

SQAP-VLA: A SYNERGISTIC QUANTIZATION-AWARE PRUNING FRAMEWORK FOR HIGH-PERFORMANCE VISION-LANGUAGE-ACTION MODELS

A PREPRINT

Hengyu Fang¹, Yijiang Liu¹, Yuan Du¹, Huanrui Yang², and Li Du¹

¹School of Electronic Science and Engineering, Nanjing University, {hengyufang, liuyijiang}@smail.nju.edu.cn, {yuandu, ldu}@nju.edu.cn

²University of Arizona, huanruiyang@arizona.edu

ABSTRACT

Vision-Language-Action (VLA) models exhibit unprecedented capabilities for embodied intelligence. However, their extensive computational and memory costs hinder their practical deployment. Existing VLA compression and acceleration approaches conduct quantization or token pruning in an ad-hoc manner but fail to enable both for a holistic efficiency improvement due to an observed incompatibility. This work introduces SQAP-VLA, the first structured, training-free VLA inference acceleration framework that simultaneously enables state-of-the-art quantization and token pruning. We overcome the incompatibility by co-designing the quantization and token pruning pipeline, where we propose new quantization-aware token pruning criteria that work on an aggressively quantized model while improving the quantizer design to enhance pruning effectiveness. When applied to standard VLA models, SQAP-VLA yields significant gains in computational efficiency and inference speed while successfully preserving core model performance, achieving a $\times 1.93$ speedup and up to a 4.5% average success rate enhancement compared to the original model.

Code: The code is available at <https://github.com/ecdine/SQAP-VLA>.

1 Introduction

Vision-Language-Action (VLA) models [Chi et al., 2024] constitute a major advancement in embodied intelligence, achieving exceptional performance on tasks that require integrated perception, language understanding, and real-world interaction. These models have catalyzed a broad spectrum of pioneering research in the field [Kim et al., 2024, Li et al., 2024, Brohan et al., 2023a, Black et al., 2024]. However, their substantial computational and memory demands stand in stark contrast to the requirements for low-latency, energy-efficient deployment in edge devices for robotic applications. Bridging this gap necessitates effective model compression strategies. In particular, quantization [Gholami et al., 2021, Nagel et al., 2021] and token pruning [Wang et al., 2021] have emerged as promising approaches for facilitating the efficient deployment of VLA models on resource-constrained hardware.

Quantization is recognized as an effective and generally applicable technique for model compression, especially at low precisions such as W4A4 (4-bit weights and 4-bit activations) [Lin et al., 2025, Liu et al., 2025a, Li et al., 2025a]. In theory, W4A4 quantization can reduce the model size to one quarter and dramatically decrease computational costs compared to full-precision models [Zhao et al., 2024, Shao et al., 2024]. Besides quantization, token pruning [Liu et al., 2023a] directly reduces the computational load and is hardware-friendly, leading to substantial improvements in inference speed [Kuzmin et al., 2024]. At first glance, quantization and token pruning appear to be naturally orthogonal and additive, suggesting that their integration should seamlessly yield highly efficient models [Liang et al., 2021]. However, a straightforward combination of these methods leads to severe performance degradation. This arises from a fundamental coupling between the two techniques: Token pruning leaves the model with limited information to work with, making the model more sensitive to quantization. In contrast, quantization profoundly alters the statistical distribution of features employed in token pruning, such as attention scores. As a result, pruning strategies developed

for high-precision models become ineffective or even invalid when naively applied to quantized networks. This intrinsic incompatibility has been largely overlooked in prior research and presents a major barrier to the deployment of compact and effective VLA models in real-world scenarios.

