

Pole-Arina: A Privacy-Preserving Dataset and Benchmark for Static Pole Tricks

Anonymous CVsports submission

Paper ID 7

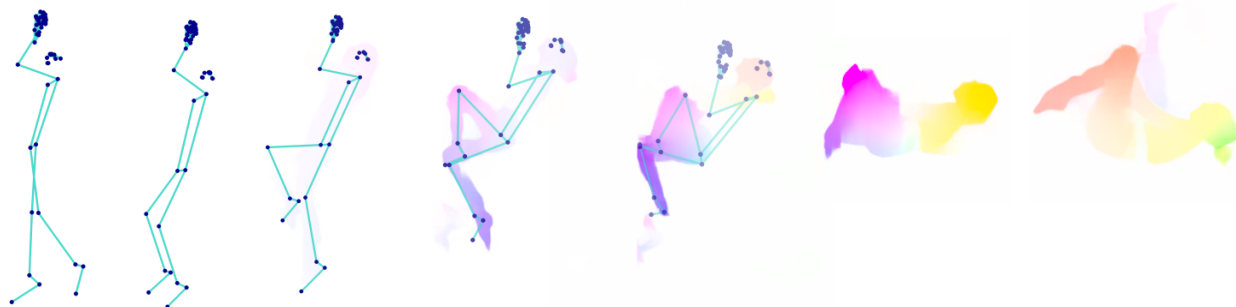


Figure 1. *Pole-Arina* replaces identifiable RGB with privacy-preserving motion representations. We extract 2D pose skeletons and dense optical flow, enabling clip-level trick recognition, frame-wise phase understanding, and coaching-oriented analysis for static pole tricks under strong privacy constraints.

Abstract

001 *Computer vision for sports lacks datasets and models that*
 002 *target pole sports, an activity marked by self-occlusion,*
 003 *body inversions, and sustained contact with the pole. We*
 004 *address this gap by introducing **Pole-Arina**, a curated,*
 005 *privacy-preserving dataset for markerless analysis of static*
 006 *pole tricks. The collection comprises 836 videos from*
 007 *58 participants and provides two appearance-suppressing*
 008 *modalities derived from the recordings: pose skeleton se-*
 009 *quences (2D joints) and dense optical flow. We anno-*
 010 *tate each clip with both per-video trick labels and per-*
 011 *frame labels that capture temporal structure (e.g., `floor`,*
 012 *`on_pole`, and trick-specific pose states), enabling evalu-*
 013 *ation of clip-level and frame-wise recognition under pri-*
 014 *vacuity constraints. As reference baselines, we benchmark*
 015 *lightweight temporal models across all benchmark settings*
 016 *(skeleton/flow \times video/frame) and analyze common con-*
 017 *fusions and imbalance effects. In addition, we provide*
 018 *a geometry-aware analysis module that measures trick-*
 019 *specific body orientation, joint alignments, and proximities*
 020 *to produce interpretable overlays and actionable feedback.*
 021 *To the best of our knowledge, this is the first dataset and*
 022 *benchmark targeting recognition and analysis of pole tricks*
 023 *from privacy-preserving motion representations. We will re-*
 024 *lease the dataset, annotations, and code upon acceptance.*
 025

1. Introduction

Video-based coaching systems increasingly combine pose
 estimation with temporal sequence models to analyze hu-
 man motion and deliver corrective feedback across fitness,
 sports, and dance [? ? ? ? ?]. Markerless approaches
 are especially compelling because they avoid wearables and
 preserve natural movement [?]. For fine-grained technique
 analysis, models must support both clip-level recognition
 and temporally localized, per-frame understanding under
 occlusions, viewpoint changes, and rapid transitions, con-
 ditions that motivate curated datasets and learning pipelines
 built on privacy-preserving representations. One such un-
 explored example is pole dancing.

As a sport, pole dancing has grown into a recognized
 athletic discipline with standardized competition rules (e.g.,
 within the International Pole and Aerial Sports Federation
 since 2009) that emphasize technique and objective evalua-
 tion [?]; while debates about hypersexualization persist, a
 sport-centered framing prioritizes inclusive, skill-based as-
 sessment [?]. Static pole tricks demand strength, flexibil-
 ity, precise alignment, and sustained skin-pole contact (ren-
 dering wearables impractical); they are practiced in studios
 and at home, where dancers commonly self-record. Brief,
 discriminative end-poses are preceded by visually similar
 approach phases, occlusions are frequent, viewpoints vary,
 and naming conventions differ across studios, complicating
 both recognition and temporal localization. To the best of
 our knowledge, no prior work provides a public dataset for

