsilon$  is the corresponding ground truth.

**Layer-skipping mechanism via MoLe router.** We propose MoLe-VLA to improve the efficiency of LLM in robotic tasks, where many transformer layers are underutilized due to the simpler reasoning demands of robotics tasks. MoLe employs a lightweight router to adaptively skip non-essential transformer layers during inference, reducing computational costs while maintaining performance.

As shown in Fig. 2, for a given MLLM  $\pi$  with  $K$  layers, the MoLe router processes the input embeddings  $\mathbf{x}_k \in \mathbb{R}^{b \times n \times d}$  and generates a binary gating vector  $G_{mol}(\mathbf{x}) = \{G_k\}_{k=1}^K$ , where  $G_k \in [0, 1]$ . To ensure efficiency, only the top- $k$  values in  $G_{mol}(\mathbf{x})$  are set to 1, determining which layers  $\pi_k$  are executed with the hidden feature  $\mathbf{h}_k$  while the rest are skipped:

$$\mathbf{h}_k = G_k \cdot \pi_k(\mathbf{h}_{k-1}) + (1 - G_k) \cdot \mathbf{h}_{k-1}. \quad (7)$$

Unlike traditional MoE routers that allocate tokens to experts, the MoLe router skips entire layers, avoiding redundant computations. This improves inference efficiency and responsiveness, making MoLe particularly suited for real-time robotic tasks like manipulation and navigation that require lightweight and adaptive processing. The complete pseudo-code of MoLe is provided in Algorithm 1.

### 3.3. Spatial-Temporal Aware Router

We propose a novel routing mechanism that synergistically leverages the spatial structure of visual inputs and the temporal dependencies in language inputs to select appropriate LLM layers for VLA tasks dynamically. Given visual features  $\mathbf{v}_t \in \mathbb{R}^{b \times n_{img} \times d}$  and textual features  $\mathbf{l} \in \mathbb{R}^{b \times n_{text} \times d}$ , both modalities are projected into a shared latent space using a learnable matrix  $\mathbf{W}_p \in \mathbb{R}^{d \times d_1}$ :

$$\mathbf{h}_{img} = \mathbf{v}_t \cdot \mathbf{W}_p, \quad \mathbf{h}_{text} = \mathbf{l} \cdot \mathbf{W}_p. \quad (8)$$

We compute spatial routing weights  $\mathbf{S} \in \mathbb{R}^{b \times N_e}$  from  $\mathbf{h}_{img}$  to capture spatial features:

$$\mathbf{S} = \mathbf{W}_s^{(2)} \cdot \varphi(\mathbf{W}_s^{(1)} \cdot \mathbf{h}_{img} + \mathbf{b}_s^{(1)}), \quad (9)$$

where  $\varphi$  is the GELU activation. Concurrently, temporal routing weights  $\mathbf{T} \in \mathbb{R}^{b \times N_e}$  are derived from  $\mathbf{h}_{text}$  using a

---

**Algorithm 1: MOLE WITH STAR AND COGKD**


---

**Input:** Observation  $o_t$ , language instruction  $l$ , ground truth  $\epsilon$ , total layers  $K$

**Output:** Student model loss  $\mathcal{L}_{MoLe}$

- 1 Obtain visual tokens:  $v \leftarrow \mathcal{E}(o_t)$ ;
- 2 Concat multimodal token:  $x \leftarrow concat(v, l)$ ;
- 3 **Step 1: Compute skip indices via STAR router**  
Compute skip indices:  
 $\{g_k\}_{k=1}^K, \mathcal{L}_{lb} \leftarrow G_{star}(v, l);$
- 4 **Step 2: Compute skip indices via STAR router for**  
 $k = 1$  to  $K$  in  $\pi^{(s)}(x)$  do  
5      $h_k \leftarrow \begin{cases} \pi_k^{(s)}(h_{k-1}^{(s)}) & \text{if } G_k = 1 \\ h_{k-1}^{(s)} & \text{otherwise} \end{cases};$   
6 end
- 7 Final student features:  $f^{(s)} \leftarrow h_K^{(s)}$ ;
- 8 Cognition tokens:  $e^{c,(s)} \leftarrow h_K^{(s)}[:, -1, :]$ ;
- 9 **Step 3: Compute skip indices via STAR router**  
Execute all layers in  $\pi^{(t)}(x)$ :  $f^{(t)} \leftarrow h_K^{(t)}$ ;
- 10 Compute CogKD loss:  $\mathcal{L}_{cog} \leftarrow \mathcal{L}_{mse} + \mathcal{L}_{reservekl}$ ;
- 11 **Step 4: Compute skip indices via STAR router**  
Compute action prediction:  $\hat{\epsilon} \leftarrow \mathcal{A}(f^{(s)})$ ;
- 12 Compute loss:  $\mathcal{L}_{MoLe} \leftarrow \mathcal{L}_{task} + \mathcal{L}_{cog} + \mathcal{L}_{lb}$
- 13 **Return**  $\mathcal{L}_{MoLe}$ ;

---

Transformer module, followed by average pooling:

$$\mathbf{T} = \mathbf{W}_t \cdot \Phi(\text{Transformer}(\mathbf{h}_{text})). \quad (10)$$

A dynamic temperature factor  $\alpha \in [0, 1]$ , computed from the [CLS] token of  $\mathbf{h}_{text}$ , modulates routing sharpness:

$$\alpha = \sigma(\mathbf{W}_\tau^\top \cdot \mathbf{h}_{text}^{[CLS]} + b_\tau), \quad (11)$$

where  $\sigma$  is the sigmoid function. The final expert gating weights  $\mathbf{G} \in \mathbb{R}^{b \times N_e}$  combine  $\mathbf{S}$  and  $\mathbf{T}$ , scaled by  $\alpha$ , and are computed via Gumbel-Softmax for differentiable selection:

$$\mathbf{G} = \tau(\alpha \cdot (\mathbf{S} + \mathbf{T}), \tau = 1.0). \quad (12)$$

By integrating spatial and temporal information, our method enables the router to select LLM layers, optimizing performance for VLA tasks adaptively. The approach is efficient, requiring only  $\mathcal{O}(N_e(d_2 + N_{text}^2))$  FLOPs per sample compared to  $\mathcal{O}(N_e d)$  in standard MoE frameworks, where  $d \gg N_{text}, d_2$ . This design ensures high adaptability and computational efficiency.

