

Captain Safari: A World Engine

Yu-Cheng Chou¹ Xingrui Wang¹ Yitong Li² Jiahao Wang¹ Hanting Liu¹
Cihang Xie³ Alan Yuille¹ Junfei Xiao^{1✉}

¹Johns Hopkins University ²Tsinghua University ³UC Santa Cruz

<https://johnson111788.github.io/open-safari/>

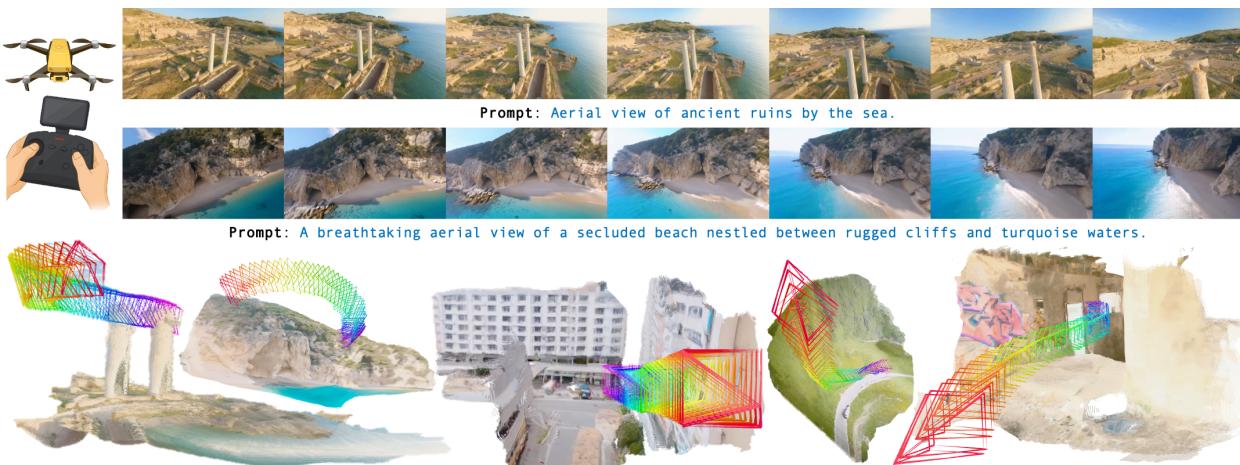


Figure 1. **Captain Safari** is a pose-aware world engine that generates long-horizon, 3D-consistent FPV videos from any user-specified camera trajectory. By retrieving pose-aligned world memory, it keeps geometry stable across large viewpoint changes and reconstructs crisp, well-formed structures while faithfully tracking aggressive 6-DoF motion.

Abstract

World engines aim to synthesize long, 3D-consistent videos that support interactive exploration of a scene under user-controlled camera motion. However, existing systems struggle under aggressive 6-DoF trajectories and complex outdoor layouts: they lose long-range geometric coherence, deviate from the target path, or collapse into overly conservative motion. To this end, we introduce Captain Safari, a pose-conditioned world engine that generates videos by retrieving from a persistent world memory. Given a camera path, our method maintains a dynamic local memory and uses a retriever to fetch pose-aligned world tokens, which then condition video generation along the trajectory. This design enables the model to maintain stable 3D structure while accurately executing challenging camera maneuvers.

To evaluate this setting, we curate OpenSafari, a new *in-the-wild* FPV dataset containing high-dynamic drone videos with verified camera trajectories, constructed through a multi-stage geometric and kinematic validation