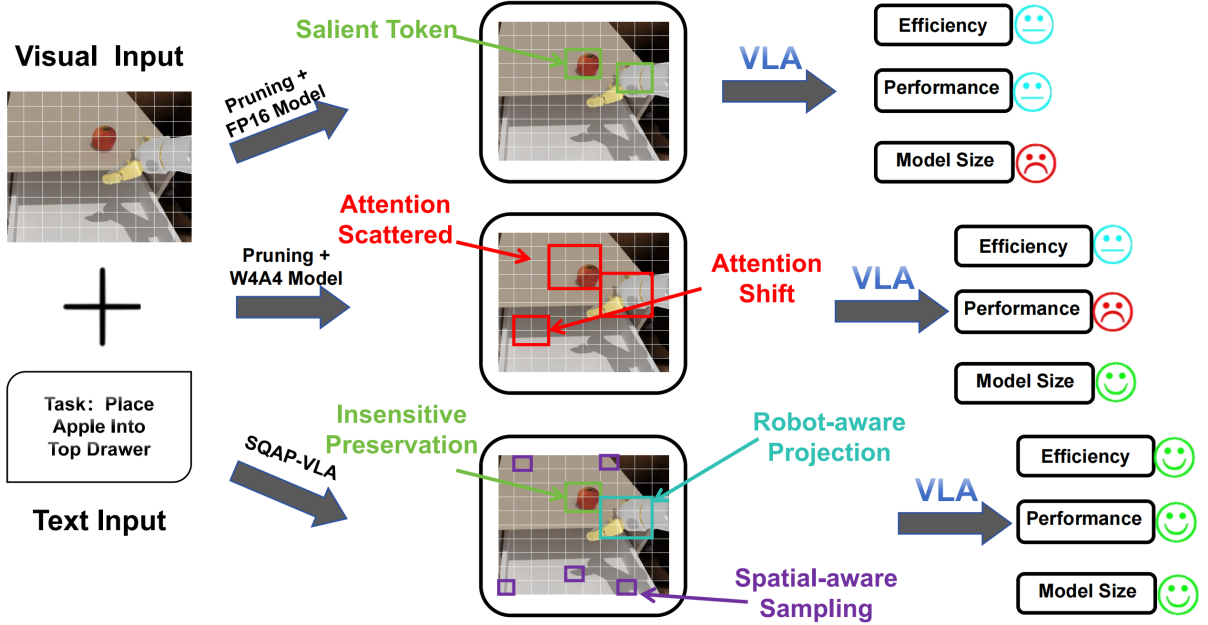


Figure 1: **Overview of SQAP-VLA framework.** SQAP-VLA resolves the incapability of token pruning on quantized VLA models via a quantization-aware pruning criteria. We propose insensitive preservation, robot-aware projection, and spatial-aware sampling to counter the scattered and shifted attention score of the quantized VLA model, enabling high performance, improved speed, and reduced model size with sparse tokens on a quantized VLA model.

To address the aforementioned incompatibility, we contend that an effective deployment of quantization and token pruning necessitates a principled co-design [Hawks et al., 2021], rather than a naive combination. We introduce a novel quantization-aware token pruning framework that jointly optimizes both compression techniques for robust performance [Li et al., 2025b]. On the pruning side, we propose three strategies to ensure adaptation to quantization effects, as shown in Figure 1. First, we identify and utilize quantization-insensitive pruning metrics by selectively retaining tokens with extreme attention scores [Ye et al., 2024], as their relative ordering is empirically robust to quantization noise and reliably preserves the most salient information. Second, we incorporate robot-aware prior protection [Kleeberger and Huber, 2020], i.e., leveraging the world coordinates of a robotic arm to ensure that its corresponding patch tokens are preserved, thereby safeguarding task-critical features. Third, we introduce spatially-aware sampling via Farthest Point Sampling [Qi et al., 2017] to maximize the spatial coverage of retained visual features and prevent information collapse. On the quantization side, we enhance the pruning-friendliness of the activation distributions by integrating Hadamard transforms [Tseng et al., 2024], which mitigate attention score distortion and yield more reliable pruning criteria. This synergistic framework outperforms conventional approaches, achieving superior performance retention under extreme compression and establishing a new standard for efficient, high-performance VLA model deployment on resource-constrained devices [Zhu et al., 2024].

Our primary contributions are:

- We identify the intrinsic incompatibility between quantization and token pruning in Vision-Language-Action (VLA) models and propose a quantization-aware token pruning co-design framework that adapts token pruning strategies to quantization-induced feature distribution shifts, while also enhancing the effectiveness of pruning through quantization techniques.
- On the pruning side, we introduce quantization-insensitive preservation, robot-aware protection, and spatially-aware sampling, thereby adapting token pruning to quantized representations. On the quantization side, we propose the per-tensor Hadamard transformation to enhance attention distributions, thereby facilitating more reliable pruning criteria.

- Comprehensive experiments demonstrate that our integrated method not only resolves the failure of traditional pruning under extreme compression but also achieves the highest degree of performance, forging a viable pathway for the efficient, low-latency deployment of high-performance VLAs on resource-constrained devices.

2 Related Work

2.1 Vision-Language-Action Models

Vision-Language-Action (VLA) models [Brohan et al., 2023b, Driess et al., 2023, Brohan et al., 2023a] constitute a major advancement in embodied artificial intelligence. These models enable end-to-end learning paradigms in which robots can interact with both visual environments and human language instructions. Typical VLA models extend pretrained Vision-Language Models (VLMs) [Liu et al., 2023b] by incorporating mechanisms for generating executable action sequences. A prominent recent trend involves employing diffusion models [Chi et al., 2024] as the action head, wherein the VLM processes visual inputs and textual commands while the diffusion process generates the corresponding action trajectories. This paradigm is exemplified by models such as CogACT [Li et al., 2024] and π_0 [Black et al., 2024]. However, the large model size and the complexity of the decoding process present substantial challenges for deployment on resource-constrained platforms, highlighting the necessity for efficient solutions. In response, this work proposes a synergistic compression framework that integrates model quantization with token pruning, offering a promising direction for enabling efficient and high-performance VLA models.

2.2 VLA Compression

Quantization is recognized as an effective and generally applicable technique for model compression. QAIL [Park et al., 2024] applies quantization-aware training to maintain robust VLA performance with 4-bit weights and activations (W4A4) yet necessitates costly retraining. The large language model (LLM) typically constitutes the dominant component of VLA architectures. Quantization techniques developed for LLMs [Frantar et al., 2023, Lin et al., 2024, Sun et al., 2025, Liu et al., 2025b] thus provide important references for VLA quantization. For instance, QUIP# [Tseng et al., 2024], which incorporates Hadamard transformations to mitigate activation outliers, is promising for VLA quantization. However, their group-wise quantization granularity poses limitations for efficient inference. In this work, we address this shortcoming by integrating the Hadamard transformation with tensor-wise quantization, thereby enhancing adaptability and efficiency in VLA scenarios.

Token Pruning represents another widely adopted approach to model compression. Previously in VLMs, mainstream methods have utilized attention scores to identify and prune insignificant tokens [Chen et al., 2024, Zhang et al., 2025a, Bolya et al., 2023]. Recent token pruning work in VLA models draws inspiration from these techniques. For example, SP-VLA [Li et al., 2025b] leverages the inherent redundancy in visual tokens. Meanwhile, EfficientVLA [Yang et al., 2025] focuses on task relevance, compressing model input by identifying tokens most relevant to the current task while also incorporating a summary of historical information. Cache-VLA [Xu et al., 2025] takes another path by discovering and leveraging position-fixed tokens that are repeatedly used in sequential decision-making, thereby reducing redundant computation through token caching across time steps. Mole-VLA [Zhang et al., 2025b] introduces a dynamic layer-skipping mechanism, which enables the model to bypass redundant layers based on input complexity. Despite their successes in full-precision models, these methods do not consider the potential impact of quantization. We observe significant performance degradation when these techniques are naively combined with quantization. To systematically address this issue, we propose SQAP-VLA, a quantization-aware token pruning co-design framework.

