

WristCompass: Kinematic Coupling as a Learnable Visual Concept for Ego-Camera Orientation

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Recovering ego-camera orientation from manipulation video*
002 *is a prerequisite for disentangling hand motion from cam-*
003 *era motion, a key step in imitation learning from egocen-*
004 *tric demonstrations. The obvious approach, inferring ori-*
005 *entation from scene geometry, fails when hands occlude*
006 *the frame: VGGT, a 1B-parameter scene reconstruction*
007 *model, scores worse than a constant predictor on the TACO*
008 *benchmark. We identify an alternative visual concept that is*
009 *present precisely when scene geometry is absent: kinematic*
010 *coupling dynamics, the structured physical relationship be-*
011 *tween wrist motion and camera orientation imposed by the*
012 *arm-shoulder-head chain. We find that this concept is com-*
013 *compact (4D inter-wrist features outperform 126D full hand*
014 *keypoints), temporal (requiring a GRU over short windows*
015 *rather than per-frame retrieval), and physically grounded*
016 *(transferring zero-shot across datasets because it is rooted*
017 *in anatomy rather than scene appearance). Trained only*
018 *on tabletop manipulation, WristCompass transfers zero-shot*
019 *to Epic Kitchens cooking video, achieving 14.3° median*
020 *geodesic error and approaching the performance of a 1B-*
021 *parameter scene model at 200K GRU parameters.*

022 1. Introduction

023 Egocentric manipulation video is an observation from a joint
024 space: ego motion (camera orientation in $SO(3)$) entangled
025 with world state (hand and object configuration). Disentan-
026 gling these components enables ego-motion-compensated
027 hand trajectories, a prerequisite for learning manipulation
028 policies from egocentric human demonstrations [5, 20, 21].
029 Identifying the right visual concept for each component is
030 critical, but the camera orientation has proven surprisingly
031 hard to recover. Scene-based approaches fail where hands
032 occlude the frame: VGGT [17], a 1B-parameter scene recon-
033 struction model, scores 23.98° median geodesic error on the
034 bimanual manipulation benchmark TACO [10], worse than
035 a constant predictor (21.22°). The scene-geometry concept

simply does not exist when the scene is occluded. 036

We identify the concept that does exist and is sufficient 037
for recovering the camera orientation: *kinematic coupling* 038
dynamics. The arm-shoulder-head chain imposes a struc- 039
tured physical relationship between wrist motion and camera 040
orientation. This suggests that egocentric video contains a 041
compact, structured representation of ego motion encoded 042
in body dynamics, a visual concept that is present precisely 043
when scene geometry is absent. 044

This concept has three properties that characterize it as a 045
structured visual concept. **Compact:** 4D inter-wrist features 046
outperform 126D full hand keypoints (17.5° vs. 13.80°), 047
consistent with the signal concentrating in the inter-wrist 048
vector. **Temporal:** a nearest-neighbor lookup on per-frame 049
features scores 26.69°, worse than constant, while a GRU 050
over 12-frame windows achieves 13.80°; the concept is not 051
accessible to static retrieval, only to temporal models. **Phys-** 052
ically grounded: trained only on TACO tabletop manipula- 053
tion, WristCompass transfers zero-shot to Epic Kitchens [3]. 054
It achieves 14.32°, approaching the performance of VGGT 055
at 1000× lower parameter count, because the concept is 056
grounded in anatomy rather than scene appearance. 057

Contributions: (1) identifying kinematic coupling dynam- 058
ics as a compact, structured, temporally-grounded visual 059
concept sufficient for camera orientation recovery in ma- 060
nipulation video; (2) WristCompass, realizing this concept 061
from bare RGB; (3) a zero-shot transfer result demonstrat- 062
ing that the concept generalizes across datasets because it is 063
grounded in anatomy rather than scene appearance. 064

065 2. Related Work

Ego-camera pose estimation. Recovering ego-camera ori- 066
entation from video has been approached primarily through 067
scene structure. Classical approaches such as SLAM sys- 068
tems [2, 15] track sparse keypoints across frames, while 069
structure-from-motion methods [14] and learning-based re- 070
construction approaches [17, 18] reconstruct camera trajec- 071
tories from scene geometry. These methods require sufficient 072
scene texture and viewpoint diversity, assumptions that break 073