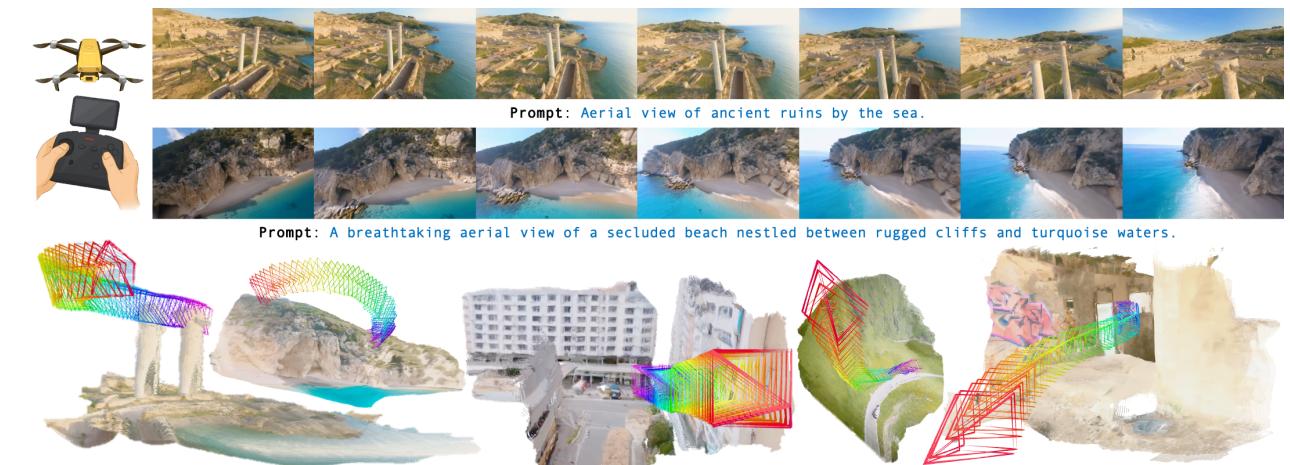
pipeline. Across video quality, 3D consistency, and trajectory following, Captain Safari substantially outperforms state-of-the-art camera-controlled generators. It reduces MET3R from 0.3703 to 0.3690, improves AUC@30 from 0.181 to 0.200, and yields substantially lower FVD than all camera-controlled baselines. More importantly, in a 50-participant, 5-way human study where annotators select the best result among five anonymized models, 67.6% of preferences favor our method across all axes. Our results demonstrate that pose-conditioned world memory is a powerful mechanism for long-horizon, controllable video generation and provide OpenSafari as a challenging new benchmark for future world-engine research.

1. Introduction

Simulating coherent 3D worlds through controllable video generation has long been a foundational challenge for augmented reality, embodied AI, and virtual agents [8–10, 13–16, 19, 20, 23, 26, 29, 40, 43–45, 50, 51, 57]. Classical game engines and physics simulators offer explicit geometry and precise control, but require heavy manual authoring and expensive computation [7, 28, 32]. More-

Captain Safari: Een Wereldmotor

Yu-ChengChou¹ Xingrui Wang¹ Yitong Li² Jiahao Wang¹ Hanting Liu¹CihangXie³ Alan Yuille¹ Junfei Xiao^{1✉}Johns Hopkins Universiteit²Tsinghua University³UCSantaCruz<https://johnson111788.github.io/open-safari/>



Figuur 1. CaptainSafari is een pose-bewuste wereldmotor die lange-horizon, 3D-consistente FPV-video's genereert vanuit elke door de gebruiker gespecificeerde cameratraject. Door pose-uitgelijnd wereldgeheugen op te halen, blijft de geometrie stabiel bij grote veranderingen in het gezichtspunt en worden scherpe, goed gevormde structuren gereconstrueerd terwijl agressieve 6-DoF bewegingen gevolgd.

Samenvatting

Wereldmotoren zijn gericht op het synthetiseren van lange, 3D-consistente video's die interactieve verkenning van een scène ondersteunen onder door de gebruiker gecontroleerde camerabewegingen. Bestaande systemen hebben echter moeite met agressieve 6-DoF trajecten en complexe buitenlayouts: ze verliezen langeafstandsgeometrische samenhang, wijken af van het doelpad of vallen terug in conservatieve bewegingen. Om deze reden introduceren we Captain Safari, een pose-geconditioneerde wereldmotor die video's genereert door te putten uit een persistent wereldgeheugen. Gegeven een camera pad, behoudt onze methode een dynamisch lokaal geheugen en gebruikt een retriever om pose-uitgelijnde wereldtokens op te halen, die vervolgens de videogeneratie langs het traject conditioneren. Dit ontwerp stelt het model in staat om een stabiele 3D-structuur te behouden terwijl uitdagende cameramanoeuvres nauwkeurig worden uitgevoerd.

