

Cue the Flow: Steering Flow-Matching Policies for Open-World Delivery Manipulation

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Open-world delivery requires mobile manipulators to follow free-form
2 user instructions and manipulate potentially novel objects. Existing dual-system
3 approaches use high-level grounding models to convert language into grounded
4 visual prompts, but their low-level controllers can remain brittle under noisy per-
5 ception, dynamic scenes, and contact-rich interactions. We instead use a pretrained
6 flow-matching vision-language-action model as the low-level control interface,
7 leveraging its reactivity and robustness to environmental changes while treating
8 the grounding output as a spatial cue for policy steering. Our key insight is that
9 the pretrained VLA already provides a strong manipulation prior, while the spatial
10 cue supplies the missing target information needed to guide actions under novel
11 language–object mappings. Concretely, we introduce a lightweight cue-conditioned
12 adapter. The adapter is first trained with contrastive objectives to produce salient
13 and spatially discriminative cue representations, and is then supervised to predict
14 a diagonal affine transformation over the generated action chunk, aligning policy
15 steering with the cued target. Across tabletop and mobile-base settings, our method
16 improves instruction following and manipulation success on both in-domain and
17 out-of-domain objects, achieving up to near $2\times$ improvement in average task
18 success rate with negligible inference overhead.

19 **Keywords:** Policy Steering, Mobile Manipulation; Generalization

20 1 Introduction

21 Consider a delivery robot receiving the request: “Give me the takeout bag with the name Joseph
22 on it.” We focus on the **manipulation phase** of this task, where the robot must turn an open-ended
23 instruction into a physical action: parsing the language, grounding the referred object in a potentially
24 cluttered scene, and executing a reliable trajectory toward it. While the underlying manipulation skill
25 may be simple, such as pick-and-place, learning a practical policy for open-world delivery remains
26 difficult. As shown in Figure 1(a), the core challenge is a substantial train–deployment mismatch:
27 the policy is trained in fixed environments with a limited vocabulary of objects and instructions, yet
28 deployed on a mobile robot in open-world settings with novel objects, unseen referring expressions,
29 and variable scene configurations.

30 Large-scale pretrained robot policies [1, 2, 3, 4] have substantially improved policy robustness
31 under visual perturbations [5], reducing the brittleness caused by changes in scene appearance and
32 environment layout. Yet a central challenge remains: how to extend a policy to novel language–object
33 mappings while preserving its previously learned manipulation capability. One promising direction is
34 a dual-system design that decouples explicit grounding from robot control [6, 7, 8, 9]. The high-level
35 grounding module, or System 2, interprets the instruction and converts it into spatial abstractions,
36 often represented as visual prompts [10, 11]. The low-level controller, or System 1, then conditions
37 on these prompts to generate executable robot actions.

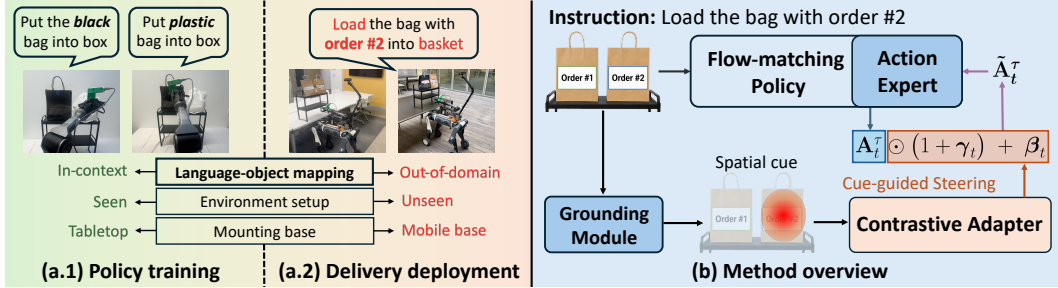


Figure 1: **(a.1)** The policy is trained on in-context demonstrations collected in a fixed environment. **(a.2)** In real-world delivery, however, the robot must execute manipulation under changing environments and novel language–object mappings. **(b)** We focus on improving policy generalization to these novel mappings. A grounding module converts the open-ended instruction into a spatial cue, which conditions a contrastive adapter to predict affine transformations over the frozen flow-matching policy’s action output, steering the generated action chunk toward the target object.

38 However, many existing dual-system methods
 39 still depend on optimization-based or geometric
 40 motion generation [12, 7, 13]. Such controllers
 41 can be fragile in open-world delivery settings:
 42 perception may be noisy, the target object may
 43 move during execution, and a fixed plan may
 44 not recover from grounding or localization errors. This motivates replacing the geometric low-
 45 level controller with a pretrained robot policy, especially a flow-matching vision-language-action
 46 model [14, 15, 16]. Since these models are trained on diverse robotic demonstrations, they provide
 47 learned priors for object interaction, contact-rich manipulation, and reactive correction. They therefore
 48 offer a practical low-level control interface for the dual-system design, while preserving the explicit
 49 grounding ability of the high-level module. Table 1 compares these paradigms.

Paradigms	VLA [3]	Dual + low-level [7]	Dual + VLA (ours)
Learned control	✓	✗	✓
Precise perception	✓	✗	✓
Explicit grounding	✗	✓	✓
OOD generalization	Limited	Strong	Strong

Table 1: **Comparison of paradigm designs.**

50 Yet integrating flow-matching VLAs into a dual-system design is not straightforward. The visual
 51 prompt produced by System 2 must be converted into a form that can reliably steer a policy trained to
 52 act from its own visual-language representations. Prior work [9, 17] explores full policy training with
 53 visual prompts as additional inputs. While effective in controlled settings, this strategy still depends
 54 on the VLA’s internal spatial reasoning ability, which is known to be limited [18, 19]. More critically,
 55 input-level prompting does not specify how the prompt should modulate the action generation process.
 56 After being encoded with the rest of the visual observation, the prompt is only implicitly represented,
 57 without a direct mechanism that encourages the generated trajectory to align with the System 2
 58 guidance. The key challenge is therefore not merely how to provide a visual prompt to a VLA, but
 59 how to turn that grounded signal into actionable guidance during flow-based action generation.

60 In this paper, we address this challenge by treating the grounding output as a spatial cue that directly
 61 guides action generation. Our key insight is that a pretrained flow-matching VLA already knows
 62 how to manipulate: it can generate feasible trajectories for objects in the workspace, but may not
 63 reliably determine which object to act on under unseen language–object mappings. Thus, the missing
 64 component is not a new manipulation policy, but a mechanism for specifying the intended target.
 65 We therefore keep the VLA frozen and introduce a lightweight cue-conditioned adapter that uses
 66 the System 2 generated spatial cue to steer the frozen policy’s action output. The adapter first
 67 learns to encode the explicit cue into a salient and spatially discriminative representation, using
 68 contrastive objectives to discourage reliance on scene-specific shortcuts. It then conditions on this cue
 69 representation to predict a diagonal affine transformation over the generated action chunk, shifting
 70 the frozen policy’s output toward the cued object while preserving its learned manipulation prior.
 71 Figure 1(b) provides an overview of the full pipeline.

72 Our contributions are threefold. First, we analyze the policy-learning challenges posed by open-world
 73 delivery and propose a dual-system framework that combines explicit language–object grounding
 74 with learned low-level action generation. Second, we develop a lightweight cue-conditioned adapter
 75 for frozen flow-matching VLAs, which steers generated action chunks through a diagonal affine

76 transformation while avoiding full policy fine-tuning. Third, we validate the proposed method across
77 three in-domain and out-of-domain object settings. Our method consistently outperforms all baselines,
78 achieving relative gains of 86% and 44% over the best-performing baselines in average task success
79 rate under tabletop and mobile-base settings, respectively.

