# Task adaptation of Vision-Language-Action model: 1st Place Solution for the 2025 BEHAVIOR Challenge

Ilia Larchenko[*]
ilya.larchenko@gmail.com

Gleb Zaring[*]
zaringleb@gmail.com

Akash Karnatak[*]
akashkarnatak.pro@gmail.com

[*]Independent Researchers

December 9, 2025

## Abstract

We present a vision-action policy that won 1st place in the 2025 BEHAVIOR Challenge—a large-scale benchmark featuring 50 diverse long-horizon household tasks in photo-realistic simulation, requiring bimanual manipulation, navigation, and context-aware decision making.

Building on the Pi0.5 architecture, we introduce several innovations. Our primary contribution is **correlated noise for flow matching**, which improves training efficiency and enables **correlation-aware inpainting** for smooth action sequences. We also apply **learnable mixed-layer attention** and **System 2 stage tracking** for ambiguity resolution. Training employs **multi-sample flow matching** to reduce variance, while inference uses **action compression** and challenge-specific **correction rules**.

Our approach achieves 26% q-score across all 50 tasks on both public and private leaderboards.

## 1  Introduction

### 1.1  The BEHAVIOR Challenge

The BEHAVIOR Challenge presents a demanding benchmark for embodied AI: performing 50 diverse, long-horizon household tasks in photo-realistic simulation. Tasks range from simple (turning on radio) to complex multi-step activities (cooking a hotdog). The challenge requires:

- **Long-horizon execution**: 6.6 minutes on average per task, with the longest tasks averaging 14 minutes

- **Bimanual manipulation**: Coordinated use of two 7-DOF arms with parallel-jaw grippers

- **Mobile navigation**: Indoor navigation through cluttered environments

- **Multi-camera perception**: Processing RGB images from head and both wrist cameras

- **Task diversity**: 50 different activities evaluated with a single policy (or small set of checkpoints)

The benchmark uses OmniGibson simulation built on NVIDIA Isaac Sim, providing realistic physics and rendering. Each task is evaluated over 10 episodes with randomized initial conditions, and performance is measured by a q-score that combines success rate with partial credit for subtask completion.

### 1.2  Key Challenges

Long-horizon household manipulation poses several fundamental challenges:

- **Compounding errors.** With episodes spanning thousands of timesteps, small prediction errors can accumulate. This demands either extremely accurate predictions or robust recovery behaviors.

- **Non-Markovian states.** Many task states are visually ambiguous—the robot holding a radio at the start of a task looks identical to holding it at the end. Without memory of past actions or explicit stage tracking, the policy cannot distinguish these states and may execute incorrect actions.

- **No recovery demonstrations.** The training data consists of successful demonstrations only. When the robot deviates from demonstrated trajectories (inevitable given compounding errors), it encounters states never seen during training. The policy must somehow generalize to recover from these out-of-distribution situations.

- **Multi-modal action distributions.** Many states admit multiple valid action sequences (e.g., which hand to use, which object to grasp first). Different episodes of the same tasks were completed at different speeds in the training data.

### 1.3  Our Approach

We build upon Pi0.5 [3], a vision-language-action (VLA) model that uses flow matching to predict action sequences.

Our modifications address the challenges above through the following novel components:

- **Modeling action structure.** Robot actions exhibit strong correlations—temporally (smooth trajectories) and across dimensions (coordinated joint movements). We model this structure explicitly by training with correlated noise sampled from $\mathcal{N}(0, \beta\Sigma + (1 - \beta)\mathbf{I})$, where $\Sigma$ is the empirical action covariance and $\beta = 0.5$. This makes training more efficient and enables principled inpainting during inference.

- **Providing non-Markovian context.** We introduce System 2 stage tracking: the model predicts the current task stage, and a voting mechanism filters noisy predictions to maintain stable stage estimates. This stage information is fused with task embeddings and fed back to the model, resolving ambiguous states.

- **Combining learning with heuristics.** Pure learning struggles with the lack of recovery data. We complement the learned policy with correction rules derived from failure analysis: simple heuristics that detect and recover from common failure modes like accidental gripper closures.

- We apply **learnable mixed-layer attention** that allows each action expert layer to attend to learned linear combinations of all VLM layers rather than arbitrarily deciding how action expert layers should attend to VLM layers.

- For training, we employ **multi-sample flow matching** (15 predictions per VLM forward pass) to reduce gradient variance while amortizing expensive vision-language computations.

- At inference, we apply **action compression** via cubic splines to speed up action execution by $1.3\times$.

- We also simplified the VLM part by removing text processing and using **trainable task embeddings instead of the text prompt**. Technically this removes "L" from "VLA" and "VLM" terms but we keep the names "VLA" and "VLM" for simplicity.

### 1.4 A Note on Evidence

This report describes a competition entry, not a research paper. As a small independent team with limited compute, we prioritized winning over rigorous ablation studies. Many design choices were guided by intuition or quick experiments rather than systematic evaluation.

We present our methods honestly, without claiming that every component is necessary or optimal. That said, first place on the leaderboard suggests the overall approach is sound—even if we cannot isolate the contribution of each piece.

### 1.5 Code and Models

Our code for training and inference is available at `https://github.com/IliaLarchenko/behavior-1k-solution`.

We also share our model weights at `https://huggingface.co/IliaLarchenko/behavior_submission`.

This code and these weights should allow anyone to reproduce our results.

## 2 Related Work

Our work builds on vision-language-action models, flow matching, action chunking, and multi-task learning for robot control.

### 2.1 Vision-Language-Action Models

VLA models fine-tune pretrained vision-language models for robot control, leveraging semantic knowledge from web-scale pretraining. Early VLAs like RT-2 [6] use autoregressive token prediction for actions, which struggles with high-frequency control. Pi0 [2] (and later Pi0.5 [3]) uses flow matching for continuous action prediction, enabling dexterous manipulation at up to 50 Hz.

Using flow matching or diffusion as an action head is a popular choice in modern VLAs, but different architectures choose different ways to combine VLM and action head:

- **Pi0.5**: Layer-to-layer expert-like attention. Layer $j$ of action expert attends to layer $j$ of VLM in a fused self-attention manner.

- **Gr00t** [4]: Attends only to the last VLM layer via standard cross-attention, like in the classic encoder-decoder transformer architecture [16].

- **SmolVLA** [7]: Keeps only half of the VLM layers, alternating self and cross attention layers.

We introduce learnable linear combinations of all VLM layers, letting the model decide optimal attention patterns rather than fixing them architecturally.

### 2.2 Flow Matching and Diffusion for Actions

Flow matching [8] and diffusion [9] provide expressive action distributions that handle multi-modality better than regression. In Pi0.5, the flow-matching expert transforms normal noise into action chunks in multiple denoising steps.

Authors of Pi0 and Pi0.5 mention that first denoising steps are more difficult for the model to solve than the rest and introduced non-uniform sampling of flow-matching time during training to resolve this problem.

Our contribution in solving the same problem is modeling action correlations explicitly in the noise distribution. Standard flow matching uses $\epsilon \sim \mathcal{N}(0, \mathbf{I})$; we use $\epsilon \sim \mathcal{N}(0, \beta\Sigma + (1 - \beta)\mathbf{I})$ where $\Sigma$ is the empirical action covariance and $\beta = 0.5$. This makes training more efficient and enables correlation-guided inpainting.

## 2.3 Action Chunking and Temporal Consistency

ACT [12] introduced predicting action sequences (chunks) rather than single actions, improving temporal consistency and enabling open-loop execution. Real-time chunking methods [13] use rolling windows with inpainting to balance consistency with reactivity—executing part of a predicted chunk while using the remainder as initial conditions for the next prediction.

Our correlation-aware inpainting uses learned action correlations to maintain smooth transitions between prediction windows, rather than treating inpainted and free actions independently.

## 2.4 Multi-Task Learning for Robotics

Training on multiple tasks can improve generalization and enable transfer learning. RT-1 [5] and RT-2 [6] and other VLA papers [2, 3, 4, 7] demonstrate that large-scale multi-task training improves robustness and generalization in real-world robots.

Our multi-task training on 50 BEHAVIOR tasks shows similar benefits: the model learns recovery behaviors (e.g., grabbing fallen objects) that aren't explicitly present in single-task demonstrations. We hypothesize that training on diverse tasks exposes the model to a wide variety of states and transitions, enabling it to generalize recovery strategies across tasks.

## 2.5 Hierarchical Reasoning for Long-Horizon Tasks

Augmenting policies with high-level reasoning improves long-horizon performance. Many methods use separate models—a VLM for semantic subtask prediction and a low-level policy for execution. Hi Robot [11] uses the same model for both high-level and low-level inference in a chain-of-thought style.

Our System 2 stage tracking is simpler: we predict task stages as an auxiliary output and use voting logic to filter noisy predictions. This provides non-Markovian context without explicit hierarchical planning or separate inference passes.

## 2.6 Action Space Design

**Delta actions** predict changes from current state rather than absolute positions, providing invariance to initial configuration, better generalization across starting states, and easier learning of smooth trajectories.

