

# Matching Every Pair to Track Every Point: PairFormer for All-Pairs Tracking and Video Trajectory Fields

Anonymous CVPR submission

## Abstract

001 *Tracking-any-point (TAP) answers query-conditioned cor-*  
002 *respondence but leaves the dense, all-pairs structure of a*  
003 *video implicit. We formulate All-Pairs Tracking (APT):*  
004 *given a video, predict dense displacement and visibility for*  
005 *every source-target frame pair, from which per-pixel trajec-*  
006 *tories can be read out. To this end, we propose PairFormer,*  
007 *a feed-forward transformer that addresses APT in a single*  
008 *pass. A spatio-temporal patch encoder computes tempo-*  
009 *rally conditioned features for all frames. CorrBank con-*  
010 *structs a learnable correlation memory for each frame pair*  
011 *to obtain pairwise motion tokens. A broadcast motion mixer*  
012 *aggregates information across space and time and refines*  
013 *these tokens with global context. A trajectory head then*  
014 *predicts full-resolution displacement for each pair, yield-*  
015 *ing a coherent all-pairs trajectory field. To support APT*  
016 *at scale, we develop PAIRender, a data platform that syn-*  
017 *thesizes photo-realistic dynamic scenes with dense annota-*  
018 *tions. From PAIRender we derive a training set ( $\pi$ -R10K)*  
019 *and a benchmark (APT-Bench) with an all-to-all evalua-*  
020 *tion protocol. Experiments show that PairFormer achieves*  
021 *strong performance on APT-Bench and competitive results*  
022 *on standard TAP benchmarks. Code and dataset will be re-*  
023 *leased upon publication.*

## 024 1. Introduction

025 Understanding motion in video is more than linking a  
026 few queried points [12, 13, 18, 22] or estimating pair-  
027 wise optical flow [20, 27, 37, 40]; it requires a represen-  
028 tation that exposes how *every* image location relates to *ev-*  
029 *ery* frame across the sequence. TAP delivers long-range,  
030 occlusion-aware but *sparse*, query-conditioned tracks [12,  
031 22], whereas optical flow focuses on adjacent-frame corre-  
032 spondences [20, 27, 37, 40]; in contrast, APT requires an  
033 explicit *all-pairs* dense field, enabling uniform reasoning  
034 and temporal consistency across the entire sequence.

035 We advocate *All-Pairs Tracking (APT)*: given a video,

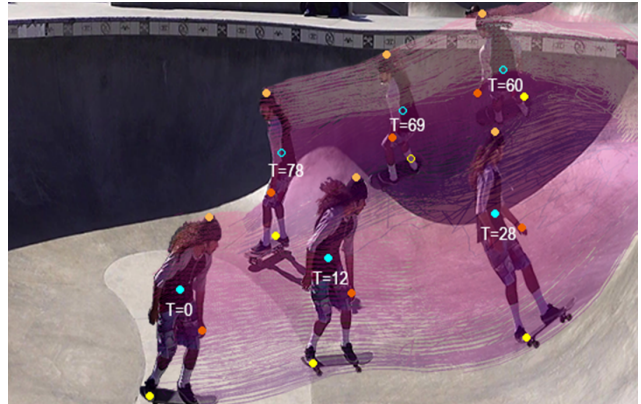


Figure 1. **All-Pairs Tracking (APT) illustration.** APT estimates a dense trajectory field for all pixels; here we visualize trajectories of pixels on the person region over the entire sequence, where hollow markers indicate occluded points.

036 predict dense correspondences for *every* ordered frame  
037 pair, together with a visibility probability that marks trace-  
038 able mappings. APT strictly generalizes two-frame optical  
039 flow [20, 27, 37, 40] and subsumes TAP-style point track-  
040 ing [12, 22]. From an explicit all-pairs field, trajectories for  
041 any pixel can be read out on demand rather than obtained  
042 by propagation from a single start frame. Framed this way,  
043 matching at any temporal offset is handled uniformly, and  
044 sequence-wide temporal consistency becomes a first-class  
045 modeling objective.

046 Realizing APT poses several challenges. First, the  
047 number of frame pairs grows quadratically with sequence  
048 length. Second, the relevant context for each correspon-  
049 dence is often spread across distant frames and spatial lo-  
050 cations. Third, repeated textures, large motions, and oc-  
051 clusions introduce substantial ambiguity. Existing hand-  
052 crafted cost volumes and local propagation schemes [5, 28,  
053 34, 37] struggle to aggregate such non-local context and to  
054 maintain consistency at the scale of an entire sequence.

055 To address these challenges, we introduce *PairFormer*, a  
056 transformer architecture guided by a simple principle: *first*  
057 *construct pairwise correspondences, then enforce global*