80 2 Related Work

81 **Flow-matching VLAs.** A central goal in robot learning is to build generalist policies that can perform
82 diverse tasks and transfer across environments, objects, and embodiments. Vision-language-action
83 (VLA) models [1, 2, 15, 14] address this goal by mapping visual observations and language instruc-
84 tions directly to robot actions, often after training on large-scale robot and human demonstration data.
85 Recent flow-matching VLAs [15, 14, 16] further model action generation as a learned continuous flow
86 conditioned on fused visual-language features, allowing them to represent complex and potentially
87 multimodal action distributions [20]. Although flow-matching VLAs have begun to be used within
88 dual-system architectures [21, 22], how to effectively steer their action generation toward novel
89 language-object mappings remains underexplored.

90 **Policy steering.** Policy steering improves pretrained generative policies by modifying their sampling
91 or refinement process while keeping the base policy fixed [23]. Existing methods typically guide
92 generated actions with external objectives or auxiliary models, steering outputs toward task, safety,
93 or user-specified constraints. Model-predictive refinement methods [24, 25] use learned dynamics
94 models to improve task performance or enforce safety, while human-in-the-loop approaches [26,
95 27, 28] refine actions using user-provided subgoals, corrections, or preferences. Classifier- and
96 dynamics-guided methods [29, 30] further leverage latent visual dynamics models [31, 32] to bias
97 actions toward desired outcomes. Unlike prior steering methods that primarily use external signals to
98 refine in-domain behavior, our method treats spatial cues as target-specification signals and learns
99 cue-conditioned transformations over frozen VLA action chunks, enabling better generalization to
100 out-of-domain language-object mappings.

101 **Visual prompting for robotic policies.** Explicit visual representations provide a useful interface for
102 guiding robotic manipulation. One line of work uses visual marks to connect a vision-language model
103 (VLM) planner with a low-level controller [7, 6, 8], where the controller translates visual predictions
104 into executable actions. However, such controllers can be brittle in open-world settings where
105 perception is noisy and reactive control is needed. Another line of work treats actions as language [33,
106 34], using the VLM to directly output executable robot actions, but these methods depend heavily on
107 the grounding precision and reasoning reliability of the VLM. Other approaches [13, 9, 35] inject
108 visual prompts into the VLA input stream, yet the prompt’s influence on action generation remains
109 implicit. In contrast, our method uses a spatial cue separated from the input image as an explicit
110 target-specification signal, steering a frozen VLA by directly modulating its generated action chunk.

111 3 Method

112 We present the problem setup and detail method design in this section. Section 3.1 formulates gener-
113 alization to novel language-object mappings as a policy-steering problem. Section 3.2 introduces the
114 overall dual-system framework. Section 3.3 describes the spatial-cue-conditioned adapter for steering
115 the frozen policy. Section 3.4 provides the implementation details.

116 3.1 Problem Definition

117 Consider two deployments of the same frozen flow-matching policy. In the first, both the instruction
118 and the referred object fall within the training distribution, and the policy can execute the task
119 successfully. In the second, the target object is replaced by a novel object and the instruction is
120 changed accordingly. The underlying motor behavior may remain the same, such as reaching,
121 grasping, and placing, but the policy must now associate the instruction and visual scene with an
122 unseen target. The central challenge is therefore not to acquire a new manipulation skill, but to

123 redirect an existing manipulation prior toward the intended object during action generation. We
 124 formulate this as a policy-steering problem: given a frozen flow-matching VLA and an explicit spatial
 125 cue produced by a grounding module, the goal is to modulate the generated action chunk so that it
 126 follows the cued target while preserving the base policy’s learned manipulation capability.

127 3.2 Dual-System Design

128 Flow-matching VLAs generate action chunks using an action expert conditioned on features from
 129 a vision-language backbone. While this end-to-end grounding is effective for in-distribution tasks,
 130 it can become brittle when the policy encounters novel instructions referring to unseen objects. To
 131 reduce this uncertainty, we adopt a dual-system design that separates language–object grounding
 132 from action generation. The high-level grounding module converts diverse free-form instructions
 133 into a unified spatial cue, allowing the downstream policy to be steered by explicit target information
 134 rather than relying solely on implicit language grounding.

135 **System 2: Grounding module.** System 2 converts open-ended user instructions into a structured
 136 **spatial cue** for the downstream VLA, specifying the 2D target location in the workspace camera view.
 137 We first use a grounding model to predict spatial references for the target object, typically in the form
 138 of bounding boxes. An external segmentation model [36, 37] then converts these references into a
 139 pixel-level target mask. From this mask, we compute the object centroid and construct a Gaussian
 140 heatmap \mathbf{h}_t . This heatmap is not injected as an image-level visual prompt but used as an explicit
 141 spatial cue. All components in System 2 are frozen and used for zero-shot inference.

142 **System 1: VLA as controller interface.** To preserve the pretrained manipulation capability of the
 143 VLA while enabling it to be steered by the spatial cue from System 2, we freeze the VLA policy and
 144 attach a lightweight adapter. The adapter predicts a cue-conditioned diagonal affine transformation
 145 that modulates the action chunk produced by the frozen policy during flow integration. Specifically,
 146 let \mathbf{A}_t^τ denote the generated intermediate action chunk at integration time $\tau \in [0, 1]$. The adapter
 147 predicts two diagonal modulation terms, γ_t and β_t , and computes the steered action chunk as

$$\tilde{\mathbf{A}}_t^\tau = \mathbf{A}_t^\tau \odot (1 + \gamma_t) + \beta_t, \quad (1)$$

148 where \odot denotes element-wise multiplication. The scaling term γ_t reweights the frozen policy’s
 149 action output, while the shift term β_t provides an additive correction. Together, they define a
 150 transformation that steers the generated action chunk toward the target object.

151 3.3 Contrastive Adapter

152 Learning an effective adapter is challenging be-
 153 cause the spatial cue is aligned with the target
 154 object in the observation. Without additional
 155 constraints, the adapter may exploit shortcut
 156 correlations from in-domain training scenes or
 157 memorize scene-specific visual patterns, rather
 158 than using the cue as the source of target spec-
 159 ification. To strengthen the correspondence be-
 160 tween cue location and action transformation,
 161 we design the adapter with two components: a
 162 cue encoder that learns salient and spatially dis-
 163 criminative features from the explicit cue, and a
 164 prediction head that maps these features to affine
 165 transformation parameters for steering the frozen
 policy output. Figure 2 illustrates the adapter
 architecture and training pipeline.

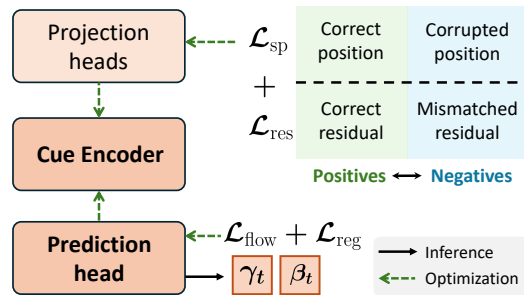


Figure 2: **Illustration of adapter learning.**

166 **Cue encoder learning.** The spatial cue specifies the target location but does not by itself encode
 167 object semantics or task context. We therefore train a cue encoder \mathcal{E} to transform the cue \mathbf{h}_t
 168 into a representation suitable for action steering. The learned representation is encouraged to
 169 satisfy two properties: **spatial discrimination**, preserving where the target lies relative to the action
 170 representation, and **cue salience**, emphasizing cue-relevant features over incidental visual patterns.