3 Methodology

Our methodology introduces a novel quantization-aware token pruning strategy aimed at maintaining robust VLA performance in resource-constrained environments. To counteract the criteria distortion caused by quantization, we propose three pruning strategies: quantization-insensitive preservation, robot-aware protection, and spatially-aware sampling. In parallel, we enhance the quantizer design to maximize pruning effectiveness via Hadamard transformation and tensor-wise quantization.

3.1 Challenges of Token Pruning on Quantization

Among existing token pruning approaches in VLMs, token importance is typically evaluated through the analysis of attention scores. In full-precision VLA models, task-relevant visual tokens generally attain high attention values. As illustrated in Figure 2a, the tokens corresponding to the “apple” and “robot arm” receive prominent attention scores for

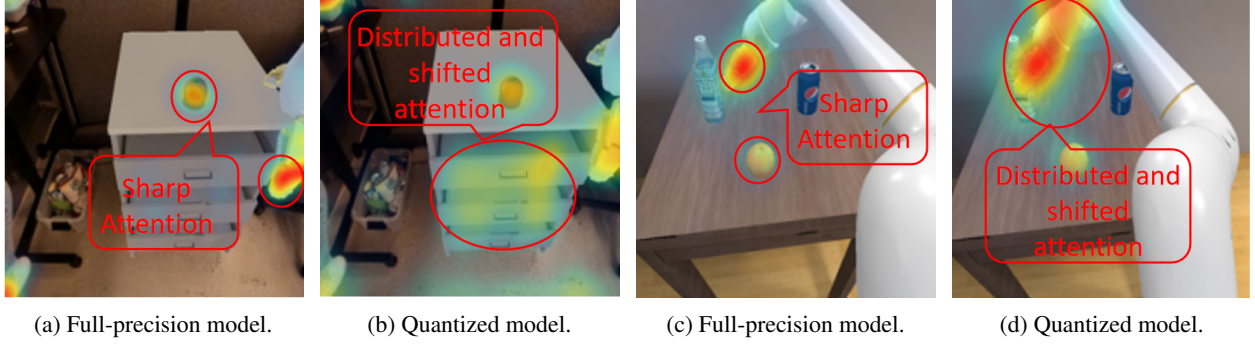


Figure 2: The attention heatmap before and after quantization. (a) and (c) Before quantization, attention is sharply focused. (b) and (d) After quantization, the attention becomes scattered and shifted.

the task “pick up the apple”. However, if weights and activations of the VLA model are quantized to low precisions, noise will be introduced to the inference process, and attention scores will be distorted, resulting in both scattered (i.e., expanded focus areas) and shifted (i.e., focus on irrelevant regions) attention maps. When token pruning methods rely on such severely degraded attention maps, their direct application produces suboptimal or even invalid performance. This incompatibility between quantization and token pruning necessitates more meticulous designs to ensure robust and reliable model operation.

3.2 Quantization-aware Pruning Strategies

To address the aforementioned challenges, we propose three pruning strategies: quantization-insensitive token preservation, robot-aware token protection, and spatially-aware token sampling.

Strategy 1: Quantization-insensitive token preservation. While quantization perturbs the overall attention landscape, this distortion does not affect all tokens uniformly. Our empirical analysis reveals that the numerical shift primarily impacts tokens with mid-range attention scores, blurring the distinction between moderately and minimally important elements. Crucially, the indices of a small number of tokens with the highest-magnitude (“top- k ”) attention scores remain remarkably stable after quantization. These specific tokens are fundamentally important in VLA models, as they consistently correspond to the most task-critical visual elements like target objects or the robot’s end-effector. Because the identity of these top- k tokens is largely invariant to quantization noise, a ‘top- k ’ selection strategy provides a direct and robust mechanism to safeguard this vital information. Accordingly, we define the set of indices for these tokens, $\mathcal{K}_{\text{attn}}$, as:

$$\mathcal{K}_{\text{attn}} = \text{Top}_k(\mathbf{a}_q, k), \quad (1)$$

where $\mathbf{a}_q \in \mathbb{R}^{N_v}$ is the attention weight vector from a task-query token to all N_v visual tokens, and $\text{Top}_k(\cdot, k)$ is an operator that returns the indices of the k largest values in the vector. With a k small enough, the top- k selection is stable under quantization.

Strategy 2: Robot-aware token protection. To compensate for the reduced k for quantization stability and create a truly quantization-invariant anchor, we leverage task-specific priors. We observed that attention scores in the original model consistently correlate with the visual tokens corresponding to the robotic arm. Instead of relying on degraded scores, we directly project the robot’s known 3D world coordinates (x_w, y_w, z_w) into 2D pixel coordinates (u, v) using the camera’s intrinsic (\mathbf{K}) and extrinsic ($[\mathbf{R}|\mathbf{t}]$) matrices:

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{K}[\mathbf{R}|\mathbf{t}] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}. \quad (2)$$

These pixel coordinates are then mapped to discrete token coordinates (t_u, t_v) based on the patch size (P_w, P_h) :

$$t_u = \left\lfloor \frac{u}{P_w} \right\rfloor, \quad t_v = \left\lfloor \frac{v}{P_h} \right\rfloor. \quad (3)$$

This allows us to form a “protected” ring of tokens around a central token \mathbf{t}_c (e.g., the robot’s end-effector). The set of tokens in this ring, $\mathcal{K}_{\text{ring}}$, is selected using the Chebyshev distance (ℓ_∞ norm) within a radius R_t :

$$\mathcal{K}_{\text{ring}} = \{\mathbf{t} \in \mathcal{T} \mid \|\mathbf{t} - \mathbf{t}_c\|_\infty \leq R_t\}. \quad (4)$$

This method, grounded in the robot’s physical state, provides a stable and reliable mechanism for preserving the model’s core visuomotor grounding, regardless of quantization errors.

