

# Multi-Modal Manipulation via Multi-Modal Policy Consensus

Anonymous CVPR submission

Paper ID \*\*\*\*

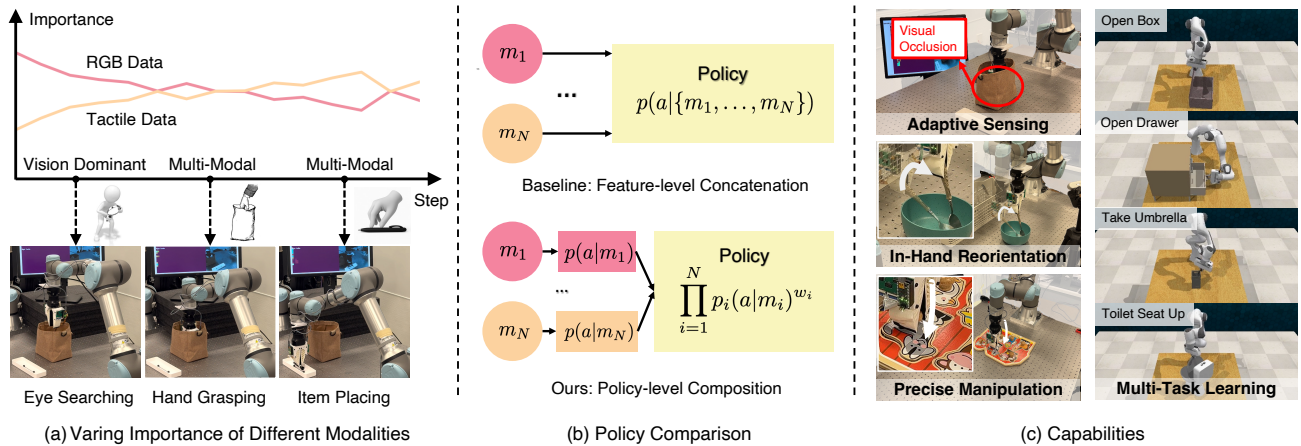


Figure 1. **Representation-Composable Policy.** (a) Perturbation-based importance analysis in the occluded marker picking task shows that vision dominates early, while tactile signals become important once occluded, demonstrating that our framework dynamically utilizes different modalities across task phases. (b) Classical feature concatenation vs. our policy-level composition, where  $m_i$  denotes a modality (e.g., RGB, point cloud, tactile, or learned visual feature). Our compositional design allows individual modality policies to be added or removed without retraining the entire network. (c) Our method unlocks key capabilities. These include *Adaptive Sensing*, retrieving an occluded marker using tactile feedback during occlusion; *In-Hand Reorientation*, reorienting a spoon within the gripper; *Precise Manipulation*, inserting a puzzle piece with fine-grained control; and *Multi-Task Learning*, consistently outperforming prior work across diverse tasks in RL Bench.

## Abstract

001 *Effectively integrating diverse sensory modalities is crucial for robotic manipulation. However, the typical approach of feature concatenation is often suboptimal: dominant modalities such as vision can overwhelm sparse but critical signals like touch in contact-rich tasks, and monolithic architectures cannot flexibly incorporate new or missing modalities without retraining. Our method factorizes the policy into a set of diffusion models, each specialized for a single representation (e.g., vision or touch), and employs a router network that learns consensus weights to adaptively combine their contributions, enabling incremental of new representations. We evaluate our approach on simulated manipulation tasks in RL Bench, as well as real-world tasks such as occluded object picking, in-hand spoon reorientation, and puzzle insertion, where it significantly outperforms feature-concatenation baselines on scenarios requiring multimodal reasoning. Our policy further demonstrates robustness to physical perturbations and sensor corruption.*

*We further conduct perturbation-based importance analysis, which reveals adaptive shifts between modalities.*

019

020

## 1. Introduction

021

022 Modern robots utilize a diverse array of modalities including RGB images, point clouds, tactile signals, and learned visual features, yet effectively integrating these data streams remains a challenge in robotics [6, 25]. These modalities can be both complementary (vision vs. touch) and overlapping (RGB-D vs. point cloud), requiring structured synergy for optimal performance. For decades, a common baseline has been the concatenation of feature-level modalities into a single high-dimensional vector, an approach that persists even in recent policies [18, 25]. Despite its popularity, this approach lacks a principled mechanism for balancing contributions across modalities and cannot easily adapt when modalities are added or missing, often resulting in suboptimal performance.

022

023

024

025

026

027

028

029

030

031

032

033

034

035

036 The weakness becomes apparent when different sensors  
037 are crucial at different task phases. Consider a robot trained  
038 to retrieve an object from an opaque container. For 90% of  
039 the trajectory, vision guides the robot to approach and position  
040 its gripper. Once the gripper enters the container, however,  
041 vision becomes useless while tactile signals become  
042 critical for success. Despite tactile information being essential  
043 during this crucial 10% of the task, feature concatenation  
044 struggles with this sparsity: the learning algorithm  
045 downweights the rarely-active tactile stream as noise, focusing  
046 instead on the statistically dominant visual features. Moreover,  
047 these monolithic approaches cannot adapt to new sensors or  
048 missing modalities without complete retraining. These limitations  
049 point to the need for a structured alternative that treats each  
050 modality as a distinct contributor rather than forcing premature  
051 fusion.

052 To address these challenges, we draw inspiration from  
053 compositional generative models [12, 13, 15, 31] and factorize  
054 robot policies into *modality-specific experts*, with a router  
055 network learning consensus among them. Each expert specializes  
056 in a single modality (e.g., vision for geometric reasoning or  
057 tactile for contact dynamics) and captures its behavioral constraints  
058 [41]. A router network learns consensus weights during training  
059 to combine these experts into a unified policy, ensuring that  
060 even rarely-active but crucial modalities (e.g., tactile) retain  
061 their influence when most needed. As illustrated in Figure 1b,  
062 our approach contrasts with classical feature-level concatenation  
063 by composing experts at the policy level, which not only  
064 mitigates sparsity issues but also allows new modalities to  
065 be added or removed without retraining the entire network.  
066 This compositional design directly resolves the limitations of  
067 prior work by preserving the contribution of each modality,  
068 adapting flexibly to sensor availability, and providing robustness  
069 to failures, while also enabling diverse manipulation skills  
070 such as those shown in Figure 1c.  
071

072 Building on this compositional structure, our framework  
073 enables context-aware shifts in modality reliance: tactile  
074 experts dominate during contact-rich phases, while vision  
075 governs geometric reasoning in free space. The unified policy,  
076 formed by weighted sums of score functions from the factorized  
077 diffusion experts, yields robustness under sensor corruption  
078 and partial observability, ensuring reliable manipulation  
079 across diverse conditions.