171 Let \mathcal{E} denote the cue encoder. To make the cue representation spatially discriminative, we require
 172 the action representation to align more closely with its corresponding spatial cue than with spatially
 173 shifted versions of the same cue. We first encode the Gaussian heatmap and project it into a
 174 cue representation, $\mathbf{z}_c = \mathcal{P}_c(\mathcal{E}(\mathbf{h}_t))$, where \mathcal{P}_c is a projection head. We also construct an action
 175 representation, $\mathbf{z}_a = \mathcal{P}_a(\mathbf{a}_t, \mathcal{E}(\mathbf{h}_t))$, where \mathbf{a}_t denotes the action-expert feature and \mathcal{P}_a is the
 176 corresponding projection head. For each positive cue representation \mathbf{z}_c^i , we construct a set of
 177 *spatial-shift negatives* $\{\mathbf{z}_{c,k}^-\}_{k=1}^K$ by perturbing the heatmap peak on the patch grid while keeping the
 178 underlying scene fixed. These negatives retain the same visual context but specify different target
 179 locations, forcing the cue representation to capture spatial position rather than appearance context.
 180 The spatial contrastive loss for sample i is

$$\mathcal{L}_{\text{sp}}^i = -\log \frac{\exp(\mathbf{z}_a^i \cdot \mathbf{z}_c^i / \tau)}{\sum_{j=1}^B \exp(\mathbf{z}_a^i \cdot \mathbf{z}_c^j / \tau) + \sum_{k=1}^K \exp(\mathbf{z}_a^i \cdot \mathbf{z}_{c,k}^- / \tau)}, \quad (2)$$

181 where B is the batch size and τ is the softmax temperature. This objective contrasts the correct cue
 182 location against both in-batch cues and spatially corrupted cues, encouraging the encoder to learn
 183 features that are sensitive to the target location rather than scene-level appearance.

184 Since the adapter also conditions on action-related features, it may learn shortcuts that are not driven
 185 by the spatial cue itself. We therefore introduce a residual cue representation that isolates the effect
 186 of the cue. For each sample, we compute two encodings: one from the real heatmap, $\mathcal{E}(\mathbf{h}_t)$, and
 187 one from a zero heatmap, $\mathcal{E}(\mathbf{0})$. Their difference captures the cue-induced residual, which is then
 188 projected into the same representation space as the action feature: $\mathbf{z}_{\text{res}} = \mathcal{P}_\delta(\mathcal{E}(\mathbf{h}_t) - \mathcal{E}(\mathbf{0}))$, where
 189 \mathcal{P}_δ is a projection head. This reinforces the adapter to learn cue-relevant action transformations.
 190 Therefore, for each sample i , we introduce the residual contrastive loss:

$$\mathcal{L}_{\text{res}}^i = -\log \frac{\exp(\mathbf{z}_a^i \cdot \mathbf{z}_{\text{res}}^i / \tau)}{\sum_{j=1}^B \exp(\mathbf{z}_a^i \cdot \mathbf{z}_{\text{res}}^j / \tau)}. \quad (3)$$

191 **Affine transformation learning.** Given the encoded cue, a prediction head outputs the modulation
 192 terms γ_t and β_t , which define the diagonal affine transformation in Eq. (1). We train this prediction
 193 head under the same flow-matching supervision as the base policy, so that the adapter learns to steer
 194 the frozen policy’s action chunks while preserving the pretrained action distribution. To discourage
 195 overly large corrections, we regularize the magnitude of the modulation terms:

$$\mathcal{L}_{\text{reg}} = \|\gamma_t\|_2^2 + \|\beta_t\|_2^2. \quad (4)$$

196 The full training objective is therefore

$$\mathcal{L} = \mathcal{L}_{\text{flow}} + \lambda_{\text{sp}} \mathcal{L}_{\text{sp}} + \lambda_{\text{res}} \mathcal{L}_{\text{res}} + \lambda_{\text{reg}} \mathcal{L}_{\text{reg}}, \quad (5)$$

197 where $\mathcal{L}_{\text{flow}}$ trains the affine transformation under the base flow-matching supervision, while the
 198 auxiliary terms align the transformation with the spatial cue and regularize its magnitude to avoid
 199 corrupting the pretrained action distribution.

200 3.4 Implementation Details

201 **Training.** We train the adapter on the same demonstrations used to learn the flow-matching policy,
 202 supplemented with ground-truth spatial cues for target-location supervision. The adapter is imple-
 203 mented as a lightweight sidecar to the action expert, modulating generated action chunks while
 204 keeping the base policy frozen. We also evaluated LoRA [38] on the action expert as an alternative
 205 for learning cue-conditioned action transformations. In our low-data setting, LoRA tended to overfit
 206 to the training cue–action pairs, likely because its updates are directly coupled to the action-expert
 207 representations, and did not yield robust cue-steering behavior.

208 **Deployment.** During deployment, the adapter uses the same residual cue formulation as in training:
 209 it takes both the System 2 heatmap \mathbf{h}_t and a zero heatmap $\mathbf{0}$, and predicts the modulation terms from
 210 their residual cue feature. For efficiency, the adapter adds about 0.02% parameters relative to the base
 211 policy and introduces negligible latency. System 2 is also invoked once at the beginning of each task
 212 rollout, unless substantial scene changes require the spatial cue to be recomputed.

Method	Tabletop				Mobile Base			
	Setting I	Setting II	Setting III	Average	Setting I	Setting II	Setting III	Average
Base [14]	70.8	45.8	0.0	38.9	33.3	50.0	8.3	30.5
Base-L [14]	37.5	8.3	29.2	25.0	8.3	0.0	8.3	5.5
MOKA [7]	20.8	25.0	8.3	18.0	16.7	8.3	0.0	8.3
VP-VLA [9]	58.3	45.8	4.2	36.1	58.3	66.7	8.3	44.4
Ours	83.3	66.7	66.7	72.2	83.3	66.7	41.7	63.9

Table 2: **Success rates (%) comparison under tabletop and mobile-base conditions across three settings.** (I) seen objects with paraphrased instructions, (II) unseen objects with in-domain instructions, and (III) unseen objects with novel instructions. Best results in each column are **bolded**.

213 4 Experiments

214 4.1 Experimental Setup

215 **Pipeline.** For policy adaptation, we collect 80 demonstration trajectories using two training objects,
 216 a black bag and a plastic bag, without any label tags. We fine-tune $\pi_{0.5}$ [14] on this dataset with
 217 a standard training recipe and use the resulting model as the base policy for our method. During
 218 deployment, we use Qwen3-VL-2B [39] and SAM2.1 [36, 37] as the grounding module. These
 219 model choices allow the full pipeline to run simultaneously on a single RTX 5080 GPU, with the
 220 VLA operating continuously at 15,Hz. Experiments are conducted with a 6-DoF robotic arm and a
 221 wheeled dog robot, which provides vibration robustness and flexible hardware integration.

222 **Evaluation.** We evaluate under three settings. **I. Seen objects with paraphrased language** uses
 223 objects from the fine-tuning dataset with perturbed prompts. **II. Unseen objects with in-domain**
 224 **language** uses target objects absent from the collected trajectories, while keeping object names within
 225 the pretrained $\pi_{0.5}$ vocabulary, such as “red bag”. **III. Unseen objects with novel language** further
 226 introduces novel referring expressions, such as “bag with order number 1”, representing the most
 227 challenging setting. For each task, we curate six prompts and two spatial configurations. Each prompt
 228 is generated from a compositional template, such as “<grasp-action> <target-object> and
 229 <place-action> into the box,” where the grasping action, target object, and placing action vary
 230 across trials. The placement location is fixed, since delivery tasks can use a deterministic handover or
 231 drop-off location. Because all training demonstrations are collected on tabletop, deployment on the
 232 mobile base introduces an additional shift in viewpoint and scene geometry. Task success is measured
 233 as a binary outcome indicating whether the robot correctly completes the user instruction. The full
 234 set of evaluation tasks, prompts, and spatial configurations is provided in the supplementary material.

