

Forecasting Animal Motion in the Wild

Anonymous CVPR submission

Paper ID 6

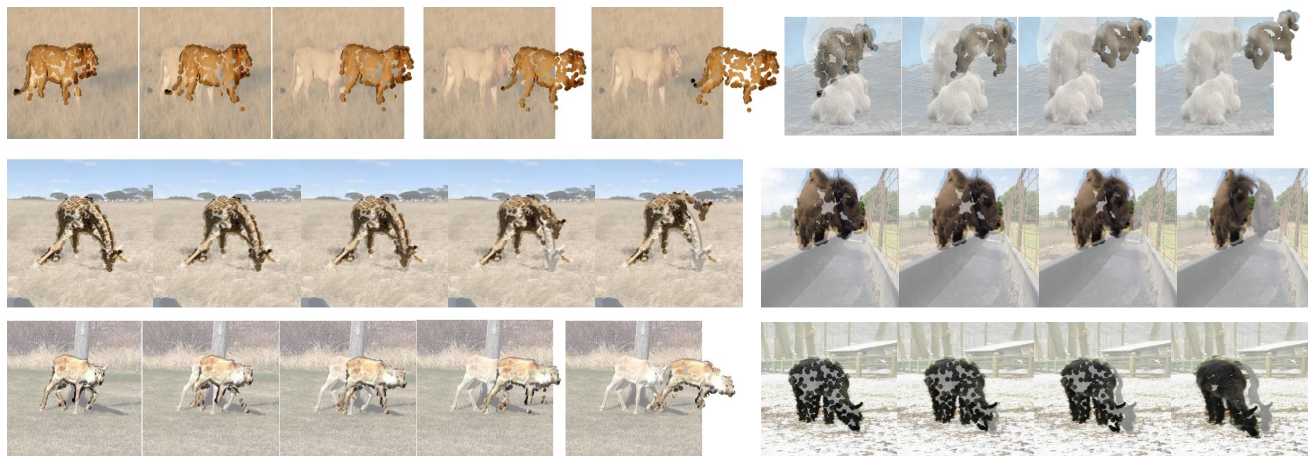


Figure 1. **Dense point trajectories act as visual tokens for behavior, enabling scalable prediction of complex motion across diverse species.** Our method takes as input a single RGB image and a short history of motion, and forecasts future animal motion in the form of point trajectories. For each predicted point trajectory, we translate a small circular patch of the input image along the motion trajectory and superimpose it on the input image (**no pixels are generated!**). Leftmost shows the start locations on the input frame; the rest is forecast by our model. Our method is capable of forecasting many different species and behaviors, even long-tail ones; the polar bear (top right) is only present in 0.31% of the training data, the caribou (bottom left) in 0.025%, and the alpaca (bottom right) in 0.50%. [Results video here.](#)

Abstract

001 Visual intelligence requires anticipating the future behavior
 002 of agents, yet vision systems lack a general representation
 003 for motion and behavior. We propose dense point trajectories
 004 as visual tokens for behavior, a structured mid-level repre-
 005 sentation that disentangles motion from appearance and
 006 generalizes across diverse non-rigid agents, such as animals
 007 in-the-wild. Building on this abstraction, we design a diffu-
 008 sion transformer that models unordered sets of trajectories
 009 and explicitly reasons about occlusion, enabling coherent
 010 forecasts of complex motion patterns. To evaluate at scale,
 011 we curate MammalMotion, 300 hours of unconstrained ani-
 012 mal video with robust shot detection and camera-motion
 013 compensation. Experiments show that forecasting trajectory
 014 tokens achieves category-agnostic, data-efficient prediction,
 015 outperforms state-of-the-art baselines, and generalizes to
 016 rare species and morphologies, providing a foundation for

predictive visual intelligence in the wild.

017

1. Introduction

018

Predicting the future motion of objects and agents is a fun-
 019 damental capability of visual intelligence. In dynamic envi-
 020 ronments, agents, from animals in the wild to humans in
 021 social settings, must anticipate the behavior of others in order
 022 to act effectively or survive. Despite major advances in vi-
 023 sual recognition and generation, predicting behavior remains
 024 one of the least understood capabilities of modern vision
 025 systems.

A key reason for this gap is the lack of an appropriate
 027 representation for behavior. In language, prediction is en-
 028 abled by discrete tokens that structure the modeling problem.
 029 Vision systems lack an analogous token for motion and be-
 030 havior. In this work, we show that dense point trajectories
 031 can serve as such tokens, enabling scalable prediction of
 032 behavior across diverse agents. To understand why such a
 033

034 representation is needed, consider the limitations of existing
035 approaches. Forecasting directly in pixel space is universal
036 but poorly structured: while recent video diffusion models
037 can generate realistic short clips, forecasting behavior di-
038 rectly in pixel space entangles appearance, lighting, and
039 camera motion with object dynamics, making the learning
040 problem unnecessarily complex and data inefficient. At the
041 opposite extreme, parameterized 3D models provide com-
042 pact and physically valid representations for forecasting, but
043 rely on strong object-specific priors and therefore apply only
044 to a small number of carefully modeled categories, such as
045 humans [43] and a handful of animals [61, 69, 78, 87–90].
046 Even in these settings they often miss fine-grained deforma-
047 tion and shape variation. Without an intermediate representa-
048 tion that captures motion structure while remaining general,
049 scalable behavior prediction remains difficult. We therefore
050 seek a representation that introduces structure without sacri-
051 ficing generality.

052 We therefore propose dense point trajectories as visual
053 tokens for behavior, providing a structured representation
054 for forecasting motion across diverse agents. While sparse
055 points carry little semantic meaning when static, their mo-
056 tion reveals rich information about 3D structure and intent,
057 as demonstrated in Johansson’s classical biological motion
058 studies [29] and subsequent work [2, 13, 17, 21, 35]. Rep-
059 resenting behavior as evolving 2D point tracks focuses pre-
060 diction directly on motion dynamics while remaining ag-
061 nostic to appearance and scene variation. This formulation
062 is significantly more data efficient than forecasting pixels
063 directly [4, 5] and naturally applies to arbitrary non-rigid
064 agents without requiring category-specific models. Point
065 trajectories therefore occupy a principled middle ground be-
066 tween raw pixels and full 3D parameterizations: structured
067 enough to constrain prediction, yet general enough to scale
068 across species, morphologies, and environments.

069 Building on this abstraction, we introduce a diffusion
070 transformer that forecasts behavior from short motion histo-
071 ries. Unlike prior trajectory-based approaches designed for
072 robotics or rigid scenes [4, 11, 77], our formulation models
073 motion for non-rigid agents in the wild. The model predicts
074 future behavior as an unordered set of point trajectories
075 (Fig. 1), treating each trajectory as a token augmented with
076 local visual context from DINOv3 features. The architecture
077 jointly models trajectories while explicitly reasoning about
078 occlusion and visibility, enabling coherent predictions of
079 complex non-rigid motion. Our model learns diverse motion
080 patterns including gait, cyclical, and linear behaviors, and
081 forecasts future motion across a wide range of species, out-
082 performing state-of-the-art baselines. Training on the broad
083 diversity of motion found in nature further enables general-
084 ization to previously unseen categories and morphologies of
085 animate agents.

086 To study long-tailed biological motion at scale, we fo-

087 cus on unconstrained video of animals in the wild. Animals
088 provide a particularly challenging testbed for behavior pre-
089 diction: they exhibit highly diverse morphologies and motion
090 patterns, and data for many species is inherently sparse. A
091 representation that succeeds in this regime must generalize
092 across categories without relying on category-specific mod-
093 els. We develop a large-scale pipeline for isolating animal
094 motion from raw video, including robust shot detection and
095 camera-motion compensation, and curate over 300 hours
096 of annotated footage for behavior forecasting which we re-
097 lease with this paper. Using this in-the-wild data, we demon-
098 strate that our approach operates on tracks extracted from
099 unconstrained video and is robust to the noise and partial
100 observability inherent in real-world tracking. This dataset
101 reveals previously unreported statistical structure in animal
102 motion, and provides a foundation for studying predictive
103 visual intelligence in natural environments.

Our contributions are:

1. **Point trajectories as visual tokens for behavior forecasting.** We introduce dense point tracks as a compact mid-level representation for modeling long-tailed natural-world behavior that disentangles motion from appearance and generalizes beyond category-specific 3D models.
2. **A diffusion transformer for trajectory forecasting.** We design a DiT-based architecture that treats trajectories as tokens and predicts diverse futures of non-rigid behavior from short histories while explicitly reasoning about occlusion in unordered track sets.
3. **MammalMotion, a large-scale dataset of animal motion.** We develop a robust pipeline for isolating animal motion in unconstrained video and release 300 hours of annotated footage for behavior forecasting in the wild.

2. Method

2.1. Diffusion-based Point Trajectory Forecasting

We model animal motion as a set of N point tracks $\mathbf{X} \in \mathbb{R}^{T \times N \times 2}$ over T timesteps. Each track \mathbf{x}_n consists of normalized coordinates (x_n^t, y_n^t) and an associated visibility state $\mathbf{O}_n^t \in [0, 1]$. We learn the conditional distribution $p(\mathbf{X}_{T_c+1:T}, \mathbf{O}_{T_c+1:T} | \mathbf{I}, \mathbf{X}_{1:T_c}, \mathbf{O}_{1:T_c}, \mathbf{d})$, where \mathbf{I} is the first frame and \mathbf{d} is an optional 2D global displacement.

To improve training dynamics and handle occlusions, we reparameterize the diffusion target as $\mathbf{Z}_0^{\text{diff}} = \{\gamma \mathbf{V}, \beta \mathbf{O}\}$, where \mathbf{V} represents velocities $\dot{x}_n^t = x_n^{t+1} - x_n^t$. For occluded points, we use linear interpolation between the nearest visible frames. Following DDPM [26], our model f_θ minimizes the L_1 denoising objective:

$$\mathcal{L} = \mathbb{E}_{\mathbf{Z}_0^{\text{diff}}, \tau, \epsilon} [\|\mathbf{Z}_0^{\text{diff}} - f_\theta(\mathbf{Z}_\tau^{\text{diff}}, \mathbf{Z}^{\text{cond}}, \tau)\|_1] \quad (1)$$

where \mathbf{Z}^{cond} encapsulates the image \mathbf{I} , motion history and initial spatial positions, and optionally \mathbf{d} . For efficient inference, we employ DDIM [67] with 100 steps.