074 on close-up manipulation video where hands occlude most
075 of the frame. As we show in Sec. 4, VGGT, a 1B-parameter
076 scene reconstruction model, inherits the same failure mode
077 despite state-of-the-art performance on standard benchmarks.
078 IMU-based approaches can recover orientation but require
079 dedicated hardware unavailable in existing RGB-only video
080 archives. WristCompass targets precisely this regime, recov-
081 ering orientation post-hoc from bare monocular RGB.

082 **Ego-body, head, and hand pose estimation.** Recovering
083 body, head, or hand pose from egocentric video represents
084 a complementary line of work, but each approach assumes
085 information that is unavailable in close-up manipulation.
086 EgoAllo [19] and EgoEgo [9] recover head or body pose
087 from egocentric video but assume a wide-field view in which
088 the body is partially visible, an assumption that fails in close-
089 up manipulation. In the hand pose domain, WiLoR [13]
090 recovers 3D hand meshes from monocular RGB (we use it
091 as our keypoint extractor) and HaWoR [21] extends this to
092 world-space trajectories via SLAM-based tracking. Both
093 treat hand pose as output requiring camera pose as input.
094 WristCompass reverses both directions: we recover cam-
095 era orientation from wrist dynamics observable in the ego
096 frame, rather than inferring body pose from camera motion
097 or requiring camera pose to estimate hand trajectories.

098 3. Method

099 **Input representation.** Given an egocentric video, we ex-
100 tract 3D hand keypoints using WiLoR [13]. From each
101 frame, we take the two wrist positions, right wrist and left
102 wrist, and compute a 4D inter-wrist feature vector:

$$103 \quad \mathbf{f}_t = \left[\|\mathbf{d}_t\|, \frac{\mathbf{d}_t}{\|\mathbf{d}_t\|} \right] \in \mathbb{R}^4 \quad (1)$$

104 where $\mathbf{d}_t = \mathbf{w}_t^L - \mathbf{w}_t^R$ is the left-minus-right wrist differ-
105 ence vector. The first component is inter-wrist distance; the
106 remaining three form a unit direction vector on S^2 . Fea-
107 tures are z-score normalized using training set statistics and
108 mean-centred per video at test time (using per-video statis-
109 tics computed from the test video itself) to remove postural
110 bias without leaking cross-video information. We use 4D
111 inter-wrist features rather than full hand keypoints. The ori-
112 entation signal concentrates in the inter-wrist vector while
113 finger articulations encode subject-specific style that hurts
114 cross-subject generalization, as we validate empirically in
115 Sec. 4.

116 **Temporal model.** We process the feature sequence with a
117 two-layer GRU:

$$118 \quad \mathbf{h}_t = \text{GRU}(\mathbf{f}_{t-W:t}; \theta), \quad W = 12 \quad (2)$$

119 where $W=12$ frames (≈ 0.4 s at 30fps), selected by sweeping
120 over window sizes (Table 1). A linear head maps $\mathbf{h}_t \in \mathbb{R}^{128}$

to a 6D rotation representation [22], projected to $SO(3)$ via
121 Gram-Schmidt orthogonalization. We train by minimizing
122 geodesic loss against TACO ground-truth rotations from
123 NOKOV optical motion capture, using Adam ($\text{lr}=5 \times 10^{-4}$)
124 with early stopping on a held-out validation split (frame-level
125 80/20, used solely for early stopping). 126

Inference. WristCompass has 200K parameters (excluding
127 WiLoR as a shared frozen feature extractor). At inference,
128 we run WiLoR on each frame and compute \mathbf{f}_t when both
129 wrists are detected; frames with only one detected wrist
130 are dropped (not interpolated). On Epic Kitchens ($\approx 50\%$
131 bimanual detection), evaluation is computed over bimanual
132 frames only; GRU windows span non-consecutive frames
133 when detections are sparse. The full pipeline runs in real
134 time on CPU from bare monocular RGB. 135