Strategy 3: Spatially-aware token sampling. Analogous to the selective attention mechanism in human vision, we maintain high fidelity for critical points of interest while representing peripheral regions with reduced detail. After securing the high-importance tokens, we process the remaining tokens $\mathcal{T}_{\text{remain}} = \mathcal{T} \setminus (\mathcal{K}_{\text{attn}} \cup \mathcal{K}_{\text{ring}})$ to reduce redundancy. We apply Farthest Point Sampling (FPS) to select a spatially diverse subset of m tokens from this remainder:

$$\mathcal{K}_{\text{fps}} = \text{FPS}(\mathcal{T}_{\text{remain}}, m). \quad (5)$$

This approach efficiently prunes redundant information while preserving broad spatial coverage, thereby leading to substantial acceleration in model inference. The final set of tokens retained is the union of these three sets: $\mathcal{K}_{\text{final}} = \mathcal{K}_{\text{attn}} \cup \mathcal{K}_{\text{ring}} \cup \mathcal{K}_{\text{fps}}$. In this allocation, $\mathcal{K}_{\text{ring}}$ represents a relatively fixed number of tokens, while the size of $\mathcal{K}_{\text{attn}}$ is the main component adjusted in proportion to the overall target pruning rate. The remaining token quota is then filled by \mathcal{K}_{fps} to ensure the final distribution of unpruned tokens is spatially balanced.

3.3 Pruning-Targeted Quantizer Enhancement

The above pruning strategies partially alleviate the failure of token pruning under quantization. However, the interpretability of the attention map, which enables effective token pruning, is still fundamentally disrupted by quantization-induced artifacts. Specifically, quantization can severely distort the internal representations required for accurate token selection by degrading the quality of attention maps. The attention mechanism depends on query and key vectors projected from input activations, i.e., $Q = W_q^T X$ and $K = W_k^T X$. Our analysis reveals that these activations exhibit a highly skewed and asymmetric distribution. As shown in Figure 3a, a small set of fixed channels consistently produce values several magnitudes larger than the rest, introducing significant outliers into the activation landscape. Traditional token-wise activation quantization is particularly susceptible to these outliers, resulting in considerable quantization errors and degraded attention representation. To address this, we propose the use of channel-wise quantization for activations, which enables finer granularity and achieves improved quantization fidelity, especially for channels with smaller activation values. However, the presence of large outlier channels continues to pose significant challenges for quantization. To further mitigate this issue, we take inspiration from LLM quantization QUIP# [Tseng et al., 2024] to use the Hadamard transformation on weights and activations of the query and key layers with the following formulation:

$$(W^T H^T)(HX) = W^T (H^T H)X = W^T X, \quad (6)$$

where H represents the Hadamard matrix. The Hadamard transform effectively redistributes the activation energy more uniformly across all channels, thereby suppressing outlier effects [Tseng et al., 2024]. As visualized in Figure 3b, this transformation substantially smooths the activation landscape, resulting in improved reliability of the attention map and thereby enhancing token pruning performance.

4 Experiments

In this section, we present the experimental results and analysis. We begin by describing the experimental setup. Next, we demonstrate our main results with a highlight on efficiency. Finally, we conduct ablation studies on the pruning ratio and our proposed compression strategies.

4.1 Experimental Setup

Model Architecture Our experiments are centered on CogAct, a state-of-the-art Vision-Language-Action (VLA) model. Its architecture consists of a Prism-DinoSigLIP-224px [Karamcheti et al., 2024] vision encoder for perception, a Large Language Model (LLM) for high-level reasoning and instruction understanding, and a diffusion model for generating continuous action sequences. We evaluate our framework on two official variants of CogAct presented in the tables, Visual Matching and Variant Aggregation, to demonstrate the general applicability of our method.

Dataset and Training-Free Setup We leverage the official CogAct models, which come pre-trained on the large-scale Open X-Embodiment (OXE) dataset [Collaboration et al., 2025]. A key advantage of our framework, SQAP-VLA, is that it is entirely training-free. We apply our post-training quantization and pruning techniques directly to the publicly available model checkpoints, completely bypassing any need for costly retraining or fine-tuning.