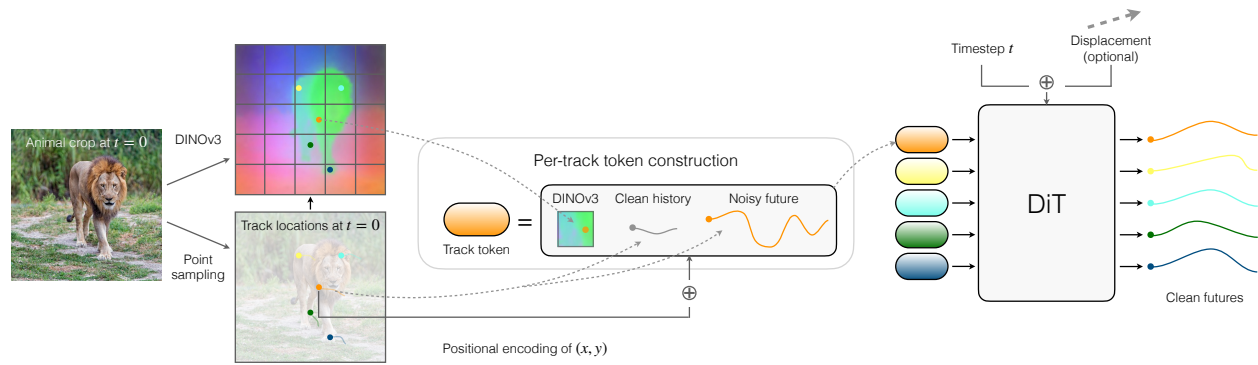


Figure 2. **Architecture.** Given an input frame and (noisy) tracks, we construct a single token for every track, which includes a DINO feature at the start location, the motion history, and the noisy track values, both with occlusion indicators. After projection, we add a position encoding for the initial point location. Tokens are stacked and fed to a transformer (DiT) to predict clean tracks (right).

137 2.2. Diffusion Transformer Architecture

138 Our denoiser f_θ is a Transformer (DiT) that treats each
 139 track as a token, ensuring permutation invariance. Each
 140 token for the n -th track is a concatenation: $\mathbf{Z}_n =$
 141 $[\mathbf{Z}_n^{\text{diff}}, \mathbf{f}_n^{\text{DINO}}, \mathbf{V}_{1:T_c}, \mathbf{O}_{1:T_c}]$.

- 142 • **Visual Context:** We extract local features $\mathbf{f}_n^{\text{DINO}}$ from a
 143 frozen DINOv3 backbone via bilinear interpolation at the
 144 track’s initial location (x_n^1, y_n^1) .
- 145 • **Motion Context:** Velocity histories are embedded using
 146 sinusoidal encodings and scaled by γ to match noise vari-
 147 ance; occlusion histories are kept as scalar and multiplied
 148 by β .

149 After concatenation, we project the tokens to the trans-
 150 former’s hidden dimension and add a sinusoidal positional
 151 encoding of the initial coordinates (x_n^1, y_n^1) to retain explicit
 152 spatial relationships. Global conditioning variables—the dif-
 153 fusion timestep τ and the optional displacement \mathbf{d} —are inte-
 154 grated directly into each DiT layer via AdaLN [54].

155 3. Data Processing and Motion Distribution

156 To isolate animal motion from camera ego-motion in the
 157 wild, we developed a data processing pipeline and camera
 158 stabilization pipeline that we apply to the MammalNet
 159 dataset [10]. We utilize VideoSAM [59] to segment animals
 160 and BootsTAPIR [15] for point tracking within segments.
 161 We then use a RANSAC-based homography estimation [15]
 162 on background points (excluding segments from VideoSAM
 163 [59]) and transform points into a camera-stabilized coordi-
 164 nate system normalized to an initial animal bounding box.
 165 This results in MammalMotion, ~ 300 hrs of animal motion,
 166 which we release.

167 **Log-Normal Distribution of Motion.** To validate our mo-
 168 tion processing, we compute a histogram of the average
 169 displacement (i.e. average distance between start and end
 170 for points that remain visible) in Figure 4. While one might

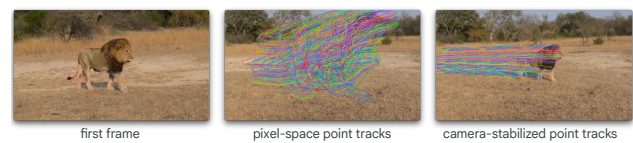


Figure 3. **Our processed data before and after camera stabilization.** Given a first frame (left), the middle image shows the point tracks in pixel space, where the motion of the animals and the camera (panning, zooming out) are entangled. On the right are our point tracks in camera-stabilized space.

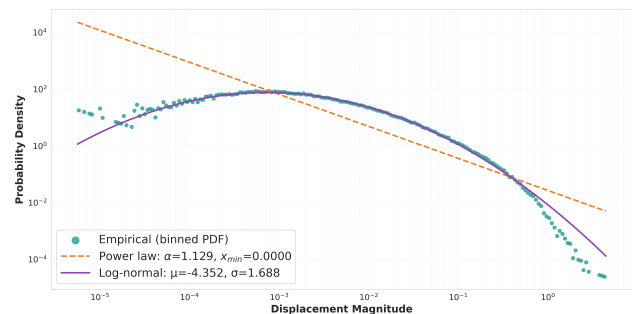


Figure 4. **Animal motion follows a log normal distribution:** We plot a histogram of animal displacement. Horizontal axis is a binned log displacement, while vertical axis is log frequency. We find that log-normal (purple) fits much better than both a power law (orange).

171 expect animal motion to follow a power law, we instead
 172 find that a log-normal distribution fits far better (i.e., the
 173 log displacements are normally distributed). Such distribu-
 174 tions have been found in other datasets of animal motions,
 175 e.g. Lévy flights [8, 23, 28], foraging decisions in rats [30],
 176 and general spontaneous behavior in animals [57], and are
 177 suggested to imply that motion magnitude is the result of a
 178 *multiplicative* interaction of independent factors. We believe
 179 this is the first time such a result has arisen from such a
 180 diverse dataset of many different species, and without any
 181 painstaking manual annotation or use of tracking devices.



(a) **Samples from our model.** Sampling from our model with different random seeds (each row) and no displacement conditioning. The frame on the left is the input state after the motion history. We see different frequencies of the grooming behavior for the jaguar and the dog’s head moves different directions.



(b) **Out-of-distribution.** Our model *generalizes* to humans and non-mammals.



(c) **Our Model vs. Stable Diffusion.** Our approach can model the behavior of less common animals in our dataset such as hares, while conventional video models struggle with these animals.

Figure 5. **Qualitative Forecasting Results.**

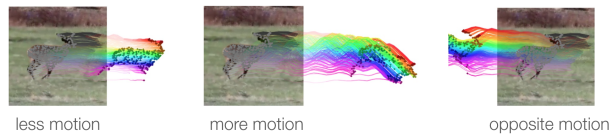


Figure 6. **Prompting our model with different levels of motion.** Grey represents motion history, colors are our predictions.

182

4. Results

183

See Sec 10 for our experimental setup: how we use the MammalNet dataset, our metrics, and our baselines.

184

185

Samples from our model. Figure 1 shows qualitative results which exhibit convincing forecasts across a variety of behaviors: e.g. the lion’s legs follows natural articulation, giraffe raises its neck, the alpaca grazes naturally. This also works for rare animal categories. Figure 5a shows the diversity that our model produces with only different seeds. Fig. 6 demonstrates prompting our model for more motion, less motion, or motion in a different direction, and the model produces plausible behaviors consistent with these motions. Static visualizations do not do justice to motion accuracy; we urge our readers to watch the [supplementary video](#).

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

Out-of-distribution examples. Fig. 5b displays qualitative results of our model on OOD data. We note that the MammalNet dataset does have videos that contain humans and other types of animals; the ostrich on the bottom right was found in our validation set, so this type of generalization

Table 1. **Quantitative comparison on MammalMotion.**

Method	ADE ↓	FDE ↓	Avg PWT ↑	FD (V) ↓	FVMD ↓
Constant Vel	0.104	0.215	41.15%	2.59	89.77
WHN	0.105	0.200	29.92%	5.34	94.70
Track2Act	0.064	0.126	43.04%	3.20	55.84
Ours	0.046	0.102	60.01%	1.96	17.0

is not surprising. However, the Lego robot (bottom left) and butterfly (top right) are unlike the expected MammalNet data distribution, but still observe physically plausible motion.

Comparison with Video Generation Models. Figure 5c displays a comparison with Stable Video Diffusion [6]. While video models often struggle with physical realism due to the overhead of modeling pixels (textures, lighting), our trajectory-token approach produces more realistic biological behavior with less compute and data. This extends the findings of [7] from rigid synthetic objects to in-the-wild nonrigid motion. E.g., our model is able to forecast the foraging behavior of a hare even though hares only constitute 0.39% of the training data. The video model struggles not only to model this behavior but even to maintain basic anatomy, morphing ears into wings.

Quantitative Results See 6 for detailed quantitative results. Results suggest that our method substantially outperforms other approaches on all metrics, and that when training on our full dataset, there is transfer between species.

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276**References**

- [1] David J Anderson and Pietro Perona. Toward a science of computational ethology. *Neuron*, 84(1):18–31, 2014. 3
- [2] Anthony P Atkinson, Winand H Dittrich, Andrew J Gemmell, and Andrew W Young. Emotion perception from dynamic and static body expressions in point-light and full-light displays. *Perception*, 33(6):717–746, 2004. 2
- [3] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical commonsense for video generation. In *The Thirteenth International Conference on Learning Representations*, 2025. 2
- [4] Homanga Bharadhwaj, Roozbeh Mottaghi, Abhinav Gupta, and Shubham Tulsiani. Track2act: Predicting point tracks from internet videos enables generalizable robot manipulation. In *European Conference on Computer Vision (ECCV)*, 2024. 2, 4, 8
- [5] Homanga Bharadhwaj, Debidatta Dwibedi, Abhinav Gupta, Shubham Tulsiani, Carl Doersch, Ted Xiao, Dhruv Shah, Fei Xia, Dorsa Sadigh, and Sean Kirmani. Gen2act: Human video generation in novel scenarios enables generalizable robot manipulation. In *Conference on Robot Learning*, pages 3936–3951. PMLR, 2025. 2
- [6] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 4, 1
- [7] Gabrijel Boduljak, Laurynas Karazija, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. What happens next? anticipating future motion by generating point trajectories. *arXiv preprint arXiv:2509.21592*, 2025. 4, 2, 8
- [8] Greg A Breed, Paul M Severns, and Andrew M Edwards. Apparent power-law distributions in animal movements can arise from intraspecific interactions. *Journal of the Royal Society Interface*, 12(103), 2015. 3
- [9] Hila Chefer, Uriel Singer, Amit Zohar, Yuval Kirstain, Adam Polyak, Yaniv Taigman, Lior Wolf, and Shelly Sheynin. Videojam: Joint appearance-motion representations for enhanced motion generation in video models. In *Forty-second International Conference on Machine Learning*. 2
- [10] Jun Chen, Ming Hu, Darren J Coker, Michael L Berumen, Blair Costelloe, Sara Beery, Anna Rohrbach, and Mohamed Elhoseiny. Mammalnet: A large-scale video benchmark for mammal recognition and behavior understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13052–13061, 2023. 3, 5, 7
- [11] Yixiang Chen, Peiyan Li, Yan Huang, Jiabing Yang, Kehan Chen, and Liang Wang. Ec-flow: Enabling versatile robotic manipulation from action-unlabeled videos via embodiment-centric flow. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11958–11968, 2025. 2
- [12] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets. *arXiv preprint arXiv:1907.06571*, 2019. 2
- [13] James E Cutting and Lynn T Kozlowski. Recognizing friends by their walk: Gait perception without familiarity cues. *Bulletin of the psychonomic society*, 9(5):353–356, 1977. 2
- [14] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10061–10072, 2023. 2, 8
- [15] Carl Doersch, Pauline Luc, Yi Yang, Dilara Gokay, Skanda Koppula, Ankush Gupta, Joseph Heyward, Ignacio Rocco, Ross Goroshin, Joao Carreira, et al. Bootstap: Bootstrapped training for tracking-any-point. In *Proceedings of the Asian Conference on Computer Vision*, pages 3257–3274, 2024. 3, 2, 6, 7
- [16] David C Dowson and B. V. Landau. The fréchet distance between multivariate normal distributions. *Journal of multivariate analysis*, 12(3):450–455, 1982. 7
- [17] Robert Fox and Cynthia McDaniel. The perception of biological motion by human infants. *Science*, 218(4571):486–487, 1982. 2
- [18] Chongkai Gao, Haozhuo Zhang, Zhixuan Xu, Zhehao Cai, and Lin Shao. Flip: Flow-centric generative planning for general-purpose manipulation tasks. *arXiv preprint arXiv:2412.08261*, 2024. 2
- [19] Google DeepMind. Veo 3 technical report. Technical report, Google DeepMind, 2025. 2
- [20] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3749–3761, 2022. 2
- [21] Emily Grossman, Michael Donnelly, R Price, D Pickens, V Morgan, G Neighbor, and Randolph Blake. Brain areas involved in perception of biological motion. *Journal of cognitive neuroscience*, 12(5):711–720, 2000. 2
- [22] Xianfan Gu, Chuan Wen, Weirui Ye, Jiaming Song, and Yang Gao. Seer: Language instructed video prediction with latent diffusion models. *arXiv preprint arXiv:2303.14897*, 2023. 2
- [23] Richard Gunner, Rory Wilson, Miguel Lurgi, Luca Borger, James Redcliffe, Emily Shepard, Mark Holton, Margaret Crofoot, Abdulaziz Alagaili, Samantha Andrzejczek, et al. High resolution data reveal fundamental steps and turning points in animal movements. 2024. 3
- [24] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2255–2264, 2018. 3
- [25] Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama. Photorealistic video generation with diffusion models. In *ECCV*, 2024. 2
- [26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2, 4

