

# Envisioning the Future, One Step at a Time

## Supplementary Material

### A. Additional Implementation Details

We provide more context on implementation details of our main model described in the paper. Please also refer to the supplementary model code, which contains extensive further comments, for reference.

#### A.1. Transformer Block

We implement our transformer [30, 107] blocks primarily following the standard Llama [101, 102]-style block architecture in a similar setup as Baumann et al. [10]. Specifically, we use pre-normalization with RMSNorm [121], omit bias terms in linear layers, and use rotary positional embeddings [97] in an axial setup with scaled cosine similarity attention following Crowson et al. [27]. Our feedforward network setup does not follow Llama’s SwiGLU [93] activation, but instead uses the more classical GELU [46], while still retaining the omission of bias terms. We observed that both choosing GELU with the typical tanh approximation as the activation and omitting the GLU-style [93] gating leads to small speed improvements without significant decreases in quality. Importantly, we implement a fully fused parallel transformer layer, where cross-attention and self-attention are combined into a single attention across both kinds of tokens, and projections are shared between the attention and feedforward network, as described in the main paper.

#### A.2. Posterior Flow Matching Head

Our flow matching posterior head follows similar high-level hyperparameters as Li et al. [60], with three layers of width 1024. Unlike them, we use a standard flow matching [65] objective instead of the DDPM [47] formulation and perform substantial architectural changes to enable efficient sampling. Each block is a standard pre-LayerNorm [4] FFN block with GELU [46] activation.

**Conditioning.** We implement conditioning such that every component can be cached. Typically, conditioning would be implemented with a local, per-layer MLP that projects a conditioning vector into channel scales, shifts, and, optionally, output gating coefficients. This causes a large number of extra kernel launches, which, as this head will perform tens to hundreds of forward passes per AR sampling step, would cause significant wall-clock overhead. Instead, we precompute all scales and shifts centrally once. Additionally, we factorize the conditioning on flow matching time  $\tau$  and the conditioning on the parameters  $\mathbf{z}_t^{(i)}$  additively, such that the time conditioning can be precomputed offline once, and the parameter conditioning can be computed once per sampling loop, further reducing computational overhead. Condition-

ing inside each block is implemented via predicted scale and shift on the output of each pre-LayerNorm [4]. We do not perform output gating.

**Input Value “Scale Cascade”.** For the posterior FM head, we use an input scale cascade to stabilize training when modeling motion. Practically, this is implemented as a logarithmically spaced set of scale coefficients

$$\mathbf{s} = \exp(\text{linspace}(\log(0.1), \log(1e5), \text{num} = 512)), \quad (1)$$

with  $\text{linspace}(\text{min}, \text{max}, \text{num})$  denoting the standard numpy/PyTorch operation, using which the features for the noisy input  $x_\tau$  are computed component-wise as

$$[\tanh(\mathbf{s} \cdot x_{\tau,0}) \parallel \tanh(\mathbf{s} \cdot x_{\tau,1})] \mathbf{W}_{in}^\top, \quad (2)$$

with  $[\cdot \parallel \cdot]$  denoting channelwise concatenation, and  $\mathbf{W}_{in}$  being the output projection to the transformer’s hidden dimension.

**Sampling.** For sampling, we solve the ODE parametrized by the FM head using an Euler solver with uniform spacing of flow matching time  $\tau$ , matching our training setting of sampling from a uniform distribution  $\tau \sim \mathcal{U}[0, 1]$ . Unless specifically noted otherwise, we use 50 sampling steps. During AR sampling, we simply sample one motion sample from the posterior, update the latest position of that trajectory, and then sample the next step defined by the AR factorization, while also conditioning on this new information. This process can be started from partial motion information, initial motion hints (pokes), or no prior motion information.

#### A.3. Hyperparameters

Tab. A provides a comprehensive list of hyperparameters that describe our training setup and model configuration. We train the open-set motion model for 400k steps with a peak learning rate of  $3 \times 10^{-5}$ . We train with a linear learning rate warmup of 5000 steps, after which we apply a linear learning rate schedule. The training setup for the Billiard simulation is similar, but trajectory positions are obtained from the Billiard physics engine [31] and thus represent ground truth motion instead of tracker annotations. Further, we focus on longer-horizon prediction in the Billiard setup. We train the model to predict 50 timesteps, where each timestep corresponds to a  $\Delta t = 0.01$  s interval for 300k iterations.

#### A.4. Training Data

We use three sources of training data for our models.

**Open-set Video Data.** To train our model for open-world motion generation, we source diverse videos from the internet, while ensuring no overlap with our evaluation data. We

Parameter	Value		
	Open-set	Open-set	Billiard
Dataset	Open-Set Videos	Open-Set Video 3→2D	Billiard Simulations
Number of clips	10M	1.5M	—
Tracker	TapNext [123]	V-DPM [98]	Ground-truth
Tracker position seeding	1024 random positions	16,641 grid positions	random ball starting positions
Flow scale	$[-1, 1]$	$[-1, 1]$	$[-1, 1]$
Image size	$512 \times 512$	$512 \times 512$	$512 \times 512$
Training track number	16	16	16
Training timesteps	16	16	50
Batch size	128	128	128
Optimizer	AdamW [67]	AdamW [67]	AdamW [67]
Betas	(0.09, 0.99)	(0.09, 0.99)	(0.09, 0.99)
Peak learning rate	$3 \times 10^{-5}$	$3 \times 10^{-5}$	$3 \times 10^{-5}$
Learning rate schedule	linear decay to $10^{-8}$	linear decay to $10^{-8}$	linear decay to $10^{-8}$
Warm-up steps	5k	5k	5k
Total steps	400k	400k	300k
Precision	bfloat16 AMP	bfloat16 AMP	bfloat16 AMP
Total Parameters	665M	665M	665M
GPUs	16 Nvidia H200	16 Nvidia H200	16 Nvidia H200s
Training Time	20 h	20 h	20 h
Depth	24	24	24
Width	1024	1024	1024
Head dim	128	128	128
Normalization	RMSNorm	RMSNorm	RMSNorm
FFN expand factor	4	4	4
Activation	GELU	GELU	GELU
Positional Encoding	see Sec. 3	see Sec. 3	see Sec. 3
Static scene conditioning	Adaptive Norm [48]	—	—
Denoyer width	1024	1024	1024
Denoyer depth	3	3	3

Table A. Hyperparameters of our main models and training setup.

then apply an off-the-shelf tracker [123] to obtain pseudo ground-truth annotations. For training, we center crop images to square resolution, cropping in both axis slightly to avoid border points for which the tracker commonly fails. We then resize frames to a uniform  $512 \times 512$  resolution.

**Reprojected 3D Data.** Large-scale open-set videos typically suffer from ego camera motion, limiting the interpretability of trajectories. We aim to train a motion model, predicting interpretable static camera trajectories on unconstrained video data for scalability. We thus apply V-DPM [98] a 3D tracker that also estimates camera motion to open-set videos. Then, we reproject tracks into the first camera view, resulting in stabilized trajectories without camera motion interference. We apply the same center crop and resize.

**Billiard Data.** Training data for the Billiard simulation is obtained using a billiard physics simulation [31]. Ball positions and velocity are sampled randomly while ensuring balls do not overlap with other balls or the border. The physics engine produces future positions of balls, which are used as tracks to train the model.

## A.5. Benchmark Construction

To create the OWM dataset, we source 95 permissively licensed videos from Pexels<sup>1</sup> that have been verified to have a static camera and cover a large variety of different kinds of motion from different kinds of entities (e.g., people, vehicles, animals, objects, ...). We prioritize structured or kinematically constrained dynamics (e.g., articulated bodies, rigid object movement) and avoid stochastic or disconnected movement (e.g., excessive background movement, exces-

sively unconstrained motion). We further manually annotate a start frame and select points of interest on moving objects. Ground truth trajectories are obtained with TAPNext [123] and the tracking quality is manually verified.