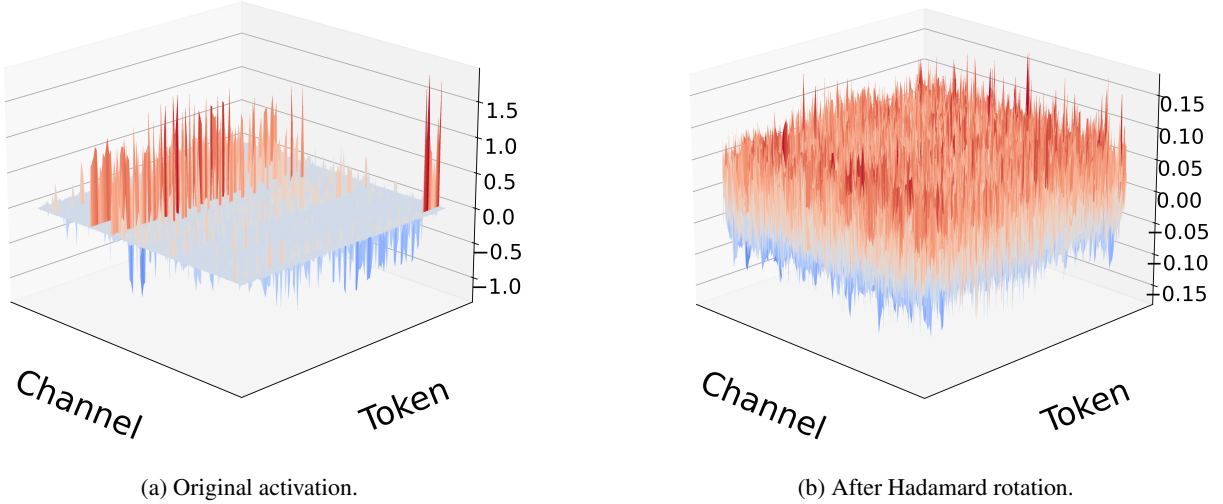


Figure 3: Visualization of activation distributions. (a) Original activations are dominated by large-magnitude outliers in specific channels. (b) After rotation, the energy from these outliers is uniformly distributed, eliminating extreme spikes.

Baselines We evaluate the effectiveness of our approach by comparing it with a comprehensive set of baselines. Specifically, we consider the original CogAct model in full precision (FP16), which serves as the primary reference for both performance and inference speed. Additionally, we benchmark our SQAP-VLA framework against several established token pruning methods applied to the FP16 model, including Random Dropping, FastV [Chen et al., 2024], VLA-Cache [Xu et al., 2025], and EfficientVLA [Yang et al., 2025]. This comparison enables a rigorous assessment of our method relative to state-of-the-art acceleration techniques.

Evaluation Environment and Tasks All evaluations are conducted in a standard robotics simulation benchmark. We report performance on four challenging and representative manipulation tasks that require precise visuomotor control: Pick Coke Can, Move Near, Open/Close Drawer, and Place Apple in Top Drawer.

Evaluation Metrics To provide a comprehensive evaluation of our approach, we employ two primary metrics. The Success Rate (SR, expressed as a percentage) quantifies task performance by measuring the proportion of successful trials. Additionally, we report the speed-up, defined as the theoretical acceleration factor relative to the FP16 CogAct baseline. Together, these metrics capture both the effectiveness and efficiency of the evaluated models.

Implementation Details All our experiments, including quantization and pruning, are conducted in a training-free manner. The efficiency metrics (latency and memory) are benchmarked on a single NVIDIA RTX 3090 GPU. For our proposed method, SQAP-VLA, we apply our synergistic quantization-aware pruning to a W4A4 (4-bit weights and 4-bit activations) quantized model. Based on our ablation studies (Table 3), we use a token pruning ratio of 0.4, which was found to yield the optimal trade-off between performance and efficiency.

4.2 Main Results

The primary task performance results are summarized in Table 1 for the visual matching scenario and Table 2 for the variant aggregation scenario. The baseline for comparison is the pretrained CogAct model in full precision and without pruning. The random dropping approach yields severely degraded performance, indicating its inadequacy for robust token pruning in VLA models. Our SQAP-VLA consistently achieves state-of-the-art success rates and speedup ratios compared to alternative pruning methods, including FastV, VLA-Cache, and EfficientVLA. Remarkably, our model operates under W4A4 quantization, while the competing methods are evaluated in full-precision settings. In the visual matching scenario, our proposed SQAP-VLA outperforms EfficientVLA by 2.9% and exceeds the baseline by 4.5%. In the variant aggregation scenario, our SQAP-VLA achieves a 3.1% success rate improvement compared to the baseline. In terms of computational efficiency, our approach delivers a $1.93\times$ speedup relative to the baseline and achieves a 36% improvement over EfficientVLA. We further observe that specialized token pruning techniques generally outperform the non-pruned baseline. For instance, both EfficientVLA and SQAP-VLA surpass the baseline in

Table 1: **Performance on the Visual Matching scenario.** We compare FP16 pruning methods against our method, which applies pruning to a W4A4 quantized model. Our approach achieves the best average performance. O/C Drawer refers to the Open/Close Drawer task, and BOPs stands for Basic OPERations. Our SQAP-VLA achieves the best average success rate and efficiency.

Method	Quant.	Visual Matching Performance					Speed-up	BOPs
		Pick Coke	Move Near	O/C Drawer	Place Apple	Average		
CogACT (Baseline)	FP16	91.3	85.0	71.8	50.9	74.8	1.0×	100.0%
Random Dropping	FP16	9.7	20.4	53.5	0.0	20.9	1.2×	58.5%
FastV	FP16	92.6	81.4	69.8	52.4	74.1	1.21×	42.0%
VLA-Cache	FP16	92.0	83.3	70.5	51.6	74.4	1.38×	80.1%
EfficientVLA	FP16	95.3	83.3	70.3	56.5	76.4	1.59×	45.1%
SQAP-VLA	W4A4	94.7	85.5	72.2	64.8	79.3	1.93×	26.3%

Table 2: **Performance on the Variant Aggregation scenario.** Our method demonstrates superior or highly competitive performance against FP16 baselines, showcasing its effectiveness even after aggressive low-bit quantization.

