

CoDA: Cosmos-driven DCT-Adaptive Sparse Attention for Efficient Robotic Trajectory Generation

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Large-scale pretrained world models—most notably*
 002 *NVIDIA Cosmos Predict 2.5—provide powerful trajectory*
 003 *generation capabilities for robotic manipulation. Yet this*
 004 *power comes at a steep computational cost: applying full*
 005 *attention across all timesteps becomes a direct bottleneck*
 006 *in real-world robot deployment. We observe that a robot*
 007 *does not need to “think hard” at every moment equally.*
 008 *Recent work on action tokenization (e.g., the FAST frame-*
 009 *work) has already demonstrated that the high-frequency*
 010 *DCT components of an action sequence effectively encode*
 011 *action complexity. We repurpose this insight into dynamic*
 012 *attention control: at moments when DCT high-frequency*
 013 *energy is strong—when the robot demands precise, com-*
 014 *plex reasoning—we retain more visual tokens; in less*
 015 *demanding intervals, we aggressively reduce them. This*
 016 *modulation is applied per-chunk, leaving open the pos-*
 017 *sibility of further acceleration in real-time deployment*
 018 *scenarios. Building on this, we propose **CoDA** (**Cosmos-***
 019 *driven **DCT-Adaptive Sparse Attention**), which applies a*
 020 *DCT-based dynamic token budget to the cross-attention*
 021 *layers of Cosmos Predict 2.5. We introduce two gating vari-*
 022 *ants: **CoDA-S**, a supervised gate that mimics DCT-derived*
 023 *complexity scores, and **CoDA-A**, an autonomous gate*
 024 *that self-optimizes the accuracy–efficiency trade-off via a*
 025 *penalty-regularized objective. Both variants are fine-tuned*
 026 *from a full-attention model (the baseline) via warm-start,*
 027 *preserving the backbone’s manipulation capability. On*
 028 *the MimicGen Square task, compared to this baseline,*
 029 *CoDA-S reduces visual tokens by **50.7%** without degrading*
 030 *trajectory accuracy, and CoDA-A reduces visual tokens by*
 031 ***25.8%** while improving GT MSE by \sim **30%** (0.51 \rightarrow 0.36).*
 032 *Eliminating unnecessary tokens suppresses attention noise,*
 033 *thereby improving trajectory accuracy itself. Our results*
 034 *demonstrate that adaptive sparsity can simultaneously*
 035 *achieve computational reduction and trajectory accuracy*
 036 *improvement in 7-DoF robotic manipulation.*

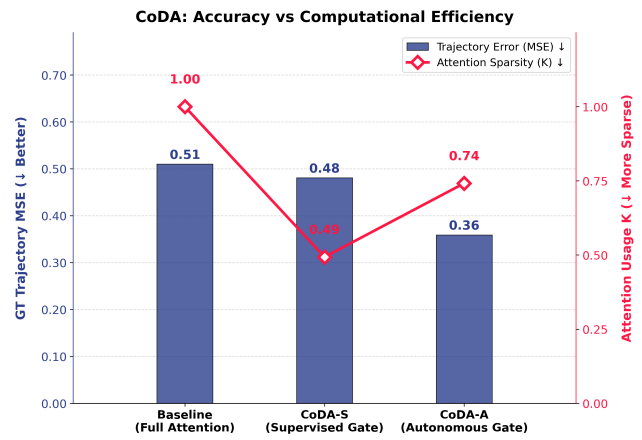


Figure 1. **CoDA achieves a Pareto improvement over the full-attention Baseline on MimicGen Square.** Blue bars (left axis, ↓) show GT Trajectory MSE; the red diamond line (right axis, ↓) shows average token usage K (lower = more sparse). CoDA-S halves the token budget ($K=0.49$) with negligible accuracy loss. CoDA-A simultaneously reduces MSE by \approx 29% and cuts token usage by 25.8%—demonstrating that task-guided sparsity suppresses attention noise and improves trajectory accuracy rather than merely saving compute.

1. Introduction

037
 038 Large-scale pretrained world models such as NVIDIA Cos-
 039 mos Predict 2.5 [1] learn causal relationships of the world
 040 from video data, providing rich visual priors for visuomotor
 041 policy learning in robotic manipulation. Unlike Vision-
 042 Language-Action (VLA) models [6, 7, 18, 25] that rely on
 043 static image-text data and must implicitly infer physical dy-
 044 namics, world models directly learn visual dynamics from
 045 video, enabling more generalizable manipulation capabil-
 046 ities. Our objective addresses a distinct bottleneck com-
 047 pared to large-scale distributed fine-tuning frameworks such
 048 as Cosmos Policy [19]: rather than scaling compute for of-
 049 fine policy generation, we focus on enabling single-GPU
 050 real-time deployment by efficiently leveraging the learned
 051 visual dynamics priors of world models. Yet this power car-
 052 ries an unresolved fundamental question: *when does a robot*

053 *need more attention?* No existing method has directly ad-
054 dressed this question.

055 In practice, only a small fraction of tokens in the atten-
056 tion matrix contribute meaningfully to the output [5, 21,
057 30]. A simple reaching motion and a precise grasping ma-
058 neuver require fundamentally different levels of computa-
059 tion. Yet full attention ignores this distinction, forcing iden-
060 tical computation at every moment. More fundamentally,
061 unnecessary tokens dilute attention scores, preventing the
062 model from focusing on the moments that truly matter.

063 Existing sparse attention methods rely on fixed patterns
064 or post-hoc token dropping, without considering the tempo-
065 ral complexity variations of robotic manipulation or lever-
066 aging domain-specific signals. They amount to reducing
067 computation blindly, without knowing when or where to fo-
068 cus.

069 We answer this question by drawing inspiration from
070 the FAST tokenizer [28], which demonstrated that high-
071 frequency DCT components of action sequences serve as
072 a reliable proxy for action complexity. We translate this
073 frequency-domain insight into an architectural mechanism:
074 a per-chunk dynamic token budget governing the cross-
075 attention layers of Cosmos Predict 2.5. Crucially, CoDA-
076 A goes beyond a fixed heuristic mapping—it treats DCT-
077 derived spectral features as learnable inputs to a gating net-
078 work optimized end-to-end against trajectory loss, enabling
079 the model to distinguish genuine action complexity from
080 high-frequency noise artifacts. We introduce two gating
081 variants: CoDA-S, a supervised gate, and CoDA-A, an au-
082 tonomous gate.

