

Robotic Manipulation by Imitating Generated Videos Without Physical Demonstrations

Anonymous CVPR submission

Paper ID ****

Abstract

001 *This work introduces Robots Imitating Generated Videos*
002 *(RIGVid), a system that enables robots to perform com-*
003 *plex manipulation tasks—such as pouring, wiping, and mix-*
004 *ing—purely by imitating AI-generated videos, without re-*
005 *quiring any physical demonstrations or robot-specific train-*
006 *ing. Given a language command and an initial scene im-*
007 *age, a video diffusion model generates potential demon-*
008 *stration videos, and a vision-language model (VLM) au-*
009 *tomatically filters out results that do not follow the com-*
010 *mand. A 6D pose tracker then extracts object trajec-*
011 *tories from the video, which are retargeted to the robot in*
012 *an embodiment-agnostic fashion. Through extensive real-*
013 *world evaluations, we show that filtered generated videos*
014 *can be as effective as real demonstrations, and that per-*
015 *formance improves with generation quality. We also show*
016 *that relying on generated videos outperforms more com-*
017 *pact alternatives such as keypoint prediction using VLMs,*
018 *and that strong 6D pose tracking outperforms other ways*
019 *to extract trajectories, such as dense feature point tracking.*
020 *These findings suggest that videos produced by a state-of-*
021 *the-art off-the-shelf model can offer a scalable and effec-*
022 *tive source of supervision for robotic manipulation. Project*
023 *page: [rigvid25.github.io](https://github.com/rigid25).*

024 1. Introduction

025 Videos offer a rich and expressive source of training data
026 for robotic manipulation, and numerous methods have suc-
027 cessfully leveraged them for supervision. Such methods
028 typically fall into two categories: (1) Learning from pub-
029 licly available large-scale datasets of real-world videos [9,
030 13, 22, 36, 106, 125], and (2) Imitation of specific demon-
031 strations captured under controlled conditions that closely
032 match the execution setting [8, 21, 55, 65, 69, 114]. Un-
033 fortunately, both of these strategies come with challenges
034 that limit scalability and broad deployment. Large-scale
035 video datasets often introduce domain gaps [36, 119, 134]
036 and require adaptation to specific robot embodiments and

tasks [9, 87]. On the other hand, video-based imitation 037
involves laborious data collection that must ensure close 038
alignment in viewpoints, morphologies, and interaction 039
modalities [7, 8, 26, 106]. 040

Motivated by recent advances in video generation, we 041
explore a potentially new paradigm: can a **single generated** 042
video, synthesized to exactly match our input environment 043
and task description, be used as the sole source of super- 044
vision for robotic manipulation? Recently released models 045
like SORA [16] and Kling [1] have demonstrated impres- 046
sive capabilities in producing realistic-seeming videos from 047
language and image inputs. At the same time, it has been 048
shown that such videos frequently suffer from distorted ob- 049
ject geometries [73, 129], physically implausible interac- 050
tions [83, 124], and unrealistic scene dynamics [11, 39]. 051
Consequently, while the idea of synthesizing supervision is 052
enticing, its usefulness in the robotics setting has not been 053
convincingly established. Prior work incorporating video 054
generation into robotics typically relies on additional super- 055
vision, such as task-specific training [30] or fine-tuning on 056
offline robot trajectory datasets [14, 15]. By contrast, we 057
ask whether a robot can perform real-world manipulation 058
tasks solely by imitating generated videos—without any ad- 059
ditional supervision or task-specific training. 060

To this end, we introduce **Robots Imitating Generated** 061
Videos (RIGVid), a framework that connects video gener- 062
ation models to real-world robotic execution. Fig. 1 shows 063
an outline of the method. Given an input RGB-D image of 064
the scene and a free-form language command (e.g., “pour 065
water on the plant”), we use a state-of-the-art video diffu- 066
sion model to generate a candidate video of the task being 067
performed. The generated video may not accurately follow 068
the language command, but we show that it is possible to 069
use a VLM to automatically filter out unsuccessful genera- 070
tions. Next, we estimate per-frame depth on the video, seg- 071
ment the manipulated object, and track its 6D pose across 072
the video using a pose tracker, FoundationPose [117], that 073
requires a pre-computed object mesh. The resulting *6D ob-* 074
ject pose trajectory serves as a high-level task representa- 075
tion that is retargeted to the robot for execution. Because 076

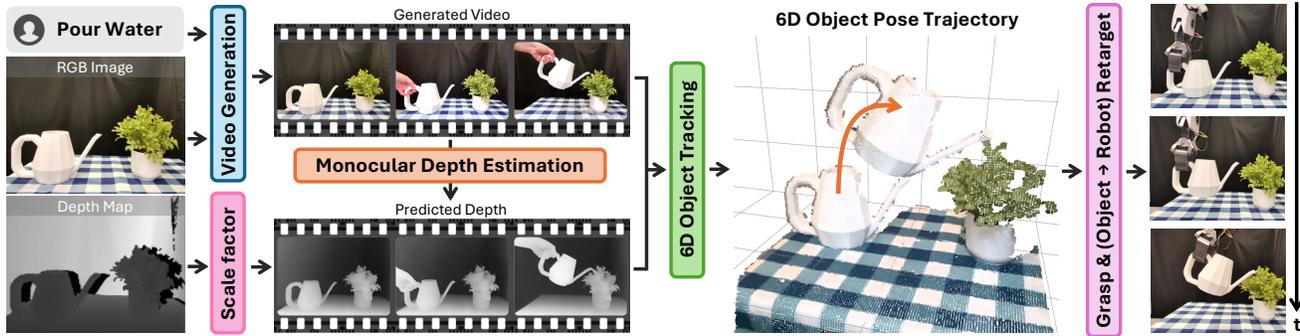


Figure 1. **Robots Imitating Generated Videos.** Given an initial scene image and depth, we generate a video conditioned on a language command. A monocular depth estimator recovers depth for each frame of the generated video, and these depth maps are combined with the corresponding RGB frames to produce 6D Object Pose Trajectory. After grasping, the trajectory is retargeted to the robot for execution.

077 this pose trajectory describes only how the object should
 078 move, rather than specifying robot-specific actions, it can
 079 be directly adapted to other robot platforms. During
 080 deployment, RIGVid also performs real-time object tracking
 081 and dynamically adjusts the robot’s actions to handle distur-
 082 bances and execution-time variations, promoting robust
 083 and adaptive behavior.

084 We evaluate RIGVid on four real-world manipulation
 085 tasks—pouring water, lifting a lid, placing a spatula, and
 086 sweeping trash. These tasks represent a wide range of ma-
 087 nipulation challenges, including minimal vs. significant
 088 depth variation (pouring vs. lifting), thin and partially oc-
 089 cluded objects (placing, sweeping), and different object ge-
 090 ometries and actions. Our results show that, when paired
 091 with our filtering mechanism, generated videos can be as
 092 effective as real human videos for visual imitation, enabling
 093 robots to act entirely from synthetic supervision. Moreover,
 094 the performance of RIGVid improves with video quality,
 095 suggesting a promising trajectory where advances in gener-
 096 ative models directly translate to stronger manipulation
 097 capabilities. While recent advances in VLMs offer a more
 098 compact alternative by predicting high-level task abstrac-
 099 tions, it has remained unclear whether these representations
 100 capture sufficient detail for robust execution, especially
 101 given the substantial computational cost of full video gener-
 102 ation. Our results indicate that generating videos is, in fact,
 103 crucial, yielding higher performance compared to SOTA
 104 VLM-based trajectory prediction method ReKep [49]. To
 105 validate our tracking approach, we compare against a broad
 106 range of existing tracking approaches that reflect the di-
 107 versity of current paradigms—sparse point tracking [15],
 108 dense optical flow [60], 3D feature-fields based pose rea-
 109 soning [58], and generated goal supervision [14]. Despite
 110 our method’s reliance on a pre-computed mesh, this require-
 111 ment consistently yields superior accuracy, confirming its
 112 practical advantage.

113 In summary, this work introduces a generative vision-
 114 driven paradigm for robotic manipulation that scales su-

pervision through synthetic videos. Our key contributions
 are: (1) We propose RIGVid, a framework that enables
 robots to perform open-world manipulation tasks using only
 generated videos, by extracting structured, embodiment-
 agnostic representations that can be directly retargeted for
 execution—without requiring any real-world demonstra-
 tions. (2) We show that RIGVid, when using filtered gen-
 erated videos, performs on par with real human videos,
 demonstrating that high-quality synthetic videos can serve
 as effective substitutes for real-world demonstrations. (3)
 We achieve state-of-the-art performance, outperforming
 representative methods based on VLMs, point tracking, op-
 tical flow, feature fields, and generated-goal supervision.

2. Related Work

Imitation from Videos. Imitation from videos seeks to ac-
 quire robotic skills directly from raw observational data,
 without requiring expert action labels or robot state in-
 formation. This paradigm has attracted significant atten-
 tion [12, 22, 32, 60, 78, 81, 91, 96, 100–102, 107, 114,
 128, 133] because it eliminates the need for precisely la-
 beled robot data. A first line of work focuses on learn-
 ing actionable affordance models from internet-scale video
 datasets [9, 10, 27, 52, 54, 67, 68, 82, 108, 127]. How-
 ever, these methods suffer from domain gap between train-
 ing videos and task-specific environments, and require ad-
 ditional mechanisms to obtain task-conditioned affordances.
 To address this, many methods adopt direct imitation from
 videos, matching visual states in demonstration videos to
 those of the learner [8, 26, 34, 46, 55, 58, 93, 103,
 104, 106, 112, 114, 121, 126]. While effective, this ap-
 proach demands paired demonstrations in the same set-
 ting. A common strategy is to leverage visual correspon-
 dences—tracks [15] or optical flow [5, 35, 122]—to infer
 object trajectories. For example, Bharadhwaj *et al.* [15]
 predicts object tracks and uses PnP to recover poses for
 closed-loop task execution. Dense descriptor learning [33,
 113, 135] has also proven powerful for handling variations
 in object geometry and appearance. Kerr *et al.* [58] recover

258 the language command, we introduce a filtering mechanism
 259 to discard inaccurate generations. We prompt GPT-4o to
 260 assess whether the generated video depicts a successful ex-
 261 ecution of the human command. We sample four evenly
 262 spaced frames in the video and concatenate them vertically
 263 into a single image to create a video summary. The VLM
 264 determines whether this summary depicts a successful ex-
 265 ecution of the task. If the response is negative, we regenerate
 266 the video and repeat the process for up to five attempts. If
 267 all attempts fail, we default to the final attempt. App. B pro-
 268 vides the full prompt used for filtering, examples of video
 269 summaries with their corresponding VLM responses, and
 270 filtering statistics.

271 The resulting filtered video is a plausible execution of
 272 the task, but it is in raw pixel space. To extract action-
 273 able information from this video, we predict its correspond-
 274 ing depth by employing the depth predictor from Ke *et*
 275 *al.* [56]. Although directly using the predicted depth is the
 276 intuitive choice here, we are faced with the challenge of
 277 scale and shift ambiguity [40], where the estimated depth
 278 values are not grounded in real-world units. Consistent
 279 with prior works adopting depth estimators in vision-based
 280 robotics [23, 37], we compute a linear scale-and-shift trans-
 281 formation that aligns the predicted depth in the first frame
 282 with the initial real depth map. This transformation is then
 283 applied to the entire predicted video to resolve scale ambi-
 284 guity.

285 3.3. Identifying Active Object Mask

286 Our next step is to identify the active object—the one being
 287 manipulated in the generated video. We require a binary
 288 mask in the initial RGB image, which is essential both for
 289 object tracking (Sec. 3.4) and for determining which object
 290 to grasp (Sec. 3.5). Given the initial image and the
 291 task command, we prompt GPT-4o [2] to identify which
 292 object in the scene is likely to be manipulated. We then
 293 ground the predicted object category into a bounding box
 294 using Grounding DINO [76], and further refine this bound-
 295 ing box into a segmentation mask using SAM-2 [99].

296 3.4. 6D Object Pose Trajectory

297 We then track the active object, localized by the binary
 298 mask, across the generated video using the scaled predicted
 299 depth. This yields the 6D pose rollout. Tracking objects in
 300 videos is a rich area of research, and we experimented with
 301 several video trackers in 6D pose space [63, 116, 117]. With
 302 the goal of real-world deployment, we found the tracker
 303 from Wen *et al.* [117] to perform the best. It requires an ob-
 304 ject mesh, which we pre-compute using BundleSDF [116].
 305 For this, we record a short RGBD video in which the object
 306 is held in front of the camera and rotated, so that it is ob-
 307 served from all sides. While this process is straightforward,
 308 it does constrain our method to objects for which a mesh
 309 can be precomputed. However, as shown in App. C, our
 310 method is compatible with mesh-free approaches, though

their inference speed is currently infeasible for real-time de-
 311 ployment. To ensure real-world feasibility, we apply an av-
 312 eraging filter to smooth abrupt pose changes, particularly in
 313 the rotational component, to prevent jerky movements. This
 314 refinement stabilizes the object pose trajectory and enables
 315 more realistic executions. App. D provides more details on
 316 pose smoothing. 317

318 3.5. Object to Robot Motion Retargeting

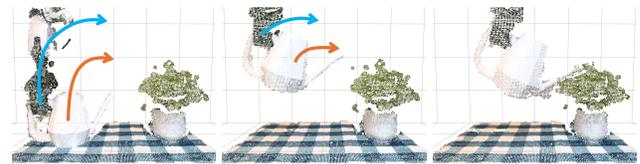


Figure 2. **Re-targeting RIGVid to a robot trajectory.** Assuming a fixed transformation between the end-effector and the object after grasping, the 6D Object Pose Trajectory (*orange arrow*) is re-targeted to the robot (*blue arrow*). This formulation is embodiment agnostic and can be transferred to a different robot.