054	pole trick recognition.	
055	We address this gap with <i>Pole-Arina</i> , a new privacy-	
056	preserving dataset and benchmark for static pole tricks. The	
057	dataset contains 836 clips spanning six foundational tricks,	
058	performed by multiple volunteers with varied skill levels	
059	and including both successful and failed attempts, totaling	
060	212,574 labeled frames. To mitigate privacy concerns while	
061	retaining kinematic fidelity, we do not release RGB video;	
062	instead, we provide two derived, appearance-suppressing	
063	modalities: (i) pose skeleton sequences (2D joints with con-	
064	fidence) and (ii) dense optical flow computed from the stan-	
065	dardized clips - Figure 1. We further provide two annota-	
066	tion granularities: per-video trick labels for clip-level action	
067	recognition and per-frame labels capturing temporal struc-	
068	ture (e.g., <code>floor</code> , <code>on_pole</code> , and trick-specific pose states),	
069	enabling temporally localized evaluation.	
070	On top of this dataset, we benchmark lightweight tempo-	
071	ral models across both modalities and both supervision	
072	settings (per-video and per-frame). Clip-level recognition	
073	is strong for skeletons (91.26% accuracy), while optical	
074	flow is more challenging (66.67%), reflecting the fact that	
075	pose extraction supplies a strong inductive bias whereas	
076	flow must recover discriminative structure from motion	
077	alone. For frame-wise recognition, overall accuracies reach	
078	86.99% (skeleton) and 81.23% (flow), but performance is	
079	affected by severe label imbalance in the per-frame taxon-	
080	omy, highlighting open challenges for fine-grained, privacy-	
081	preserving recognition. Beyond recognition, we include an	
082	interpretable geometric scoring module that encodes techni-	
083	que cues (angles, alignments, distances) on the skeleton	
084	to produce actionable, per-rule feedback, and we evaluate a	
085	coaching-oriented prototype in a controlled user study.	
086	Our contributions are as follows:	
087	• Pole Dancing Multimodal Privacy-Preserving Dataset.	
088	We introduce a dataset of 836 clips spanning six fun-	
089	damental static pole tricks, totaling 212,574 labeled	
090	frames, and release privacy-preserving representations	
091	(pose skeleton sequences and dense optical flow) together	
092	with per-video and per-frame labels.	
093	• Unified Benchmark Across Modalities and Label	
094	Granularities. We define and benchmark four recogni-	
095	tion settings (skeleton/flow \times video/frame), enabling sys-	
096	tematic evaluation of clip-level trick recognition and tem-	
097	porally localized per-frame classification under appear-	
098	ance suppression.	
099	• Baseline Models, Analysis, and Coaching-Oriented	
100	Feedback. We establish lightweight baseline results for	
101	all settings, analyze confusions and imbalance effects,	
102	and provide an interpretable geometric scoring module	
103	with a user-study evaluation in a coaching prototype.	
104	Together, the dataset, benchmark, and baselines establish	
105	pole dance as a new computer vision problem setting for	
106	privacy-aware, fine-grained athletic technique understand-	
	ing, leaving clear headroom for improved temporal model-	107
	ing, robustness, and minority-class recognition.	108
	2. Related Work	109
	Motion capture for athletic analysis. Optical marker-	110
	based mocap remains the gold standard for accuracy but is	111
	costly, intrusive, and impractical outside controlled labs [?	112
	?]. Wearable inertial systems improve portability yet suffer	113
	from drift and magnetic disturbances and can impede natu-	114
	ral movement in contact sports [?]. Markerless, vision-	115
	based pose estimation offers the best trade-off for in-situ	116
	training, preserving unconstrained motion while requiring	117
	only commodity cameras [?].	118
	Pose estimation. Modern methods estimate body key-	119
	points directly from RGB, from DeepPose regression [?] and	120
	multi-person Part Affinity Fields in OpenPose [?] to	121
	mobile or edge pipelines such as BlazePose or MediaPipe	122
	Pose with monocular estimation [? ?]. Complementary	123
	open-source implementations and toolkits, including Open-	124
	PifPaf (a composite-field, bottom-up formulation) and MM-	125
	Pose, provide strong baselines and practical training and in-	126
	ference pipelines across a range of backbones and deploy-	127
	ment settings [? ? ?]. These approaches enable non-	128
	invasive capture in studios and at home, but remain sensi-	129
	tive to occlusions, fast motion, and viewpoint changes, con-	130
	ditions common in technical sports and dance.	131
	Optical flow and privacy. Optical flow encodes motion	132
	as per-pixel displacement fields and suppresses most ap-	133
	pearance cues, making it a common input for action recog-	134
	nition when texture or identity should be minimized [?	135
]. Recent work explicitly studies appearance-free (motion-	136
	only) recognition, showing that competitive performance	137
	is possible even when static visual cues are removed, and	138
	motivating architectures that recover and leverage explicit	139
	flow [?]. In privacy-preserving action recognition more	140
	broadly, motion-centric encodings are often used as prag-	141
	matic compromises: they reduce identifiability while retain-	142
	ing discriminative temporal structure [?].	143
	Sports recognition and temporal modeling. For se-	144
	quence understanding, recurrent models remain strong on	145
	modest datasets [? ?], while Transformers capture long-	146
	range dependencies given sufficient data [? ?]. Literature	147
	reports a common pipeline: pose estimation followed by	148
	temporal modeling or rule-based assessment for technique	149
	analysis, counting, phase segmentation, and feedback [?	150
]. Beyond recognition, pose-based scoring correlates with	151
	expert judgments in judged sports [?], and computer vi-	152
	sion supports biomechanical analysis in both individual and	153

154 team settings [? ?]. In dance, pose-based evaluation at-
155 tains high agreement with expert ratings despite occlusions
156 and rapid motion [?].

157 **Sports datasets.** Large-scale action sets (e.g., Sports-
158 1M [?], Kinetics [?]) initiated video recognition but are
159 coarse-grained and RGB-centric, offering limited pose or
160 phase supervision. Fine-grained, judged-sport benchmarks
161 add temporal structure and skill labels, e.g., Diving48 [?]
162 and FineGym [?], for subtle technique discrimination and
163 phase localization. Pose-annotated sets such as Penn Ac-
164 tion [?] and PoseTrack [?] enable skeleton-based analysis
165 and tracking, while domain datasets target tactics (Soccer-
166 Net [?], VNL-STES [?]) to study multi-person dynam-
167 ics. Despite this breadth, we find no public dataset for *pole*
168 *dancing*. Existing resources neither capture its contact-rich
169 setting, frequent self-occlusions, and brief discriminative
170 end-poses, nor address the privacy constraints that make
171 appearance-based RGB release problematic. In contrast to
172 sports benchmarks centered on RGB appearance, Pole-
173 Arina emphasizes privacy-preserving motion representa-
174 tions for contact-heavy static pole poses, where sustained
175 body-pole interaction and fine-grained end-pose structure
176 are central.

177 3. Dataset

178 We introduce a privacy-preserving, skeleton and optical
179 flow dataset for fine-grained recognition of *static* pole
180 tricks. Our data consists of per-frame pose sequences (ob-
181 tained as MediaPipe poses [?], having 75 joints with con-
182 fidence values), computed optical flow using RAFT, and
183 phase-aware annotations suitable for trick recognition and
184 end-pose evaluation.

185 3.1. Data Collection & Labeling

186 The data collection goal was to create an ethical, realistic
187 representation of fundamental static pole tricks. To that
188 end, the dataset targets six foundational tricks and balances
189 feasibility for novices with sufficient discriminability for
190 modeling. Annotations structure each clip into semantically
191 meaningful temporal states.