235 **Baselines.** We compare our method against the following baselines:

- 236 • **Base** [14]. We fine-tune pretrained $\pi_{0.5}$ on our demonstrations using one fixed language prompt
 237 per task, corresponding to a standard single-VLA setting.
- 238 • **Base-L.** We fine-tune pretrained $\pi_{0.5}$ with multiple language prompts per trajectory, using the
 239 same paraphrased prompts as in Setting I. This aligns the training prompt distribution more closely
 240 with the evaluation prompts and tests whether language augmentation improves generalization.
- 241 • **MOKA** [7]. MOKA represents a dual-system approach with a high-level planner and a geometric
 242 low-level controller. It relies on the planner to produce accurate 2D keypoints and waypoints as
 243 affordance representations, which are then deprojected into 3D space for controller execution.
- 244 • **VP-VLA** [9]. VP-VLA is also a dual-system approach that uses a VLA as the controller interface.
 245 Unlike our method, it injects the generated visual prompt as additional policy input and trains the
 246 VLA with an auxiliary grounding objective to focus on the prompted region.

247 4.2 Results

248 **Quantitative results.** Table 2 summarizes the quantitative results. Moving the robot arm from
 249 a stable tabletop setup to a mobile base degrades performance for all methods, due to changes

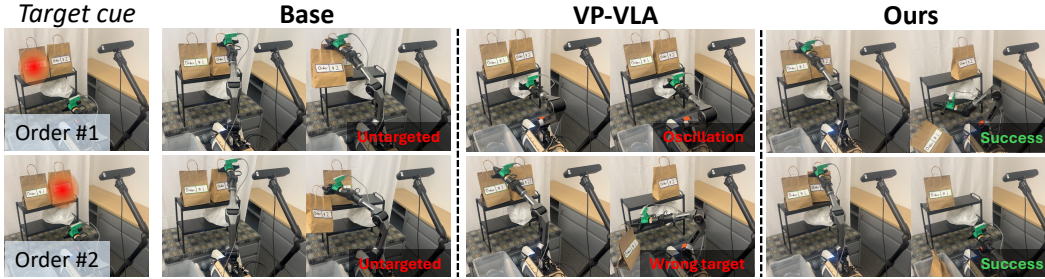


Figure 3: **Qualitative rollouts under the unseen-object with novel-language setting (Setting III).** The leftmost column shows the System 2 spatial cue, where the red heatmap indicates the target object. For each method, two temporal keyframes are shown. The Base policy produces untargeted behavior, while VP-VLA either oscillates or reaches the wrong object. Our method follows the spatial cue and successfully manipulates the specified target in both task examples.

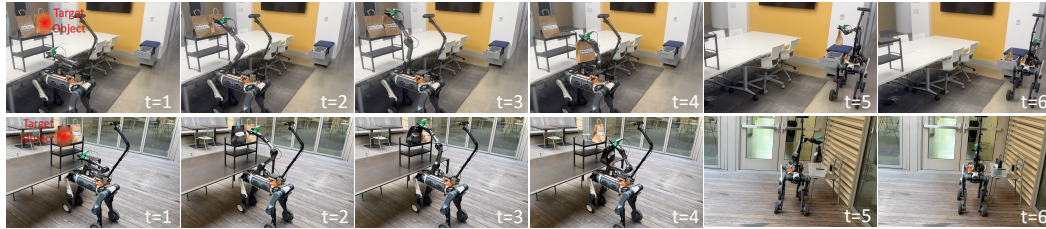


Figure 4: **Extension to a complete delivery pipeline.** We show example rollouts that combine manual navigation with autonomous manipulation. The red heatmap indicates the correctly grounded target object. After the mobile base is manually navigated to the workspace, our policy performs autonomous cue-conditioned manipulation to complete the task.

250 in viewpoint and base instability. resembling real-world delivery and often producing near-zero
 251 success for several baselines. The Base policy performs well on in-domain tasks and shows some
 252 generalization to object-level out-of-domain cases. We attribute this to the pretrained VLA prior,
 253 where the relevant visual-language patterns are likely represented in large-scale pretraining data.
 254 However, this prior is insufficient when both the target object and the referring expression are out of
 255 distribution. For Base-L, adding diverse language prompts during fine-tuning introduces a tradeoff: it
 256 reduces in-domain success but improves over Base on object-language-level out-of-domain tasks.
 257 Yet, we observe that this gain mainly comes from a biased motion trajectory that occasionally reaches
 258 the correct target, rather than from reliable language-conditioned grounding.

259 For the dual-system baselines, MOKA fails when the System 2 planner produces inaccurate visual
 260 prompts. This reflects the difficulty of generating dense spatial guidance for the full execution
 261 trajectory from a high-level planner alone. In our setup, the challenge is further amplified by the
 262 camera configuration: the forward-facing camera observes a relatively large workspace with the robot
 263 also in view, increasing perspective variation and visual clutter, and making accurate localization
 264 more difficult. VP-VLA performs comparably to the Base policy in Settings I and II, but still struggles
 265 in Setting III. This suggests that auxiliary grounding supervision alone does not reliably align the
 266 generated action trajectory with the visual prompt under stronger distribution shifts. In contrast, our
 267 method achieves better performance across all settings. We attribute this improvement to two design
 268 choices: our framework requires only sparse spatial cues from System 2, which are easier to obtain
 269 than full trajectory-level reasoning; and the learned adapter uses these cues to steer the frozen policy
 270 output while preserving the pretrained VLA prior and fine-tuned manipulation skill.

271 **Qualitative results.** Figure 3 compares qualitative rollouts in the unseen-object with novel-language
 272 setting on the mobile base. The Base policy produces untargeted trajectories, often reaching toward
 273 both bags on the shelf. VP-VLA shows stronger visual conditioning, but still lacks stable cue-action
 274 alignment, leading to oscillatory behavior or reaching toward the wrong object. This indicates that
 275 visual-prompt conditioning alone provides limited improvement under this stronger distribution shift
 276 and may also degrade manipulation stability. In contrast, our method consistently guides the frozen
 277 policy toward the specified target and enables successful manipulation. Figure 4 further demonstrates

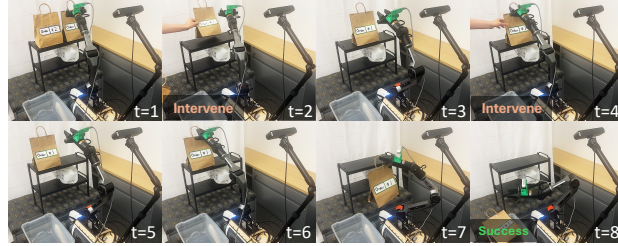
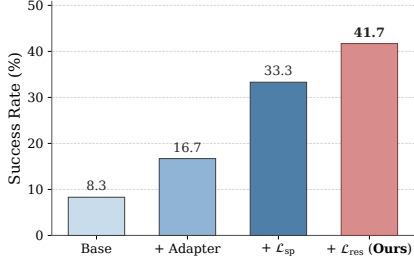


Figure 5: **Effectiveness of components.** Figure 6: **Reactivity to dynamically moving objects.**

278 that the cue-conditioned manipulation skill transfers across diverse scenes. The policy remains robust
 279 to environmental changes during execution and can be integrated with mobile-base navigation. Full
 280 video demonstrations and policy execution details are provided in the supplementary material.

281 4.3 Analysis

282 **Component effectiveness.** Figure 5 evaluates the effectiveness of the proposed learning components
 283 by comparing their performance under Setting III. Training the adapter only with the flow-matching
 284 loss and spatial cue features yields limited gains, indicating that cue conditioning alone is insufficient.
 285 Adding the spatial and residual contrastive losses consistently improves performance, confirming the
 286 importance of learning cue representations that are both spatially discriminative and salient.