We complement our dataset with samples from existing solid mechanics benchmarks with known high complexity. For this purpose, we obtain 97 samples from Physics-IQ [69] (subset “solid mechanics”) and 134 samples from Physion [11] (excluding the “Drape” subset because of its focus on soft-body collisions). We manually verify the correctness of motion in the Physion subset, as we observed some examples with unrealistic physical simulation. We, again, manually select starting frames and query points, and verify the correctness of motion annotations for all the additional samples.

## A.6. Metrics

**Open-World Motion Prediction.** For the open-world and physical motion prediction benchmark, we rely on a simple MSE objective between the ground truth trajectory points  $\mathbf{p}_{gt}$  and the predicted trajectory  $\mathbf{p}_{pred}$  by evaluated methods, where  $\mathbf{p}$  is a sequence of  $T$  2D points  $\mathbf{p} \in \mathbb{R}^{T \times 2}$ . The ground truth is obtained by applying TapNext [123] to the full original video. As in a given initial configuration, multiple outcomes could be reasonable, we give each method the chance to produce an ensemble of predictions, whereby the ensemble size is  $N_{ens} = 5$  for the **Best-of-5** setting and  $N_{ens}$  depends on the throughput of each method in the **Best-in-5min** setting. Throughput is calculated using best effort, meaning we utilize optimized implementations and lower-precision calculations when possible.

**Billiard Planning.** We calculate throughput similarly under optimized settings. To calculate the planning accuracy, we use the best action found during rollouts using the principle from Eq. (8). Then, we perform rollouts of the true Billiard simulation using the found action as the initial motion, while all balls except the action ball are initialized as stationary. A selected action is counted as correct if the target ball at least touches or covers the predefined goal position within the allocated time frame. If not the selected action is counted as incorrect. The accuracy is then calculated by dividing the number of correct actions  $N_{correct}$  by the total number of trials  $N_{total}$ .

## A.7. Baselines

**Open-World Baselines.** For the open-world and physics evaluation, we compare against five state-of-the-art video generation models: MAGI-1 4.5B [100], Wan2.2 I2V-A14B [112], CogVideo-X 1.5 5B-I2V [118], SkyReels V2 DF 1.3B 540P [24], and Stable Video Diffusion 1.1 (SVD) [16]. We utilize the implementation provided in the diffusers [109, 114] library for Wan, CogVideo-X, SkyReels,

<sup>1</sup><https://www.pexels.com/>

and SVD. For MAGI, no diffusers implementation is available as of the writing of the paper, therefore we instead adopt the official repository and checkpoint and use the provided 4.5B distill+quant variant. All models except MAGI are run in I2V mode. Thus, they receive the last known image as conditioning and are tasked to simulate the video rollout. As multiple continuations are possible, we sample the Best-out-of-5 and Best-in-5min motion, respectively, giving the models the chance to explore multiple possible outcomes under uncertainty. For MAGI-1, we run the model in video-2-video mode and provide frames preceding the last known frame as hint conditioning. We subsequently apply Tap-Next [123] tracking to generated videos to obtain predicted trajectories, which we use to compute metrics.

**Billiard Baselines.** We compare billiard action search performance against four video generation baselines and two trajectory prediction baselines, which we implement and train from scratch to ensure fair comparison. We match the training setup as closely as possible to the setup for our model.

Video generation models are implemented as image-conditioned spatio-temporal Diffusion Transformers [76]. For efficiency, we utilize latent diffusion [82] and perform diffusion in the latent space of the pretrained VAE from Stable Diffusion-XL [80]. Image-conditioning is achieved by cross-attending to the VAE-produced tokens of the start image. We train four variants of video diffusion models, differing along two axes to cover a variety of previous approaches. Our video diffusion models either use auto-regressive generation or full sequence diffusion. In the former setting, the image conditioning is auto-regressively updated to include the prior  $N_{hist}$  images. The auto-regressive video generation model then generates the single next frame, conditioned on the history of previous images. The full sequence diffusion approach, on the other hand, is conditioned solely on the initial image and generates the full video from a single noise sample  $x_1 \in \mathbb{R}^{T \times H \times W \times C}$ . The models further differ in how they are informed about motion prompts. The *Images to Video* variants receive an additional second conditioning image to which they cross-attend. Note that this is natively supported by AR video generation models, while full sequence diffusion requires modification. Therefore, these models can infer the initial motion from visual cues. The *poke-cond.* models receive the instantaneous flow as an additional conditioning similar to our method. The flow and positions are first embedded using Fourier Embeddings and then passed through a small-scale MLP before being pooled into a fixed-size vector with a linear layer for multiple trajectories. The model is then conditioned on the flow embedding using Adaptive Layer Normalization [4, 48]. We use L-sized DiT backbones [76] for our experiments and train the video diffusion models until convergence.

For the full trajectory diffusion baseline, we ensure a fair

comparison by reusing our motion models’ backbone, but replacing the auto-regressive point-wise diffusion head with a DiT [76]. The training setup and motion model hyperparameters are consistent with our standard setup; however, we ensure that the model always receives only the first step flow.

For the FPT [10] baseline, we utilize the official implementation. Note that all other models predict step-wise motion, while FPT samples future positions in a single step. We align the horizon of the FPT baseline with that of the step-wise models and predict the final position of the balls at the end of the prediction window.

## B. Additional Ablations

In the following, we elaborate further on design choices in our implementation.

### B.1. Number of Function Evaluations

We test the impact of using more evaluations of the denoising flow matching head on the endpoint error (EPE) in the Billiard setting. Results in Tab. B show that our approach yields lower endpoint error with more function evaluations. Beyond 10 function evaluations, the benefits begin to diminish. Therefore, for our main evaluations in Sec. 5 we use 50 evaluations to balance quality and speed.

NFEs	Mean-best-of-5-EPE
1	0.00361
5	0.00143
10	0.00140
25	0.00139
50	<b>0.00138</b>

Table B. **Inference Time Scaling:** Our approach achieves lower End-Point-Error in the Billiard simulation with more function evaluations of the diffusion head.

### B.2. Trajectory ID Embedding

As outlined in Sec. 3 we draw random, (nearly) orthogonal trajectory embeddings  $\text{id}_{\text{traj}}^{(i)} \sim \mathcal{U}(\mathbb{S}^{d-1})$  to indicate trajectory correspondence to the model. Other, more common approaches would be to use no explicit embedding and instead only rely on positional embeddings, or to use a learnable trajectory embedding with a fixed-size codebook.

We compare these options in Tab. C on the Billiard simulation data. We find that our randomized embeddings outperform both learnable embeddings (likely attributable to a reduction in the likelihood of the model learning position-related biases) and the setting with no extra embeddings. Importantly, unlike learnable embeddings, the model is capable of zero-shot trajectory number extrapolation from 16 (the number observed during training) to larger and smaller numbers, with minimal performance degradation.

Traj. Emb.	Num. Traj.	Mean-best-of-5-EPE
No Emb.	8	0.00116
	16	0.00150
	24	0.00277
Learnable	8	0.00112
	16	0.00149
	24	not possible
Ours	8	<b>0.00108</b>
	16	<b>0.00141</b>
	24	<b>0.00263</b>

Table C. **Trajectory ID Embedding:** Our trajectory ID embeddings provide lower end-point-error in billiard simulations and enable zero-shot generalization to *both* increased and reduced number of trajectories.