080 Our contributions are as follows:

- 081 1. We introduce a framework that composes modality-specific  
082 experts through learned consensus weights, offering a principled  
083 and extensible alternative to monolithic feature concatenation.  
084 This design naturally supports incremental learning, allowing  
085 new experts to be integrated without retraining existing  
086 policies.
- 087 2. We demonstrate strong performance on the multi-task  
088 RL-Bench [22] simulation benchmark compared to vari-

ous baselines and validate our approach on complex real-world  
tasks, including occluded marker picking, spoon reorientation  
in hand, and puzzle insertion. We also show robustness to  
physical perturbations, runtime disturbances, and sensor  
corruption.

3. We provide comprehensive analyses, including  
perturbation-based importance studies, that quantitatively  
demonstrate how our policy learns to shift reliance between  
modalities in response to changing task context.

## 2. Related Works

**Multimodal Fusion in Robotics.** The integration of vision  
and touch is critical for robust robotic manipulation, allowing  
robots to handle occlusions and make precise contact [3, 21, 24].  
Vision provides global scene understanding while tactile sensing  
offers high-fidelity local information critical for tasks like  
grasping and insertion [11, 38]. Existing fusion approaches  
range from simple feature concatenation [25, 26] to sophisticated  
architectures including Visuo-Tactile Transformers that  
dynamically weigh modality features [9], stage-guided fusion  
[16], cross-modal attention mechanisms, and unified modality  
methods [20, 42]. However, all these approaches share a  
common limitation: they perform *feature-level* fusion to create  
a single conditional input for a monolithic policy. This makes  
them vulnerable to sparsity, where learning becomes biased  
toward statistically dominant but contextually irrelevant  
modalities (e.g., vision) while discarding essential but sparse  
signals (e.g., touch). Our work addresses this issue by  
factorizing policies into modality-specific modules and  
composing complete action distributions at the *policy level*.

**Compositional Generation and Energy-Based Models.**  
Our framework builds on compositional modeling, particularly  
the principle of combining simple models as a product of  
distributions to form a more expressive joint distribution [12,  
30, 31, 39]. This idea is central to Energy-Based Models  
(EBMs), where summing energy functions corresponds to  
multiplying probability distributions [13, 17]. Diffusion  
models can be viewed as score-matching EBMs [36], making  
their score functions naturally composable in the same way.  
Beyond image generation, compositionality has been applied  
to trajectory modeling [1, 23], language models [14], robotic  
planning with constraints [32, 41], and hierarchical planning  
[2]. More broadly, compositional generative models extend to  
domains such as traffic generation [29] and human motion  
[35, 37], highlighting the versatility of compositionality  
as a modeling principle. In robotics, related work has  
explored combining policies across heterogeneous sources  
such as simulation and real-world data [40]. In contrast,  
our approach focuses on compositional learning across  
diverse modalities, enabling adaptive integration

141 of modalities such as vision, touch, and semantic features  
142 within a unified policy.

143 **Diffusion Models in Robot Learning.** Diffusion models  
144 have emerged as state-of-the-art policy representations  
145 in robotics, capable of modeling complex, multimodal  
146 action spaces. They have been successfully applied to  
147 imitation learning for single-arm [10, 33, 43] and bimanual  
148 manipulation [7], tool use [8], motion planning [23, 34],  
149 and reinforcement learning [1]. Extensions of diffusion  
150 planning incorporate hierarchical structures [27], and  
151 domain-specific adaptations such as autonomous driv-  
152 ing [28]. Concurrent work by Zhang et al. [4] explores  
153 composition of RGB and point cloud policies using fixed,  
154 manually-tuned weights, but is limited to visual modalities  
155 without adaptive weighting. Our framework introduces a  
156 structured approach to compositional learning that accom-  
157 modates arbitrary modalities, including tactile feedback  
158 and semantic features, through learned context-dependent  
159 routing for adaptive modality weighting.

### 160 3. Approach

161 Our work introduces a compositional approach to multi-  
162 modal robot learning that addresses sparsity in sensory  
163 modalities through policy factorization and consensus-  
164 based composition (see Figure 2). We ground our  
165 framework in energy-based composition principles and  
166 instantiate them via policy consensus across modalities for  
167 robotic manipulation.

#### 168 3.1. Problem Formulation

169 We consider the imitation learning setting for robot manip-  
170 ulation. Given a dataset of expert demonstrations  $\mathcal{D} =$   
171  $\{(s_t, \mathbf{a}_t)\}_{t=1}^T$ , where  $s_t$  is the state and  $\mathbf{a}_t$  the action at  
172 timestep  $t$ , the state comprises  $N$  sensory modalities:

$$173 \quad \mathbf{s}_t = \{m_{1,t}, m_{2,t}, \dots, m_{N,t}\}.$$

174 Our goal is to learn a policy  $\pi(\mathbf{a}_t | \mathbf{s}_t)$  that leverages  
175 heterogeneous information within  $\mathbf{s}_t$ . Each raw modality  
176  $m_{i,t}$  is encoded into a latent embedding  $\mathbf{e}_{i,t}$  through a  
177 modality-specific encoder. Here, an embedding  $\mathbf{e}_{i,t}$  refers  
178 to the modality encoding together with relevant robot state  
179 information (e.g., joint angles, gripper status), providing  
180 a richer context for action prediction. We denote  $\mathbf{a}_t$  as  
181 the ground-truth action in the dataset,  $a$  as a generic  
182 action candidate, and  $a^k$  as its noised version at diffusion  
183 timestep  $k$ . For each modality  $i$ , we train  $K_i$  sub-policies  
184  $p_{i,j}$ , parameterized by  $\theta_{i,j}$ , which define energy functions  
185  $E_{\theta_{i,j}}(a, m_{i,t})$  and corresponding diffusion scores  $\epsilon_{\theta_{i,j}}$ . In  
186 our experiments, we set  $K_i = 2$  to capture complementary  
187 behavioral modes (e.g., geometry vs. fine detail for vision,  
188 contact onset vs. sustained force for tactile), although  
189 the formulation allows general  $K_i$ . A router network  $R_\psi$