287 **Effect of spatial cue format.** We further study the effect of different spatial cue formats. In addition
 288 to the soft Gaussian heatmap, we evaluate a binary target mask, defined as a segmentation silhouette
 289 at patch resolution, and an axis-aligned bounding box, which captures coarse object extent. The
 290 Gaussian heatmap and binary mask yield similar behavior and performance, indicating that localized
 291 target cues are sufficient for effective steering. In contrast, the bounding box performs worse, likely
 292 because its coarse boundaries can direct the robot toward invalid or non-manipulable object regions.
 293 This suggests that fine-grained spatial cues are more effective than the coarse box-level one.

294 **Reactivity to dynamically changing environments.** Figure 6 demonstrates the benefit of using a
 295 large-scale VLA as the System 1 controller by perturbing the target object’s position during execution.
 296 With the spatial cue grounded only once at the beginning of the rollout, our method still reacts to
 297 environmental changes without re-querying the high-level planner. This reactivity is inherited from
 298 the underlying VLA and preserved by our adapter, highlighting an advantage over conventional
 299 geometric controllers that rely on fixed trajectories or repeated replanning.

300 5 Conclusion

301 We introduced a spatial-cue-steered dual-system framework for open-world delivery manipulation.
 302 By replacing the geometric low-level controller with a pretrained flow-matching VLA, our framework
 303 improves robustness to visual perturbations and reactivity to dynamic environments. It converts
 304 System 2 outputs into action-steering signals via a lightweight adapter while keeping the VLA
 305 frozen, preserving its manipulation prior and improving generalization to novel language–object
 306 mappings. Across tabletop and mobile manipulation evaluations, our method improves in-domain
 307 and out-of-domain performance with negligible additional inference cost, showing that spatial-cue
 308 steering provides a practical interface between open-vocabulary grounding and reactive VLA control.

309 6 Limitations

310 Our method depends on the underlying flow-matching policy and remains sensitive to out-of-
 311 distribution affordance regions and object positions, constraining the mobile base’s stopping location
 312 and orientation relative to the target. Future work will address more complex delivery scenarios,
 313 including human interaction and handover, item selection from densely clustered bags, and tighter
 314 integration with navigation policies to improve manipulation reliability.

References

- 315
- 316 [1] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster,
317 G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine,
318 P. Liang, and C. Finn. Openvla: An open-source vision-language-action model, 2024. URL
319 <https://arxiv.org/abs/2406.09246>.
- 320 [2] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess,
321 A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman,
322 A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal,
323 L. Lee, T.-W. E. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao,
324 K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut,
325 H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, J. Wu, F. Xia, T. Xiao,
326 P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-2: Vision-language-action models transfer web
327 knowledge to robotic control, 2023. URL <https://arxiv.org/abs/2307.15818>.
- 328 [3] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman,
329 B. Ichter, S. Jakubczak, T. Jones, L. Ke, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch,
330 L. X. Shi, J. Tanner, Q. Vuong, A. Walling, H. Wang, and U. Zhilinsky. π_0 : A vision-language-
331 action flow model for general robot control, 2026. URL [https://arxiv.org/abs/2410.](https://arxiv.org/abs/2410.24164)
332 [24164](https://arxiv.org/abs/2410.24164).
- 333 [4] X. Chen, Y. Chen, Y. Fu, N. Gao, J. Jia, W. Jin, H. Li, Y. Mu, J. Pang, Y. Qiao, Y. Tian,
334 B. Wang, B. Wang, F. Wang, H. Wang, T. Wang, Z. Wang, X. Wei, C. Wu, S. Yang, J. Ye,
335 J. Yu, J. Zeng, J. Zhang, J. Zhang, S. Zhang, F. Zheng, B. Zhou, and Y. Zhu. Internvla-m1:
336 A spatially guided vision-language-action framework for generalist robot policy, 2025. URL
337 <https://arxiv.org/abs/2510.13778>.
- 338 [5] W. Pumacay, I. Singh, J. Duan, R. Krishna, J. Thomason, and D. Fox. The colosseum:
339 A benchmark for evaluating generalization for robotic manipulation, 2024. URL [https://](https://arxiv.org/abs/2402.08191)
340 arxiv.org/abs/2402.08191.
- 341 [6] W. Huang, C. Wang, Y. Li, R. Zhang, and L. Fei-Fei. Rekep: Spatio-temporal reasoning of
342 relational keypoint constraints for robotic manipulation, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2409.01652)
343 [abs/2409.01652](https://arxiv.org/abs/2409.01652).
- 344 [7] F. Liu, K. Fang, P. Abbeel, and S. Levine. Moka: Open-world robotic manipulation through
345 mark-based visual prompting, 2024. URL <https://arxiv.org/abs/2403.03174>.
- 346 [8] W. Yuan, J. Duan, V. Blukis, W. Pumacay, R. Krishna, A. Murali, A. Mousavian, and D. Fox.
347 Robopoint: A vision-language model for spatial affordance prediction for robotics, 2024. URL
348 <https://arxiv.org/abs/2406.10721>.
- 349 [9] Z. Wang, Y. Chen, Y. Liu, J. Ye, P. Chen, C. Lu, S. Liu, B. Yu, and J. Jia. Vp-vla: Visual
350 prompting as an interface for vision-language-action models, 2026. URL [https://arxiv.](https://arxiv.org/abs/2603.22003)
351 [org/abs/2603.22003](https://arxiv.org/abs/2603.22003).
- 352 [10] A. Bar, Y. Gandelsman, T. Darrell, A. Globerson, and A. Efros. Visual prompting via image
353 inpainting. *Advances in neural information processing systems*, 35:25005–25017, 2022.
- 354 [11] J. Yang, H. Zhang, F. Li, X. Zou, C. Li, and J. Gao. Set-of-mark prompting unleashes extraordi-
355 nary visual grounding in gpt-4v, 2023. URL <https://arxiv.org/abs/2310.11441>.
- 356 [12] L. Manuelli, W. Gao, P. Florence, and R. Tedrake. kpm: Keypoint affordances for category-
357 level robotic manipulation, 2019. URL <https://arxiv.org/abs/1903.06684>.
- 358 [13] R. Zheng, Y. Liang, S. Huang, J. Gao, H. D. III, A. Kolobov, F. Huang, and J. Yang. Tracevla:
359 Visual trace prompting enhances spatial-temporal awareness for generalist robotic policies,
360 2025. URL <https://arxiv.org/abs/2412.10345>.

