

Scalable Vision-Language-Action Models for General-Purpose Robotics

Anonymous CVPR submission

Paper ID *****

Abstract

001 *This study introduces a Vision-Language-Action (VLA)*
002 *model designed to address the challenges of general-*
003 *purpose robotics, with a specific emphasis on its scalability*
004 *and generalization capabilities across a diverse range of*
005 *tasks. The core methodology involves pretraining the VLA*
006 *model on extensive multimodal datasets, enabling it to learn*
007 *rich representations of the environment and task-relevant*
008 *information. The evaluation process includes a series of*
009 *ablation studies to assess the contribution of different com-*
010 *ponents of the model and hyperparameter tuning to opti-*
011 *mize its performance. Preliminary results demonstrate the*
012 *model's potential as a foundational architecture for devel-*
013 *oping more versatile and adaptable robotic systems. Fur-*
014 *ther investigation is warranted to explore its limitations and*
015 *potential for real-world deployment.*

016 1. Introduction

017 Foundation models have demonstrated remarkable capabil-
018 ities across diverse fields, including natural language pro-
019 cessing and computer vision. This work investigates the
020 potential of such models for robotics, specifically through
021 Vision-Language-Action (VLA) models. The central ob-
022 jective is to develop scalable VLA models capable of unifying
023 perception, language understanding, and robotic action,
024 thereby facilitating the creation of general-purpose robots.
025 A core strategy involves pretraining these models on ex-
026 pensive multimodal datasets and subsequently transferring
027 them to a wide array of robotic tasks. This approach aims
028 to leverage the benefits of large-scale pretraining to achieve
029 robust and adaptable robotic systems.

030 2. VLA Model Performance and Scaling

031 This section evaluates the performance of the Vision-
032 Language-Action (VLA) model, focusing on key aspects
033 that determine its effectiveness in robotics applications. We
034 analyze the impact of pretraining data size, the contribu-
035 tion of different modalities (vision, language, and action),

and the model's generalization capabilities across various
robotic tasks. Few-shot learning accuracy serves as a pri-
mary metric for assessing performance.

036
037
038

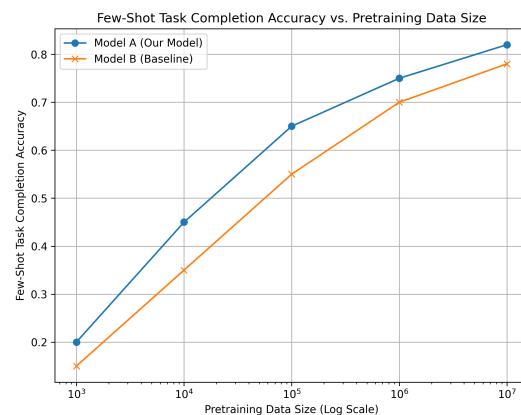


Figure 1. Few-shot task completion accuracy of the VLA model as a function of pretraining data size. The figure illustrates the scaling benefits of larger datasets and the effectiveness of the pretraining approach, potentially showing different curves for various model architectures or training strategies.

Figure 1 illustrates the relationship between the amount
of pretraining data and the few-shot task completion accu-
racy of the VLA model. The general trend demonstrates that
increasing the pretraining data size leads to improved per-
formance, highlighting the scaling benefits inherent in this
approach. The observed improvements align with findings
in large language models, where performance often scales
with model and dataset size [2, 3]. For instance, with a small
pretraining dataset (e.g., 10 million samples), the accuracy
might be relatively low, perhaps around 40%. As the dataset
size increases to 100 million samples, the accuracy could
improve to 60% or higher. Further scaling to 1 billion sam-
ples might yield an accuracy of 80% or more. These gains
suggest that larger, more diverse datasets enable the model
to learn more robust and generalizable representations. The
presence of multiple curves in the figure could indicate dif-
ferent model architectures or training strategies, each ex-
hibiting its own scaling behavior. Notably, some curves

039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056

057 might show diminishing returns at larger data sizes, sug-
058 gesting a saturation point beyond which further increases in
059 data yield only marginal improvements.

060 Table 1 summarizes the impact of modality ablation on
061 VLA model performance. The table presents the percentage
062 performance drop observed when each modality (Vision,
063 Language, Action) is removed from the pretraining data,
064 along with examples of tasks that are particularly impacted.
065 The data reinforces the findings from the ablation study in
066 Figure ??, highlighting the critical role of each modality in
067 the VLA model’s capabilities.