192 **Scope and terminology.** Pole dance comprises *dance*
193 *moves* (at least one foot stays on the floor), *spins* (airborne
194 rotations), *floor work* (near-ground movement), and *tricks*
195 (static/semi-static shapes). Poles operate in *static* or *spin-*
196 *ning* mode. We adopt terminology from Spin City’s Pole
197 Bible and IPSF [? ?]. We target *static pole tricks* to obtain
198 consistent camera-facing orientation, fewer self-occlusions,
199 and reduced motion blur, conditions that improve pose es-
200 timation and geometric measurement on commodity video.
201 Each trick follows three phases: *entry* (approach/contact),

transition (movement on pole), and *end pose* (held target
202 shape). 203

204 **Selection criteria.** To balance feasibility, safety, and dis-
205 criminability, we included six foundational tricks across
206 two posture types: upright (*Layout*, *Pin-Up*, *Wrist Seat*) and
207 inverted (*Straddle Invert*, *Gemini*, *Crucifix*) - Figure 2. Up-
208 right tricks share similar entries but end in distinct shapes;
209 *Layout* vs. *Pin-Up* form intentional near-neighbors (Fig-
210 ure 2a). Inverted tricks begin from a basic invert grip and
211 diverge to clearly different end poses (Figure 2b), adding
212 biomechanical complexity while remaining achievable for
213 athletic beginners to lower-intermediate dancers.

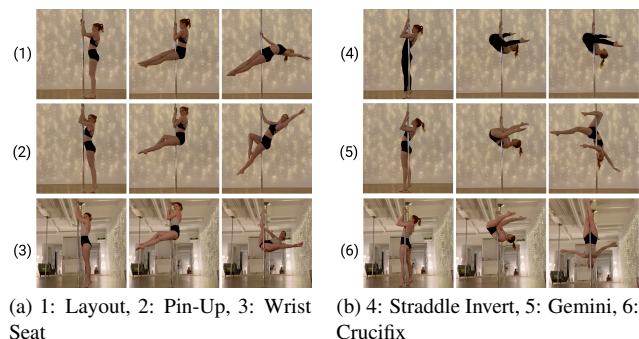


Figure 2. Progression of each trick, highlighting similar entries and transitions before the final pose.

3.1.1. Data Acquisition

214 Collecting high-quality and representative data was a cru-
215 cial step, as together with the annotations, it provides the
216 ground truth for our dataset generation and model evalua-
217 tion. The aim was to assemble a large enough collection
218 of clips, covering a balanced distribution of the six selected
219 tricks and including both successful and failed attempts. We
220 sought balanced sampling across tricks, experience levels,
221 gender, and age; nevertheless, residual demographic skew
222 and estimator-specific biases may persist. Beyond the scale
223 and diversity of the data, particular attention was given to
224 the ethical aspects of data collection. 225

226 We acquired videos via two routes: (i) *online submis-*
227 *sions* through a form with a clearly-defined consent pro-
228 tocol, standardized filming guidance (smartphone, frontal
229 full-body framing, even lighting, minimal background), and
230 direct upload; and (ii) *in-class recordings* in dedicated ses-
231 sions where an instructor demonstrated and supervised each
232 trick with a fixed camera, ensuring consistent viewpoints,
233 safer spotting for beginners, and better balance across the
234 six tricks. Participants were encouraged to submit multi-
235 ple attempts (including incomplete or failed tries) to capture
236 natural variability.

237 The final dataset comprises 836 clips from $N = 58$

238 participants. The protocol deliberately encouraged multi-
239 ple attempts per participant to capture natural variability in
240 the progression of a trick. Section 3.3 presents a detailed
241 breakdown of the final dataset. Table 3 compiles concise,
242 syllabus-aligned instructions and categories for the six selected
243 tricks.

244 **Ethical considerations.** All participants were fully in-
245 formed about the project, data handling, and withdrawal
246 rights, and provided written consent; collection and pro-
247 cessing complied with local data-protection policies. In
248 pole sports, participants typically wear short shorts and
249 a sports bra/tank top to maximize friction with the pole,
250 so identity cues extend to the full body. Combined
251 with mirrors and distinctive studio backgrounds, credible
252 de-identification would require near full-body obfuscation
253 and extensive background handling, substantially reducing
254 video utility. Hence, to protect privacy, we only use and re-
255 lease skeleton joint coordinates, optical flow, and per-video
256 and per-frame labels; raw videos are used solely for creat-
257 ing the dataset and are not shared. No personal identifiers,
258 faces, audio, or background content are stored.

259 3.1.2. Annotation Protocol

260 Our dataset provides two annotation granularities: per-
261 video labels and per-frame labels. The per-video labels
262 capture the performed trick (L, P, W, S, G, I), while per-
263 frame labels capture the broad temporal structure of each
264 clip through states such as `floor`, `on_pole`, and trick-
265 specific pose labels (e.g., `*_pose`), enabling supervision of
266 temporal progression from setup to execution. As expected
267 for unconstrained performance videos, the frame-level label
268 distribution is imbalanced, with support states (`floor` and
269 `on_pole`) occupying a larger fraction of frames than termi-
270 nal trick poses; to support transparent benchmarking, we
271 provide the full label definitions and class frequencies with
272 the dataset. The annotation protocol could be extended in
273 future work toward finer phase-aware per-frame labels for
274 the support states, e.g., by separating entry types and inter-
275 mediate transition states.

276 Frame labels were annotated with expert guidance at
277 frame-accurate temporal boundaries. We do not provide
278 subjective per-phase quality scores, as these were found
279 to be less consistent across clips; instead, form assessment
280 is handled separately through a geometric rule-based mod-
281 ule. Annotation efficiency and consistency across the full
282 dataset were supported by a lightweight custom annotation
283 tool. The dataset further includes per-video metadata such
284 as the trick label and the performer’s experience level.

285 3.2. Dataset Generation

286 Our pipeline converts raw training videos into stable,
287 privacy-preserving representations for learning and evalu-

Method	Missed Frames	Avg. # Joints	Avg. Conf with Face	Avg. Conf. w/o Face
MediaPipe	0.12%	60/75	0.76	0.76
MMPose	0.09%	129/133	0.77	0.69
OpenPifPaf	38.9%	52/133	0.27	0.15
OpenPose	2.9%	37/67	0.39	0.39

Table 1. Quantitative results of multiple skeleton extractors run on the entire video collection.