### 3.4. Cognition self-Knowledge Distillation

While achieving an efficient layer-skipping mechanism, we also design a self-distillation strategy to compensate for the

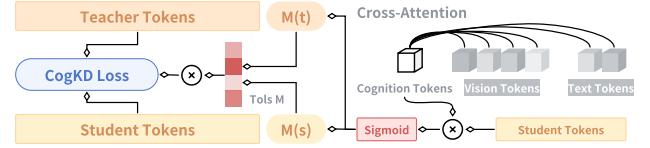


Figure 3. Detailed illustration of our proposed CogKD loss.

cognition loss in the sparse LLM as shown in Fig. 3. Here, we take the original model as the teacher and the MoLe model as the student. To distill the tokens, one common approach is to mimic the tensor token-wisely [7, 41, 49]. Formally, with the tokens  $\mathbf{f}^{(t)} \in \mathbb{R}^{n \times d}$  and  $\mathbf{f}^{(s)} \in \mathbb{R}^{n \times d_s}$  of teacher and student networks, the mimicking can be fulfilled via token reconstruction as

$$\mathcal{L}_{\text{mimic}} = \frac{1}{N} \left\| \mathbf{f}^{(t)} - \mu(\mathbf{f}^{(s)}) \right\|_2^2, \quad (13)$$

However, the Eq.13 treats and distills each token equally, which is inappropriate. For instance, the visual tokens related to the text description should receive more attention [42, 48].

Therefore, we introduce a learnable embedding  $\mathbf{e}_t^c \in \mathbb{R}^{1 \times d}$  dubbed cognition token to distill adaptively. Specifically, it is inserted in the bottom layer, effectively integrating vision tokens and language instruction to understand task requirements better and generate action sequences relevant to the task. The teacher and student model each have their own  $\mathbf{e}_t^{c,(t)}$  and  $\mathbf{e}_t^{c,(s)}$ , respectively. During the distillation, we get the tokens of interests (ToIs)  $\mathbf{M}$  by calculating the similarity between the cognition token and the student tokens:

$$\mathbf{M}^{(i)} = \eta(\mathbf{e}^{c,(i)} \mathbf{f}^{(s)}), i \in \{s, t\}, \quad (14)$$

where  $\eta$  denotes the Sigmoid function. Next, we utilize the intersection of ToIs generated by the teacher and student cognition tokens to decide the distillation degree of each token, where  $\mathbf{M} = \mathbf{M}^{(t)} \odot \mathbf{M}^{(s)}$ , because the distillation tokens should consist of the ones both important to the teacher and student. Therefore, the Eq.13 can be updated as:

$$\mathcal{L}_{\text{cog-mimic}} = \frac{1}{N} \left\| \mathbf{M} \odot \mathbf{f}^{(t)} - \mu(\mathbf{M} \odot \mathbf{f}^{(s)}) \right\|_2^2. \quad (15)$$

Furthermore, we introduce the Reverse-KL [12] paired with our cognition token as the before manner to obtain  $\mathcal{L}_{\text{cog-reversekl}}$  to enhance distribution constraint:

$$\mathcal{L}_{\text{cog-reversekl}} = (\mathbf{M} \odot \mathbf{f}^{(s)}) \log \left( \frac{\mathbf{M} \odot \mathbf{f}^{(s)}}{\mathbf{M} \odot \mathbf{f}^{(t)}} \right). \quad (16)$$

Finally, our eventual CogKD loss can be formulated as

$$\mathcal{L}_{\text{cog}} = (1 - \lambda_1) \mathcal{L}_{\text{cog-mimic}} + \lambda_1 \mathcal{L}_{\text{cog-reversekl}}, \quad (17)$$

where  $\lambda_1$  is the factor and set to 0.5 for balancing the losses.

Table 1. **Performance comparison with existing VLA models across ten tasks in RLBench settings.** We colour-coded the results **red** (1st) and **blue** (2nd) and the row colour reflects the baseline type. The five efficiency methods operate with only 50% LLM layers.

Methods	Action Head	Backbone	Put Rubbish in Bin	Close Box	Close Laptop Lid	Take Umbrella out of Stand	Close Fridge
<b>RLBench</b>							
OpenVLA [16] (CoRL'24)	MLP	LLaMA2-7B	8.0%	72.0%	64.0%	28.0%	88.0%
CogAct [19] (Arxiv'24)	Diffusion	LLaMA2-7B	<b>60.0%</b>	64.0%	76.0%	32.0%	48.0%
RoboMamba [25] (NeruIPS'24)	MLP	Mamba-2.8B	<b>36.0%</b>	60.0%	52.0%	32.0%	68.0%
Random-skip-CogAct	MLP	LLaMA2-7B	16.0%	80.0%	<b>80.0%</b>	32.0%	84.0%
MoD-CogAct [34] (Arxiv'24)	Diffusion	LLaMA2-7B	<b>56.0%</b>	80.0%	68.0%	<b>40.0%</b>	92.0%
DeeR-CogAct [43] (NeruIPS'24)	Diffusion	LLaMA2-7B	52.0%	72.0%	60.0%	36.0%	76.0%
MoLe-OpenVLA (Ours)	MLP	LLaMA2-7B	12.0%	80.0%	76.0%	<b>40.0%</b>	<b>96.0%</b>
MoLe-CogAct (Ours)	Diffusion	LLaMA2-7B	24.0%	<b>84.0%</b>	<b>80.0%</b>	36.0%	88.0%
Methods	Sweep to Dustpan	Phone on Base	Change Clock	Toilet Seat Down	Take Frame off Hanger	Mean Acc.% $\uparrow$	FLOPs (G) $\downarrow$
OpenVLA [16] (CoRL'24)	68.0%	20.0%	16.0%	76.0%	12.0%	45.4%	1930.0
CogAct [19] (Arxiv'24)	44.0%	56.0%	12.0%	<b>100.0%</b>	60.0%	57.2%	1935.8
RoboMamba [25] (NeruIPS'24)	32.0%	44.0%	16.0%	64.0%	32.0%	43.6%	<b>826.3</b>
Random-skip-CogAct	64.0%	24.0%	8.0%	92.0%	32.0%	51.2%(-6.0%)	984.3
MoD-CogAct [34] (Arxiv'24)	4.0%	36.0%	<b>20.0%</b>	96.0%	<b>72.0%</b>	56.4%(-0.8%)	985.8
DeeR-CogAct [43] (NeruIPS'24)	36.0%	<b>68.0%</b>	<b>20.0%</b>	96.0%	68.0%	<b>59.2%</b> (+2.0%)	997.4
MoLe-OpenVLA (Ours)	<b>72.0%</b>	20.0%	12.0%	<b>100.0%</b>	44.0%	55.6%(+10.2%)	<b>981.5</b>
MoLe-CogAct (Ours)	68.0%	36.0%	<b>20.0%</b>	<b>100.0%</b>	<b>72.0%</b>	<b>60.8%</b> (+3.6%)	985.8

### 3.5. Optimization Objective

For the update of the teacher model, we initialize both models with pre-trained parameters and use the exponential moving average (EMA) to update the teacher model  $\pi^{(t)}$ :

$$\pi_t^{(t)} = \alpha \cdot \pi_{t-1}^{(t)} + (1 - \alpha) \cdot \pi_t^{(s)}. \quad (18)$$

In this setup,  $t$  indicates the time step, and we set the update weight  $\alpha = 0.999$  [38].

Our final training objective can be formulated with the combination of  $\mathcal{L}_{task}$ ,  $\mathcal{L}_{cog}$  and  $\mathcal{L}_{lb}$ :

$$\mathcal{L}_{MoLe} = \mathcal{L}_{task} + \lambda_2 \mathcal{L}_{cog} + \lambda_3 \mathcal{L}_{lb}, \quad (19)$$

where  $\lambda_2$  and  $\lambda_3$  are two hyperparameters which are set to 0.5 and 0.1 by default. A more detailed discussion about the hyperparameters can be found in the Appendix.