- 334 [27] Tobias Höppe, Arash Mehrjou, Stefan Bauer, Didrik Nielsen, 392
335 and Andrea Dittadi. Diffusion models for video prediction 393
336 and infilling. *arXiv preprint arXiv:2206.07696*, 2022. 2 394
337 [28] Nicolas E Humphries, Nuno Queiroz, Jennifer RM Dyer, 395
338 Nicolas G Pade, Michael K Musyl, Kurt M Schaefer, 396
339 Daniel W Fuller, Juerg M Brunnschweiler, Thomas K Doyle, 397
340 Jonathan DR Houghton, et al. Environmental context explains 398
341 lévy and brownian movement patterns of marine predators. 399
342 *Nature*, 465(7301):1066–1069, 2010. 3 400
343 [29] Gunnar Johansson. Visual perception of biological motion 401
344 and a model for its analysis. *Perception & psychophysics*, 14 402
345 (2):201–211, 1973. 2 403
346 [30] Kanghoon Jung, Hyeran Jang, Jerald D Kralik, and Jaeseung 404
347 Jeong. Bursts and heavy tails in temporal and sequential 405
348 dynamics of foraging decisions. *PLoS computational biology*, 406
349 10(8):e1003759, 2014. 3 407
350 [31] Bingyi Kang, Yang Yue, Rui Lu, Zhijie Lin, Yang Zhao, 408
351 Kaixin Wang, Gao Huang, and Jiashi Feng. How far is video 409
352 generation from world model: A physical law perspective. In 410
353 *International Conference on Machine Learning*, pages 28991– 411
354 29017. PMLR, 2025. 2 412
355 [32] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia 413
356 Neverova, Andrea Vedaldi, and Christian Rupprecht. Co- 414
357 tracker: It is better to track together. In *European conference 415*
358 *on computer vision*, pages 18–35. Springer, 2024. 2 416
359 [33] Nikita Karaev, Yuri Makarov, Jianyuan Wang, Natalia 417
360 Neverova, Andrea Vedaldi, and Christian Rupprecht. Co- 418
361 tracker3: Simpler and better point tracking by pseudo- 419
362 labelling real videos. In *Proceedings of the IEEE/CVF Inter- 420*
363 *national Conference on Computer Vision*, pages 6013–6022, 421
364 2025. 2 422
365 [34] Kris M Kitani, Brian D Ziebart, James Andrew Bagnell, and 423
366 Martial Hebert. Activity forecasting. In *European conference 424*
367 *on computer vision*, pages 201–214. Springer, 2012. 3 425
368 [35] Lynn T Kozlowski and James E Cutting. Recognizing the sex of 426
369 a walker from a dynamic point-light display. *Perception & 427*
370 *psychophysics*, 21(6):575–580, 1977. 2 428
371 [36] Jessy Lauer, Mu Zhou, Shaokai Ye, William Menegas, Steffen 429
372 Schneider, Tanmay Nath, Mohammed Mostafizur Rahman, 430
373 Valentina Di Santo, Daniel Soberanes, Guoping Feng, et al. 431
374 Multi-animal pose estimation, identification and tracking with 432
375 deeplabcut. *Nature Methods*, 19(4):496–504, 2022. 3 433
376 [37] Alex X Lee, Richard Zhang, Frederik Ebert, Pieter Abbeel, 434
377 Chelsea Finn, and Sergey Levine. Stochastic adversarial video 435
378 prediction. *arXiv preprint arXiv:1804.01523*, 2018. 2 436
379 [38] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B 437
380 Choy, Philip HS Torr, and Manmohan Chandraker. Desire: 438
381 Distant future prediction in dynamic scenes with interacting 439
382 agents. In *Proceedings of the IEEE conference on computer 440*
383 *vision and pattern recognition*, pages 336–345, 2017. 3 441
384 [39] Shijie Li, Chunyu Liu, Xun Xu, Si Yong Yeo, and Xulei Yang. 442
385 Future-aware interaction network for motion forecasting. In 443
386 *Proceedings of the IEEE/CVF International Conference on 444*
387 *Computer Vision (ICCV)*, pages 7505–7515, 2025. 3 445
388 [40] Jiaye Liu, Youran Qu, Qi Yan, Xiaohui Zeng, Lele Wang, 446
389 and Renjie Liao. Fr\`echet video motion distance: A metric 447
390 for evaluating motion consistency in videos. *arXiv preprint 448*
391 *arXiv:2407.16124*, 2024. 7 449
- [41] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao 392
Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun 393
Zhu, et al. Grounding dino: Marrying dino with grounded 394
pre-training for open-set object detection. *arXiv preprint 395*
arXiv:2303.05499, 2023. 6 396
- [42] Zhenguang Liu, Pengxiang Su, Shuang Wu, Xuanjing Shen, 397
Haipeng Chen, Yanbin Hao, and Meng Wang. Motion predic- 398
tion using trajectory cues. In *Proceedings of the IEEE/CVF 399*
International Conference on Computer Vision (ICCV), pages 400
13299–13308, 2021. 3 401
- [43] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard 402
Pons-Moll, and Michael J Black. Smpl: A skinned multi- 403
person linear model. In *Seminal Graphics Papers: Pushing 404*
the Boundaries, Volume 2, pages 851–866. 2023. 2 405
- [44] Konrad Lorenz. Der kumpan in der umwelt des vogels. 406
der artgenosse als auslösendes moment sozialer verhal- 407
tungsweisen. *Journal für Ornithologie. Beiblatt.(Leipzig)*, 408
1935. 3 409
- [45] Konrad Lorenz and Nikolaas Tinbergen. Taxis und instink- 410
thandlung in der eirollbewegung der graugans. *Zeitschrift für 411*
Tierpsychologie, 1938. 3 412
- [46] Kartikeya Mangalam, Harshayu Girase, Shreyas Agarwal, 413
Kuan-Hui Lee, Ehsan Adeli, Jitendra Malik, and Adrien 414
Gaidon. It is not the journey but the destination: Endpoint 415
conditioned trajectory prediction. In *European conference on 416*
computer vision, pages 759–776. Springer, 2020. 3 417
- [47] Alexander Mathis, Pranav Mamidanna, Kevin M Cury, Taiga 418
Abe, Venkatesh N Murthy, Mackenzie Weygandt Mathis, and 419
Matthias Bethge. Deeplabcut: markerless pose estimation of 420
user-defined body parts with deep learning. *Nature neuro- 421*
science, 21(9):1281–1289, 2018. 3 422
- [48] Seokha Moon, Hyun Woo, Hongbeen Park, Haeji Jung, Reza 423
Mahjourian, Hyung-gun Chi, Hyerin Lim, Sangpil Kim, and 424
Jinkyu Kim. Visiontrap: Vision-augmented trajectory predic- 425
tion guided by textual descriptions. In *European Conference 426*
on Computer Vision, pages 361–379. Springer, 2024. 3 427
- [49] Tanmay Nath, Alexander Mathis, An Chi Chen, Amir Patel, 428
Matthias Bethge, and Mackenzie Weygandt Mathis. Using 429
deeplabcut for 3d markerless pose estimation across species 430
and behaviors. *Nature protocols*, 14(7):2152–2176, 2019. 3 431
- [50] Dantong Niu, Yuvan Sharma, Haoru Xue, Giscard Biamby, 432
Junyi Zhang, Ziteng Ji, Trevor Darrell, and Roei Herzig. Pre- 433
training auto-regressive robotic models with 4d representa- 434
tions. *arXiv preprint arXiv:2502.13142*, 2025. 2 435
- [51] Ian Noronha, Aneeq Chowdhury, Saru Bharti, and Upinder 436
Kaur. Quadforecaster: Diffusion-based quadruped pose pre- 437
diction for animal communication analysis. In *The Thirty- 438*
Ninth Annual Conference on Neural Information Processing 439
Systems workshop: AI for non-human animal communication. 440
3 441
- [52] OpenAI. Sora, 2024. 2 442
- [53] Sergiu Oprea, Pablo Martinez-Gonzalez, Alberto Garcia- 443
Garcia, John Alejandro Castro-Vargas, Sergio Orts-Escolano, 444
Jose Garcia-Rodriguez, and Antonis Argyros. A review on 445
deep learning techniques for video prediction. *IEEE Trans- 446*
actions on Pattern Analysis and Machine Intelligence, 44(6): 447
2806–2826, 2022. 2 448