136 4. Experiments

137 4.1. Datasets and Evaluation

TACO [10] is our primary training and evaluation bench-
138 mark, providing 17 subjects performing 15 tool-action-object
139 activities (5,210 frames) with a helmet-mounted RealSense
140 L515 and NOKOV optical motion-capture ground truth for
141 camera rotation. Evaluation uses Procrustes alignment, a
142 single global rotation applied to all predictions within a
143 video to minimize mean geodesic error, measuring relative
144 orientation structure rather than absolute pose. We report
145 median geodesic error, averaged over 5 random seeds on a
146 fixed train/val split (80/20, stratified by subject), with early
147 stopping on the held-out portion. 148

Epic Kitchens [3] is a large-scale egocentric cooking
149 dataset with a chest-mounted GoPro that serves as our zero-
150 shot transfer target. We use EPIC Fields [16] COLMAP
151 poses as ground truth and evaluate on 36 participants, 62
152 videos, 16,609 frames (minimum thresholds: 70% COLMAP
153 coverage, 10° constant baseline per video). COLMAP poses
154 are a proxy for ground truth and may exhibit drift in texture-
155 less or heavily occluded regions. 156

157 4.2. TACO In-Distribution Results

Table 1 and Figure 2 compare WristCompass against all
158 baselines and report architecture ablation results on TACO.
159 WristCompass achieves $13.80^\circ \pm 0.14^\circ$, outperforming all
160 baselines. 161

Three findings stand out. First, VGGT (23.98°) is worse
162 than the constant predictor (21.22°), confirming that scene-
163 based approaches fail on close-up manipulation video. Sec-
164 ond, the 126D full-keypoint MLP (17.5°) is worse than 4D
165 wrist geometry. A controlled comparison with an identical
166 GRU architecture on 126D input confirms this: 4D wins on
167 all 5 folds of subject-level cross-validation, ruling out archi-
168 tecture and capacity as confounds. Third, NN retrieval on the
169 same 4D features scores 26.69° , worse than constant, show-
170

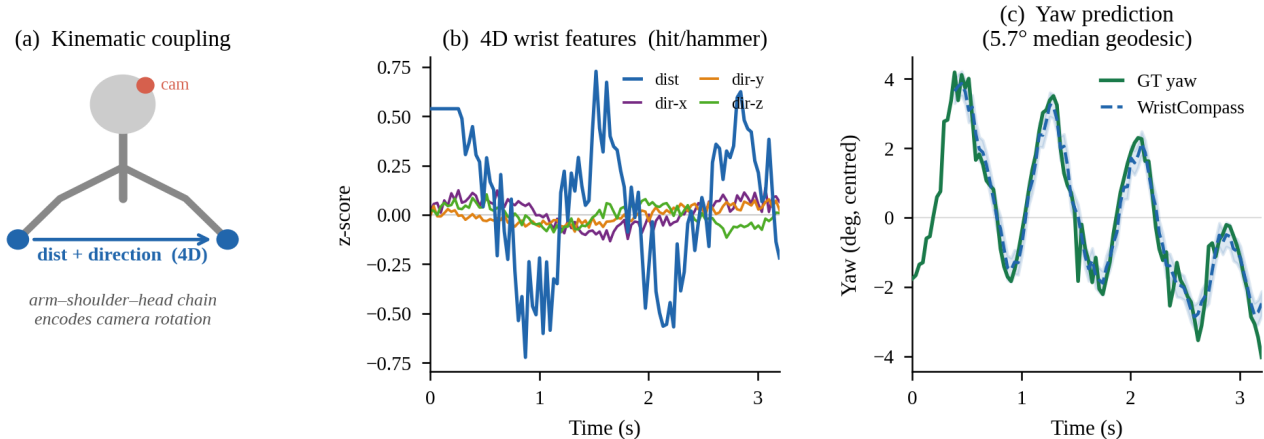


Figure 1. **WristCompass overview.** (a) Kinematic coupling: the arm-shoulder-head chain couples wrist motion dynamics to ego-camera rotation. (b) 4D inter-wrist features for a representative session: only the distance feature carries temporal variation at this timescale, direction components encode relative hand geometry. (c) WristCompass predicted yaw vs. ground-truth (GT) yaw (5.7° median geodesic). Shaded region indicates prediction uncertainty. The model tracks head rotation from wrist dynamics alone, with no scene information.