**Per-timestamp normalization** [17] addresses the fact that action distributions change over the prediction horizon—early actions in a chunk have small deltas while later actions are more varied. Normalizing per-timestamp makes the learning problem more uniform.

## 2.7 BEHAVIOR Benchmark

BEHAVIOR-1K [1] is a benchmark featuring 1,000 everyday household activities in photo-realistic simulation (OmniGibson/Isaac Sim). The 2025 challenge includes 50 tasks with long-horizon execution (average 6.6 minutes, up to 14 minutes), complex object interactions, diverse activities (rearrangement, cooking, cleaning), and 10,000 expert demonstrations via teleoperation.

# 3 Method Overview

We organize our contributions into three main parts: **model architecture**, **training procedures**, and **inference optimizations**.

## 3.1 High-Level Architecture

Our model processes multi-camera RGB observations along with robot proprioceptive state to predict action sequences using flow matching (see Figure 1). The architecture consists of:

1. **Vision Backbone**: SigLIP-So400m/14 [14] encoder processes images from three cameras (head, left wrist, right wrist)

2. **Language Model Backbone**: PaliGemma-based [15] transformer processes visual tokens, task and stage information and robot state.

3. **Action Expert**: Separate transformer that predicts action sequences via flow matching

See Appendix B for more details.

## 3.2 Our method description structure

**Part 1: Model Architecture (Section 4–5)** We modify the Pi0.5 architecture in several ways:

1. **Task Embeddings**: Replace language model processing with 50 trainable task embeddings (one per task)

2. **KV Cache Transformation**: Learn linear combinations of VLM layer outputs for flexible mixed-layer attention

3. **System 2 Stage Tracking**: Predict and track task stages to provide non-Markovian context

4. **Custom Attention Masks**: Hierarchical attention pattern isolating reliable inputs from noisy ones

5. **Correlated Noise Generation**: Sample noise from $\mathcal{N}(0, \beta\Sigma + (1-\beta)\mathbf{I})$ instead of independent noise

**Part 2: Training Procedures (Section 6)** Our training innovations focus on reducing variance and improving action space learning:

1. **Multi-Sample Flow Matching**: Compute 15 flow predictions per VLM forward pass with different noise/time samples

2. **Delta Action Space**: Predict action deltas with per-timestamp normalization

3. **Multi-Task Training**: Joint training on all 50 tasks with stage prediction as auxiliary task

4. **Loss Composition**: Weighted combination of action loss, stage prediction loss, and FAST auxiliary loss
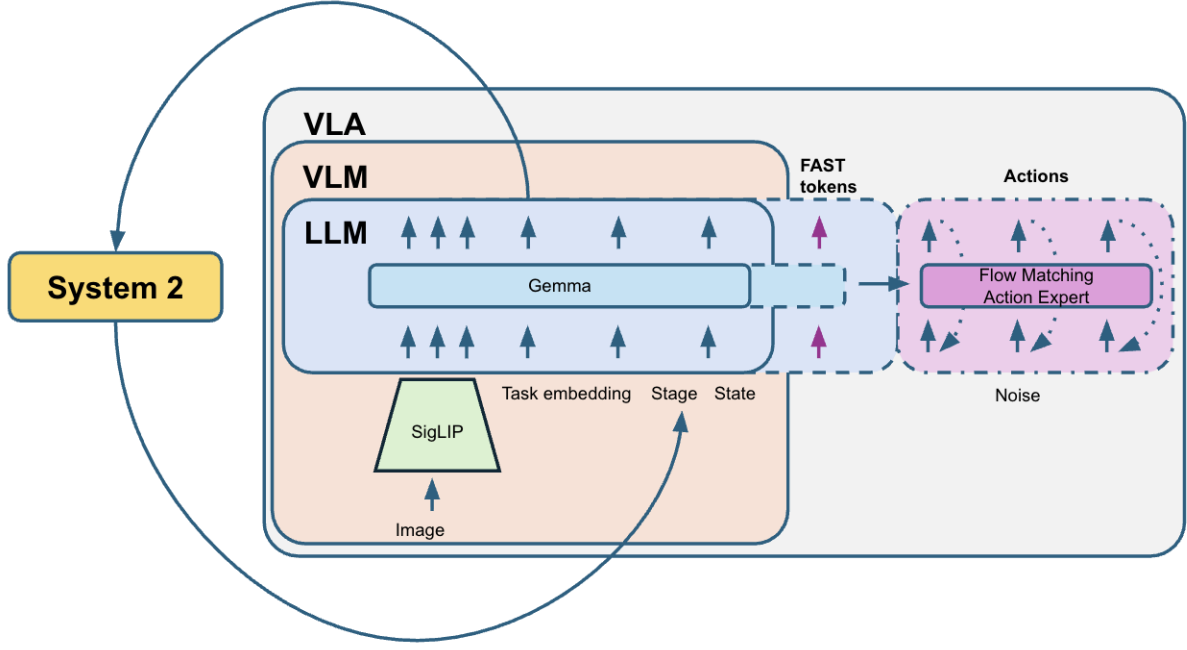
**Figure 1:** Overall architecture of our system. The vision backbone (SigLIP) processes images from three cameras, PaliGemma VLM processes visual tokens along with task embeddings and robot state, and the Action Expert predicts actions via flow matching. Task embeddings replace language processing, and System 2 provides stage information for non-Markovian context.

**Part 3: Inference Optimizations (Section 7)** We introduce several inference-time techniques to improve action quality and efficiency:

1. **Soft Correlation-Aware Inpainting**: Guide free dimensions during inpainting to maintain action correlation

2. **Action Compression**: Interpolate 26 predicted actions to 20 execution steps using cubic splines (1.3× speedup)

3. **Stage Voting**: Track stage transitions using majority voting over prediction history

4. **Correction Rules**: Task-specific fixes for common failure modes (e.g., gripper recovery)

### 3.3 Mathematical Notation

Throughout this report, we use the following notation:

- $\mathbf{x}_t \in \mathbb{R}^{H \times D}$: Action sequence at flow time $t$, where $H = 30$ is the action horizon and $D = 23$ is the action dimension

- $\mathbf{a} \in \mathbb{R}^{H \times D}$: Target action sequence (ground truth)

- $\mathbf{v}_t \in \mathbb{R}^{H \times D}$: Predicted velocity (flow direction) at flow time $t$

- $\boldsymbol{\epsilon} \in \mathbb{R}^{H \times D}$: Noise vector

- $\boldsymbol{\Sigma} \in \mathbb{R}^{HD \times HD}$: Action covariance matrix (equivalent to correlation matrix since actions are normalized to zero mean and unit variance)

- $t \in [0, 1]$: Flow matching time parameter ($t = 1$ is pure noise, $t = 0$ is target)

- $\ell$: Timestep index in a trajectory (distinct from flow time $t$)

- $\tau \in \{0, 1, \ldots, 49\}$: Task ID

- $s \in \{0, 1, \ldots, s_{\max}(\tau)\}$: Task stage/substage

- $\mathbf{K}, \mathbf{V}$: Key and value caches from VLM layers

- $\mathcal{O}, \mathcal{U}$: Sets of inpainted (observed) and free (unobserved) action dimensions

## 4 Model Architecture

This section describes the architectural modifications we made to the Pi0.5 base model to adapt it for BEHAVIOR-1K tasks.

### 4.1 Task Embeddings

The original Pi0.5 uses language embeddings for task specification. We replace this with task-specific embeddings for the structured BEHAVIOR-1K tasks.

The BEHAVIOR challenge requires very limited generalization. There are only 50 tasks that are present in both training and evaluation data. This means there is no explicit requirement for the policy to generalize to new tasks described in natural language.

Instead of processing natural language prompts, we use trainable task embeddings—one 2048-dimensional embedding for each of 50 tasks, trained from scratch.
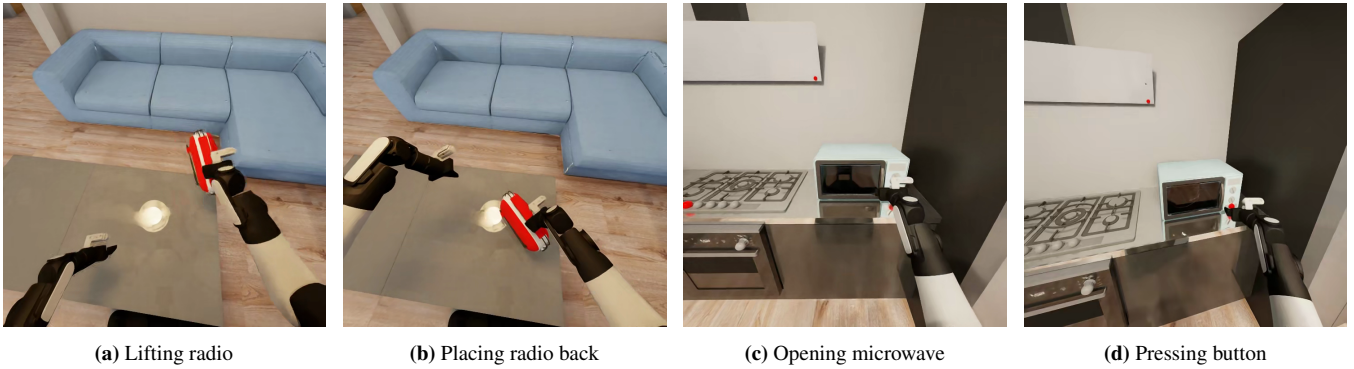
**(a)** Lifting radio              **(b)** Placing radio back            **(c)** Opening microwave            **(d)** Pressing button

**Figure 2:** Examples of visually similar but semantically different states. (a,b) Radio task: lifting vs. placing back. (c,d) Popcorn task: opening microwave vs. pressing button after loading. Without stage tracking, the model confuses these states.