Om deze instelling te evalueren, hebben we OpenSafari samengesteld, een nieuw FPV-dataset in het wild met hoog-dynamische dronevideo's en geverifieerde cameratrajecten, geconstrueerd door middel van een meertraps geometrische en kinematische validatie.

[✉] Correspondentieauteur: Junfei Xiao (xiaojf97@gmail.com)
Preprint, werk in uitvoering.

pijplijn. Op het gebied van videokwaliteit, 3D consistentie en trajectvolgeling presteert CaptainSafari aanzienlijk beter dan de meest geavanceerde camera-gestuurde generatoren. Het vermindert MET3R van 0,3703 naar 0,3690, verbetert AUC@30 van 0,181 naar 0,200 en levert aanzienlijk lagere FVD op dan alle camera-gestuurde baselines. Belangrijker nog, in een menselijke studie met 50 deelnemers, waarbij annotatoren het beste resultaat kiezen uit vijf geanonimiseerde modellen, geeft 67,6% van de voorkeuren de voorkeur aan onze methode over alle assen. Onze resultaten tonen aan dat pose-geconditioneerde wereldgeheugen een krachtig mechanisme is voor lange-termijn, controleerbare videogeneratie en bieden OpenSafari als een uitdagende nieuwe benchmark voor toekomstig wereld-engine onderzoek.

1. Inleiding

Het simuleren van coherente 3D-werelden door middel van controleerbare videogeneratie is al lang een fundamentele uitdaging voor augmented reality, belichaamde AI en virtuele agenten [8–10, 13–16, 19, 20, 23, 26, 29, 40, 43–45, 50, 51, 57]. Klassieke game-engines en fysica-simulatoren bieden expliciete geometrie en nauwkeurige controle, maar vereisen veel handmatige bewerking en dure berekeningen [7, 28, 32]. Daar-

[✉] Correspondentieauteur: Junfei Xiao (xiaojf97@gmail.com)
Preprint, werk in progress.

over, while they may achieve visual realism in specialized domains, they still fall short in capturing the richness and diversity characteristic of real world, such as natural scenes [27, 35, 58]. In contrast, modern video diffusion models synthesize high-fidelity, diverse videos from text or images, yet typically operate as feed-forward clip generators without persistent world state: *they struggle with long-range 3D consistency, complex trajectory following, and faithful reconstruction of diverse scenes* [18, 24]. In this work, we move toward bridging this gap with *Captain Safari*, a world engine that enables pose-conditioned modeling of 3D-consistent and diverse environments, surpassing the limitations of traditional game engines in terms of generality, diversity, and interactivity.

Contemporary video world models face three intertwined challenges. First, *long-horizon consistency* is limited by the temporal window of context frames; models often “forget” distant scenery or violate spatial coherence, leading to abrupt appearance changes that break the realism and continuity of the generated environment [8, 14, 45]. Second, achieving *complex camera maneuvers under strict 3D consistency* remains difficult: existing pose- or trajectory-conditioned methods typically work well only for slow, near-forward motions [16, 34, 49]. When the path involves fast 6-DoF movement, strong parallax, or sharp turns, models exhibit a trade-off—either dampening motion and restricting viewpoint changes to preserve geometry, or committing to the requested path at the cost of distortions, flicker, and structural drift. Third, current approaches underrepresent *complex outdoor layouts*. Much of the works focuses on structured, constrained settings (e.g., indoor tours, driving scenes, or real-estate videos), and models are seldom stress-tested in in-the-wild FPV scenarios where the camera weaves around buildings, vegetation, and varied terrain with substantial parallax [6, 22, 59, 60]. As a result, methods that look competitive in simplified environments often fail to preserve geometry and appearance when confronted with truly diverse, complex outdoor scenes.