Table 1. **TACO results and ablation.** Top: baselines including scene-based (VGGT), static (NN retrieval), and full-keypoint (MLP 126D) approaches. A controlled comparison (identical GRU on 126D input) confirms 4D wins on all 5 folds of subject-level CV. Bottom: step-by-step architecture ablation, all GRU rows use 4D inter-wrist features; LR and val split were tuned jointly. All values are median geodesic error in degrees so lower is better.

Method	Geodesic ($^\circ$) \downarrow
Constant-R baseline	21.2
VGGT 1B [17]	24.0
NN retrieval (4D, static)	26.7
MLP, 42 joints (126D)	17.5
GRU $W=32$ (init)	16.9
→ $W=16$	15.8
→ LR 5×10^{-4} + val split	13.8
→ $W=12$	13.7
WristCompass (5 seeds)	13.8 ± 0.14

171 ing the orientation signal is not accessible to static retrieval
172 and requires temporal context.

173 **Per-activity analysis.** Supplementary Figure 4 shows per-
174 activity results. WristCompass beats the constant baseline
175 on 11/15 activities. It performs best on cyclic motions:
176 stir/spoon (5.9°), brush/brush (6.5°). It does worst on activ-
177 ities where head orientation decouples from wrist dynam-
178 ics: smear/glue-gun (37.3°), measure/ruler (24.3°). In these
179 cases, the subject’s wrists remain nearly stationary while the
180 head rotates to inspect different workspace regions, breaking
181 the kinematic coupling assumption.

182 4.3. Zero-Shot Transfer to Epic Kitchens

183 Trained exclusively on TACO, WristCompass achieves
184 14.32° zero-shot on Epic Kitchens, with 33 of 36 participants
185 beating the constant baseline. Despite different ground-truth

Table 2. **Epic Kitchens comparison** (36 participants, zero-shot for WristCompass). Parameter count for WristCompass refers to the GRU only; WiLoR keypoint extractor is a shared prerequisite not counted in this comparison.

Method	Params	Geodesic ($^\circ$) \downarrow
Constant-R	—	18.5
VGGT [17]	1B	12.8
WristCompass	200K	14.3

sources (mocap vs. COLMAP), camera mountings (helmet
vs. chest), and activity domains, performance is comparable
to in-distribution TACO (13.80°). Figure 3 shows the per-
participant breakdown. P14 (constant $66.1^\circ \rightarrow$ GRU 5.3°)
demonstrates strong signal extraction when head movement
is rich. The three failures (P10, P13, P31) are single-video
participants with limited COLMAP pose quality.

As shown in Table 2, compared to VGGT (12.83°), Wrist-
Compass (14.32°) trails by 1.5° at $1000\times$ lower GRU pa-
rameter count. COLMAP ground truth may favor VGGT,
which shares its scene-geometry assumptions. Both substan-
tially beat the constant baseline (18.54°). The methods have
complementary failure modes: VGGT wins on discrete ma-
nipulation tasks where scene structure changes predictably,
while WristCompass wins on cyclic tasks where inter-wrist
dynamics are more stable. On TACO, where scene features
are largely occluded, WristCompass wins by over 10° .