Method	Quant.	Variant Aggregation Performance					Speed-up	BOPs
		Pick Coke	Move Near	O/C Drawer	Place Apple	Average		
CogACT (Baseline)	FP16	89.6	80.8	28.3	46.6	61.3	1.0×	100.0%
Random Dropping	FP16	4.0	16.1	15.6	0.0	8.9	1.20×	58.5%
FastV	FP16	91.4	78.6	27.6	50.6	62.1	1.19×	42.0%
VLA-Cache	FP16	91.7	79.3	32.5	45.8	62.3	1.37×	82.6%
EfficientVLA	FP16	94.8	77.6	28.4	51.9	63.2	1.57×	45.1%
SQAP-VLA	W4A4	92.8	80.6	27.2	57.0	64.4	1.93×	26.3%

the visual matching scenario, and all the methods except random dropping achieve improved performance in the variant aggregation scenario.

4.3 Efficiency Analysis

We conduct our efficiency analysis on a single NVIDIA RTX 3090 GPU using the Simpler simulator, focusing on the primary computational bottleneck in Vision-Language-Action (VLA) models: the Large Language Model (LLM) processing during the prefill stage. As demonstrated in Figure 4a, our SQAP-VLA method achieves a notable $1.93\times$ end-to-end system speedup over the FP16 baseline. This system-level gain is driven by a $2.56\times$ acceleration within the LLM backbone, which results from a synergistic combination of W4A4 quantization ($2.09\times$ speedup) and token pruning ($1.21\times$ speedup). In contrast to existing methods such as VLA-Cache and EfficientVLA, which focus exclusively on pruning to reduce the number of operations, our dual approach concurrently reduces both the input sequence length (number of operations) and the computational cost per operation. This integrated strategy enables a more effective reduction of the total computational load, resulting in superior acceleration for large-scale parallel inference. Beyond inference speedup, SQAP-VLA substantially mitigates GPU memory consumption. As illustrated in Figure 4b, the peak GPU memory usage decreases from 14.3 GB with the baseline model to 7.6 GB. This reduction in memory footprint is critical for resource-constrained edge devices, effectively overcoming the primary obstacle to on-device deployment of advanced VLA models.

4.4 Ablation Study

Table 3: **Ablation Study on Pruning Ratios.** We report the average success rate (%) over all tasks. The best performance is highlighted in bold.

Model Variant	Pruning Ratio					Baseline (FP16)
	N/A (W4A4)	0.3	0.4	0.5	0.6	
Visual Matching	71.78	76.55	79.30	78.43	76.85	74.80
Variant Aggregation	58.18	61.35	64.40	63.63	61.18	61.30

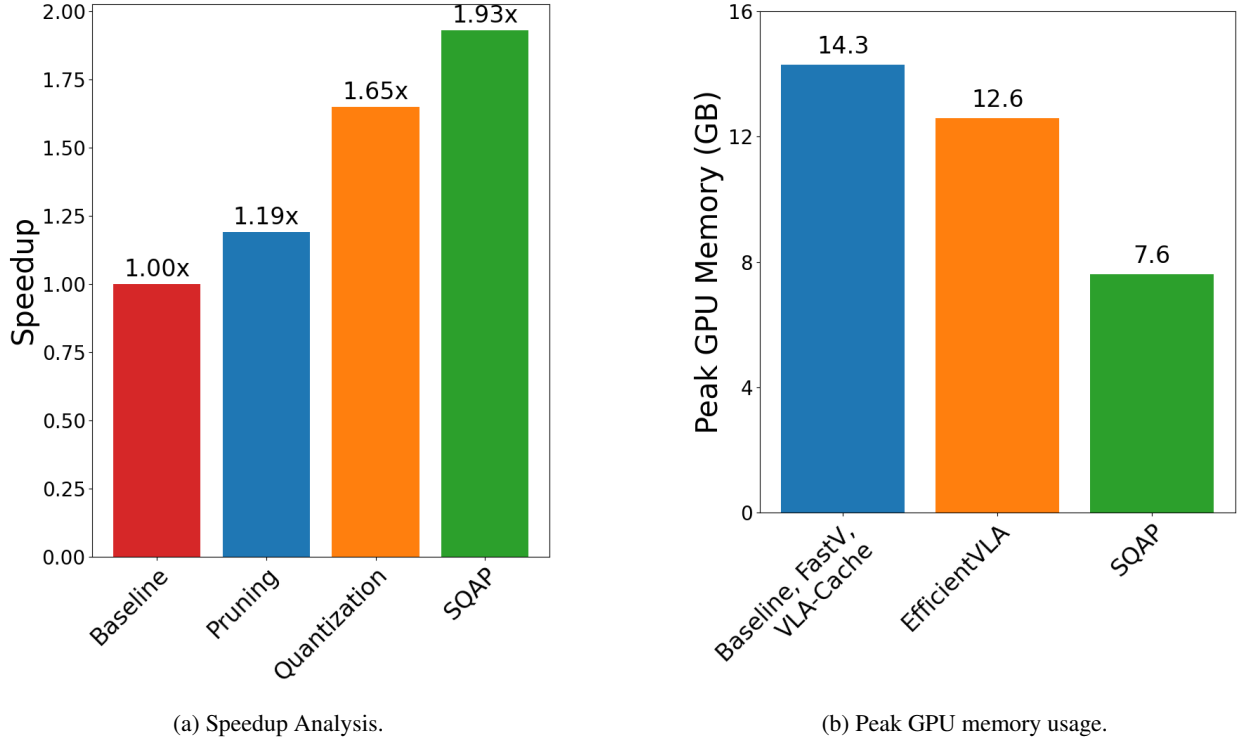


Figure 4: **Latency and memory experiments.** (a) Ablation study on pruning and quantization of SQAP-VLA. (b) GPU memory comparison of the baseline, EfficientVLA, FastV, VLA-Cache, and SQAP-VLA.