## 4. Experiments

### 4.1. Implementation details

**Simulation and real-world deployment.** To evaluate our approach and demonstrate its generalization ability, we conduct experiments on both RLBench [13] in the CoppeliaSim simulator and real-world environments with :

1) *RLBench* includes 10 diverse tabletop tasks performed with a Franka Panda robot and a front-view camera. These tasks range from object manipulation to environment interaction, such as: *Close box*, *Close laptop lid*, *Toilet seat down*,

*Put rubbish in bin*, *Sweep to dustpan*, *Close fridge*, *Phone on base*, *Take umbrella out of stand*, *Frame off hanger*, and *Change clock*. Task data are generated using predefined waypoints and the Open Motion Planning Library [36]. Following prior work [14], each task includes 100 training trajectories sampled using a frame-based approach and evaluated in 25 trials per task within the training workspace.

2) *Real-world deployment* is evaluated on the Franka Research 3 (FR3) robot equipped with a 3D-printed UMI gripper [8] across three tasks. A GoPro 9 camera mounted on the wrist captures real-world visual observations. We collect 50 demonstrations for each task, including *detach charger*, *pull drawer*, and *pour water*, using a hand-held UMI gripper within a defined workspace range. A single agent is trained across all tasks and evaluated in 10 trials per task within the training workspace. The success rate is determined through human assessment and serves as the evaluation metric.

**Baselines** The innovation of *MoLe-VLA* lies in its novel, plug-in MoLe architecture, which accelerates VLA inference while improving the robot's success rate. To evaluate its effectiveness, we compare *MoLe* with three state-of-the-art VLA methods across two action generation paradigms: 1) *Autoregressive models*, including *OpenVLA* [16], which uses LLaMA for discrete action prediction, and 2) *Diffusion-based models*, such as *CogAct* [19], which predicts action chunks via a diffusion head. Additionally, we evaluate several VLA efficiency baselines: *RoboMamba* [25], which replaces transformer-based LLMs with a lightweight Mamba model; *DeeR* [43], which enables early exits in LLMs;

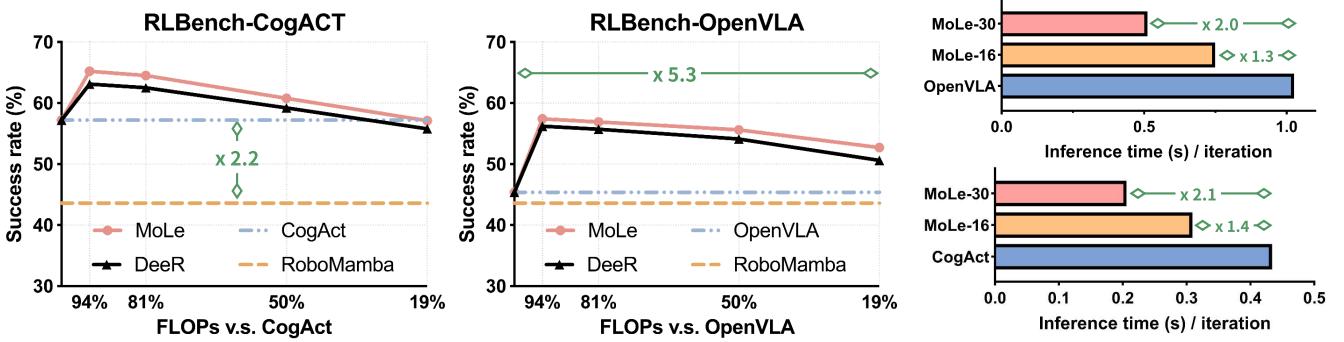


Figure 4. **Efficiency analysis compared with state-of-the-art baselines with FLOPs and inference time.** (Left) Success rate v.s. the FLOPs reduction compared to model backbone. (Right) Inference time per iteration for different layers of MoLe and model backbones.

Table 2. **Inference analysis** evaluated on RLBench simulation environment with FLOPs, inference time, and mean success rate.

Methods	Inference time↓	FLOPs (G)↓	Mean↑
CogAct	0.434 s	1935.8	57.2%
DeeR	0.337 s	997.4	59.2%
MoLe	0.309 s	985.8	60.8%

Table 3. **Effective of model quantization** evaluated on RLBench simulation environment with NVIDIA 4090D GPU.

Methods	Precision	Frequency↑	GPU memory↓	Mean↑
CogAct	FP16	9.8 Hz	16055 MB	57.2%
MoLe	INT8	15.7 Hz	8887 MB	58.8%

*MoD* [34], which allocates input tokens dynamically across layers; and *Random-skip*, which skips LLM layers randomly. For a fair comparison, the latter three baselines are implemented on CogAct with the same setting, with *DeeR* using single-phase training and full model loading. We integrate MoLe with two VLA models, forming *MoLe-OpenVLA* and *MoLe-CogAct*, both using a default **50% layer-skip**.

**Training and evaluation details.** All baselines are trained using the same task configuration for fair comparison. Each method’s official pre-trained parameters are loaded, following their respective training settings. For MoLe-VLA, the single-view RGB input is resized to  $224 \times 224$ , and the robot state is aligned with the predicted actions (7-DOF end-effector poses). The model is trained with a batch size of 64 and 8 diffusion steps per sample, using pre-trained weights for the vision and language modules. The vision module incorporates *DINO-v2* and *SigLIP*, while the language module *LLAMA-2* and the action module *DiT-Base* are trained end-to-end with a constant learning rate of  $2 \times 10^{-5}$  for 1k iterations. Training is conducted on 8 NVIDIA A800 GPUs in approximately 1.5 hours using PyTorch’s Fully Sharded Data Parallel (FSDP) framework.

Table 4. **Scalability analysis** with mean success rate evaluated on RLBench simulation environment with different model sizes.

Methods	CogAct-Small	CogAct-Base	CogAct-Large
CogAct	47.2%	57.2%	70.0%
MoLe	49.9%(+2.7%)	60.8%(+3.6%)	71.5%(+1.5%)

## 4.2. Quantitative results in simulation.

**Performance enhancement** We compare the performance of our proposed MoLe method with state-of-the-art VLA models across ten RLBench tasks, utilizing only half of the LLM layers for efficiency, as shown in Tab. 1. MoLe, implemented with OpenVLA and CogAct backbones, achieves superior success rates and efficiency. Notably, MoLe-CogAct achieves the highest mean success rate of 60.8%, outperforming competing efficiency methods like DeeR of 59.2% and MoD of 56.4% as they overlook the most semantic layers and result in token-wise perception inconsistency, with significant improvements in tasks such as *Close Fridge* and *Sweep to Dustpan*. Similarly, MoLe-OpenVLA demonstrates a 10.2% improvement over the original OpenVLA. Despite requiring only 981.5 and 985.8 GFLOPs, MoLe surpasses DeeR and MoD in efficiency and success rate, highlighting its ability to balance computational cost and task performance. These results underscore MoLe’s effectiveness as a plug-in LLM architecture for robotic manipulation.