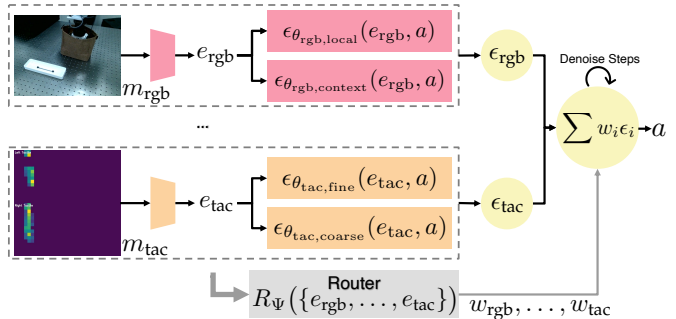


Figure 2. **Overview of Our Compositional Policy Framework.** Raw sensory modalities ( $m_{\text{rgb}}, m_{\text{tac}}$ ) are encoded into embeddings ( $\mathbf{e}_{\text{rgb}}, \mathbf{e}_{\text{tac}}$ ). Each modality is factorized into complementary sub-policies (e.g.,  $\epsilon_{\theta_{\text{rgb,context}}}(e_{\text{rgb}}, a)$ ,  $\epsilon_{\theta_{\text{rgb,local}}}(e_{\text{rgb}}, a)$ ,  $\epsilon_{\theta_{\text{tac,coarse}}}(e_{\text{tac}}, a)$ ,  $\epsilon_{\theta_{\text{tac,fine}}}(e_{\text{tac}}, a)$ ), which produce score predictions that are averaged into a modality-specific score. A router network  $R_\psi(\{\mathbf{e}_{\text{rgb}}, \dots, \mathbf{e}_{\text{tac}}\})$  then predicts consensus weights  $\{w_i\}$  to reconcile these modality-specific scores into the final composed score  $\sum_i w_i \epsilon_i$ , which defines the policy for action generation.

190 maps the embeddings  $\{\mathbf{e}_{i,t}\}$  to consensus weights  $\{w_i\}$ ,  
191 normalized by softmax so that  $\sum_i w_i = 1$ , indicating the  
192 relative influence of each modality.

#### 193 3.2. Energy-Based Policy Composition

194 Our approach builds on energy-based composition princi-  
195 ples [13, 15, 31, 40], where policies are viewed as energy  
196 functions that assign low energy to preferred actions and  
197 high energy otherwise. Each base policy  $p_{i,j}(a | m_{i,t})$  im-  
198 plicitly defines an energy function through

$$199 \quad p_{i,j}(a | m_{i,t}) \propto \exp(-E_{\theta_{i,j}}(a, m_{i,t})).$$

200 Composing such policies corresponds to multiplying dis-  
201 tributions, or equivalently summing their energies:

$$202 \quad p(a | \mathbf{s}_t) \propto \prod_{i=1}^N \left[ \prod_{j=1}^{K_i} \exp(-E_{\theta_{i,j}}(a, m_{i,t})) \right]^{w_i},$$

$$203 \quad = \exp \left( - \sum_{i=1}^N w_i \sum_{j=1}^{K_i} E_{\theta_{i,j}}(a, m_{i,t}) \right). \quad 204$$

205 This reveals that our composition performs a weighted  
206 sum of energy functions, where the router weights  $\{w_i\}$   
207 determine each modality's influence based on the current  
208 state. Unlike feature concatenation, which forces all modal-  
209 ities through a shared network that tends to suppress statisti-  
210 cally rare signals, our formulation preserves separate energy  
211 functions for each modality, ensuring that sparse but critical  
212 signals remain influential.

### 213 3.3. Compositional Policy Factorization

214 We factorize the policy at two levels to capture both inter-  
215 and intra-modality structure:

$$216 p(a|s_t) \propto \prod_{i=1}^N p_i(a|m_{i,t})^{w_i},$$

$$217 p_i(a|m_{i,t}) \propto \prod_{j=1}^{K_i} p_{i,j}(a|m_{i,t}).$$

218 Here  $p_i$  is the composite policy for modality  $i$ , while  $p_{i,j}$   
219 are complementary sub-policies (e.g., vision experts for ge-  
220 ometry vs. fine detail, tactile experts for contact onset vs.  
221 sustained force).

222 This product-of-distributions view admits a constraint  
223 satisfaction interpretation [41], where each modality-  
224 specific policy imposes behavioral constraints on the final  
225 action (e.g., geometry from vision, contact dynamics  
226 from touch). Importantly, unlike feature concatenation  
227 where sparse signals compete with dominant ones for  
228 network capacity, each  $p_i$  is trained independently on its  
229 own modality stream. This ensures that rarely active but  
230 task critical modalities, such as tactile inputs that appear  
231 only during contact, retain their influence through the  
232 learned consensus weights  $w_i$ , which provide a consistent  
233 balancing of modality influence across the policy.  
234

### 235 3.4. Score-Based Implementation via Diffusion 236 Models

237 We implement each base policy  $p_{i,j}$  using Denoising Diffu-  
238 sion Probabilistic Models (DDPMs) [19]. Diffusion mod-  
239 els can be interpreted as score-matching energy-based mod-  
240 els [36], and following compositional generation princi-  
241 ples [13, 31], sampling from a product of distributions cor-  
242 responds to summing their score functions. This leads to a  
243 two-step aggregation process at each denoising step  $k$ :

244 *Intra-Modality Composition.* The composed score for  
245 modality  $i$  is obtained by averaging scores of its  $K_i$  factor-  
246 ized sub-policies:

$$247 \epsilon_{i,\text{comp}}(a^k, m_{i,t}, k) = \frac{1}{K_i} \sum_{j=1}^{K_i} \epsilon_{\theta_{i,j}}(a^k, m_{i,t}, k).$$

248 *Inter-Modality Composition.* The final composed score  
249 is then computed using router-weighted combination:

$$250 \epsilon_{\text{comp}}(a^k, s_t, k) = \sum_{i=1}^N w_i(s_t) \cdot \epsilon_{i,\text{comp}}(a^k, m_{i,t}, k).$$

251 The modality-specific scores  $\epsilon_{i,\text{comp}}$  define gradient  
252 fields that encode each modality’s behavioral constraints.  
253 The router assigns consensus weights  $w_i(s_t)$  to combine

254 these fields into a unified score, balancing their contri-  
255 butions. This weighted composition connects energy  
256 based composition with score based diffusion, establishing  
257 a direct theoretical link between compositional energy  
258 models and diffusion-based policies.

### 259 3.5. Router Network

260 The router  $R_\psi$  maps modality-specific embeddings  $\{e_i\}$   
261 to consensus weights  $\{w_i\}$ , normalized with a softmax so  
262 that they are positive and sum to one, making them inter-  
263 pretable as relative modality influence. Conceptually, the  
264 router reconciles the action proposals of modality-specific  
265 experts into a unified policy. The consensus weights are  
266 learned during training and fixed at execution, providing a  
267 interpretable mechanism for balancing modalities.

### 268 3.6. Advantages Over Existing Fusion

269 Our framework offers (1) **Robustness to sparsity**, as rarely  
270 active modalities such as tactile remain represented by sep-  
271 arate experts rather than being suppressed. (2) **Modularity**,  
272 as experts can be trained and extended independently with-  
273 out retraining the entire policy. (3) **Interpretability**, with  
274 consensus weights  $\{w_i\}$  directly revealing the influence of  
275 each modality. (4) **Principled consensus**, as the router pro-  
276 vides a consistent weighting scheme grounded in energy-  
277 based composition, rather than blind feature fusion.

278 These properties make our framework a theoretically  
279 grounded and extensible alternative to monolithic feature  
280 concatenation for multi-modal robot learning.