058	<i>consistency across the sequence.</i> The model comprises four	109
059	components. A <b>Spatio-Temporal Patch Encoder</b> first pro-	110
060	duces temporally conditioned patch representations for all	111
061	frames. <b>CorrBank</b> then converts each ordered frame pair	112
062	into a learnable correlation memory and performs memory-	113
063	augmented matching to yield <i>pairwise motion tokens</i> . Next,	114
064	a <b>Broadcast Motion Mixer</b> aggregates information along	115
065	trajectories (per-trajectory), within frames (per-frame), and	116
066	across the sequence (all-frames), and at selected depths	117
067	<i>broadcasts</i> the aggregated context back to the pairwise to-	118
068	kens to resolve ambiguities, including occlusion and re-	119
069	appearance. Finally, a <b>Trajectory Field Decoder</b> maps the	120
070	refined tokens to per-pair displacement and visibility, yield-	121
071	ing an explicit all-pairs trajectory field.	122
072	To support large-scale training and evaluation, we de-	123
073	velop a Blender-based platform ( <i>PAIRender</i> ) featuring di-	124
074	verse environments, moving characters, and natural cam-	125
075	era trajectories, in the spirit of recent synthetic data gen-	126
076	erators for correspondence tasks [17, 43]. The platform	127
077	generates synchronized RGB frames, per-pixel 2D/3D tra-	128
078	jectories, depth, semantic labels, visibility, and camera pa-	129
079	rameters, enabling dense supervision and controlled exper-	130
080	imentation. From PAIRender we derive a training dataset	131
081	( $\pi$ - <i>R10K</i> ) and a held-out benchmark ( <i>APT-Bench</i> ) with an	132
082	<i>all-to-all</i> evaluation protocol that queries points sampled	133
083	from every frame, encouraging sequence-wide reasoning	134
084	rather than first-frame propagation as in TAP-style bench-	135
085	marks [12].	136
086	In summary, our contributions are:	137
087	• We introduce <b>APT</b> as an explicit, dense, visibility-aware	138
088	correspondence formulation that generalizes optical flow	139
089	and TAP while targeting sequence-wide consistency.	140
090	• We propose <b>PairFormer</b> , a feed-forward transformer ar-	141
091	chitecture that constructs a globally consistent all-pairs	142
092	trajectory field in a single forward pass.	143
093	• We develop <b>PAIRender</b> and release $\pi$ - <i>R10K</i> and <i>APT-</i>	144
094	<i>Bench</i> , providing dense all-pairs supervision and an all-	145
095	to-all evaluation protocol tailored to APT.	146
096	<b>2. Related Work</b>	147
097	<b>Optical flow.</b> The notion of optic flow in perception dates	148
098	back to Gibson’s work [16], and we refer to Niehorster et	149
099	al. [30] for a historical overview. In computer vision, op-	150
100	tical flow typically denotes the dense 2D motion field re-	151
101	lating two consecutive video frames, mapping each pixel in	152
102	the first frame to a location in the second [20, 27]. Most	153
103	classical and modern approaches share a similar optimiza-	154
104	tion pattern: an initial flow estimate is obtained (for ex-	155
105	ample, assuming zero motion or a simple prior), a match-	156
106	ing cost is computed from appearance features, this cost	157
107	is regularized by smoothness assumptions or learned pri-	158
108	ors, and the flow is iteratively refined in a local neighbor-	159
	hood [4, 21, 34, 37, 42]. Recent learning-based methods	160
	such as RAFT and SEA-RAFT push this paradigm fur-	161
	ther [37, 40]. They construct a correlation volume between	
	feature maps extracted by a deep backbone [19], maintain a	
	low-resolution flow field that is updated recurrently by con-	
	volutional networks using both flow and correlations, and	
	finally upsample the prediction to full resolution using spe-	
	cialized upsampling modules [15, 21, 33]. Our formulation	
	of All-Pairs Tracking (APT) is related in spirit—dense cor-	
	respondences remain central—but differs in two key ways:	
	(i) we target <i>all</i> ordered frame pairs rather than only adja-	
	cent ones, and (ii) we replace handcrafted correlation vol-	
	umes and recurrent conv refinements with a transformer-	
	based architecture that builds a learnable correlation mem-	
	ory (CorrBank) and a global motion field over the entire	
	sequence.	
	<b>Flow-based point tracking.</b> Dense optical flow is often	
	used as a primitive for constructing longer trajectories. A	
	common strategy is to link per-frame flow estimates over	
	time, forming tracks by chaining displacement vectors and	
	detecting occlusions to prevent drift [32, 35]. More recent	
	systems extend this idea by using multi-step flows [7–9, 28]	
	or by incorporating temporal priors learned from data [5],	
	allowing trajectories to bridge over occlusions instead of	
	terminating. These approaches can yield dense multi-frame	
	tracks, but they rely on a separate optical flow model as a	
	pre-processing stage and treat tracking as a post-hoc opera-	
	tion on its outputs.	
	<b>Point trackers without flow.</b> A second line of work tracks	
	points directly, without explicitly predicting dense optical	
	flow. Classical methods such as the Lucas–Kanade fam-	
	ily [27, 38] track sparse features using local image align-	
	ment. More recently, Harley et al. [18] proposed a deep	
	point tracker with a sliding multi-frame window that al-	
	lows tracks to pass through occlusions, and Doersch et	
	al. [12] introduced the “Tracking Any Point” (TAP) bench-	
	mark, which catalyzed a wave of follow-up methods. Sub-	
	sequent work has explored multi-point interactions [22,	
	23], improved initialization and re-initialization [13], ex-	
	panded correlation neighborhoods [2, 6], transformer-based	
	designs [25, 26], and post-hoc densification from sparse	
	tracks [24]. Among these, CoTracker3 [22] represents a	
	strong state of the art for TAP-style tracking. However,	
	CoTracker-type models operate on a <i>sparse</i> set of user-	
	chosen queries, and performance depends on how these	
	queries are distributed [23]. Memory and compute con-	
	straints limit the number of simultaneously tracked points	
	to a few thousand.	
	<b>Training dataset.</b> For point tracking, common training	
	sources include Kubric [17], FlyingThings++ [18], and	
	PointOdyssey [43]. Performance can be sensitive to how	
	these datasets are combined and scheduled [40], and differ-	
	ent works adopt different curricula.	

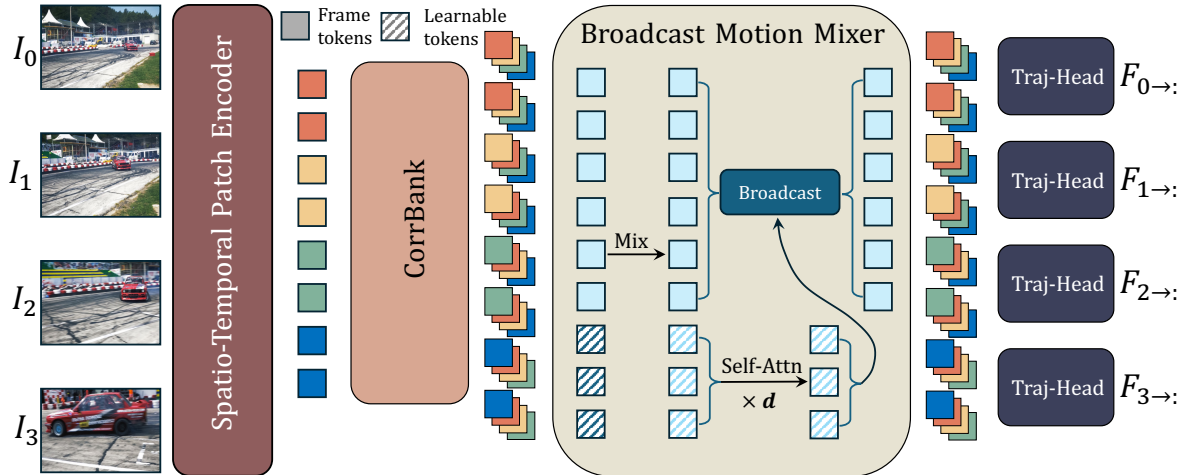


Figure 2. **PairFormer** architecture for **All-Pairs Tracking (APT)**. Given a video, the **Spatio-Temporal Patch Encoder** maps all frames to temporally conditioned patch features on a regular grid. **CorrBank** performs memory-augmented matching to produce *pairwise motion tokens*. The **Broadcast Motion Mixer** applies trajectory-wise mixing, followed by per-frame and all-frames context aggregation and a broadcast step that conditions each motion token on global context. Finally, the **Trajectory Field Decoder** (Traj-Head) applies a DPT-style head to each pair and outputs full-resolution displacement.  $F_{s \rightarrow \cdot}$  denotes displacement field from frame  $s$  to all other frames.

### 3. Method

#### 3.1. Problem Formulation

Let  $\mathcal{V} = \{I_t\}_{t=1}^T$  be a video, where each frame  $I_t$  is defined on an image domain  $\Omega \subset \mathbb{R}^2$ . *All-Pairs Tracking (APT)* aims to estimate, for every ordered frame pair  $(s, t)$ , a dense forward correspondence from frame  $s$  to frame  $t$ .

We represent the correspondence from a pixel  $x \in \Omega$  in frame  $s$  to its location in frame  $t$  by

$$\phi_{s \rightarrow t}(x) = x + F_{s \rightarrow t}(x), \quad (1)$$

where  $F_{s \rightarrow t} : \Omega \rightarrow \mathbb{R}^2$  is a displacement field. To handle occlusion and out-of-view cases, we further predict a *visibility map*

$$V_{s \rightarrow t} : \Omega \rightarrow [0, 1], \quad (2)$$

which we interpret as the probability that  $\phi_{s \rightarrow t}(x)$  is visible and corresponds to the same physical point in frame  $t$ . Given the family of mappings  $\{\phi_{s \rightarrow t}\}_{t=1}^T$  for a fixed source frame  $s$ , the trajectory of a pixel  $x$  is the sequence  $\{\phi_{s \rightarrow t}(x)\}_{t=1}^T$ .