5. Discussion and Limitations

Failure modes and scope. WristCompass works best when
head movement is rich (constant baseline $> 15^\circ$) and both
hands are consistently visible. P14 (constant 66° , GRU 5°)
shows the potential ceiling. Cyclic tasks (stir, brush) are
particularly well-suited. The model cannot improve on the
constant baseline when head orientation is near-static: on

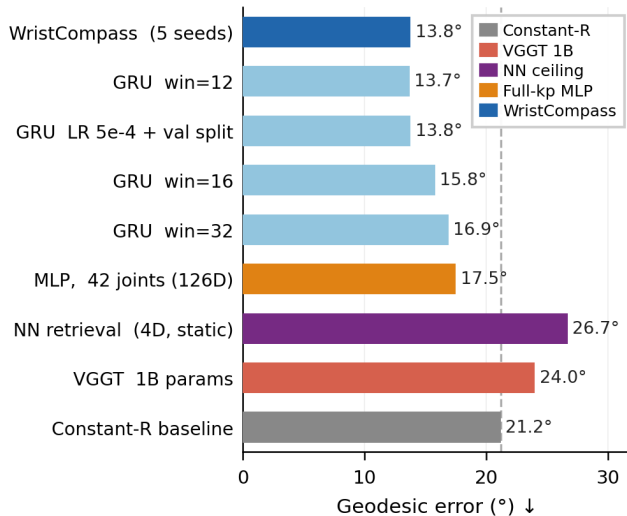


Figure 2. **TACO results and ablation.** NN retrieval (26.7°) is worse than the constant baseline (21.2°); the orientation signal is not accessible to static retrieval. WristCompass (13.8°) outperforms VGGT 1B (24.0°) at 200K GRU parameters.

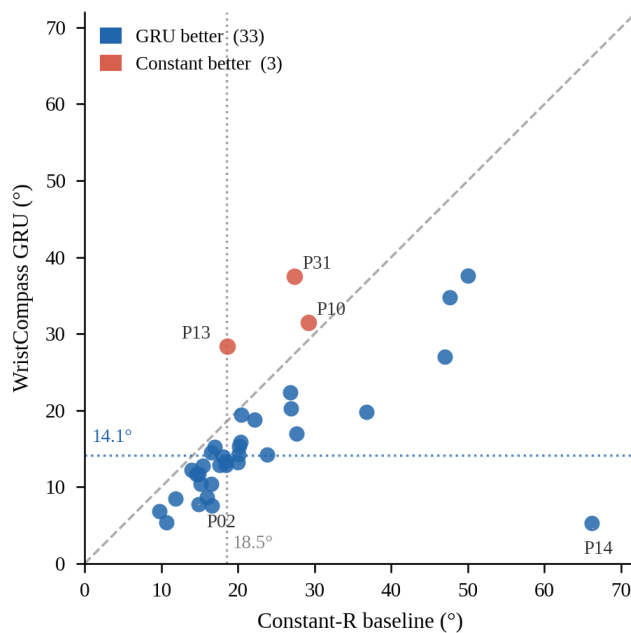


Figure 3. **Epic Kitchens zero-shot (36 participants).** Points below diagonal: WristCompass beats constant baseline (33/36, blue). Points above: failures (3/36, red). P14 (constant 66°, GRU 5°) is a striking outlier.

close-up manipulation datasets with limited head movement such as HOT3D (5.1°) and ARCTIC (5.5°), there is insufficient rotational signal to recover [1, 4]. Activity-level failures (smear 37.3°, measure 24.3°) occur when wrist position is constrained independently of head orientation, breaking kinematic coupling. Both of the datasets we evaluate, TACO and Epic Kitchens, involve standing manipulation, transfer to seated or full-body activities remains untested.

Kinematic coupling as a physical visual concept. Most concept learning methods, including β -VAE [7], Slot Attention [12], and Concept Bottleneck Models [8], discover structure from co-occurrence statistics in training data, so their generalization is bounded by training distribution diversity. Inductive-bias approaches impose architectural priors such as competition and capacity limits [11], but these are still learned constraints rather than physical ones. Kinematic coupling occupies a third position: its 4D inter-wrist representation is not learned from data but dictated by biomechanics—the physical structure of the human arm-shoulder-head chain. This makes it *physically grounded*: it generalizes across environments and activity domains because human anatomy is shared, not because the training distribution covers the test distribution. It is also *intrinsically temporal*—the NN retrieval ablation shows the concept is invisible in any single frame and only emerges over time. The zero-shot transfer from TACO to Epic Kitchens provides direct empirical evidence for both properties.