## 281 4. Experiment

282 We evaluate our modality-composable policy framework  
283 by addressing three key research questions: (1) How does a  
284 compositional architecture compare to feature-level fusion  
285 in tasks with modality sparsity? (2) Does the compositional  
286 model capture context-dependent and modified reliance on  
287 different modalities, allowing new experts to be composed  
288 without retraining? (3) Does the policy maintain robustness  
289 under physical perturbations and sensor corruption?

### 290 4.1. Experimental Setup

291 **Tasks.** We evaluate our method on both simulation and  
292 real-world manipulation tasks. In simulation, we use RL-  
293 Bench [22] as a multi-task benchmark with four tasks: *open*  
294 *box*, *open drawer*, *take umbrella out of stand*, and *toilet seat*  
295 *up*. The policy is trained on a total of 200 demonstration  
296 episodes and evaluated on 200 unseen configurations.

297 For real-world experiments, we employ a UR5e ma-  
298 nipulator equipped with dual cameras and tactile sensors,  
299 as shown in Figure 3(a). We evaluate three challenging  
300 manipulation tasks, with Figure 3(b–d) showing overlays  
301 of their initial testing conditions: (i) *occluded marker*  
302 *picking*; (ii) *spoon reorientation*; and (iii) *puzzle insertion*.

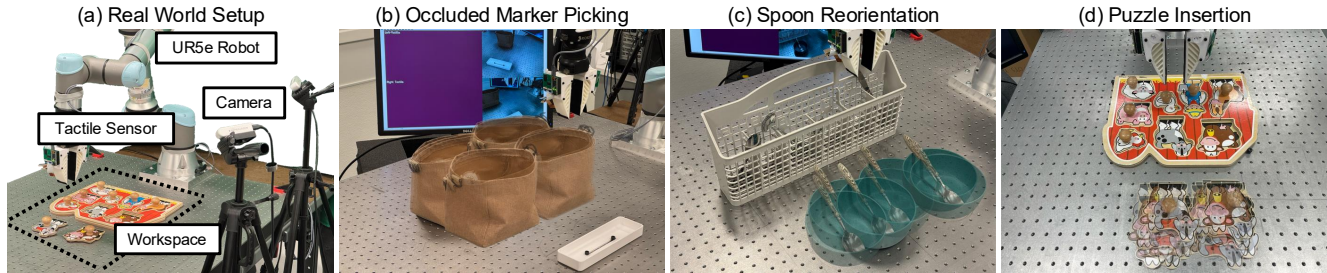


Figure 3. **Real-World Experimental Setup.** (a) UR5e manipulator equipped with dual cameras and tactile sensors. (b–d) Overlays of initial conditions for the evaluation tasks: occluded marker picking, spoon reorientation, and puzzle insertion.

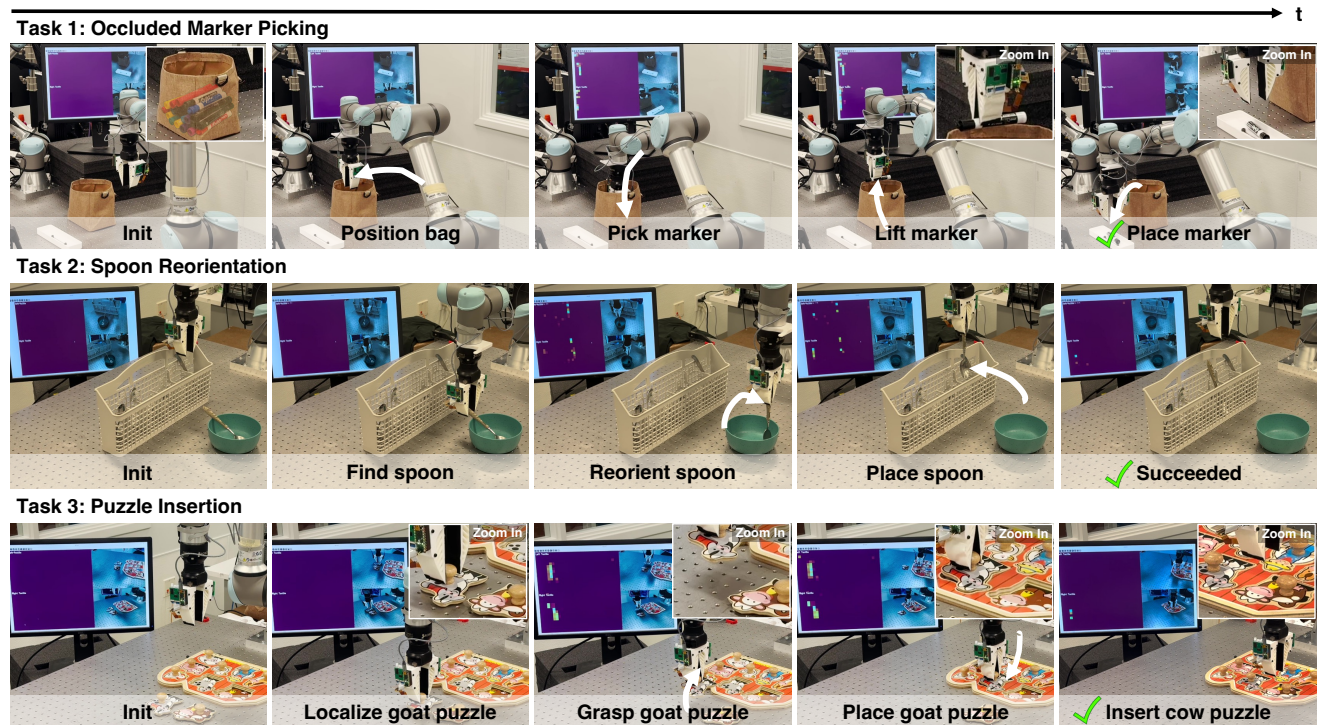


Figure 4. **Qualitative Policy Rollouts.** Representative execution traces from three tasks: Task 1 occluded marker picking, where tactile feedback guides manipulation when vision is unavailable; Task 2 spoon reorientation, demonstrating dexterous in-hand manipulation; Task 3 puzzle insertion, requiring high-precision alignment at millimeter accuracy.

303 We collect 80, 60, and 50 teleoperated demonstrations for  
304 these tasks, respectively.

305 **Sensory Modalities.** In simulation, we utilize RGB  
306 images from two cameras, point cloud (PCD) data from  
307 depth sensors, and 3D semantic features extracted using a  
308 pretrained DINO model [5]. For real-world experiments,  
309 we use RGB images from dual side-view Intel RealSense  
310 D415 cameras with a resolution of  $96 \times 128$ , along with  
311 dense tactile arrays from *FlexiTac* sensors [20]. Each finger  
312 is equipped with a tactile pad consisting of  $12 \times 32$  sensing  
313 units, each with a spatial resolution of 2 mm.