- 361 [14] P. Intelligence, K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail, M. Equi,
362 C. Finn, N. Fusai, M. Y. Galliker, D. Ghosh, L. Groom, K. Hausman, B. Ichter, S. Jakubczak,
363 T. Jones, L. Ke, D. LeBlanc, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, A. Z.
364 Ren, L. X. Shi, L. Smith, J. T. Springenberg, K. Stachowicz, J. Tanner, Q. Vuong, H. Walke,
365 A. Walling, H. Wang, L. Yu, and U. Zhilinsky. $\pi_{0.5}$: a vision-language-action model with
366 open-world generalization, 2025. URL <https://arxiv.org/abs/2504.16054>.
- 367 [15] NVIDIA, :, J. Bjorck, F. Castañeda, N. Cherniadev, X. Da, R. Ding, L. J. Fan, Y. Fang, D. Fox,
368 F. Hu, S. Huang, J. Jang, Z. Jiang, J. Kautz, K. Kundalia, L. Lao, Z. Li, Z. Lin, K. Lin, G. Liu,
369 E. Llontop, L. Magne, A. Mandlekar, A. Narayan, S. Nasiriany, S. Reed, Y. L. Tan, G. Wang,
370 Z. Wang, J. Wang, Q. Wang, J. Xiang, Y. Xie, Y. Xu, Z. Xu, S. Ye, Z. Yu, A. Zhang, H. Zhang,
371 Y. Zhao, R. Zheng, and Y. Zhu. Gr00t n1: An open foundation model for generalist humanoid
372 robots, 2025. URL <https://arxiv.org/abs/2503.14734>.
- 373 [16] M. Shukor, D. Aubakirova, F. Capuano, P. Kooijmans, S. Palma, A. Zouitine, M. Aractingi,
374 C. Pascal, M. Russi, A. Marafioti, S. Alibert, M. Cord, T. Wolf, and R. Cadene. Smolvla:
375 A vision-language-action model for affordable and efficient robotics, 2025. URL <https://arxiv.org/abs/2506.01844>.
- 377 [17] D. Kim, S. Jan, H. Park, and D. Lim. Vca: Vision-click-action framework for precise
378 manipulation of segmented objects in target ambiguous environments, 2026. URL <https://arxiv.org/abs/2602.23583>.
- 380 [18] D. Qu, H. Song, Q. Chen, Y. Yao, X. Ye, Y. Ding, Z. Wang, J. Gu, B. Zhao, D. Wang, and X. Li.
381 Spatialvla: Exploring spatial representations for visual-language-action model, 2025. URL
382 <https://arxiv.org/abs/2501.15830>.
- 383 [19] Y. Feng, W. Zhang, Y. Wang, H. Luo, H. Yuan, S. Zheng, and Z. Lu. Spatial-aware vla
384 pretraining through visual-physical alignment from human videos, 2025. URL <https://arxiv.org/abs/2512.13080>.
- 386 [20] C. Pan, G. Anantharaman, N.-C. Huang, C. Jin, D. Pfrommer, C. Yuan, F. Permenter, G. Qu,
387 N. Boffi, G. Shi, and M. Simchowitz. Much ado about noising: Dispelling the myths of
388 generative robotic control, 2026. URL <https://arxiv.org/abs/2512.01809>.
- 389 [21] L. X. Shi, B. Ichter, M. Equi, L. Ke, K. Pertsch, Q. Vuong, J. Tanner, A. Walling, H. Wang,
390 N. Fusai, A. Li-Bell, D. Driess, L. Groom, S. Levine, and C. Finn. Hi robot: Open-ended
391 instruction following with hierarchical vision-language-action models, 2025. URL <https://arxiv.org/abs/2502.19417>.
- 393 [22] C. Li, J. Wen, Y. Peng, Y. Peng, F. Feng, and Y. Zhu. Pointvla: Injecting the 3d world into
394 vision-language-action models, 2025. URL <https://arxiv.org/abs/2503.07511>.
- 395 [23] M. Nakamoto, O. Mees, A. Kumar, and S. Levine. Steering your generalists: Improving robotic
396 foundation models via value guidance, 2025. URL <https://arxiv.org/abs/2410.13816>.
- 397 [24] K. Nakamura, L. Peters, and A. Bajcsy. Generalizing safety beyond collision-avoidance via
398 latent-space reachability analysis. In *Robotics: Science and Systems XXI*. Robotics: Science
399 and Systems Foundation, June 2025. doi:10.15607/rss.2025.xxi.113. URL <http://dx.doi.org/10.15607/RSS.2025.XXI.113>.
- 401 [25] Y. Song, L. Le, Y.-H. Park, J. Wang, J. Shi, L. Liu, J. Gu, E. Eaton, D. Jayaraman, and
402 K. Daniilidis. Omniguide: Universal guidance fields for enhancing generalist robot policies,
403 2026. URL <https://arxiv.org/abs/2603.10052>.
- 404 [26] H. Cai, Z. Peng, and B. Zhou. Predictive preference learning from human interventions, 2025.
405 URL <https://arxiv.org/abs/2510.01545>.

- 406 [27] Y. Wu, R. Tian, G. Swamy, and A. Bajcsy. From foresight to forethought: Vlm-in-the-loop
407 policy steering via latent alignment, 2025. URL <https://arxiv.org/abs/2502.01828>.
- 408 [28] Z. Peng, W. Mo, C. Duan, Q. Li, and B. Zhou. Learning from active human involvement through
409 proxy value propagation, 2025. URL <https://arxiv.org/abs/2502.03369>.
- 410 [29] M. Du and S. Song. Dynaguide: Steering diffusion polices with active dynamic guidance, 2025.
411 URL <https://arxiv.org/abs/2506.13922>.
- 412 [30] Z. Sun and S. Song. Latent policy barrier: Learning robust visuomotor policies by staying
413 in-distribution, 2025. URL <https://arxiv.org/abs/2508.05941>.
- 414 [31] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent
415 dynamics for planning from pixels, 2019. URL <https://arxiv.org/abs/1811.04551>.
- 416 [32] P. Wu, A. Escontrela, D. Hafner, K. Goldberg, and P. Abbeel. Daydreamer: World models for
417 physical robot learning, 2022. URL <https://arxiv.org/abs/2206.14176>.
- 418 [33] X. Li, L. Xu, M. Zhang, J. Liu, Y. Shen, I. Ponomarenko, J. Xu, L. Heng, S. Huang, S. Zhang,
419 and H. Dong. Crayonrobo: Object-centric prompt-driven vision-language-action model for
420 robotic manipulation, 2025. URL <https://arxiv.org/abs/2505.02166>.
- 421 [34] A. J. Hancock, X. Wu, L. Zha, O. Russakovsky, and A. Majumdar. Actions as language:
422 Fine-tuning vlms into vlas without catastrophic forgetting, 2025. URL <https://arxiv.org/abs/2509.22195>.
- 424 [35] S. Lee, S. Mo, and W.-S. Han. Bring my cup! personalizing vision-language-action models
425 with visual attentive prompting, 2026. URL <https://arxiv.org/abs/2512.20014>.
- 426 [36] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland,
427 L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár,
428 and C. Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint*
429 *arXiv:2408.00714*, 2024. URL <https://arxiv.org/abs/2408.00714>.
- 430 [37] N. Carion, L. Gustafson, Y.-T. Hu, S. Debnath, R. Hu, D. Suris, C. Ryali, K. V. Alwala,
431 H. Khedr, A. Huang, J. Lei, T. Ma, B. Guo, A. Kalla, M. Marks, J. Greer, M. Wang, P. Sun,
432 R. Rädle, T. Afouras, E. Mavroudi, K. Xu, T.-H. Wu, Y. Zhou, L. Momeni, R. Hazra, S. Ding,
433 S. Vaze, F. Porcher, F. Li, S. Li, A. Kamath, H. K. Cheng, P. Dollár, N. Ravi, K. Saenko,
434 P. Zhang, and C. Feichtenhofer. Sam 3: Segment anything with concepts, 2026. URL <https://arxiv.org/abs/2511.16719>.
- 436 [38] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora:
437 Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- 438
- 439 [39] S. Bai, Y. Cai, R. Chen, K. Chen, X. Chen, Z. Cheng, L. Deng, W. Ding, C. Gao, C. Ge, W. Ge,
440 Z. Guo, Q. Huang, J. Huang, F. Huang, B. Hui, S. Jiang, Z. Li, M. Li, M. Li, K. Li, Z. Lin,
441 J. Lin, X. Liu, J. Liu, C. Liu, Y. Liu, D. Liu, S. Liu, D. Lu, R. Luo, C. Lv, R. Men, L. Meng,
442 X. Ren, X. Ren, S. Song, Y. Sun, J. Tang, J. Tu, J. Wan, P. Wang, P. Wang, Q. Wang, Y. Wang,
443 T. Xie, Y. Xu, H. Xu, J. Xu, Z. Yang, M. Yang, J. Yang, A. Yang, B. Yu, F. Zhang, H. Zhang,
444 X. Zhang, B. Zheng, H. Zhong, J. Zhou, F. Zhou, J. Zhou, Y. Zhu, and K. Zhu. Qwen3-vl
445 technical report, 2025. URL <https://arxiv.org/abs/2511.21631>.