Limitations. Several limitations of the current approach point toward future directions. WristCompass recovers relative camera orientation, which is how the camera rotates over time, rather than absolute pose. This is sufficient for downstream tasks such as ego-motion-compensated hand trajectory extraction, where action representations depend on frame-to-frame rotation rather than global heading, but absolute calibration remains an open problem. Single-hand frames ($\approx 48\%$ of Epic Kitchens frames) currently fall back to a constant predictor; improving single-hand fallback is a clear direction for future work. Epic Kitchens evaluation uses raw WiLoR-mini keypoints without smoothing ($\approx 50\%$ bimanual detection); Kalman-RTS smoothing improves TACO ($\approx 100\%$ detection) but degrades Epic Kitchens by replacing temporal signal with near-constant interpolations, suggesting that detection-aware smoothing strategies warrant further investigation.

6. Conclusion

We identify kinematic coupling dynamics, the temporal relationship between bimanual wrist motion and ego-camera orientation imposed by the arm-shoulder-head chain, as a compact, physically-grounded visual concept for recovering the camera orientation of manipulation video. WristCompass realizes this concept with 200K GRU parameters from bare monocular RGB, outperforming a 1B-parameter scene model on close-up manipulation and generalizing zero-shot to kitchen video. Future directions include depth integration, explicit ego-world disentanglement [6], and downstream robot policy learning.

References

- [1] Prithviraj Banerjee, Sindi Shkodrani, Pierre Moulon, Shreyas Hampali, Fan Zhang, Jade Fountain, Edward Miller, Selen Ba-