314 **Baselines.** We compare against comprehensive baselines  
315 tailored to each domain:

- **Simulation:** (i) Single-modality policies trained on  
RGB-only, PCD-only, or DINO-only inputs; (ii) Feature  
concatenation combining all modality embeddings; (iii)  
Factorized MoE fusion using soft routing.
- **Real-world:** (i) RGB-only policy; (ii) RGB+Tactile fea-  
ture concatenation baseline.

**Metrics.** We report success rate as the primary metric  
and completion time as a secondary metric for successful  
trials.

## 4.2. Main Results

**Simulation Performance.** Table 1 presents our simula-  
tion results. Our method achieves the highest average suc-

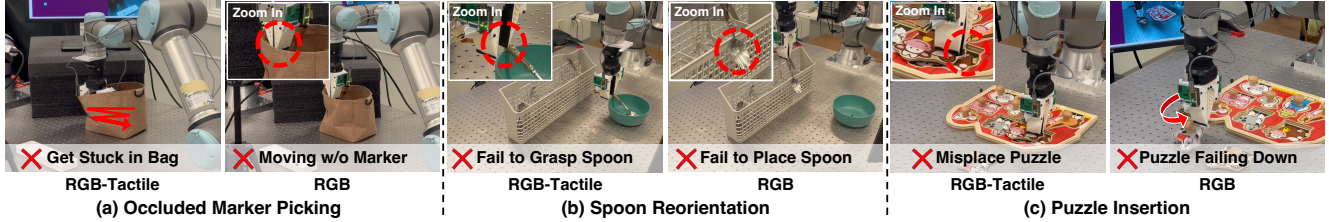


Figure 5. **Typical Failure Cases of Baseline Methods.** We show failure cases of an RGB-only policy compared with an RGB+Tactile concatenation baseline. Each task highlights the complementary roles of the two modalities: vision provides global spatial and geometric information, while tactile sensing provides contact awareness and fine-grained grasp feedback. (a) In occluded marker picking, the concatenation baseline becomes trapped without grasping, while RGB-only lacks awareness of the grasp state once occluded. (b) In spoon reorientation, the concatenation baseline fails at initial grasping, while RGB-only fails at precise placement. (c) In puzzle insertion, the concatenation baseline causes misalignment, while RGB-only suffers frequent grasp failures.

328 cess rate (66%), significantly outperforming both single-  
 329 modality policies and the feature concatenation baseline  
 330 (56%). This 18% relative improvement is achieved with  
 331 minimal parameter overhead (+0.7M parameters, where M  
 332 denotes millions, corresponding to only a 0.3% increase),  
 333 demonstrating the efficiency of our compositional approach  
 334 compared to naive fusion strategies.

Table 1. Policy performance and parameter count on RL-Bench tasks. Our method achieves 18% relative improvement over concatenation baseline with negligible parameter increase.

Method	Success Rate	Params (M)
<i>Single-Modality</i>		
RGB only	0.54	257.3
Point Cloud only	0.49	251.9
3D DINO only	0.45	251.9
<i>Multi-Modality</i>		
Concatenation	0.56	262.9
<b>Ours</b>	<b>0.66</b>	<b>263.6</b>

335 **Real-World Performance.** Tables 2–4 demonstrate  
 336 consistent superiority across all real-world tasks. In con-  
 337 trast to baselines, our method achieves the highest success  
 338 rates (65% for occluded picking, 75% for spoon reorienta-  
 339 tion, and 52% for puzzle insertion) while maintaining the  
 340 lowest average completion times. Figure 4 provides quali-  
 341 tative evidence of these performance differences.

342 Notably, the RGB+Tactile concatenation baseline ex-  
 343 hibits catastrophic failure: in *Occluded Marker Picking*, its  
 344 success rate is only 5% compared to 35% for RGB-only, and  
 345 in *Spoon Reorientation*, it achieves just 21% compared to  
 346 67% for RGB-only. These results confirm that naive fusion  
 347 can be actively detrimental when modalities have varying  
 348 informativeness. Figure 5 illustrates these systematic fail-  
 349 ure modes: the concatenation baseline gets trapped without  
 350 successful grasping in occluded picking, fails at initial  
 351 grasping in spoon reorientation, and causes misalignment

in puzzle insertion, while RGB-only policies suffer from  
 lack of tactile feedback during critical manipulation phases.

Table 2. Decomposed Success Rates on Real-World Puzzle Insertion Task.

Sub-Task	RGB	RGB+Tac.	Ours
Pick up goat	0.90	0.90	<b>1.00</b>
Place goat	0.90	0.90	<b>0.95</b>
Task success (goat)	0.25	0.35	<b>0.58</b>
Pick up cow	0.85	0.80	<b>0.95</b>
Place cow	0.75	0.80	<b>0.95</b>
Task success (cow)	0.30	<b>0.45</b>	<b>0.45</b>

Table 3. Performance on Occluded Marker Picking. Failed trials counted as 120s timeout.

Metric	RGB	RGB+Tac.	Ours
Success Rate	0.35	0.05	<b>0.65</b>
Avg. Time (s)	107.8	117.9	<b>96.5</b>

Table 4. Performance on Spoon Reorientation. Failed trials counted as 300s timeout.

Metric	RGB	RGB+Tac.	Ours
Success Rate	0.67	0.21	<b>0.75</b>
Avg. Time (s)	117.1	221.5	<b>94.8</b>

### 4.3. Analysis of Learned Modality Dependencies

To assess whether the policy leverages different modalities depending on task context, we conduct a perturbation-based importance analysis. Specifically, we measure modality dependency by injecting calibrated Gaussian noise  $\mathcal{N}(0, \sigma^2)$  into each modality and computing the normalized L2 distance between perturbed and original action outputs. For stability, temporal smoothing is applied using an exponential moving average (EMA,  $\alpha = 0.1$ ).