This simplification is justified for BEHAVIOR-1K since:

1. Only 50 distinct tasks (fixed set)

2. Task semantics are implicit in demonstrations

3. Removes language model processing overhead

4. Allows the model to learn task-specific features directly

## 4.2 System 2: Stage Prediction and Fusion

One of the big challenges that we faced across different problems is the existence of non-Markovian states. This means that the current state of the task is not enough to predict the correct next action. For example, the robot can see nearly identical images at the beginning and end of the same task (Figure 2).

This often confuses the model, which then acts incorrectly. To fix this issue we added a simple System 2 that predicts the current stage of the task based on the images and task embedding, applies voting logic to filter out incorrect predictions, and uses it as additional input to the model on following steps.

### 4.2.1 Stage Prediction

Each task is divided into 5–15 stages based on demonstration length. We predict the current stage using a linear classifier on the VLM output:

$$\mathbf{logits} = \mathbf{W}_{\text{stage}} \cdot \text{VLM}(\text{images}, \text{task\_id}) + \mathbf{b}_{\text{stage}} \quad (1)$$

where $\mathbf{W}_{\text{stage}} \in \mathbb{R}^{15 \times 2048}$ (15 is the maximum number of stages across all tasks). Invalid stages for each task are masked with $-\infty$ before softmax.

Stage prediction achieves $\sim$99% accuracy on training data, providing reliable context for action prediction.

### 4.2.2 Stage-Task Fusion

We fuse the task embedding with stage information using multiple learned representations (sine-cosine encoding, task-specific learned embeddings, and gated combinations). This provides 5 task-related tokens to the model. Details are in Appendix A.

## 4.3 KV Cache Transformation for Mixed-Layer Attention

We don't see a clear winner among different ways how the flow-matching/diffusion action head attends to the VLM part in different VLA models, so we decided to let the model decide which layers to attend to and how.

### 4.3.1 Learnable Layer Combination

During both training and inference we first compute Key-Value cache for all layers of the VLM part. Then we transform it using learnable weights and biases.

For each action expert layer $j$, we compute transformed keys and values as linear combinations of all VLM layers:

$$\mathbf{K}_j^{\text{new}} = \sum_{i=1}^{L} w_{ij}^{(K)} \mathbf{K}_i + \mathbf{b}_j^{(K)} \quad (2)$$

$$\mathbf{V}_j^{\text{new}} = \sum_{i=1}^{L} w_{ij}^{(V)} \mathbf{V}_i + \mathbf{b}_j^{(V)} \quad (3)$$
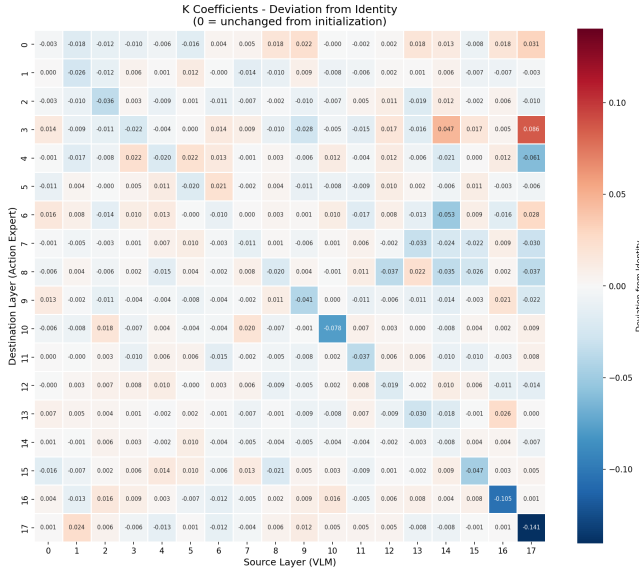
where $\mathbf{K}_i, \mathbf{V}_i \in \mathbb{R}^{B \times S \times N_h \times d_h}$ are key/value caches from VLM layer $i$ (with batch size $B$, sequence length $S$, $N_h$ attention heads, head dimension $d_h$), $w_{ij}^{(K)}, w_{ij}^{(V)} \in \mathbb{R}$ are learnable scalar weights, $\mathbf{b}_j^{(K)}, \mathbf{b}_j^{(V)} \in \mathbb{R}^{N_h \times d_h}$ are learnable bias terms, and $L$ is the number of VLM layers.

**Initialization**: $w_{ij} = \delta_{ij}$ (identity), $\mathbf{b}_j = \mathbf{0}$, so the model starts with Pi0.5's layer-to-layer attention.
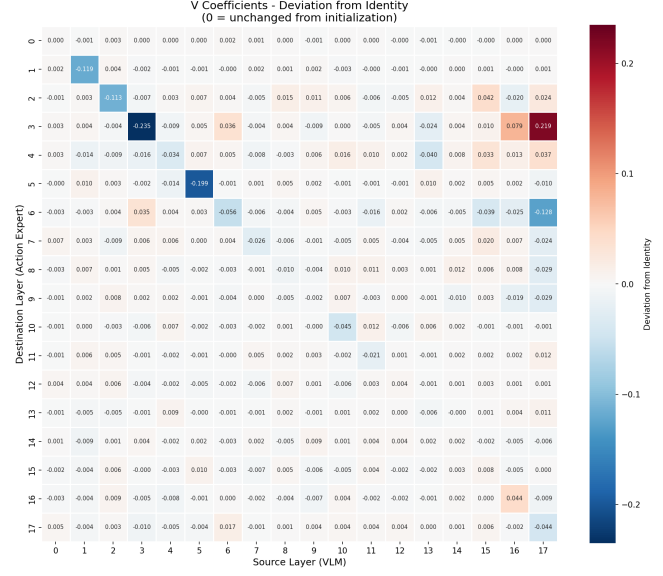
### 4.3.2 Properties

This design allows the model to: (1) attend to any VLM layer—weights $w_{ij}$ can select early, middle, or late layers; (2) form smooth combinations by attending to weighted averages of multiple layers; (3) learn from data without manual architecture search. The approach is parameter-efficient: for each of 18 action expert layers, we learn only 18 scalar coefficients plus one bias tensor $\in \mathbb{R}^{N_h \times d_h}$, separately for keys and values.

We use different coefficients for keys and values, as there is no reason for them to be identical. Figure 3 shows the learned deviations from identity initialization. Since we initialized from pretrained Pi0.5 weights after significant fine-tuning, the

**(a)** Key coefficients deviation from identity



**(b)** Value coefficients deviation from identity

**Figure 3:** Learned KV transformation weights showing deviation from initialization (identity). Each row corresponds to an action expert layer, each column to a VLM layer. Brighter colors indicate higher deviation.

identity initialization was already a good prior. We observe some tendency to attend more to the last VLM layer, though this could be noise. We expect this approach to have larger effects for models trained from scratch or initialized from non-robotics VLMs.

### 4.4 Custom Attention Masks



**Figure 4:** Attention mask structure during training. Blue indicates bidirectional attention, white indicates no attention, and diagonal patterns indicate causal attention.

We use a hierarchical attention pattern that isolates reliable inputs from noisy ones:

1. **Image tokens**: Bidirectional attention among themselves and task token

2. **Task token**: Bidirectional with images

3. **Stage tokens**: Attend to images, task, and proprio state

4. **State tokens**: Attend to images, task, stage, and other state tokens

5. **FAST tokens** (if used): Attend to all prefix tokens and causally to each other

6. **Action expert tokens**: Bidirectional among themselves, attend to all other prefix tokens (except FAST)

**Rationale**:

1. Images and task embeddings are the most reliable inputs (deterministic from observations). We prevent them from attending to noisier inputs like robot state (which can drift during inference) or predicted stage (which could be incorrect). This keeps visual features clean.

2. Only images and task embeddings are used in System 2 to predict the current stage; to avoid target leakage, they should not attend to stage tokens.

3. FAST tokens predict actions autoregressively, so they attend to all prefix tokens and causally to each other (used only during training).

4. Action expert tokens predict the whole chunk simultaneously, so they use bidirectional attention among themselves while attending to all prefix tokens except FAST. See Figure 4.