To address these issues, we introduce *Captain Safari*, a pose-aware world engine that explicitly maintains a persistent notion of world state to uphold *long-horizon 3D consistency* across strong parallax. Because storing and propagating a full long-term state is computationally prohibitive, we develop a retrieval mechanism that *selects and aggregates* only the most informative scene cues, thereby providing strong geometric guidance without incurring prohibitive cost. Crucially, this retrieval is *pose-aware*: given the target camera pose, it assembles a pose-aligned world prior that steers the generation process, enabling accurate tracking of *aggressive camera maneuvers* while preserving 3D-consistent structure in complex environments.

Furthermore, to close the gap in *complex outdoor layouts* and *aggressive camera motion*, we curate *OpenSafari*,

a large-scale dataset of high-dynamic FPV drone videos with verified camera poses. Much of the literature targets structured, constrained settings (e.g., indoor tours, driving or real-estate videos), and even outdoor datasets typically feature slow, near-forward motion. In contrast, *OpenSafari* comprises in-the-wild FPV flights that weave around buildings and vegetation across uneven terrain, exhibiting large parallax, rapid 6-DoF maneuvers, and sharp viewpoint changes. Paired with verified camera trajectories, these videos present diverse, cluttered outdoor scenes and long-range motion, challenging models to maintain 3D consistency while faithfully tracking complex maneuvers.

We evaluate *Captain Safari* along three axes: *video quality*, *3D consistency*, and *trajectory following*. Across these criteria, our method consistently outperforms contemporary camera-controlled video generators on *OpenSafari*: Table 1 reports clear gains in 3D consistency and accurate tracking under complex maneuvers, while maintaining strong perceptual quality. Importantly, a large-scale human study (Table 2) shows that *Captain Safari* receives **67%** of votes in five-way comparisons, indicating that the improvements are perceptually salient. Qualitative comparisons (Fig. 4 and Fig. 5) further demonstrate stable geometry under long-range path and faithful adherence to sharp 6-DoF camera turns in cluttered outdoor scenes.

In summary, our contributions are:

1. We present *Captain Safari*, the first camera-controlled video generation method to enforce long-horizon 3D consistency while tracking aggressive FPV maneuvers.
2. We propose a *pose-guided, long-horizon retrieval* that efficiently reconciles strict 3D consistency with accurate tracking of complex maneuvers.
3. We curate *OpenSafari*, a large-scale in-the-wild FPV dataset with verified camera poses, featuring diverse, cluttered outdoor scenes and rapid 6-DoF motion that stress-test geometry-consistent camera control.
4. In *OpenSafari*, our pose-aware retrieval notably improves video quality, 3D consistency, and trajectory alignment, also achieving a **67%** human preference rate.

2. Related Work

2.1. 3D-Consistent World Models

Early image-to-3D approaches reconstruct geometry indirectly via multi-view consistency or implicit fields, but often fail to maintain coherent structure across large view changes [13, 21, 43, 50]. Recent efforts integrate 3D reasoning into the generative process. DiffusionGS [5] injects Gaussian Splatting into the diffusion denoiser, enforcing view consistency and enabling single-stage, scalable 3D generation. GenEx [23] and EvoWorld [40] extends this idea from static reconstruction to dynamic world creation, generating explorable 360° panoramic environments

naast, hoewel ze visueel realisme kunnen bereiken in gespecialiseerde domeinen, schieten ze nog steeds tekort in het vastleggen van de rijkdom en diversiteit die kenmerkend zijn voor de echte wereld, zoals natuurlijke scènes [27, 35, 58]. In tegenstelling hiermee synthetiseren moderne videodiffusiemodellen hoogwaardige, diverse video's van tekst of afbeeldingen, maar functioneren ze meestal als feed-forward clipgeneratoren zonder blijvende wereldstatus: *ze hebben moeite met langeafstands 3D-consistentie, complexe trajectvolg en getrouwe reconstructie van diverse scènes* [18, 24]. In dit werk streven we ernaar deze kloof te overbruggen met *Captain Safari*, een wereldengine die pose-geconditioneerde modellering van 3D-consistente en diverse omgevingen mogelijk maakt, en de beperkingen van traditionele game-engines overtreft op het gebied van algemeenheid, diversiteit en interactiviteit.