**Efficiency analysis** To demonstrate the efficiency of MoLe-VLA, we analyze success rate changes with increasing skipped layers in Fig. 4. MoLe achieves similar success rates compared to the full-layer backbone while only 19% of the FLOPs and delivering  $\times 2$  faster inference. Notably, MoLe-OpenVLA significantly outperforms the original OpenVLA by a large margin. Furthermore, detailed statistics on model efficiency are provided in Tab. 2. MoLe achieves the highest efficiency, requiring only 0.309 seconds per iteration during inference while maintaining the highest mean success rate of 60.8%. These results highlight the superiority of MoLe in balancing efficiency and performance.

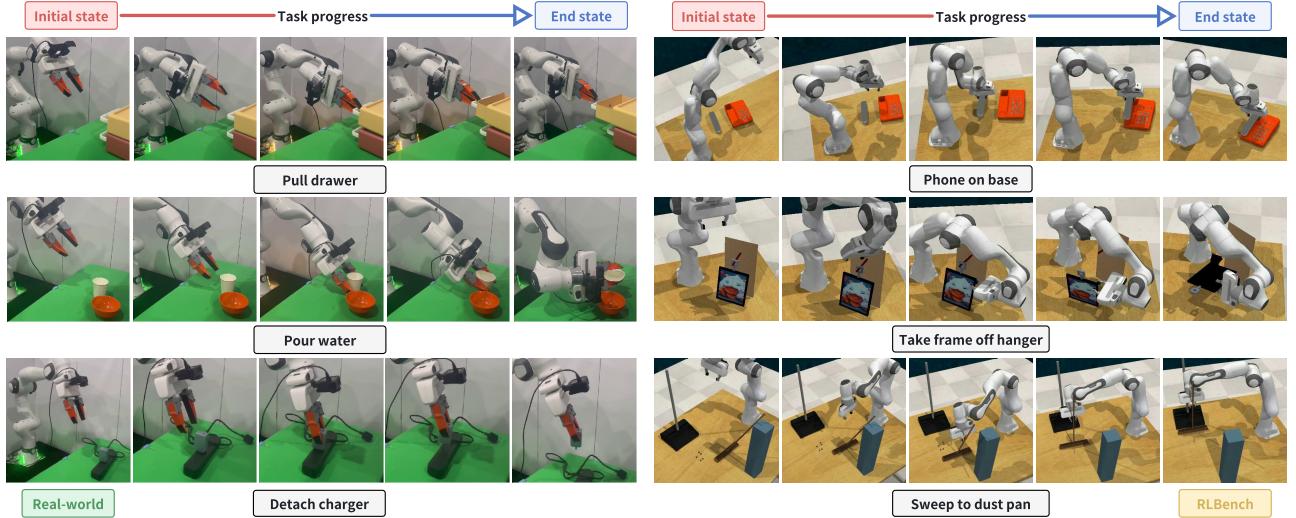


Figure 5. **The qualitative results of MoLe-VLA in both RLBench and real-world**, including the manipulation progress and the task completion end state for both simulation and real-world environments, are shown. More visualizations can be found in the Appendix.

Table 5. **Ablation study** on STAR Router and CogKD loss with its variants on RLBench simulation environment.

Methods	STAR	Cognition	CogKD Loss			Mean↑
			MSE	KL	Reserve KL	
<b>RLBench</b>						
<i>Ex<sub>0</sub></i>	✗	✗	✗	✗	✗	57.2%
<i>Ex<sub>1-1</sub></i>	✓	✗	✗	✗	✗	56.3%
<i>Ex<sub>1-2</sub></i>	✓	✗	✓	✗	✗	54.8%
<i>Ex<sub>2-1</sub></i>	✓	✓	✓	✗	✗	58.3%
<i>Ex<sub>2-2</sub></i>	✓	✓	✗	✓	✗	57.7%
<i>Ex<sub>2-3</sub></i>	✓	✓	✗	✗	✓	59.4%
<i>Ex<sub>2-4</sub></i>	✓	✓	✓	✗	✓	<b>60.8%</b>

**MoLe with quantization analysis** We highlight the efficiency of MoLe under 8-bit quantization, which is more representative of real-world deployment scenarios, compared to FP16 CogAct as shown in Tab. 3 on a commercial-grade RTX 4090D. MoLe achieves a higher success rate of 58.8% with an inference frequency of 15.7 Hz, while utilizing only 55% of the GPU memory compared to CogAct, which achieves just 9.8 Hz. This demonstrates MoLe’s ability to maintain superior performance with significantly lower computational costs after quantization.

**Scalability evaluation** Table 4 highlights the scalability of our proposed MoLe compared to full-layer CogAct across different model sizes evaluated on the RLBench. MoLe consistently achieves higher mean success rates, with improvements of +2.7%, +3.6%, and +1.5% for Small, Base, and Large models, respectively. Notably, MoLe-Large achieves a mean success rate of 71.5%, demonstrating its ability to leverage increased model capacity effectively. These results validate the robustness and adaptability of MoLe across diverse computational budgets and model scales.

Table 6. **Success rate for real-world** evaluated on the FR3 robot equipped with a 3D-printed UMI gripper.

Methods	Detach charger	Pull drawer	Pour water	Mean↑
MoLe	70.0%	60.0%	80.0%	70.0%
CogAct	60.0%	60.0%	80.0%	66.7%

**Ablation study** Table 5 demonstrates the effectiveness of our *STAR* and *CogKD* in the RLBench simulation environment. The baseline CogAct (*Ex<sub>0</sub>*) achieves a mean success rate of 57.2%, while integrating *STAR* with cognition tokens (*Ex<sub>2-1</sub>*) boosts performance to 58.3%, showcasing their synergy. Further improvements are observed with tailored CogKD loss variants, where combining *STAR*, cognition tokens, and Reserve KL loss (*Ex<sub>2-3</sub>*) achieves 59.4%, and the best performance of 60.8% is achieved by adding both MSE and Reserve KL losses (*Ex<sub>2-4</sub>*), a +3.6% gain over the baseline. These results highlight the strength of *STAR* in capturing spatial-temporal dependencies and the importance of cognition tokens for self-knowledge distillation.

### 4.3. Evaluation with real-world tasks.

We conducted experiments involving interactions with various real-world objects, as summarized in Tab. 6. The results show that MoLe consistently delivers strong performance across three tasks. Notably, in the challenging *pour water* task, which demands precise 3D position and rotation predictions, MoLe achieved an impressive success rate of 80%. These results highlight that MoLe preserves the ability to understand 3D spatial scenes and make accurate predictions with a 50% reduction in LLM computational cost.