083 Our main contributions are summarized as follows:

- 084 • We propose the **CoDA** framework, which leverages DCT-
085 based action complexity as a dynamic token budget to al-
086 leviate the computational bottleneck of world models in
087 real-time robotic control.
- 088 • We introduce two gating variants—**CoDA-S** (supervised)
089 and **CoDA-A** (autonomous)—and show that CoDA-A’s
090 data-driven optimization not only surpasses the super-
091 vised variant but also suppresses attention noise, improv-
092 ing trajectory accuracy beyond the full-attention baseline.
- 093 • We empirically demonstrate on the MimicGen Square
094 task that adaptive sparsity can simultaneously reduce
095 visual token usage and improve trajectory accuracy—a
096 **Pareto improvement** over the full-attention baseline.

097 2. Related Work

098 2.1. Visuomotor Control via World Models

099 The paradigm of robotic manipulation has recently shifted
100 toward large-scale generative architectures. Beyond tradi-
101 tional VLAs, generative policies such as Diffusion Policy
102 [10] and video generation models [2, 14, 20] are increas-
103 ingly repurposed as universal policies. Generative world

104 models such as Genie [8] and large-scale video generators
105 [26] further demonstrate that video generation can serve as
106 an interactive simulator of physical dynamics, motivating
107 their use as robot policy backbones. Flow-matching-based
108 action policies such as π_0 [4] achieve remarkable gener-
109 alization across dexterous, high-frequency tasks. World-
110 model-based approaches such as Cosmos Policy [19], fine-
111 tuned from the NVIDIA Cosmos Predict 2.5 foundation
112 model [1], demonstrate powerful planning and manipula-
113 tion capabilities by modeling causal visual dynamics. How-
114 ever, these large-scale models typically assume multi-GPU
115 or cloud computing environments. For real-world deploy-
116 ment, especially closed-loop control on a single GPU, the
117 inference latency caused by computing full attention across
118 all timesteps remains a critical bottleneck. While recent
119 works such as DiffuserLite [12] address inference speed by
120 reducing diffusion sampling steps, and transformer-based
121 policies such as Q-Transformer [9] highlight the computa-
122 tional demands of sequential attention in robotic control, the
123 fundamental bottleneck of the cross-attention mechanism it-
124 self remains unaddressed. Our work seeks to bridge this
125 gap, enabling the deployment of world-model-level intelli-
126 gence within strict real-time constraints.

127 2.2. Sparse Attention and Dynamic Routing

128 Since the advent of Vision Transformers [13], the quadratic
129 complexity of self-attention has spurred numerous effi-
130 ciency methods, ranging from hardware-aware exact atten-
131 tion [11] to structural sparsity. Fixed pattern approaches
132 like Longformer [3] and Swin Transformer [23] restrict at-
133 tention to local windows. Token merging methods such
134 as ToMe [5] and adaptive token sampling methods such as
135 ATS [15] reduce FLOPs by combining or dropping tokens
136 based on visual similarity or attention weights. While effec-
137 tive in static vision tasks, these methods rely on fixed pat-
138 terns or visual similarities, critically ignoring the temporal
139 complexity and action-specific demands inherent in robotic
140 control, such as action chunking [16, 31].

141 Another line of research employs learnable sparsity via
142 dynamic token routing and penalty regularization. Methods
143 such as DynamicViT [30] and EViT [21] learn to progres-
144 sively drop inattentive tokens by optimizing a computation-
145 performance trade-off with sparsity penalties. However,
146 these general-purpose penalties optimize for generic atten-
147 tion entropy rather than task-specific accuracy: tokens may
148 be dropped precisely during the most demanding manipula-
149 tion phases, where computational savings are most costly.
150 CoDA overcomes this by introducing an accuracy-aware
151 penalty directly tied to the trajectory loss \mathcal{L}_{action} , ensuring
152 sparsity is pursued only when trajectory accuracy is suffi-
153 ciently high.

154 2.3. Frequency Analysis for Action Representation

155 Analyzing physical actions in the frequency domain of-
 156 fers a robust way to disentangle structural trends from fine-
 157 grained details [17, 29]. The FAST [28] framework demon-
 158 strated that applying the Discrete Cosine Transform (DCT)
 159 to action sequences effectively encodes action complex-
 160 ity, proving that high-frequency components are crucial for
 161 dexterous control tasks. However, FAST uses this insight
 162 solely for offline tokenization during VLA training. In
 163 contrast, CoDA reinterprets DCT high-frequency energy as
 164 a dynamic per-chunk token budget for the cross-attention
 165 mechanism—directly linking action complexity to visual
 166 attention allocation at inference time. Critically, CoDA-
 167 A goes further: rather than using DCT features as a fixed
 168 heuristic, it treats them as learnable inputs optimized end-
 169 to-end against the trajectory loss, enabling the model to
 170 distinguish genuine action complexity from high-frequency
 171 noise artifacts.

172 3. Methodology

173 To address the computational bottleneck of applying full-
 174 attention world models to real-time robotic control, we pro-
 175 pose **CoDA** (**Cosmos-driven DCT-Adaptive Sparse Attention**).
 176 Our core philosophy is that a robot’s computational
 177 load should dynamically scale with the physical complex-
 178 ity of its actions. In this section, we first mathematically
 179 formulate action complexity using Discrete Cosine Trans-
 180 form (DCT) (Sec. 3.1), introduce the dynamic token bud-
 181 geting mechanism in cross-attention (Sec. 3.2), and present
 182 our two adaptive gating variants: CoDA-S and CoDA-A
 183 (Sec. 3.3). Finally, we describe our progressive warm-start
 184 fine-tuning strategy (Sec. 3.4).

185 3.1. Quantifying Action Complexity via DCT

186 Not all segments of a robotic trajectory require the same
 187 level of visual reasoning. A steady translational motion rep-
 188 represents a simple, low-frequency signal, whereas a precise
 189 contact-rich manipulation (e.g., insertion or grasping) ex-
 190 hibits rapid variations and high-frequency patterns. Draw-
 191 ing inspiration from FAST [28], we employ the unnormal-
 192 ized 1D Discrete Cosine Transform (DCT) to mathemati-
 193 cally quantify this “action complexity.”