#### 3.2. Network Architecture

We propose *PairFormer*, a transformer architecture tailored to APT and guided by a simple principle: *first construct pairwise correspondences, then enforce global consistency across the sequence*. As illustrated in Fig. 2, the model comprises four components. A **Spatio-Temporal Patch Encoder** (ST-Patch Encoder) maps each frame to temporally conditioned patch representations. **CorrBank** converts each ordered frame pair into a learnable *correlation memory* and performs memory-augmented matching

to produce *pairwise motion tokens*. The **Broadcast Motion Mixer** (BMM) performs trajectory-wise and spatio-temporal context mixing, and *broadcasts* the aggregated context back to refine the pairwise motion tokens. Finally, a **Trajectory Field Decoder** applies a DPT-style head [31] independently to each pair, mapping refined tokens to full-resolution displacement  $F_{s \rightarrow t}$ , visibility  $V_{s \rightarrow t}$ , and confidence  $C_{s \rightarrow t}$ , forming an explicit all-pairs trajectory field.

**ST-Patch Encoder.** Given video frames  $\{I_t\}_{t=1}^T$ , the ST-Patch Encoder splits each frame into a regular grid of patches and maps each patch to a feature vector. We then add a temporal embedding (shared across patches within a frame) and a 2D positional encoding on the patch grid, and apply a stack of self-attention blocks over all frames jointly. This yields temporally conditioned patch representations for all frames, which serve as the input to CorrBank and the subsequent motion modeling modules.

**CorrBank.** Let  $Z_s, Z_t \in \mathbb{R}^{P \times D}$  denote ST-Patch Encoder features for frames  $s$  and  $t$  on a patch grid with  $P$  locations. For each ordered pair  $(s, t)$ , we first form a residual feature map

$$R_{s,t} = h(Z_t - Z_s), \quad (3)$$

where  $h(\cdot)$  is a small convolutional network that emphasizes local motion cues. CorrBank maintains a set of learnable *correlation tokens*  $A \in \mathbb{R}^{K \times D}$ , shared across all pairs. We use cross-attention to let the residual features query these tokens and obtain a *pair-specific correlation memory* on the

217 patch grid  $M_{s,t} \in \mathbb{R}^{P \times D}$ :

218 
$$M_{s,t} = \text{Attn}(Q = R_{s,t}, K = A, V = A). \quad (4)$$

219 In a second step, we perform *memory-augmented matching*:  
 220 source features  $Z_s$  act as queries, target features  $Z_t$  provide  
 221 keys, and the correlation memory  $M_{s,t}$  supplies values to  
 222 produce *pairwise motion tokens*  $Y_{s,t} \in \mathbb{R}^{P \times D}$ :

223 
$$Y_{s,t} = \text{Attn}(Q = Z_s, K = Z_t, V = M_{s,t}). \quad (5)$$

224 Thus, CorrBank replaces an explicit 4D cost volume with a  
 225 learnable, token-based correlation module that can leverage  
 226 efficient kernels such as FlashAttention [10].

227 **Broadcast Motion Mixer (BMM).** Let  $T$  denote the  
 228 number of frames and  $P$  the number of patches per frame.  
 229 Given pairwise motion tokens  $Y_{s,t} \in \mathbb{R}^{P \times D}$  for each source  
 230 frame  $s$  and target frame  $t \in \{1, \dots, T\}$ , BMM refines  
 231 them in two stages.

232 *Trajectory-wise mixing.* For each source frame  $s$  and  
 233 patch index  $p$ , we consider the temporal sequence of motion  
 234 tokens  $\{Y_{s,\tau}(p)\}_{\tau=1}^T \in \mathbb{R}^{T \times D}$ . BMM introduces  $K$   
 235 learnable context tokens  $U \in \mathbb{R}^{K \times D}$ , shared across all tra-  
 236 jectories, and concatenates them with this sequence:

237 
$$X_{s,p} = [U; Y_{s,1}(p), \dots, Y_{s,T}(p)] \in \mathbb{R}^{(K+T) \times D}. \quad (6)$$

238 Several self-attention blocks are applied to  $X_{s,p}$ . We keep  
 239 the first  $K$  output positions as *trajectory context tokens*  
 240  $H_{s,p} \in \mathbb{R}^{K \times D}$ , and the remaining  $T$  positions as updated  
 241 motion tokens  $\tilde{Y}_{s,\tau}(p)$ . In this way, the learnable tokens  
 242  $H_{s,p}$  absorb information from all time steps and form a  
 243 compact summary of the motion along the trajectory for  
 244 patch  $(s, p)$ .

245 *Context mixing (per-frame and all-frames).* For each  
 246 source frame  $s$ , we feed the  $P$  trajectory context tokens  
 247  $\{H_{s,p}\}_{p=1}^P$  into a stack of  $d$  per-frame self-attention blocks  
 248 (with  $s$  as the batch index), yielding per-frame features  
 249  $F_s \in \mathbb{R}^{P \times D}$  that encode spatial structure and context within  
 250 frame  $s$ . We then concatenate  $\{F_s\}_{s=1}^T$  along the token di-  
 251 mension and pass the resulting  $T \times P$  spatio-temporal to-  
 252 kens through another stack of  $d$  self-attention blocks, ob-  
 253 taining scene-level features  $G \in \mathbb{R}^{T \times P \times D}$  that mix infor-  
 254 mation across both space and time.

255 *Broadcast refinement.* At selected depths, we construct,  
 256 for each  $(s, p)$ , a small context set  $S_{s,p} = [F_s(p), G_{s,p}]$ ,  
 257 where  $F_s(p)$  and  $G_{s,p}$  denote the per-frame and scene-level  
 258 features at that patch, respectively. We then refine the mo-  
 259 tion tokens at patch  $(s, p)$  via cross-attention:

260 
$$\tilde{Y}_{s,\cdot}^*(p) = \text{Attn}(Q = \tilde{Y}_{s,\cdot}(p), K = S_{s,p}, V = S_{s,p}). \quad (7)$$

261 The updated sequence  $\tilde{Y}_{s,\cdot}^*(p) \in \mathbb{R}^{T \times D}$  is finally reshaped  
 262 back to refined pairwise motion tokens  $Y_{s,t}^* \in \mathbb{R}^{P \times D}$ .

This broadcast step conditions each pairwise motion token  
 on both its own trajectory summary and the global spatio-  
 temporal context, while preserving its dependence on the  
 specific frame pair  $(s, t)$ .

**Trajectory Field Decoder.** The final stage maps the re-  
 fined pairwise motion tokens to dense per-pixel outputs.  
 Given  $Y_{s,t}^* \in \mathbb{R}^{P \times D}$  on the patch grid for each ordered  
 pair  $(s, t)$ , we apply a DPT-style decoder head [31] with  
 learned upsampling independently to each pair. This head  
 takes  $Y_{s,t}^*$  and predicts a four-channel full-resolution map,  
 comprising 2D displacement  $F_{s \rightarrow t} \in \mathbb{R}^{H_{\text{im}} \times W_{\text{im}} \times 2}$ , a vis-  
 ibility field  $V_{s \rightarrow t} \in [0, 1]^{H_{\text{im}} \times W_{\text{im}}}$ , and a confidence field  
 $C_{s \rightarrow t} \in \mathbb{R}^{H_{\text{im}} \times W_{\text{im}}}$ .