- 269 sol, Richard A. Newcombe, Robert Wang, Jakob Julian Engel, 325
270 and Tomas Hodan. Introducing hot3d: An egocentric dataset 326
271 for 3d hand and object tracking. *ArXiv*, abs/2406.09598, 2024. 327
272 4
- 273 [2] Carlos Campos, Richard Elvira, Juan J. Gómez Rodríguez, 328
274 José M. M. Montiel, and Juan D. Tardós. ORB-SLAM3: An 329
275 accurate open-source library for visual, visual-inertial and 330
276 multi-map SLAM. *IEEE Transactions on Robotics*, 2021. 1 332
277 [3] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, 333
278 Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Da- 334
279 vide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, 335
280 and Michael Wray. Scaling egocentric vision: The EPIC- 336
281 KITCHENS dataset. In *ECCV*, 2018. 1, 2
- 282 [4] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed 337
283 Kocabas, Manuel Kaufmann, Michael J. Black, and Otmar 338
284 Hilliges. Arctic: A dataset for dexterous bimanual hand- 339
285 object manipulation. *2023 IEEE/CVF Conference on Com- 340
286 puter Vision and Pattern Recognition (CVPR)*, pages 12943– 341
287 12954, 2022. 4
- 288 [5] Hongming Fu, Wenjia Wang, Xiaozhen Qiao, Shuo Yang, 342
289 Zheng Liu, and Bo Zhao. Egograsp: World-space hand- 343
290 object interaction estimation from egocentric videos. *ArXiv*, 344
291 abs/2601.01050, 2026. 1 345
- 292 [6] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, 346
293 Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar 347
294 Ashutosh, et al. Ego-Exo4D: Understanding skilled human 348
295 activity from first- and third-person perspectives. In *CVPR*, 349
296 2024. 4 350
- 297 [7] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, 351
298 Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and 352
299 Alexander Lerchner. β -VAE: Learning basic visual concepts 353
300 with a constrained variational framework. In *ICLR*, 2017. 4 354
301 [8] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Muss- 355
302 mann, Emma Pierson, Been Kim, and Percy Liang. Concept 356
303 bottleneck models. In *ICML*, 2020. 4 357
- 304 [9] Jiaman Li, Karen Liu, and Jiajun Wu. Ego-body pose estima- 358
305 tion via ego-head pose estimation. In *CVPR*, 2023. 2
- 306 [10] Yun Liu, Haolin Yang, Xu Xu, Mingsheng Ding, Weicheng Li, 359
307 Yuxin Li, Ziyu Liu, Jianyu Luo, Jing Cheng, and Li Cheng. 360
308 TACO: Benchmarking generalizable bimanual tool-action- 361
309 object understanding. In *CVPR*, 2024. 1, 2
- 310 [11] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar 362
311 Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier 363
312 Bachem. Challenging common assumptions in the unsu- 364
313 pervised learning of disentangled representations. In *ICML*, 365
314 2019. 4
- 315 [12] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, 366
316 Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, 367
317 Alexey Dosovitskiy, and Thomas Kipf. Object-centric learn- 368
318 ing with slot attention. In *NeurIPS*, 2020. 4
- 319 [13] Rolandos Alexandros Potamias, Jinglei Shu, German Bar- 369
320 quero, Cristina Palmero, Sergio Escalera, and Stefanos 370
321 Zafeiriou. WiLoR: End-to-end 3D hand localization and 371
322 reconstruction in-the-wild. In *ECCV*, 2024. 2
- 323 [14] Johannes L. Schönberger and Jan-Michael Frahm. Structure- 372
324 from-motion revisited. In *CVPR*, 2016. 1
- [15] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for 325
monocular, stereo, and rgb-d cameras. In *Neural Information 326
Processing Systems*, 2021. 1 327
- [16] Vadim Tschernezki, Ahmad Darkhalil, Zhifan Zhu, David 328
Fouhey, Iro Laina, Diane Sheratt, Mike Brookes, Roberto 329
Cipolla, Spyridon Leonardos, and Dima Damen. EPIC Fields: 330
Marrying 3D geometry and video understanding. *NeurIPS*, 331
2023. 2 332
- [17] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea 333
Vedaldi, Christian Rupprecht, and David Novotny. VGGT: 334
Visual geometry grounded deep structure from motion. In 335
CVPR, 2025. 1, 3 336
- [18] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris 337
Chidlovskii, and Jérôme Revaud. Dust3r: Geometric 3d 338
vision made easy. *2024 IEEE/CVF Conference on Computer 339
Vision and Pattern Recognition (CVPR)*, pages 20697–20709, 340
2023. 1 341
- [19] Jiaman Ye, Yiye Ye, Manolis Savva, Angel X. Chang, and 342
Li Yi. EgoAllo: Egocentric human motion estimation via 343
explicit alignment with a grounded world coordinate system. 344
In *ECCV*, 2024. 2 345
- [20] Zhengdi Yu, Stefanos Zafeiriou, and Tolga Birdal. Dyn- 346
hamr: Recovering 4d interacting hand motion from a dynamic 347
camera. *2025 IEEE/CVF Conference on Computer Vision 348
and Pattern Recognition (CVPR)*, pages 27716–27726, 2024. 349
1 350
- [21] Jinglei Zhang, Jiankang Deng, Chao Ma, and Rolan- 351
dos Alexandros Potamias. Hawor: World-space hand motion 352
reconstruction from egocentric videos. *2025 IEEE/CVF Con- 353
ference on Computer Vision and Pattern Recognition (CVPR)*, 354
pages 1805–1815, 2025. 1, 2 355
- [22] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao 356
Li. On the continuity of rotation representations in neural 357
networks. In *CVPR*, 2019. 2 358

359 **Supplementary Material**

360 **A. Controlled 4D vs. 126D comparison.** Table 3 reports a
361 head-to-head comparison using identical GRU architectures
362 ($W=12$, hidden=128, 2 layers) on 4D inter-wrist vs. 126D
363 full-keypoint input, evaluated via 5-fold subject-level cross-
364 validation. 4D outperforms 126D on all 5 folds, ruling out
365 architecture and capacity as confounds. The 126D model
366 overfits to subject-specific finger articulations that do not
367 transfer to held-out subjects.

Table 3. **4D vs. 126D GRU** (5-fold subject-level CV, 5 seeds per fold). 4D wins on every fold.