## 4.5 Delta Action Space with Per-Timestamp Normalization

### 4.5.1 Delta Actions

Instead of predicting absolute joint positions, we predict **delta actions**:

$$\mathbf{a}_{\text{delta}}[\ell + i] = \mathbf{a}_{\text{abs}}[\ell + i] - \mathbf{q}_{\text{current}}[\ell] \qquad (4)$$

where $\mathbf{q}$ represents joint positions, $i \in \{0, \ldots, H-1\}$ is the index within the action chunk, and $\ell$ is the current timestep.

### 4.5.2 Per-Timestamp Normalization

For each action dimension $d$ and index in the chunk $i$:

$$\mu_{di} = \text{mean}(\mathbf{a}_{\text{delta}}[i, d]), \quad \sigma_{di} = \text{std}(\mathbf{a}_{\text{delta}}[i, d]) \qquad (5)$$

Normalized actions:

$$\tilde{\mathbf{a}}[i, d] = \frac{\mathbf{a}_{\text{delta}}[i, d] - \mu_{di}}{\sigma_{di}} \qquad (6)$$

**Why per-timestamp?** Action distributions change over time within a trajectory: initial actions in the chunk are very close to the current state (small delta), while later actions are more varied. Per-timestamp normalization makes the learning problem more uniform across the horizon.

**Note:** velocities and gripper positions are excluded from per-timestamp normalization.

## 4.6 Action Expert Architecture

The action expert uses the Gemma 300M architecture ($\sim$300M parameters, 18 layers), identical to Pi0.5. It attends to the transformed VLM KV cache and outputs velocity predictions $\mathbf{v}_t = \text{ActionExpert}(\mathbf{x}_t, t, \text{KV}_{\text{VLM}})$. The training loss is $\mathcal{L}_{\text{action}} = \mathbb{E}_{t,\epsilon}[\|\mathbf{v}_t - (\epsilon - \mathbf{a})\|^2]$. See Appendix C for architecture details.

## 4.7 FAST Auxiliary Training

We implemented FAST [10] auxiliary training with loss weight 0.05. Since we removed text input, we directly train FAST token embeddings mapping to PaliGemma input. FAST uses global quantile normalization for better DCT compression because per-timestamp normalization breaks temporal smoothness of actions. See Appendix D for details.

## 5 Correlated Noise for Flow Matching

One of our key innovations is modeling action correlations explicitly during flow matching training and inference.

### 5.1 Motivation: Action Correlations in Robot Control

Robot actions exhibit strong correlations in two ways (see Figure 5):

1. **Temporal correlation**: Actions at nearby timesteps are similar (smooth trajectories)

2. **Cross-dimensional correlation**: Joint velocities are co-ordinated (e.g., trunk joints move together)
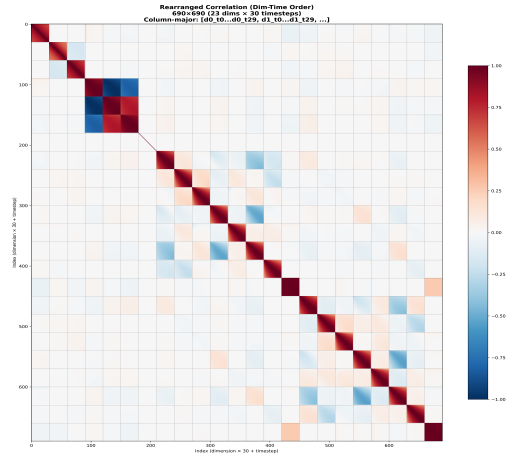


**Figure 5:** Action correlation matrix $\mathbf{\Sigma} \in \mathbb{R}^{690 \times 690}$ computed from training data (30 timesteps $\times$ 23 action dimensions = 690). The matrix shows strong block-diagonal structure indicating high temporal correlation and significant cross-dimensional correlation.

Standard flow matching uses independent Gaussian noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$. This creates a training problem: the early denoising steps ($t \approx 1$) are difficult, but once a few denoising steps are made and the model learns correlation structure, later prediction becomes trivial.

**Our solution**: Generate correlated noise $\epsilon \sim \mathcal{N}(0, \mathbf{\Sigma})$ that already matches the action structure. This makes all denoising steps more balanced in difficulty.

### 5.2 Correlation Matrix Estimation

We estimate the full correlation matrix from the training set. Let $\mathbf{a}^{(n)} \in \mathbb{R}^{H \times D}$ be the $n$-th normalized action sequence (flattened to $\mathbb{R}^{HD}$). The sample covariance is:

$$\hat{\mathbf{\Sigma}} = \frac{1}{N} \sum_{n=1}^{N} \text{vec}(\mathbf{a}^{(n)}) \cdot \text{vec}(\mathbf{a}^{(n)})^T \qquad (7)$$

where $\text{vec}(\cdot)$ flattens the $H \times D$ matrix to a $HD$-dimensional vector.

### 5.3 Beta Shrinkage for Robustness

Using pure $\mathbf{\Sigma}$ can be unstable. We apply **shrinkage regularization**:

$$\mathbf{\Sigma}_{\text{reg}} = \beta \mathbf{\Sigma} + (1 - \beta)\mathbf{I} \qquad (8)$$

where $\beta \in [0, 1]$ is the shrinkage parameter. We use $\beta = 0.5$ as a balanced choice.

### 5.4 Correlated Noise Generation

To sample $\epsilon \sim \mathcal{N}(0, \mathbf{\Sigma}_{\text{reg}})$, we use the Cholesky decomposition:

$$\mathbf{\Sigma}_{\text{reg}} = \mathbf{L}\mathbf{L}^T \qquad (9)$$

where $\mathbf{L}$ is lower triangular. Then:

$$\epsilon = \mathbf{L}\mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}) \qquad (10)$$

### 5.5 Impact on Training

With correlated noise, the flow matching interpolation becomes:

$$\mathbf{x}_t = t \cdot \boldsymbol{\epsilon} + (1 - t) \cdot \mathbf{a}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{\Sigma}_{\text{reg}}) \qquad (11)$$

**At $t = 1$ (pure noise)**: $\mathbf{x}_1 = \boldsymbol{\epsilon}$ has the same correlation structure as real actions. The model sees plausible action patterns even at the noisiest step.

**At $t \in (0, 1)$**: The interpolation maintains correlation structure throughout the denoising process.

This makes the training task more uniform and keeps difficulties of different denoising steps more balanced.

## 6 Training

### 6.1 Training Data

We train on the BEHAVIOR-1K demonstration dataset: 10,000 expert demonstrations across 50 tasks (1200+ hours total). Each demonstration is divided into 5–15 stages based on temporal position for stage prediction training.

### 6.2 Multi-Sample Flow Matching

Standard flow matching computes one action prediction per observation, sampling $(t, \boldsymbol{\epsilon})$ randomly for each batch element. This introduces significant variance in the training signal.

#### 6.2.1 Motivation

The flow matching loss has two sources of randomness:

1. **Time sampling**: $t \sim \text{Beta}(1.5, 1)$

2. **Noise sampling**: $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{\Sigma}_{\text{reg}})$

We can amortize the expensive VLM forward pass across multiple flow samples and reduce the randomness of the resulting gradient.

#### 6.2.2 Algorithm

For each training batch:

1. **VLM forward pass** (once): Compute KV cache for all prefix tokens

2. **Multi-sample action prediction** ($N = 15$ times): Sample different $(t_n, \boldsymbol{\epsilon}_n)$ for each sample, compute noisy actions, run action expert

3. **Backward pass**: Gradients flow back through all $N$ samples

### 6.3 Loss Function

Our total loss is a weighted combination of three components:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{action}} + \lambda_s \mathcal{L}_{\text{stage}} + \lambda_f \mathcal{L}_{\text{FAST}} \qquad (12)$$

#### 6.3.1 Action Loss

The flow matching loss averaged over $N$ samples:

$$\mathcal{L}_{\text{action}} = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{HD} \sum_{h=1}^{H} \sum_{d=1}^{D} (v_{t_n}[h, d] - u_{t_n}[h, d])^2 \qquad (13)$$

#### 6.3.2 Stage Prediction Loss

Cross-entropy loss for stage classification with weight $\lambda_s = 0.1$.

#### 6.3.3 FAST Auxiliary Loss

Weight $\lambda_f = 0.05$ (reduced from 0.1 after initial training).

### 6.4 Multi-Task vs. Task-Specific Training

We train models in two steps:

**1. Multi-task (initial)**: Train on all 50 tasks simultaneously. This phase took 15 days of non-stop training on $8 \times$H200 GPUs.

**2. Task-group-specific (fine-tuning)**: We split tasks based on validation results: best (highest success rate), good (score $> 0$), bad (score $\sim 0$) and trained them separately. This phase took $\sim$1 week per group.

The final submission uses 4 task-specific checkpoints, automatically switching based on task ID.

### 6.5 Training Budget

Our total competition budget was $\sim$\$13k. We spent $\sim$\$3k of personal money on experiments and evaluation. \$10k was sponsored by Nebius and we used it for the main training runs on $8 \times$H200 GPUs.

## 7 Inference Optimizations

### 7.1 Correlation-Aware Inpainting

To ensure smooth action sequences and resolve local multi-modality, we use a rolling inpainting strategy.