## 4.4. Qualitative results

As shown in Fig. 5, we visualize the manipulation process for three real-world and three RLBench simulation tasks. Our method accurately predicts continuous 7-DoF end-effector poses, enabling precise task execution along planned trajectories. For instance, in the *pour water* task, MoLe-VLA successfully grasps the cup, lifts the can, positions it above the bowl, and smoothly rotates the gripper to control water flow. Detailed demonstrations are provided in the supplementary video, with failure cases analyzed in the appendix.

## 5. Conclusion

We proposed **MoLe-VLA**, a framework inspired by the Shallow Brain Hypothesis, to optimize VLA models for robotics. MoLe dynamically activates key LLM layers with a specially devised STAR router, reducing redundancy while preserving essential information. To address performance loss from layer skipping, we developed CogKD to enhance efficiency and cognitive capacity. Experiments on real-world and RL-Bench environments show that MoLe reduces computational costs, enabling efficient and adaptable robotic systems.

## 6. Acknowledgment

This work was supported by the National Natural Science Foundation of China (62476011).

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. [3](#)
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022. [2](#)
- [3] Ruichuan An, Sihan Yang, Ming Lu, Kai Zeng, Yulin Luo, Ying Chen, Jiajun Cao, Hao Liang, Qi She, Shanghang Zhang, et al. Mc-llava: Multi-concept personalized vision-language model. *arXiv preprint arXiv:2411.11706*, 2024. [2](#)
- [4] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for training large autoregressive vision-language models, 2023. [3](#)
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pan-nag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricu, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023. [3](#)
- [6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023. [2](#)
- [7] Jiajun Cao, Yuan Zhang, Tao Huang, Ming Lu, Qizhe Zhang, Ruichuan An, Ningning Ma, and Shanghang Zhang. Movekd: Knowledge distillation for vlms with mixture of visual encoders. *arXiv preprint arXiv:2501.01709*, 2025. [5](#)
- [8] Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. *arXiv preprint arXiv:2402.10329*, 2024. [6, 1](#)
- [9] Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, 2024. [3](#)
- [10] Luciano Del Corro, Allie Del Giorno, Sahaj Agarwal, Bin Yu, Ahmed Awadallah, and Subhabrata Mukherjee. Skipdecode: Autoregressive skip decoding with batching and caching for efficient llm inference. *arXiv preprint arXiv:2307.02628*, 2023. [2](#)
- [11] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024. [3](#)
- [12] Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. *arXiv preprint arXiv:2306.08543*, 2023. [5](#)
- [13] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020. [2, 6](#)
- [14] Yueru Jia, Jiaming Liu, Sixiang Chen, Chenyang Gu, Zhilue Wang, Longzan Luo, Lily Lee, Pengwei Wang, Zhongyuan Wang, Renrui Zhang, et al. Lift3d foundation policy: Lift-ing 2d large-scale pretrained models for robust 3d robotic manipulation. *arXiv preprint arXiv:2411.18623*, 2024. [6](#)
- [15] Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: In-

- vestigating the design space of visually-conditioned language models, 2024. 3
- [16] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model, 2024. 2, 3, 6
- [17] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR, 2022. 2
- [18] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation, 2022. 3
- [19] Qixiu Li, Yaobo Liang, Zeyu Wang, Lin Luo, Xi Chen, Mozheng Liao, Fangyun Wei, Yu Deng, Sicheng Xu, Yizhong Zhang, Xiaofan Wang, Bei Liu, Jianlong Fu, Jianmin Bao, Dong Chen, Yuanchun Shi, Jiaolong Yang, and Baining Guo. Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation, 2024. 2, 3, 4, 6
- [20] Tianqin Li, Ziqi Wen, Yangfan Li, and Tai Sing Lee. Emergence of shape bias in convolutional neural networks through activation sparsity. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [21] Xinghang Li, Minghuan Liu, Hanbo Zhang, Cunjun Yu, Jie Xu, Hongtao Wu, Chilam Cheang, Ya Jing, Weinan Zhang, Huaping Liu, et al. Vision-language foundation models as effective robot imitators. *arXiv preprint arXiv:2311.01378*, 2023. 2
- [22] Xinghang Li, Minghuan Liu, Hanbo Zhang, Cunjun Yu, Jie Xu, Hongtao Wu, Chilam Cheang, Ya Jing, Weinan Zhang, Huaping Liu, Hang Li, and Tao Kong. Vision-language foundation models as effective robot imitators, 2024. 3
- [23] Bin Lin, Zhenyu Tang, Yang Ye, Jinfu Huang, Junwu Zhang, Yatian Pang, Peng Jin, Munan Ning, Jiebo Luo, and Li Yuan. Moe-llava: Mixture of experts for large vision-language models, 2024. 2, 3
- [24] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 3
- [25] Jiaming Liu, Mengzhen Liu, Zhenyu Wang, Pengju An, Xiaoli Li, Kaichen Zhou, Senqiao Yang, Renrui Zhang, Yandong Guo, and Shanghang Zhang. Robomamba: Efficient vision-language-action model for robotic reasoning and manipulation, 2024. 2, 3, 6
- [26] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019. 3
- [27] Yulin Luo, Ruichuan An, Bocheng Zou, Yiming Tang, Jiaming Liu, and Shanghang Zhang. Llm as dataset analyst: Subpopulation structure discovery with large language model. In *European Conference on Computer Vision*, pages 235–252. Springer, 2024. 2
- [28] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and Automation Letters*, 7(3):7327–7334, 2022. 2
- [29] OpenAI et al. Gpt-4 technical report, 2024. 3
- [30] Maxime Oquab, Timothée Darzet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 4
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 3
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [33] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. 3
- [34] David Raposo, Sam Ritter, Blake Richards, Timothy Lillicrap, Peter Conway Humphreys, and Adam Santoro. Mixture-of-depths: Dynamically allocating compute in transformer-based language models, 2024. 2, 3, 6, 7
- [35] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer, 2017. 3
- [36] Ioan A Sucan, Mark Moll, and Lydia E Kavraki. The open motion planning library. *IEEE Robotics & Automation Magazine*, 19(4):72–82, 2012. 6
- [37] Mototaka Suzuki, Cyriel MA Pennartz, and Jaan Aru. How deep is the brain? the shallow brain hypothesis. *Nature Reviews Neuroscience*, 24(12):778–791, 2023. 2
- [38] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Learning*, 2017. 6
- [39] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. 3
- [40] Tiannan Wang, Wangchunshu Zhou, Yan Zeng, and Xinsong Zhang. Efficientvlm: Fast and accurate vision-language models via knowledge distillation and modal-adaptive pruning, 2022. 3
- [41] Zhendong Yang, Zhe Li, Ailing Zeng, Zexian Li, Chun Yuan, and Yu Li. Vitkd: Feature-based knowledge distillation for vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1379–1388, 2024. 5