### B.3. Multi-Step Reasoning

T	$\Delta t$	Num. Steps	Mean-best-of-5-EPE
0.5	0.01	50	<b>0.00141</b>
0.5	0.05	10	<u>0.00999</u>
0.5	0.5	1	0.02823

Table D. **Reasoning in multiple steps.** We compare predicting 0.5 s into the future using models trained with different step sizes. Our standard method integrates 50 steps, while the other models perform fewer steps. Therefore, these models require fewer auto-regressive steps, yet have to model more of the dynamics internally.

Our approach predicts the motion over a short time horizon  $\Delta t$  in one step and then auto-regressively samples movement to predict motion over the entire time horizon  $T$ , thus factorizing the full motion prediction over  $\Delta T$  into a sequence of small-step predictions. In theory, a motion model with infinite capacity should be able to predict the final position of all scene elements in a single step by internally accounting for all potential interactions. However, we argue that predicting step-wise motion is a substantially more practically viable task when not assuming abundant model capacity. We investigate this assumption by comparing model variants in the Billiard setting.

We compare our standard model, predicting  $\Delta t = 0.01$  s into the future per step, against variants predicting over a larger  $\Delta t$ . We perform a 0.5 s rollout (making the largest-step model perform predictions over the full horizon  $T$  in a single step, as in [10]) and evaluate the end-point-error of each model. The results in Tab. D show that multi-step motion prediction improves modeling performance, with overall improved performance for smaller step intervals. The single-step model performs significantly worse than both multi-step variants, mirroring the planning results in Tab. 2. We attribute this failure to the complexity of internally modeling and accounting for all interactions in a large  $\Delta t$  timeframe.

### B.4. Classic Trajectory Forecasting Setting

We explore our approach’s efficacy in classic trajectory forecasting settings in a *zero-shot* setting. We compare on the canonical ETH-UCY [58, 77] benchmark following the setting of Trajectron++ [87]. All baselines are trained exclusively on in-domain data, while we apply our model zero-shot. The baselines directly operate on tracked abstract agents in a 2D top-down view space (obtained from the camera-space tracks via projection), while we operate directly on the original images, as our model uses that as the input. Since the given homographies are not accurate for reprojecting back into the camera space, we manually annotate correspondences and fit homographies, obtaining the equivalent tracks in camera space, which serve as input and output space for our model. ETH generally annotates people’s heads, while UCY seems to rely on people’s feet. This does not matter in a top-down view, as the head will typically be in the same 2D position as the feet, but it matters in camera space. We annotate homographies to follow the ETH convention. Metrics are computed in the original 2D top-down space directly following the baselines.

We show results in Tab. E. Despite not being trained for this setting, our method achieves competitive results with canonical task-specific baselines. This demonstrates that our much more generic approach can still perform well even in specific settings. With additional finetuning on sufficiently large-scale in-domain data, results should further improve significantly.

### B.5. OWM Breakdown

We report metrics for subsets of OWM focusing on specific kinds of motion in Tab. F. Our model is competitive with substantially larger video baselines for all subsets, including intricate multi-agent interactions. In the time constraint setting our method achieves the best results across all subsets as it’s fast inference allows to explore a much larger variety of potential futures.

## C. Additional Qualitative Samples

**Benchmark samples.** We provide qualitative samples from the OWM benchmark (Fig. A), Physics-IQ subset (Fig. B), and Physion subset (Fig. C) with the **Best-out-of-5** motion annotation for our approach and all baseline methods.

Qualitatively, our approach predicts motion that is on par with state-of-the-art video generation approaches in open-world settings, found in the OWM benchmark. Comparing on Physics-IQ [69], video generation approaches tend to predict overly simplified, physically implausible trajectories, whereas our method is able to capture the complexity of real-world physical interactions. For Physion [11], state-of-the-art video generation models hallucinate overly complex motion, whereas our approach is able to capture the rigid



Method	ETH		Hotel		Zara01		Zara02	
	Deterministic	Best-of-20	Deterministic	Best-of-20	Deterministic	Best-of-20	Deterministic	Best-of-20
SocialLSTM [2]	1.09/2.35	–	0.79/1.76	–	<u>0.47/1.00</u>	–	<u>0.56/1.17</u>	–
SocialGAN [39]	–	0.81/1.52	–	0.72/1.61	–	<u>0.34/0.69</u>	–	<u>0.34/0.69</u>
Trajectron [49]	–	0.59/1.14	–	0.35/0.66	–	<u>0.43/0.83</u>	–	<u>0.43/0.83</u>
Trajectron++ [87]	<b>0.71/1.66</b>	<b>0.39/0.83</b>	<b>0.22/0.46</b>	<b>0.12/0.19</b>	<b>0.39/0.77</b>	<b>0.15/0.33</b>	<b>0.23/0.59</b>	<b>0.11/0.25</b>
Ours (zero-shot)	<u>0.81/1.50</u>	<b>0.31/0.80</b>	<u>0.30/0.54</u>	<u>0.17/0.30</u>	0.79/1.75	0.53/1.21	0.58/1.30	0.40/0.91

Table E. **Zero-shot Comparison with Closed-Domain Trajectory Forecasting on ETH-UCY [58, 77].** All numbers (except ours) are sourced from Trajectron++ [87]. Note that some sequences from UCY are missing due to missing RGB videos. Values are (following Trajectron++) (min-){ADE/FDE}, “–” means not reported by Trajectron++. In the “deterministic” setting, we sample from our model once with fixed seed.

Method	Rigidity				Number of Agents				Agents with Free Will				Avg. Rank		Throughput SAMPLES/MIN↑	Params↓
	Rigid		Non-rigid		Single-Agent		Multi-Agent		w/ Free Will		w/o Free Will					
	N=5↓ T=5MIN↓		N=5↓ T=5MIN↓		N=5↓ T=5MIN↓		N=5↓ T=5MIN↓		N=5↓ T=5MIN↓		N=5↓ T=5MIN↓		N=5↓ T=5MIN↓			
MAGI-1	0.032	0.058	0.039	0.069	0.020	0.044	0.048	0.080	0.040	0.066	0.030	0.065	2.00	3.00	0.303	4.5B
Wan2.2 [112]	0.042	DNF	0.038	DNF	0.039	DNF	0.039	DNF	0.036	DNF	0.045	DNF	2.33	DNF	0.141	14B
CogVideo-X 1.5 [118]	0.051	DNF	0.051	DNF	0.041	DNF	0.052	DNF	0.049	DNF	0.054	DNF	4.50	DNF	0.051	5B
SkyReels V2 [24]	0.061	0.071	0.056	0.066	0.048	0.056	0.064	0.075	0.054	0.063	0.065	0.076	5.50	3.00	0.304	1.3B
SVD 1.1 [16]	0.048	0.055	0.057	0.073	0.037	0.053	0.065	0.077	0.060	0.069	0.042	0.064	4.66	3.00	0.714	1.5B
Ours	0.031	0.007	0.039	0.016	0.036	0.008	0.044	0.017	0.037	0.014	0.044	0.011	2.00	1.00	2200	0.6B

Table F. **OWM Subset-wise Metrics.** Breakdown of Tab. 1 results. While orders of magnitude faster, our method is consistently competitive with state-of-the-art video models across not only the overall benchmark, but also when split across multiple properties (rigidity of motion, number of agents, presence of agents with free will).

body physics of the benchmark setting. Therefore, our approach is able to balance complexity and simplicity better than previous approaches making it applicable to a wider range of inference contexts.

**Open-set samples.** We provide samples for a variety of open-set conditioning images, sourced from the internet. Fig. D shows that our approach predicts motion informed by the context provided through the starting image. We provide two examples where we edit the image using nano banana and use the same motion hint. The qualitative samples show that for a person on a trampoline, more bouncy motion is predicted compared to a person jumping on a wooden floor. Further, a ball rolling across a table has a more straight trajectory compared to an egg rolling across the same table.

Fig. E shows samples generated with initial motion hints. The samples show that unrolling hinted trajectories results in consistent motion across entities and realistic long-term behaviour.