194 Given an action chunk $A \in \mathbb{R}^{T \times D}$ where T is the time
 195 horizon and D is the action dimension—computed from
 196 ground-truth action chunks during training and from the
 197 previous chunk’s predicted actions (or teacher-forced tar-
 198 gets) during inference-time gating—we transform the tem-
 199 poral sequence of each dimension d into the frequency do-
 200 main:

$$201 \quad C_d(k) = \sum_{n=0}^{T-1} A_d(n) \cos\left(\frac{\pi(2n+1)k}{2T}\right), \quad (1)$$

where $k \in \{0, \dots, T-1\}$ denotes the frequency index. 202

203 We hypothesize that the high-frequency components
 204 contain the critical variations that demand acute visual at-
 205 tention. Thus, we define the action complexity score ϕ as
 206 the ratio of high-frequency spectral magnitude to the total
 207 spectral magnitude, summed across all action dimensions:

$$208 \quad \phi = \frac{\sum_{d=1}^D \sum_{k=\lfloor T/2 \rfloor}^{T-1} |C_d(k)|}{\sum_{d=1}^D \sum_{k=0}^{T-1} |C_d(k)| + \epsilon}, \quad (2)$$

209 Following common frequency-band partitioning in prior
 210 action-frequency analyses, we use the Nyquist midpoint
 211 as the boundary and treat indices $k \geq \lfloor T/2 \rfloor$ as high-
 212 frequency components, which capture rapid temporal vari-
 213 ations in action trajectories. Here, ϵ is a small constant for
 214 numerical stability. Since both numerator and denominator
 215 share the same per-dimension scale, the ratio remains in-
 216 variant to the absolute magnitude of each action dimension.
 217 A higher ϕ indicates a complex maneuver requiring denser
 218 visual context, while a lower ϕ implies a predictable motion
 219 where visual attention can be safely sparsified. Although ϕ
 220 is computed from ground-truth actions during training and
 221 from model predictions at inference, this train–inference
 222 distribution shift is mitigated by CoDA-A’s end-to-end op-
 223 timization, which directly trains the gating network on pre-
 224 dicted spectral features.

225 3.2. Dynamic Cross-Attention via Token Budgeting

226 We build our policy upon the Cosmos Predict 2.5 DiT archi-
 227 tecture [1, 27]. The fundamental computational bottleneck
 228 lies in the cross-attention mechanism, where action query
 229 tokens Q attend to a dense set of visual key tokens K_v, V_v :

$$230 \quad \text{Attn}(Q, K_v, V_v) = \text{softmax}\left(\frac{QK_v^\top}{\sqrt{d}}\right) V_v \quad (3)$$

231 Instead of aggregating over all N visual tokens, CoDA in-
 232 troduces a dynamic *token budget* $K_{budget} \in [0, 1]$. Af-
 233 ter computing attention scores, we retain only the top-
 234 $\lceil K_{budget} \times N \rceil$ visual tokens by masking the rest to $-\infty$
 235 before the softmax. This confines value aggregation to the
 236 most relevant visual tokens, reducing the effective context
 237 each action query must integrate. The degree of sparsity is
 238 determined per-chunk by the complexity score ϕ , directly
 239 linking action complexity to visual attention allocation.
 240 While this masking theoretically reduces the computation
 241 required for value aggregation, realizing strict hardware-
 242 level acceleration requires specialized sparse kernels; this
 243 limitation is further discussed in Sec. 5.

244 3.3. Adaptive Gating: CoDA-S and CoDA-A

245 To determine K_{budget} dynamically during inference, we de-
 246 sign a lightweight **Gating Head** \mathcal{G} . It is a 3-layer MLP
 247 that extracts four DCT-derived spectral features from the