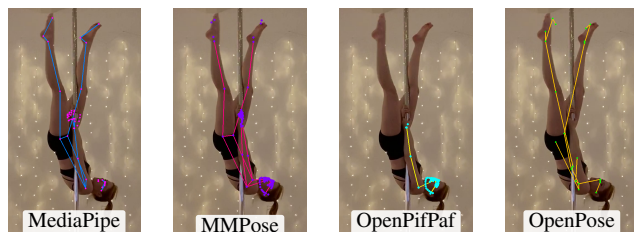


Figure 3. Qualitative example of all extractors on one single frame.

288 Instead of operating on pixels, we extract body skele-
289 tons and optical flow vectors of all the recorded videos.

290 **Skeleton Extraction.** Accurate assessment of pole tricks
291 requires reliable detection of both fingers and toes. We
292 therefore evaluate four full-body pose estimators (Figure 3):
293 OpenPose BODY_25 with hand keypoints [?], Medi-
294 aPipe Holistic (model complexity 2, confidence thresh-
295 old 0.3) [?], MMPose (RTMPose-L, 384×288 [?])
296 with COCO WholeBody keypoints [?], and OpenPifPaf
297 (ShuffleNetV2K30-WholeBody [?]) with COCO Whole-
298 Body keypoints [?]. Some methods predict more keypoints
299 than needed; for example, MMPose and OpenPifPaf include
300 68 facial landmarks that are not relevant to our use case.

301 For all four methods, we compute extraction statistics
302 over the full dataset, including the percentage of missed
303 frames (defined as all joints having confidence below 0.3),
304 the average number of detected joints per frame, and the av-
305 erage confidence score (Table 1). MMPose yields the lowest
306 number of missed frames and the highest overall confidence
307 on the full keypoint set. However, after excluding facial
308 keypoints, which are not useful in our setting, its average
309 confidence drops significantly, and the number of missed
310 frames increases by 31. MediaPipe, by contrast, remains
311 only slightly behind in total missed frames (262 vs. 201)
312 and shows a more favorable error pattern: most of its missed
313 frames (156) occur at the end of videos, when no person is
314 visible, whereas MMPose failures are more often scattered
315 throughout the central parts of clips, including contiguous
316 intervals. For this reason, we select MediaPipe both as the
317 released skeleton representation and as input to our example
318 action-recognition pipeline.

Trick	Layout	Pin-Up	Wrist Seat	Straddle Invert	Gemini	Inverted Crucifix
#Videos	188	167	131	144	88	94

Table 2. Number of videos per trick end-pose.

319 **Optical Flow Extraction.** We extracted optical flow with
 320 RAFT [?] using the official PyTorch implementation with
 321 the RAFT-Large model and default inference settings.
 322 All videos were first standardized to 30 fps and resized by
 323 setting the shorter image edge to 300 px while preserving
 324 aspect ratio to reduce computation and ease data sharing.
 325 We then computed dense forward flow between consecutive
 326 frame pairs only ($t \rightarrow t+1$).

327 3.3. Final Dataset Statistics

328 **Per-video class balance.** Table 2 reports per-trick video
 329 counts with at least one end-pose. Upright tricks (Layout,
 330 Pin-Up) are more frequent than inverted ones, reflecting
 331 their lower difficulty. In total, 812/836 clips include one
 332 end-pose, 9 of which contain a second attempt, and 24/836
 333 are failed attempts.

334 **Phase distribution.** Regarding the per-frame phase distri-
 335 bution, phase `on_pole` immediately stands out, as it holds
 336 the majority share of 60.0% of all frames. It is followed
 337 by `floor` with 20.5%, while the trick-specific labels con-
 338 tribute only a small fraction individually, with around 2-4%
 339 and a combined value of 19.5%. By design, each record-
 340 ing typically shows a single target trick but always includes
 341 background, idle, and transition segments. As a result,
 342 `on_pole` and `floor` dominate frame counts.

343 **Demographics.** To reflect real studio populations while
 344 maximizing variability from failed to perfect attempts,
 345 we prioritized novice-friendly, beginner-intermediate tricks
 346 and encouraged multiple tries. The dataset includes $N = 58$
 347 participants, deliberately emphasizing limited prior experi-
 348 ence to diversify entries, transitions, and final shapes: 0=to-
 349 tal beginner, 1=intermediate, 2=advanced. Gender distri-
 350 bution mirrors typical class composition: Female $n=43$
 351 (74.1%), Male $n=15$ (25.9%), and ages span 20–58 years
 352 (median 28). We provide the experience labels as part of
 353 our dataset.

354 4. Trick Recognition & Analysis

355 Our technical contributions demonstrate the value of the
 356 dataset while exposing the remaining challenges of privacy-
 357 preserving pole-trick understanding. Specifically, we
 358 benchmark skeleton-based and optical-flow-based action
 359 recognition under both per-video and per-frame labeling.
 360 As a unifying design choice, all action recognition models





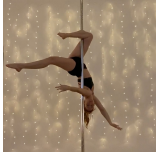

Trick	Description	Category / Level
 Layout	<ol style="list-style-type: none"> 1. Stand on the right side of the pole 2. Pull up, cross legs at ankles 3. Lean back, arch, push hips up 	Upright Beginner
 Pin-Up	<ol style="list-style-type: none"> 1. Stand on the right side of the pole 2. Pull up, right toes to left knee 3. Lean slightly back and arch 	Upright Beginner
 Wrist Seat	<ol style="list-style-type: none"> 1. Stand on the right side of the pole 2. Pull up, right toes to left knee 3. Place left hand underneath the thigh 4. Lean back, open legs into a V shape 	Upright Beginner
 Straddle Invert	<ol style="list-style-type: none"> 1. Stand behind the pole, facing left 2. Stronghold grip, pull up, lean back 3. Open legs into a V shape 	Inverted Intermediate
 Gemini	<ol style="list-style-type: none"> 1. Start with Straddle Invert 2. Hook the outside leg, other leg down 3. Chest up, arch and release hands 	Inverted Intermediate
 Inverted Crucifix	<ol style="list-style-type: none"> 1. Start with Straddle Invert 2. Place legs into crucifix hold 3. Upper body low, release arms 	Inverted Intermediate

Table 3. Compact summary of the selected tricks; terminology aligned with IPSF and Spin City [? ?].

use a lightweight GRU-based architecture to provide a con- 361
 sistent baseline across modalities and to emphasize dataset 362
 difficulty rather than architecture-specific gains. All mod- 363
 els use the same fixed train/validation/test split (70/15/15) 364
 to ensure fair comparison across modalities and supervision 365
 settings. The coaching application enabled by this pipeline 366

367 is evaluated separately in Section 5.