- 449 [54] William Peebles and Saining Xie. Scalable diffusion models
450 with transformers. In *Proceedings of the IEEE/CVF inter-*
451 *national conference on computer vision*, pages 4195–4205,
452 2023. 3, 5
- 453 [55] Talmo D Pereira, Nathaniel Tabris, Arie Matsliah, David M
454 Turner, Junyu Li, Shruthi Ravindranath, Eleni S Papadoyannis,
455 Edna Normand, David S Deutsch, Z Yan Wang, et al. Sleep: A
456 deep learning system for multi-animal pose tracking. *Nature*
457 *methods*, 19(4):486–495, 2022. 3
- 458 [56] Jernej Polajnar, Elizaveta Kvinikadze, Adam W Harley, and
459 Igor Malenovsky. Wing buzzing as a mechanism for generat-
460 ing vibrational signals in psyllids (hemiptera: Psylloidea).
461 *Insect science*, 31(5):1466–1476, 2024. 3
- 462 [57] Alex Proekt, Jayanth R Banavar, Amos Maritan, and Don-
463 ald W Pfaff. Scale invariance in the dynamics of spontaneous
464 behavior. *Proceedings of the National Academy of Sciences*,
465 109(26):10564–10569, 2012. 3
- 466 [58] MarcAurelio Ranzato, Arthur Szlam, Joan Bruna, Michael
467 Mathieu, Ronan Collobert, and Sumit Chopra. Video (lan-
468 guage) modeling: a baseline for generative models of natural
469 videos. *arXiv preprint arXiv:1412.6604*, 2014. 2
- 470 [59] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu,
471 Chaitanya Ryali, Tengyu Ma, Haiham Khedr, Roman Rädle,
472 Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment any-
473 thing in images and videos. *arXiv preprint arXiv:2408.00714*,
474 2024. 3, 6
- 475 [60] Andrey Rudenko, Luigi Palmieri, Michael Herman, Kris M
476 Kitani, Dariu M Gavrilu, and Kai O Arras. Human motion
477 trajectory prediction: A survey. *The International Journal of*
478 *Robotics Research*, 39(8):895–935, 2020. 3
- 479 [61] Nadine Rüegg, Shashank Tripathi, Konrad Schindler,
480 Michael J Black, and Silvia Zuffi. Bite: Beyond priors for
481 improved three-d dog pose estimation. In *Proceedings of*
482 *the IEEE/CVF Conference on Computer Vision and Pattern*
483 *Recognition*, pages 8867–8876, 2023. 2, 3
- 484 [62] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and
485 Marco Pavone. Trajectron++: Dynamically-feasible trajectory
486 forecasting with heterogeneous data. In *European conference*
487 *on computer vision*, pages 683–700. Springer, 2020. 3
- 488 [63] Tim Salzmann, Hao-Tien Lewis Chiang, Markus Ryll, Dorsa
489 Sadigh, Carolina Parada, and Alex Bewley. Robots that can
490 see: Leveraging human pose for trajectory prediction. *IEEE*
491 *Robotics and Automation Letters*, 8(11):7090–7097, 2023. 3
- 492 [64] Leandro A Scholz, Tessa Mancienne, Sarah J Stednitz,
493 Ethan K Scott, and Conrad CY Lee. Plug-and-play auto-
494 mated behavioral tracking of zebrafish larvae with deeplabcut
495 and sleep: pre-trained networks and datasets of annotated
496 poses. *bioRxiv*, 2025. 3
- 497 [65] Ari Seff, Brian Cera, Dian Chen, Mason Ng, Aurick Zhou,
498 Nigamaa Nayakanti, Khaled S Refaat, Rami Al-Rfou, and
499 Benjamin Sapp. Motionlm: Multi-agent motion forecasting
500 as language modeling. In *Proceedings of the IEEE/CVF*
501 *International Conference on Computer Vision*, pages 8579–
502 8590, 2023. 3
- 503 [66] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico
504 Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc
505 Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Di-
506 nov3. *arXiv preprint arXiv:2508.10104*, 2025. 5
- [67] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising
diffusion implicit models. In *International Conference on*
Learning Representations, 2021. 2, 5
- [68] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudin-
ov. Unsupervised learning of video representations using
lstms. In *International conference on machine learning*, pages
843–852. PMLR, 2015. 2
- [69] Keqiang Sun, Dor Litvak, Yunzhi Zhang, Hongsheng Li, Jia-
jun Wu, and Shangzhe Wu. Ponymation: Learning articulated
3d animal motions from unlabeled online videos. In *European*
Conference on Computer Vision, pages 100–119. Springer,
2024. 2, 3
- [70] Neerja Thakkar, Karttikeya Mangalam, Andrea Bajcsy, and
Jitendra Malik. Adaptive human trajectory prediction via
latent corridors. In *European Conference on Computer Vision*,
pages 297–314. Springer, 2024. 3
- [71] Niko Tinbergen. On aims and methods of ethology. *Zeitschrift*
für tierpsychologie, 20(4):410–433, 1963. 3
- [72] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan
Kautz. Mocogan: Decomposing motion and content for video
generation. In *Proceedings of the IEEE conference on com-
puter vision and pattern recognition*, pages 1526–1535, 2018.
2
- [73] Anirudh Vemula, Katharina Muelling, and Jean Oh. Social
attention: Modeling attention in human crowds. In *2018*
IEEE international Conference on Robotics and Automation
(ICRA), pages 4601–4607. IEEE, 2018. 3
- [74] Karl Von Frisch. *The dancing bees*. A Harvest, 1953. 3
- [75] Jacob C Walker, Pedro Vélez, Luisa Polania Cabrera,
Guangyao Zhou, Rishabh Kabra, Carl Doersch, Maks Ovs-
janikov, João Carreira, and Shiry Ginosar. Generalist fore-
casting with frozen video models via latent diffusion. 2025.
2, 7
- [76] Yaohui Wang, Piotr Bilinski, Francois Bremond, and Antitza
Dantcheva. Imaginator: Conditional spatio-temporal gan for
video generation. In *Proceedings of the IEEE/CVF winter*
conference on applications of computer vision, pages 1160–
1169, 2020. 2
- [77] Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang
Gao, and Pieter Abbeel. Any-point trajectory modeling for
policy learning, 2023. 2, 4, 8
- [78] Shangzhe Wu, Ruining Li, Tomas Jakab, Christian Rupprecht,
and Andrea Vedaldi. Magicpony: Learning articulated 3d
animals in the wild. In *Proceedings of the IEEE/CVF Con-
ference on Computer Vision and Pattern Recognition*, pages
8792–8802, 2023. 2, 3
- [79] Zhen Xing, Qi Dai, Zejia Weng, Zuxuan Wu, and Yu-
Gang Jiang. Aid: Adapting image2video diffusion mod-
els for instruction-guided video prediction. *arXiv preprint*
arXiv:2406.06465, 2024. 2
- [80] Mengda Xu, Zhenjia Xu, Yinghao Xu, Cheng Chi, Gor-
don Wetzstein, Manuela Veloso, and Shuran Song. Flow
as the cross-domain manipulation interface. *arXiv preprint*
arXiv:2407.15208, 2024. 2
- [81] Jiange Yang, Haoyi Zhu, Yating Wang, Gangshan Wu, Tong
He, and Limin Wang. Tra-moe: Learning trajectory predic-
tion model from multiple domains for adaptive policy con-
ditioning. In *Proceedings of the IEEE/CVF Conference on*

- 565 *Computer Vision and Pattern Recognition*, pages 6960–6970,
566 2025. 2
- 567 [82] Siyuan Yang, Lu Zhang, Yu Liu, Zhizhuo Jiang, and You He.
568 Video diffusion models with local-global context guidance.
569 *arXiv preprint arXiv:2306.02562*, 2023. 2
- 570 [83] Shaokai Ye, Anastasiia Filippova, Jessy Lauer, Steffen Schnei-
571 der, Maxime Vidal, Tian Qiu, Alexander Mathis, and Macken-
572 zie Weygandt Mathis. Superanimal pretrained pose estimation
573 models for behavioral analysis. *Nature communications*, 15
574 (1):5165, 2024. 3
- 575 [84] Xi Ye and Guillaume-Alexandre Bilodeau. Stdif: Spatio-
576 temporal diffusion for continuous stochastic video prediction.
577 In *Proceedings of the AAAI Conference on Artificial Intelli-*
578 *gence*, pages 6666–6674, 2024. 2
- 579 [85] Chengbo Yuan, Chuan Wen, Tong Zhang, and Yang Gao. Gen-
580 eral flow as foundation affordance for scalable robot learning.
581 *arXiv preprint arXiv:2401.11439*, 2024. 2
- 582 [86] Artem Zholus, Carl Doersch, Yi Yang, Skanda Koppula, Vior-
583 ica Patraucean, Xu Owen He, Ignacio Rocco, Mehdi SM
584 Sajjadi, Sarath Chandar, and Ross Goroshin. Tapnext: Track-
585 ing any point (tap) as next token prediction. In *Proceedings of*
586 *the IEEE/CVF International Conference on Computer Vision*,
587 pages 9693–9703, 2025. 2
- 588 [87] Silvia Zuffi, Angjoo Kanazawa, David W Jacobs, and
589 Michael J Black. 3d menagerie: Modeling the 3d shape and
590 pose of animals. In *Proceedings of the IEEE conference on*
591 *computer vision and pattern recognition*, pages 6365–6373,
592 2017. 2, 3
- 593 [88] Silvia Zuffi, Angjoo Kanazawa, and Michael J Black. Lions
594 and tigers and bears: Capturing non-rigid, 3d, articulated
595 shape from images. In *Proceedings of the IEEE conference*
596 *on Computer Vision and Pattern Recognition*, pages 3955–
597 3963, 2018. 3
- 598 [89] Silvia Zuffi, Angjoo Kanazawa, Tanya Berger-Wolf, and
599 Michael J Black. Three-d safari: Learning to estimate ze-
600 bra pose, shape, and texture from images” in the wild”. In
601 *Proceedings of the IEEE/CVF International Conference on*
602 *Computer Vision*, pages 5359–5368, 2019. 3
- 603 [90] Silvia Zuffi, Ylva Mellbin, Ci Li, Markus Hoeschle, Hedvig
604 Kjellström, Senya Polikovsky, Elin Hernlund, and Michael J.
605 Black. Varen: Very accurate and realistic equine network.
606 In *Proceedings of the IEEE/CVF Conference on Computer*
607 *Vision and Pattern Recognition (CVPR)*, pages 5374–5383,
608 2024. 2, 3

Forecasting Animal Motion in the Wild

Supplementary Material

609 5. Qualitative Results: Supplementary Video

610 We provide a [supplementary video](#) to further demonstrate
611 the qualitative performance of our method. Each example
612 first displays the ground truth video segment used to extract
613 point tracks via our data processing pipeline, followed by our
614 model’s predicted motions. The four conditioning timesteps
615 (sampled at 15 FPS) are indicated by a grey border, while all
616 subsequent frames are model predictions. Points predicted
617 as occluded by our method are not rendered.

618 The video is organized into the following sections:

619 1) **Diverse Species and Behaviors:** We showcase re-
620 sults across a wide range of behaviors—including walking,
621 mating, eating, fighting, and grooming—and across vari-
622 ous species. Notably, our model demonstrates robust perfor-
623 mance on rare species that are significantly underrepresented
624 in the training set (e.g., fossa at 0.038%, tapir at 0.22%, and
625 the caribou and eskimo dog at 0.025%). For context, even
626 the most frequent species in our dataset (squirrel, giraffe, ele-
627 phant, hamster, and deer) each comprise only approximately
628 3% of the total data.

629 2) **Stochastic Motion Generation:** By varying the ran-
630 dom seed while keeping the input image and motion history
631 fixed, we demonstrate the model’s ability to generate diverse,
632 physically plausible motion trajectories from the same initial
633 context.