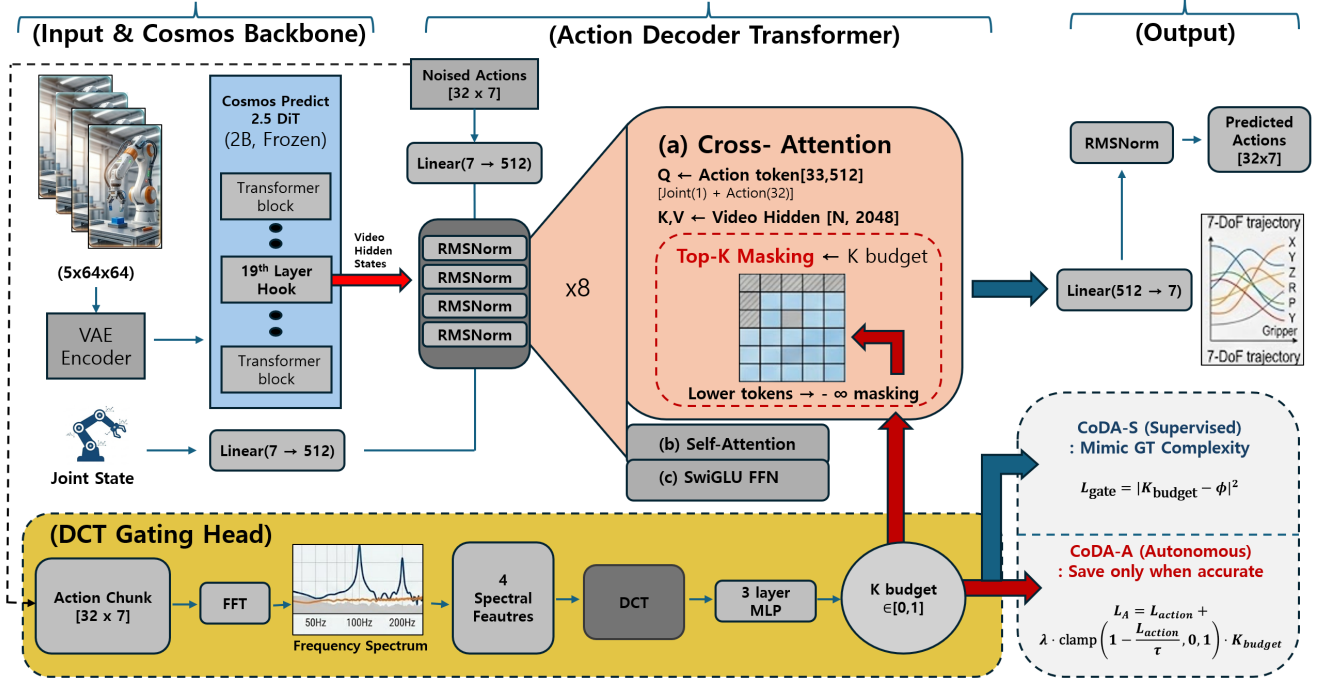


Figure 2. **Overview of the CoDA framework.** (Left) The frozen Cosmos Predict 2.5 backbone extracts video hidden states from 5-frame RGB observations via its 19th-layer hook. Joint state is embedded and concatenated with noised action tokens as queries to the Action Decoder. (Center) The Action Decoder Transformer (8 layers) applies (a) Cross-Attention with dynamic Top-K token masking, (b) Self-Attention, and (c) SwiGLU FFN at each layer. In cross-attention, action query tokens attend to visual key-value tokens; CoDA retains only the top- $\lceil K_{budget} \times N \rceil$ tokens by masking the rest to $-\infty$, concentrating attention on task-relevant visual features. (Bottom) The DCT Gating Head extracts four spectral features from the current action chunk via FFT and DCT, and predicts $K_{budget} \in [0, 1]$ through a 3-layer MLP. CoDA-S (supervised) trains the gate to mimic ground-truth DCT complexity ϕ ; CoDA-A (autonomous) optimizes K_{budget} directly via an accuracy-aware sparsity penalty, reducing tokens only when trajectory accuracy is guaranteed. (Right) The decoder predicts a 7-DoF action trajectory (32×7).

248 available action chunk signal at each step (ground truth in
249 training; model-generated chunk signal in inference)—total
250 spectral magnitude, high-frequency ratio (i.e., ϕ), dominant
251 frequency index, and spectral entropy—and predicts the token
252 budget $K_{budget} = \mathcal{G}(f_{DCT}) \in [0, 1]$. We propose two
253 distinct learning paradigms for this gating mechanism.

254 **CoDA-S (Supervised Gating).** In the supervised vari-
255 ant, we enforce the gating head to explicitly mimic the
256 human-understandable physical complexity. Using the of-
257 fline dataset, we pre-calculate the ground-truth complexity
258 ϕ for every action chunk. The gate is then trained using a
259 standard Mean Squared Error (MSE) objective:

$$260 \quad \mathcal{L}_{gate} = \|K_{budget} - \phi\|^2 \quad (4)$$

261 This ensures the model reliably reduces computation dur-
262 ing low-frequency motions, grounding the sparsity in pure
263 physical kinematics. In practice, the predicted K_{budget} is
264 discretized to one of three hardware-aligned token counts
265 $\{16, 32, 64\}$ to satisfy the block-alignment requirements of

the Triton sparse attention kernel (thresholds: $\phi < 0.15 \rightarrow$ 266
267 16 , $\phi < 0.25 \rightarrow 32$, otherwise $\rightarrow 64$). These thresh-
268 olds were selected to partition the empirical ϕ distribution
269 of the training set into three roughly equal-frequency bins,
270 ensuring balanced coverage across low-, mid-, and high-
271 complexity chunks.

272 **CoDA-A (Autonomous Gating).** While CoDA-S relies
273 on fixed kinematic heuristics, CoDA-A optimizes the
274 accuracy–efficiency trade-off autonomously. We introduce
275 an *Accuracy-aware Sparsity Penalty* that forces the model
276 to reduce visual tokens only when the trajectory accuracy is
277 guaranteed:

$$278 \quad \mathcal{L}_A = \mathcal{L}_{action} + \lambda \cdot \text{clamp}\left(1 - \frac{\mathcal{L}_{action}}{\tau}, 0, 1\right) \cdot K_{budget} \quad (5)$$

279 Here, \mathcal{L}_{action} is the flow-matching trajectory loss [4, 22].
280 We set $\tau=0.3$ based on the converged Stage 0 validation-
281 loss scale, and use a small $\lambda=0.01$ so that the sparsity
282 term regularizes rather than dominates the primary trajec-
283 tory objective; both were selected empirically based on val-

284 idation stability and accuracy–efficiency trade-offs. When
 285 $\mathcal{L}_{action} > \tau$, the gate term vanishes and the model fo-
 286 cuses solely on accuracy recovery. When $\mathcal{L}_{action} < \tau$,
 287 the penalty is activated and gradients flow directly through
 288 K_{budget} to the Gating Head, rewarding reduced token usage.
 289 The small λ and the clamped gate prevent abrupt oscil-
 290 lations, yielding stable convergence. As demonstrated
 291 in our experiments, this data-driven optimization not only
 292 saves computation but also suppresses attention noise, im-
 293 proving trajectory accuracy.

294 3.4. Progressive Warm-Start Fine-tuning

295 Training a 2B parameter foundation model with dynamic
 296 sparsity from scratch is highly susceptible to representation
 297 collapse. To stabilize the learning process, we propose a
 298 progressive warm-start strategy. In **Stage 0**, we fine-tune
 299 the pretrained Cosmos Predict 2.5 weights using full at-
 300 tention ($K_{budget} = 1.0$) to establish a robust visuomotor
 301 policy baseline. In **Stage 1**, we initialize from the Stage 0
 302 checkpoint and attach the Gating Head. The Cosmos back-
 303 bone remains frozen throughout; only the action decoder
 304 and Gating Head are fine-tuned. This curriculum ensures
 305 the model first acquires the dense visuomotor priors neces-
 306 sary for manipulation before learning to sparsify them.

307 4. Experiments

308 4.1. Experimental Setup

309 **Task and Dataset.** We evaluate CoDA on the MimicGen
 310 Square task [24], a 7-DoF robotic manipulation benchmark
 311 requiring a Franka Panda arm to pick up a square nut and
 312 place it onto a fixed peg. We use 3,000 automatically gen-
 313 erated demonstrations for training and 100 held-out demon-
 314 strations for evaluation, with a 90/10 train/validation split.