368 4.1. Skeleton-Based Trick Recognition

369 **Data preprocessing.** For skeleton-based recognition,
370 each clip is represented as a variable-length sequence of
371 frames, each consisting of 75 joints with (x, y, conf) values.
372 Missing detections at sequence boundaries are trimmed,
373 short internal missing frames are linearly interpolated, and
374 a light temporal smoothing is applied to reduce jitter.
375 Training-only augmentation includes mild spatial perturba-
376 tions (rotations and Gaussian noise) and temporal warping.
377 For the per-frame setting, the same preprocessing is used,
378 while frame-wise labels are aligned to the processed se-
379 quence, and unlabeled or padded positions are masked out
380 during training. To mitigate the strong class imbalance in
381 frame-wise labels, we additionally apply class-aware aug-
382 mentation during training, with stronger augmentation for
383 underrepresented classes.

384 **Model architecture.** We use a lightweight 2-layer bidi-
385 rectional GRU backbone (hidden size 64, dropout 0.2) for
386 both skeleton-based settings, taking the flattened per-frame
387 joint representation as input. For video-level classification,
388 the final forward and backward hidden states are concate-
389 nated and passed to a linear classifier. For per-frame pre-
390 diction, the recurrent output at each timestep is decoded by
391 a shared linear head to produce frame-wise logits. Models
392 are trained with cross-entropy, using masking for padded or
393 unlabeled frames in the per-frame setting.

394 4.2. Optical Flow-Based Trick Recognition

395 **Data preprocessing.** Optical-flow recognition operates
396 on precomputed dense flow fields, represented as variable-
397 length sequences of two-channel motion frames. Flow se-
398 quences are temporally subsampled (take 16, skip 5), re-
399 sized to a common processing resolution (maximum size
400 192), normalized using statistics computed on the training
401 split, and augmented only during training. For the per-frame
402 setting, frame-wise labels are assigned from the annota-
403 tions, and padded or unlabeled positions are masked out.

404 **Model architecture.** For optical flow, each frame is first
405 encoded by a 2-layer compact CNN and then aggregated
406 temporally with a 2-layer GRU. In the video-level setting,
407 the final GRU state is used for clip classification. In the per-
408 frame setting, the GRU output at each timestep is decoded
409 by a shared linear layer to obtain frame-wise predictions.
410 This design keeps the temporal backbone comparable to the
411 skeleton-based models while introducing only a modality-
412 specific spatial encoder.

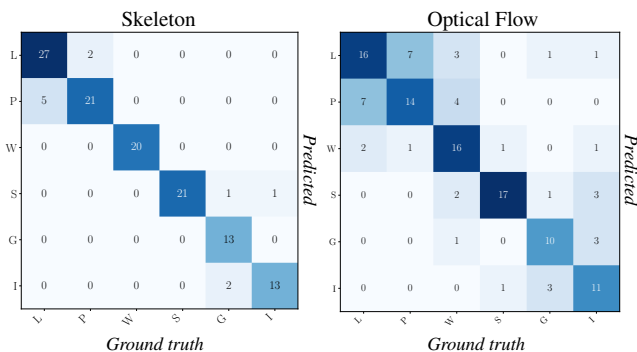


Figure 4. Confusion matrices for our per-video models on both modalities.

4.3. Results and Discussion

We report results for two supervision granularities: (i) *per-video* action recognition and (ii) *per-frame* state/trick recognition. For each task, we evaluate the two privacy-preserving modalities introduced in Pole-Arina: skeleton sequences and dense optical flow. Figures 4 and 5 show the corresponding confusion matrices, while Table 5 provides an overview of our results.

Per-video recognition (clip labels). On the shared test split, the skeleton-based video model achieves **91.26%** accuracy, while the optical-flow model reaches **66.67%** accuracy (Fig. 4). This gap is expected: skeleton extraction already provides a structured representation of body pose and motion, supplying a strong inductive bias for classification. In contrast, optical flow is a weaker and noisier proxy for articulation, particularly in unconstrained recordings with camera motion, background clutter, mirrors, and self-occlusions, leaving greater headroom for improved modeling on the flow stream.

Across both modalities, the dominant confusions occur between Layout and Pin-Up, which are visually and kinematically similar and thus represent a natural ambiguity in the label space. Additionally, errors tend to remain within the same broad regime: upright tricks are more commonly confused with other upright tricks, and inverted tricks with other inverted tricks. Mixing upright with inverted does not happen for the skeleton model, indicating that the pose representation strongly disambiguates inversion state; it is less rare (though still uncommon) for optical flow, consistent with flow providing weaker explicit cues about body orientation and inversion.

Per-frame recognition (frame labels). For frame-wise classification, the skeleton model achieves **86.99%** accuracy with a **77.71%** macro F1-score, whereas optical flow reaches **81.23%** accuracy but only **45.84%** macro F1 (Fig. 5). This discrepancy confirms that overall frame accu-

Modality	floor	on_pole	L	W	P	S	I	G
Skeleton	87.98	90.35	77.19	81.54	76.43	76.07	86.24	45.89
Optical Flow	88.08	88.03	52.33	68.22	7.99	62.02	0.00	0.00

Table 4. Per-class F1-scores (%) for per-frame recognition, highlighting the stronger effect of class imbalance on optical flow than on skeleton-based recognition.