In Fig. F we provide samples generated without an initial motion hint. While the trajectories tend to be more simplistic, realistic motion is obtained based solely on the input image. Query points without motion hint are marked in grey.

Fig. G illustrates samples where only a single query point on the object received a hint, and motion for other queries has to be inferred from appearance alone. The results highlight that sampled motion is coherent across objects. Further, our model is able to capture multi-modal behaviour if two outcomes are realistic given the same input. Query points

without motion hint are marked in grey while queries with motion hint are colored black.

**Billiard samples.** We show qualitative predictions from our model trained on billiard simulation data in Fig. H. We show the predicted simulation given an initial impulse in the upper row, and the ground truth simulation overlaid with the prediction in the lower row for each respective sample. Simulation time increases linearly from left to right. Our model is able to accurately predict the ground truth motion, up to stochastic uncertainties.

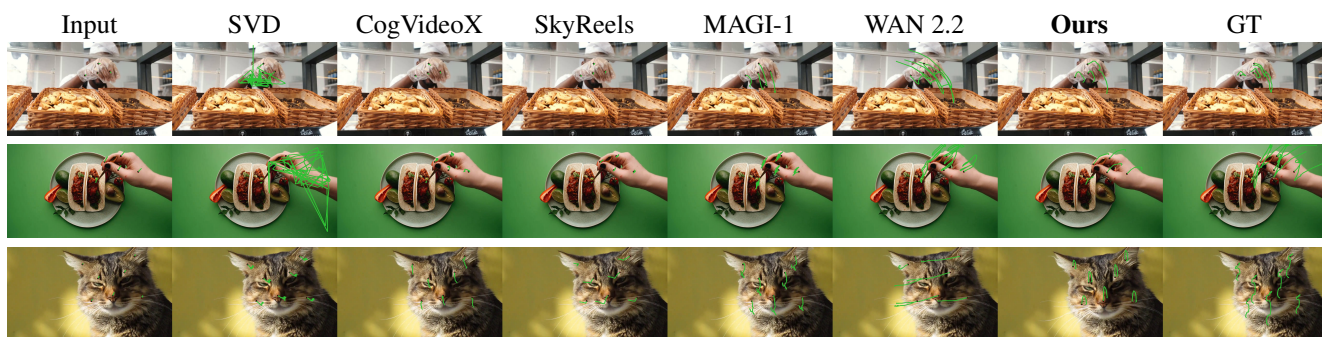


Figure A. **Qualitative comparison on OWM:** Our model produces motion samples that are qualitatively on par with much larger models such as WAN2.2 and MAGI-1.

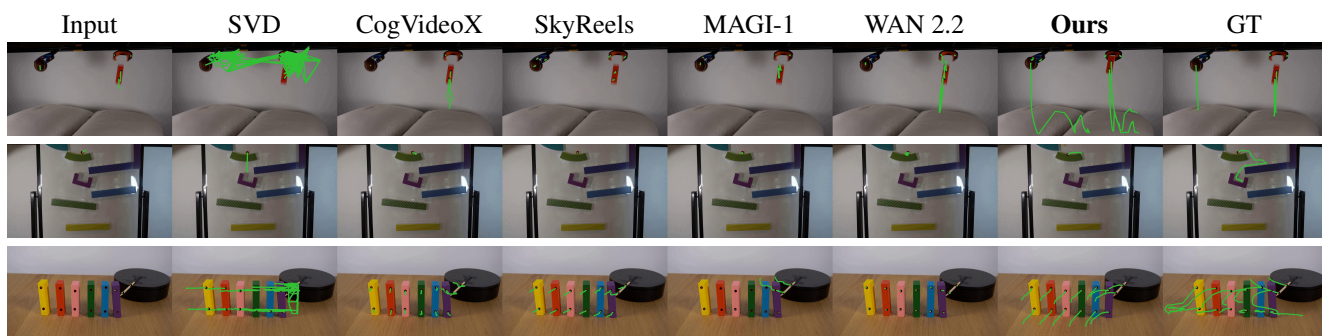


Figure B. **Qualitative comparison on Physics-IQ:** While video generation models fail to capture the complexity of object interactions and predict simplified or no motion, our approach captures realistic physical interactions.

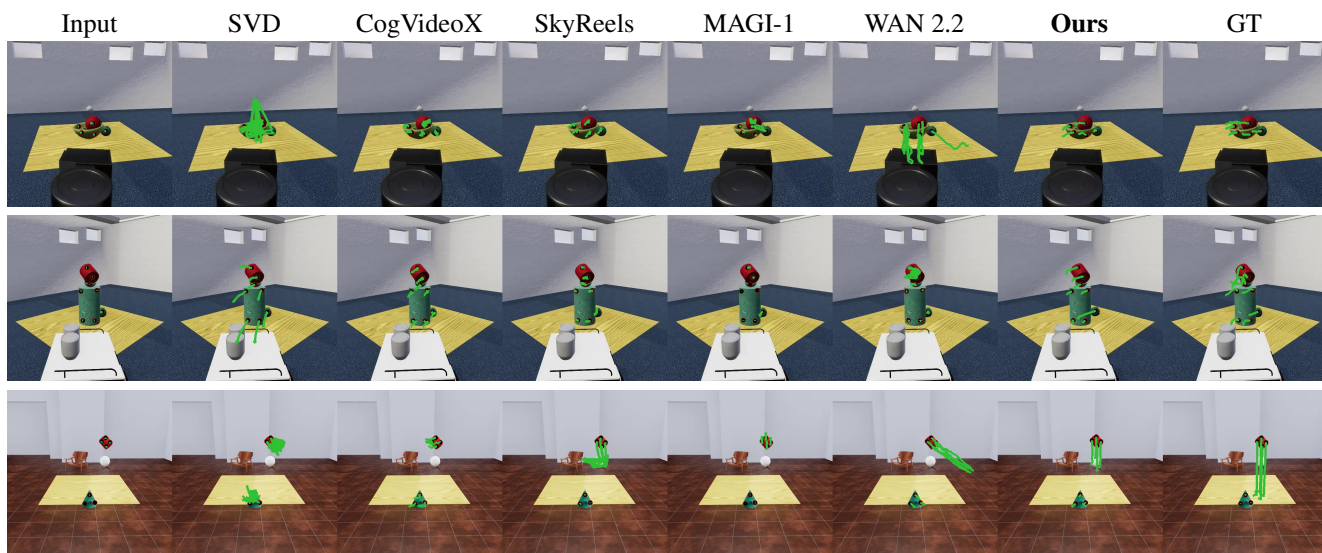


Figure C. **Qualitative comparison on Physion:** For simplified rigid body settings in Physion, video generation models hallucinate overly complex motion, while our approach is able to capture physical dynamics.



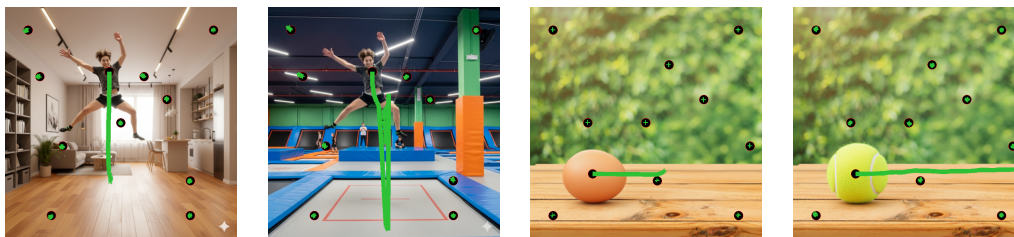


Figure D. **Context informed samples:** The above samples show our model’s ability to take appearance information into account when predicting motion. In both comparisons, images were sampled with the same initial poke. Images where edited using nano banana for high similarity in appearance.