634 3) **Controllable Generation via Displacement Vectors:**
635 We illustrate the model’s responsiveness to an optional 2D
636 displacement vector. All results before these were generat-
637 ed without this prompting. Each set of results holds the
638 input and random seed constant, but uses a different 2D
639 displacement vector. The displacement vectors used, where
640 $d = [d_x, d_y]$ is the ground truth displacement, are, from left
641 to right, d , $-d$, $\frac{d}{2}$, and $2d$.

642 4) **Out-of-distribution generalization:** We evaluate our
643 model’s zero-shot capabilities by prompting it with non-
644 mammal animals, humans, and other objects.

645 5) **Baseline Comparisons:** We provide side-by-side visu-
646 alizations against the “Oracle Velocity” (our strongest non-
647 learned baseline) and Track2Act trained on our full dataset.
648 Comparisons with Track2Act use identical random seeds
649 and motion history. Note that Track2Act and oracle velocity
650 cannot handle occlusions, so all points are treated as visible.

651 6) **Comparison with Stable Video Diffusion [6]:** While
652 SVD produces high-quality results for common species
653 (e.g., horses), it often struggles with rare species, frequently
654 “shape-shifting” them into more common animals or fail-
655 ing to capture realistic behavioral patterns. We highlight
656 these failure modes in species such as the hare (0.39%

657 in our training dataset), elk (1.2%), bison (0.89%), and
658 black rhino (0.20%). We specifically use the Stable Dif-
659 fusion XL model available through the interface available at
660 <https://stablediffusionweb.com/>.

661 7) **Data Preprocessing and Camera Stabilization:** We
662 visualize results from our data preprocessing pipeline, show-
663 casing both raw outputs and results after camera stabilization.
664 We observe that while many animals are detected, some are
665 missed; furthermore, while the segmentation masks from
666 VideoSAM are highly accurate, they are not perfect on this
667 challenging data. Crucially, the camera stabilization of point
668 tracks allows us to effectively disentangle animal motion
669 from camera motion.

670 6. Quantitative Results

671 Results comparing our method with baselines can be seen in
672 table 2 and table 3 for the Panthera genus data. Simple base-
673 lines like no-motion and constant-velocity can sometimes
674 perform well on the combined data due to the large amount
675 of low-motion data, but fail for higher motion, and particu-
676 larly for FVMD which accurately scores motion statistics.
677 Interestingly, WHN gives an accurate acceleration distribu-
678 tion despite not being trained on this data, yet fails to esti-
679 mate overall velocity and other statistics well (qualitatively
680 it gives low-motion, jittery predictions that don’t match ani-
681 mal skeletons). ATM and Track2Act, which we retrained on
682 Panthera data, give predictions that are somewhat closer in
683 terms of the final endpoint error and velocity statistics, but
684 actually perform worse in terms of acceleration and point
685 level accuracy, suggesting they learn overall motion but miss
686 motion details, perhaps in part because the overall losses are
687 on displacement rather than velocity. Our method—trained
688 exclusively on Panthera data—substantially outperforms oth-
689 ers in prediction accuracy on every metric. Furthermore, our
690 method can take the true velocity as conditioning to improve
691 results even further, even though for many metrics simply
692 using the oracle velocity provides little boost.

693 Tables 4 and 5 give analogous results for our model (and
694 Track2Act) trained on all species in our dataset. Results
695 follow similar trends overall, but our model trained on the
696 full data is substantially better, e.g. FVMD for high motion
697 examples falls from 84.8 to 49.3 and PWT rises from 20.6
698 to 26.0. This isn’t because the full dataset is easier than the
699 Panthera subset; other baselines actually perform similarly
700 or worse on these metrics. Instead, this suggests that training
701 on the full dataset improves performance due to transfer
702 between species.

Table 2. Quantitative results on **Panthera Data**, distribution level. FD values are multiplied by 10^3 ; Variance values are multiplied by 10^5 ; FVMD values are divided by 10^3 . Best results in **bold**, second best underlined. \uparrow indicates higher is better; \downarrow indicates lower is better.

Selection	Method	FD (V) \downarrow	FD (A) \downarrow	Var (V)	Var (A)	FVMD \downarrow
High motion	GT	-	-	29.5	10.8	-
	No motion	16.6	5.61	0	0	335.406
	Constant vel	7.49	5.61	37.3	0	149.518
	WHN	15.2	3.27	1.37	4.11	247.56
	ATM	6.52	6.18	10	6.99	112.71
	Track2Act	<u>6.32</u>	5.06	8.01	0.446	<u>104.85</u>
	Ours (uncond)	3.71	<u>4.3</u>	12.8	1.02	84.79
	Oracle vel	5.7	5.61	20.9	0	218.73
	Ours (cond)	2.82	4.19	16.6	1.18	79.38
Medium motion	GT	-	-	1.26	0.931	-
	No motion	0.681	0.484	0	0	91.93
	Constant vel	0.822	0.484	0.959	0	54.51
	WHN	0.726	1.16	1.45	4.35	46.47
	ATM	0.494	0.384	0.28	0.475	<u>36.94</u>
	Track2Act	<u>0.421</u>	<u>0.416</u>	0.345	0.04	43.62
	Ours (uncond)	0.405	0.417	0.179	0.027	26.85
	Oracle vel	0.614	0.484	0.0708	0	107.49
	Ours (cond)	0.389	0.414	0.184	0.0285	28.96
Low motion	GT	-	-	0.142	0.173	-
	No motion	0.077	0.09	0	0	46.0
	Constant vel	0.0814	0.09	0.093	0	<u>15.20</u>
	WHN	0.527	1.56	1.44	4.33	19.89
	ATM	0.0746	0.116	0.123	0.254	19.05
	Track2Act	<u>0.0517</u>	0.0731	0.0748	0.0294	24.36
	Ours (uncond)	0.0444	<u>0.0776</u>	0.026	0.00488	7.54
	Oracle vel	0.0679	0.09	0.00643	0	119.12
	Ours (cond)	0.0431	0.0774	0.0258	0.00499	9.95
Combined	GT	-	-	6.93	2.66	-
	No motion	3.77	1.38	0	0	149.53
	Constant vel	1.86	1.38	8.58	0	62.51
	WHN	3.37	<u>1.12</u>	1.43	4.29	86.89
	ATM	1.49	1.4	2.42	1.75	<u>35.50</u>
	Track2Act	<u>1.43</u>	1.21	1.89	0.121	38.44
	Ours (uncond)	0.874	1.05	2.94	0.226	24.82
	Oracle vel	1.43	1.38	4.73	0	118.61
	Ours (cond)	0.679	1.02	3.73	0.262	24.90

Table 3. Quantitative evaluation on **Panthera**, example-level metrics. Best results in **bold**, second best underlined. For non-learned baselines and ATM (single output), we report single-sample metrics; for WHN, Track2Act, and Ours we report best of $K = 5$.

Selection	Method	ADE \downarrow	FDE \downarrow	VMD \downarrow	Avg PWT \uparrow
High motion	No motion	0.211	0.393	5.51	13.22%
	Constant vel	0.193	0.413	<u>4.91</u>	<u>16.73%</u>
	WHN	0.215	0.393	5.82	10.24%
	ATM	0.143	0.262	5.95	16.31%
	Track2Act	<u>0.135</u>	<u>0.245</u>	5.04	16.72%
	Ours (uncond)	0.107	0.209	4.77	20.68%
	Oracle vel	0.082	0.095	6.03	17.05%
	Ours (cond)	0.067	0.097	4.61	27.31%
	Medium motion	No motion	<u>0.022</u>	<u>0.030</u>	4.18
Constant vel		0.044	0.080	4.78	44.33%
WHN		0.032	0.040	5.03	36.42%
ATM		0.025	0.037	4.51	51.24%
Track2Act		0.024	0.032	<u>3.99</u>	53.69%
Ours (uncond)		0.020	0.027	3.82	60.91%
Oracle vel		0.022	0.027	4.37	52.53%
Ours (cond)		0.016	0.019	3.75	63.66%
Low motion		No motion	<u>0.007</u>	<u>0.010</u>	2.71
	Constant vel	0.013	0.024	3.55	70.99%
	WHN	0.022	0.023	4.45	42.40%
	ATM	0.010	0.016	3.29	72.06%
	Track2Act	0.008	0.012	<u>2.70</u>	76.87%
	Ours (uncond)	0.006	0.009	2.57	86.10%
	Oracle vel	0.007	0.009	3.40	82.43%
	Ours (cond)	0.005	0.007	2.56	87.76%
	Combined	No motion	0.076	0.138	4.02
Constant vel		0.079	0.164	4.33	46.16%
WHN		0.086	0.146	5.05	30.43%
ATM		0.057	0.101	4.48	48.39%
Track2Act		<u>0.053</u>	<u>0.092</u>	<u>3.81</u>	51.13%
Ours (uncond)		0.042	0.078	3.62	58.11%
Oracle vel		0.035	0.042	4.51	53.14%
Ours (cond)		0.028	0.039	3.55	61.67%

703

7. Related Work

704

705

706

707

708

709

Pixel Forecasting. When it comes to forecasting visual information, pixels have been the natural choice for several years. Early approaches predicted future pixels deterministically, as a regression problem [53, 58, 68], which is exceedingly challenging, since the problem is ambiguous, and leads to blurry predictions.

710

711

712

713

714

715

716

While GANs [12, 72, 76] and variational models [37] were once promising, many modern approaches use diffusion models [26] which produce sharp videos [22, 25, 27, 79, 82, 84] – and have brought on a creative video revolution [19, 52]. However, training models directly on video is expensive and data-inefficient and models still struggle with hallucinations and basic physical interactions [3, 9, 31].

717

718

719

720

721

Point Track Forecasting. Several works have pushed the frontier in high-quality point-tracking [14, 15, 32, 33, 86], with broad applications across different computer vision tasks. When it comes to forecasting point tracks, the most significant advancements have come from the robotics

domain. Any-point Trajectory Modeling [77] introduced the paradigm of first training a regression model to predict point tracks from an image and language instruction, and learning a robot policy on top of the track prediction model. Several approaches have followed in this direction [4, 11, 18, 50, 80, 81, 85]. These works have explored different architectures for forecasting point tracks such as conditional diffusion transformers [4, 11] and latent diffusion models [80], all with the end-goal of learning good robotic manipulation policies. Similarly, [75] applies DiTs to forecast frozen video encodings along with future decoded point tracks. We draw inspiration from these conditional DiT architectures but focus on a different application, forecasting motion in the complex domain of in-the-wild animal data.

Most recently, [7] showed that point-track forecasting outperforms pixel generation for simple Kubric [20] object motions. Our work provides further evidence that point tracks can be a more data-efficient representation for motion, by expanding their scope to more challenging and non-rigid domain of in-the-wild animal data.

Table 4. Quantitative results on **All Data**, distribution level. FD values are multiplied by 10^3 ; Variance values are multiplied by 10^5 ; FVMD values are divided by 10^3 . Best results in **bold**, second best underlined. \uparrow indicates higher is better; \downarrow indicates lower is better.