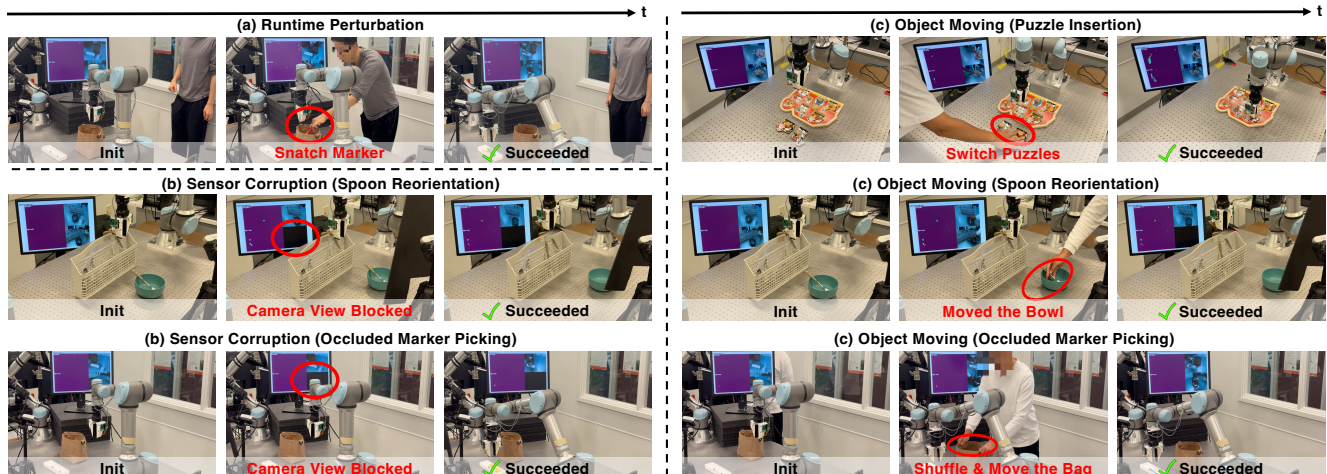


Figure 6. **Policy Robustness under Diverse Perturbations.** We evaluate three types of interventions: (a) runtime perturbation, where the marker is suddenly snatched away during execution; (b) sensor corruption, where a camera is occluded to simulate partial sensor failure; and (c) object repositioning, where task-relevant objects are reset and moved to new positions between executions. Our method maintains reliable performance across all scenarios.

363 As shown in Figure 1a, this analysis reveals context  
 364 aware modality usage during the *Occluded Marker Picking*  
 365 task, which unfolds in two distinct phases. In the first  
 366 stage ( $t = 0$  to 5 steps), vision dominates: the policy shows  
 367 high sensitivity to visual perturbations, relying on visual  
 368 feedback for spatial navigation and approach. In the second  
 369 stage (after 5 steps), the policy shifts to multi modal coordi-  
 370 nation. As the gripper enters the occluded container, tactile  
 371 sensitivity rises sharply upon contact while vision remains  
 372 important, highlighting their complementary roles. This  
 373 emergent two stage behavior demonstrates that our composi-  
 374 tional architecture can transition from single modality  
 375 reliance to multi modal integration based on task demands.

#### 376 4.4. Robustness and Adaptation Evaluation

377 We evaluate the robustness and adaptation of our ap-  
 378 proach through systematic perturbation experiments across  
 379 physical and sensory domains.

380 **Physical Perturbations.** We introduce two categories of  
 381 physical interventions to assess policy resilience under envi-  
 382 ronmental disturbances (Figure 6a, c): (a) *Runtime pertur-*  
 383 *bation*: a mid-execution intervention in which the marker is  
 384 suddenly snatched away during the occluded marker pick-  
 385 ing task, introducing a dynamic disturbance. (c) *Object*  
 386 *repositioning*: objects are deliberately reset and moved to  
 387 new positions between task executions (e.g., relocating the  
 388 bag in the occluded-marker task, displacing the bowl and  
 389 spoon in the spoon-reorientation task, or swapping the po-  
 390 sitions of the goat and cow puzzle pieces).

391 **Sensor Corruption.** We simulate partial sensor failure  
 392 by occluding one camera with an opaque card, eliminat-  
 393 ing visual input (Figure 6b). This corruption is applied to

both spoon-reorientation and occluded-marker tasks. Our  
 method maintains consistent success despite the degraded  
 sensory input, demonstrating resilience to partial sensor  
 failure.

394  
395  
396  
397  
398 **Incremental Learning.** A key benefit of our composi-  
 399 tional framework is the ability to combine independently  
 400 trained policies without retraining. We demonstrate this  
 401 by composing an RGB policy and a tactile policy using  
 402 fixed consensus weights (0.5, 0.5). As shown in Figure 7,  
 403 while the RGB-only policy fails under occlusion, the  
 404 composed policy succeeds in grasping despite never being  
 405 jointly trained. This highlights that our method supports  
 406 incremental integration of new sensory modalities, enabling  
 407 efficient extension of existing skills.

#### 408 4.5. Ablation Studies

409 To validate our key architectural choices, we conduct com-  
 410 prehensive ablation studies examining the impact of differ-  
 411 ent routing and fusion strategies. Table 5 summarizes the  
 412 experimental results.

413 **Learned vs. Fixed Consensus.** We first investigate the  
 414 contribution of our learned consensus mechanism by replac-  
 415 ing it with a fixed equal-weight strategy. This ablation re-  
 416 sults in a significant performance degradation of 7.6% (from  
 417 66% to 61%), demonstrating that adaptive, learned consen-  
 418 sus weighting is important for effective policy composition.

419 **Policy-Level vs. Feature-Level Fusion.** We compare  
 420 our policy-level consensus-based composition against a  
 421 soft-routing MoE baseline that performs fusion at the  
 422 feature level. Our method achieves a substantial improve-  
 423 ment of 15.8% over the MoE baseline (66% vs. 57%; see  
 424 Table 5). This performance gap suggests that composing

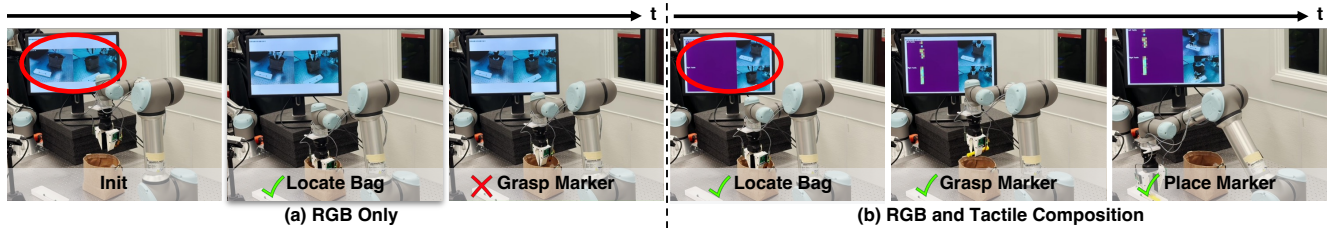


Figure 7. **Incremental Learning.** (a) An RGB-only policy fails to grasp the marker without tactile feedback. (b) By composing a pre-trained RGB policy with a tactile policy using manually set consensus weights (0.5, 0.5), the combined policy successfully grasps the marker under occlusion without requiring retraining.

425 modalities at the action-output level better preserves  
426 modality-specific information compared to early feature  
427 selection, which may prematurely discard relevant signals  
428 needed for downstream decision-making.

Table 5. Ablation study comparing routing and fusion strategies. Results show average success rates across all evaluation tasks in simulation.