### 3.3. Training Scheme

We denote ground-truth displacement and visibility by  
 $F_{s \rightarrow t}^*$  and  $V_{s \rightarrow t}^*$  when available. For a pixel  $x \in \Omega$  in source  
 frame  $s$  and a target frame  $t$ , we define the per-pixel residual

$$r_{s \rightarrow t}(x) = \|F_{s \rightarrow t}(x) - F_{s \rightarrow t}^*(x)\|_1. \quad (8)$$

The overall training objective is a weighted sum

$$\mathcal{L} = \lambda_{\text{traj}} \mathcal{L}_{\text{traj}} + \lambda_{\text{conf}} \mathcal{L}_{\text{conf}} + \lambda_{\text{vis}} \mathcal{L}_{\text{vis}} + \lambda_{\text{corr}} \mathcal{L}_{\text{corr}}. \quad (9)$$

**Trajectory loss.** We supervise displacements with a stan-  
 dard per-pixel  $\ell_1$  loss:

$$\mathcal{L}_{\text{traj}} = \mathbb{E}_{(s,t),x} [r_{s \rightarrow t}(x)], \quad (10)$$

where the expectation is taken over all supervised frame  
 pairs and all valid pixels.

**Confidence adjustment.** PairFormer predicts an auxil-  
 iary confidence field  $C_{s \rightarrow t}(x) > 0$  alongside the displace-  
 ment  $F_{s \rightarrow t}(x)$ . We treat  $C_{s \rightarrow t}(x)$  as a learned weight on  
 the residual, regularized by a log term. For a single pixel  
 $(s, t, x)$ , the confidence-adjusted penalty is

$$\ell_{\text{conf}}(s, t, x) = r_{s \rightarrow t}(x) C_{s \rightarrow t}(x) - \alpha \log C_{s \rightarrow t}(x), \quad (11)$$

with  $\alpha = 0.3$  controlling the strength of regularization. The  
 overall confidence loss is then

$$\mathcal{L}_{\text{conf}} = \mathbb{E}_{(s,t),x} [\ell_{\text{conf}}(s, t, x)]. \quad (12)$$

The first term encourages larger confidence on pixels with  
 small residuals, while the negative log term prevents the  
 model from driving  $C_{s \rightarrow t}(x)$  to zero.



Figure 3. Sample image and dense flow field rendered from our PAIRender data platform.

300 **Visibility loss.** In addition to displacement, PairFormer  
 301 predicts a visibility field  $V_{s \rightarrow t}(x)$  that should match the  
 302 ground-truth visibility  $V_{s \rightarrow t}^*(x)$ . We supervise visibility  
 303 with a standard pixel-wise binary cross-entropy loss:

$$304 \quad \mathcal{L}_{\text{vis}} = \mathbb{E}_{(s,t),x} [\ell_{\text{BCE}}(V_{s \rightarrow t}(x), V_{s \rightarrow t}^*(x))], \quad (13)$$

305 where  $\ell_{\text{BCE}}$  denotes the binary cross-entropy between a  
 306 predicted probability and a binary label.

307 **Correspondence regularization.** Synthetic supervision  
 308 from PAIRender provides temporally consistent trajectories,  
 309 which allows us to identify pixels that correspond to  
 310 the same physical point across frames. Let  $\Omega_{\text{corr}}$  denote a  
 311 set of matched pixel pairs  $((s, x), (t, x'))$  lying on the same  
 312 ground-truth trajectory. For such pairs, the predicted all-  
 313 pairs trajectory fields starting from  $x$  and  $x'$  should agree.  
 314 We define a trajectory descriptor  $\mathbf{q}_s(x) \in \mathbb{R}^{T \times 2}$  as the stack  
 315 of forward-mapped positions

$$316 \quad \mathbf{q}_s(x) = [\phi_{s \rightarrow 1}(x), \dots, \phi_{s \rightarrow T}(x)], \quad (14)$$

317 and similarly  $\mathbf{q}_t(x')$  for frame  $t$ . The correspondence regu-  
 318 larizer is

$$319 \quad \mathcal{L}_{\text{corr}} = \mathbb{E}_{((s,x),(t,x')) \in \Omega_{\text{corr}}} [\|\mathbf{q}_s(x) - \mathbf{q}_t(x')\|_1]. \quad (15)$$

320 This term encourages points on the same trajectory to share  
 321 a consistent all-pairs prediction, beyond what is enforced by  
 322 per-pair supervision alone.

323 Together, Eqs. (10)–(15) provide complementary super-  
 324 vision:  $\mathcal{L}_{\text{traj}}$  anchors displacements to ground-truth flow,  
 325  $\mathcal{L}_{\text{conf}}$  adaptively reweights and regularizes per-pixel resid-  
 326 uals,  $\mathcal{L}_{\text{vis}}$  trains the visibility field with a probabilistic in-  
 327 terpretation, and  $\mathcal{L}_{\text{corr}}$  ties together predictions for pixels  
 328 along the same physical trajectory. Their weighted combi-  
 329 nation in Eq. (9) is used for all experiments. Empirically,  
 330 we set  $(\lambda_{\text{traj}}, \lambda_{\text{conf}}, \lambda_{\text{vis}}, \lambda_{\text{corr}}) = (1, 1, 1, 0.3)$  in all exper-  
 331 iments, which we find yields stable training.

## 4. Data Platform: PAIRender 332

333 Data-driven modeling of dynamic scenes is limited by the  
 334 lack of large-scale datasets that provide *dense*, *reliable*,  
 335 and *long-range* supervision. Existing synthetic corpora are  
 336 modest in scale and biased toward rigid or short-term mo-  
 337 tion, with sparse annotations that are insufficient for learn-  
 338 ing all-pairs correspondence at the sequence level. To ad-  
 339 dress this, we develop *PAIRender*, a scalable Blender-based  
 340 platform that generates photorealistic dynamic videos with  
 341 per-pixel ground truth tailored to APT.

342  **$\pi$ -R10K dataset.** Using PAIRender, we assemble  $\pi$ -  
 343 *R10K*, a training dataset for all-pairs trajectory field estima-  
 344 tion. The collection (Fig. 3) contains 10,000 scenes, each a  
 345 sequence of 60 frames at  $512 \times 512$  resolution. Diversity is  
 346 characterized along three aspects: (i) **Environment:** indoor  
 347 and outdoor layouts from public asset libraries and procedu-  
 348 ral generation; (ii) **Dynamics:** articulated humans and ev-  
 349 eryday objects undergoing both rigid and non-rigid motion,  
 350 with frequent occlusions; (iii) **Camera:** smooth trajectories  
 351 centered on regions with significant motion, obtained by re-  
 352 targeting real-world camera motions. Each sequence pro-  
 353 vides synchronized RGB, dense 2D/3D trajectories, depth,  
 354 visibility masks, semantic labels, and camera parameters,  
 355 enabling the supervision scheme in Sec. 3. PAIRender is  
 356 fully programmable, so scene count and motion patterns can  
 357 be easily scaled.