Selection	Method	FD (V) \downarrow	FD (A) \downarrow	Var (V)	Var (A)	FVMD \downarrow
High motion	GT	-	-	31.5	8.94	-
	No motion	27.1	7.51	0	0	481.99
	Constant vel	<u>13.7</u>	7.51	23.8	0	210.47
	WHN	25.2	3.19	1.1	3.34	280.77
	Track2Act	14.4	5.54	6.37	1.13	<u>114.30</u>
	Ours (uncond)	8.96	<u>3.74</u>	13.1	1.68	49.30
	Oracle vel	12.1	7.51	19.8	0	326.80
	Ours (cond)	4.86	3.33	28.3	2.14	40.24
Medium motion	GT	-	-	1.32	1.07	-
	No motion	1.14	0.897	0	0	139.91
	Constant vel	1.43	0.897	1.03	0	89.23
	WHN	0.559	0.679	1.21	3.65	<u>33.86</u>
	Track2Act	<u>0.511</u>	<u>0.454</u>	0.193	0.297	43.63
	Ours (uncond)	0.257	0.298	0.396	0.251	12.90
	Oracle vel	1.03	0.897	0.0825	0	163.67
	Ours (cond)	0.197	0.28	0.613	0.314	12.13
Low motion	GT	-	-	0.111	0.157	-
	No motion	0.0957	0.132	0	0	80.13
	Constant vel	0.124	0.132	0.0891	0	34.31
	WHN	0.46	1.39	1.29	3.89	<u>16.68</u>
	Track2Act	<u>0.0652</u>	<u>0.115</u>	0.128	0.254	40.10
	Ours (uncond)	0.016	0.0309	0.0382	0.0365	4.11
	Oracle vel	0.0886	0.132	0.00404	0	212.13
	Ours (cond)	0.0148	0.0304	0.0416	0.0383	4.51
Combined	GT	-	-	5.41	1.82	-
	No motion	4.66	1.53	0	0	204.14
	Constant vel	<u>2.59</u>	1.53	4.11	0	89.77
	WHN	5.34	0.691	1.23	3.7	94.7
	Track2Act	3.2	1.31	1.48	0.454	<u>55.84</u>
	Ours (uncond)	1.96	<u>0.877</u>	2.94	0.453	17.0
	Oracle vel	2.26	1.53	3.15	0	185.62
	Ours (cond)	1.07	0.778	6.26	0.57	14.38

Table 5. Quantitative evaluation on **All Data**, example-level metrics. Best results in **bold**, second best underlined. For non-learned baselines, we report single-sample metrics; for WHN, and ours we report best of $K = 5$.

Selection	Method	ADE \downarrow	FDE \downarrow	VMD \downarrow	Avg PWT \uparrow
High motion	No motion	0.325	0.596	6.50	12.44%
	Constant vel	0.286	0.591	5.02	11.94%
	WHN	0.262	0.538	5.74	11.62%
	Track2Act	<u>0.157</u>	<u>0.332</u>	<u>4.62</u>	<u>18.11%</u>
	Ours (uncond)	0.119	0.275	4.33	26.01%
	Oracle vel	0.110	0.156	7.04	14.70%
	Ours (cond)	0.068	0.103	4.25	31.50%
	Medium motion	No motion	0.032	0.057	5.29
Constant vel		0.068	0.142	5.70	33.28%
WHN		0.035	0.049	4.75	34.49%
Track2Act		<u>0.027</u>	<u>0.044</u>	<u>3.90</u>	46.44%
Ours (uncond)		0.020	0.035	3.57	59.45%
Oracle vel		0.030	0.042	5.73	43.37%
Ours (cond)		0.016	0.021	3.51	63.05%
Low motion		No motion	<u>0.007</u>	<u>0.011</u>	3.44
	Constant vel	0.018	0.034	4.51	65.13%
	WHN	0.023	0.024	4.19	41.93%
	Track2Act	0.013	0.016	<u>2.88</u>	61.17%
	Ours (uncond)	0.005	0.008	2.27	88.48%
	Oracle vel	0.008	0.010	4.54	80.95%
	Ours (cond)	0.004	0.006	2.26	90.19%
	Combined	No motion	0.099	0.180	4.82
Constant vel		0.104	0.215	5.02	41.15%
WHN		0.105	0.200	4.85	29.92%
Track2Act		<u>0.064</u>	<u>0.126</u>	<u>3.73</u>	43.04%
Ours (uncond)		0.046	0.102	3.31	60.01%
Oracle vel		0.042	0.058	5.57	51.74%
Ours (cond)		0.028	0.042	3.26	63.48%

742 Behavioral Forecasting in Computer Vision.

743 Beyond pixels and tracks, there has also been work focusing
744 on forecasting the behavior of intelligent entities as well
745 as their interactions. For example, human trajectory predic-
746 tion has a long history with a variety of approaches [34, 60].
747 For direct trajectory prediction, these range from RNN based
748 approaches [62, 73] to VAEs [46] and GANS [24], leverag-
749 ing generative modeling of future human trajectories. Many
750 recent papers also focus on utilizing scene context [63, 70]
751 for human trajectory forecasting. Behavioral forecasting
752 has also been extensively explored in the context of au-
753 tonomous driving [38, 39, 48, 65]. Relatively few vision
754 papers have focused on forecasting animal motion. Quad-
755 Forecaster [51] predicted the poses of animals in constrained
756 contexts while [42] demonstrated a proof of concept of their
757 approach on fish and mice. In contrast, our approach lever-
758 ages large and diverse datasets and forecasts animal motion
759 on a general level.

760 **Animal Pose, Motion and Behavior.** Ethology, the study of
761 animal behavior, has a long history [44, 45, 71, 74]. Recent
762 advances in computing and machine learning show promise
763 in aiding discoveries – e.g. the emerging field of Computa-

tional Ethology [1] where computer vision and automated
764 motion analysis plays a major role. For example, work such
765 as DeepLabCut [36, 47, 49] and SLEAP [55] have acceler-
766 ated annotating poses of animals in video. Video analysis has
767 aided the ethology of a wide range of animals, from jump-
768 ing plant lice [56], mice [83] and even zebrafish larvae [64].
769 However, as [1] notes, for most of these approaches, humans
770 still need to manually annotate behaviors in training data
771 – which can be subjective due to varied spatial and tempo-
772 ral scales, and limited by human perception and difficulties
773 in discovering new behaviors. Our work leverages massive
774 datasets of unlabeled videos and is a step towards automatic
775 motion understanding of animals.
776

777 There has also been a line of work on reconstructing
778 animal pose in 3D [87, 88], moving towards accurate recon-
779 structions of individual species [61, 78, 89, 90], or creating
780 species-specific models to generate 3D animal motion [69].
781 These works provide insights into individual animal species,
782 but our work focuses on developing an approach that is data-
783 efficient and can generalize to many species including long
784 tail ones.

8. Method Details: Forecasting Point Trajectories with a Diffusion Model

We present a diffusion-based approach for generating animal motion as a sequence of point tracks. Unlike video generation models that predict RGB pixels, our method operates directly on point trajectories. Given a single observation frame and optional conditioning information like motion history or desired velocity, our model generates plausible future trajectories.

8.1. Problem Formulation

We represent the motion of a single animal as a set of N point tracks, where each track describes the 2D trajectory of a single surface point over a time horizon of T timesteps. Formally, we aim to predict a set of tracks $\mathbf{X} \in \mathbb{R}^{T \times N \times 2}$. Each point track $\mathbf{x}_n = [(x_n^1, y_n^1), (x_n^2, y_n^2), \dots, (x_n^T, y_n^T)]$, consists of a sequence of normalized coordinates (x_n^t, y_n^t) where t indexes time. Points may become occluded, in which case we assume the location is unknown: we represent the occlusion state as $\mathbf{O} \in \mathbb{R}^{T \times N}$, where $\mathbf{O}_n^t \in [0, 1]$ indicates that it the n 'th point is visible (1) or occluded (0) at time t .

Our forecasting model learns a conditional generative distribution:

$$p(\mathbf{X}_{T_c+1:T}, \mathbf{O}_{T_c+1:T} | \mathbf{I}, \mathbf{X}_{1:T_c}, \mathbf{O}_{1:T_c}, \mathbf{d})$$

Where \mathbf{I} is the first frame, $\mathbf{X}_{1:T_c}$ and $\mathbf{O}_{1:T_c}$ are the observed conditioning motion history tracks and occlusion states over the first T_c timesteps, and a single optional 2D displacement vector $\mathbf{d} \in \mathbb{R}^2$ describing the average motion of tracks from the last frame: $\mathbf{d} = \sum_{n=1}^N \mathbf{O}_n^T [(x_n^T, y_n^T) - (x_n^1, y_n^1)] / \sum_{n=1}^N \mathbf{O}_n^T$. The model generates future trajectories $\mathbf{X}_{T_c+1:T}$ and occlusion states $\mathbf{O}_{T_c+1:T}$ conditioned on this observed history and the optional conditioning. Because the main challenge is to predict the track positions $\mathbf{X}_{T_c+1:T}$, which are a high-dimensional and continuous value, we draw inspiration from prior work [4] and model distribution with a diffusion process.

Parameterization of the diffusion target. Diffusion involves adding Gaussian noise to the inputs (tracks and occlusions) and training a network to denoise them. While we could directly denoise \mathbf{X} and \mathbf{O} , there are two problems. First, \mathbf{X} has missing values for occluded points (prior work, e.g. [77], assumes there are no missing points, which is untenable for longer horizons). Second, \mathbf{X} values are extremely correlated, and most of the variance is due to the initial point that is tracked rather than due to the motion itself. We therefore reparameterize the tracks to improve training dynamics. Specifically, we construct the diffusion target $\mathbf{Z}_0^{\text{diff}} = \{\gamma \mathbf{V}, \beta \mathbf{O}\}$ where the n 'th row of $\mathbf{V} \in \mathbb{R}^{N \times T \times 2}$ is $[(\dot{x}_n^1, \dot{y}_n^1), (\dot{x}_n^2, \dot{y}_n^2), \dots, (\dot{x}_n^T, \dot{y}_n^T)]$, and γ and β are scaling parameters so the overall variance roughly matches the noise

distribution. Here $\dot{x}_n^t = (x_n^{t+1} - x_n^t)$, and $\dot{y}_n^t = (y_n^{t+1} - y_n^t)$. We interpolate occluded values $\dot{x}_n^t = (x_n^i - x_n^j) / (i - j)$ where i and j are the next and previous visible points (for occluded points at the end of the sequence, which don't have any such j , we simply use 0). We don't do any special preprocessing for the occlusion indicator; even though it's discrete, we find that the model can still denoise to the discrete values provided that they are scaled appropriately.

8.2. Diffusion Process

Following DDPM [26], we define a forward diffusion process that gradually corrupts the diffusion targets $\mathbf{Z}_0^{\text{diff}}$ with Gaussian noise. The forward process over $\tau = 1, 2, \dots, S$ diffusion steps is:

$$q(\mathbf{Z}_\tau^{\text{diff}} | \mathbf{Z}_0^{\text{diff}}) = \mathcal{N}(\mathbf{Z}_\tau^{\text{diff}}; \sqrt{\bar{\alpha}_\tau} \mathbf{Z}_0^{\text{diff}}, (1 - \bar{\alpha}_\tau) \mathbf{I}), \quad (2)$$

where $\bar{\alpha}_\tau = \prod_{s=1}^{\tau} \alpha_s$ with $\alpha_s = 1 - \beta_s$ and $\{\beta_s\}_{s=1}^S$ is a linear noise schedule from $\beta_1 = 0.0001$ to $\beta_S = 0.02$.