315 **Architecture and Training.** Our backbone is Cosmos
 316 Predict 2.5 [1], a 2B-parameter Diffusion Transformer
 317 (DiT) pretrained on large-scale video data. In **Stage 0**,
 318 we fine-tune the pretrained backbone with full attention
 319 ($K_{budget} = 1.0$) for 200 epochs using a learning rate of
 320 2×10^{-5} with AdamW optimizer and cosine annealing, es-
 321 tablishing the *Baseline* (our Stage 0 full-attention model;
 322 see Sec. 3.4). In **Stage 1**, we initialize from the Stage
 323 0 checkpoint, attach the Gating Head, freeze the Cosmos
 324 backbone, and fine-tune only the action decoder and Gating
 325 Head for an additional 200 epochs. All experiments are con-
 326 ducted on a single NVIDIA RTX 6000 GPU using bfloat16
 327 precision.

328 **Evaluation Metrics.** We report four complementary met-
 329 rics:

Table 1. Quantitative comparison on MimicGen Square. GT MSE measures trajectory accuracy against ground truth (\downarrow). S0 Fidelity MSE measures deviation from the full-attention baseline (\downarrow). Avg K indicates the average token budget ratio (\downarrow = more sparse). Token Reduction is computed as $(1 - \text{Avg K}) \times 100\%$.

Model	GT MSE \downarrow	S0 Fidelity \downarrow	Avg K \downarrow	Token Red.
Baseline	0.51	0.00	1.00	—
CoDA-S	0.48	0.21	0.49	50.7%
CoDA-A	0.36	0.20	0.74	25.8%

- **GT MSE (\downarrow):** Mean Squared Error between the predicted and ground-truth action trajectories, measuring physical trajectory accuracy. 330
331
332
- **S0 Fidelity MSE (\downarrow):** MSE between the sparse model’s output and the full-attention baseline’s output, quantifying how much intelligence is lost due to sparsification. 333
334
335
- **Avg K (\downarrow):** Average token budget ratio across all evaluation chunks. Lower values indicate greater sparsity and computational savings. 336
337
338
- **Per-DoF MSE:** MSE decomposed across each of the 7 action dimensions (X, Y, Z, Roll, Pitch, Yaw, Gripper), enabling fine-grained error analysis. 339
340
341

All MSE values are reported in per-dimension Z-score normalized action space (zero mean, unit standard deviation computed from the training set). Reported values therefore reflect relative model comparisons rather than absolute physical error magnitudes. 342
343
344
345
346

4.2. Main Results 347

Table 1 and Figure 1 present our main quantitative results. 348
We highlight three key findings. 349

CoDA-A achieves a Pareto improvement. CoDA-A reduces GT MSE from 0.51 to **0.36** while also lowering Avg K to 0.74 in this setting, indicating a favorable accuracy–efficiency trade-off. In Figure 1, both the bar and the line move downward relative to the Baseline, confirming that the accuracy gain is not achieved at the cost of efficiency. This demonstrates that adaptive sparsity is not merely a cost-saving measure: eliminating low-relevance tokens actively suppresses attention noise, enabling the model to concentrate on task-critical visual features and thereby improving trajectory accuracy. 350
351
352
353
354
355
356
357
358
359
360

CoDA-S maximizes token reduction. CoDA-S achieves the most aggressive sparsification, reducing the average token budget to $K=0.49$ (50.7% reduction) while maintaining GT MSE at 0.48—comparable to the full-attention Baseline (0.51). The supervised gate reliably maps low-frequency, kinematically simple motions to high sparsity, grounding token reduction in interpretable physical signals. 361
362
363
364
365
366
367

330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368

369 **Sparse models preserve structural fidelity.** Both
370 CoDA-S and CoDA-A exhibit an S0 Fidelity MSE of
371 ≈ 0.20 , indicating that despite operating with fewer tokens,
372 their output trajectories remain structurally consistent
373 with the full-attention Baseline. The warm-start strategy
374 (Sec. 3.4) is critical: by inheriting dense visuomotor
375 priors from Stage 0, neither sparse model experiences
376 representation collapse during Stage 1 fine-tuning.

377 4.3. Per-DoF Error Analysis

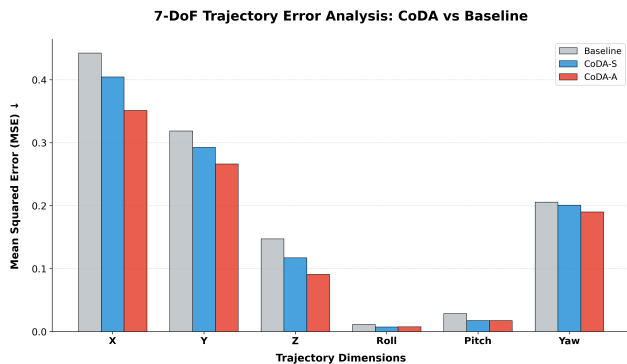


Figure 3. **Per-DoF MSE decomposition across 6 continuous action dimensions.** CoDA-A (red) reduces error over the Baseline (gray) across all translational (X, Y, Z) and rotational (Roll, Pitch, Yaw) axes. The improvement is most pronounced in the Z-axis ($0.15 \rightarrow 0.09$, $\approx 40\%$), which governs vertical alignment during nut insertion—the most contact-rich phase. Roll and Pitch show smaller absolute gains consistent with their inherently low variance. The Gripper dimension is excluded: its binary open/close signal makes continuous MSE an uninformative metric.

378 Figure 3 decomposes trajectory error across the six con-
379 tinuous action dimensions. Two patterns stand out.

380 **Consistent improvement across all axes.** CoDA-A
381 reduces MSE over the Baseline in every dimension
382 without exception (X: $0.44 \rightarrow 0.35$, Y: $0.32 \rightarrow 0.27$, Z:
383 $0.15 \rightarrow 0.09$, Roll: $0.011 \rightarrow 0.008$, Pitch: $0.030 \rightarrow 0.017$,
384 Yaw: $0.21 \rightarrow 0.19$). This consistency rules out the possi-
385 bility that the aggregate MSE gain in Sec. 4.2 is driven by a
386 single outlier dimension.