Figure E. **Hinted Samples** Our model is capable of producing complex, coherent, and appearance-informed motion given an initial motion hint.

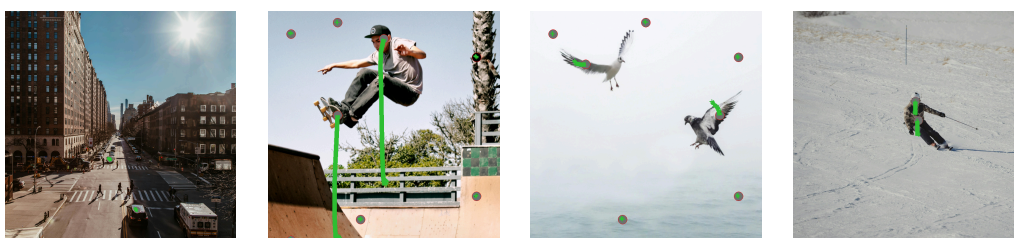


Figure F. **Un-hinted Samples** Given only appearance conditioning, our approach is able to produce physically correct and coherent motion, also showing more complex understanding, such as that cars at an intersection should *not* move when pedestrians are blocking their path.

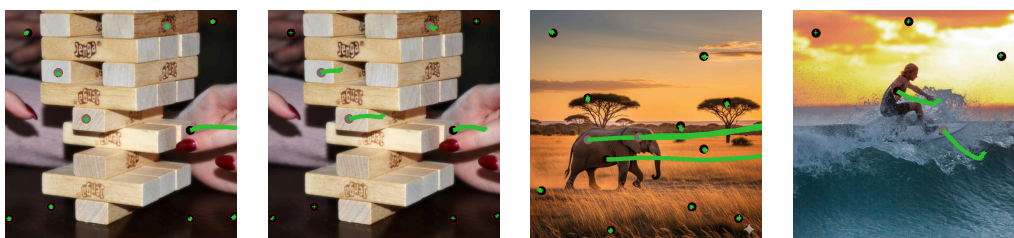


Figure G. **Partially-hinted Samples** Given only a single poke conditioning per example, our model produces coherent motion for queries on the same or linked objects. The Jenga example highlights that our model is able to capture multi-modality if two outcomes are possible given the same initial motion hint.

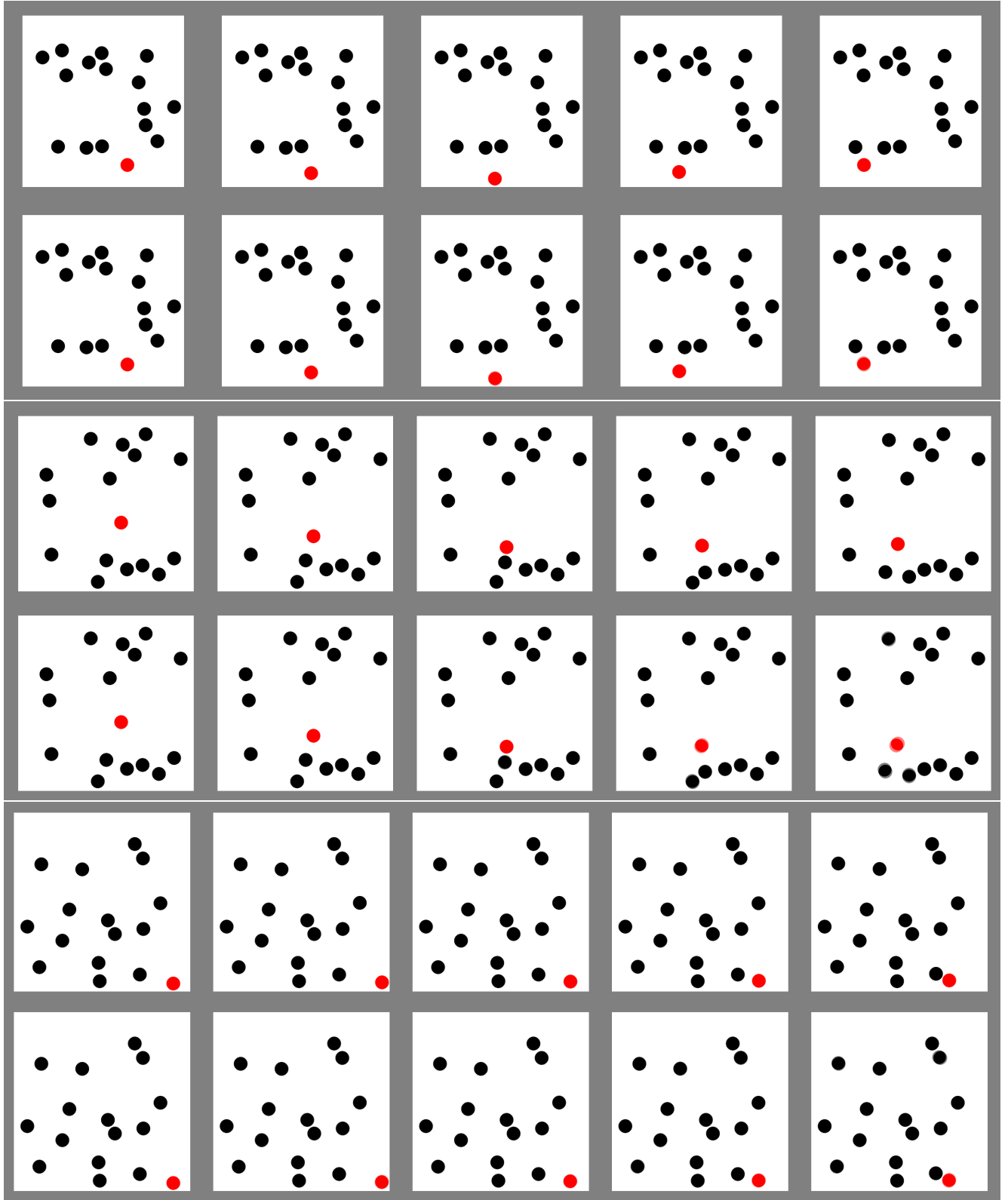


Figure H. **Qualitative samples on our billiard simulation.** The respective top row shows our model's prediction given an initial impulse for the ball marked in red, where we visualize the predicted trajectory state using a frame-wise renderer; the lower row shows an overlay of the ground truth simulation with the prediction to enable comparisons. Our model can successfully predict the observed motion up to minor stochastic details.