446 **A Method Details**

447 **A.1 MOKA**

448 MOKA was originally designed for tabletop manipulation with top-down observations and 2D planar
449 affordance prediction, which does not directly match our third-person shelf-picking setting with
450 deformable bag objects and non-planar interactions. Rather than fully re-engineering MOKA for
451 this domain, we evaluate its high-level planning capability under an oracle execution assumption: a
452 prediction is considered successful if the generated 2D affordance points correspond to semantically
453 correct and physically feasible interaction locations.

454 For consistency across methods, we standardize the perception stack with Qwen3-VL [39] and
455 SAM2.1 [36], and adapt the prompting strategy to our shelf-manipulation tasks. We also evaluate
456 MOKA under its 3-shot in-context configuration, corresponding to the strongest setting reported in the
457 original work. Despite these adaptations and the oracle execution assumption, performance remains
458 limited due to target misclassification, incorrect grasp-point prediction, and infeasible proposed
459 motions. These results indicate that MOKA’s 2D affordance representation has limited transferability
460 to our shelf-based delivery manipulation setting.

461 **A.2 VP-VLA**

462 The original VP-VLA method is designed for autoregressive VLAs and is not directly formulated for
463 flow-matching policies. We adapt its visual-prompting mechanism by providing the generated visual
464 prompt as an additional visual input to the VLA and training the policy with auxiliary grounding
465 supervision. This baseline differs from our method primarily in how the System 2 visual signal is
466 used: VP-VLA treats it as an input-level prompt, whereas our method converts the spatial cue into an
467 action-steering signal that modulates the frozen policy output. For consistency across methods, we
468 use the same perception stack with Qwen3-VL [39] and SAM2.1 [36].

469 **A.3 Ours**

470 We rewrite each user instruction into a standardized command format compatible with the VLA’s
471 training distribution, removing task-irrelevant details while keeping the target specification. For
472 demonstrations involving mobile navigation, we divide execution into three stages. First, the mobile
473 base remains stationary while the robot arm performs autonomous manipulation. Once the gripper
474 grasps the target object, policy inference is paused and a hard-coded arm motion moves the object to
475 a predefined holding pose for stable transport. The grasping moment is detected when the gripper jaw
476 width falls below a fixed threshold. Second, we manually teleoperate or use map-based autonomous
477 navigation to guide the mobile base toward the target placement location. Third, policy inference is
478 resumed to allow the robot to autonomously complete the placement stage. The demonstrations are
479 recorded in both indoor and outdoor environments, with diverse backgrounds and lighting conditions.
480 We observe that policy performance and steering effectiveness decrease under these out-of-domain
481 conditions, but the model still shows a consistent steering trend and retains the ability to identify
482 target objects and complete the tasks.

483 **B Experiment Details**

484 **B.1 Hardware Setup**

485 Our experimental platform consists of a DEEP Robotics Lynx M20 Pro wheeled quadruped integrated
486 with an AgileX Robotics PiPER robotic arm. We use a Stereolabs ZED 2 stereo camera for the
487 third-person view and an Intel RealSense D435i depth camera as a wrist-mounted camera on the
488 manipulator, although depth and stereo measurements are not used in this paper. A Jetson AGX
489 Orin mounted on the robot coordinates communication between the robot hardware and a remote
490 workstation running the policy on an RTX 5080 GPU. For outdoor experiments, network access is

491 provided by USB tethering the Jetson to a mobile phone, which serves as a 5G uplink. The final
492 hardware layout is shown in Figure 7.

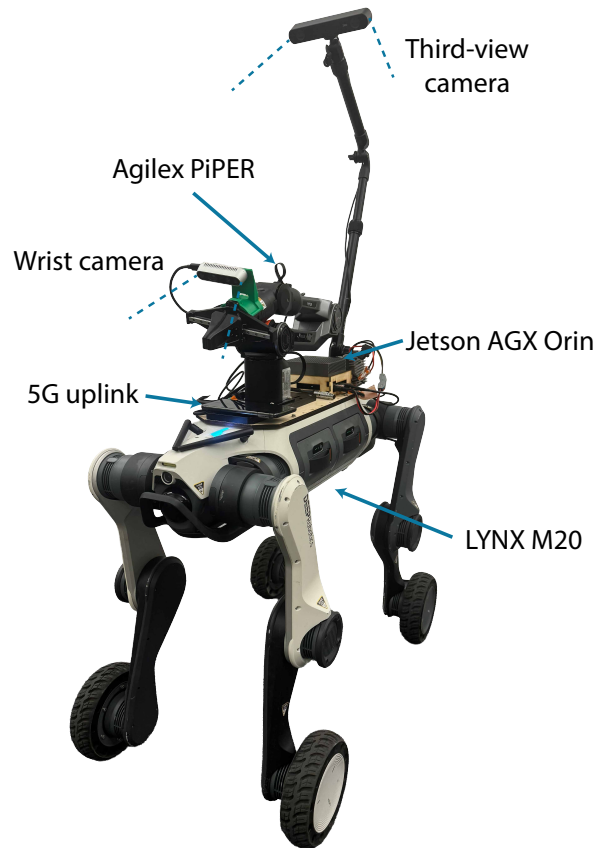


Figure 7: Robot hardware setup.

492

493 The Jetson AGX Orin serves as a lightweight communication bridge rather than an inference device.
494 It captures both camera streams and transmits them to the remote workstation through local network
495 communication. The workstation receives the video streams, runs the policy, and sends action
496 commands back through the same network path to the Jetson, which forwards them to the robot arm.
497 This setup incurs approximately 80ms round-trip latency, modestly higher than wired operation.

498 B.2 Evaluation Setup

499 To evaluate robustness to linguistic variation and compositional generalization, we construct a fixed
500 set of prompt templates covering the three evaluation settings described in the main paper. Prompts are
501 generated compositionally by varying grasp verbs, object descriptors, and placement phrasing while
502 preserving the same underlying pick-and-place task semantics. Representative prompt templates are
503 shown in Table 3.

504 In addition to language variation, we evaluate robustness under diverse spatial configurations and
505 object layouts. Figure 8 shows representative scenes from each evaluation setting. Across config-
506 urations, we vary object placement, relative spacing, distractor positions, and camera viewpoints
507 while preserving the same high-level pick-and-place objective. The seen-object setting reuses training
508 objects in novel arrangements, whereas the unseen-object settings introduce novel targets and
509 increasingly challenging referring expressions. All configurations are held fixed across compared
510 methods to ensure consistent evaluation conditions.

Evaluation Setting	Prompt Templates
Setting I & II	Pick the <descriptor> bag and put it into the box Grab <descriptor> bag and place it in the box Take the <descriptor> bag and drop it into the box Move the <descriptor> bag into the box Put <descriptor> bag into the box Place the <descriptor> bag in the box
Setting III	Put bag with order number <value> into the box Grab bag with number <value> into the box Pick order number <value> bag and place it into the box

Table 3: **List of prompts used for evaluation.** For seen-object evaluations, <descriptor> refers to objects observed during training, i.e., “black” or “plastic.” For unseen-object evaluations, it refers to novel descriptors such as “brown” or “red.” In the novel-language setting, <value> denotes ordinal identifiers, i.e., 1, 2, or 3, used to evaluate compositional reference grounding under previously unseen instruction patterns.