Our diffusion model, f_θ , learns to reverse this process by predicting the clean diffusable data $\mathbf{Z}_0^{\text{diff}}$ directly. The training objective minimizes the L1 loss:

$$\mathcal{L} = \mathbb{E}_{\mathbf{Z}_0^{\text{diff}}, \tau, \epsilon} [\|\mathbf{Z}_0^{\text{diff}} - f_\theta(\mathbf{Z}_\tau^{\text{diff}}, \mathbf{Z}_\tau^{\text{cond}}, \tau)\|_1], \quad (3)$$

where $\mathbf{Z}_\tau^{\text{diff}} = \sqrt{\bar{\alpha}_\tau} \mathbf{Z}_0^{\text{diff}} + \sqrt{1 - \bar{\alpha}_\tau} \epsilon$ with $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and $\mathbf{Z}_\tau^{\text{cond}} = \{\mathbf{I}, \mathbf{X}_1, \mathbf{V}_{1:T_c}, \mathbf{O}_{1:T_c}, \mathbf{d}\}$ is conditioning information including the image \mathbf{I} , as well as motion and occlusion history and desired displacement, if available.

Diffusion Transformer Architecture. We now turn to the description of f_θ , which predicts the clean tracks given noisy tracks and conditioning information. We do not assume that tracks are given in any meaningful order or on any grid. However, similar to [4], we note that a transformer model, where each token corresponds to a track, can handle the permutation invariance, as long as we include relevant conditioning information within each token that encodes what the track corresponds to. This design means that the model can easily reason about the full motion forecast for a single point (since everything about a point is encoded within the same point), and yet it can also easily compare and contrast nearby points via attention. It also means that we can make our network is invariant to the input ordering of the tracks.

Figure 2 shows our overall architecture. Each input token corresponds to a full point trajectory; that is, we construct a token for each track before stacking them into a matrix to pass to the transformer. Each token contains all per-track conditioning information: image features, and clean history of conditioning velocities and occlusions $\{\mathbf{I}, \mathbf{X}_1, \mathbf{V}_{1:T_c}, \mathbf{O}_{1:T_c}\}$, as well as the noisy diffusion target for the track. We can then predict the clean data for each track via simple linear projection from the transformer's output.

We construct a token for the n 'th point track in the following way. We start with a visual feature derived from I , the image frame at time $t = 1$. We extract the full bounding box around the animal plus a 50% margin, and compute image features from a frozen DINOv3 [66], which should capture priors about animal parts. We then extract a feature for the track's initial location (x_n^1, y_n^1) using bilinear interpolation. Next, we encode the velocity and occlusion history $(x_n^{1:T_c-1}, y_n^{1:T_c-1}, \mathbf{O}_n^{1:T_c})$; we embed the velocities $x_n^{1:T_c-1}$ and $y_n^{1:T_c-1}$ using a sinusoidal embedding and scale by γ ; we keep the occlusions $\mathbf{O}_n^{1:T_c}$ as scalar and multiply by β . This component of the token is set to zero in the case where the conditioning is not provided. Finally, we add the noisy velocities and occlusion values $\mathbf{Z}_\tau^{\text{diff}} = \{\hat{\mathbf{V}}, \hat{\mathbf{O}}\}$. The full token construction is the concatenation of the clean conditioning DINOv3 features, the clean conditioning velocity history embedding and the occlusion history, the noisy velocities, and the noisy occlusions, along the channel dimension: $\mathbf{Z}_n = [\mathbf{Z}_n^{\text{diff}}, \mathbf{f}_n^{\text{DINO}}, \mathbf{V}_{1:T_c}, \mathbf{O}_{1:T_c}]$.

We project each token to the transformed dimension D_T and add a position encoding. Unlike sequence models, where the added position encoding is derived from the sequence index, we derive our position encoding from the initial location. (x_n^1, y_n^1) . We use a simple sinusoidal position encoding with length D_T and add it to the track token embedding. Finally we apply a standard DiT transformer [54], before linearly projecting the final layer to the dimension of each track in Z^{diff} .

The final conditioning information is global, rather than per-track: the diffusion timestep τ and optionally the desired total displacement d . We embed these values via a linear embedding, zeroing out the embedding for d in the cases where it is not given, and use adaptive layer norm [54] as input directly at each layer of the diffusion model, as is typical for encoding the diffusion timestep in a diffusion transformer.

8.3. Sampling with DDIM

For efficient inference, we use the DDIM sampling algorithm [67], which enables deterministic sampling with fewer steps than the training diffusion process. DDIM defines a non-Markovian forward process that preserves the same marginals $q(\mathbf{Z}_\tau | \mathbf{Z}_0)$ but allows skipping diffusion timesteps during sampling.

Given the model's prediction $\hat{\mathbf{Z}}_0 = f_\theta(\mathbf{Z}_\tau, \tau, \mathbf{d})$ at diffusion timestep τ , we compute the next state $\mathbf{Z}_{\tau-\Delta}$ as:

$$\epsilon_\theta = \frac{\mathbf{Z}_\tau - \sqrt{\bar{\alpha}_\tau} \hat{\mathbf{Z}}_0}{\sqrt{1 - \bar{\alpha}_\tau}}, \quad (4)$$

$$\mathbf{Z}_{\tau-\Delta} = \sqrt{\bar{\alpha}_{\tau-\Delta}} \hat{\mathbf{Z}}_0 + \sqrt{1 - \bar{\alpha}_{\tau-\Delta} - \sigma_\tau^2} \epsilon_\theta + \sigma_\tau \epsilon, \quad (5)$$

where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and $\sigma_\tau = \eta \sqrt{(1 - \bar{\alpha}_{\tau-\Delta}) / (1 - \bar{\alpha}_\tau) \sqrt{1 - \bar{\alpha}_\tau / \bar{\alpha}_{\tau-\Delta}}}$ controls

stochasticity. We use deterministic sampling ($\eta = 0$) with 100 diffusion steps instead of the full 1000 training steps, yielding 10 \times speedup with minimal quality degradation.

After sampling in velocity space, we convert back to absolute coordinates via cumulative summation: $x_n^t = x_n^1 + \sum_{s=1}^{t-1} v_n^{x,s}$ and $y_n^t = y_n^1 + \sum_{s=1}^{t-1} v_n^{y,s}$ for each point n and trajectory time t .

8.4. Implementation Details

Architecture Our model uses a DiT-B configuration with 12 transformer blocks, hidden dimension of 768, and 12 attention heads. $\sigma_v = 12.0$, $\sigma_o = 0.1$.

Training. We train with Adam optimizer with learning rate 5×10^{-4} , cosine annealing schedule with 5-epoch warmup, and batch size 64 distributed across 16 GPUs. We apply gradient clipping with norm 5.0. Training runs for 140 epochs.

Exponential Moving Average (EMA). We maintain an exponential moving average of model parameters with decay $\gamma = 0.9997$, a standard technique in diffusion models that stabilizes sample quality:

$$\theta_{\text{EMA}} \leftarrow \gamma \theta_{\text{EMA}} + (1 - \gamma) \theta. \quad (6)$$

The EMA weights are used for all evaluation and inference.

Data representation. Each training example consists of $N = 320$ point tracks over a trajectory horizon of $T = 32$ timesteps (sampled at 15 FPS), conditioned on the first $T_c = 4$ timesteps. Tracks longer than T timesteps are subsampled with stride 8. Tracks are normalized to $[0, 1]$ image coordinates within the animal's bounding box and stabilized via homography transformation. We handle variable numbers of valid points ($32 \leq N_{\text{valid}} \leq 320$) using attention masking to ignore padded points.

9. Data Processing details

Here we present an in-depth overview of our data processing pipeline, resulting in the MammalMotion dataset.

9.1. Data Filtering

Our pipeline begins with an initial quality filtering stage applied to the full (untrimmed) 539-hour MammalNet [10] dataset. Videos were excluded if they did not meet our minimum requirements for temporal and spatial resolution: a frame rate of at least 29.9 FPS and a total resolution of 200,000 pixels.

We also remove all videos with a low dynamic range. The dynamic range of each video is computed by analyzing the pixel intensity distribution across all frames. For each frame, we first convert the image to grayscale and then calculate the dynamic range ratio using percentile-based thresholds to account for potential outliers. The dynamic range ratio R for a frame is defined as: $R = \frac{P_{99} - P_1}{I_{\text{max}} - I_{\text{min}}}$, where P_{99}

978 and P_1 are the 99th and 1st percentiles of the pixel intensity
979 distribution respectively, and I_{max} and I_{min} are the theoret-
980 ical maximum and minimum intensity values possible for
981 the image’s data type. The final dynamic range measure for
982 a video is computed as the mean of the frame-wise ratios.
983 This metric provides a normalized measure between 0 and
984 1, where values closer to 1 indicate a wider effective dy-
985 namic range in the video content. We removed videos with a
986 dynamic range value below 0.55.

987 After filtering according to the above criteria, we were
988 left with 280 hours of video data.

989 9.2. Shot Detection via Point Tracking

990 After filtering at the video level, we divided the remaining
991 videos into shots. Seeing that popular open-source libraries
992 such as PySceneDetect fail to detect accurate shot boundaries
993 on the difficult animal data, we developed a novel method
994 for detecting shots based on the same point tracker that we
995 used for obtaining point track training data.

996 Our algorithm works as follows: We use point-tracking
997 to identify temporal discontinuities in video sequences that
998 indicate shot boundaries. Our algorithm operates by greed-
999 ily dividing input videos into contiguous segments of up
1000 to 100 frames and systematically analyzing the temporal
1001 consistency of sparse point correspondences within each
1002 segment. For each video segment, the system samples 50
1003 random query points at the first frame. These query points
1004 are then tracked forward in time using BootsTAPIR, which
1005 outputs point trajectories for the whole segment as well as
1006 visibility booleans.

1007 The shot change detection criterion is based on mon-
1008 itoring the percentage of visible points across all frames
1009 within each segment—when the visibility percentage drops
1010 below 6% (less than 3 points are able to be tracked) for any
1011 frame, the algorithm identifies this as a shot change bound-
1012 ary, under the assumption that abrupt scene transitions cause
1013 widespread tracking failures due to the disappearance or sig-
1014 nificant transformation of visual features. When a segment
1015 contains multiple frames below the visibility threshold, only
1016 the earliest is recorded as the boundary. The segment win-
1017 dows then restarts at that boundary frame t' , where new query
1018 points are sampled, allowing subsequent boundaries to be
1019 discovered in successive passes without any post-hoc merg-
1020 ing. When no boundary is detected, the window advances by
1021 100 frames, ensuring complete, gap-free coverage.

1022 Using this algorithm for shot detection has the added
1023 advantage that shots returned are ones where we will be able
1024 to track points.

1025 9.3. Detection and Segmentation

1026 We now get a segmentation of every animal within each shot.
1027 Our pipeline begins with an initial animal detection stage for
1028 each video shot. We employ Grounding-DINO [41] on every

frame, using the text prompt “animal” and a confidence
threshold of 0.35. Any shots without a single successful
detection are discarded from the dataset.