387 **Largest gains in contact-rich axes.** The Z-axis—which
388 controls vertical positioning during the critical nut-insertion
389 phase—shows the largest relative improvement ($\approx 40\%$).
390 We attribute this to CoDA-A dynamically allocating more
391 tokens during high-frequency, contact-rich moments (as
392 further evidenced by the attention maps in Sec. 4.6), where
393 precise depth control matters most. Roll and Pitch exhibit
394 smaller absolute errors overall, and their improvements,
395 while consistent in direction, are accordingly smaller in
396 magnitude. The Gripper channel is omitted from this fig-
397 ure: its binary open/close control signal renders continuous
398 MSE an uninformative proxy for grasping performance.

Table 2. Ablation on gating strategy and temporal granularity under a fixed cosine LR schedule for fair comparison (Table 1 reports best-performing checkpoints). **Fixed-K**: uniform static sparsity at $K=0.5$ with no dynamic adjustment. **Episode-dynamic**: a single token budget assigned per full episode based on episode-level DCT complexity. **Chunk-dynamic**: per-chunk dynamic budget using the pre-computed GT DCT k -score directly as K (no learned gating head). **CoDA-S/A**: our proposed DCT-guided gating variants (Sec. 3.3).

Model	GT MSE ↓	Avg K ↓
Baseline (Full Attention)	0.652	1.000
Fixed-K ($K=0.5$, static)	0.800	0.500
Episode-dynamic	1.059	0.172
Chunk-dynamic (no learned gate)	0.747	0.172
CoDA-S (Ours)	0.443	0.501
CoDA-A (Ours)	0.492	0.698

4.4. Ablation Study: Architectural Choices for Dynamic Gating

While Table 1 reports best-performing checkpoints to demonstrate peak capability, this ablation evaluates all variants under a fixed, identical training schedule to isolate the effect of each architectural choice. Table 2 reveals a clear three-step progression.

Step 1: Necessity of dynamic gating. Applying static sparsification (Fixed-K, $K=0.5$) severely degrades trajectory accuracy (GT MSE=0.800) relative to the full-attention Baseline (0.652), despite using the same token budget as CoDA-S. Uniformly dropping tokens destroys task-critical visual priors at precisely the moments demanding precise reasoning.

Step 2: Necessity of chunk-level temporal granularity. Episode-dynamic gating—which assigns a single token budget per full episode—performs worst of all (GT MSE=1.059, exceeding even the Baseline), because a globally underestimated budget forces the model to operate with severe sparsity throughout contact-rich insertion phases. Finer-grained per-chunk budgeting (Chunk-dynamic) reduces this error to 0.747, confirming that action complexity fluctuates rapidly and requires localized temporal resolution.

Step 3: Effectiveness of DCT-guided gating signals. Even with chunk-level granularity, the naive Chunk-dynamic variant remains inferior to the Baseline. Replacing the naive signal with our DCT-derived complexity score (CoDA-S) dramatically improves accuracy to GT MSE=0.443—below the Baseline—at a comparable token budget (Avg $K \approx 0.50$). This confirms that the high-frequency DCT signal is the critical factor enabling the

431 model to retain tokens precisely when and where precise
432 reasoning is required.

433 4.5. Training Dynamics

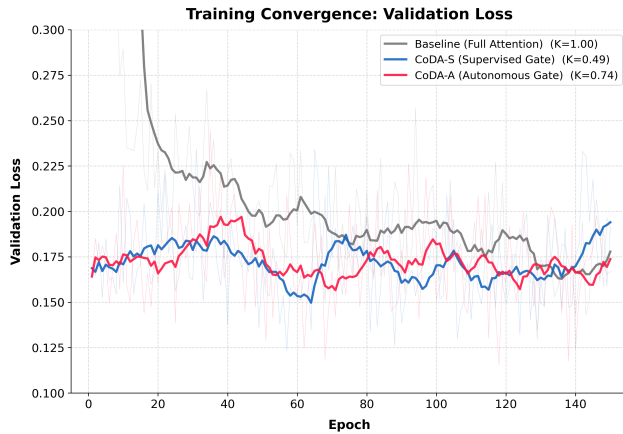


Figure 4. **Validation loss convergence for all three models.** The Baseline (gray) starts from pretrained Cosmos weights at loss ≈ 0.30 and requires ~ 60 epochs to settle. CoDA-S (blue) and CoDA-A (red) inherit the Stage 0 checkpoint and begin Stage 1 already near loss ≈ 0.17 —a direct consequence of the warm-start strategy. Both sparse models converge stably without loss spikes, confirming that freezing the Cosmos backbone and fine-tuning only the action decoder and Gating Head prevents representation collapse.

434 Figure 4 compares the validation loss trajectories of all
435 three models.

436 **Warm-start eliminates cold-start cost.** The Baseline
437 (Stage 0) must build visuomotor priors from scratch atop
438 the video-pretrained Cosmos weights, starting at a high
439 initial loss (≈ 0.30) and requiring roughly 60 epochs
440 before settling. CoDA-S and CoDA-A bypass this cold-start
441 phase entirely: inheriting the Stage 0 checkpoint, both
442 models enter Stage 1 with an already-low validation loss
443 (≈ 0.17) and maintain stable convergence throughout.

444 **No representation collapse under sparsification.** De-
445 spite operating with reduced token budgets from the very
446 first epoch of Stage 1, neither sparse model exhibits loss
447 spikes or divergence. This stability arises from our archi-
448 tectural choice to freeze the Cosmos backbone and fine-tune
449 only the lightweight action decoder and Gating Head, ensur-
450 ing that the dense visuomotor priors established in Stage 0
451 are preserved intact while the gating mechanism is learned
452 on top.