319 Once the object trajectory is obtained, the first step is
 320 to grasp the object. We use an off-the-shelf grasper, Any-
 321 Grasp [31], to identify and execute the highest-scoring
 322 grasp within a defined boundary around the active ob-
 323 ject mask. After grasping, we retarget its trajectory to
 324 the robot’s end-effector. Since the object remains firmly
 325 grasped, we assume a fixed transformation between the
 326 robot’s end-effector and the object. This transformation is
 327 composed of two components: (1) the pose of the object
 328 relative to the gripper at the moment it is grasped and (2)
 329 the offset between the gripper and the robot’s end-effector.
 330 By combining these two components, we obtain a single
 331 transformation from the end-effector to the object.

332 The corresponding end-effector trajectory is obtained by
 333 applying the fixed end-effector-to-object transformation to
 334 the object’s pose along the entire trajectory. This formu-
 335 lation ensures that the re-targeted 6D pose rollout follows
 336 the object’s motion while maintaining a stable grasp. These
 337 are executed on the physical robot, enabling it to reproduce
 338 the object’s movement as observed in the generated video.
 339 A key strength of this approach is that it is robot-agnostic.
 340 Specifically, to accommodate a different robot or gripper,
 341 only the end-effector to the object transformation needs to
 342 be updated to reflect the new end-effector configuration.

343 4. Experiments

344 We describe our experimental setup, evaluation metrics, and
 345 overview of the baselines. We provide empirical evidence to
 346 answer key research questions about RIGVid: (1) How does
 347 the choice of video generation model impact performance,
 348 and how does the robot perform with real videos of humans
 349 demonstrating these tasks? (2) How does RIGVid compare
 350 with VLM-based trajectory prediction methods that allow

351 zero-shot robot executions? (3) How does RIGVid compare
352 to the most relevant visual imitation methods?

353 4.1. Robot Setup, Tasks, and Evaluation

354 We conduct experiments on an xArm7 robot arm with a
355 stationary Orbbec Femto Bolt camera, positioned next to
356 the robot to capture RGBD observations. We evaluate our
357 method on four everyday manipulation tasks that together
358 span a diverse range of robotic challenges:

- 359 1. **Pouring water** requires the robot to position and tilt a
360 watering can above a plant. While the depth of the can
361 from the camera remains largely constant, successful
362 execution demands a smooth trajectory spanning the
363 pick-up, transport, and pouring phases. A trial is con-
364 sidered successful if the spout of the watering can is
365 positioned above the plant at the end of the execution.
- 366 2. **Lifting a lid** requires the robot to lift a pot lid. Unlike
367 pouring, this task involves significant variation in ob-
368 ject depth, as the lid moves closer to the camera during
369 execution. This task assesses the method’s adaptability
370 to changing object-to-camera distances and the robust-
371 ness of trajectory extraction under substantial depth
372 shifts. Success is achieved if the lid is no longer in
373 contact with the pot at the end of the trial.
- 374 3. **Placing a spatula on a pan** requires the robot to
375 place the head of a spatula into a pan. The spatula
376 presents a thin, elongated geometry and is often par-
377 tially occluded during manipulation, which presents a
378 challenge for object tracking. This task evaluates the
379 method’s ability to handle objects with small surface
380 area and persistent occlusion, both of which are par-
381 ticularly difficult for non-mesh-based approaches. The
382 task is considered successful if the spatula’s head is in
383 the pan at the end of execution.
- 384 4. **Sweeping trash** requires the robot to sweep trash into
385 a dustpan. This task is especially challenging as it
386 combines the need for precise positioning to align the
387 head of the sweeping brush with the trash, along with
388 all the challenges encountered from the previous task.
389 A trial is successful if the trash is touching the base of
390 the dustpan at the end of the execution.

391 Task success is determined via human judgment based
392 on these criteria, though the procedure could be readily au-
393 tomated with a VLM. The initial setup is fixed across tri-
394 als of the same task and each trial has a different gener-
395 ated video. All baselines use the same videos for consist-
396 ent comparison and reporting. We run all experiments on a
397 single Nvidia TitanX GPU machine with 32 GB RAM.

398 4.2. Quality and Filtering of Generated Videos

399 As mentioned in Sec. 3.2, we experimented with Sora,
400 Kling v1.5, and Kling v1.6 and compared different video
401 filtering strategies. Next, we summarize our key empirical
402 findings.

403 *How do different video generation models compare for*

robotic imitation? Sora is known for creating visually im-
404 pressive and cinematic videos. However, these videos often
405 prioritize aesthetics over following the human command.
406 For example, as seen in the top row of Fig. 3, Sora fre-
407 quently alters the camera viewpoint, changes object posi-
408 tions, or even swaps out objects mid-sequence. This lack
409 of scene and object consistency makes Sora poorly suited
410 for imitation. Kling v1.5 places more emphasis on follow-
411 ing language instructions. In our evaluations, videos from
412 Kling v1.5 generally preserved the original scene and cor-
413 rectly depicted the target object. Nonetheless, it still ex-
414 hibited physically implausible behaviors, such as objects
415 moving in unnatural ways or actions that defy basic phys-
416 ical constraints. These issues, although less frequent than
417 in Sora, still prevent successful downstream robot execu-
418 tion. In addition, we frequently observed that Kling v1.5
419 would fail to follow the commanded instruction at all, noth-
420 ing happens in the video, and the intended manipulation is
421 simply not attempted. Kling v1.6 further improved com-
422 mand following and physical plausibility. Videos gener-
423 ated by it were the most consistent with the initial scene and the
424 intended task. As shown in the bottom row of Fig. 3, Kling
425 v1.6 avoids altering the scene layout, maintains the posi-
426 tions of all objects, and depicts motions that are physically
427 reasonable and closely aligned with the human command.
428 Hence, Kling v1.6 proved to be the most reliable video gen-
429 erator for us.

Does higher video quality lead to better robot perfor-
430 *mance?* To quantify this, Fig. 4 plots RIGVid’s task success
431 across five video sources: unfiltered Sora, unfiltered Kling
432 v1.5, unfiltered Kling v1.6, filtered Kling v1.6, and real hu-
433 man demonstration videos. For each source, we used 10
434 videos per task. We observe a clear trend: as video qual-
435 ity improves, so does success rate. Sora’s videos led to the
436 lowest success, Kling v1.5 performed better, and Kling v1.6
437 gave the highest results among all generated videos. Filter-
438 ing further improved reliability: by discarding failed gen-
439 erations using our automatic approach, performance with
440 filtered Kling v1.6 videos matched that obtained with real
441 demonstration videos.

How effective is automatic video filtering? Filtering suc-
442 cess rates varied by task: 83% for pouring, 66% for lift-
443 ing, 55% for placing, and 45% for sweeping. Most errors
444 were false negatives—i.e., the filter occasionally discarded
445 some usable videos, but almost never passed an incorrect
446 one. Compared to standard metrics like video-text consis-
447 tency and subject consistency from VBench++ [50, 115],
448 our VLM-based filter correlated much more strongly with
449 true task completion (Tab. 1).

Can generated videos replace real videos for imitation?
450 The results in Fig. 4 indicate that, when using filtered
451 Kling v1.6 videos, RIGVid’s performance is similar to that
452 achieved with real human demonstration videos. This find-
453
454
455
456



Figure 3. **Qualitative comparison of video generation.** Sora (top) drastically alters the scene layout and objects. Kling v1.5 (middle) is better but exhibits physically-implausible interactions. Kling v1.6 (bottom) produces the most consistent and realistic videos.

457 ing suggests that, at current model quality, generated videos
458 are already sufficient for visual imitation, substantially re-
459 ducing the need for manual data collection.

460 *What causes failure?* Aside from one case where the ob-
461 ject slipped out of the gripper, all failures are attributed to
462 errors in monocular depth estimation. These errors result in
463 inaccurate 6D trajectories and lead to tracking failures. No-
464 tably, similar depth estimation issues are also observed in
465 real videos, suggesting that the limitation lies in the depth
466 model itself. App. I provides a detailed analysis with quali-
467 tative examples.

468 4.3. RIGVid vs. VLM-based Trajectory Prediction

469 Video generation is computationally expensive, prompt-
470 ing the question of whether more efficient alternatives

471 can enable robot manipulation without any demonstrations.
472 VLMs offer one by predicting simplified task abstrac-
473 tions—goal states [48], constraints [49], or reward func-
474 tions [92]—without generating full visual sequences, mak-
475 ing them cheaper in computation and inference time. We
476 compare against the state-of-the-art VLM-based method
477 ReKep [49] in Fig. 5, where RIGVid achieves 85% vs. 50%
478 success over four tasks. App. F illustrates ReKep’s fail-
479 ures, which we attribute to inaccurate keypoint predictions.
480 This comparison suggests that while VLM-generated abstrac-
481 tions are compact, they may lack the rich, necessary
482 details. Thus, despite its higher cost, video generation pro-
483 vides crucial supervision rather than being a wasteful ex-
484 pense.

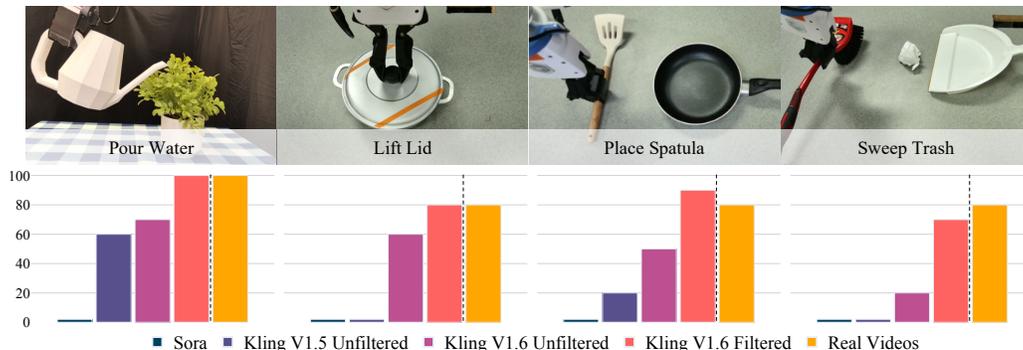


Figure 4. **RIGVid performance vs. video quality.** The dashed lines separate performance on generated videos from real videos. Kling V1.6 produces most reliable videos and leads to highest RIGVid success. Filtered videos yields performance on par with real videos.

Filtering Metrics	Description	Pour Water	Lift Lid	Place Spatula	Sweep Trash	Average
Video-text Consistency	Text-video match	0.06	0.47	0.70	0.11	0.34
I2V Subject Consistency	Image-video subject match	0.35	0.58	-0.09	0.63	0.37
Querying GPT o1	VLM-based filtering	0.91	0.91	0.91	0.66	0.84

Table 1. **Comparison of video filtering metrics.** Pearson correlation coefficients measure each metric’s effectiveness in assessing whether a generated video follows the language command. Each task has 10 success and failure cases. Querying GPT o1 proves to be most effective.

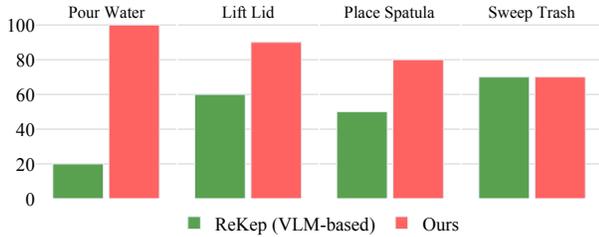


Figure 5. **RIGVid vs. VLM-Based Trajectory Prediction.** By leveraging rich spatial and temporal structure, RIGVid provides accurate execution for challenging manipulation tasks.

485

4.4. Baseline Tracking Methods

486

We adapted several well-studied and competitive baselines that use different types of tracking for visual imitation to be able to work *without any demonstrations*. For each baseline, we describe its inputs and outputs, core approach, our modifications, and the motivation for its inclusion. We describe these baselines below and defer further details to App. E.

487

488

489

490

491

492

Track2Act [15] (Tracks-Based). It takes as input RGBD image of the initial scene, together with a goal image that specifies the desired final configuration. Since we have no way to get the goal configuration, we generate its goal image by using the last frame of the generated video. The method predicts a dense grid of 2D point tracks between the initial and goal image to estimate pixel-wise correspondences. These tracks are then lifted to 3D using the depth map from the initial frame and converted into a sequence of 3D object poses via the Perspective-n-Point (PnP) algorithm. All other components of the method are kept unchanged. Track2Act is notable for its use of a dedicated track prediction network, which is conceptually similar to approaches that predict affordances, but here enables object motion to be inferred directly from observation pairs.

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

AVDC [60] (Flow-Based). Given an initial image and a task description, AVDC predicts object motion by first generating a task-conditioned video and then computing optical flow between frames. This optical flow is used in an opti-

508

509

510

mization process to reconstruct the object trajectory. In our adaptation, we substitute AVDC’s original video generator with our improved model, while preserving all downstream processing. This method is compelling because it leverages dense correspondences across all points on the object, providing more correspondences for tracking.