Method	Avg. Success Rate
<b>Ours (Learned Router)</b>	<b>0.66</b>
<i>Ablation Variants</i>	
Fixed Equal Weights	0.61
Factorized MoE Fusion	0.57

## 429 5. Conclusion

430 We presented a compositional framework for multimodal  
431 robot manipulation that factorizes policies at the modality  
432 level. The action distribution conditioned on each input  
433 modality is modeled by a separate expert policy, and their  
434 outputs are integrated through learned consensus weights  
435 from a router network. This consensus-based composition  
436 allows the policy to adaptively balance sensing modalities,  
437 preserving their expertise while enabling robust coordination.  
438 Experiments on a multi-task simulation benchmark and contact-rich  
439 real-world tasks demonstrate that our method consistently  
440 outperforms conventional feature-fusion baselines. Beyond these  
441 performance gains, our importance analysis reveals that policies  
442 dynamically shift reliance between modalities based on context,  
443 with vision handling geometric reasoning and tactile managing  
444 contact-rich phases. These compositional advantages enable  
445 incremental sensor deployment without retraining and provide  
446 natural robustness to sensor failures.  
447

448 While promising, our framework opens several directions  
449 for future research. First, the current router adapts  
450 consensus weights at the policy level but does not provide  
451 fine-grained or temporally dynamic adjustments; developing  
452 more expressive consensus mechanisms could enhance

adaptability. Second, our experiments focus primarily on  
vision and tactile inputs in controlled real-world environments,  
leaving extensions to additional modalities (e.g., audio, force,  
language) and deployment in more diverse settings as important  
next steps.

## References

- [1] Anurag Ajay, Yilun Du, Abhi Gupta, Joshua Tenenbaum, Tommi Jaakkola, and Pulkit Agrawal. Is conditional generative modeling all you need for decision-making? *arXiv preprint arXiv:2211.15657*, 2022. 2, 3
- [2] Anurag Ajay, Song Han, Yilun Du, Shuang Li, Abhinav Gupta, Tommi Jaakkola, Joshua Tenenbaum, Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models for hierarchical planning. *arXiv preprint arXiv:2309.08587*, 2023. 2
- [3] Antonio Bicchi and Vijay R. Kumar. Robotic grasping and contact: a review. *Proceedings 2000 ICRA. Millennium Conference. IEEE International Conference on Robotics and Automation. Symposia Proceedings (Cat. No.00CH37065)*, 1: 348–353 vol.1, 2000. 2
- [4] Jiahang Cao, Qiang Zhang, Hanzhong Guo, Jiaxu Wang, Hao Cheng, and Renjing Xu. Modality-composable diffusion policy via inference-time distribution-level composition, 2025. 3
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 5
- [6] Peixin Chang, Shuijing Liu, Haonan Chen, and Katherine Driggs-Campbell. Robot sound interpretation: Combining sight and sound in learning-based control. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, page 5580–5587. IEEE Press, 2020. 1
- [7] Haonan Chen, Jiaming Xu, Lily Sheng, Tianchen Ji, Shuijing Liu, Yunzhu Li, and Katherine Driggs-Campbell. Learning coordinated bimanual manipulation policies using state diffusion and inverse dynamics models. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025. 3
- [8] Haonan Chen, Cheng Zhu, Shuijing Liu, Yunzhu Li, and Katherine Driggs-Campbell. Tool-as-interface: Learning