453 4.6. Qualitative Analysis: Attention Maps

454 Figure 5 shows Layer 3 cross-attention heatmaps across key
455 manipulation stages.

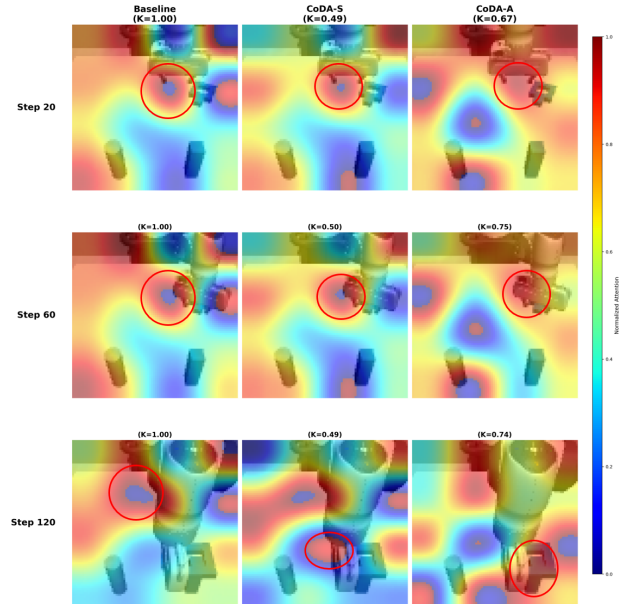


Figure 5. **Cross-attention heatmaps (Layer 3) at three manipulation stages.** Each row corresponds to a timestep (Step 20: approach, Step 60: contact, Step 120: insertion). Each column shows the attention distribution for Baseline ($K=1.0$), CoDA-S ($K=0.49$), and CoDA-A ($K=0.74$). Red circles mark the task-critical nut-peg contact zone. **Left (Baseline):** Attention is broadly diffused across the visual field throughout all stages, including task-irrelevant background regions. **Center (CoDA-S):** Attention is more concentrated near the manipulation workspace, though some diffuse patterns remain. **Right (CoDA-A):** At Step 120 (insertion), attention sharpens most strongly onto the nut-peg contact zone, while K rises dynamically to accommodate the increased complexity—demonstrating that CoDA-A allocates more tokens precisely when and where they matter.

Baseline exhibits attention diffusion. With full token access ($K=1.0$), the Baseline distributes attention broadly across the scene throughout all stages, including background regions irrelevant to the task. This attention dilution is the mechanism behind its higher trajectory error: softmax normalization forces the model to aggregate signal from noisy, task-irrelevant tokens alongside meaningful ones.

CoDA-A acts as an implicit attention regularizer. By pruning low-relevance tokens, CoDA-A constrains the attention distribution to focus on the task-critical zone. Critically, the dynamic nature of CoDA-A’s gating is visible here: at the demanding insertion step (Step 120), the gate raises K to retain more tokens, producing the sharpest heatmap concentration of all three models. This is the mechanism underlying the Pareto improvement observed in Sec. 4.2: removing noisy tokens does not merely save compute—it forces the attention mechanism to find sharper, more informative alignments, directly improving trajectory

475 prediction.

476 5. Conclusion and Future Work

477 We presented **CoDA**, a DCT-adaptive sparse attention
478 framework that dynamically allocates visual token budgets
479 in proportion to robotic action complexity. By repurpos-
480 ing DCT high-frequency energy—originally proposed by
481 FAST for offline tokenization—as a per-chunk gating signal
482 in the cross-attention layers of Cosmos Predict 2.5, CoDA
483 bridges the gap between world-model-level intelligence and
484 real-time deployment constraints. Our two gating variants
485 reveal a fundamental insight: sparsity, when guided by task-
486 relevant signals, is not merely a computational shortcut but
487 an implicit regularizer. CoDA-S reduces token usage by
488 50.7% while preserving trajectory accuracy, and CoDA-
489 A achieves a Pareto improvement—simultaneously cutting
490 tokens by 25.8% and improving GT MSE by $\sim 30\%$ —by
491 learning to suppress attention noise from task-irrelevant vi-
492 sual regions.

493 **Limitations.** Current experiments are limited to a single
494 simulated manipulation task (MimicGen Square). Due to
495 the substantial computational cost of fine-tuning the 2B-
496 parameter backbone, our evaluation relies on a single train-
497 ing run; multi-task generalization and multi-seed variance
498 are left for future work. Furthermore, while CoDA sig-
499 nificantly reduces visual tokens, strict wall-clock inference
500 speedup has not yet been demonstrated. Although CoDA-
501 S already discretizes token budgets into hardware-aligned
502 block sizes ($\{16, 32, 64\}$) to satisfy Triton kernel require-
503 ments, end-to-end kernel integration and wall-clock latency
504 validation remain an immediate next step.

505 **Future Work.** Immediate next steps include validating
506 Triton kernel integration on RTX 6000 Ada (SM_89) to
507 demonstrate concrete latency reduction, and extending
508 CoDA to multi-task and real-robot settings. Longer-term
509 directions include exploring per-query token budgets for
510 finer-grained adaptive attention, and applying the CoDA
511 framework to other large-scale world models beyond Cos-
512 mos Predict 2.5.

513 References

- 514 [1] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji,
515 Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin
516 Chen, et al. Cosmos world foundation model platform for
517 physical AI. *arXiv preprint arXiv:2501.03575*, 2025. 1, 2, 3,
518 5
- 519 [2] Anurag Ajay et al. Is conditional video generation a good
520 object-centric policy? In *ICLR*, 2024. 2
- 521 [3] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Long-
522 former: The long-document transformer. *arXiv preprint*
523 *arXiv:2004.05150*, 2020. 2