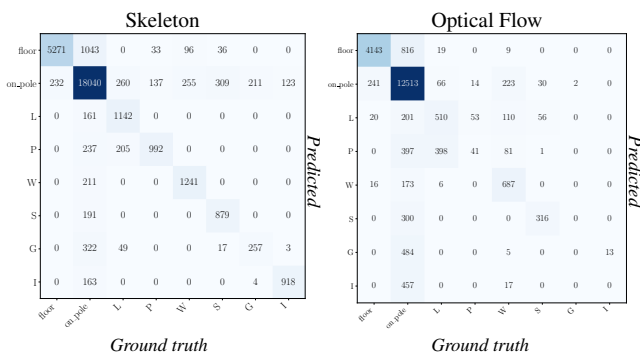


Figure 5. Confusion matrices for our per-frame models on both modalities.

449 racy is strongly inflated by class imbalance: the dominant
450 on_pole and floor states account for most frames and
451 are recognized comparatively well in both modalities, while
452 minority trick classes remain substantially harder.

453 The class-wise skeleton results show that several trick
454 classes remain meaningfully distinguishable despite this
455 imbalance, including I (86.24% F1), W (81.54% F1), L
456 (77.19% F1), P (76.43% F1), and S (76.07% F1). The main
457 exception is G (gemini), which is substantially harder at
458 only 45.89% F1. In contrast, the optical-flow model per-
459 forms reasonably on the dominant states (floor: 88.08%
460 F1, on_pole: 88.03% F1) and on some trick classes such
461 as W (68.22% F1) and S (62.02% F1), but fails on G and
462 I, both of which obtain 0.00% F1. This indicates that flow
463 alone is much less reliable for fine-grained pose discrimina-
464 tion under severe imbalance, as reflected in Table 4.

465 Overall, skeletons provide a strong and compact repre-
466 sentation for frame-level recognition, with substantially bet-
467 ter class-sensitive performance across most trick classes.
468 Optical flow remains useful for dominant motion states,
469 but currently underperforms on fine-grained minority poses.
470 These results support the dual-modality design of Pole-
471 Arina: skeletons offer a strong baseline for privacy-
472 preserving pose recognition, while optical flow defines a
473 more challenging appearance-suppressed setting with clear
474 room for methodological improvement.

475 4.4. Limitations

476 We intentionally adopt compact GRU-based models to sup-
477 port low-latency coaching use, and therefore present them
478 as practical baselines rather than capacity-saturating archi-
479 tectures. Skeleton-based results are also *conditioned on*

Modality	Task	Acc. (%)	Macro F1 (%)
Skeleton	Per-video	91.26	91.50
Optical Flow	Per-video	66.67	67.11
Skeleton	Per-frame	86.99	77.71
Optical Flow	Per-frame	81.23	45.84

Table 5. Summary of benchmark results across modalities and supervision settings. For the imbalanced per-frame task, macro F1 is more informative than overall accuracy alone.

pose estimation: extracting skeletons already introduces a
strong inductive bias, and errors or biases of the selected es-
timator, especially under occlusion and sustained pole con-
tact, directly affect both recognition and geometric scoring.
While optical flow reduces appearance leakage, it remains
sensitive to camera motion and mirror artifacts. These de-
pendencies are part of the challenge captured by Pole-Arina
and point to clear directions for future work in more robust
privacy-preserving sports analysis.

5. Coaching Feedback Case Study and User Study

Beyond benchmark recognition, we use Pole-Arina to study
a practical downstream application: coaching-oriented
feedback for static pole tricks. Our goal is not to replace ex-
pert instruction, but to examine whether privacy-preserving
motion representations and interpretable pose analysis can
support a useful feedback workflow. To this end, we first
define a transparent geometric scoring system for end-pose
assessment and then evaluate its integration into a prototype
coaching application through a user study.

5.1. Geometric Scoring System

We design an interpretable scoring module for skeleton-
based end-pose assessment. We considered an autoencoder-
based approach that would learn an *ideal* pose and score
deviations through reconstruction error [?], but such a for-
mulation would require many near-perfect exemplars and
would be highly sensitive to dataset bias. Instead, we adopt
a geometric, rule-based scorer: it is transparent, aligns with
instructor practice through angles, alignments, and relative
joint positions, and avoids the subjectivity often discussed
in dance evaluation [?].

Each trick is scored by 5–7 geometric rules that capture
its defining cues: (i) *body orientation* with respect to the
pole, (ii) *limb alignment* through joint angles, and (iii) *joint
proximity* through characteristic distances, such as ankles
being close together or toe–knee contact.

Each rule specifies a target value and tolerance for the
required joint angle or distance. A rule is evaluated only
if the corresponding joint confidences satisfy $\min(v_i) \geq \tau$,

519 with $\tau = 0.75$. We then compute the relevant angles and
520 distances from the joint coordinates and assign a *pass* or *fail*
521 decision according to the target and tolerance. The end-pose
522 score is defined as the fraction of passed evaluable rules,

$$523 \text{ score} = \frac{\# \text{passed}}{\# \text{evaluable}}.$$

524 For interpretability, the system reports per-rule feedback to-
525 gether with the measured values and target ranges. We will
526 provide the complete rule catalog for all six tricks, includ-
527 ing targets and tolerances, as part of the benchmark.

528 5.2. User Study

529 To assess the practical utility of our pipeline, we integrated
530 the per-frame skeleton-based recognition model and geo-
531 metric scorer into an end-to-end coaching-guide web ap-
532 plication. Given a recorded training video, the system ex-
533 tracts privacy-preserving pose skeletons, predicts the trick
534 label, and computes a geometric end-pose score, return-
535 ing structured feedback. Computation is centralized on a
536 FastAPI+PyTorch server, while a React/Next.js client runs
537 on consumer laptops and supports a coach-like workflow:
538 upload \rightarrow trick gallery \rightarrow per-rule feedback (Figure 6).