358 **APT-Bench benchmark.** To evaluate APT models, we  
 359 construct *APT-Bench*, a held-out test set of 100 sequences  
 360 curated from PAIRender. Each sequence contains 120  
 361 frames at  $1920 \times 1080$  resolution. Unlike established  
 362 point-tracking benchmarks that evaluate only points sam-  
 363 pled from the first frame, APT-Bench uses an *all-to-all* pro-  
 364 tocol that samples points from every frame and evaluates  
 365 trajectories across the entire sequence.

Table 1. TAP benchmark results in terms of  $\delta_{\text{avg}}$  (higher is better).

Method	Bad.	Cro.	Dav.	Dri.	Ego.	Kin.	Rgb.	Rob.	Avg.
RAFT	23.7	29.3	48.5	44.8	41.0	64.3	82.8	72.2	50.8
SEA-RAFT	23.9	21.9	48.7	49.4	44.0	64.3	85.7	67.6	50.7
AccFlow	10.3	22.2	23.5	26.4	4.0	38.8	63.2	57.9	30.8
PIPs++	34.1	27.5	62.5	51.3	38.5	64.2	70.4	73.4	52.7
LocoTrack	41.4	43.1	68.0	66.5	58.4	70.0	80.3	76.9	63.1
BootsTAPIR	42.7	34.9	67.9	66.9	56.8	70.6	81.0	78.2	62.4
DELTA	44.6	42.9	75.3	67.8	40.3	66.5	83.0	74.8	61.9
CoTracker2	40.0	31.7	70.9	67.8	43.2	65.8	73.4	73.0	58.2
CoTracker3	48.3	44.5	77.1	69.8	60.4	71.8	84.2	81.6	67.2
PairFormer	49.7	44.1	75.8	66.9	63.2	72.1	89.1	83.8	68.1

Table 2. Average Jaccard (AJ, higher is better).

Method	Dav.	Kin.	Rgb.	Rob.	Avg.
CoTracker3	62.9	53.5	69.5	66.3	63.1
PairFormer	63.1	56.1	79.0	69.9	65.9

Table 3. Occlusion Accuracy (OA, higher is better).

Method	Dav.	Kin.	Rgb.	Rob.	Avg.
CoTracker3	90.1	85.6	91.6	89.8	89.3
PairFormer	90.2	89.5	92.5	90.2	91.1

## 366 5. Experiments

### 367 5.1. Training Details

368 We train PairFormer on a 1:1 mixture of  $\pi$ -R10K and  
 369 Kubric [17] training dataset. Video clips have random se-  
 370 quence length (uniformly sampled from 30 to 60 frames) in  
 371 fixed size of  $384 \times 512$ . We adopt standard data augmen-  
 372 tations used in point tracking and optical flow, including  
 373 random cropping and scaling, horizontal flips, color jitter,  
 374 and synthetic occluders [22]. We train the model with 8  
 375 NVIDIA H100-80G GPUs for 50,000 iterations. We opti-  
 376 mize with AdamW (weight decay 0.01); the learning rate is  
 377  $2 \times 10^{-4}$  with cosine decay and 1,000-step linear warmup.

### 378 5.2. Evaluation on Benchmarks

379 **Baseline methods.** We compare our PairFormer with both  
 380 point trackers and optical flow models: RAFT [37], SEA-  
 381 RAFT [40], AccFlow [41], PIPs++ [43], LocoTrack [6],  
 382 DELTA [29], BootsTAPIR [14], DOT [24], and Co-  
 383 Tracker [22, 23]. For fair-comparison, we use their publicly  
 384 released checkpoints for evaluation.

385 **Evaluation on TAP benchmarks.** We report results  
 386 on public point-tracking datasets spanning varied do-  
 387 mains: *BADJA* [3] (animals), *CroHD* [36] (surveillance),  
 388 *TAPVid-DAVIS* [12] (YouTube-like), *DriveTrack* [1] (driv-  
 389 ing), *EgoPoints* [11] (egocentric), *TAPVid-Kinetics* [12]  
 390 (YouTube), *RGB-Stacking* [12] (robotic manipulation), and  
 391 *RoboTAP* [39] (robotics). For long videos we cap the length  
 392 to at most 600 frames for tractability, and follow the official  
 393 TAP-Vid protocol for sampling query trajectories and vis-  
 394 ibility labels. On these benchmarks we use standard TAP  
 395 metrics: (1)  $\delta$ -**accuracy**, averaged over pixel thresholds  
 396  $k \in \{1, 2, 4, 8, 16\}$  at a canonical  $256 \times 256$  scale, measur-  
 397 ing the fraction of query points whose estimated positions  
 398 fall within  $k$  pixels of ground truth; (2) **Average Jaccard**  
 399 (**AJ**), which jointly evaluates localization and visibility on  
 400 TAP-style splits; and (3) **Occlusion Accuracy (OA)**, the

fraction of correctly predicted visible/occluded states. All  
 quantities are computed at official query points and reported  
 per dataset; AJ, and OA are additionally summarized on the  
 TAP-Vid subsets (DAVIS, Kinetics, RGB-Stacking, Robo-  
 TAP).

In Tab. 1–3, we report performance on each TAP-style  
 benchmark. Specifically, we use offline version of Co-  
 Tracker3 model for comparison. From Tab. 1 we observe  
 that PairFormer attains the best  $\delta_{\text{avg}}$  on 6/8 datasets, with  
 particularly strong gains on Rgb. and Rob., where dense  
 spatial context and long-range reasoning are crucial. On  
 Bad., Ego., and Kin. our method also slightly outperforms  
 CoTracker3. Nevertheless, when averaged over all bench-  
 marks, PairFormer achieves the highest  $\delta_{\text{avg}}$  (68.1 vs. 67.2),  
 indicating better overall robustness across diverse domains.  
 Tables 2 and 3 summarize AJ and OA on TAP-Vid datasets.  
 Here PairFormer consistently improves over CoTracker3.  
 These results show that our model provides more precise  
 trajectories while also yielding more reliable occlusion de-  
 cisions.

**Evaluation on APT benchmarks.** For evaluating APT,  
 we utilize three datasets with dense pairwise annotations:  
*CVO* [41], *CVO-Extended* [24], and *APT-Bench*. The dense  
 scenario can be formulated as densely sampled query points  
 and their corresponding trajectories, thus enabling the direct  
 adoption of metrics designed for TAP-style benchmarks.  
 Additionally, we introduce *strided-EPE* for a comprehen-  
 sive analysis. Specifically, for temporal gaps  $g \in [1, T - 1]$ ,  
 we compute average EPE (Endpoint Error) of displacement  
 field  $F_{s \rightarrow t}$  when  $|s - t| = g$ , denoted as  $\Delta_g^{\text{EPE}}$ . Notably,  
 CVO and CVO-Extended only provide first-to-any dense  
 flow, while APT-Bench supports all-pairs evaluation.

In Table 4, we summarize  $\delta_{\text{avg}}$ , EPE, AJ and strided-  
 EPE with stride  $g = 1, 5, 25$  across three datasets. Our  
 PairFormer attains the best performance on all reported  
 metrics, outperforming DOT and the CoTracker3 variants  
 in terms of both accuracy and calibration. The perfor-

Table 4. APT benchmark results and Strided-EPE on APT benchmarks.  $\Delta_g^{\text{epe}}$  denotes average EPE on frame stride  $g$  (lower is better).