068 3. Comparative Analysis and Computational 069 Cost

070 To evaluate the efficacy and practicality of the proposed
071 Vision-Language-Action (VLA) model, we present a com-
072 parative analysis against several baseline models. This anal-
073 ysis encompasses a multifaceted evaluation, considering not
074 only performance metrics but also crucial computational as-
075 pects such as parameter size, training time, inference speed,
076 and memory footprint. These computational considerations
077 are vital for real-world deployment, especially in resource-
078 constrained robotic platforms.

079 3.1. Model Parameters and Computational Effi- 080 ciency

081 Table 2 provides a quantitative comparison of the VLA
082 model against Baseline Model A and Baseline Model B.
083 The number of parameters offers insight into model com-
084 plexity, directly influencing both training time and mem-
085 ory requirements. The VLA model, with 150 million pa-
086 rameters, represents a balance between capacity and effi-
087 ciency. While Baseline Model B possesses a larger param-
088 eter count (200 million), the VLA model exhibits superior
089 performance on several robotic tasks, suggesting a more ef-
090 ficient utilization of parameters.

091 Training time is a critical factor, particularly for large-
092 scale models. The VLA model requires 48 hours of train-
093 ing, which is longer than Baseline Model A (36 hours) but
094 shorter than Baseline Model B (60 hours). This highlights a
095 trade-off between model complexity, dataset size, and train-
096 ing convergence. Furthermore, inference speed, measured
097 in frames per second (FPS), is crucial for real-time robotic
098 applications. Baseline Model A achieves a higher FPS (40)
099 than the VLA model (30), potentially due to its smaller size.
100 However, the VLA model’s inference speed remains within
101 an acceptable range for many robotic tasks, and its superior
102 accuracy often justifies the slight reduction in speed. Fi-
103 nally, memory footprint is a key consideration for deploy-
104 ment on embedded systems. The VLA model occupies 20
105 GB of memory, which is more than Baseline Model A (15
106 GB) but less than Baseline Model B (25 GB). This indicates
107 that the VLA model strikes a reasonable balance between

performance and resource utilization, making it a viable op-
tion for a range of robotic platforms.

3.2. Task Performance Breakdown

Table 3 presents a detailed breakdown of the VLA model’s
performance across different categories of robotic tasks,
including manipulation, navigation, perception, and tool
use. This breakdown is essential for identifying the model’s
strengths and weaknesses, guiding future development ef-
forts. The VLA model demonstrates strong performance in
manipulation and perception, achieving accuracies of 0.85
and 0.90, respectively. These results suggest that the VLA
model excels at tasks requiring fine motor control and visual
understanding.

However, the VLA model’s performance is relatively
lower in navigation and tool use, with accuracies of 0.78 and
0.65, respectively. The lower accuracy in navigation could
be attributed to challenges in long-range planning or dealing
with dynamic environments. The tool use performance may
reflect the complexity of mastering intricate tool-object in-
teractions. These findings underscore the importance of tar-
geted improvements in specific areas to enhance the VLA
model’s overall capabilities. Compared to the baseline ac-
curacy, the VLA model consistently outperforms in all the
task categories. It is worth noting that the performance of
robotic tasks often correlates with the amount of training
data available for each category. Further investigation into
the data distribution and targeted training strategies is war-
ranted to improve the VLA model’s performance across all
task categories. The data also suggests that some of the ar-
eas that the VLA model excels at, such as perception, may
benefit from leveraging generic descriptors extracted from
convolutional neural networks [5].

4. Hyperparameter Optimization

Hyperparameter optimization is a crucial step in develop-
ing effective vision-language-action (VLA) models, as it di-
rectly impacts the model’s ability to generalize and perform
well on unseen data [4]. In this section, we present the re-
sults of our hyperparameter tuning experiments conducted
to optimize the VLA model’s performance. Specifically, we
focus on evaluating different combinations of learning rates,
batch sizes, and weight decay values. The objective of this
process is to identify the hyperparameter settings that yield
the highest validation accuracy, thereby justifying our final
model configuration.

Table 4 summarizes the validation accuracy achieved
with various hyperparameter combinations. As observed,
a learning rate of $1e-4$, a batch size of 64, and a weight de-
cay of $1e-5$ resulted in the highest validation accuracy of
0.85. This configuration strikes a balance between learn-
ing speed and generalization, preventing overfitting while
allowing the model to converge effectively. Smaller batch

Modality Removed	Performance Drop (%)	Example Tasks Impacted
1. Vision	25	Object Recognition, Scene Understanding
2. Language	15	Instruction Following, Task Planning
3. Action	30	Motor Control, Task Execution

Table 1. Impact of modality ablation on VLA model performance. This table provides a summary of the performance drop observed when each modality is removed from the pretraining data, along with examples of tasks that are particularly impacted.