We conducted ablation studies on token pruning ratios varying from 0.3 to 0.6 and compression strategies including quantization, quantization-insensitive preservation, robot-aware protection, and spatially-aware sampling.

Pruning Ratios Table 3 presents the results of an ablation study evaluating various token pruning ratios. The baseline model indicates the original model without quantization and pruning. The “N/A” ratio corresponds to the W4A4 model without token pruning. We investigate pruning ratios ranging from 0.3 to 0.6, where the value denotes the proportion of tokens pruned. The W4A4 quantization harms the success rate by 3.02% and 3.12% respectively in the visual matching and variant aggregation scenarios. Notably, models employing our pruning methods consistently outperform the original, unpruned model. SQAP-VLA achieves the highest success rate at a pruning ratio of 0.4. Thereby, we selected a ratio of 0.4 for our main method as it represents an excellent trade-off between high performance and the efficiency gains from more aggressive token pruning.

Compression Strategies Table 4 presents an ablation study of our quantization and pruning strategies, designed to meticulously evaluate the contribution of each component. Our analysis commences with the uncompressed CogACT model as a baseline. As anticipated, applying a uniform W4A4 quantization results in a performance degradation, underscoring the inherent challenge of aggressive compression. While the subsequent integration of a generic quantization-insensitive token pruning affirms the basic compatibility of these techniques, it fails to fully recover the performance loss. The pivotal improvement stems from our context-aware strategies. The introduction of robot-aware token protection effectively reverses the performance decline. Ultimately, the addition of spatially-aware token sampling culminates in a final success rate of 79.30% in visual matching scenarios, not only mitigating the initial quantization-induced deficit but decisively outperforming the full-precision baseline by a significant margin of 4.5%. This empirically validates our core hypothesis that co-designing quantization with intelligent, task-centric pruning is crucial for developing highly efficient models that also achieve superior performance.

5 Conclusion

This paper tackles the deployment of large-scale Vision-Language-Action (VLA) models on resource-constrained platforms via joint token pruning and model quantization. We first identify a fundamental conflict between the two

Table 4: **Ablation Study on Compression Strategies.** “Baseline” denotes the original full-precision model. We use 4-bit quantization on weights and activations. Three token pruning strategies are sequentially introduced to demonstrate the contribution of each component.

Compression Strategy	Model Variant	Success Rate (%)				Average
		Pick Coke Can	Move Near	Open/Close Drawer	Place Apple	
Baseline	Visual Matching	91.3	85.0	71.8	50.9	74.80
	Variant Aggregation	89.6	80.8	28.3	46.6	61.30
+ Quantization	Visual Matching	93.3	83.9	69.2	40.7	71.78
	Variant Aggregation	87.5	75.6	26.7	42.9	58.18
+ Quant-Insensitive	Visual Matching	91.7	81.1	65.4	45.4	70.90
	Variant Aggregation	87.6	74.3	21.0	46.0	57.23
+ Robot-Aware	Visual Matching	92.0	85.0	69.2	49.7	73.98
	Variant Aggregation	90.6	75.6	26.7	48.7	60.40
+ Spatially-Aware	Visual Matching	94.7	85.5	72.2	64.8	79.30
	Variant Aggregation	92.8	80.6	27.2	57.0	64.40

techniques: low-precision quantization catastrophically distorts the attention distributions that traditional token pruning methods rely upon, rendering them ineffective. To resolve this tension, we introduce a novel, training-free framework for the Synergistic Quantization-Aware Pruning of VLA models. Our approach is distinguished by a pruning strategy explicitly co-designed to be compatible with quantization that operates on quantization-insensitive signals to ensure effective token selections. Extensive experiments on the challenging ManiSkill2 benchmark validate the superiority of our method. It achieves a $1.93\times$ speedup and reduces the GPU memory by over 73% while maintaining or even surpassing the task performance of the full-precision baseline. This work presents a principled and effective solution that bridges the gap between high-performance VLA models and the stringent requirements of on-device deployment, offering a viable path toward more capable and responsive robotic systems. Future work could extend this synergistic framework to diverse architectures and explore dynamic, task-adaptive pruning policies.

References

- Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion, 2024. URL <https://arxiv.org/abs/2303.04137>.
- 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. URL <https://arxiv.org/abs/2406.09246>.
- 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, Jialong Yang, and Bain-ing Guo. Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation, 2024. URL <https://arxiv.org/abs/2411.19650>.
- 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, Pannag Sanketi, Pierre Sermanet, Jaspier Singh, Anikait Singh, Radu Soricut, 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, 2023a. URL <https://arxiv.org/abs/2307.15818>.
- Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky. π_0 : A vision-language-action flow model for general robot control, 2024. URL <https://arxiv.org/abs/2410.24164>.

- Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W. Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference, 2021. URL <https://arxiv.org/abs/2103.13630>.
- Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tijmen Blankevoort. A white paper on neural network quantization, 2021. URL <https://arxiv.org/abs/2106.08295>.
- Hanrui Wang, Zhekai Zhang, and Song Han. Spatten: Efficient sparse attention architecture with cascade token and head pruning. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, February 2021. doi:10.1109/hpca51647.2021.00018. URL <http://dx.doi.org/10.1109/HPCA51647.2021.00018>.
- Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song Han. Qserve: W4a8kv4 quantization and system co-design for efficient llm serving, 2025. URL <https://arxiv.org/abs/2405.04532>.
- Yijiang Liu, Hengyu Fang, Liulu He, Rongyu Zhang, Yichuan Bai, Yuan Du, and Li Du. Fbquant: Feedback quantization for large language models, 2025a. URL <https://arxiv.org/abs/2501.16385>.
- Muyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Xiuyu Li, Junxian Guo, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. Svdquant: Absorbing outliers by low-rank components for 4-bit diffusion models, 2025a. URL <https://arxiv.org/abs/2411.05007>.
- Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. Atom: Low-bit quantization for efficient and accurate llm serving, 2024. URL <https://arxiv.org/abs/2310.19102>.
- Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large language models, 2024. URL <https://arxiv.org/abs/2308.13137>.
- Yifei Liu, Mathias Gehrig, Nico Messikommer, Marco Cannici, and Davide Scaramuzza. Revisiting token pruning for object detection and instance segmentation, 2023a. URL <https://arxiv.org/abs/2306.07050>.
- Andrey Kuzmin, Markus Nagel, Mart van Baalen, Arash Behboodi, and Tijmen Blankevoort. Pruning vs quantization: Which is better?, 2024. URL <https://arxiv.org/abs/2307.02973>.
- Tailin Liang, John Glossner, Lei Wang, Shaobo Shi, and Xiaotong Zhang. Pruning and quantization for deep neural network acceleration: A survey, 2021. URL <https://arxiv.org/abs/2101.09671>.
- Benjamin Hawks, Javier Duarte, Nicholas J. Fraser, Alessandro Pappalardo, Nhan Tran, and Yaman Umuroglu. Ps and qs: Quantization-aware pruning for efficient low latency neural network inference. *Frontiers in Artificial Intelligence*, 4, July 2021. ISSN 2624-8212. doi:10.3389/frai.2021.676564. URL <http://dx.doi.org/10.3389/frai.2021.676564>.
- Ye Li, Yuan Meng, Zewen Sun, Kangye Ji, Chen Tang, Jiajun Fan, Xinzhu Ma, Shutao Xia, Zhi Wang, and Wenwu Zhu. Sp-vla: A joint model scheduling and token pruning approach for vla model acceleration, 2025b. URL <https://arxiv.org/abs/2506.12723>.
- 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, 2024. URL <https://arxiv.org/abs/2409.10197>.
- Kilian Kleeberger and Marco F. Huber. Single shot 6d object pose estimation, 2020. URL <https://arxiv.org/abs/2004.12729>.
- Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space, 2017. URL <https://arxiv.org/abs/1706.02413>.
- Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks, 2024. URL <https://arxiv.org/abs/2402.04396>.
- Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models, 2024. URL <https://arxiv.org/abs/2308.07633>.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale, 2023b. URL <https://arxiv.org/abs/2212.06817>.

- Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model, 2023. URL <https://arxiv.org/abs/2303.03378>.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b. URL <https://arxiv.org/abs/2304.08485>.
- Seongmin Park, Hyungmin Kim, Wonseok Jeon, Juyoung Yang, Byeongwook Jeon, Yoonseon Oh, and Jungwook Choi. Quantization-aware imitation-learning for resource-efficient robotic control, 2024. URL <https://arxiv.org/abs/2412.01034>.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers, 2023. URL <https://arxiv.org/abs/2210.17323>.
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration, 2024. URL <https://arxiv.org/abs/2306.00978>.
- Yuxuan Sun, Ruikang Liu, Haoli Bai, Han Bao, Kang Zhao, Yuening Li, Jiaxin Hu, Xianzhi Yu, Lu Hou, Chun Yuan, Xin Jiang, Wulong Liu, and Jun Yao. Flatquant: Flatness matters for llm quantization, 2025. URL <https://arxiv.org/abs/2410.09426>.
- Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. Spinqant: Llm quantization with learned rotations, 2025b. URL <https://arxiv.org/abs/2405.16406>.
- Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models, 2024. URL <https://arxiv.org/abs/2403.06764>.
- Qizhe Zhang, Aosong Cheng, Ming Lu, Renrui Zhang, Zhiyong Zhuo, Jiajun Cao, Shaobo Guo, Qi She, and Shanghang Zhang. Beyond text-visual attention: Exploiting visual cues for effective token pruning in vlms, 2025a. URL <https://arxiv.org/abs/2412.01818>.
- Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster, 2023. URL <https://arxiv.org/abs/2210.09461>.
- Yantai Yang, Yuhao Wang, Zichen Wen, Luo Zhongwei, Chang Zou, Zhipeng Zhang, Chuan Wen, and Linfeng Zhang. Efficientvla: Training-free acceleration and compression for vision-language-action models, 2025. URL <https://arxiv.org/abs/2506.10100>.
- Siyu Xu, Yunke Wang, Chenghao Xia, Dihao Zhu, Tao Huang, and Chang Xu. Vla-cache: Towards efficient vision-language-action model via adaptive token caching in robotic manipulation, 2025. URL <https://arxiv.org/abs/2502.02175>.
- Rongyu Zhang, Menghang Dong, Yuan Zhang, Liang Heng, Xiaowei Chi, Gaole Dai, Li Du, Yuan Du, and Shanghang Zhang. Mole-vla: Dynamic layer-skipping vision language action model via mixture-of-layers for efficient robot manipulation, 2025b. URL <https://arxiv.org/abs/2503.20384>.
- Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models, 2024. URL <https://arxiv.org/abs/2402.07865>.
- Embodiment Collaboration, Abby O’Neill, Abdul Rehman, Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frueger, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu,

Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minh Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Muhammad Zubair Irshad, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundareshan, Qiuyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Martín-Martín, Rohan Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant Halder, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent Vanhoucke, Vitor Guizilini, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open x-embodiment: Robotic learning datasets and rt-x models, 2025. URL <https://arxiv.org/abs/2310.08864>.