#### 7.1.1 Soft Inpainting Strategy

Instead of executing all predicted actions directly we:

1. **Predict** 30 actions: $\{\mathbf{a}_0, \mathbf{a}_1, \ldots, \mathbf{a}_{29}\}$

2. **Execute** first 26 actions: $\{\mathbf{a}_0, \ldots, \mathbf{a}_{25}\}$

3. **Save** last 4 actions: $\{\mathbf{a}_{26}, \ldots, \mathbf{a}_{29}\}$ as initial conditions

4. **Next prediction**: Generate 30 new actions such that first 4 almost match the saved actions
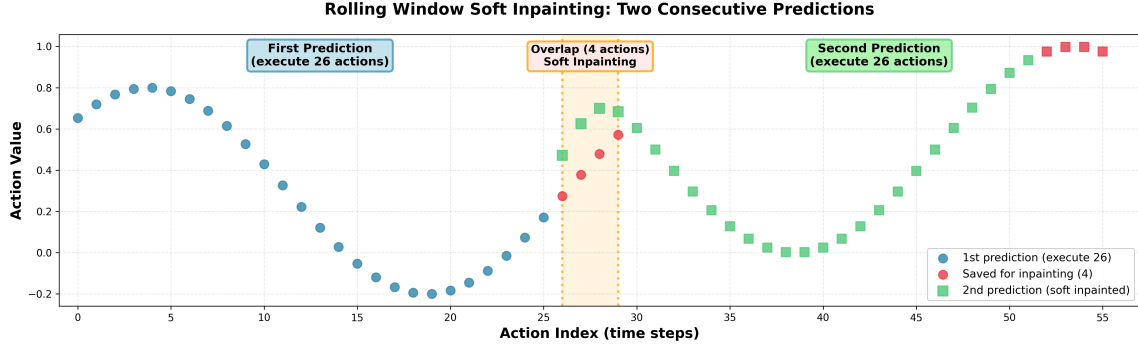
5. **Repeat**

**Figure 6:** Rolling window soft inpainting across two consecutive predictions. The first prediction generates 30 actions (blue circles), executes 26, and saves the last 4 (red circles) for inpainting. The second prediction (green squares) starts close to the saved actions through soft inpainting.
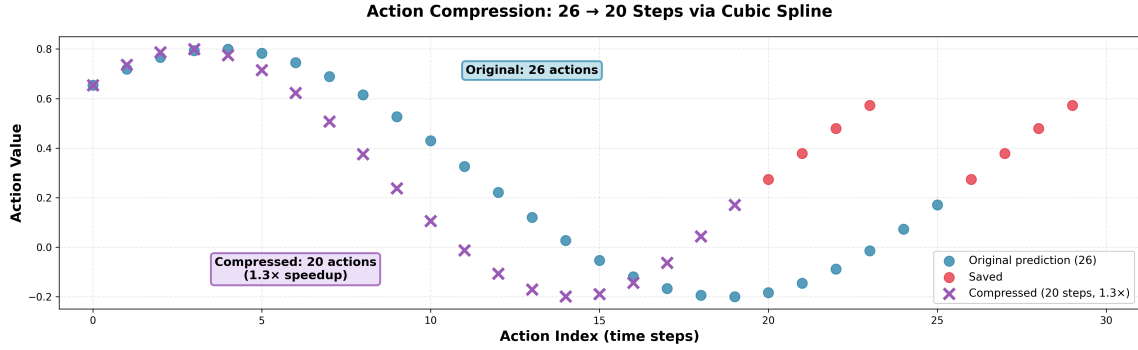


**Figure 7:** Action compression via cubic spline interpolation. The original 26 actions (blue circles) are compressed to 20 steps (purple crosses), achieving $1.3\times$ speedup.

### 7.1.2 *Correlation-Aware Correction Propagation*

The key challenge is: how to constrain the initial actions while respecting the correlation structure? A naive approach applies hard constraints on the first 4 actions with no adjustment for the rest. This creates discontinuity at the boundary between actions 4 and 5—the model predictions don't respect correlation between inpainted and free actions, and the input at every flow-matching step becomes out-of-distribution.

Our approach propagates corrections using the learned correlation structure. Let $\mathcal{O}$ = indices of inpainted actions (first 4 timesteps, $|\mathcal{O}| = 4 \times 23 = 92$ dimensions), $\mathcal{U}$ = indices of free actions ($|\mathcal{U}| = 26 \times 23 = 598$ dimensions). Partition the correlation matrix:

$$\boldsymbol{\Sigma}_{\mathrm{reg}} = \begin{bmatrix} \boldsymbol{\Sigma}_{\mathcal{OO}} & \boldsymbol{\Sigma}_{\mathcal{OU}} \\ \boldsymbol{\Sigma}_{\mathcal{UO}} & \boldsymbol{\Sigma}_{\mathcal{UU}} \end{bmatrix} \tag{14}$$

At each denoising step with time $t$, after the model predicts $\mathbf{x}_t^{\mathrm{pred}}$:

1. Compute desired state on inpainted dimensions using saved initial actions $\mathbf{x}_0^{\mathcal{O}}$ and fixed noise $\mathbf{z}^{\mathcal{O}}$: $\mathbf{x}_t^{\mathrm{desired},\mathcal{O}} = (1-t)\mathbf{x}_0^{\mathcal{O}} + t\mathbf{z}^{\mathcal{O}}$

2. Compute additive correction: $\boldsymbol{\delta}^{\mathcal{O}} = \mathbf{x}_t^{\mathrm{desired},\mathcal{O}} - \mathbf{x}_t^{\mathrm{pred},\mathcal{O}}$

3. Apply hard constraint on inpainted dimensions

4. Propagate correction to free dimensions: $\boldsymbol{\delta}^{\mathcal{U}} = \mathbf{M}_{\mathrm{corr}}\boldsymbol{\delta}^{\mathcal{O}}$ where $\mathbf{M}_{\mathrm{corr}} = \boldsymbol{\Sigma}_{\mathcal{UO}}\boldsymbol{\Sigma}_{\mathcal{OO}}^{-1}$

5. Apply the correlated correction: $\mathbf{x}_t^{\mathcal{U}} = \mathbf{x}_t^{\mathrm{pred},\mathcal{U}} + \boldsymbol{\delta}^{\mathcal{U}}$

The correction matrix $\mathbf{M}_{\mathrm{corr}}$ encodes how much each free dimension should adjust given a correction on inpainted dimensions, ensuring smooth transitions. We precompute this matrix once using a numerically stable solver. See Figure 6.

### 7.1.3 *Soft Inpainting via Time Threshold*

We apply inpainting correction only during early denoising steps ($t > 0.3$). At early $t$, maintaining constraints and correlation is crucial for smooth trajectories. At late $t$ (close to the target), the model should have full freedom to adapt to current observations. This "soft" inpainting allows deviation from the initial plan if observations change.

## 7.2 Action Compression via Cubic Spline Interpolation

Speeding up action execution relative to demonstration data can improve task completion rates by allowing more prediction cycles per episode and more attempts to recover from failures. This idea appeared in prior work [18] ($2\times$ speedup) and Figure AI's "sport mode" [19] (linear resampling of action chunks). We apply a similar principle with cubic spline interpolation (see Figure 7):

- **Predicted**: 26 actions at 30Hz

- **Executed**: 20 steps at 30Hz

- **Speedup**: $26/20 = 1.3\times$

We use cubic spline interpolation to generate smooth intermediate actions rather than linear resampling, which can introduce jitter.

**Velocity scaling**: We scale base velocity dimensions by $1.3\times$ to account for faster execution: $\mathbf{a}_{scaled}[0 : 3] = \mathbf{a}_{interpolated}[0 : 3] \times 1.3$. Joint velocities (arms, trunk) remain unchanged since they are already normalized and the controller handles timing.

**Adaptive compression**: We disable compression when gripper states change significantly. Many failures are related to grasping, so we slow down and give the policy more time when the robot tries to grasp an object.

### 7.3 Stage Tracking with Voting Logic

The model predicts the current stage at each inference step. Since individual predictions can be noisy (see Figure 8), we employ a majority voting scheme to ensure stable stage transitions.

We maintain a sliding window of the three most recent stage predictions. At each inference step, the model outputs stage logits, and we take the argmax to obtain the predicted stage $\hat{s}$. This prediction is appended to the history buffer. Stage transitions follow three rules based on the prediction history:

**Forward transition**: If at least two of three predictions indicate the next stage $s_{curr} + 1$, we advance to that stage and clear history. This majority voting prevents premature transitions from single noisy predictions while still allowing responsive progression.

**Skip detection**: If all three predictions consistently indicate stage $s_{curr} + 2$, this suggests the robot has completed a stage faster than expected or the stage was already satisfied. We advance by one stage to catch up, then clear history.

**Rollback**: If all three predictions unanimously indicate the previous stage $s_{curr} - 1$, we roll back by one stage. This handles cases where a subtask needs to be re-attempted. Requiring unanimous agreement makes rollback more conservative than forward transitions.

After any stage transition, the prediction history is cleared to prevent stale predictions from influencing future transitions.