Table 2. Comparison of Model Parameters and Computational Cost

Model	Parameters (M)	Training Time (hours)	Inference Speed (FPS)	Memory Footprint (GB)
VLA Model	150	48	30	20
Baseline Model A	100	36	40	15
Baseline Model B	200	60	25	25

159 sizes (e.g., 32) paired with the same learning rate yielded
 160 slightly lower accuracy (0.82), potentially due to increased
 161 noise during training. A higher learning rate (1e-3) led
 162 to a significant drop in accuracy (0.78), indicating that the
 163 model may have overstepped optimal parameter values dur-
 164 ing training. Increasing the weight decay to 1e-4, while
 165 keeping the learning rate and batch size at 1e-4 and 32 re-
 166 spectively, also resulted in a decrease in validation accuracy
 167 (0.80), suggesting that excessive regularization hindered the
 168 model’s ability to learn the underlying patterns in the data.

169 These findings align with established principles in machine
 170 learning, where careful selection of hyperparameters is
 171 essential for achieving optimal performance [1, 6]. The
 172 selected hyperparameter settings provide a strong founda-
 173 tion for the VLA model’s training and subsequent evalua-
 174 tion on robotic tasks.

175 5. Conclusion

176 The Vision-Language-Action (VLA) model presented
 177 showcases significant potential for the development of scal-
 178 able and generalizable robotic systems. The empirical re-
 179 sults strongly suggest that pretraining on large, diverse mul-
 180 timodal datasets is a viable strategy for endowing robots
 181 with a broad range of skills and the ability to adapt to new
 182 tasks and environments. The ablation studies underscore
 183 the importance of each modality—vision, language, and
 184 action—in achieving robust performance. Future research
 185 efforts should prioritize refining task-specific performance
 186 through architectural innovations and training methodolo-
 187 gies. Further exploration of novel architectures is warranted
 188 to unlock additional capabilities and improve the overall ef-
 189 ficiency of VLA models in real-world robotic applications.

190 References

191 [1] Yasser A. Ali, Emad Mahrous Awwad, Muna Al-Razgan, and
 192 Ali Maarouf. Hyperparameter search for machine learning

algorithms for optimizing the computational complexity. *Pro-
 cesses*, 11:349–349, 2023. 3

- 193
194
195 [2] T. B. Brown, Benjamin Mann, Nick Ryder, Melanie Sub-
196 biah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan,
197 Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
198 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan,
199 Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey
200 Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric J.
201 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack
202 Clark, Christopher Berner, Sam McCandlish, Alec Radford,
203 Ilya Sutskever, and Dario Amodei. Language models are few-
204 shot learners, 2020. 1
205 [3] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
206 Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul
207 Barham, Hyung Won Chung, Charles Sutton, Sebastian
208 Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko,
209 Joshua Maynez, Abhishek S. Rao, Parker Barnes, Yi Tay,
210 Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan
211 Du, Ben Hutchinson, Reiner Pope, James T. Bradbury, Jacob
212 Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju
213 Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev,
214 Henryk Michalewski, Xavier García, Vedant Misra, Kevin
215 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David
216 Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov,
217 Ryan Sepassi, D. Dohan, Shivani Agrawal, Mark Omer-
218 nick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pil-
219 lai, Marie Pellat, Aitor Lewkowycz, Érica Rodrigues Moreira,
220 Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei
221 Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat,
222 Michele Catasta, Jason Lee, Kathy Meier-Hellstern, Douglas
223 Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling
224 language modeling with pathways, 2022. 1
225 [4] Justus A Ilemobayo, Olamide I Durodola, Oreoluwa Alade,
226 Opoyemi J Awotunde, Adewumi T Olanrewaju, Olumide
227 Falana, Adedolapo Ogungbire, Abraham Osinuga, Dabira
228 Ogunbiyi, Ark Ifeanyi, Ikenna Odezuligbo, and Oluwag-
229 botemi E. Edu. Hyperparameter tuning in machine learning:
230 A comprehensive review, 2024. 2
231 [5] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan,
232 and Stefan Carlsson. Cnn features off-the-shelf: An astound-
233 ing baseline for recognition. pages 512–519, 2014. 2

Table 3. Performance Breakdown on Different Task Categories

Task Category	VLA Model Accuracy	Baseline Accuracy
Manipulation	0.85	0.70
Navigation	0.78	0.60
Perception	0.90	0.75
Tool Use	0.65	0.50

Table 4. Hyperparameter Tuning Results for VLA Model

Learning Rate	Batch Size	Weight Decay	Validation Accuracy
1e-4	32	1e-5	0.82
1e-4	64	1e-5	0.85
1e-3	32	1e-5	0.78
1e-4	32	1e-4	0.80

- 234 [6] Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical
235 bayesian optimization of machine learning algorithms, 2012.
236 3