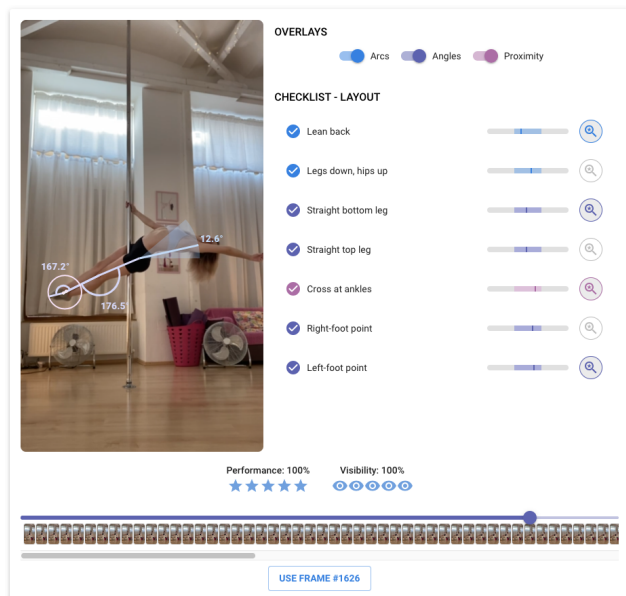


Figure 6. Application interface showing the frame viewer, overlay controls, and per-rule feedback for the predicted trick.

539 We conducted a between-subjects study with 33 partic-
540 ipants comparing a *Control* condition (smartphone video
541 self-review) against an *Experimental* condition (Pole-Arına
542 feedback). Participants were randomly assigned to one
543 condition, performed five practice trials with review after
544 each trial using experience-matched tricks (*Layout* and
545 *Pin-Up*), and then completed post-session Likert ratings

and the System Usability Scale (SUS) [?], together with
optional qualitative feedback and debriefing.

Results. Normality tests (Kolmogorov–Smirnov and
Shapiro–Wilk) indicated non-normal distributions, so we
used Mann–Whitney U tests for between-condition com-
parisons. In our between-subjects study ($N = 33$; 17 Ex-
perimental, 16 Control), the Pole-Arına condition yielded
higher ratings for trust/adoption, clarity, and usability, in-
cluding a higher SUS score (95.44 ± 4.07 vs. 86.41 ± 11.14 ;
 $U = 72.00, p = .020$), while efficiency-related mea-
sures showed no significant difference ($p = .485$). Partic-
ipants’ self-ratings generally trended below the system
scores. Qualitative feedback highlighted the usefulness of
the overlays and expressed interest in future additions such
as reference poses and real-time feedback.

These results suggest that interpretable geometry-based
feedback can be useful in practice as a downstream ap-
plication of our benchmark. While this user study does
not constitute a benchmark evaluation in itself, it indicates
that combining accurate skeleton-based recognition with
transparent rule-based explanations can support a coaching-
oriented workflow.

6. Conclusion

We introduce the first privacy-preserving benchmark for
fine-grained pole dance analysis, an underrepresented yet
increasingly popular sport that poses distinctive challenges
(contact constraints, self-occlusions, mirrors, and appear-
ance obfuscation). We provide skeleton and optical-flow
modalities with expert annotations, enabling learning with-
out releasing identifiable RGB. Strong baseline results from
our lightweight skeleton- and flow-based models demon-
strate that meaningful recognition is possible from these
representations, while leaving substantial headroom for im-
proved temporal modeling, robustness, and feedback qual-
ity. Finally, we show how these components can be in-
tegrated into a coaching-oriented prototype and evaluated
with a user study, connecting benchmark performance to
practical feedback. Together, the dataset, baselines, and
evaluation establish pole dance as a new, well-defined com-
puter vision problem setting for privacy-aware, fine-grained
athletic technique understanding.

References

- [] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. PoseTrack: A benchmark for Human Pose Estimation and Tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5167–5176, 2018.
- [] Valentin Bazarevsky, Ivan Grishchenko, Karthik Raveendran, Tyler Zhu, Fan Zhang, and Matthias Grundmann.