## References

- [1] Veo: a text-to-video generation system (veo 3 tech report). Technical report, Google DeepMind, 2025. Technical report. [2](#)
- [2] Alexandre Alahi, Kratharth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–971, 2016. [3](#), [5](#)
- [3] Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos J Storkey, Tim Pearce, and François Fleuret. Diffusion for world modeling: Visual details matter in atari. *Advances in Neural Information Processing Systems*, 37: 58757–58791, 2024. [2](#)
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. [1](#), [3](#)
- [5] Federico Baldassarre, Marc Szafraniec, Basile Terver, Vasil Khalidov, Francisco Massa, Yann LeCun, Patrick Labatut, Maximilian Seitzer, and Piotr Bojanowski. Back to the features: Dino as a foundation for video world models, 2025. [3](#)
- [6] Philip J. Ball, Jakob Bauer, Frank Belletti, Bethanie Brownfield, Ariel Ephrat, Shlomi Fruchter, Agrim Gupta, Kristian Holsheimer, Aleksander Holynski, Jiri Hron, Christos Kaplanis, Marjorie Limont, Matt McGill, Yanko Oliveira, Jack Parker-Holder, Frank Perbet, Guy Scully, Jeremy Shar, Stephen Spencer, Omer Tov, Ruben Villegas, Emma Wang, Jessica Yung, Cip Baetu, Jordi Berbel, David Bridson, Jake Bruce, Gavin Buttmore, Sarah Chakera, Bilva Chandra, Paul Collins, Alex Cullum, Bogdan Damoc, Vibha Dasagi, Maxime Gazeau, Charles Gbadamosi, Woohyun Han, Ed Hirst, Ashyana Kachra, Lucie Kerley, Kristian Kjems, Eva Knoepfel, Vika Koriakin, Jessica Lo, Cong Lu, Zeb Mehring, Alex Moufarek, Henna Nandwani, Valeria Oliveira, Fabio Pardo, Jane Park, Andrew Pierson, Ben Poole, Helen Ran, Tim Salimans, Manuel Sanchez, Igor Saprykin, Amy Shen, Sailesh Sidhwani, Duncan Smith, Joe Stanton, Hamish Tomlinson, Dimple Vijaykumar, Luyu Wang, Piers Wingfield, Nat Wong, Keyang Xu, Christopher Yew, Nick Young, Vadim Zubov, Douglas Eck, Dumitru Erhan, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Raia Hadsell, Aaron van den Oord, Inbar Mosseri, Adrian Bolton, Satinder Singh, and Tim Rocktäschel. Genie 3: A new frontier for world models. 2025. [2](#)
- [7] Federico Barbero, Alex Vitvitskyi, Christos Perivolaropoulos, Razvan Pascanu, and Petar Veličković. Round and round we go! what makes rotary positional encodings useful? In *The Thirteenth International Conference on Learning Representations*, 2025. [4](#)
- [8] Florent Bartoccioni, Elias Ramzi, Victor Besnier, Shashanka Venkataramanan, Tuan-Hung Vu, Yihong Xu, Loick Chambron, Spyros Gidaris, Serkan Odabas, David Hurych, Renaud Marlet, Alexandre Boulch, Mickael Chen, Éloi Zablocki, Andrei Bursuc, Eduardo Valle, and Matthieu Cord. Vavim and vavam: Autonomous driving through video generative modeling, 2025. [2](#)
- [9] Peter W Battaglia, Jessica B Hamrick, and Joshua B Tenenbaum. Simulation as an engine of physical scene understanding. *Proceedings of the national academy of sciences*, 110 (45):18327–18332, 2013. [1](#), [3](#)
- [10] Stefan Andreas Baumann, Nick Stracke, Timy Phan, and Björn Ommer. What if: Understanding motion through sparse interactions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2025. [2](#), [3](#), [4](#), [5](#), [7](#), [8](#), [1](#)
- [11] Daniel Bear, Elias Wang, Damian Mrowca, Felix Jedidja Binder, Hsiao-Yu Tung, RT Pramod, Cameron Holdaway, Sirui Tao, Kevin A. Smith, Fan-Yun Sun, Li Fei-Fei, Nancy Kanwisher, Joshua B. Tenenbaum, Daniel LK Yamins, and Judith E Fan. Physion: Evaluating physical prediction from vision in humans and machines. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. [5](#), [6](#), [7](#), [2](#), [4](#)
- [12] Homanga Bharadhwaj, Roozbeh Mottaghi, Abhinav Gupta, and Shubham Tulsiani. Track2act: Predicting point tracks from internet videos enables generalizable robot manipulation, 2024. [3](#), [7](#)
- [13] Randolph Blake and Maggie Shiffrar. Perception of human motion. *Annu. Rev. Psychol.*, 58(1):47–73, 2007. [2](#)
- [14] Andreas Blattmann, Timo Milbich, Michael Dorkenwald, and Björn Ommer. ipoke: Poking a still image for controlled stochastic video synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14707–14717, 2021. [2](#)
- [15] Andreas Blattmann, Timo Milbich, Michael Dorkenwald, and Björn Ommer. Understanding object dynamics for interactive image-to-video synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5171–5181, 2021. [2](#)
- [16] 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. [2](#), [7](#), [5](#)
- [17] Gabrijel Boduljak, Laurynas Karazija, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. What happens next? anticipating future motion by generating point trajectories. In *The Fourteenth International Conference on Learning Representations*, 2026. [3](#)
- [18] Leo Bringer, Joey Wilson, Kira Barton, and Maani Ghafari. Mdmp: Multi-modal diffusion for supervised motion predictions with uncertainty. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 2889–2899, 2025. [3](#)
- [19] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators, 2024. [2](#), [7](#)
- [20] Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024. [3](#)

- [21] Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In *Conference on Robot Learning*, pages 86–99. PMLR, 2020. 3
- [22] Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in Neural Information Processing Systems*, 37:24081–24125, 2024. 7
- [23] Boyuan Chen, Hanxiao Jiang, Shaowei Liu, Saurabh Gupta, Yunzhu Li, Hao Zhao, and Shenlong Wang. Physgen3d: Crafting a miniature interactive world from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6178–6189, 2025. 3
- [24] Guibin Chen, Dixuan Lin, Jiangping Yang, Chunze Lin, Junchen Zhu, Mingyuan Fan, Hao Zhang, Sheng Chen, Zheng Chen, Chengcheng Ma, Weiming Xiong, Wei Wang, Nuo Pang, Kang Kang, Zhiheng Xu, Yuzhe Jin, Yupeng Liang, Yubing Song, Peng Zhao, Boyuan Xu, Di Qiu, Debang Li, Zhengcong Fei, Yang Li, and Yahui Zhou. Skyreels-v2: Infinite-length film generative model, 2025. 2, 7, 5
- [25] Junsong Chen, Jincheng YU, Chongjian GE, Lewei Yao, Enze Xie, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. 4
- [26] Xinle Cheng, Tianyu He, Jiayi Xu, Junliang Guo, Di He, and Jiang Bian. Playing with transformer at 30+ fps via next-frame diffusion. *arXiv preprint arXiv:2506.01380*, 2025. 2
- [27] Katherine Crowson, Stefan Andreas Baumann, Alex Birch, Tanishq Mathew Abraham, Daniel Z Kaplan, and Enrico Shippole. Scalable high-resolution pixel-space image synthesis with hourglass diffusion transformers. In *Proceedings of the 41st International Conference on Machine Learning*, pages 9550–9575. PMLR, 2024. 4, 1
- [28] Etched Decart, Quinn McIntyre, Spruce Campbell, Xinlei Chen, and Robert Wachen. Oasis: A universe in a transformer. URL: <https://oasis-model.github.io>, 2024. 2
- [29] Michael Dorkenwald, Timo Milbich, Andreas Blattmann, Robin Rombach, Konstantinos G Derpanis, and Bjorn Ommer. Stochastic image-to-video synthesis using cinns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3742–3753, 2021. 2
- [30] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 4, 1
- [31] Markus Ebke. python-billiards. <https://github.com/markus-ebke/python-billiards>, 2025. 6, 7, 1, 2
- [32] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 4
- [33] Lijie Fan, Tianhong Li, Siyang Qin, Yuanzhen Li, Chen Sun, Michael Rubinstein, Deqing Sun, Kaiming He, and Yonglong Tian. Fluid: Scaling autoregressive text-to-image generative models with continuous tokens. In *The Thirteenth International Conference on Learning Representations*, 2025. 5
- [34] Birte U Forstmann, Roger Ratcliff, and E-J Wagenmakers. Sequential sampling models in cognitive neuroscience: Advantages, applications, and extensions. *Annual review of psychology*, 67(1):641–666, 2016. 3
- [35] Katerina Fragkiadaki, Pulkit Agrawal, Sergey Levine, and Jitendra Malik. Learning visual predictive models of physics for playing billiards, 2016. 3
- [36] Jiyang Gao, Chen Sun, Hang Zhao, Yi Shen, Dragomir Anguelov, Congcong Li, and Cordelia Schmid. Vectornet: Encoding hd maps and agent dynamics from vectorized representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11525–11533, 2020. 3
- [37] Ruohan Gao, Bo Xiong, and Kristen Grauman. Im2flow: Motion hallucination from static images for action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5937–5947, 2018. 3
- [38] Junliang Guo, Yang Ye, Tianyu He, Haoyu Wu, Yushu Jiang, Tim Pearce, and Jiang Bian. Mineworld: a real-time and open-source interactive world model on minecraft. *arXiv preprint arXiv:2504.08388*, 2025. 2
- [39] 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, 5
- [40] David Ha and Jürgen Schmidhuber. World models. In *Advances in Neural Information Processing Systems 31*, pages 2451–2463. Curran Associates, Inc., 2018. 3
- [41] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning*, pages 2555–2565. PMLR, 2019.
- [42] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- [43] Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. In *International Conference on Learning Representations*, 2021.
- [44] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse control tasks through world models. *Nature*, 640(8059):647–653, 2025. 3
- [45] Danijar Hafner, Wilson Yan, and Timothy Lillicrap. Training agents inside of scalable world models. *arXiv preprint arXiv:2509.24527*, 2025. 2