511

512

513

514

515

516

4D-DPM [58] (Feature Field-Based). We modify this method from tracking object part poses to tracking single objects. It takes as input a 3D Gaussian splatting field of the object and a video of the task, and outputs estimated object trajectories over time. To build the field, this requires a separate video where object is captured from all sides. In our adaptation, since 4D-DPM typically expects a real human demonstration video, we instead use a generated video as the task input video. This approach is compelling because it applies geometric, feature-based reasoning to track objects, capturing entire object structure from video, without relying on explicit correspondences.

517

518

519

520

521

522

523

524

525

526

527

528

Gen2Act [14] (Generated Goal-Based). Gen2Act takes as input an RGBD image of the scene and a task description, and outputs robot actions predicted by a learned policy. The method first generates a video of the task as an intermediate step, then extracts object tracks from this video. These tracks, together with offline robot demonstrations, are used to train a policy for robot execution. In our adaptation, we remove the fine-tuning on real robot demonstrations and instead directly estimate object poses from the video tracks, eliminating any dependence on expert demonstration data. Gen2Act is notable for leveraging sparse correspondences extracted from the generated video, enabling task-relevant object motion to be tracked and retargeted without requiring explicit actions.

529

530

531

532

533

534

535

536

537

538

539

540

541

542

4.5. RIGVid vs Other Robotics Tracking Methods

543

In the following, we include takeaways based on the results in Fig. 6 and Fig. 7.

544

545

How does RIGVid compare to prior works that use other tracking methods, and what accounts for its advantage? Fig. 6 reveal that RIGVid achieves a success rate of 85.0%, compared to 67.5% for Gen2Act and considerably lower

546

547

548

549

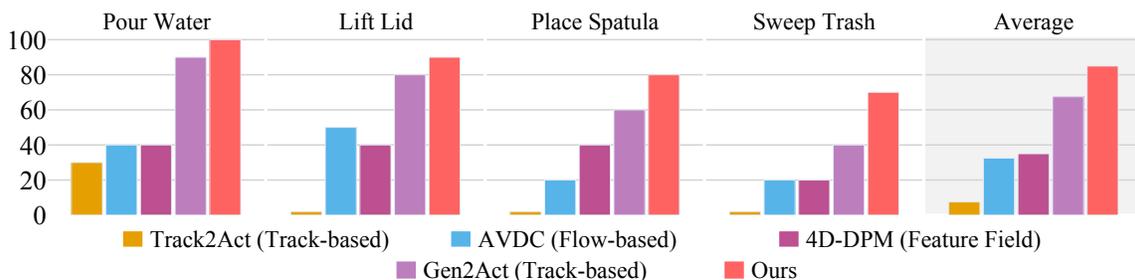


Figure 6. **Performance on everyday manipulation tasks.** Shown in the table, RIGVid, which uses 6D Object Pose Trajectory, consistently achieves higher success rates (*bottom*) across all four tasks (*top*). Relative improvements are higher as tasks become harder (*i.e.*, from left to right in the bar plot). All results are on valid video generations (*i.e.*, human filtered; detailed study in Sec. 4.2) and ten episodes for each reported metric.

550 rates for all other baselines. This margin grows on the
 551 more complex tasks. Crucially, all approaches are evalu-
 552 ated using the same set of generated videos, isolating the ef-
 553 fect of the trajectory representation itself. Methods such as
 554 Track2Act (7.5%), AVDC (32.5%), and 4D-DPM (35.0%)
 555 rely on point tracks or optical flow, but their performance re-
 556 mains limited—especially as object rotations or occlusions
 557 become severe. Gen2Act, which combines video genera-
 558 tion with point-based tracking, closes part of the gap but
 559 consistently struggles when large portions of the object be-
 560 come untrackable. In contrast, RIGVid’s use of a structured
 561 6D object pose trajectory enables robust execution across all
 562 tasks, accounting for the observed 17.5% absolute improve-
 563 ment over Gen2Act. This advantage persists even when
 564 more powerful tracking models like CoTracker3 [53] are
 565 used, as shown in App. G. These results indicate that it is
 566 the accuracy and stability of the 6D object pose trajectory
 567 that is key to RIGVid’s stronger performance across a range
 568 of manipulation tasks.

569 *How does RIGVid perform on tasks involving depth vari-*
 570 *ation, small objects, and occlusion?* Looking at the task-
 571 wise breakdown in Fig. 6, we find that RIGVid consistently

572 maintains high success rates even as object depth varies sig-
 573 nificantly (such as in the lifting task) or when the manipu-
 574 lated objects are thin, small, or partially occluded (such as
 575 in placing a spatula or sweeping trash). Other methods fre-
 576 quently struggle in these settings, often failing to recover ac-
 577 curate object trajectories when objects become partially hid-
 578 den or change distance rapidly. The advantage of RIGVid
 579 is most pronounced on the most challenging tasks: for both
 580 spatula placement and sweeping, RIGVid achieves success
 581 rates 20–25% higher than the next best baseline. These re-
 582 sults suggest that the structured 6D pose trajectory not only
 583 enables robust tracking under depth changes and occlusion,
 584 but also scales to manipulation scenarios where traditional
 585 correspondence-based methods break down.

586 *What do the intermediate predictions reveal about these*
 587 *methods?* Visualizing the outputs in Fig. 7 for the same gen-
 588 erated video, we observe the intermediate predictions and
 589 resulting robot executions produced by each method. For
 590 Track2Act, the predicted tracks diverge from the true object
 591 path, leading to failed execution. AVDC generates reason-
 592 able optical flow in individual frames, but when summed
 593 across the entire video, the resulting trajectory is often phys-

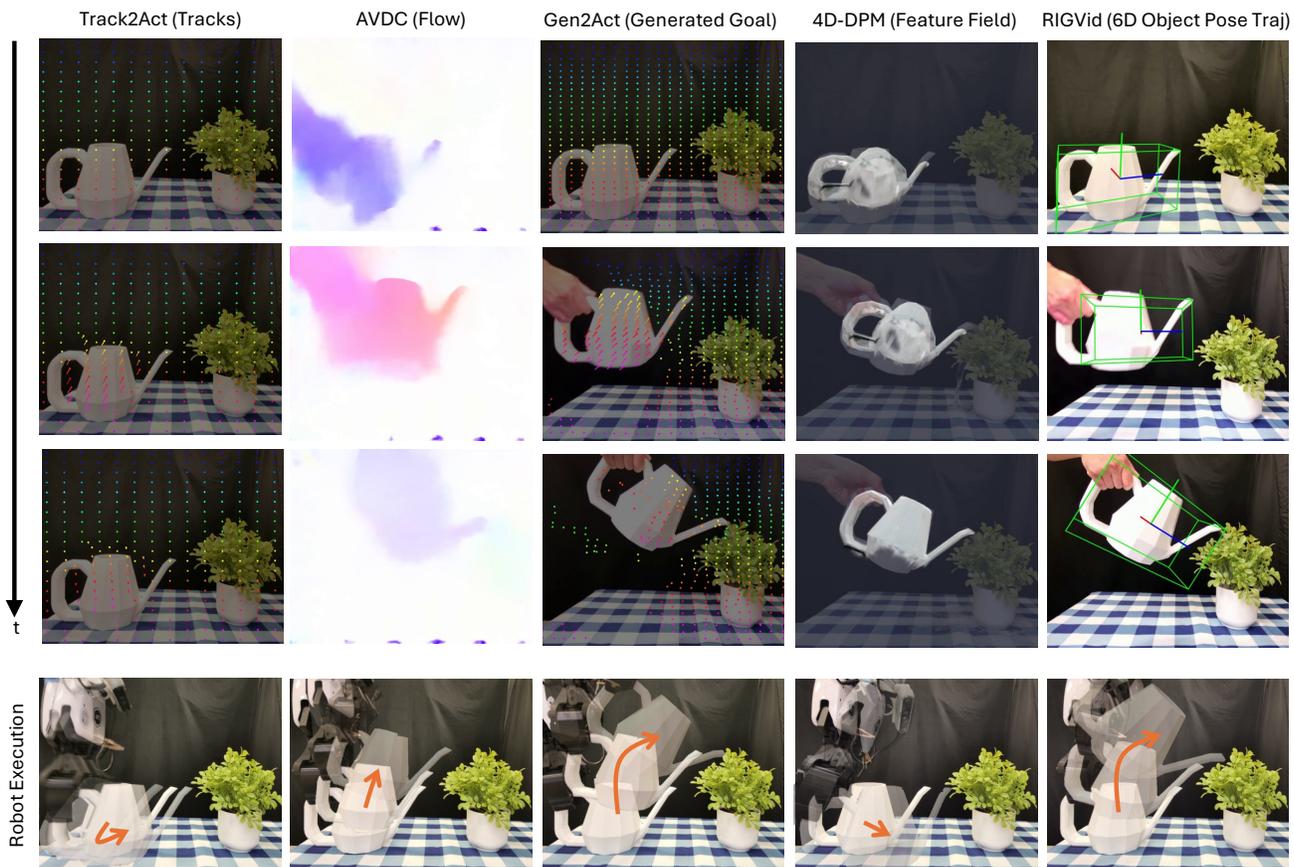


Figure 7. **Analyzing intermediate visual representations.** Our 6D Object Pose Trajectory can correctly track the position and rotation of the watering can (rightmost column), leading to a successful execution.

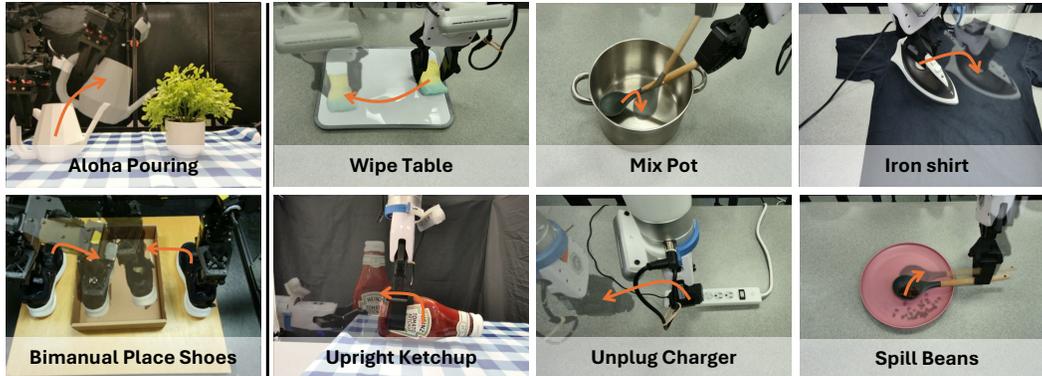


Figure 8. **RIGVid’s embodiment-agnostic capabilities and examples on solving complex, open-world tasks.** RIGVid can readily work on ALOHA setup [131] as shown on top left. On the bottom left, RIGVid is retargeted to the ALOHA setup. On the right, it generates trajectories for diverse manipulation tasks—including wiping, mixing, and ironing—without using any physical demonstrations.

594 ically implausible and the execution fails. Gen2Act yields
595 plausible tracks and leads to successful manipulation. 4D-
596 DPM exhibits inconsistent tracking performance. While it
597 accurately follows the object in certain segments, the ex-
598 ample shown reveals incorrect tracking during the first half
599 of the episode, which ultimately causes the rollout to fail.
600 In contrast, the 6D object pose trajectories produced by
601 RIGVid remain stable throughout the episode and closely
602 match the actual object motion, resulting in successful ex-
603 ecution.

604 4.6. Testing Generalization and Robustness

605 **Embodiment-Agnostic Transfer.** We test RIGVid’s gen-
606 eralizability to another embodiment by deploying it on the
607 ALOHA robot for the pouring task (Fig. 8, top left). On this
608 setup, it achieves 80% success, compared to 100% on our
609 default xArm setup.¹ RIGVid also generalizes to a biman-
610 ual setup, simultaneously placing a pair of shoes into a box
611 using both arms (Fig. 8, bottom left).

612 **Extensions to New Tasks.** RIGVid enables zero-shot
613 execution of diverse and challenging manipulation tasks
614 that involve complex and unconstrained trajectories. As
615 shown in Fig. 8 (right), it completes tasks such as wip-
616 ing, mixing, and ironing. It also handles physically intricate
617 scenarios, including uprighting a ketchup bottle, rotating a
618 spoon to spill beans, and unplugging a charger, despite these
619 tasks involving extreme rotations.

620 **Recovery from Perturbations.** A key strength of
621 RIGVid is its robustness to external disturbances during ex-
622 ecution, as shown in Fig. 9. The system continuously tracks
623 the object’s position using f_{track} and updates the robot’s end-
624 effector trajectory in real time. To detect a deviation (im-
625 age 1), the current object pose is compared to the precom-
626 puted motion plan. If it strays beyond 3 cm in position or 20
627 degrees in orientation (image 2), the system classifies it as a
628 disturbance. The robot then backtracks to the last success-

fully executed trajectory point (image 3) and resumes the
629 planned motion (image 4). This realignment allows RIGVid
630 to maintain stable task execution even under physical per-
631 turbations. Additional robustness examples are discussed in
632 App. H. 633

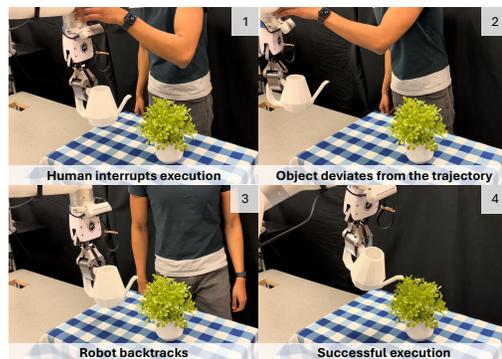


Figure 9. **RIGVid is robust to perturbations.** The robot detects and recovers from external disturbances by backtracking to the last achieved pose before resuming execution.