Method	CVO [41]					CVO-Ext [24]					APT-Bench						
	$\delta_{\text{avg}} \uparrow$	EPE $\downarrow$	AJ $\uparrow$	$\Delta_1^{\text{epe}} \downarrow$	$\Delta_5^{\text{epe}} \downarrow$	$\delta_{\text{avg}} \uparrow$	EPE $\downarrow$	AJ $\uparrow$	$\Delta_1^{\text{epe}} \downarrow$	$\Delta_5^{\text{epe}} \downarrow$	$\Delta_{25}^{\text{epe}} \downarrow$	$\delta_{\text{avg}} \uparrow$	EPE $\downarrow$	AJ $\uparrow$	$\Delta_1^{\text{epe}} \downarrow$	$\Delta_5^{\text{epe}} \downarrow$	$\Delta_{25}^{\text{epe}} \downarrow$
CoTracker3_Offline	78.9	1.53	75.5	1.23	1.63	71.5	5.30	70.4	3.66	5.59	8.95	75.2	4.25	72.5	2.47	4.18	6.48
CoTracker3_Online	77.2	1.55	74.3	1.27	1.71	70.8	5.35	69.2	3.79	5.82	9.21	73.9	4.31	73.2	2.59	4.32	6.73
DOT	<u>81.2</u>	<u>1.40</u>	<u>79.8</u>	<u>1.12</u>	<u>1.39</u>	<u>72.7</u>	<u>5.28</u>	<u>70.8</u>	<u>3.52</u>	<u>5.38</u>	<u>8.77</u>	<u>76.7</u>	<u>3.68</u>	<u>75.4</u>	<u>2.18</u>	<u>3.69</u>	<u>6.02</u>
PairFormer	<b>81.8</b>	<b>1.32</b>	<b>80.3</b>	<b>1.04</b>	<b>1.31</b>	<b>73.2</b>	<b>5.13</b>	<b>71.2</b>	<b>3.34</b>	<b>5.09</b>	<b>8.26</b>	<b>77.5</b>	<b>3.26</b>	<b>76.3</b>	<b>2.03</b>	<b>3.21</b>	<b>5.39</b>

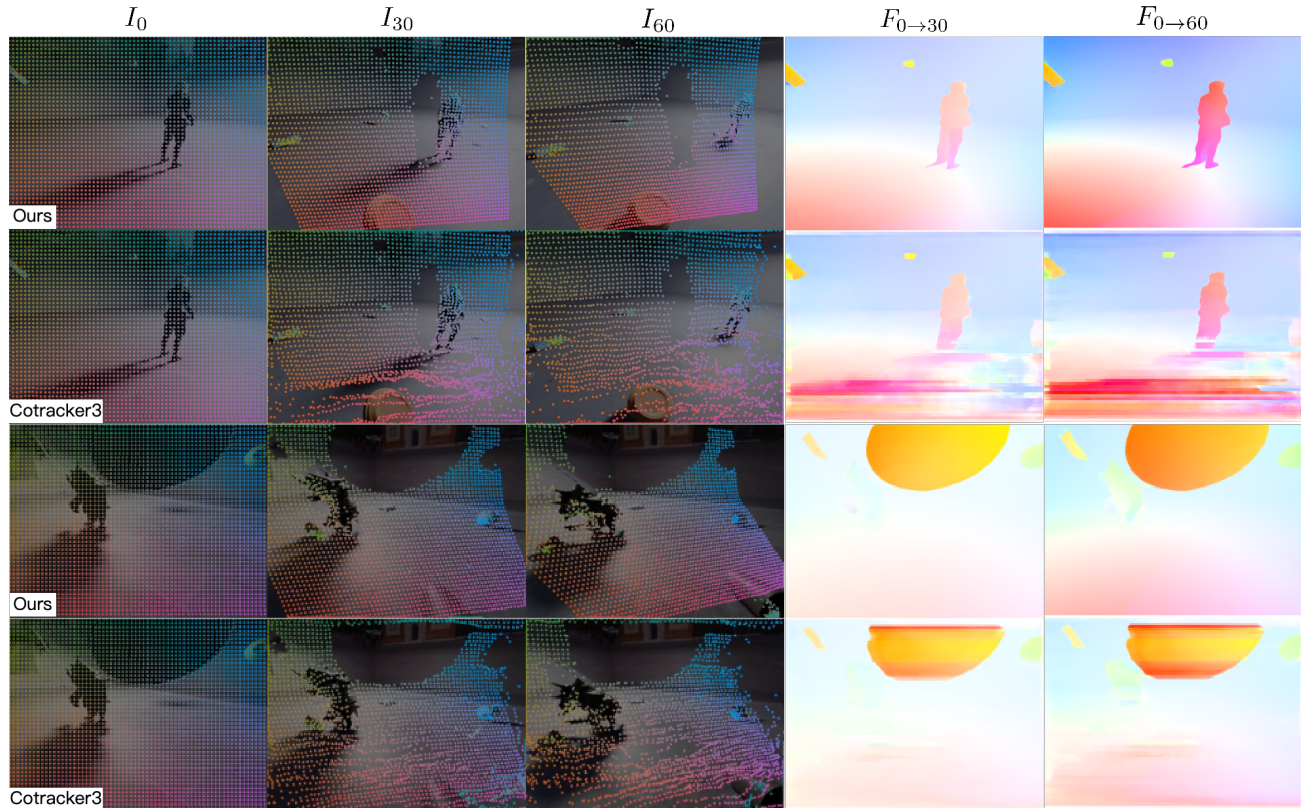


Figure 4. Flow and trajectory visualization across long temporal gaps. We visualize  $F_{s \rightarrow t}$  for  $s = 0$  and  $t = 30, 60$  on a representative APT-Bench sequence, comparing CoTracker3\_Offline and PairFormer.

438 mance gains are particularly notable on APT-Bench, which  
 439 consists of 100-frame-long videos. On CVO and CVO-  
 440 Extended (comprising  $< 50$ -frame video clips) the perfor-  
 441 mance margins are smaller yet consistent. This indicates  
 442 that our method achieves more substantial improvements on  
 443 long-sequence video benchmarks while maintaining steady  
 444 advantages over shorter video clips. This is further evi-  
 445 denced by strided-EPE values with stride  $g = 1, 5, 25$   
 446 on three dataset. Since CVO only contains 7-frame video  
 447 clips,  $\Delta_{25}^{\text{epe}}$  values are omitted for this dataset. For all  
 448 methods and datasets, EPE increases monotonically as  $g$   
 449 grows from 1 to 25, which matches the intuition that  
 450 longer temporal distances lead to larger displacements,  
 451 more occlusions, and higher ambiguity. However, the  
 452 growth rate for PairFormer is consistently lower. These

453 trends indicate that our model accumulates less error as  
 454 temporal distance increases, yielding more stable long-  
 455 range trajectories and fewer catastrophic failures under  
 456 prolonged occlusions and large viewpoint changes.