- [46] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023. 1
- [47] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1
- [48] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 4, 5, 2, 3
- [49] Boris Ivanovic and Marco Pavone. The trajectron: Probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2375–2384, 2019. 5
- [50] Miguel Jaques, Michael Burke, and Timothy Hospedales. Physics-as-inverse-graphics: Unsupervised physical parameter estimation from video. In *International Conference on Learning Representations*, 2020. 3
- [51] Yuxin Jiang, Shengcong Chen, Siyuan Huang, Liliang Chen, Pengfei Zhou, Yue Liao, Xindong HE, Chiming Liu, Hongsheng Li, Maoqing Yao, and Guanghui Ren. Enerverse-AC: Envisioning embodied environments with action condition. 2025. 2
- [52] Gunnar Johansson. Visual perception of biological motion and a model for its analysis. *Perception & psychophysics*, 14(2):201–211, 1973. 2
- [53] Nikita Karaev, Yuri Makarov, Jianyuan Wang, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Co-tracker3: Simpler and better point tracking by pseudo-labelling real videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6013–6022, 2025. 5
- [54] Efstathios Karypidis, Ioannis Kakogeorgiou, Spyros Gidaris, and Nikos Komodakis. DINO-foresight: Looking into the future with DINO. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. 2, 3
- [55] Georg B Keller and Thomas D Mrsic-Flogel. Predictive processing: a canonical cortical computation. *Neuron*, 100(2):424–435, 2018. 1
- [56] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6
- [57] Kuaishou Technology. Kling: Kuaishou’s proprietary text-to-video generation model. <https://ir.kuaishou.com/news-releases/news-release-details/kuaishou-unveils-proprietary-video-generation-model-kling>, 2024. Press release. 2
- [58] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. Crowds by example. In *Computer graphics forum*, pages 655–664. Wiley Online Library, 2007. 4, 5
- [59] Ruining Li, Chuanxia Zheng, Christian Rupprecht, and Andrea Vedaldi. Puppet-master: Scaling interactive video generation as a motion prior for part-level dynamics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13405–13415, 2025. 2
- [60] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *NeurIPS 2024*, 2024. 4, 5, 1
- [61] Zhengqi Li, Richard Tucker, Noah Snaveley, and Aleksander Holynski. Generative image dynamics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24142–24153, 2024. 3
- [62] Zizhang Li, Hong-Xing Yu, Wei Liu, Yin Yang, Charles Herrmann, Gordon Wetzstein, and Jiajun Wu. Wonderplay: Dynamic 3d scene generation from a single image and actions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9080–9090, 2025. 3
- [63] Jingyun Liang, Yuchen Fan, Kai Zhang, Radu Timofte, Luc Van Gool, and Rakesh Ranjan. Movidio: Motion-aware video generation with diffusion models. In *European Conference on Computer Vision*, 2024. 3, 7
- [64] Ming Liang, Bin Yang, Rui Hu, Yun Chen, Renjie Liao, Song Feng, and Raquel Urtasun. Learning lane graph representations for motion forecasting. In *European Conference on Computer Vision*, pages 541–556. Springer, 2020. 3
- [65] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. 4, 5, 1
- [66] Shaowei Liu, Zhongzheng Ren, Saurabh Gupta, and Shenglong Wang. Physgen: Rigid-body physics-grounded image-to-video generation. In *European Conference on Computer Vision (ECCV)*, 2024. 3
- [67] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 6, 2
- [68] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, pages 405–421, 2020. 4
- [69] Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do generative video models understand physical principles? In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 948–958, 2026. 5, 6, 7, 2, 4
- [70] Roozbeh Mottaghi, Hessam Bagherinezhad, Mohammad Rastegari, and Ali Farhadi. Newtonian scene understanding: Unfolding the dynamics of objects in static images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3521–3529, 2016. 3
- [71] Roozbeh Mottaghi, Mohammad Rastegari, Abhinav Gupta, and Ali Farhadi. “what happens if...” learning to predict the effect of forces in images. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 269–285. Springer, 2016. 3
- [72] Nigamaa Nayakanti, Rami Al-Rfou, Aurick Zhou, Kratharth Goel, Khaled S. Refaat, and Benjamin Sapp. Wayformer: Motion forecasting via simple & efficient attention networks. *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2980–2987, 2022. 3
- [73] Jiquan Ngiam, Vijay Vasudevan, Benjamin Caine, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal,



David J Weiss, Benjamin Sapp, Zhifeng Chen, and Jonathon Shlens. Scene transformer: A unified architecture for predicting future trajectories of multiple agents. In *International Conference on Learning Representations*, 2022. 3