634 5. Conclusions and Limitations

635 We introduced Robots Imitating Generated Videos
636 (RIGVid), the first method for robotic manipulation that
637 works without demonstrations—no teleoperation, no hu-
638 man demonstration, or expert policy rollouts. By leveraging
639 recent advances in generative vision models and 6D pose
640 estimation, RIGVid enables robots to execute complex
641 tasks entirely from generated video. We extract 6D Object
642 Pose Trajectory from the generated videos and retarget
643 it to the robot, demonstrating a scalable, data-efficient
644 approach to robotic skill acquisition. Our analysis shows
645 a clear correlation between video quality and task success:
646 as generation improves, RIGVid approaches real demo
647 performance. Additionally, our comparisons with SOTA
648 VLM-based zero-shot manipulation methods confirm that
649 leveraging detailed visual and temporal cues from gener-
650 ated videos yields substantially superior performance. We

¹The slight performance drop stems primarily from camera calibration challenges, as ALOHA’s arms yield less accurate pose estimates.

651 also show that RIGVid significantly outperforms baselines
652 across a diverse set of visual imitation tasks, and demon-
653 strate the robustness of our approach to environmental
654 disturbances. Our work represents a step toward enabling
655 robots to learn from the vast visual knowledge in generative
656 models, reducing reliance on costly and time-consuming
657 real-world data collection.

658 Despite the advancements, our method has certain limi-
659 tations. We need a precomputed mesh of the objects. In the
660 future, it can be simplified by using single-image to mesh
661 reconstruction methods [74, 75]. Additionally, our ap-
662 proach depends on the video generation quality, which may
663 struggle with complex prompts or scenes. As video gener-
664 ation improves, we anticipate this limitation will become
665 less significant. Our work aims to democratize robotics by
666 removing the need for demonstrations, which could enable
667 broader accessibility of robotic capabilities. However, we
668 also acknowledge potential risks. Highly generalizable
669 manipulation methods could be misused in applications
670 such as automated weapons or unsafe industrial automation.
671

672 References

673 [1] Kling ai. <https://www.klingai.com/>. Accessed:
674 2024-02-10. 1, 3

675 [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
676 Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo
677 Almeida, Janko Altenschmidt, Sam Altman, Shyamal
678 Anadkat, et al. Gpt-4 technical report. *arXiv preprint*
679 *arXiv:2303.08774*, 2023. 4

680 [3] Anurag Ajay, Seungwook Han, Yilun Du, Shuang Li, Abhi
681 Gupta, Tommi Jaakkola, Josh Tenenbaum, Leslie Kael-
682 bling, Akash Srivastava, and Pulkit Agrawal. Composi-
683 tional foundation models for hierarchical planning. *Ad-
684 vances in Neural Information Processing Systems*, 36:
685 22304–22325, 2023. 3

686 [4] Mert Albaba, Chenhao Li, Markos Diomataris, Omid
687 Taheri, Andreas Krause, and Michael Black. Nil: No-data
688 imitation learning by leveraging pre-trained video diffusion
689 models. *arXiv preprint arXiv:2503.10626*, 2025. 3

690 [5] Max Argus, Lukas Hermann, Jon Long, and Thomas Brox.
691 Flowcontrol: Optical flow based visual servoing. In *2020*
692 *IEEE/RSJ International Conference on Intelligent Robots*
693 *and Systems (IROS)*, pages 7534–7541. IEEE, 2020. 2

694 [6] Philipp Ausserlechner, David Habegger, Stefan Thalham-
695 mer, Jean-Baptiste Weibel, and Markus Vincze. Zs6d:
696 Zero-shot 6d object pose estimation using vision transfor-
697 mers. In *2024 IEEE International Conference on Robotics*
698 *and Automation (ICRA)*, pages 463–469. IEEE, 2024. 3

699 [7] Arpit Bahety, Priyanka Mandikal, Ben Abbatematteo, and
700 Roberto Martín-Martín. Screwmimic: Bimanual imitation
701 from human videos with screw space projection. *arXiv*
702 *preprint arXiv:2405.03666*, 2024. 1

703 [8] Shikhar Bahl, Abhinav Gupta, and Deepak Pathak.
704 Human-to-robot imitation in the wild. *arXiv preprint*
705 *arXiv:2207.09450*, 2022. 1, 2

[9] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, 706
and Deepak Pathak. Affordances from human videos as 707
a versatile representation for robotics. In *Proceedings of* 708
the IEEE/CVF Conference on Computer Vision and Pattern 709
Recognition, pages 13778–13790, 2023. 1, 2 710

[10] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, 711
Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampe- 712
dro, and Jeff Clune. Video pretraining (vpt): Learning to act 713
by watching unlabeled online videos. *Advances in Neural* 714
Information Processing Systems, 35:24639–24654, 2022. 2 715

[11] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, 716
Michal Yarom, Yonatan Bitton, Chenfanfu Jiang, Yizhou 717
Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evalu- 718
ating physical commonsense for video generation. *arXiv* 719
preprint arXiv:2406.03520, 2024. 1 720

[12] Leonardo Barcellona, Andrii Zadaianchuk, Davide Allegro, 721
Samuele Papa, Stefano Ghidoni, and Efstratios Gavves. 722
Dream to manipulate: Compositional world models em- 723
powering robot imitation learning with imagination. *arXiv* 724
preprint arXiv:2412.14957, 2024. 2 725

[13] Homanga Bharadhwaj, Abhinav Gupta, Shubham Tulsiani, 726
and Vikash Kumar. Zero-shot robot manipulation from pas- 727
sive human videos. *arXiv preprint arXiv:2302.02011*, 2023. 728
1 729

[14] Homanga Bharadhwaj, Debidatta Dwibedi, Abhinav 730
Gupta, Shubham Tulsiani, Carl Doersch, Ted Xiao, 731
Dhruv Shah, Fei Xia, Dorsa Sadigh, and Sean Kirmani. 732
Gen2act: Human video generation in novel scenarios en- 733
ables generalizable robot manipulation. *arXiv preprint* 734
arXiv:2409.16283, 2024. 1, 2, 3, 7, 17 735

[15] Homanga Bharadhwaj, Roozbeh Mottaghi, Abhinav Gupta, 736
and Shubham Tulsiani. Track2act: Predicting point tracks 737
from internet videos enables diverse zero-shot robot manip- 738
ulation, 2024. 1, 2, 7, 16 739

[16] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, 740
Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luh- 741
man, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya 742
Ramesh. Video generation models as world simulators. 743
2024. 1, 3 744

[17] Ming Cai and Ian Reid. Reconstruct locally, localize glob- 745
ally: A model free method for object pose estimation. In 746
Proceedings of the IEEE/CVF Conference on Computer Vi- 747
sion and Pattern Recognition, pages 3153–3163, 2020. 3 748

[18] Sylvain Calinon. A tutorial on task-parameterized move- 749
ment learning and retrieval. *Intelligent service robotics*, 9:
1–29, 2016. 3 750

[19] Andrea Caraffa, Davide Boscaini, Amir Hamza, and Fabio 751
Poiesi. Freeze: Training-free zero-shot 6d pose estimation 752
with geometric and vision foundation models. *European* 753
Conference on Computer Vision (ECCV), 2024. 3 754

[20] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, 755
Julien Mairal, Piotr Bojanowski, and Armand Joulin. 756
Emerging properties in self-supervised vision transformers. 757
In *Proceedings of the International Conference on Com-* 758
puter Vision (ICCV), 2021. 20 759

[21] Elliot Chane-Sane, Cordelia Schmid, and Ivan Laptev. 760
Learning video-conditioned policies for unseen manipu- 761
lation tasks. In *2023 IEEE International Conference on* 762
Computer Vision (ICCV), 2023. 763

- 764 *Robotics and Automation (ICRA)*, pages 909–916. IEEE, 2023. 1
- 765
- 766 [22] Matthew Chang, Arjun Gupta, and Saurabh Gupta. Semantic visual navigation by watching youtube videos. In *NeurIPS*, 2020. 1, 2
- 767
- 768
- 769 [23] Matthew Chang, Theophile Gervet, Mukul Khanna, Sriram Yenamandra, Dhruv Shah, So Yeon Min, Kavitha Shah, Chris Paxton, Saurabh Gupta, Dhruv Batra, Roozbeh Mottaghi, Jitendra Malik, and Devendra Singh Chaplot. Goat: Go to any thing. *arXiv preprint arXiv:2311.06430*, 2023. 4
- 770
- 771
- 772
- 773
- 774 [24] Xuxin Cheng, Yandong Ji, Junming Chen, Ruihan Yang, Ge Yang, and Xiaolong Wang. Expressive whole-body control for humanoid robots. *arXiv preprint arXiv:2402.16796*, 2024. 3
- 775
- 776
- 777
- 778 [25] Sungjoon Choi, Matthew KXJ Pan, and Joohyung Kim. Nonparametric motion retargeting for humanoid robots on shared latent space. In *Robotics: science and systems*, 2020. 3
- 779
- 780
- 781
- 782 [26] Sudeep Dasari and Abhinav Gupta. Transformers for one-shot visual imitation. In *Conference on Robot Learning*, pages 2071–2084. PMLR, 2021. 1, 2
- 783
- 784
- 785 [27] Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Abhinav Gupta. An unbiased look at datasets for visuomotor pre-training. In *Conference on Robot Learning*, pages 1183–1198. PMLR, 2023. 2
- 786
- 787
- 788
- 789 [28] Carl Doersch, Yi Yang, Dilara Gokay, Pauline Luc, Skanda Koppula, Ankush Gupta, Joseph Heyward, Ross Goroshin, João Carreira, and Andrew Zisserman. Bootstap: Bootstrapped training for tracking-any-point. *arXiv preprint arXiv:2402.00847*, 2024. 17
- 790
- 791
- 792
- 793
- 794 [29] Yilun Du, Mengjiao Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet, Tianhe Yu, Pieter Abbeel, Joshua B Tenenbaum, et al. Video language planning. *arXiv preprint arXiv:2310.10625*, 2023. 3
- 795
- 796
- 797
- 798 [30] Yilun Du, Sherry Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 3
- 799
- 800
- 801
- 802
- 803 [31] Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhui Liu, Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and temporal domains. *IEEE Transactions on Robotics*, 2023. 4
- 804
- 805
- 806
- 807
- 808 [32] Chelsea Finn, Tianhe Yu, T. Zhang, P. Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. In *CoRL*, 2017. 2
- 809
- 810
- 811 [33] Peter R Florence, Lucas Manuelli, and Russ Tedrake. Dense object nets: Learning dense visual object descriptors by and for robotic manipulation. *arXiv preprint arXiv:1806.08756*, 2018. 2
- 812
- 813
- 814
- 815 [34] Zipeng Fu, Qingqing Zhao, Qi Wu, Gordon Wetzstein, and Chelsea Finn. Humanplus: Humanoid shadowing and imitation from humans. *arXiv preprint arXiv:2406.10454*, 2024. 2
- 816
- 817
- 818
- 819 [35] Chongkai Gao, Haozhuo Zhang, Zhixuan Xu, Zhehao Cai, and Lin Shao. Flip: Flow-centric generative planning for general-purpose manipulation tasks. *arXiv preprint arXiv:2412.08261*, 2024. 2
- 820
- 821
- 822
- 823 [36] Shenyan Gao, Siyuan Zhou, Yilun Du, Jun Zhang, and Chuang Gan. Adaworld: Learning adaptable world models with latent actions. *arXiv preprint arXiv:2503.18938*, 2025. 1
- 824
- 825
- 826
- 827 [37] Theophile Gervet, Soumith Chintala, Dhruv Batra, Jitendra Malik, and Devendra Singh Chaplot. Navigating to objects in the real world. *Science Robotics*, 2023. 4
- 828
- 829
- 830 [38] Michael Gleicher. Retargeting motion to new characters. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, pages 33–42, 1998. 3
- 831
- 832
- 833 [39] Xuyang Guo, Jiayan Huo, Zhenmei Shi, Zhao Song, Jiahao Zhang, and Jiale Zhao. T2vphysbench: A first-principles benchmark for physical consistency in text-to-video generation. *arXiv preprint arXiv:2505.00337*, 2025. 1
- 834
- 835
- 836
- 837 [40] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 4
- 838
- 839
- 840 [41] Tairan He, Zhengyi Luo, Wenli Xiao, Chong Zhang, Kris Kitani, Changliu Liu, and Guanya Shi. Learning human-to-humanoid real-time whole-body teleoperation. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8944–8951. IEEE, 2024. 3
- 841
- 842
- 843
- 844
- 845 [42] Xingyi He, Jiaming Sun, Yuang Wang, Di Huang, Hujun Bao, and Xiaowei Zhou. Onepose++: Keypoint-free one-shot object pose estimation without cad models. *Advances in Neural Information Processing Systems*, 35: 35103–35115, 2022. 3
- 846
- 847
- 848
- 849
- 850 [43] Yisheng He, Wei Sun, Haibin Huang, Jianran Liu, Haoqiang Fan, and Jian Sun. Pvn3d: A deep point-wise 3d keypoints voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11632–11641, 2020. 3
- 851
- 852
- 853
- 854
- 855 [44] Yisheng He, Haibin Huang, Haoqiang Fan, Qifeng Chen, and Jian Sun. Ffb6d: A full flow bidirectional fusion network for 6d pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3003–3013, 2021. 3
- 856
- 857
- 858
- 859
- 860 [45] Yisheng He, Yao Wang, Haoqiang Fan, Jian Sun, and Qifeng Chen. Fs6d: Few-shot 6d pose estimation of novel objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6814–6824, 2022. 3
- 861
- 862
- 863
- 864
- 865 [46] Cheng-Chun Hsu, Bowen Wen, Jie Xu, Yashraj Narang, Xiaolong Wang, Yuke Zhu, Joydeep Biswas, and Stan Birchfield. Spot: Se (3) pose trajectory diffusion for object-centric manipulation. *arXiv preprint arXiv:2411.00965*, 2024. 2
- 866
- 867
- 868
- 869
- 870 [47] Kai Hu, Christian Ott, and Dongheui Lee. Online human walking imitation in task and joint space based on quadratic programming. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3458–3464. IEEE, 2014. 3
- 871
- 872
- 873
- 874
- 875 [48] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023. 6
- 876
- 877
- 878