### 7.4 Correction Rules

#### 7.4.1 General Gripper Correction

BEHAVIOR-1K dataset is very clean and doesn't contain recovery demonstrations. In practice, if the policy fails any action, there is a very high chance of ending up in an out-of-distribution state and getting completely stuck.

One of the most common failures across all tasks is to fail a grasp and close the gripper on empty air. There are almost no training data where the robot opens the gripper after closing

it—this leads to complete failure as the robot is stuck and can't perform any actions.

**Note**: After training on 50 tasks, we observed that recovery behavior actually emerged in some cases—the robot learned to open the gripper even without explicit training data for it. However, these cases were rare and didn't improve the total success rate substantially.

To resolve this problem, we implemented a simple rule: if the gripper is closed and it was never closed in training data for this particular task at the same stage, we treat it as a failed grasp and completely open the gripper.

**Impact**: This correction rule alone approximately **doubled our success rate** in selected tasks where grabbing objects was a common failure mode.

#### 7.4.2 Task-Specific Rules

We started experimenting with more fine-grained task-specific rules but eventually added only one simple rule for the "turning_on_radio" task (roll back by 2 stages if we reach the final stage but didn't succeed). We believe this approach has potential for improving results but is not scalable or generalizable, so we didn't pursue it further.

## 8 Evaluation

### 8.1 Evaluation Setup

#### 8.1.1 Challenge Protocol

The BEHAVIOR-1K challenge uses a standardized evaluation protocol:

- **Tasks**: 50 household activities (same as in the demo dataset)

- **Episodes per task**: 10 evaluation episodes (fixed instances with randomized initial conditions)

- **Success metric**: Goal condition satisfaction (binary and partial)

- **Time limit**: Task-specific ($2\times$ average human task completion time in the demo dataset)

#### 8.1.2 Evaluation Instances

Each task has 10 pre-defined evaluation instances with randomized object positions (task-relevant objects) and robot initial pose.

Public leaderboard scores use 10 published instances; private scores are computed on an additional held-out set of 10 instances per task known only to the organizers.

Evaluation instances are distinct from training demonstrations (different seeds but drawn from the same distribution), testing generalization to new spatial configurations. The required sequence of actions stays mostly the same as in the demo dataset; there is no test of unseen object categories, language goal understanding, or entirely new tasks.

**Note**: During train data collection there was an unintentional bias toward simpler instances. The public and private
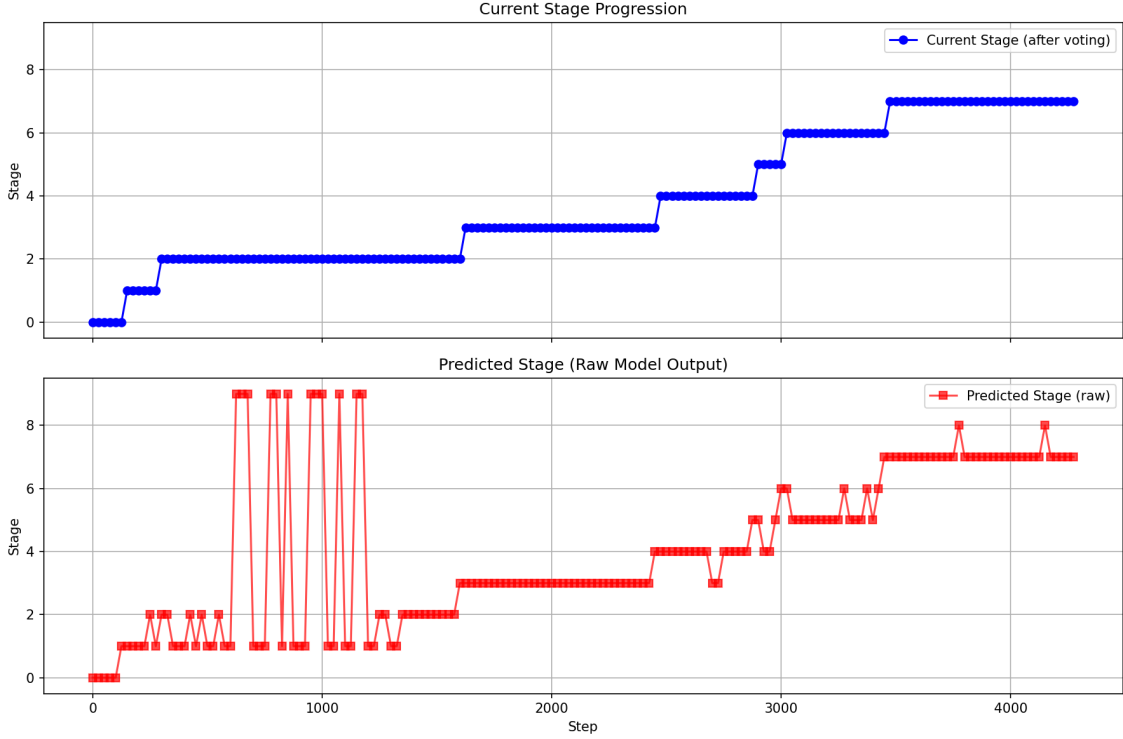
**Figure 8:** Stage tracking with voting logic resolves non-Markovian states. **Top**: Current stage after voting shows smooth, stable progression. **Bottom**: Raw model predictions are noisy and sometimes jump ahead (e.g., predicting stage 9 when still in early stages). In the radio task example, the robot grasping the radio at the beginning (stage 1-2) looks visually identical to placing it back later (stage 9), but System 2 correctly maintains stage 2 context despite model confusion.

evaluation instances do not share this bias and are on average more challenging.

## 8.2 Metrics

### 8.2.1 Primary Metric: Averaged Partial Success

Each task has a list of goal conditions. Partial success is the fraction of those conditions satisfied at the end of an episode—how far the robot progressed toward full completion.

- **Episode-level partial success**: fraction of goal conditions satisfied in that episode

- **Task-level score**: average episode partial success across the 10 evaluation episodes for that task

- **Overall score**: average task-level partial success across the 50 tasks (q_score)

Binary success (all goal conditions satisfied) is measured and reported for analysis, but it is **not** the ranking metric.

### 8.2.2 Secondary Metrics: Efficiency

The challenge also tracks efficiency metrics (not used for ranking):

- **Simulated time**: Total simulation steps × time per step

- **Distance navigated**: Cumulative base movement distance

- **Hand displacement**: Cumulative gripper movement distance

## 8.3 Evaluation Infrastructure

Running all 500 episodes (tasks × instances) on a single machine would take ~10 days, so we parallelized.

Parallelization options considered:

- Multiple environments in one process (experimental in OmniGibson; avoided to keep the official eval script unchanged)

- Many simulator instances on one machine (impractical due to high RAM use and a memory leak)

- Multiple machines (chosen): elastic access to inexpensive, RTX-capable GPUs with stable setup; we used RunPod and scaled to 20 RTX 4090s in parallel

As an independent team without existing infrastructure, we built a lightweight infrastructure that fit our budget and kept iteration speed high.

Resulting evaluation infrastructure:

- On-demand cluster to avoid paying for idle compute; full evaluation finishes in under 2 days

- Heavy assets (simulation code, model checkpoints) on a shared network volume. Using a default Docker container allows us to spin up new machines in minutes

11

- Evaluation outputs (videos, metrics, logs) automatically synced to an S3 bucket (we chose Cloudflare R2) with a simple UI for browsing models' evaluation results

- A lightweight offline balancer spreads tasks evenly and reruns episodes that crash (e.g., simulator segfaults)

## 9 Results

### 9.1 Competition Results

Top 5 teams evaluated on held-out instances (see Table 1).

We achieved a 26% q_score with negligible difference between public and private evaluation. In our case, partial successes contribute roughly half of the total score (see Figure 9).