Hedendaagse videowereldmodellen staan voor drie onderling verbonden uitdagingen. Ten eerste is *lange-termijn consistentie* beperkt door het temporele venster van contextframes; modellen “vergeten” vaak verre landschappen of schenden ruimtelijke samenhang, wat leidt tot abrupte veranderingen in uiterlijk die de realisme en continuïteit van de gegenererde omgeving doorbreken [8, 14, 45]. Ten tweede blijft het moeilijk om *complexere camerabewegingen onder strikte 3D consistentie* te bereiken: bestaande methoden die afhankelijk zijn van pose- of trajectconditionering werken meestal goed alleen voor langzame, bijna voorwaartse bewegingen [16, 34, 49]. Wanneer het pad snel 6-DoF-beweging, sterke parallax of scherpe bochten omvat, vertonen modellen een afweging—of wel het dempen van beweging en het beperken van veranderingen in het gezichtspunt om de geometrie te behouden, ofwel het volgen van het gevraagde pad ten koste van vervormingen, flikkeringen en structurele afwijkingen. Ten derde worden *complexere buitenindelingen* momenteel ondervertegenwoordigd. Veel van het werk richt zich op gestructureerde, beperkte omgevingen (bijv. indoor tours, rijssituaties of vastgoedvideo's), en modellen worden zelden onderworpen aan stress-tests in wilde FPV-scenario's waar de camera zich tussen gebouwen, vegetatie en gevarieerd terrein met aanzienlijke parallax beweegt [6, 22, 59, 60]. Als gevolg daarvan falen methoden die competitief lijken in vereenvoudigde omgevingen vaak om geometrie en uiterlijk te behouden wanneer ze worden geconfronteerd met echt diverse, complexe buitenomgevingen.

Om deze problemen aan te pakken, introduceren we *Captain Safari*, een pose-bewuste wereldmotor die expliciet een blijvend begrip van de wereldtoestand handhaeft om *langetermijn 3D consistentie* te behouden bij sterke parallax. Omdat het opslaan en doorgeven van een volledige langetermijnstaat computationeel onhaalbaar is, ontwikkelen we een ophaalmechanisme dat *alleen de meest informatieve scène-aanwijzingen selecteert en aggregateert*, waardoor sterke geometrische begeleiding wordt geboden zonder onaantrekkelijke kosten. Cruciaal is dat deze ophaalactie *pose-bewust* is: gegeven de doelcamerapositie, stelt het een pose-uitgelijnde wereldvoorbeelding samen die het generatieproces stuurt, waardoor nauwkeurige tracking van *agressieve camerabewegingen* mogelijk wordt gemaakt, terwijl 3D-consistente structuren in complexe omgevingen behouden blijven.

Bovendien, om de kloof te dichten in *complexere buitenomgevingen* en *agressieve camerabewegingen*, hebben we *OpenSafari* samengesteld,

een grootschalig dataset van hoog-dynamische FPV dronevideo's met geverifieerde cameraposities. Veel van de literatuur richt zich op gestructureerde, beperkte omgevingen (bijv. indoor tours, rij- of vastgoedvideo's), en zelfs buitendatasets bevatten meestal langzame, bijna voorwaartse bewegingen. Daarentegen omvat *OpenSafari* FPV-vluchten in de vrije natuur die zich om gebouwen en vegetatie heen bewegen over oneffen terrein, met grote parallax, snelle 6-DoF manœuvres en scherpe veranderingen in gezichtspunt. In combinatie met geverifieerde cameratrajecten presenteren deze video's diverse, rommelige buitenscènes en langeafstandsbewegingen, wat modellen uitdaagt om 3D-consistentie te behouden terwijl ze complexe manœuvres nauwkeurig volgen.