- 879 [49] Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang,
880 and Li Fei-Fei. Rekep: Spatio-temporal reasoning of rela-
881 tional keypoint constraints for robotic manipulation. *arXiv*
882 *preprint arXiv:2409.01652*, 2024. 2, 6
- 883 [50] Ziqi Huang, Fan Zhang, Xiaojie Xu, Yinan He, Jiashuo Yu,
884 Ziyue Dong, Qianli Ma, Nattapol Chanpaisit, Chenyang
885 Si, Yuming Jiang, Yaohui Wang, Xinyuan Chen, Ying-
886 Cong Chen, Limin Wang, Dahua Lin, Yu Qiao, and Zi-
887 wei Liu. Vbench++: Comprehensive and versatile bench-
888 mark suite for video generative models. *arXiv preprint*
889 *arXiv:2411.13503*, 2024. 5, 20, 21
- 890 [51] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and
891 Tao Chen. Motiongpt: Human motion as a foreign lan-
892 guage. *Advances in Neural Information Processing Sys-*
893 *tems*, 36:20067–20079, 2023. 3
- 894 [52] Yuanchen Ju, Kaizhe Hu, Guowei Zhang, Gu Zhang, Min-
895 grun Jiang, and Huazhe Xu. Robo-abc: Affordance general-
896 ization beyond categories via semantic correspondence for
897 robot manipulation. In *European Conference on Computer*
898 *Vision*, pages 222–239. Springer, 2024. 2
- 899 [53] Nikita Karaev, Iurii Makarov, Jianyuan Wang, Natalia
900 Neverova, Andrea Vedaldi, and Christian Rupprecht. Co-
901 tracker3: Simpler and better point tracking by pseudo-
902 labelling real videos. *arXiv preprint arXiv:2410.11831*,
903 2024. 8
- 904 [54] Siddharth Karamcheti, Suraj Nair, Annie S Chen, Thomas
905 Kollar, Chelsea Finn, Dorsa Sadigh, and Percy Liang.
906 Language-driven representation learning for robotics. *arXiv*
907 *preprint arXiv:2302.12766*, 2023. 2
- 908 [55] Simar Kareer, Dhruv Patel, Ryan Punamiya, Pranay
909 Mathur, Shuo Cheng, Chen Wang, Judy Hoffman, and Dan-
910 fei Xu. Egomimic: Scaling imitation learning via egocen-
911 tric video. *arXiv preprint arXiv:2410.24221*, 2024. 1, 2
- 912 [56] Bingxin Ke, Dominik Narnhofer, Shengyu Huang, Lei Ke,
913 Torben Peters, Katerina Fragkiadaki, Anton Obukhov, and
914 Konrad Schindler. Video depth without video models,
915 2024. 4
- 916 [57] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler,
917 and George Drettakis. 3d gaussian splatting for real-time
918 radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1,
919 2023. 17
- 920 [58] Justin Kerr, Chung Min Kim, Mingxuan Wu, Brent Yi,
921 Qianqian Wang, Ken Goldberg, and Angjoo Kanazawa.
922 Robot see robot do: Imitating articulated object manipu-
923 lation with monocular 4d reconstruction. *arXiv preprint*
924 *arXiv:2409.18121*, 2024. 2, 3, 7, 17
- 925 [59] Chung Min Kim, Mingxuan Wu, Justin Kerr, Ken Gold-
926 berg, Matthew Tancik, and Angjoo Kanazawa. Garfield:
927 Group anything with radiance fields. In *Proceedings of*
928 *the IEEE/CVF Conference on Computer Vision and Pattern*
929 *Recognition*, pages 21530–21539, 2024. 17
- 930 [60] Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and
931 Joshua B Tenenbaum. Learning to act from actionless
932 videos through dense correspondences. *arXiv preprint*
933 *arXiv:2310.08576*, 2023. 2, 7, 17
- 934 [61] Scott Kuindersma, Robin Deits, Maurice Fallon, Andrés
935 Valenzuela, Hongkai Dai, Frank Permenter, Twan Koolen,
Pat Marion, and Russ Tedrake. Optimization-based lo-
936 comotion planning, estimation, and control design for the
937 atlas humanoid robot. *Autonomous robots*, 40:429–455,
938 2016. 3
- 939 [62] Yann Labbé, Justin Carpentier, Mathieu Aubry, and Josef
940 Sivic. Cosypose: Consistent multi-view multi-object 6d
941 pose estimation. In *Computer Vision—ECCV 2020: 16th*
942 *European Conference, Glasgow, UK, August 23–28, 2020,*
943 *Proceedings, Part XVII 16*, pages 574–591. Springer, 2020.
944 3
- 945 [63] Yann Labbé, Lucas Manuelli, Arsalan Mousavian, Stephen
946 Tyree, Stan Birchfield, Jonathan Tremblay, Justin Carpen-
947 tier, Mathieu Aubry, Dieter Fox, and Josef Sivic. Mega-
948 pose: 6d pose estimation of novel objects via render &
949 compare. In *Proceedings of the 6th Conference on Robot*
950 *Learning (CoRL)*, 2022. 3, 4, 19
- 951 [64] Arjun S Lakshminpathy, Jessica K Hodgins, and Nancy S
952 Pollard. Kinematic motion retargeting for contact-
953 rich anthropomorphic manipulations. *arXiv preprint*
954 *arXiv:2402.04820*, 2024. 3
- 955 [65] Marion Lepert, Jiaying Fang, and Jeannette Bohg. Phan-
956 tom: Training robots without robots using only human
957 videos. *arXiv preprint arXiv:2503.00779*, 2025. 1
- 958 [66] Fu Li, Shishir Reddy Vutukur, Hao Yu, Ivan Shugurov,
959 Benjamin Busam, Shaowu Yang, and Slobodan Ilic. Nerf-
960 pose: A first-reconstruct-then-regress approach for weakly-
961 supervised 6d object pose estimation. In *Proceedings of the*
962 *IEEE/CVF International Conference on Computer Vision*,
963 pages 2123–2133, 2023. 3
- 964 [67] Gen Li, Deqing Sun, Laura Sevilla-Lara, and Varun Jam-
965 pani. One-shot open affordance learning with foundation
966 models. In *Proceedings of the IEEE/CVF Conference on*
967 *Computer Vision and Pattern Recognition*, pages 3086–
968 3096, 2024. 2
- 969 [68] Gen Li, Nikolaos Tsagkas, Jifei Song, Ruaridh Mon-
970 Williams, Sethu Vijayakumar, Kun Shao, and Laura
971 Sevilla-Lara. Learning precise affordances from ego-
972 centric videos for robotic manipulation. *arXiv preprint*
973 *arXiv:2408.10123*, 2024. 2
- 974 [69] Jinhan Li, Yifeng Zhu, Yuqi Xie, Zhenyu Jiang, Mingyo
975 Seo, Georgios Pavlakos, and Yuke Zhu. Okami: Teaching
976 humanoid robots manipulation skills through single video
977 imitation. In *8th Annual Conference on Robot Learning*,
978 2024. 1, 3
- 979 [70] Zhen Li, Zuo-Liang Zhu, Ling-Hao Han, Qibin Hou, Chun-
980 Le Guo, and Ming-Ming Cheng. Amt: All-pairs multi-
981 field transforms for efficient frame interpolation. In *IEEE*
982 *Conference on Computer Vision and Pattern Recognition*
983 *(CVPR)*, 2023. 20
- 984 [71] Junbang Liang, Ruoshi Liu, Ege Ozguroglu, Sruthi Sud-
985 hakar, Achal Dave, Pavel Tokmakov, Shuran Song, and Carl
986 Vondrick. Dreamitate: Real-world visuomotor policy learn-
987 ing via video generation. *arXiv preprint arXiv:2406.16862*,
988 2024. 3
- 989 [72] Yuwei Liang, Weijie Li, Yue Wang, Rong Xiong, Yichao
990 Mao, and Jiafan Zhang. Dynamic movement primitive
991 based motion retargeting for dual-arm sign language mo-
992

- 993 tions. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8195–8201. IEEE, 2021. 3
- 994
- 995 [73] Fangfu Liu, Wenqiang Sun, Hanyang Wang, Yikai Wang, Haowen Sun, Junliang Ye, Jun Zhang, and Yueqi Duan. Reconx: Reconstruct any scene from sparse views with video diffusion model. *arXiv preprint arXiv:2408.16767*, 2024. 1
- 996
- 997
- 998
- 999 [74] Minghua Liu, Chao Xu, Haiyan Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in Neural Information Processing Systems*, 36:22226–22246, 2023. 10
- 1000
- 1001
- 1002
- 1003
- 1004 [75] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9298–9309, 2023. 10
- 1005
- 1006
- 1007
- 1008
- 1009 [76] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 4
- 1010
- 1011
- 1012
- 1013
- 1014 [77] Xingyu Liu, Gu Wang, Ruida Zhang, Chenyangguang Zhang, Federico Tombari, and Xiangyang Ji. Unopose: Unseen object pose estimation with an unposed rgb-d reference image. *arXiv preprint arXiv:2411.16106*, 2024. 3
- 1015
- 1016
- 1017
- 1018 [78] YuXuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation from observation: Learning to imitate behaviors from raw video via context translation. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1118–1125. IEEE, 2018. 2
- 1019
- 1020
- 1021
- 1022
- 1023 [79] Yuan Liu, Yilin Wen, Sida Peng, Cheng Lin, Xiaoxiao Long, Taku Komura, and Wenping Wang. Gen6d: Generalizable model-free 6-dof object pose estimation from rgb images. In *European Conference on Computer Vision*, pages 298–315. Springer, 2022. 3
- 1024
- 1025
- 1026
- 1027
- 1028 [80] Zhengyi Luo, Jinkun Cao, Kris Kitani, Weipeng Xu, et al. Perpetual humanoid control for real-time simulated avatars. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10895–10904, 2023. 3
- 1029
- 1030
- 1031
- 1032 [81] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emons, Anchit Gupta, Emre Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In *Conference on Robot Learning*, pages 879–893. PMLR, 2018. 2
- 1033
- 1034
- 1035
- 1036
- 1037
- 1038 [82] Russell Mendonca, Shikhar Bahl, and Deepak Pathak. Structured world models from human videos. *arXiv preprint arXiv:2308.10901*, 2023. 2
- 1039
- 1040
- 1041 [83] Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do generative video models learn physical principles from watching videos? *arXiv preprint arXiv:2501.09038*, 2025. 1
- 1042
- 1043
- 1044
- 1045 [84] Shinichiro Nakaoka, Atsushi Nakazawa, Fumio Kanehiro, Kenji Kaneko, Mitsuharu Morisawa, and Katsushi Ikeuchi. Task model of lower body motion for a biped humanoid robot to imitate human dances. In *2005 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3157–3162. IEEE, 2005. 3
- 1046
- 1047
- 1048
- 1049
- 1050
- [85] Van Nguyen Nguyen, Thibault Groueix, Mathieu Salzmann, and Vincent Lepetit. Gigapose: Fast and robust novel object pose estimation via one correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9903–9913, 2024. 3
- [86] Scott Niekum, Sarah Osentoski, George Konidaris, and Andrew G Barto. Learning and generalization of complex tasks from unstructured demonstrations. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5239–5246. IEEE, 2012. 3
- [87] Abby O’Neill, Abdul Rehman, Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023. 1
- [88] Evin Pinar Örnek, Yann Labbé, Bugra Tekin, Lingni Ma, Cem Keskin, Christian Forster, and Tomáš Hodaň. Foundpose: Unseen object pose estimation with foundation features. *European Conference on Computer Vision (ECCV)*, 2024. 3
- [89] Kiru Park, Timothy Patten, and Markus Vincze. Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7668–7677, 2019. 3
- [90] Keunhong Park, Arsalan Mousavian, Yu Xiang, and Dieter Fox. Latentfusion: End-to-end differentiable reconstruction and rendering for unseen object pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10710–10719, 2020. 3
- [91] Austin Patel, Andrew Wang, Ilija Radosavovic, and Jitendra Malik. Learning to imitate object interactions from internet videos. *arXiv preprint arXiv:2211.13225*, 2022. 2
- [92] Shivansh Patel, Xinchun Yin, Wenlong Huang, Shubham Garg, Hooshang Nayyeri, Li Fei-Fei, Svetlana Lazebnik, and Yunzhu Li. A real-to-sim-to-real approach to robotic manipulation with vlm-generated iterative keypoint rewards. *arXiv preprint arXiv:2502.08643*, 2025. 6
- [93] Deepak Pathak, Parsa Mahmoudieh, Guanghao Luo, Pulkit Agrawal, Dian Chen, Yide Shentu, Evan Shelhamer, Jitendra Malik, Alexei A Efros, and Trevor Darrell. Zero-shot visual imitation. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 2050–2053, 2018. 2
- [94] Luigi Penco, Nicola Scianca, Valerio Modugno, Leonardo Lanari, Giuseppe Oriolo, and Serena Ivaldi. A multimode teleoperation framework for humanoid loco-manipulation: An application for the icub robot. *IEEE Robotics & Automation Magazine*, 26(4):73–82, 2019. 3
- [95] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Transactions on Graphics (ToG)*, 40(4):1–20, 2021. 3
- [96] Georgy Ponimatkin, Martin Cifka, Tomáš Souček, Mederic Fourmy, Yann Labbé, Vladimir Petrik, and Josef Sivic. 6d object pose tracking in internet videos for robotic manipulation. *arXiv preprint arXiv:2503.10307*, 2025. 2
- [97] Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong Wang. Dexmv: Im-