### 9.2 Analysis

#### 9.2.1 General Observations

- Some tasks are almost solved, except under particularly tricky initial conditions

- For tasks with 0 success, we do not observe that they are generally impossible; instead, they usually contain one tricky step that involves very high-precision manipulation (with low success rate even for human teleoperators) or a carefully followed sequence that is slightly beyond the current model's limits

- Task duration does not appear to be a fundamental obstacle: longer tasks simply have many more steps, which makes full success harder, but partial success remains very achievable

#### 9.2.2 Failure Mode Analysis

To understand what goes wrong when the robot fails, we labeled a subset of tasks (15/50) with multiple-choice failure reasons (see Figure 10):

- **Dexterity**: Clumsiness—the robot cannot reliably pick up or release items. This is the dominant problem, accounting for around one-third of failures

- **Order**: Progress errors—wrong action sequence. Many tasks require a particular order. Other common issue is deciding to finish early, exacerbated by "tails" on passive actions in the demo dataset

- **Confusion**: Weird behavior probably caused by the robot entering an out-of-distribution state

- **Robot fell**: While attempting to squat to pick an item from the floor, the robot sometimes starts to fall backward

- **Reasoning**: The robot should choose a locally non-obvious action (e.g., go around the table instead of trying to reach through it)

- **Search**: Randomness in the denoising process works surprisingly well for continued exploration, but can lead to repeatedly traversing the same area
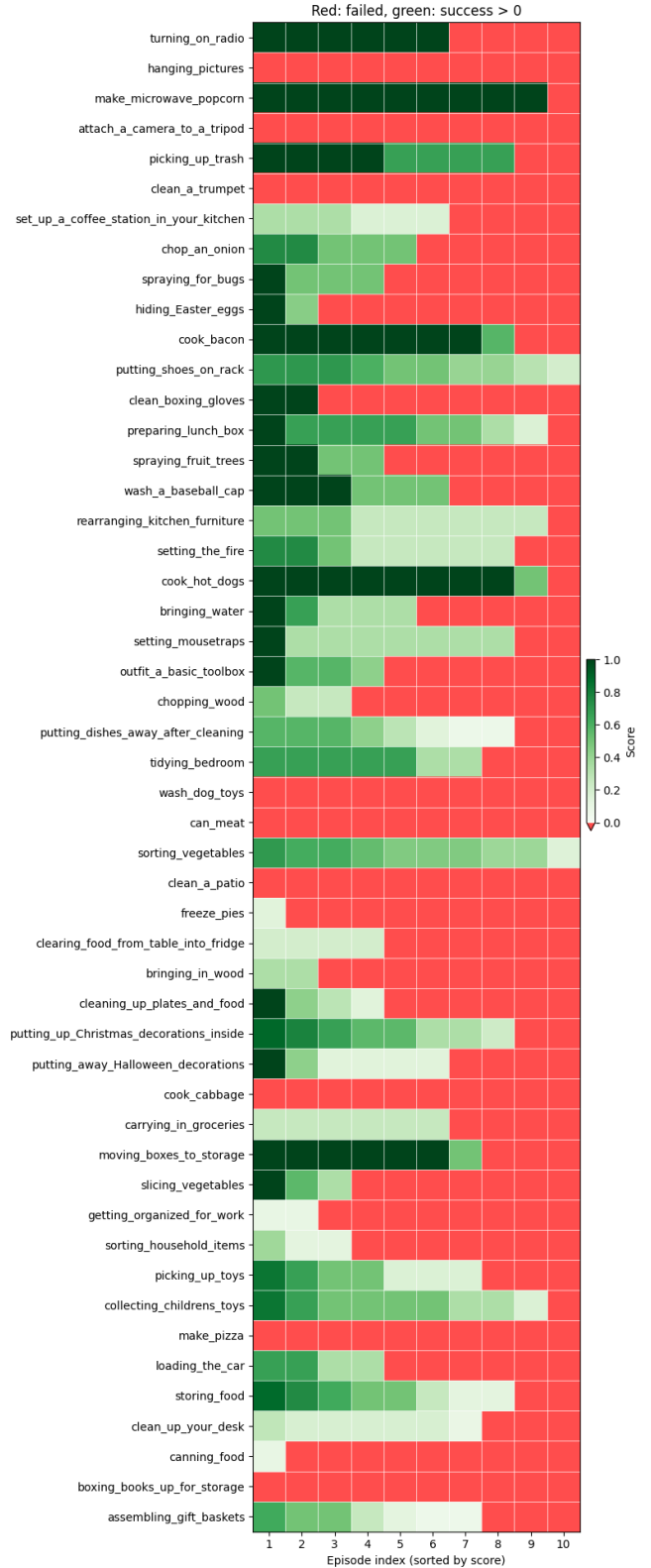


**Figure 9:** Per-episode scores on public evaluation. Each row is a task (sorted by average duration in demo dataset), each column an episode (sorted by score). Dark green indicates success, light green - partial success, red indicates failure.

| Rank | Team | Affiliation | Binary Success | | q_score | |
|------|------|-------------|------|------|------|------|
| | | | public | private | public | private |
| 1 | **Robot Learning Collective (ours)** | Independent | 0.1120 | **0.1240** | 0.2605 | **0.2599** |
| 2 | Comet | NVIDIA Research | 0.1440 | 0.1140 | 0.1830 | 0.2514 |
| 3 | SimpleAI Robot | Beijing Simple AI Technology Co Ltd | 0.1400 | 0.1080 | 0.1943 | 0.1591 |
| 4 | The North Star | Huawei CRI EAI Team | 0.1280 | 0.0760 | 0.1702 | 0.1204 |
| 5 | Embodied Intelligence | Independent | 0.0620 | 0.0520 | 0.1110 | 0.0947 |

**Table 1:** 2025 BEHAVIOR Challenge leaderboard results. Full leaderboard available at `https://behavior.stanford.edu/challe nge/leaderboard.html`.
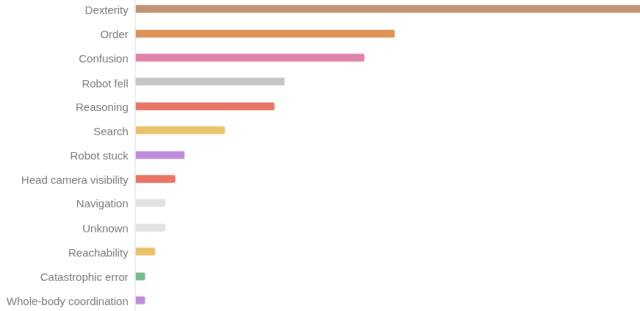


**Figure 10:** Distribution of failure reasons across labeled evaluation episodes. Dexterity issues (clumsiness in grasping/releasing) dominate, followed by progress/order errors and confusion from out-of-distribution states.

- **Robot stuck**: The model sometimes keeps trying to move while the robot base contact prevents it

- **Head camera visibility**: The head camera does not cover the floor in front of the robot or shelves above the head

- **Navigation**: Navigating to fixed locations (same as in the demo dataset) is generally not a problem, but in rare cases the model confuses the task and tries to go to the location of another task.

- **Unknown**: Failures that we could not attribute to any specific problem.

- **Reachability**: The robot is realistically limited in reach; sometimes—especially after unsuccessful grasps—this makes task completion impossible.

- **Catastrophic error**: The robot makes a mess of the environment, rendering the task unachievable.

- **Whole-body coordination**: The robot lacks awareness that an arm may be stuck, which prevents further safe navigation.

Many problems are amplified by the specifics of the data-collection procedure: operating a humanoid-like robot in simulation is challenging for human operators. Tasks were divided into steps; if the operator failed at a given step, the simulation was rolled back to the start of that substep. This enabled collection of complex, long-horizon demonstrations, but the resulting trajectories are very idealized and cover only a narrow manifold—small deviations from the ideal trajectory often cause the model to output random actions or freeze.

### 9.2.3 Additional Observations

**Emerging recovery behavior from cross-task learning**: This was a key factor in improving the model. Single-task models show no recovery behavior. The same architecture trained on all 50 tasks exhibits a wide range of recovery behaviors, such as picking up fallen objects.

In general, more training improves success rates across tasks, but for certain tasks the multi-task model performs worse. We hypothesize that this may be due to short task duration and thus low relative weight in the dataset, or because of the confusion between different tasks with similar visual features.

During the main training run we periodically branched off checkpoints and fine-tuned them on 1–2 tasks. Early in training this gave a significant performance boost, but at later stages the main run reached comparable performance, suggesting that undertraining was the main limiting factor.

### 9.2.4 Summary

Our analysis highlights problems that are currently a major focus for practitioners building real-world systems: VLA models for dexterous manipulation, System-2-style components to guide IL policies, and diverse pre-training datasets to widen the manifold where the model produces meaningful actions. This suggests that work on this challenge is relevant to real-world problems.

### 9.3 Ablations

We did not have the budget to run full ablations, but our small-scale experiments show:

- The model is surprisingly tolerant to image quality. Comparing 224×224 generation vs. 720×720 downscaled did not lead to meaningful changes

- More surprisingly, machines in our cloud provider used for evaluation did not support NGX (due to Docker settings), which led to easily noticeable image-quality degradation, but this had very small impact on success rate. As the private score was calculated with the correct setting, we decided to do nothing about it.

- Advantage-weighting dataset rows based on subtask duration (following the intuition that if a step took a long time for the operator, demonstration quality tends to be lower) did not show a success-rate improvement on the one task we tried (although we noticed a significant difference in task duration, likely corresponding to different operators)

- Small variations in inference parameters (number of actions to execute, execution speedup, voting history length) did not show significant change. Extreme changes (like execution of only 10 steps for 30 steps chunk, or speed up 2 times) led to the score degradation

- Gripper opening correction rule showed $2.2\times$ increase in q-score for a subset of 13 tasks with 3 episodes each (39 episodes total)

## 10    Discussion and Conclusion

We demonstrated how to adapt Pi0.5 for the BEHAVIOR Challenge, achieving 1st place with 26% q-score across 50 diverse household tasks. Our approach combines architectural innovations (correlated noise, learnable mixed-layer attention, System 2 stage tracking) with practical optimizations (action compression, correction rules).

**The result is far from optimal.** None of our models fully converged—despite training nearly continuously for a month, only $\sim$2 epochs passed through the dataset. Longer training would almost certainly improve results. The challenge timeline, not model capacity, was the limiting factor.