We evalueren *Captain Safari* langs drie assen: *videokwaliteit*, *3D consistentie*, *entrajectvolg*. Volgens deze criteria presteert onze methode consequent beter dan hedendaagse camerastuurde videogeneratoren op *OpenSafari*: Tabel 1 rapporteert duidelijke verbeteringen in 3D consistentie en nauwkeurige tracking onder complexe manœuvres, terwijl sterke perceptuele kwaliteit behouden blijft. Belangrijk is dat een grootschalige menselijke studie (Tabel 2) aantoon dat *Captain Safari* **67% van de stemmen ontvangt in vergelijkingen met vijf opties**, wat aangeeft dat de verbeteringen perceptueel opvallend zijn. Kwalitatieve vergelijkingen (Fig. 4 en Fig. 5) tonen verder stabiele geometrie aan bij langeafstandspaden en getrouwde nadering van scherpe 6-DoF camerabewegingen in rommelige buitenomgevingen.

Samenvattend zijn onze bijdragen:

1. We presenteren *Captain Safari*, de eerste camerastuurde videogeneratiemethode die lange-termijn 3D consistentie afdwingt terwijl agressieve FPV-bewegingen worden gevolgd.
2. We stellen een *pose-geleide, lange-termijn retrieval* voor die efficiënt strikte 3D consistentie verzoent met nauwkeurige tracking van complexe manœuvres.
3. We stellen *OpenSafari* samen, een grootschalige in-the-wild FPV-dataset met geverifieerde cameraposities, met diverse, rommelige buitenscènes en snelle 6-DoF bewegingen die de geometrie-consistente camerabesturing op de proef stellen.
4. In *OpenSafari* verbetert onze pose-bewuste retrieval opmerkelijk de videokwaliteit, 3D consistentie en trajectuitlijning, en behaalt ook een **67%** menselijke voorkeurscore.

2. Gerelateerd Werk

2.1. 3D-consistente wereldmodellen

Vroege benaderingen voor beeld-naar-3D reconstrueren geometrie indirect via multi-view consistentie of impliciete velden, maar slagen er vaak niet in om een coherente structuur te behouden bij grote veranderingen in het zicht [13, 21, 43, 50]. Recente inspanningen integreren 3D-redenering in het generatieve proces. DiffusionGS [5] injecteert Gaussian S platting in de diffusie denoiser, waardoor zichtconsistentie wordt afgedwongen en schaalbare 3D-generatie in één stap mogelijk wordt. GenEx [23] en EvoWorld [40] breiden dit idee uit van statische reconstructie naar dynamische wereldcreatie, waarbij verkenbare 360° panoramische omgevingen worden gegenereerd.