- [42] Weihao Ye, Qiong Wu, Wenhao Lin, and Yiyi Zhou. Fit and prune: Fast and training-free visual token pruning for multi-modal large language models. *arXiv preprint arXiv:2409.10197*, 2024. 5
- [43] Yang Yue, Yulin Wang, Bingyi Kang, Yizeng Han, Shenzhi Wang, Shiji Song, Jiashi Feng, and Gao Huang. Deer-vla: Dynamic inference of multimodal large language models for efficient robot execution, 2024. 2, 3, 6
- [44] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. 4
- [45] Rongyu Zhang, Yun Chen, Chenrui Wu, and Fangxin Wang. Multi-level personalized federated learning on heterogeneous and long-tailed data. *IEEE Transactions on Mobile Computing*, 2024. 3
- [46] Rongyu Zhang, Aosong Cheng, Yulin Luo, Gaole Dai, Huanrui Yang, Jiaming Liu, Ran Xu, Li Du, Yuan Du, Yanbing Jiang, et al. Decomposing the neurons: Activation sparsity via mixture of experts for continual test time adaptation. *arXiv preprint arXiv:2405.16486*, 2024. 2, 3
- [47] Rongyu Zhang, Yulin Luo, Jiaming Liu, Huanrui Yang, Zhen Dong, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, Yuan Du, et al. Efficient deweather mixture-of-experts with uncertainty-aware feature-wise linear modulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 16812–16820, 2024. 2, 3
- [48] Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm: Visual token sparsification for efficient vision-language model inference. *arXiv preprint arXiv:2410.04417*, 2024. 5
- [49] Yuan Zhang, Tao Huang, Jiaming Liu, Tao Jiang, Kuan Cheng, and Shanghang Zhang. Freekd: Knowledge distillation via semantic frequency prompt. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15931–15940, 2024. 5



# MoLe-VLA: Dynamic Layer-skipping Vision Language Action Model via Mixture-of-Layers for Efficient Robot Manipulation

## Supplementary Material

The supplementary materials accompanying this paper provide an extensive quantitative and qualitative analysis of the proposed method. First, we show the deployment of real-world robots in Appendix A. Then, we present the complete set of hyperparameters used in our experiments, detailed in Appendix B, to ensure reproducibility and facilitate further exploration by the research community. In Appendix C, we examine the training scalability of our proposed method, highlighting its performance across varying scales of data and model configurations. Furthermore, we provide complete experiment results on RLBench, exploring the impact of different hyperparameters, including  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$ . These experiments follow the same settings described in the main manuscript, and the results are summarized in Appendix D. We also investigate the impact of skipping different numbers of layers in the RLBench environment, providing a complete evaluation of the trade-offs between efficiency and performance, as detailed in Appendix E. In addition, we include further qualitative analyses in Appendix F, offering visual and descriptive insights into our method's performance enhancements and capabilities. This section highlights specific examples where our approach excels, emphasizing its ability to effectively handle diverse scenarios and complex tasks. Finally, we analyze the failure cases of our proposed methods in real-world environments in Appendix G. These supplementary materials aim to provide a deeper understanding of the proposed method, supporting its robustness and applicability to various tasks.

### A. Real-world Franka robot setup

For our real-world experiments, we utilize the Franka Research 3 (FR3) robotic arm as the hardware platform. To overcome the limitations of the FR3's default gripper, which has relatively short fingers and struggles with certain complex tasks, we 3D-printed and replaced it with a UMI gripper [8]. A GoPro 9 camera is positioned to the right of the setup to capture high-quality RGB images, providing visual input for the pipeline. We conduct experiments on three tasks: *detach charger*, *pull drawer*, and *pour water*. Keyframes are extracted to construct the training set for each task, with 10 frames used for each. Figure 6 illustrates the experimental setup and assets. During the evaluation, task success is determined through human assessment. All actions are performed within the robot's coordinate system to ensure precision and consistency throughout the process. The successful outcomes of the three tasks are shown in the



Figure 6. Franka robot setup.

Table 7. Training hyper-parameters for RLBench.

Hyper-parameters	Values
batch size	64*8
optimizer	AdamW
MLLM learning rate	2e-5
action head learning rate	2e-5
learning rate schedule	constant
warmup steps	2500
LSTM dropout	0.3
MLP dropout	0.4
training epochs	100
$\lambda$	0.05
LSTM window size	12

"End State" images in Figure 5 of the main text.

### B. Training details

We conducted experiments on the RLBench benchmark using the hyperparameters summarized in Table 7. The model was trained with a batch size of  $64 \times 8$ , and the AdamW optimizer was utilized for optimization. The learning rate for the MLLM was set to  $2 \times 10^{-5}$ , while the action head learning rate was configured as  $2 \times 10^{-5}$ . A constant learning rate schedule was adopted, with 2500 warmup steps to stabilize training at the initial stages. To prevent overfitting and enhance generalization, we applied a dropout rate of 0.3 for the LSTM layers and 0.4 for the MLP layers. The LSTM window size was configured to 12 to effectively capture tem-

Table 8. **Performance comparison with existing VLA models across ten tasks in RLBench settings.** We colour-coded the results red (1st) and blue (2nd) and the row colour reflects the baseline type. The four efficiency methods operate with only 50% LLM layers.

Methods	Backbone	Close Fridge	Put Rubbish in Bin	Sweep to Dustpan	Phone on Base	Change Clock
<b>RLBench</b>						
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 0.5, \alpha = 0.999$	MoLe-CogAct	80.0%	16.0%	60.0%	36.0%	24.0%
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 0.5, \alpha = 0.9999$	MoLe-CogAct	84.0%	20.0%	52.0%	28.0%	12.0%
$\lambda_1 = 0.5, \lambda_2 = 0.8, \lambda_3 = 0.5, \alpha = 0.999$	MoLe-CogAct	76.0%	24.0%	32.0%	24.0%	36.0%
$\lambda_1 = 0.5, \lambda_2 = 0.8, \lambda_3 = 0.5, \alpha = 0.9999$	MoLe-CogAct	68.0%	8.0%	44.0%	40.0%	44.0%
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 0.1, \alpha = 0.999$	MoLe-CogAct	60.0%	20.0%	60.0%	48.0%	28.0%
$\lambda_1 = 0.8, \lambda_2 = 0.5, \lambda_3 = 0.1, \alpha = 0.999$	MoLe-CogAct	96.0%	36.0%	8.0%	52.0%	16.0%
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 1.0, \alpha = 0.999$	MoLe-CogAct	96.0%	64.0%	28.0%	52.0%	20.0%
$\lambda_1 = 0.5, \lambda_2 = 0.1, \lambda_3 = 0.5, \alpha = 0.999$	MoLe-CogAct	88.0%	24.0%	68.0%	36.0%	20.0%
Methods	Take Umbrella out of Stand	Take Frame off Hanger	Close Box	Close Laptop Lid	Toilet Seat Down	Mean Acc.% ↑
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 0.5, \alpha = 0.999$	48.0%	72.0%	76.0%	68.0%	84.0%	53.8%
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 0.5, \alpha = 0.9999$	52.0%	68.0%	80.0%	56.0%	80.0%	53.2%
$\lambda_1 = 0.5, \lambda_2 = 0.8, \lambda_3 = 0.5, \alpha = 0.999$	60.0%	64.0%	76.0%	60.0%	92.0%	54.4%
$\lambda_1 = 0.5, \lambda_2 = 0.8, \lambda_3 = 0.5, \alpha = 0.9999$	48.0%	68.0%	80.0%	56.0%	84.0%	54.0%
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 0.1, \alpha = 0.999$	40.0%	60.0%	96.0%	68.0%	92.0%	57.2%
$\lambda_1 = 0.8, \lambda_2 = 0.5, \lambda_3 = 0.1, \alpha = 0.999$	52.0%	56.0%	84.0%	60.0%	100.0%	56.0%
$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 1.0, \alpha = 0.999$	68.0%	24.0%	80.0%	72.0%	92.0%	59.6%
$\lambda_1 = 0.5, \lambda_2 = 0.1, \lambda_3 = 0.5, \alpha = 0.999$	36.0%	72.0%	84.0%	80.0%	100.0%	60.8%