- 1109 imitation learning for dexterous manipulation from human
1110 videos. In *European Conference on Computer Vision*, pages
1111 570–587. Springer, 2022. 3
- 1112 [98] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
1113 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
1114 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen
1115 Krueger, and Ilya Sutskever. Learning transferable visual
1116 models from natural language supervision, 2021. 20
- 1117 [99] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang
1118 Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman
1119 Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junt-
1120 yan Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-
1121 Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Fe-
1122 ichtenhofer. Sam 2: Segment anything in images and
1123 videos, 2024. 4
- 1124 [100] Juntao Ren, Priya Sundareshan, Dorsa Sadigh, Sanjiban
1125 Choudhury, and Jeannette Bohg. Motion tracks: A unified
1126 representation for human-robot transfer in few-shot imita-
1127 tion learning. *arXiv preprint arXiv:2501.06994*, 2025. 2
- 1128 [101] Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine
1129 Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google
1130 Brain. Time-contrastive networks: Self-supervised learn-
1131 ing from video. In *2018 IEEE international conference on*
1132 *robotics and automation (ICRA)*, pages 1134–1141. IEEE,
1133 2018.
- 1134 [102] Pratyusha Sharma, Lekha Mohan, Lerrel Pinto, and Abhinav
1135 Gupta. Multiple interactions made easy (mime): Large
1136 scale demonstrations data for imitation. *arXiv:1810.07121*,
1137 2018. 2
- 1138 [103] Pratyusha Sharma, Deepak Pathak, and Abhinav Gupta.
1139 Third-person visual imitation learning via decoupled hier-
1140 archical controller. *Advances in Neural Information Pro-*
1141 *cessing Systems*, 32, 2019. 2
- 1142 [104] Junyao Shi, Zhuolun Zhao, Tianyou Wang, Ian Pedroza,
1143 Amy Luo, Jie Wang, Jason Ma, and Dinesh Jayaraman. Ze-
1144 romimic: Distilling robotic manipulation skills from web
1145 videos. *arXiv preprint arXiv:2503.23877*, 2025. 2
- 1146 [105] Ivan Shugurov, Fu Li, Benjamin Busam, and Slobodan
1147 Ilic. Osop: A multi-stage one shot object pose estimation
1148 framework. In *Proceedings of the IEEE/CVF Conference*
1149 *on Computer Vision and Pattern Recognition*, pages 6835–
1150 6844, 2022. 3
- 1151 [106] Aravind Sivakumar, Kenneth Shaw, and Deepak Pathak.
1152 Robotic telekinesis: Learning a robotic hand imita-
1153 tor by watching humans on youtube. *arXiv preprint*
1154 *arXiv:2202.10448*, 2022. 1, 2
- 1155 [107] Laura Smith, Nikita Dhawan, Marvin Zhang, Pieter Abbeel,
1156 and Sergey Levine. Avid: Learning multi-stage tasks via
1157 pixel-level translation of human videos. *arXiv preprint*
1158 *arXiv:1912.04443*, 2019. 2
- 1159 [108] Mohan Kumar Srirama, Sudeep Dasari, Shikhar Bahl, and
1160 Abhinav Gupta. Hrp: Human affordances for robotic pre-
1161 training. *arXiv preprint arXiv:2407.18911*, 2024. 2
- 1162 [109] Jiaming Sun, Zihao Wang, Siyu Zhang, Xingyi He,
1163 Hongcheng Zhao, Guofeng Zhang, and Xiaowei Zhou.
1164 Onepose: One-shot object pose estimation without cad
1165 models. In *Proceedings of the IEEE/CVF Conference on*
Computer Vision and Pattern Recognition, pages 6825–
6834, 2022. 3
- [110] Yihong Sun, Hao Zhou, Liangzhe Yuan, Jennifer J Sun,
Yandong Li, Xuhui Jia, Hartwig Adam, Bharath Hariharan,
Long Zhao, and Ting Liu. Video creation by demonstration.
arXiv preprint arXiv:2412.09551, 2024. 3
- [111] Jonathan Tremblay, Thang To, Balakumar Sundaralingam,
Yu Xiang, Dieter Fox, and Stan Birchfield. Deep object
pose estimation for semantic robotic grasping of household
objects. *arXiv preprint arXiv:1809.10790*, 2018. 3
- [112] Eugene Valassakis, Georgios Papagiannis, Norman Di Palo,
and Edward Johns. Demonstrate once, imitate immediately
(dome): Learning visual servoing for one-shot imitation
learning. In *2022 IEEE/RSJ International Conference on*
Intelligent Robots and Systems (IROS), pages 8614–8621.
IEEE, 2022. 2
- [113] Mel Vecerik, Carl Doersch, Yi Yang, Todor Davchev, Yusuf
Aytar, Guangyao Zhou, Raia Hadsell, Lourdes Agapito, and
Jon Scholz. Robotap: Tracking arbitrary points for few-shot
visual imitation. In *2024 IEEE International Conference on*
Robotics and Automation (ICRA), pages 5397–5403. IEEE,
2024. 2
- [114] Chen Wang, Linxi Fan, Jiankai Sun, Ruohan Zhang, Li Fei-
Fei, Danfei Xu, Yuke Zhu, and Anima Anandkumar. Mim-
icplay: Long-horizon imitation learning by watching hu-
man play. *arXiv preprint arXiv:2302.12422*, 2023. 1, 2
- [115] Yi Wang, Yanan He, Yizhuo Li, Kunchang Li, Jiashuo Yu,
Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui
Wang, et al. Internvid: A large-scale video-text dataset
for multimodal understanding and generation. In *The*
Twelfth International Conference on Learning Representa-
tions, 2023. 5, 20
- [116] Bowen Wen, Jonathan Tremblay, Valts Blukis, Stephen
Tyree, Thomas Muller, Alex Evans, Dieter Fox, Jan Kautz,
and Stan Birchfield. Bundlesdf: Neural 6-dof tracking and
3d reconstruction of unknown objects. *CVPR*, 2023. 4, 16
- [117] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield.
Foundationpose: Unified 6d pose estimation and tracking
of novel objects. *arXiv preprint arXiv:2312.08344*, 2023.
1, 3, 4, 19
- [118] Albert Wu, Ruocheng Wang, Sirui Chen, Clemens Eppner,
and C Karen Liu. One-shot transfer of long-horizon ex-
trinsic manipulation through contact retargeting. In *2024*
IEEE/RSJ International Conference on Intelligent Robots
and Systems (IROS), pages 13891–13898. IEEE, 2024. 3
- [119] Annie Xie, Lisa Lee, Ted Xiao, and Chelsea Finn. Decom-
posing the generalization gap in imitation learning for vi-
sual robotic manipulation. In *2024 IEEE International Con-*
ference on Robotics and Automation (ICRA), pages 3153–
3160. IEEE, 2024. 1
- [120] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Reza Tofighi,
and Dacheng Tao. Gmflow: Learning optical flow via global
matching. In *Proceedings of the IEEE/CVF conference on*
computer vision and pattern recognition, pages 8121–8130,
2022. 17
- [121] Mengda Xu, Zhenjia Xu, Cheng Chi, Manuela Veloso, and
Shuran Song. Xskill: Cross embodiment skill discovery. In

- 1223 *Conference on robot learning*, pages 3536–3555. PMLR,
1224 2023. 2
- 1225 [122] Mengda Xu, Zhenjia Xu, Yinghao Xu, Cheng Chi, Gor-
1226 don Wetzstein, Manuela Veloso, and Shuran Song. Flow
1227 as the cross-domain manipulation interface. *arXiv preprint*
1228 *arXiv:2407.15208*, 2024. 2
- 1229 [123] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan
1230 Tompson, Dale Schuurmans, and Pieter Abbeel. Learn-
1231 ing interactive real-world simulators. *arXiv preprint*
1232 *arXiv:2310.06114*, 1(2):6, 2023. 3
- 1233 [124] Xindi Yang, Baolu Li, Yiming Zhang, Zhenfei Yin, Lei Bai,
1234 Liqian Ma, Zhiyong Wang, Jianfei Cai, Tien-Tsin Wong,
1235 Huchuan Lu, et al. Vlpp: Towards physically plausible
1236 video generation with vision and language informed physi-
1237 cal prior. *arXiv e-prints*, pages arXiv–2503, 2025. 1
- 1238 [125] Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Sejune Joo,
1239 Jianwei Yang, Baolin Peng, Ajay Mandlekar, Reuben Tan,
1240 Yu-Wei Chao, Bill Yuchen Lin, et al. Latent action pretrain-
1241 ing from videos. *arXiv preprint arXiv:2410.11758*, 2024. 1
- 1242 [126] Tianhe Yu, Chelsea Finn, Annie Xie, Sudeep Dasari, Tian-
1243 hao Zhang, Pieter Abbeel, and Sergey Levine. One-shot im-
1244 itation from observing humans via domain-adaptive meta-
1245 learning. *arXiv preprint arXiv:1802.01557*, 2018. 2
- 1246 [127] Chengbo Yuan, Chuan Wen, Tong Zhang, and Yang Gao.
1247 General flow as foundation affordance for scalable robot
1248 learning. *arXiv preprint arXiv:2401.11439*, 2024. 2
- 1249 [128] Kevin Zakka, Andy Zeng, Pete Florence, Jonathan Tomp-
1250 son, Jeannette Bohg, and Debidatta Dwibedi. Xirl: Cross-
1251 embodiment inverse reinforcement learning. In *Conference*
1252 *on Robot Learning*, pages 537–546. PMLR, 2022. 2
- 1253 [129] Qihang Zhang, Shuangfei Zhai, Miguel Angel Bautista,
1254 Kevin Miao, Alexander Toshev, Joshua Susskind, and Ji-
1255 atao Gu. World-consistent video diffusion with explicit 3d
1256 modeling. *arXiv preprint arXiv:2412.01821*, 2024. 1
- 1257 [130] Zhengyou Zhang. A flexible new technique for camera cal-
1258 ibration. *IEEE Transactions on pattern analysis and ma-
1259 chine intelligence*, 22(11):1330–1334, 2000. 16, 17
- 1260 [131] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea
1261 Finn. Learning fine-grained bimanual manipulation with
1262 low-cost hardware. *arXiv preprint arXiv:2304.13705*,
1263 2023. 9
- 1264 [132] Haoyu Zhen, Qiao Sun, Hongxin Zhang, Junyan Li, Siyuan
1265 Zhou, Yilun Du, and Chuang Gan. Tesseract: Learning 4d
1266 embodied world models. 2025. 3
- 1267 [133] Huayi Zhou, Ruixiang Wang, Yunxin Tai, Yueci Deng,
1268 Guiliang Liu, and Kui Jia. You only teach once: Learn one-
1269 shot bimanual robotic manipulation from video demonstra-
1270 tions. *arXiv preprint arXiv:2501.14208*, 2025. 2
- 1271 [134] Jiaming Zhou, Teli Ma, Kun-Yu Lin, Zifan Wang, Ronghe
1272 Qiu, and Junwei Liang. Mitigating the human-robot domain
1273 discrepancy in visual pre-training for robotic manipulation.
1274 *arXiv preprint arXiv:2406.14235*, 2024. 1
- 1275 [135] Junzhe Zhu, Yuanchen Ju, Junyi Zhang, Muhan Wang,
1276 Zhecheng Yuan, Kaizhe Hu, and Huazhe Xu. Dense-
1277 matcher: Learning 3d semantic correspondence for
1278 category-level manipulation from a single demo. *arXiv*
1279 *preprint arXiv:2412.05268*, 2024. 2

Appendix

1280 A. Best Practices for Video Generation

1281 We found that the following practices lead to reliable video
1282 generation: (1) having a clean background without visual
1283 distractions, (2) minimizing the number of distractor objects
1284 in the scene, (3) ensuring objects are reasonably large and
1285 viewed from a natural, human-like perspective, (4) ensuring
1286 there is one clearly identifiable task that can be performed,
1287 (5) using simple and concise text prompts, and (6) setting
1288 the relevance factor to 0.7 with the negative prompt “fast
1289 motion” led to the most reliable video generations.

1290 B. Prompting for Video Filtering and Filtering 1291 Statistics

1292 The prompt for GPT o1-based filtering is shown in Figure
1293 10. We provide GPT o1 with the prompt, the video sum-
1294 mary—constructed by vertically concatenating evenly sam-
1295 pled frames from the video—and the language command
1296 (e.g., “pour water”). GPT o1 then responds with “Yes” or
1297 “No” to indicate whether the specified task is successfully
1298 performed in the video. The filtering success rates are: 83%
1299 for pouring, 66% for lifting, 55% for placing, and 45% for
1300 sweeping.