- 495 robot policies from observing human tool use. In *Confer-*  
496 *ence on Robot Learning (CoRL)*, 2025. 3
- 497 [9] Yizhou Chen, Mark Van der Merwe, Andrea Sipos, and  
498 Nima Fazeli. Visuo-tactile transformers for manipulation.  
499 In *6th Annual Conference on Robot Learning*, 2022. 2
- 500 [10] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric  
501 Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion  
502 policy: Visuomotor policy learning via action diffusion. In  
503 *Proceedings of Robotics: Science and Systems (RSS)*, 2023.  
504 3
- 505 [11] Vedant Dave, Fotios Lygerakis, and Elmar Rueckert.  
506 Multimodal visual-tactile representation learning through  
507 self-supervised contrastive pre-training. *arXiv preprint*  
508 *arXiv:2401.12024*, 2024. 2
- 509 [12] Yilun Du and Leslie Kaelbling. Compositional generative  
510 modeling: A single model is not all you need. *arXiv preprint*  
511 *arXiv:2402.01103*, 2024. 2
- 512 [13] Yilun Du, Shuang Li, and Igor Mordatch. Compositional  
513 visual generation with energy based models. In *Advances in*  
514 *Neural Information Processing Systems*, 2020. 2, 3, 4
- 515 [14] Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-  
516 baum, and Igor Mordatch. Improving factuality and reason-  
517 ing in language models through multiagent debate. *arXiv*  
518 *preprint arXiv:2305.14325*, 2023. 2
- 519 [15] Yilun Du, Conor Durkan, Robin Strudel, Joshua B. Tenen-  
520 baum, Sander Dieleman, Rob Fergus, Jascha Sohl-Dickstein,  
521 Arnaud Doucet, and Will Grathwohl. Reduce, reuse, recycle:  
522 Compositional generation with energy-based diffusion mod-  
523 els and mcmc, 2024. 2, 3
- 524 [16] Ruoxuan Feng, Di Hu, Wenke Ma, and Xuelong Li. Play  
525 to the score: Stage-guided dynamic multi-sensory fusion for  
526 robotic manipulation. In *8th Annual Conference on Robot*  
527 *Learning*, 2024. 2
- 528 [17] Nikolaos Gkanatsios, Ayush Jain, Zhou Xian, Yunchu  
529 Zhang, Christopher Atkeson, and Katerina Fragkiadaki.  
530 Energy-based Models are Zero-Shot Planners for Composi-  
531 tional Scene Rearrangement. In *Robotics: Science and Sys-*  
532 *tems*, 2023. 2
- 533 [18] Johanna Hansen, Francois Hogan, Dmitriy Rivkin, David  
534 Meger, Michael Jenkin, and Gregory Dudek. Visuotactile-  
535 rl: Learning multimodal manipulation policies with deep re-  
536 inforcement learning. In *2022 International Conference on*  
537 *Robotics and Automation (ICRA)*, pages 8298–8304, 2022. 1
- 538 [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising dif-  
539 fusion probabilistic models. In *Advances in Neural Informa-*  
540 *tion Processing Systems*, 2020. 4
- 541 [20] Binghao Huang, Yixuan Wang, Xinyi Yang, Yiyue Luo, and  
542 Yunzhu Li. 3d vitac: learning fine-grained manipulation with  
543 visuo-tactile sensing. In *Proceedings of Robotics: Confer-*  
544 *ence on Robot Learning (CoRL)*, 2024. 2, 5
- 545 [21] Zhe Huang, Ye-Ji Mun, Haonan Chen, Yiqing Xie, Yilong  
546 Niu, Xiang Li, Ninghan Zhong, Ha-II You, David Livingston  
547 McPherson, and K. Driggs-Campbell. Towards safe multi-  
548 level human-robot interaction in industrial tasks. *ArXiv*,  
549 *abs/2308.03222*, 2023. 2
- 550 [22] Stephen James, Zicong Ma, David Rovick Arrojo, and An-  
551 drew J. Davison. Rlbench: The robot learning benchmark &  
552 learning environment, 2019. 2, 4
- [23] Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey  
Levine. Planning with diffusion for flexible behavior synthe-  
sis. In *International Conference on Machine Learning*, 2022.  
2, 3
- [24] Roland S. Johansson and Göran Westling. Roles of glabrous  
skin receptors and sensorimotor memory in automatic con-  
trol of precision grip when lifting rougher or more slippery  
objects. *Experimental Brain Research*, 56:550–564, 2004. 2
- [25] Michelle A. Lee, Yuke Zhu, Krishnan Srinivasan, Parth  
Shah, Silvio Savarese, Li Fei-Fei, Animesh Garg, and  
Jeannette Bohg. Making sense of vision and touch:  
Self-supervised learning of multimodal representations for  
contact-rich tasks. In *2019 International Conference on*  
*Robotics and Automation (ICRA)*, page 8943–8950. IEEE  
Press, 2019. 1, 2
- [26] Hao Li, Yizhi Zhang, Junzhe Zhu, Shaoxiong Wang,  
Michelle A Lee, Huazhe Xu, Edward Adelson, Li Fei-Fei,  
Ruohan Gao, and Jiajun Wu. See, hear, and feel: Smart  
sensory fusion for robotic manipulation. In *Proceedings of*  
*The 6th Conference on Robot Learning*, pages 1368–1378.  
PMLR, 2023. 2
- [27] Wenhao Li, Xiangfeng Wang, Bo Jin, and Hongyuan Zha.  
Hierarchical diffusion for offline decision making. In *In-*  
*ternational Conference on Machine Learning*, pages 20035–  
20064, 2023. 3
- [28] Bencheng Liao, Shaoyu Chen, Haoran Yin, Bo Jiang, Cheng  
Wang, Sixu Yan, Xinbang Zhang, Xiangyu Li, Ying Zhang,  
Qian Zhang, et al. Diffusiondrive: Truncated diffusion  
model for end-to-end autonomous driving. *arXiv preprint*  
*arXiv:2411.15139*, 2024. 3
- [29] Haohong Lin, Xin Huang, Tung Phan-Minh, David S Hay-  
den, Huan Zhang, Ding Zhao, Siddhartha Srinivasa, Eric M  
Wolff, and Hongge Chen. Causal composition diffusion  
model for closed-loop traffic generation. *arXiv preprint*  
*arXiv:2412.17920*, 2024. 2
- [30] Nan Liu, Shuang Li, Yilun Du, Joshua B. Tenenbaum, and  
Antonio Torralba. Learning to compose visual relations. In  
*Advances in Neural Information Processing Systems*, 2021.  
2
- [31] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and  
Joshua B Tenenbaum. Compositional visual genera-  
tion with composable diffusion models. *arXiv preprint*  
*arXiv:2206.01714*, 2022. 2, 3, 4
- [32] Udit Arora Mishra, Shuchen Xue, Yifan Chen, and Danfei  
Xu. Generative skill chaining: Long-horizon skill planning  
with diffusion models. In *Conference on Robot Learning*,  
pages 2905–2925. PMLR, 2023. 2
- [33] Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Li-  
outikov. Goal-conditioned imitation learning using score-  
based diffusion policies. In *Proceedings of Robotics: Sci-*  
*ence and Systems (RSS)*, 2023. 3
- [34] Kallol Saha, Vishal Mandadi, Jayaram Reddy, Ajit Srikanth,  
Aditya Agarwal, Bipasha Sen, Arun Singh, and Madhava Kr-  
ishna. Edmp: Ensemble-of-costs-guided diffusion for mo-  
tion planning. *arXiv*, 2023. 3
- [35] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H  
Bermano. Human motion diffusion as a generative prior.  
*arXiv preprint arXiv:2303.01418*, 2023. 2

- 611 [36] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Ab-  
612 hishek Kumar, Stefano Ermon, and Ben Poole. Score-based  
613 generative modeling through stochastic differential equa-  
614 tions. In *International Conference on Learning Representations*, 2021. 2, 4  
615
- 616 [37] Shanlin Sun, Gabriel De Araujo, Jiaqi Xu, Shenghan Zhou,  
617 Hanwen Zhang, Ziheng Huang, Chenyu You, and Xiaohui  
618 Xie. Coma: Compositional human motion generation with  
619 multi-modal agents. *arXiv preprint arXiv:2412.07320*, 2024.  
620 2
- 621 [38] Sudharshan Suresh, Haozhi Qi, Tingfan Wu, Taosha Fan,  
622 Luis Pineda, Mike Lambeta, Jitendra Malik, Mrinal Kalakr-  
623 ishnan, Roberto Calandra, Michael Kaess, Joseph Ortiz,  
624 and Mustafa Mukadam. Neural feels with neural fields:  
625 Visuo-tactile perception for in-hand manipulation. *Science*  
626 *Robotics*, page adl0628, 2024. 2
- 627 [39] Julen Urain, Anqi Li, Puze Liu, Carlo D’Eramo, and Jan Pe-  
628 ters. Composable energy policies for reactive motion gener-  
629 ation and reinforcement learning. *The International Journal*  
630 *of Robotics Research*, 42(10):827–858, 2023. 2
- 631 [40] Lirui Wang, Jialiang Zhao, Yilun Du, Edward H Adelson,  
632 and Russ Tedrake. Poco: Policy composition from  
633 and for heterogeneous robot learning. *arXiv preprint*  
634 *arXiv:2402.02511*, 2024. 2, 3
- 635 [41] Zhutian Yang, Jiayuan Mao, Yilun Du, Jiajun Wu, Joshua B.  
636 Tenenbaum, Tomás Lozano-Pérez, and Leslie Pack Kael-  
637 bling. Compositional diffusion-based continuous constraint  
638 solvers. In *Proceedings of The 7th Conference on Robot*  
639 *Learning*, pages 3242–3265. PMLR, 2023. 2, 4
- 640 [42] Ying Yuan, Haichuan Che, Yuzhe Qin, Binghao Huang,  
641 Zhao-Heng Yin, Kang-Won Lee, Yi Wu, Soo-Chul Lim, and  
642 Xiaolong Wang. Robot synesthesia: In-hand manipulation  
643 with visuotactile sensing. *arXiv preprint arXiv:2312.01853*,  
644 2023. 2
- 645 [43] Xinyue Zhu, Binghao Huang, and Yunzhu Li. Touch in the  
646 wild: Learning fine-grained manipulation with a portable  
647 visuo-tactile gripper. *arXiv preprint arXiv:2507.15062*,  
648 2025. 3