### 5.3. Qualitative Results 457

458 Figure 4 presents a qualitative comparison between Co-  
 459 Tracker3\_Offline and PairFormer on a representative se-  
 460 quence from the APT-Bench benchmark. For each method,  
 461 we first recover dense per-pixel trajectories by following  
 462 the predicted correspondences starting from the reference  
 463 frame, and then derive long-range flow maps from these  
 464 trajectories. The visualized results are in Figure 4. Com-  
 465 pared with CoTracker3\_Offline, PairFormer produces vis-  
 466 ibly smoother and more coherent trajectories that adhere

467 more faithfully to object geometry and remain stable over  
 468 time. The corresponding flow maps from PairFormer ex-  
 469 hibit sharper motion boundaries and are better aligned  
 470 with edges and structures in the starting frame, indicating  
 471 stronger spatial consistency between the estimated motion  
 472 and the underlying appearance. These qualitative observa-  
 473 tions are consistent with the quantitative gains in EPE and  
 474 strided-EPE, and suggest that PairFormer can propagate in-  
 475 formation through long sequences without sacrificing local  
 476 geometric fidelity.

#### 477 5.4. Ablations

478 We conduct three ablation studies under a fixed training  
 479 budget and report mean validation performance on a held-  
 480 out split of APT-Bench. Unless otherwise noted, all vari-  
 481 ants adopt the same optimization settings, data schedule,  
 482 and model width as the default PairFormer. We focus on  
 483 two summary metrics: the TAP-style score  $\delta_{\text{avg}}$  (higher is  
 484 better) and endpoint error (EPE, lower is better), averaged  
 485 over all validation videos.

486 **Depth allocation.** We first study how to allocate trans-  
 487 former depth between the ST-Patch Encoder and the Broad-  
 488 cast Motion Mixer (BMM). Concretely, we compare a  
 489 configuration with deeper encoding and shallower mixing  
 490 (24/9 blocks in ST-Patch Encoder/BMM) against a config-  
 491 uration that reverses this ratio (9/24), while keeping the to-  
 492 tal number of blocks comparable (Table 5, columns 1–2).  
 493 Both variants use the same CorrBank and decoder.

494 The results show that shifting capacity from the encoder  
 495 to BMM consistently improves both  $\delta_{\text{avg}}$  and EPE. With a  
 496 deeper encoder, the model can already extract strong local  
 497 features, but additional encoder layers mainly refine per-  
 498 frame descriptors. In contrast, allocating more depth to  
 499 BMM increases the capacity for long-range spatio-temporal  
 500 mixing over pairwise motion tokens, which directly benefits  
 501 APT: the 9/24 variant improves  $\delta_{\text{avg}}$  and reduces EPE, in-  
 502 dicating that extra computation is better used on trajectory-  
 503 and sequence-level reasoning than on further deepening the  
 504 initial cross-frame encoding.

505 **With and without  $\mathcal{L}_{\text{corr}}$ .** Next, we evaluate the effect of the  
 506 correspondence regularization term  $\mathcal{L}_{\text{corr}}$  from Section 3.3.  
 507 This term encourages points that belong to the same phys-  
 508 ical trajectory to share a consistent all-pairs prediction by  
 509 enforcing agreement between direct and multi-hop corre-  
 510 spondences along the trajectory (beyond what is enforced  
 511 by per-pair supervision alone).

512 Table 5, columns 3–4 compare models trained without  
 513 and with  $\mathcal{L}_{\text{corr}}$  under identical settings. Adding  $\mathcal{L}_{\text{corr}}$  yields  
 514 a moderate but consistent gain in  $\delta_{\text{avg}}$  and a reduction in  
 515 EPE. In practice, we observe that models trained with this  
 516 constraint produce smoother trajectories over long temporal  
 517 ranges and fewer small, inconsistent deviations when revis-

Table 5. Joint ablation of ST-Patch Encoder / BMM depth, corre-  
 spondence regularization, and CorrBank capacity.

Metric	Depth Alloc.		$\mathcal{L}_{\text{corr}}$		# Corr-tokens		
	24 / 9	9 / 24	w/o.	w/.	32	64	128
$\delta_{\text{avg}} \uparrow$	66.9	<b>67.6</b>	67.2	<b>67.8</b>	67.1	<b>67.6</b>	67.5
EPE $\downarrow$	4.4	<b>4.2</b>	4.3	<b>4.1</b>	4.3	<b>4.2</b>	4.2

518 iting the same region multiple times. This suggests that ex-  
 519 plicitly tying together correspondences along the same tra-  
 520 jectory helps the model resolve ambiguous matches and im-  
 521 proves temporal coherence, even though the loss is applied  
 522 only as a secondary regularizer.

523 **CorrBank capacity.** Finally, we study the capacity of Cor-  
 524 rBank by varying the number of learnable correlation to-  
 525 kens from 32 to 128 (Table 5, columns 5–7). Each token  
 526 participates in the construction of the pair-specific correla-  
 527 tion memory, so increasing this number enlarges the space  
 528 in which pairwise similarity patterns can be represented.

529 The ablation shows that using too few correlation tokens  
 530 (32) slightly hurts both  $\delta_{\text{avg}}$  and EPE, indicating that the  
 531 correlation memory becomes under-parameterized and can-  
 532 not fully capture complex appearance-motion relationships.  
 533 Increasing to 64 tokens yields a noticeable improvement  
 534 and is used as our default setting. Further expanding to 128  
 535 tokens does not bring additional gains and only increases  
 536 memory and compute. Overall, these results suggest that a  
 537 moderately sized CorrBank is sufficient to encode the rele-  
 538 vant pairwise structure, and that most of the modeling bene-  
 539 fit comes from how this memory is used (via CorrBank and  
 540 BMM) rather than from unbounded capacity.

## 541 6. Conclusion

542 This work formalizes All-Pairs Tracking (APT), a paradigm  
 543 that explicitly models dense displacement and visibility  
 544 across all frame pairs in a video, extending beyond the  
 545 query-conditioned focus of existing TAP methods. We  
 546 propose PairFormer, a feed-forward transformer integrat-  
 547 ing spatio-temporal encoding, CorrBank, and global con-  
 548 text aggregation, alongside the PAIRender platform for  
 549 synthetic annotated data (supporting  $\pi$ -R10K and APT-  
 550 Bench). Experiments confirm PairFormer’s state-of-the-art  
 551 performance on APT-Bench and competitive results on TAP  
 552 benchmarks, advancing dense spatiotemporal tracking via a  
 553 unified task, model, and data ecosystem. We believe that  
 554 making all-pairs correspondence a first-class representation  
 555 also opens up opportunities for downstream applications  
 556 such as 4D reconstruction and motion editing.

557 **Limitations.** PairFormer faces quadratic computational  
 558 costs for ultra-long videos, and its reliance on synthetic  
 559 PAIRender data might introduces a domain gap with real-  
 560 world scenarios.