- [74] OpenAI. Sora 2 system card. [https://cdn.openai.com/pdf/50d5973c-c4ff-4c2d-986f-c72b5d0ff069/sora\\_2\\_system\\_card.pdf](https://cdn.openai.com/pdf/50d5973c-c4ff-4c2d-986f-c72b5d0ff069/sora_2_system_card.pdf), 2025. System card. 2
- [75] Jack Parker-Holder, Philip Ball, Jake Bruce, Vibhavari Dasagi, Kristian Holsheimer, Christos Kaplanis, Alexandre Moufarek, Guy Scully, Jeremy Shar, Jimmy Shi, Stephen Spencer, Jessica Yung, Michael Dennis, Sultan Kenjeyev, Shangbang Long, Vlad Mnih, Harris Chan, Maxime Gazeau, Bonnie Li, Fabio Pardo, Luyu Wang, Lei Zhang, Frederic Besse, Tim Harley, Anna Mitenkova, Jane Wang, Jeff Clune, Demis Hassabis, Raia Hadsell, Adrian Bolton, Satinder Singh, and Tim Rocktäschel. Genie 2: A large-scale foundation world model. 2024. 2
- [76] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 3
- [77] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You’ll never walk alone: Modeling social behavior for multi-target tracking. In *2009 IEEE 12th International Conference on Computer Vision*, pages 261–268. IEEE, 2009. 4, 5
- [78] Pika Labs. Pika 2.1. <https://pika.art/faq>, 2025. Product documentation/FAQ. 2
- [79] Silvia L Pinteá, Jan C van Gemert, and Arnold WM Smeulders. Déjà vu: Motion prediction in static images. In *European Conference on Computer Vision*, pages 172–187. Springer, 2014. 3
- [80] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. 3
- [81] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. 4
- [82] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 3
- [83] Pol Rosello. Predicting future optical flow from static video frames. Retrieved on: Jul, 18:2, 2016. 3
- [84] David Ruhe, Jonathan Heek, Tim Salimans, and Emiel Hoogetboom. Rolling diffusion models. In *International Conference on Machine Learning*, pages 42818–42835. PMLR, 2024. 7
- [85] Pablo Ruiz-Ponce, German Barquero, Cristina Palmero, Sergio Escalera, and José García-Rodríguez. Mixermdm: Learnable composition of human motion diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12380–12390, 2025. 3
- [86] Runway Research. Introducing runway gen-4. <https://runwayml.com/research/introducing-runway-gen-4>, 2025. 2
- [87] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data. In *European Conference on Computer Vision*, pages 683–700. Springer, 2020. 3, 4, 5
- [88] Daniel L Schacter, Donna Rose Addis, and Randy L Buckner. Remembering the past to imagine the future: the prospective brain. *Nature reviews neuroscience*, 8(9):657–661, 2007. 1
- [89] Daniel L Schacter, Donna Rose Addis, and Randy L Buckner. Episodic simulation of future events: Concepts, data, and applications. *Annals of the New York Academy of Sciences*, 1124(1):39–60, 2008.
- [90] Daniel L Schacter, Roland G Benoit, and Karl K Szpunar. Episodic future thinking: Mechanisms and functions. *Current opinion in behavioral sciences*, 17:41–50, 2017. 1
- [91] Rami Seid. Lucid v1. <https://ramimo.substack.com/p/lucid-v1-a-world-model-that-does>, 2024. 2
- [92] Martin EP Seligman, Peter Railton, Roy F Baumeister, and Chandra Sripada. Navigating into the future or driven by the past. *Perspectives on psychological science*, 8(2):119–141, 2013. 1
- [93] Noam Shazeer. Glue variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020. 1
- [94] Xiaoyu Shi, Zhaoyang Huang, Fu-Yun Wang, Weikang Bian, Dasong Li, Yi Zhang, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, et al. Motion-i2v: Consistent and controllable image-to-video generation with explicit motion modeling. *SIGGRAPH 2024*, 2024. 3, 7
- [95] Joonghyuk Shin, Daehyeon Choi, and Jaesik Park. Instant-drag: Improving interactivity in drag-based image editing. In *SIGGRAPH Asia 2024 conference papers*, pages 1–10, 2024. 3
- [96] Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darcet, Théo Moutakanni, Leonel Sentana, Claire Roberts, Andrea Vedaldi, Jamie Tolan, John Brandt, Camille Couprie, Julien Mairal, Hervé Jégou, Patrick Labatut, and Piotr Bojanowski. DINOv3, 2025. 6
- [97] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 4, 1
- [98] Edgar Sucar, Eldar Insafutdinov, Zihang Lai, and Andrea Vedaldi. V-dpm: 4d video reconstruction with dynamic point maps, 2026. 6, 2
- [99] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier



- features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. 4
- [100] Hansi Teng, Hongyu Jia, Lei Sun, Lingzhi Li, Maolin Li, Mingqiu Tang, Shuai Han, Tianning Zhang, WQ Zhang, Weifeng Luo, et al. Magi-1: Autoregressive video generation at scale. *arXiv preprint arXiv:2505.13211*, 2025. 2, 7
- [101] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 1
- [102] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 1
- [103] Michael Tschannen, Cian Eastwood, and Fabian Mentzer. Givt: Generative infinite-vocabulary transformers. In *European Conference on Computer Vision*, pages 292–309. Springer, 2024. 5, 7, 8
- [104] Tomer D Ullman, Elizabeth Spelke, Peter Battaglia, and Joshua B Tenenbaum. Mind games: Game engines as an architecture for intuitive physics. *Trends in cognitive sciences*, 21(9):649–665, 2017. 1
- [105] Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. In *The Thirteenth International Conference on Learning Representations*, 2025. 2
- [106] Balakrishnan Varadarajan, Ahmed Hefny, Avikalp Srivastava, Khaled S Refaat, Nigamaa Nayakanti, Andre Cornman, Kan Chen, Bertrand Douillard, Chi Pang Lam, Dragomir Anguelov, et al. Multipath++: Efficient information fusion and trajectory aggregation for behavior prediction. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 7814–7821. IEEE, 2022. 3
- [107] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 3, 6, 1
- [108] Rahul Venkatesh, Honglin Chen, Kevin Feigels, Daniel M Bear, Khaled Jedoui, Klemen Kotar, Felix Binder, Wanhee Lee, Sherry Liu, Kevin A Smith, et al. Understanding physical dynamics with counterfactual world modeling. In *European Conference on Computer Vision*, pages 368–387. Springer, 2024. 2
- [109] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>, 2022. 2
- [110] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Anticipating visual representations from unlabeled video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 98–106, 2016. 3
- [111] Jacob Walker, Abhinav Gupta, and Martial Hebert. Dense optical flow prediction from a static image. In *Proceedings of the IEEE international conference on computer vision*, pages 2443–2451, 2015. 3
- [112] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wenten Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 7, 2, 5
- [113] Ben Wang and Aran Komatsuzaki. Gpt-j-6b: A 6 billion parameter autoregressive language model. <https://github.com/kingoflolz/mesh-transformer-jax>, 2021. 4
- [114] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface’s transformers: State-of-the-art natural language processing, 2020. 2
- [115] Jiajun Wu, Ilker Yildirim, Joseph J Lim, Bill Freeman, and Josh Tenenbaum. Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. *Advances in neural information processing systems*, 28, 2015. 3
- [116] Jiajun Wu, Erika Lu, Pushmeet Kohli, Bill Freeman, and Josh Tenenbaum. Learning to see physics via visual deanimation. *Advances in neural information processing systems*, 30, 2017.
- [117] Tianyi Xie, Zeshun Zong, Yuxing Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. Physgaussian: Physics-integrated 3d gaussians for generative dynamics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4389–4398, 2024. 3
- [118] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Yuxuan Zhang, Weihang Wang, Yean Cheng, Bin Xu, Xiaotao Gu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video diffusion models with an expert transformer. In *The Thirteenth International Conference on Learning Representations*, 2025. 2, 7, 5
- [119] Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human motion diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 3
- [120] Jeffrey M Zacks and Khena M Swallow. Event segmentation. *Current directions in psychological science*, 16(2):80–84, 2007. 3

- [121] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019. [1](#)
- [122] Hang Zhao, Jiyang Gao, Tian Lan, Chen Sun, Ben Sapp, Balakrishnan Varadarajan, Yue Shen, Yi Shen, Yuning Chai, Cordelia Schmid, et al. Tnt: Target-driven trajectory prediction. In *Conference on robot learning*, pages 895–904. PMLR, 2021. [3](#)
- [123] Artem Zholus, Carl Doersch, Yi Yang, Skanda Koppula, Viorica Patraucean, Xu Owen He, Ignacio Rocco, Mehdi SM Sajjadi, Sarath Chandar, and Ross Goroshin. Tapnext: Tracking any point (tap) as next token prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9693–9703, 2025. [5](#), [6](#), [2](#), [3](#)
- [124] Gaoyue Zhou, Hengkai Pan, Yann LeCun, and Lerrel Pinto. Dino-wm: World models on pre-trained visual features enable zero-shot planning. In *Forty-second International Conference on Machine Learning*. [2](#), [3](#)
- [125] Mo Zhou, Jianwei Wang, Xuanmeng Zhang, Dylan Campbell, Kai Wang, Long Yuan, Wenjie Zhang, and Xuemin Lin. Probdiffflow: an efficient learning-free framework for probabilistic single-image optical flow estimation. *Frontiers of Computer Science*, 20(8):2008342, 2026. [3](#)
- [126] Chuning Zhu, Raymond Yu, Siyuan Feng, Benjamin Burchfiel, Paarth Shah, and Abhishek Gupta. Unified world models: Coupling video and action diffusion for pretraining on large robotic datasets. *arXiv preprint arXiv:2504.02792*, 2025. [2](#)