- BlazePose: On-device Real-time Body Pose tracking. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [] John Brooke. SUS-A quick and dirty usability scale. *Usability Evaluation in Industry*, 189(194):4–7, 1996.
- [] Yuanhao Cai, Zhicheng Wang, Zhengxiong Luo, Binyi Yin, Angang Du, Haoqian Wang, Xiangyu Zhang, Xinyu Zhou, Erjin Zhou, and Jian Sun. Learning Delicate Local Representations for Multi-person Pose Estimation. In *Computer Vision – ECCV 2020*, pages 455–472, Cham, 2020. Springer International Publishing.
- [] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Openpose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):172–186, 2019.
- [] Joao Carreira and Andrew Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [] Spin City. *The Ultimate Pole Bible*. Spin City Aerial Fitness Ltd, 2025.
- [] Steffi L Colyer, Murray Evans, Darren P Cosker, and Aki IT Salo. A Review of the Evolution of Vision-Based Motion Analysis and the Integration of Advanced Computer Vision Methods Towards Developing a Markerless System. *Sports Medicine - Open*, 4:1–15, 2018.
- [] MIMPose Contributors. OpenMMLab Pose Estimation Toolbox and Benchmark. <https://github.com/open-mmlab/mmpose>, 2020.
- [] Adrien Delière, Anthony Cioppa, Silvio Giancola, Meisam J. Seikavandi, Jacob V. Dueholm, Kamal Nasrollahi, Bernard Ghanem, Thomas B. Moeslund, and Marc Van Droogenbroeck. SoccerNet-v2: A Dataset and Benchmarks for Holistic Understanding of Broadcast Soccer Videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2021.
- [] Zhao Du, Shan Wang, Ziyan Deng, and Fang Wang. Unveiling the Power of AI Fitness Apps: A Uses and Gratifications Perspective. *Journal of Global Information Management (JGIM)*, 33(1):1–28, 2025.
- [] Aysu Ezen-Can. A Comparison of LSTM and BERT for Small Corpus. *arXiv preprint arXiv:2009.05451*, 2020.
- [] International Pole Sports Federation. Code of Points 2025 – 2027. https://ipsfsports.org/downloads/Uncategorised/ipsf_pole_sports_code_of_points_2025-2027_final_070120240.pdf, 2025. Accessed: 2025-08-16.
- [] Indrajeet Ghosh, Sreenivasan Ramasamy Ramamurthy, Avijoy Chakma, and Nirmalya Roy. Sports analytics review: Artificial intelligence applications, emerging technologies, and algorithmic perspective. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 13(5):e1496, 2023.
- [] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [] Filip Ilic, Thomas Pock, and Richard P. Wildes. Is Appearance Free Action Recognition Possible? In *European Conference on Computer Vision*, 2022.
- 596[] Filip Ilic, He Zhao, Thomas Pock, and Richard P. Wildes. 654
597 Selective, Interpretable and Motion Consistent Privacy At- 655
598 tribute Obfuscation for Action Recognition. In *Proceedings 656*
599 *of the IEEE Conference on Computer Vision and Pattern 657*
600 *Recognition*, 2024. 658
- 601[] Sheng Jin, Lumin Xu, Jin Xu, Can Wang, Wentao Liu, Chen 659
602 Qian, Wanli Ouyang, and Ping Luo. Whole-Body Human 660
603 Pose Estimation in the Wild. In *Computer Vision – ECCV 661*
604 *2020*, pages 196–214, Cham, 2020. Springer International 662
605 Publishing. 663
- 606[] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas 664
607 Leung, Rahul Sukthankar, and Li Fei-Fei. Large-Scale Video 665
608 Classification with Convolutional Neural Networks. In *Pro- 666*
609 *ceedings of the IEEE Conference on Computer Vision and 667*
610 *Pattern Recognition*, 2014. 668
- 611[] Yeonho Kim and Daijin Kim. Real-time dance evaluation by 669
612 markerless human pose estimation. *Multimedia Tools and 670*
613 *Applications*, 77:31199–31220, 2018. 671
- 614[] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. PifPaf: 672
615 Composite Fields for Human Pose Estimation. In *Proce- 673*
616 *edings of the IEEE Conference on Computer Vision and Pattern 674*
617 *Recognition*, 2019. 675
- 618[] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Open- 676
619 PifPaf: Composite Fields for Semantic Keypoint Detection 677
620 and Spatio-Temporal Association. *IEEE Transactions on In- 678*
621 *telligent Transportation Systems*, pages 1–14, 2021. 679
- 622[] Yingwei Li, Yi Li, and Nuno Vasconcelos. RESOUND: To- 680
623 wards Action Recognition without Representation Bias. In 681
624 *Proceedings of the European Conference on Computer Vi- 682*
625 *sion*, pages 513–528, 2018. 683
- 626[] Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris Mc- 684
627 Clanahan, Esha Uboweja, Michael Hays, Fan Zhang, Chuo- 685
628 Ling Chang, Ming Yong, Juhyun Lee, Wan-Teh Chang, Wei 686
629 Hua, Manfred Georg, and Matthias Grundmann. Mediapipe: 687
630 A framework for perceiving and processing reality. In *Third 688*
631 *Workshop on Computer Vision for AR/VR at IEEE Computer 689*
632 *Vision and Pattern Recognition*, 2019. 690
- 633[] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. 691
634 ShuffleNet V2: Practical Guidelines for Efficient CNN Ar- 692
635 chitecture Design, 2018. 693
- 636[] Marina Mikami and Noriyuki Kida. Categorizing Rhyth- 694
637 mic Jumping Motion Using Motion Capture without Mark- 695
638 ers. *Advances in Physical Education*, 13(2):93–105, 2023. 696
- 639[] Hoang Nguyen, Ankhzaya Jamsrandorj, Vanyi Chao, 697
640 Yin May Oo, Muhammad Amrulloh Robbani, Kyung-Ryoul 698
641 Mun, and Jinwook Kim. VNL-STES: A Benchmark Dataset 699
642 and Model for Spatiotemporal Event Spotting in Volleyball 700
643 Analytics. In *Proceedings of the Computer Vision and Pat- 701*
644 *tern Recognition Conference*, pages 5862–5871, 2025. 702
- 645[] Paritosh Parmar and Brendan Tran Morris. Learning to Score 703
646 Olympic Events. In *Proceedings of the IEEE Conference on 704*
647 *Computer Vision and Pattern Recognition Workshops*, pages 705
648 20–28, 2017. 706
- 649[] Zhiqiang Pu, Yi Pan, Shijie Wang, Boyin Liu, Min Chen, 707
650 Hao Ma, and Yixiong Cui. Orientation and Decision-Making 708
651 for Soccer Based on Sports Analytics and AI: A Systematic 709
652 Review. *IEEE/CAA Journal of Automatica Sinica*, 11(1):37– 710
653 57, 2024. 711

- [] Jiping Qu. A dance movement quality evaluation model using transformer encoder and convolutional neural network. *Scientific Reports*, 14(1):32058, 2024. 712
713
714
- [] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. FineGym: A Hierarchical Video Dataset for Fine-grained Action Understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 715
716
717
718
- [] Karen Simonyan and Andrew Zisserman. Two-Stream Convolutional Networks for Action Recognition in Videos. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*, page 568–576, Cambridge, MA, USA, 2014. MIT Press. 719
720
721
722
723
- [] Xiang Suo, Weidi Tang, and Zhen Li. Motion Capture Technology in Sports Scenarios: A Survey. *Sensors*, 24(9):2947, 2024. 724
725
726
- [] Zachary Teed and Jia Deng. RAFT: Recurrent All-Pairs Field Transforms for Optical Flow. In *European Conference on Computer Vision*, pages 402–419. Springer, 2020. 727
728
729
- [] Atima Tharatipyakul, Thanawat Srikaewsiew, and Suporn Pongnumkul. Deep learning-based pose estimation in providing feedback for physical movement: A review. *Heliyon*, 10(17):e36589, 2024. 730
731
732
733
- [] Alexander Toshev and Christian Szegedy. DeepPose: Human Pose Estimation via Deep Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1653–1660, 2014. 734
735
736
737
- [] Michael Tschannen, Olivier Bachem, and Mario Lucic. Recent Advances in Autoencoder-Based Representation Learning. *arXiv preprint arXiv:1812.05069*, 2018. 738
739
740
- [] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All You Need. *Advances in Neural Information Processing Systems*, 30, 2017. 741
742
743
744
- [] Charlene Weaving. Sliding Up and Down a Golden Glory Pole: Pole Dancing and the Olympic Games. *Sport, Ethics and Philosophy*, 14(4):525–536, 2020. 745
746
747
- [] Weiyu Zhang, Menglong Zhu, and Konstantinos G Derpanis. From Actemes to Action: A Strongly-supervised Representation for Detailed Action Understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2248–2255, 2013. 748
749
750
751
752