We next identify frames within each shot that can be used
to initialize a video segmenter on every animal in the shot.
To ensure tractability, shots longer than 1000 frames are first
partitioned into 1000-frame segments. We then developed a
multi-stage heuristic to identify a frame where all animals
are clearly visible and spatially distinct.

We first estimate the number of animals in the shot, N ,
by averaging the number of detections across all frames and
rounding to the nearest integer. We then form a candidate
pool of all frames containing exactly N detections. From
this pool, we isolate the top 10% of frames with the lowest
average Intersection over Union (IoU) among their bounding
boxes. This step prioritizes frames where the animals exhibit
minimal overlap. From this refined subset, we select the
single frame with the highest mean detection confidence to
serve as the definitive query frame.

Finally, we initialize VideoSAM [59] with the bounding
boxes from the selected query frame. The resulting segmen-
tation masks are then propagated bi-directionally to cover
the entire shot.

9.4. Point Tracking

Once we have shots with animals segmented and tracked, we
can track points within each animal. As point trackers are
somewhat unreliable over long timeframes, we break each
shot into sub-shots of length up to 8 seconds (240 frames).
For each animal segmentation mask, we sample 500 points
across each sub-shot. To sample each point, we first sample
uniformly in time (random frame indices within the shot).
Then, we sample from the mask.

Our sampling strategy constrains query points to lie
within animal segmentation masks and employs a distance
transform-based weighting scheme to allow for sampling of
thinner structures such as legs, tails, and heads. Specifically,
75% of points are drawn according to an inverse distance
transform distribution. Let $D(\mathbf{p})$ denote the Euclidean dis-
tance transform, i.e. the distance from pixel $\mathbf{p} \in M$ to the
nearest boundary of segmentation mask M . The sampling
probability is: $P(\mathbf{p}) = \frac{1/(D(\mathbf{p})+\epsilon)}{\sum_{\mathbf{q} \in M} 1/(D(\mathbf{q})+\epsilon)}$, where $\epsilon = 10^{-6}$
ensures numerical stability. This assigns higher probability
to pixels closer to mask boundaries, encouraging coverage
of thin structures. The remaining 25% of points are sampled
uniformly within the mask to ensure coverage of interior
regions.

Once query points are sampled, we track across the shot
(up to 8 seconds) using BootsTAPIR [15].

1077 9.5. Camera Stabilization

1078 While the tracked points are faithful to the animal pixels, the
1079 motion of the tracked points in pixel space confounds the
1080 motion of animals and the camera. Therefore, we employ a
1081 stabilization algorithm to disentangle the animal and camera
1082 motion, and train models on "stabilized" point tracks that
1083 only reflect the motion of animals.

1084 Our approach first samples approximately 300 back-
1085 ground points from regions outside dilated animal segmenta-
1086 tion masks, applying a 32-pixel dilation buffer to ensure ade-
1087 quate separation from foreground motion. These background
1088 query points are evenly distributed across video frames
1089 and tracked using BootsTAPIR to establish correspondence
1090 across the temporal sequence. The resulting background
1091 point trajectories are then used to estimate inter-frame cam-
1092 era transformations through a robust RANSAC-based op-
1093 timization process using publicly available code [15] that
1094 estimates a full homography (8 degrees of freedom). The
1095 camera motion estimation employs a reference frame ap-
1096 proach where transformations are computed relative to a
1097 canonical middle frame, with iterative refinement passes to
1098 improve accuracy. To ensure high-quality transformations,
1099 we require a $> 50\%$ average inlier ratio, and for the trans-
1100 formation matrix to be well-conditioned. We fail to stabilize
1101 7% of the data and discard this before training.

1102 Once the homographies for each frame in a shot relative
1103 to a reference frame are computed, we can stabilize the
1104 point tracks at each timestep relative to the start of a time
1105 horizon. This enables us to understand how an animal moves,
1106 irrespective of camera motion.

1107 9.6. Training Example Construction

1108 We construct training examples by selecting an animal and
1109 a particular starting frame t . We extract the input image by
1110 taking a bounding box around the segment and expanding
1111 it by 50% on each side. We transform all other points with
1112 respect to this bounding box using the homographies (i.e.
1113 multiply each point on frame t' by $H_t H_{t'}^{-1}$). We then nor-
1114 malize all coordinates with respect to the first bounding box,
1115 so that $(0, 0)$ corresponds to the upper-left corner and $(1, 1)$
1116 the bottom right.

1117 10. Experimental Setup

1118 10.1. Experimental Dataset

1119 Before processing the data to create MammalMotion, we
1120 filter the full 539-hour MammalNet dataset [10], cutting
1121 it down to 280 hours. Videos were excluded if they failed
1122 to meet minimum requirements for temporal and spatial
1123 resolution or displayed a low dynamic range.

1124 We evaluate our approach on our filtered *all-species*
1125 dataset spanning the entire MammalNet taxonomy, as well
1126 as a *Panthera*-only subset comprising lions, tigers, and leop-

ards. For each configuration, we construct evaluation sets by
1127 randomly sampling from the validation split with different
1128 levels of motion. In the all-species setting, random samples
1129 are also drawn using stratified sampling across species \times
1130 behavior classes to ensure balanced representation of rare cat-
1131 egories. In contrast, the *Panthera*-only setting uses uniform
1132 random sampling due to its more homogeneous taxonomy.
1133 In both cases, we draw even amounts of samples where the
1134 animal averages the following amounts of frame-to-frame
1135 absolute motion: less than half a pixel, half to 1.5 pixels, and
1136 greater than 1.5 pixels. 1137

1138 10.2. Metrics

1139 We evaluate our model's performance using a suite of met-
1140 rics that assess both example-level trajectory accuracy and
1141 distribution-level motion. All metrics are computed on pre-
1142 dicted trajectories compared against our ground truth. 1142

Distribution-Level Motion Statistics: we apply several
1143 metrics to the overall distributions of predicted trajectories. 1144

Fréchet Distance (FD). To assess whether our model cap-
1145 tures the statistical properties of animal motion, we compute
1146 the Fréchet distance [16] between predicted and ground truth
1147 trajectory distributions. It fits multivariate Gaussian distri-
1148 butions to a set of vectors and compares them. We compute
1149 FD on two representations: first-order differences (veloc-
1150 ities), and second-order differences (accelerations), capturing
1151 motion dynamics, and motion smoothness, respectively. Fol-
1152 lowing prior work [75], we restrict this analysis to individual
1153 tracks visible in all predicted frames to ensure complete
1154 motion sequences. 1155

Trajectory Variance. We measure the temporal variance
1156 of predicted trajectories $\text{Var}_{\text{pred}} = \text{Var}(\text{flat}(\mathbf{P}^{\text{pred}}))$ where
1157 $\mathbf{P}^{\text{pred}} \in \mathbb{R}^{N_{\text{samples}} \times T \times 2}$ is the matrix of all predicted track
1158 samples. This captures the diversity and magnitude of motion
1159 in generated trajectories. We report this alongside ground
1160 truth variance Var_{gt} to assess whether the model reproduces
1161 natural motion magnitudes. 1162

Fréchet Video Motion Distance (FVMD). To evaluate tem-
1163 poral coherence, we use the Fréchet Video Motion Distance
1164 (FVMD) [40]. FVMD quantifies the discrepancy between the
1165 distributions of motion feature vectors, where the features
1166 are local histograms of motion orientation and magnitude. 1167

Example-Level Metrics: Since diffusion models are
1168 stochastic, we follow common practice and report best-of- K
1169 metrics by sampling $K = 5$ predictions with different ran-
1170 dom seeds for each test example. For metrics where lower
1171 is better (ADE, FDE, VMD), we compute $\min_k \text{metric}_k$ for
1172 each example, and average across examples. When higher is
1173 better (PWT) we use max. 1174

1175 **Displacement Error (ADE and FDE).** Following standard
1176 protocols, we evaluate trajectory accuracy using Average
1177 Displacement Error (ADE) and Final Displacement Error
1178 (FDE). ADE is the mean squared Euclidean distance between
1179 predicted and ground truth trajectories for all visible points
1180 across the predicted timesteps. FDE measures endpoint ac-
1181 curacy at the terminal timestep T . **Points Within Threshold**
1182 **(PWT).** As established in point tracking literature [14], we
1183 report the fraction of predicted points within pixel-wise dis-
1184 tance thresholds of the ground truth $\delta \in \{1, 2, 4, 8, 16\}$, in
1185 pixel space, where the input bounding boxes are all resized
1186 to (256,256).
1187 **Video Motion Distance (VMD).** This is a straightforward
1188 extension of FVMD to an example-level metric: we compute
1189 the same feature vector used for FVMD for both the sample
1190 and ground-truth, and report the average Euclidean distance.

1191 10.3. Baselines

1192 We first compare our approach against three non-learned
1193 baselines. We then also compare our approach with learned
1194 baselines ATM and Track2Act. All baselines and our model
1195 use $N_{\text{cond}} = 4$ and predict 28 timesteps at 15 FPS.

1196 **No-Motion Baseline.** The simplest prediction strategy, re-
1197 peating the last conditioning position for all future timesteps:
1198 $\hat{\mathbf{p}}_t = \mathbf{p}_{N_{\text{cond}}-1}$ for $t \geq N_{\text{cond}}$.

1199 **Constant Velocity Baseline.** We estimate a per-point-track
1200 velocity from conditioning frames as $\mathbf{v} = (\mathbf{p}_{N_{\text{cond}}-1} -$
1201 $\mathbf{p}_0)/(N_{\text{cond}} - 1)$ and linearly extrapolate future positions:
1202 $\hat{\mathbf{p}}_t = \mathbf{p}_0 + t \cdot \mathbf{v}$. This provides a simple physics-based pre-
1203 dictor assuming constant motion dynamics.

1204 **Oracle Velocity Baseline.** Uses ground truth average veloc-
1205 ity computed from all of the points on the animal, giving
1206 a fair lower bound for the setting of our model that takes
1207 ground-truth displacement.

1208 **What Happens Next (WHN).** WHN [7] aims for general-
1209 purpose point track forecasting, but the model architecture
1210 has a grid constraint that makes it difficult to train on our
1211 non-constrained data. Therefore we apply it zero-shot.

1212 **Any Trajectory Modeling (ATM).** ATM’s [77] Track
1213 Transformer is a regression-based method. Similarly to our
1214 method, it treats each point track over time as a token. It
1215 masks out future timesteps and learns to regress these co-
1216 ordinates. ATM does not handle visibility, regresses on ab-
1217 solute xy-coordinates, and can only predict one plausible
1218 future. We train this baseline using our Panthera subset, using
1219 $N_{\text{cond}} = 4$.

1220 **Track2Act** [4]. Most similar to our model, using a diffu-
1221 sion backbone, point track as tokens, and point conditioning
1222 setup. We use public Track2Act code and diffuse directly
1223 on absolute XY-coordinates, without any positional encod-
1224 ing following the original implementation. We use the[4]’s
1225 learned ResNet visual features integrated through AdaLN,

and use a standard L2 loss. We omit the goal image (unavail- 1226
able in our setting), condition the model solely on the initial 1227
image, and train the model using $N_{\text{cond}} = 4$. 1228