Fold	4D ($^{\circ}$)	126D ($^{\circ}$)	Δ
0	4.67	6.81	-2.14
1	6.44	7.98	-1.54
2	3.68	7.82	-4.14
3	8.56	12.14	-3.58
4	2.68	6.66	-3.98
Mean	4.77 ± 0.24	7.84 ± 0.25	-3.07

368 **B. Per-axis error decomposition.** On TACO, WristCompass
369 achieves per-axis median errors of yaw 8.44° , pitch
370 5.28° , roll 3.87° . Yaw (left-right head rotation) carries the
371 most error, consistent with the inter-wrist vector being most
372 informative about lateral head movement. Pitch and roll
373 are better constrained by the arm-shoulder-head kinematic
374 chain.

375 **C. All-frames Epic Kitchens evaluation.** The main paper
376 reports 14.32° on bimanual frames only ($\approx 52\%$ of frames).
377 For a deployment-realistic estimate, we blend GRU predic-
378 tions on bimanual frames with the constant-R baseline on
379 single-hand/no-hand frames:

Table 4. **Epic Kitchens blended evaluation** (62 videos, 16,609 total frames).

Metric	Geodesic ($^{\circ}$) \downarrow
Constant-R (all frames)	18.54
GRU-only (bimanual, 52%)	14.32
Blended median (per-video)	16.56

380 The 2.0° gap between GRU-only and blended reflects the
381 48% of frames where only one or no hands are detected.
382 Improving single-hand fallback is a clear direction for future
383 work.

384 **D. Per-activity breakdown.** Figure 4 shows per-activity
385 TACO results.

386 **E. Failure case: kinematic decoupling.** Figure 5 shows
387 a representative failure (empty/bowl/plate, session 8, 21.3°
388 geodesic). Between $t=4-6$ s, the subject’s head rotates $\approx 14^{\circ}$
389 in yaw while wrists remain stationary — the subject looks

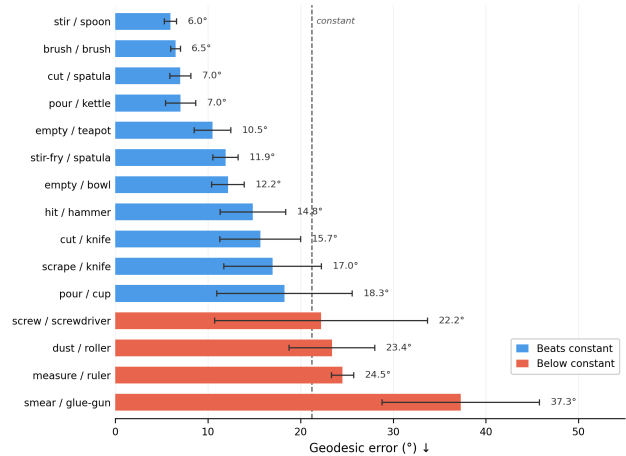


Figure 4. **Per-activity TACO results** (5 seeds, IQR error bars). Blue: beats constant (11/15). Red: below constant (4/15). Cyclic motions (stir, brush) are easiest; decoupled motions (smear, measure) are hardest.

away from their hands to inspect the workspace. Wrist-Compass cannot track this head rotation because the kinematic coupling between wrist motion and head orientation is broken. This failure mode is systematic: it accounts for the worst per-activity results (smear/glue-gun 37.3° , measure/ruler 24.3°) where wrist position is constrained independently of gaze direction.

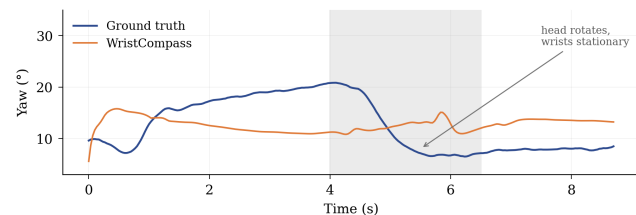


Figure 5. **Failure case: kinematic decoupling.** GT yaw (blue) drops 14° between $t=4-6$ s while WristCompass prediction (orange) remains flat. The subject’s head rotates independently of their stationary wrists.