**Task-specific checkpoints helped but are they necessary?** We used 4 checkpoints primarily to optimize our remaining time before the deadline (even though we have seen positive effect from the task confusion reduction). We believe comparable results are achievable with a single multi-task model, which would be more elegant and practical for deployment. Comparing single model vs. multiple task-specific fine-tuned models is an interesting research question.

**Gripper correction rules were surprisingly effective.** This simple heuristic—detecting and recovering from accidental gripper closures—contributed significantly to our final score. More task-specific rules could likely boost performance further, though such rules are competition-specific and not generalizable research contributions.

**We propose several novel ideas but lack proper ablations.** Due to resource constraints, we could not isolate the contribution of each component. Rigorous ablation studies would be valuable to identify which innovations actually matter and which are redundant.

**Pretrained VLA models are essential.** Our attempts to train smaller models from scratch failed completely. Leveraging Pi0.5's pretrained weights was critical—the VLA paradigm of "pretrain on internet scale, fine-tune for robotics" remains the most practical path forward.

Long-horizon household manipulation remains challenging. Even with extensive engineering, we achieve only 26% success.

### 10.1    Future Directions

The following directions seem promising to us:

- **Collecting corrective actions**: We believe that even a comparatively small number of hours of demonstrations collected from states where the model fails can significantly improve its robustness

- **Offline RL with advantage reweighting**: This seems a promising approach to train the model to avoid unfavorable states

- **VLM as System 2**: Modern VLMs exhibit good video understanding. With no strict budget on using cloud APIs, powerful VLMs prompted with task-specific "checklists" may offer a strong effort/results trade-off

- **RL in simulation using privileged information with teacher–student distillation**: This should work, but is not directly applicable to real-world problems

- **Robot control constraints to prevent falling**: Implementing a safety envelope so that actions produced by the model do not cause the robot to fall

- **Fish-eye head camera**: As well as an additional camera on the back of the robot should help

## Acknowledgments

## References

[1] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, et al. Behavior-1k: A human-centered, embodied ai benchmark with 1,000 everyday activities and realistic simulation. *arXiv preprint arXiv:2403.09227*, 2024.

[2] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, et al. $\pi_0$: A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.

[3] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, et al. $\pi_{0.5}$: A vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.

[4] Johan Björck, Fernando Castaneda, Nikita Cherniadev, et al. GR00T N1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.

[5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, et al. RT-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.

[6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.

[7] Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn Kooijmans, Steven Palma, Adil Zouitine, et al. SmolVLA: A vision-language-action model for affordable and efficient robotics. *arXiv preprint arXiv:2506.01844*, 2025.

[8] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2023.

[9] 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. *The International Journal of Robotics Research*, 2024.

[10] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. FAST: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025.

[11] Lucy Xiaoyang Shi, Brian Ichter, Michael Equi, Liyiming Ke, et al. Hi robot: Open-ended instruction following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*, 2025.

[12] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.

[13] Kevin Black, Manuel Y. Galliker, and Sergey Levine. Real-time execution of action chunking flow policies. *arXiv preprint arXiv:2506.07339*, 2025.

[14] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. *arXiv preprint arXiv:2303.15343*, 2023.

[15] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, et al. PaliGemma: A versatile 3B VLM for transfer. *arXiv preprint arXiv:2407.07726*, 2024.

[16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.

[17] TRI LBM Team, Jose Barreiros, Andrew Beaulieu, Aditya Bhat, Rick Cory, Eric Cousineau, et al. A Careful Examination of Large Behavior Models for Multitask Dexterous Manipulation. *arXiv preprint arXiv:2507.05331*, 2025.

[18] Ilia Larchenko. Decoder Only Transformer Policy – simple but powerful model for behavioral cloning. Zenodo, 2025. doi:10.5281/zenodo.14789259. `https://github.com/IliaLarchenko/dot_policy`

[19] Figure AI. Helix Accelerating Real-World Logistics. Technical Blog, February 2025. `https://www.figure.ai/news/helix-logistics`

# Appendix

## A  Stage-Task Fusion Details

We fuse the task embedding with stage information using multiple learned representations. We don't advocate for this specific approach and didn't conduct proper ablation studies due to time and resource constraints—we simply added different representations hoping some would work well.

**Input representations**:

1. **Task embedding**: Base task vector $\mathbf{e}_\tau \in \mathbb{R}^{2048}$

2. **Sine-cosine stage encoding**: $\mathbf{s}_{\text{sincos}} \in \mathbb{R}^{1024}$ (positional encoding of normalized stage)

3. **Learned stage embedding**: $\mathbf{s}_{\text{learned}} \in \mathbb{R}^{1024}$ (task-specific trainable embeddings)

**Stage normalization** (task-specific):

$$s_{\text{norm}} = \frac{s}{\max(s_{\max}(\tau) - 1, 1)} \tag{15}$$

where $s_{\max}(\tau) \in [5, 15]$ is the number of stages for task $\tau$. This ensures stage $0 \to 0.0$ and final stage $\to 1.0$ for all tasks.

**Gated fusion mechanism**: We concatenate all inputs and learn sigmoid gates to modulate each component:

$$\mathbf{x}_{\text{all}} = [\mathbf{e}_\tau \,;\, \mathbf{s}_{\text{sincos}} \,;\, \mathbf{s}_{\text{learned}}] \in \mathbb{R}^{4096} \tag{16}$$

**Four fusion strategies**:

1. **Task-gated**: $\mathbf{f}_1 = \mathbf{e}_\tau \odot \mathbf{g}_{\text{task}}$

2. **Balanced fusion**: $\mathbf{f}_2 = \mathbf{W}_5 \cdot \text{ReLU}(\mathbf{W}_4 \mathbf{x}_{\text{all}} + \mathbf{b}_4) + \mathbf{b}_5$

3. **Stage-dominant**: $\mathbf{f}_3 = \mathbf{W}_6 \cdot [(\mathbf{s}_{\text{sincos}} \odot \mathbf{g}_{\text{sincos}}) \,;\, (\mathbf{s}_{\text{learned}} \odot \mathbf{g}_{\text{learned}})] + \mathbf{b}_6$

4. **Pure stage**: $\mathbf{f}_4 = [\mathbf{s}_{\text{sincos}} \,;\, \mathbf{s}_{\text{learned}}]$

All four representations have dimension 2048. The final output is $[\mathbf{f}_1 \,;\, \mathbf{f}_2 \,;\, \mathbf{f}_3 \,;\, \mathbf{f}_4] \in \mathbb{R}^{4 \times 2048}$. These 4 stage-conditioned tokens are passed to the model along with the base task embedding (total: 5 task-related tokens).

## B    Implementation Details

- **Action Space**: 23D continuous actions (3D base velocity, 4D trunk, 7D left arm, 1D left gripper, 7D right arm, 1D right gripper)

- **Action Horizon**: 30 timesteps

- **Image Resolution**: 224×224 RGB

- **Vision Backbone**: Frozen SigLIP-So400m/14

- **Framework**: JAX

- **Training**: FSDP model parallelism on 8×H200 GPUs

- **Inference**: Single GPU (tested on RTX 4090, 24GB VRAM)

## C    Action Expert Architecture

The action expert uses the Gemma 300M architecture (identical to Pi0.5):

- **Layers**: 18

- **Hidden dimension**: 1024

- **Feedforward dimension**: 4096

- **Attention heads**: 8 (with 1 KV head using grouped-query attention)

- **Head dimension**: 256

- **Parameters**: $\sim$311M (action expert only)

Input sequence consists of noisy actions $\mathbf{x}_t = t \cdot \epsilon + (1 - t) \cdot \mathbf{a}$ projected to 1024D, plus time embedding posemb_sincos$(t, 1024)$ processed by 2-layer MLP (AdaRMS conditioning).

## D    FAST Auxiliary Training

Since we removed text input and the Gemma tokenizer (replaced with task embeddings), we directly train FAST token embeddings rather than reusing the language model's vocabulary:

- **FAST token embedding layer**: Maps discrete FAST tokens (vocab size 1024) to PaliGemma width (2048D)

- **FAST projection head**: Projects VLM outputs back to FAST vocabulary for next-token prediction

- **Action encoding**: We encode 22 action dimensions (dimensions 0:6, 7:23) using a trained FAST tokenizer, excluding the padded dimensions and one dimension with 0 variance

- **Normalization**: FAST uses global quantile normalization to preserve temporal smoothness for better DCT compression, even when the rest of the model uses per-timestamp normalization

- **Sequence processing**: FAST tokens are processed autoregressively with causal attention, attending to all prefix tokens (images, task embeddings, stages, state)

The auxiliary loss is computed as standard cross-entropy on the predicted token sequence. The final model uses a FAST loss weight of **0.05**.

**Note**: Overall, we believe that FAST may be more beneficial for VLA pre-training and have a limited impact on fine-tuning. But we speculate (without proper ablations) that it helps with training the model to complete simpler tasks like navigation faster and prevents the model from "forgetting" them while the action expert focuses on dexterous manipulation tasks.