1301 C. Mesh-Free Object Tracking

1302 We experiment with a mesh-free object tracking version of
1303 our method. Specifically, we use BundleSDF [116], which
1304 jointly performs 6-DoF object tracking and reconstruction
1305 from RGBD observations. For the *pouring* task, we evalu-
1306 ate our method using trajectories obtained via BundleSDF
1307 over 10 trials and observe a success rate of (90%), match-
1308 ing our default tracking setup. While the BundleSDF paper
1309 reports real-time capabilities, we found that its official im-
1310 plementation takes approximately 30 minutes to process each
1311 video in practice, which limits its applicability for real-time
1312 deployment. In contrast, our default tracker operates in real-
1313 time, enabling closed-loop execution and recovery from dis-
1314 turbances as discussed in Sec. 4.6. While the BundleSDF
1315 paper reports real-time capabilities, we observed signifi-
1316 cantly higher runtimes in practice with the official imple-
1317 mentation. We expect that future advances in model-free
1318 tracking will address these efficiency bottlenecks, allowing
1319 for real-time mesh-free deployment.

1320 D. Smoothing Object Trajectories

1321 To reduce noise and jitter in the estimated object poses, we
1322 apply a moving average filter with a fixed sliding window
1323 (centered on each point), separately to the position and ori-

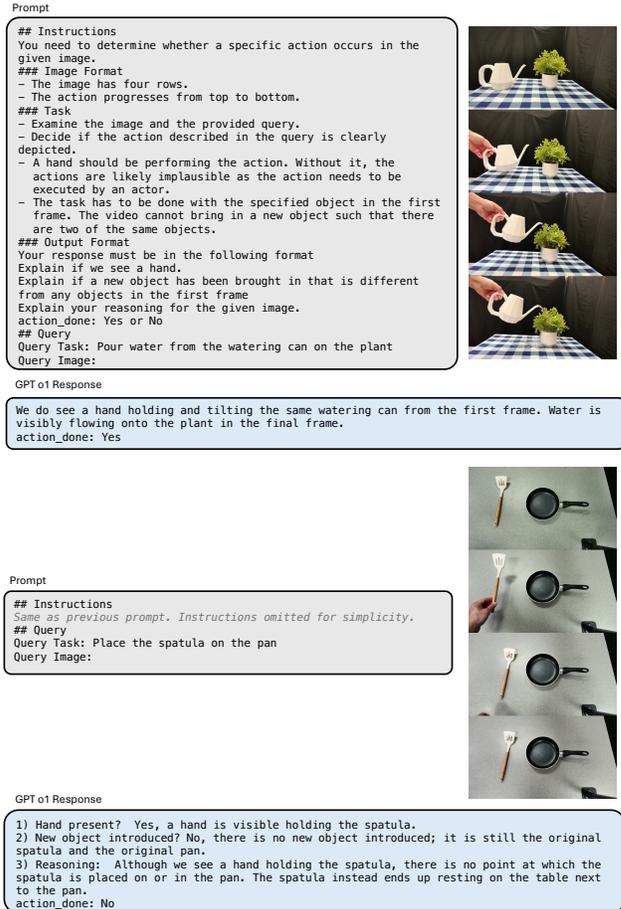


Figure 10. Examples of prompting GPT o1 to filter generated videos. We sample generated video frames and prompt GPT o1 to assess whether the specified task is performed successfully in the video.

entation components. Translations are smoothed independ- 1324
 1325
 1326
 1327
 1328
 1329

E. Description of Baselines 1330

Track2Act [15]: We adapt Track2Act’s procedure to our 1331
 1332
 1333
 1334
 1335
 1336
 1337

To integrate this into our pipeline, we use their published checkpoint but modify the input formulation—while the initial image remains identical to our real camera’s view, the goal image is taken from the last frame of a generated video rather than being physically captured. We then use PnP on the predicted point tracks along with the initial depth image to estimate the object’s rigid motion across frames, thereby defining the end-effector trajectory. We use interpolation between consecutive poses because Track2Act generates only a sparse set of frames, and denser sampling is needed for smooth trajectory estimation and execution. However, we do not include Track2Act’s closed-loop residual policy correction, focusing solely on open-loop 6D object-pose estimation and execution. This adaptation allows us to directly evaluate how well a vision-based, open-loop approach generalizes to our setting without additional corrections.

Gen2Act [14]: Gen2Act introduces a video-conditioned policy learning framework that first generates a human video using a video generation model from a scene image and a task description. The system then extracts object tracks using BootsTAP [28], and trains a policy using behavior cloning with an auxiliary track prediction loss and offline robot demonstrations. At inference, Gen2Act only uses the generated video and the learned policy to predict robot actions.

Our approach presents a simplified adaptation of this framework that removes the need for behavior cloning, and offline demonstrations. Instead of using the extracted tracks as an auxiliary loss, we directly process them for pose estimation. To recover 3D object positions, we leverage an initial depth image corresponding to the scene image, allowing us to obtain depth values for the extracted 2D tracks. We apply RANSAC filtering to remove outlier track points and then use the Perspective-n-Point (PnP) [130] to estimate the object’s 6DoF pose. This adaptation preserves the core idea of leveraging video and track-based signals while eliminating the need for supervised policy learning.

AVDC [60]: The AVDC approach models action trajectories by synthesizing a task-driven video (using a trained text-conditioned video generation model) and using optical flow from GMFlow [120] to estimate dense pixel correspondences. It then reconstructs 3D object motion using an optimization step that refines pose estimates based on the tracked flow and depth information. To improve robustness, AVDC also includes a replanning mechanism that re-executes the pipeline when predicted motion stagnates.

Since the trained text-conditioned video generation model did not generalize well to our setup, we instead use the same generated video as in other experiments to ensure a fair comparison. While we do not employ AVDC’s replanning strategy, we predict object poses using a similar optimization framework based on flow and depth informa-

tion.

4D-DPM [58]: 4D-DPM is designed to track the 3D motion of articulated object parts from a single video. It first constructs a 3D Gaussian splatting [57] representation of the scene to capture object features, then applies GARField [59] to cluster the Gaussians into discrete object components. In our adaptation, we modify this approach to operate on entire objects rather than individual parts. Specifically, we set the clustering parameters to treat the object as a single entity, ensuring that motion estimation is performed at the object level rather than segmenting it into multiple parts. This allows us to track and execute trajectories for the whole object.

F. ReKep Predictions and Executions

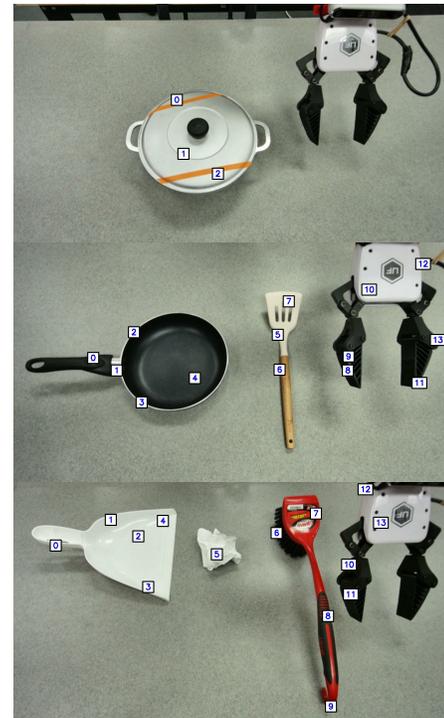


Figure 11. **Examples of ReKep’s Keypoint Locations.** The keypoint placements are often suboptimal, except for sweeping task, where the keypoints are reasonable.

A detailed example of ReKep’s keypoint and VLM predictions for pouring task is shown in Fig. 12. The VLM first predicts to grasp the watering can at keypoint 1. For the transport phase, it instructs moving keypoint 8 above keypoint 15, while keeping its height above keypoint 7. For the pouring action, keypoint 8 remains above 15 (to place the spout over the plant) and above keypoint 4 (to induce tilting). The resulting robot execution fails. We attribute most ReKep failures to inaccurate keypoint predictions, as shown

1414 in Fig. 11. In the lid image, there is no keypoint at the handle
 1415 of the lid. In the placing task, keypoints cluster around
 1416 pan corners. For the sweeping task, the keypoints are generally
 1417 well-placed, and executions succeeded. Because the initial keypoints
 1418 are suboptimal, downstream VLM predictions are also inaccurate.
 1419



Figure 12. ReKep’s output for the pouring task and the resulting robot execution (top-right). The VLM predictions on the generated keypoints lead to failed execution.

1420 G. Limitation of Tracking with Point Tracks

1421 All point tracks fail under extreme rotations, as initially visible
 1422 points often become occluded. This is a fundamental limitation of any
 1423 correspondence-based tracking method that relies solely on visible surface
 1424 features. We show this failure in Fig. 13. As the object rotates, most initial
 1425 points are lost, resulting in insufficient 2D-3D correspondences to
 1426 solve a stable PnP problem. This degrades pose estimation quality, leading
 1427 to large drift or abrupt jumps in estimated object motion. Such instability
 1428 cascades into robot
 1429

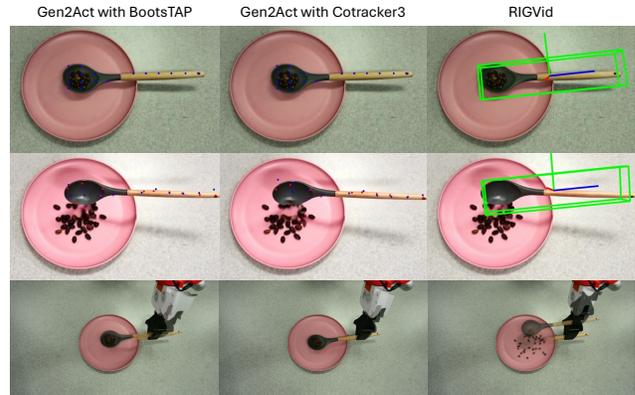


Figure 13. Gen2Act with BootsTAP, CoTracker, and RIGVid. Blue points denote the tracked points used for PnP; red points represent the reprojected 3D points. For a good PnP solution, these should align, as seen in the first frame. For Gen2Act, the blue points drift significantly from the red ones in later frames, indicating failure in pose estimation due to tracking loss, which leads to failed robot execution.

1430 execution errors, often causing the robot to fail at the task
 1431 altogether. As a result, both variants of Gen2Act—despite
 1432 stronger tracking backbones like CoTracker—still fail under
 1433 large out-of-plane rotations. In contrast, RIGVid’s model-based 6D
 1434 tracking handles these situations more robustly, as it uses full-object
 1435 geometry and SE(3) filtering to maintain stable trajectories.
 1436

H. Additional Robustness Examples

1437

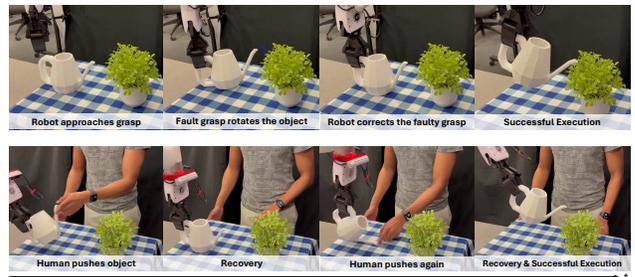


Figure 14. Additional examples of RIGVid’s robustness. In the top row, RIGVid recovers from a faulty initial grasp by reorienting the object before continuing execution. In the bottom row, it corrects for external disturbances on the object when a human pushes the object mid-execution, realigning and successfully completing the task.

1438 Examples of RIGVid’s robustness are shown in Fig. 14.
 1439 In the first row, the robot initially grasps the object, but due
 1440 to a misaligned grasp, the object rotates unexpectedly. The
 1441 robot compensates by rotating the object back to the cor-

1442 rect orientation and then resumes the planned trajectory, ul-
 1443 timately completing the task successfully. In the bottom
 1444 row, a human perturbs the object during execution while it
 1445 is held by the robot. RIGVid detects the resulting change in
 1446 the relative transformation and automatically re-aligns the
 1447 object before continuing. When the human intervenes a sec-
 1448 ond time, RIGVid again corrects the deviation, ultimately
 1449 leading to successful task completion.

1450 I. Errors from Depth Estimation

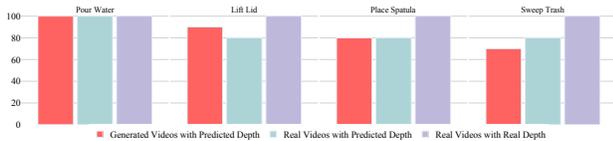


Figure 15. **Impact of Depth Estimation Errors on RIGVid performance.** Errors in monocular depth estimation result in worse performance of generated and real videos. RIGVid achieves perfect success across all tasks with real videos and real depth.