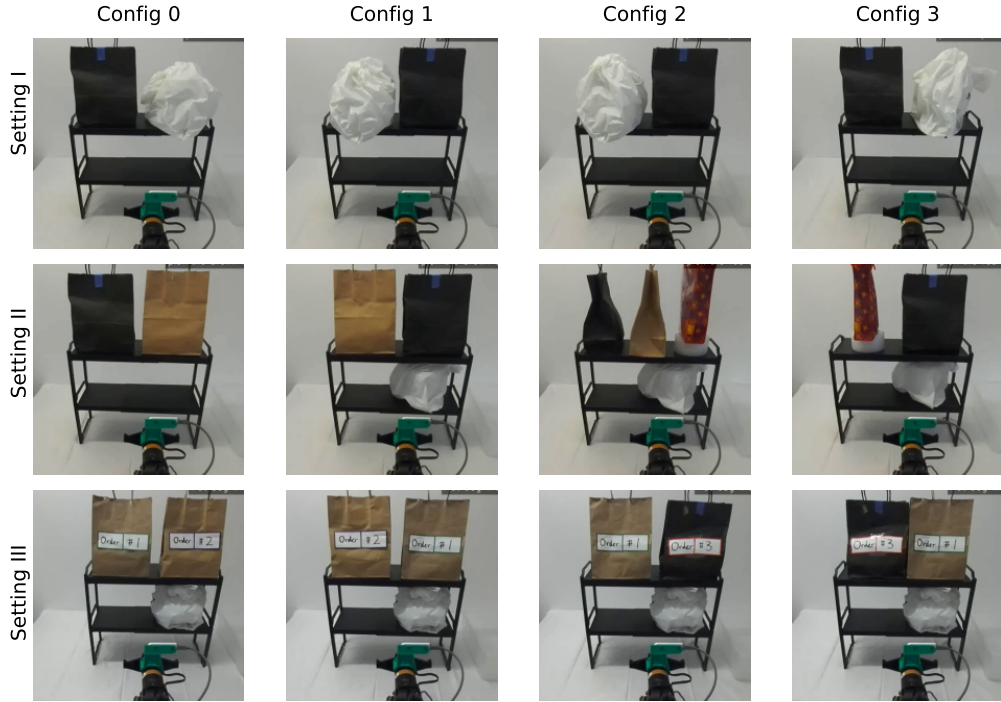
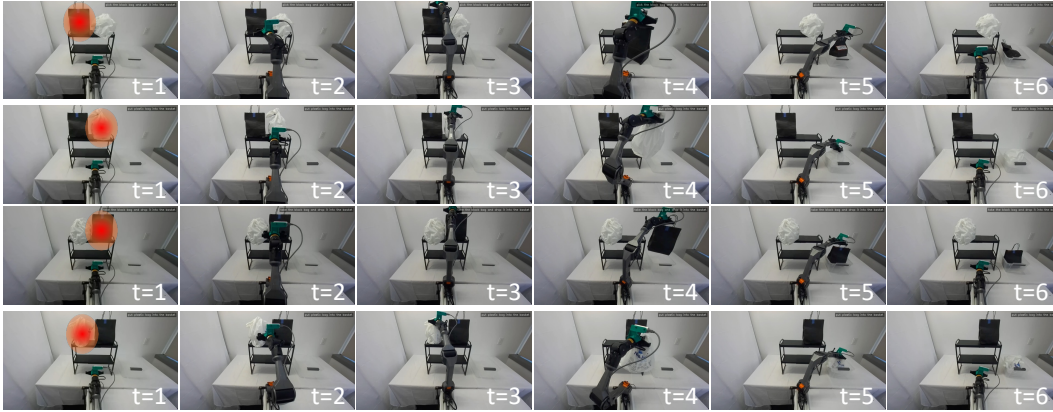


Figure 8: **Representative evaluation configurations across the three experimental settings.** Rows correspond to the evaluation regimes described in the main paper: Setting I uses seen objects with paraphrased language, Setting II uses unseen objects with known language, and Setting III uses unseen objects with novel language. Columns show different spatial layouts used during evaluation, varying object placement, distractor arrangement, and scene geometry. In Setting I, Configurations 0–1 correspond to the black-bag task, while Configurations 2–3 correspond to the plastic-bag task. All methods are evaluated on the same prompts and spatial configurations.

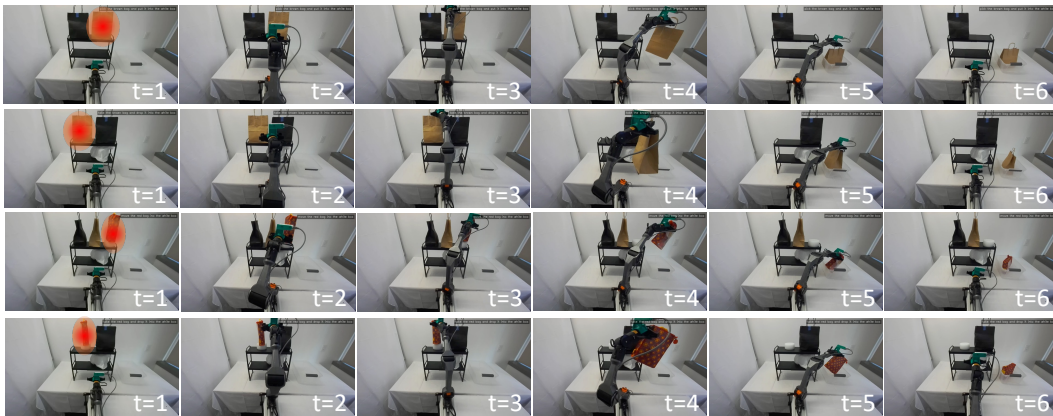
511 C Additional Results

512 Figure 9 provides additional visual comparisons, where the red heatmap indicates the correctly
 513 grounded target object. Under the same spatial object layouts, our method consistently steers the
 514 policy toward the spatial cue across all three evaluation settings. We observe that the policy often
 515 produces similar actions in the early execution steps, but as the end-effector approaches the objects,
 516 the generated trajectory is progressively redirected toward the cued target.

Setting I



Setting II



Setting III

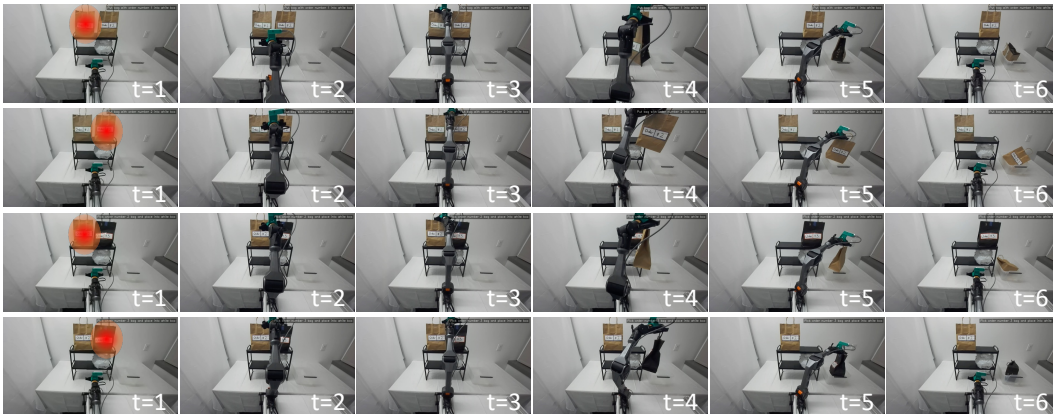


Figure 9: **Example demonstrations of policy steering across the three evaluation settings.** The red heatmap indicates the correctly grounded target object.

517 D Failure Cases

518 Our method can fail in two scenarios. First, it cannot recover when the underlying manipulation prior
519 is unreliable. Noisy rollouts, accumulated execution errors, or environmental changes can still cause
520 the frozen policy to deviate from a feasible manipulation trajectory. Second, the steering strength is
521 applied consistently throughout the rollout, which can over-modulate the policy output when precise
522 actions are required. This may lead to manipulation errors such as overshooting or unstable grasping.

523 **E Video Demonstrations**

524 The supplementary material includes a paper overview video and the corresponding full policy rollout
525 demonstrations.