- [4] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, 524
Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, 525
Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim 526
Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith 527
Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, 528
James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, 529
and Ury Zhilinsky. π_0 : A vision-language-action flow model 530
for general robot control. *arXiv preprint arXiv:2410.24164*, 531
2024. 2, 4 532
- [5] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao 533
Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token 534
merging: Your ViT but faster. In *The Eleventh International* 535
Conference on Learning Representations, 2023. 2 536
- [6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen 537
Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakr- 538
ishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. 539
RT-1: Robotics transformer for real-world control at scale. 540
In *Robotics: Science and Systems*, 2022. 1 541
- [7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen 542
Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, 543
Danny Driess, Avinava Dubey, Chelsea Finn, et al. RT-2: 544
Vision-language-action models transfer web knowledge to 545
robotic control. In *Conference on Robot Learning*, 2023. 1 546
- [8] Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker- 547
Holder, Yuge Shi, Edward Hughes, et al. Genie: Generative 548
interactive environments. In *Proceedings of the 41st Inter-* 549
national Conference on Machine Learning, 2024. 2 550
- [9] Yevgen Chebotar, Quan Vuong, Karol Hausman, Fei Xia, 551
Yao Lu, Alex Irpan, Aviral Kumar, Tianhe Yu, Alexan- 552
der Herzog, Karl Pertsch, Keerthana Gopalakrishnan, Ju- 553
lian Ibarz, Sean Kirmani, Kanishka Rao, Paul Xu, Dmitry 554
Kalashnikov, Chelsea Finn, and Sergey Levine. Q- 555
Transformer: Scalable offline reinforcement learning via au- 556
toregressive Q-functions. In *Conference on Robot Learning*, 557
2023. 2 558
- [10] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric 559
Cousineau, Benjamin Burchfiel, and Shuran Song. Diffu- 560
sion policy: Visuomotor policy learning via action diffusion. 561
In *Robotics: Science and Systems*, 2023. 2 562
- [11] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and 563
Christopher Ré. FlashAttention: Fast and memory-efficient 564
exact attention with IO-awareness. In *Advances in Neural* 565
Information Processing Systems, pages 16344–16359, 2022. 566
2 567
- [12] Zibin Dong, Jianye Hao, Yifu Yuan, Fei Ni, Yitian Mu, Yan 568
Li, and Yujing Zheng. Diffuserlite: Towards real-time diffu- 569
sion planning. In *Advances in Neural Information Process-* 570
ing Systems, 2024. 2 571
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, 572
Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, 573
Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl- 574
vain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is 575
worth 16x16 words: Transformers for image recognition at 576
scale. In *International Conference on Learning Representa-* 577
tions, 2021. 2 578
- [14] Y Du et al. Learning universal policies via text-guided video 579
generation. In *Advances in Neural Information Processing* 580
Systems, 2023. 2 581

- 582 [15] Mohsen Fayyaz, Soroush Abbasi Koohpayegani, 640
583 Farnoush Rezaei Jafari, Sunando Sengupta, Hamid 641
584 Reza Vaezi Joze, Eric Sommerlade, Hamed Pirsiavash, and 642
585 Juergen Gall. ATS: Adaptive token sampling for efficient 643
586 vision transformers. In *European Conference on Computer 644*
587 *Vision*, 2022. 2 645
588 [16] Zipeng Fu, Tony Z. Zhao, and Chelsea Finn. Mobile 646
589 ALOHA: Learning bimanual mobile manipulation with low- 647
590 cost whole-body teleoperation. In *Conference on Robot 648*
591 *Learning*, 2024. 2 649
592 [17] Lionel Guéguen, Alex Sergeev, Ben Kadlec, Rosanne Liu, 650
593 and Jason Yosinski. Faster neural networks straight from 651
594 JPEG. In *Advances in Neural Information Processing Sys- 652*
595 *tems*, 2018. 3
596 [18] Moo Jin Kim, Karl Fang, Ted Xiao, et al. OpenVLA: An 640
597 open-source vision-language-action model. *arXiv preprint 641*
598 *arXiv:2406.09246*, 2024. 1 642
599 [19] Moo Jin Kim, Yihuai Gao, Tsung-Yi Lin, Yen-Chen Lin, 643
600 Yunhao Ge, Grace Lam, Percy Liang, Shuran Song, Ming- 644
601 Yu Liu, Chelsea Finn, and Jinwei Gu. Cosmos policy: Fine- 645
602 tuning video models for visuomotor control and planning. 646
603 *arXiv preprint arXiv:2601.16163*, 2026. 1, 2 647
604 [20] Junbang Liang, Pavel Tokmakov, Ruoshi Liu, Sruthi Sud- 648
605 hakar, Paarth Shah, Rares Ambrus, and Carl Vondrick. 649
606 Video generators are robot policies. *arXiv preprint 650*
607 *arXiv:2508.00795*, 2025. 2 651
608 [21] Youwei Liang, Chongjian Ge, Zhan Tong, Yang Song, Jian- 652
609 ping Wang, and Pengtao Xie. Not all patches are what you 640
610 need: Expediting vision transformers via token reorganiza- 641
611 tions. In *International Conference on Learning Representa- 642*
612 *tions*, 2022. 2 643
613 [22] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximil- 644
614 ian Nickel, and Matt Le. Flow matching for generative mod- 645
615 eling. *arXiv preprint arXiv:2210.02747*, 2022. 4 646
616 [23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng 647
617 Zhang, Stephen Lin, and Baining Guo. Swin transformer: 648
618 Hierarchical vision transformer using shifted windows. In 649
619 *Proceedings of the IEEE/CVF International Conference on 650*
620 *Computer Vision*, pages 10012–10022, 2021. 2 651
621 [24] Ajay Mandlekar, Soroush Nasiriany, Bowen Wen, Iretoiyo 652
622 Akinola, Yashraj Narang, Linxi Fan, Yuke Zhu, and Dieter 640
623 Fox. MimicGen: A data generation system for scalable robot 641
624 learning using human demonstrations. In *Proceedings of The 642*
625 *7th Conference on Robot Learning*, pages 1820–1864, 2023. 643
626 5 644
627 [25] Octo Model Team et al. Octo: An open-source generalist 645
628 robot policy. *arXiv preprint arXiv:2405.12213*, 2024. 1 646
629 [26] OpenAI. Video generation models as world simulators. *Ope- 647*
630 *nAI Technical Report*, 2024. 2 648
631 [27] William Peebles and Saining Xie. Scalable diffusion models 649
632 with transformers. In *Proceedings of the IEEE/CVF Inter- 650*
633 *national Conference on Computer Vision*, pages 4199–4209, 651
634 2023. 3 652
635 [28] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, 640
636 Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, 641
637 and Sergey Levine. FAST: Efficient action tokeniza- 642
638 tion for vision-language-action models. *arXiv preprint 643*
639 *arXiv:2501.09747*, 2025. 2, 3 644
- [29] Zequn Qin, Pengyi Zhang, Fei Wu, and Xi Li. FcaNet: Fre- 640
quency channel attention networks. In *Proceedings of the 641*
IEEE/CVF International Conference on Computer Vision, 642
pages 783–792, 2021. 3 643
- [30] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie 644
Zhou, and Cho-Jui Hsieh. DynamicViT: Efficient vision 645
transformers with dynamic token sparsification. In *Advances 646*
in Neural Information Processing Systems, pages 13937– 647
13949, 2021. 2 648
- [31] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea 649
Finn. Learning fine-grained bimanual manipulation with 650
low-cost hardware. In *Robotics: Science and Systems*, 2023. 651
2 652