1451 In Fig. 15, we isolate the impact of depth estimation
 1452 errors. Robot executions on real videos with real depth
 1453 (captured using an RGBD camera) achieve a 100% suc-
 1454 cess rate, whereas executions from real videos with gen-
 1455 erated depth result in 85% average success. Similarly, ex-
 1456 ecutions from Kling V1.6-generated videos with gener-
 1457 ated depth also achieve 85% success, suggesting that the primary
 1458 source of error lies in monocular depth estimation. Upon
 1459 inspection, we observe two common undesirable behaviors
 1460 in the predicted depth: inaccurate depth values and tempo-
 1461 ral flickering. An example of inaccurate depth is shown in
 1462 Fig. 16. In the generated video, when the spatula is brought
 1463 close to the camera, the depth changes by only 6.8 cm,
 1464 which is visibly inconsistent with the video and likely much
 1465 smaller than the real-world change. Inaccuracies also oc-
 1466 cur in real videos, as shown in the figure—the head of the
 1467 spatula is estimated to be far from the camera, despite ap-
 1468 pearing close, revealing another failure mode in monocular
 1469 depth estimation. An illustration of flickering is shown in
 1470 Fig. 17. Although the position of the watering can relative
 1471 to the camera remains nearly unchanged across three con-
 1472 secutive frames, the estimated depth varies significantly. In
 1473 particular, the zoomed-in region on the right shows the can
 1474 appearing much whiter than on the left, indicating a sub-
 1475 stantial change in predicted depth. The average depth of the
 1476 can changes from 40.1 cm to 38.2 cm—a 1.9 cm difference
 1477 over just 0.066 seconds—which is physically implausible for
 1478 the generated video. We find similar flickering behavior in
 1479 real videos as well, where the depth changes from 43.2 cm
 1480 to 40.9 cm in the given example—a 2.3 cm difference.

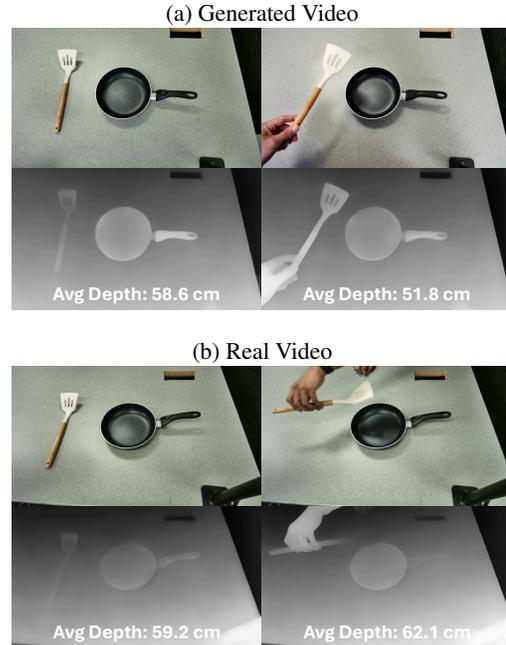


Figure 16. **Errors in Monocular Depth Estimation.** In the generated video (top), the depth of the spatula changes only slightly despite a large visual change. In the real video (bottom), the spatula’s head is predicted to lie farther away, contradicting the visual appearance.

J. Choice between MegaPose and Foundation- Pose

We compare the stability of trajectories obtained from MegaPose [63] and FoundationPose [117] by computing the translational and rotational RMS jitter. For each method, we apply a Gaussian smoothing filter ($\sigma = 2$ frames) to the raw SE(3) pose sequences, compute the residual between the original and smoothed trajectories, and then calculate:

$$\text{jitter}_{\text{trans}} = \sqrt{\frac{1}{N} \sum_{t=1}^N \|\Delta \mathbf{t}_t\|^2}, \quad \text{jitter}_{\text{rot}} = \sqrt{\frac{1}{N} \sum_{t=1}^N \theta_t^2},$$

where $\Delta \mathbf{t}_t$ is the translational residual at frame t , and θ_t is the angular magnitude (in radians) of the relative rotation $R_{\text{smooth}}^{-1} R_{\text{raw}}$, converted to degrees. We average these metrics over ten pouring trajectories extracted from generated videos.

MegaPose yields an average translational RMS jitter of 0.0045m and rotational RMS jitter of 37.47°, whereas FoundationPose achieves 0.0029m translational and 14.31° rotational jitter. These results demonstrate that FoundationPose produces significantly smoother and more stable trajectories. Additionally, it allows for real-time tracking during the execution, allowing us to make RIGVid robust to external disturbances.

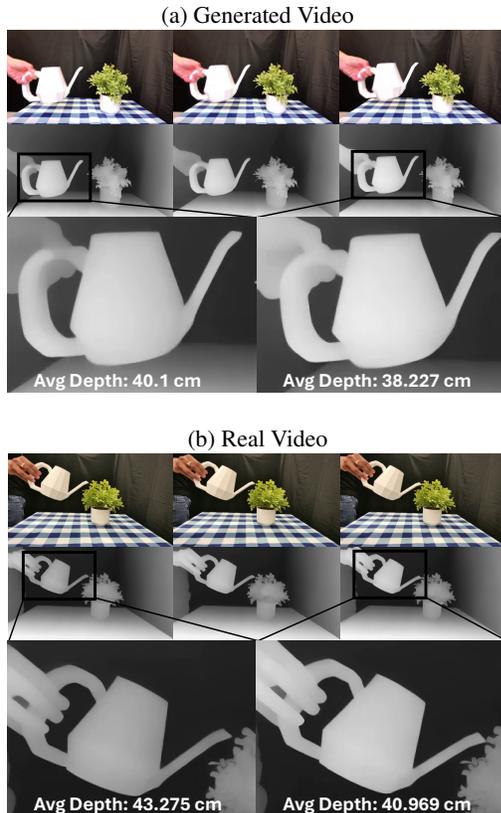


Figure 17. **Flickering in Depth Prediction.** We show three consecutive frames of the video and its corresponding predicted depth. The depth of the watering can change noticeably across frames—appearing significantly whiter in the third frame despite minimal actual motion. We observe this behavior in both generated and real videos.

1503

K. Comparing Video Generative Models

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

1518

1519

To further assess video quality, we report VBench++ [50] metrics in Table 2 and explain them below. The numbers in the table are scaled $100\times$ for easier interpretation. We collect these metrics on 40 randomly selected and unfiltered videos per model, 10 for each of the four tasks. Kling v1.6 outperformed the other models on most metrics but performed similarly or worse in video-text consistency and dynamic degree. Human evaluations discussed in Sec. 4.2 suggest that the video-text consistency and I2V subject consistency are not reliable indicators of whether a generated video correctly follows a given command. Sora scored high on dynamic degree, likely due to its tendency to drastically alter the scene, resulting in exceptionally large motions. Generated videos from these models and their corresponding metrics are shown in Fig. 18 and further details on these metrics can be found the next section.

1520

VBench++ Metric Definitions:

- **Subject Consistency.** Subject consistency describes whether subjects’ appearance remain consistent, which is computed by DINOv1 [20] similarities across video frames. 1521
- **Background Consistency.** Background temporal consistency by CLIP [98] similarities across frames. 1522
- **Motion Smoothness.** Evaluates smoothness of videos by utilizing video frame interpolation model AMT [70]. 1523
- **Dynamic Degree.** Describes whether the video contains large motions as a binary metric. 1524
- **Aesthetic Quality.** Human perceived artistic and beauty value such as photo-realism, layout and color harmony. 1525
- **Imaging Quality.** Assesses the presence of distortion in a video, such as noisiness, blurriness, and over-exposure. 1526
- **Video-Text Consistency.** Text-to-video alignment score calculated by ViCLIP [115]. 1527
- **I2V Subject Consistency.** Similarity between subjects in input image and each video frame, as well as similarity between consecutive frames. Features are extracted from DINOv1 [20]. 1528

Metrics	Video Generation Models			Human Demos
	Kling V1.6	Kling V1.5	Sora	
Subject Consistency	96.34	91.66	83.09	94.91
Background Consistency	96.64	93.97	89.34	95.00
Motion Smoothness	99.68	99.57	99.06	99.51
Dynamic Degree	52.5	57.5	70.0	80.0
Aesthetic Quality	51.75	49.77	46.22	49.30
Imaging Quality	72.80	71.48	68.68	72.52
Video-Text Consistency	22.01	22.61	21.42	21.57
I2V Subject Consistency	97.88	95.96	89.09	97.89

Table 2. **Video generation quality metrics for real human demonstration videos and different models.** Higher values indicate better quality. Kling v1.6 performs comparably to or surpasses other models on most metrics.

SORA



Kling AI v1.5



Kling AI v1.6

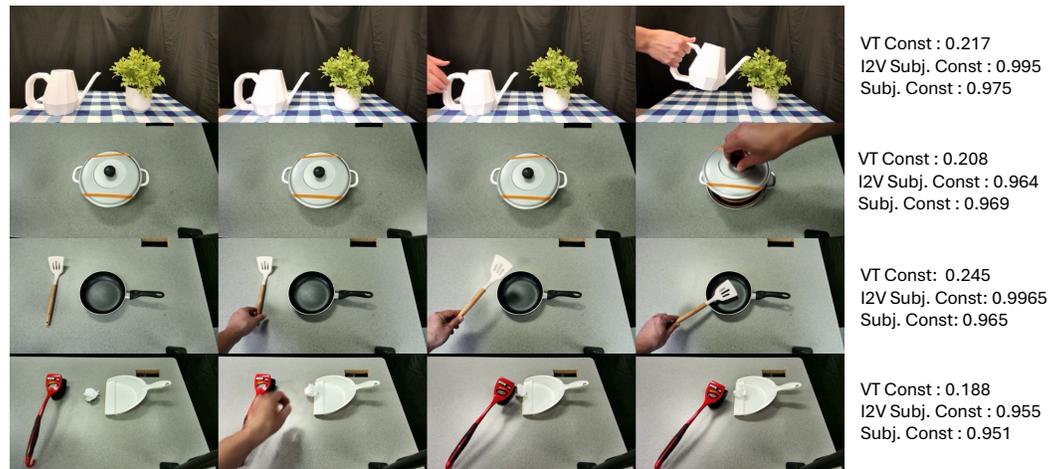


Figure 18. **Qualitative Comparison of Different Video Generative Models.** Videos generated by three models are shown in evenly sampled frames. We show VBench++ [50] metrics including video-text consistency, image-to-video subject consistency, and subject consistency.

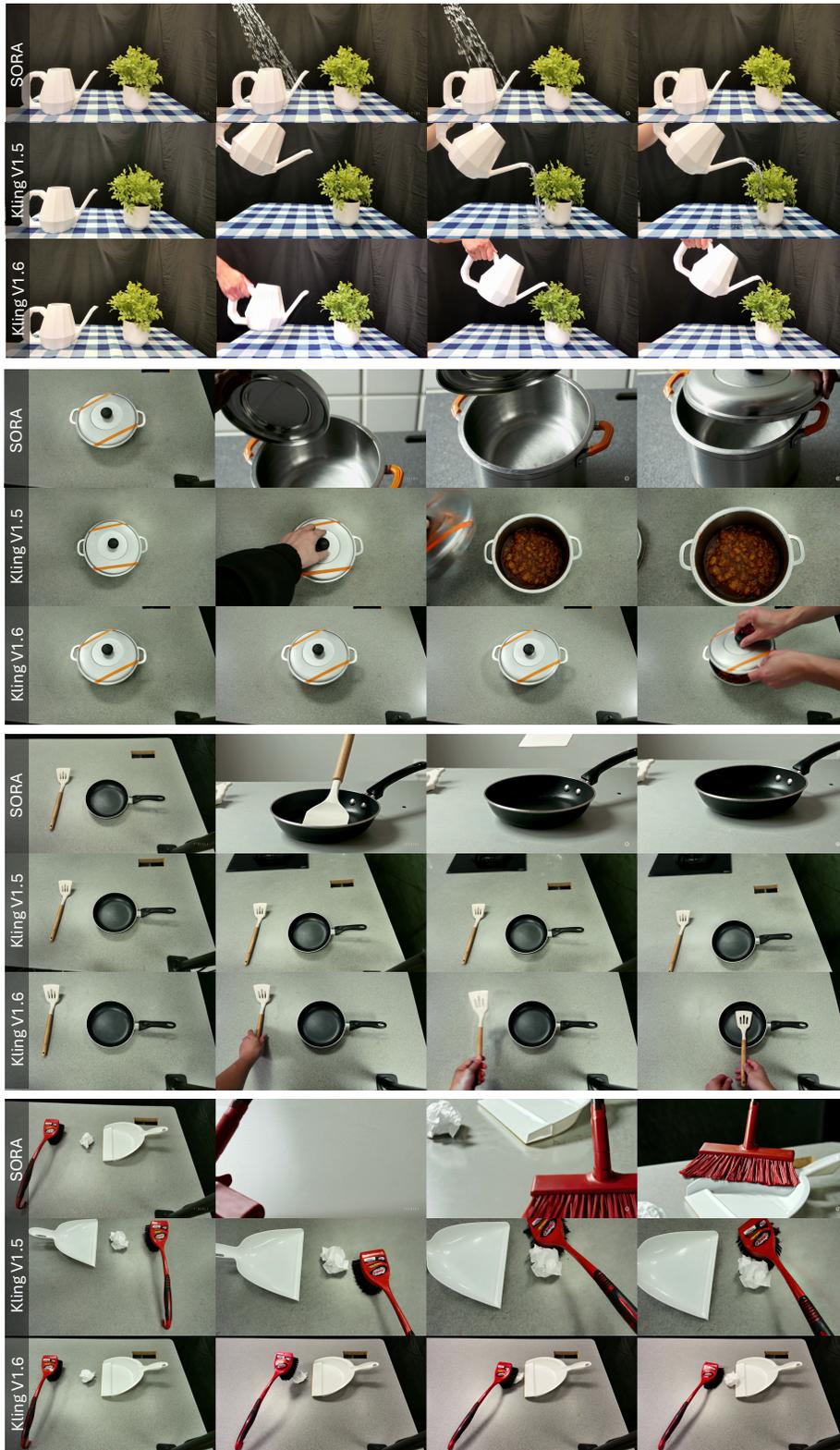


Figure 19. Qualitative comparison of video generation.