561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616

## References

- [1] Arjun Balasingam, Joseph Chandler, Chenning Li, Zhoutong Zhang, and Hari Balakrishnan. Drivetrack: A benchmark for long-range point tracking in real-world videos. In *CVPR*, 2024. 6
- [2] Weikang Bian, Zhaoyang Huang, Xiaoyu Shi, Yitong Dong, Yijin Li, and Hongsheng Li. Context-pips: Persistent independent particles demands context features. *NeurIPS*, 2024. 2
- [3] Benjamin Biggs, Thomas Roddick, Andrew Fitzgibbon, and Roberto Cipolla. Creatures great and SMAL: Recovering the shape and motion of animals from video. In *ACCV*, 2018. 6
- [4] Thomas Brox and Jitendra Malik. Large displacement optical flow: descriptor matching in variational motion estimation. *TPAMI*, 33(3):500–513, 2010. 2
- [5] Seokju Cho, Jiahui Huang, Seungryong Kim, and Joon-Young Lee. Flowtrack: Revisiting optical flow for long-range dense tracking. In *CVPR*, 2024. 1, 2
- [6] Seokju Cho, Jiahui Huang, Jisu Nam, Honggyu An, Seungryong Kim, and Joon-Young Lee. Local all-pair correspondence for point tracking. In *ECCV*, 2025. 2, 6
- [7] Pierre-Henri Conze, Philippe Robert, Tomas Crivelli, and Luce Morin. Dense long-term motion estimation via statistical multi-step flow. In *VISAPP*, 2014. 2
- [8] Pierre-Henri Conze, Philippe Robert, Tomas Crivelli, and Luce Morin. Multi-reference combinatorial strategy towards longer long-term dense motion estimation. *Computer Vision and Image Understanding*, 150:66–80, 2016.
- [9] Tomas Crivelli, Pierre-Henri Conze, Philippe Robert, and Patrick Pérez. From optical flow to dense long term correspondences. In *International Conference on Image Processing*, 2012. 2
- [10] Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations (ICLR)*, 2024. 4
- [11] Ahmad Darkhalil, Rhodri Guerrier, Adam W Harley, and Dima Damen. Egopoints: Advancing point tracking for ego-centric videos. *arXiv:2412.04592*, 2024. 6
- [12] Carl Doersch, Ankush Gupta, Larisa Markeeva, Adrià Recasens, Lucas Smaira, Yusuf Aytar, João Carreira, Andrew Zisserman, and Yi Yang. TAP-Vid: A benchmark for tracking any point in a video. In *NeurIPS Datasets and Benchmarks*, 2022. 1, 2, 6
- [13] Carl Doersch, Yi Yang, Mel Vecerik, Dilara Gokay, Ankush Gupta, Yusuf Aytar, Joao Carreira, and Andrew Zisserman. TAPIR: Tracking any point with per-frame initialization and temporal refinement. In *ICCV*, 2023. 1, 2
- [14] Carl Doersch, Pauline Luc, Yi Yang, Dilara Gokay, Skanda Koppula, Ankush Gupta, Joseph Heyward, Ignacio Rocco, Ross Goroshin, João Carreira, et al. Bootstap: Bootstrapped training for tracking-any-point. In *ACCV*, 2024. 6
- [15] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, 2015. 2
- [16] James J Gibson. The perception of visual surfaces. *The American journal of psychology*, 63(3):367–384, 1950. 2
- [17] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanaprasam, Florian Golemo, Charles Herrmann, et al. Kubric: A scalable dataset generator. In *CVPR*, 2022. 2, 6
- [18] Adam W Harley, Zhaoyuan Fang, and Katerina Fragkiadaki. Particle video revisited: Tracking through occlusions using point trajectories. In *ECCV*, 2022. 1, 2
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 2
- [20] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 1, 2
- [21] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, 2017. 2
- [22] Nikita Karaev, Iurii Makarov, Jianyuan Wang, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. CoTracker3: Simpler and better point tracking by pseudo-labelling real videos. *arXiv:2410.11831*, 2024. 1, 2, 6
- [23] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. CoTracker: It is better to track together. In *ECCV*, 2024. 2, 6
- [24] Guillaume Le Moing, Jean Ponce, and Cordelia Schmid. Dense optical tracking: connecting the dots. In *CVPR*, 2024. 2, 6, 7
- [25] Hongyang Li, Hao Zhang, Shilong Liu, Zhaoyang Zeng, Feng Li, Tianhe Ren, Bohan Li, and Lei Zhang. TAPTRv2: Attention-based position update improves tracking any point. In *NeurIPS*, 2024. 2
- [26] Hongyang Li, Hao Zhang, Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, and Lei Zhang. TAPTR: Tracking any point with transformers as detection. In *ECCV*, 2024. 2
- [27] Bruce D Lucas, Takeo Kanade, et al. An iterative image registration technique with an application to stereo vision. In *IJCAI*, 1981. 1, 2
- [28] Michal Neoral, Jonáš Šerých, and Jiří Matas. Mft: Long-term tracking of every pixel. In *WACV*, 2024. 1, 2
- [29] Tuan Duc Ngo, Peiye Zhuang, Chuang Gan, Evangelos Kalogerakis, Sergey Tulyakov, Hsin-Ying Lee, and Chaoyang Wang. DELTA: Dense efficient long-range 3d tracking for any video. In *ICLR*, 2025. 6
- [30] Diederick C Niehorster. Optic flow: a history. *i-Perception*, 12(6):20416695211055766, 2021. 2
- [31] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020. 3, 4
- [32] P. Sand and S. Teller. Particle video: Long-range motion estimation using point trajectories. In *CVPR*, 2006. 2
- [33] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan

674 Wang. Real-time single image and video super-resolution  
675 using an efficient sub-pixel convolutional neural network. In  
676 *CVPR*, 2016. 2

677 [34] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz.  
678 PWC-Net: CNNs for optical flow using pyramid, warping,  
679 and cost volume. In *CVPR*, 2018. 1, 2

680 [35] Narayanan Sundaram, Thomas Brox, and Kurt Keutzer.  
681 Dense point trajectories by GPU-accelerated large displace-  
682 ment optical flow. In *ECCV*, 2010. 2

683 [36] Ramana Sundararaman, Cedric De Almeida Braga, Eric  
684 Marchand, and Julien Pettre. Tracking pedestrian heads in  
685 dense crowd. In *CVPR*, 2021. 6

686 [37] Zachary Teed and Jia Deng. RAFT: Recurrent all-pairs field  
687 transforms for optical flow. In *ECCV*, 2020. 1, 2, 6

688 [38] Carlo Tomasi and Takeo Kanade. Detection and tracking of  
689 point. *IJCV*, 9:137–154, 1991. 2

690 [39] Mel Vecerik, Carl Doersch, Yi Yang, Todor Davchev, Yusuf  
691 Aytar, Guangyao Zhou, Raia Hadsell, Lourdes Agapito, and  
692 Jon Scholz. RoboTAP: Tracking arbitrary points for few-shot  
693 visual imitation. In *ICRA*, 2024. 6

694 [40] Yihan Wang, Lahav Lipson, and Jia Deng. Sea-raft: Simple,  
695 efficient, accurate raft for optical flow. In *ECCV*, 2024. 1, 2,  
696 6

697 [41] Guangyang Wu, Xiaohong Liu, Kunming Luo, Xi Liu,  
698 Qingqing Zheng, Shuaicheng Liu, Xinyang Jiang, Guangtao  
699 Zhai, and Wenyi Wang. Accflow: Backward accumulation  
700 for long-range optical flow. In *ICCV*, 2023. 6, 7

701 [42] Jia Xu, René Ranftl, and Vladlen Koltun. Accurate optical  
702 flow via direct cost volume processing. In *CVPR*, 2017. 2

703 [43] Yang Zheng, Adam W. Harley, Bokui Shen, Gordon Wet-  
704 zstein, and Leonidas J. Guibas. Pointodyssey: A large-scale  
705 synthetic dataset for long-term point tracking. In *ICCV*,  
706 2023. 2, 6