poral dependencies in sequential data. The training process is set to 100 epochs. Additionally, we set the regularization parameter  $\lambda$  to 0.05 to balance different loss components. This setup was chosen to ensure stable and efficient training while maximizing performance on the RLBench tasks. These hyperparameters were fine-tuned based on preliminary experiments to achieve optimal results.

## C. Data scalability

To evaluate the data scalability of our MoLe model, we conducted experiments on a reduced dataset comprising only three tasks: *Close box*, *Close laptop lid*, and *Toilet seat down*. The results, summarized in Tab. 9, demonstrate that MoLe consistently achieves the highest success rate with 82.7% across all tasks, outperforming other methods such as CogAct with 71.0%, Random-skip with 64.1%, and DeeR with 78.6%. Notably, despite being trained on fewer tasks, MoLe maintains superior performance while utilizing only 50% of the computational resources compared to the baseline models. These findings highlight the strong data scalability and computational efficiency of MoLe, making it particularly effective in scenarios with limited training data or constrained computational budgets.

## D. Hyperparameter analysis

We conducted additional experiments on RLBench with different combinations of the hyperparameters. Table 8 demon-

Table 9. **Data scalability analysis** evaluated on RLBench simulation environment with only three tasks.

Methods	Close box↓	Close laptop↓	Seat down↓	Mean↑
CogAct	96.0%	84.0%	32.0%	71.0%
Random-skip	84.0%	60.0%	52.0%	64.1%
DeeR	100.0%	76.0%	60.0%	78.6%
MoLe	100.0%	84.0%	64.0%	82.7%

strates the results of our parameter ablation study, evaluating different configurations of MoLe-CogAct across ten RLBench tasks under a 50% layer-skip setting. The analysis highlights the impact of key hyperparameters ( $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$ , and  $\alpha$ ) on task performance. The results show that increasing the weight of  $\lambda_3$ , which emphasizes the role of specific layers, consistently improves performance, while smaller  $\alpha$  values lead to better optimization stability. Additionally, the influence of  $\lambda_2$  varies depending on the task, indicating its role in balancing intermediate layer contributions. These findings underline the adaptability of MoLe-CogAct and its ability to achieve strong performance and efficiency through careful parameter tuning.

## E. Layer skip analysis

We provide complete experiments on RLBench with different numbers of skipping layers from 2 to 30. As shown in Tab. 10, our proposed MoLe-CogAct demonstrates re-

Table 10. **Performance comparison with existing VLA models across ten tasks in RLBench settings.** We colour-coded the results **red** (1st) and **blue** (2nd) and the row colour reflects the baseline type. The four efficiency methods operate with only 50% LLM layers.

Methods	Backbone	Close Fridge	Put Rubbish in Bin	Sweep to Dustpan	Phone on Base	Change Clock
<b>RLBench</b>						
Skip 2 layers	MoLe-CogAct	92.0%	48.0%	56.0%	56.0%	20.0%
Skip 6 layers	MoLe-CogAct	88.0%	52.0%	48.0%	60.0%	24.0%
Skip 8 layers	MoLe-CogAct	96.0%	64.0%	36.0%	56.0%	56.0%
Skip 12 layers	MoLe-CogAct	100.0%	56.0%	40.0%	60.0%	48.0%
Skip 20 layers	MoLe-CogAct	80.0%	24.0%	36.0%	48.0%	16.0%
Skip 24 layers	MoLe-CogAct	72.0%	40.0%	52.0%	40.0%	16.0%
Skip 26 layers	MoLe-CogAct	80.0%	44.0%	48.0%	36.0%	24.0%
Skip 30 layers	MoLe-CogAct	56.0%	0.0%	0.0%	32.0%	40.0%
Methods	Take Umbrella out of Stand	Take Frame off Hanger	Close Box	Close Laptop Lid	Toilet Seat Down	Mean Acc.% $\uparrow$
Skip 2 layers	48.0%	76.0%	96.0%	68.0%	92.0%	65.2%
Skip 6 layers	52.0%	64.0%	92.0%	56.0%	96.0%	63.2%
Skip 8 layers	32.0%	44.0%	80.0%	68.0%	84.0%	61.6%
Skip 12 layers	52.0%	52.0%	60.0%	56.0%	100.0%	62.4%
Skip 20 layers	56.0%	68.0%	76.0%	48.0%	100.0%	55.2%
Skip 24 layers	52.0%	48.0%	80.0%	52.0%	80.0%	53.2%
Skip 26 layers	48.0%	56.0%	72.0%	48.0%	84.0%	54.0%
Skip 30 layers	44.0%	40.0%	68.0%	28.0%	76.0%	38.4%

markable robustness and efficiency as the number of skipped LLM layers increases. Even with substantial layer skipping, the performance remains stable across most tasks, with only a slight decline in success rates up to 24 skipping layers. Notably, it is only when skipping 30 layers, resulting in an almost 95% reduction in FLOPs, that a significant drop in performance is observed. This highlights the exceptional efficiency of our method, which maintains strong task success rates while drastically reducing computational costs. These results underscore the robustness and adaptability of MoLe, making it a highly effective solution for efficient embodied intelligence tasks.

## F. Additional qualitative results

This section presents visualizations of the manipulation processes for seven RLbench simulation tasks not covered in the main text. As shown in Fig. 7, these visualizations illustrate task executions performed by our proposed MoLe-CogAct model. Each task is specifically designed to evaluate different capabilities of the efficient layer-skipping architecture. Our method accurately predicts 7-DoF end-effector poses, enabling smooth and precise task completion along the defined trajectories. For example, in the *put rubbish in bin* task, MoLe demonstrates a robust ability to grasp the rubbish accurately, lift it smoothly, and drop it precisely into the bin. This task exemplifies the model’s spatial reasoning capabilities, requiring precise perception of both the rubbish and bin positions, as well as the ability to distinguish the rubbish from other objects in the scene. These results

highlight MoLe’s effectiveness in solving tasks that demand both spatial understanding and precise control. Demonstration videos of these tasks are provided in the supplementary material for further reference.

## G. Failure case analysis

As shown in Fig. 8, through comprehensive real-world testing, we identified four key categories of failure cases that hinder MoLe’s performance. The first category is **loss of control**, which often occurs during interactions with target objects, such as in *pull drawer*. These failures are marked by improper force application when handling objects of different weights or by the gripper slipping unexpectedly on smooth surfaces. The second category involves **rotational prediction errors**, which are most evident in tasks requiring precise rotational control, such as *pour water*. Failures in this group include incorrect angles during object interactions and cumulative errors in multi-step rotational motions. The third category pertains to **pose predictions that exceed the robot’s physical limits**. Here, the model occasionally predicts poses beyond the mechanical capabilities of the Franka robotic arm or generates unreachable target positions due to workspace constraints, as observed in tasks like *detach charger*. These failure cases suggest that while our method achieves significant efficiency gains through LLM layer skipping, it comes at the cost of reduced expressiveness and reasoning capacity in the LLM, particularly for tasks requiring fine-grained control and precise spatial understanding.

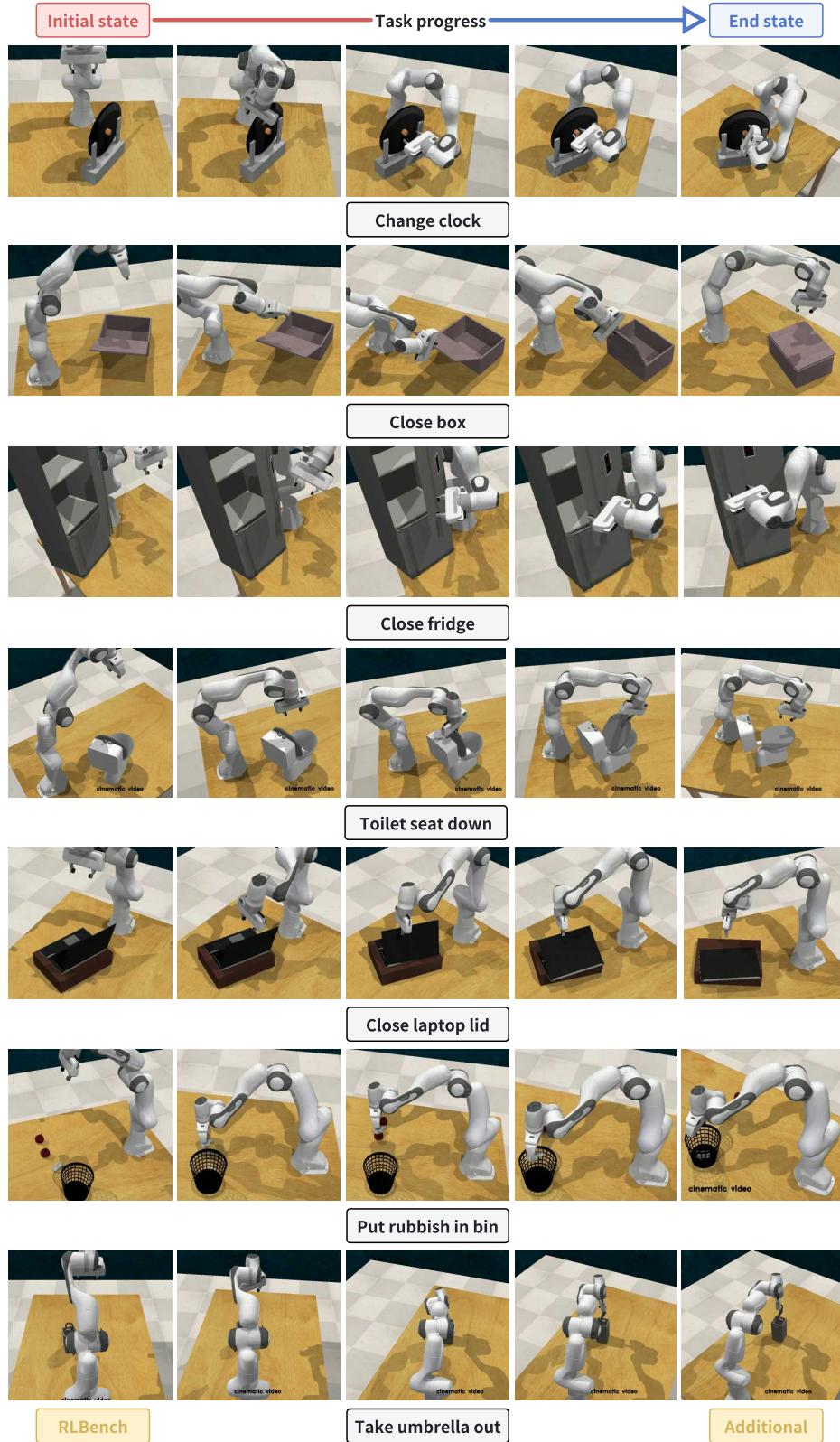


Figure 7. The qualitative results of MoLe-VLA in RLBench simulation environment, including the manipulation progress and the task completion end state.

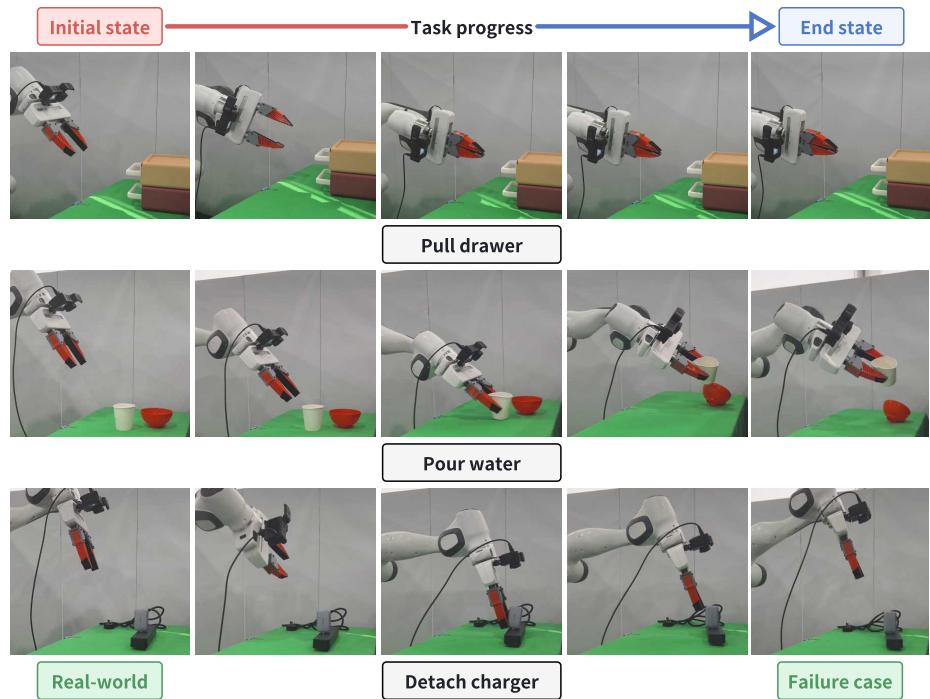


Figure 8. The failure case analysis of MoLe-VLA in real-world environment, including the manipulation progress and the task completion end state.