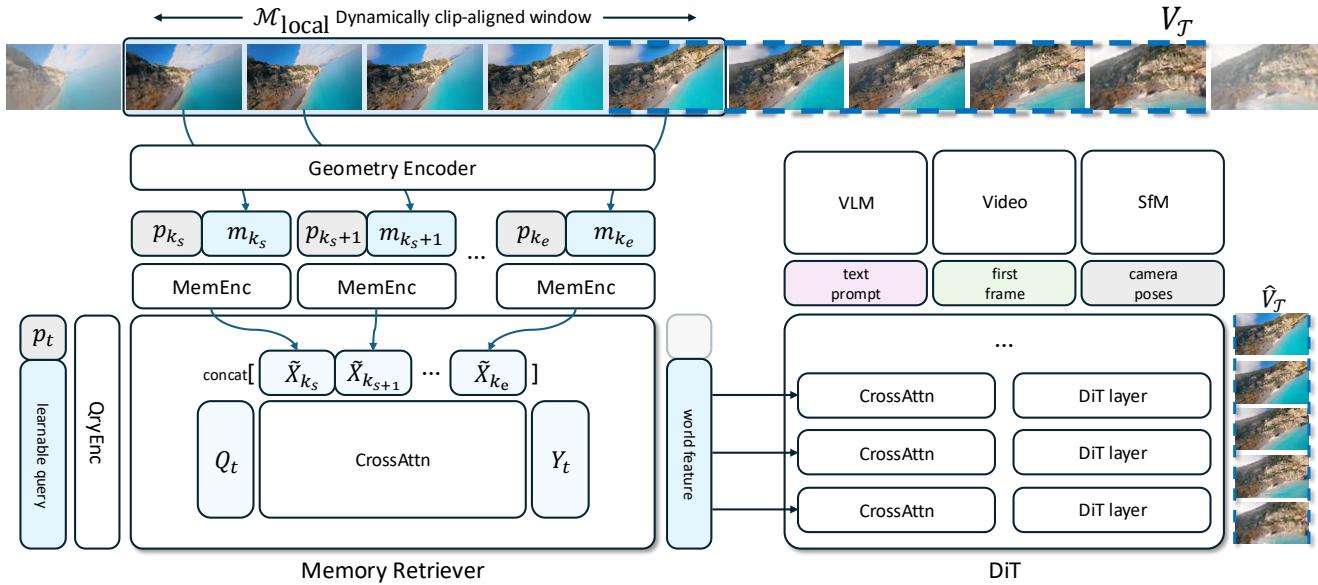


Figure 2. **Method overview.** *Captain Safari* builds a local world memory and, given a query camera pose, retrieves pose-aligned tokens that summarize the scene. These tokens then condition video generation along the user-specified trajectory, preserving a stable 3D layout.

grounded in physical priors. Complementary to these generative reconstructions, Geometry Forcing [44] and Memory Forcing [14] explicitly couple training signals with geometric supervision and spatio-temporal memory, ensuring consistency during long rollouts. Meanwhile, open-world models such as Wonderland [20], WonderWorld [51], Wonder-Turbo [26], and EvoWorld [40] further integrate geometry-indexed or adaptive memories to maintain persistent world states across interactions. However, these approaches still use implicit, clip-bound memories, whereas we introduce an explicit pose-indexed world memory retrieved on demand for camera-controlled generation.

2.2. Camera Controlled Video Generation

Early T2V/I2V models learned camera motion implicitly and struggle to reliably repeat explicit trajectories [11, 42, 61]. Recent work such as CameraCtrl [10] treats camera parameters as explicit conditions, encoding camera extrinsics and trajectories or enforcing path constraints—to improve controllability and accuracy [2, 25, 41, 47, 55]. Motion-Prompting [9] implement compositional control by point-track conditioning and MotionPro [56] use path-alignment losses that lower rotational and translational error; training-free control is also achieved by fitting a lightweight point-cloud and using a noise-layout prior to steer denoising [11]. Scene-preserving geometric priors further strengthen clip-level consistency. Cami2V[57] treats camera pose as a physical prior and exploits epipolar and multiview constraints; RealCam-I2V [19] recovers metric depth with DepthAnything v2 [48] to reconstruct a scale-stable scene; PoseTraj[16] employs pose-aware pretraining to obtain

rotation-aligned motion. Compared with parameter-only conditioning, these priors reduce within-clip layout drift and better preserve local geometry under view changes. Further, recent work links camera control with world modeling. CVD [17], Cavia [46], and WoVoGen [22] jointly synthesize multi-view and multi-trajectory videos from a shared scene representation, enforcing cross-path consistency. Meanwhile, methods that conditioning from explicit renderable 3D representations (e.g., 3D Gaussians) can anchor geometry, improve cross-view 3D consistency and path adherence [15, 29, 33, 52, 53]. However, these approaches typically build one-off 3D scenes, whereas we unify long-horizon camera control with a persistent pose-indexed world memory shared across trajectories.

3. Captain Safari

We introduce *Captain Safari*, a memory-guided video generation framework. Sec. 3.1 presents an implicit world memory for stable scene representation, while Sec. 3.2 describes a pose-conditioned retrieval system that maps camera views to world tokens, guiding a DiT-based generator for coherent outputs along arbitrary trajectories.

3.1. Implicit Memory of World Geometry

Problem setup. We represent a video as $V = \{I_t\}_{t=0}^T$, where I_t is the frame at time step t . On the same time axis we define camera poses $\mathcal{C} = \{(R_t, T_t)\}_{t=0}^T$ and obtain a 3D-aware memory feature m_t at each time step t using a pretrained geometry encoder. All memory features form a global memory bank $\mathcal{M} = \{m_t\}_{t=0}^T$.

Given a text prompt p , the camera poses \mathcal{C} , and a target

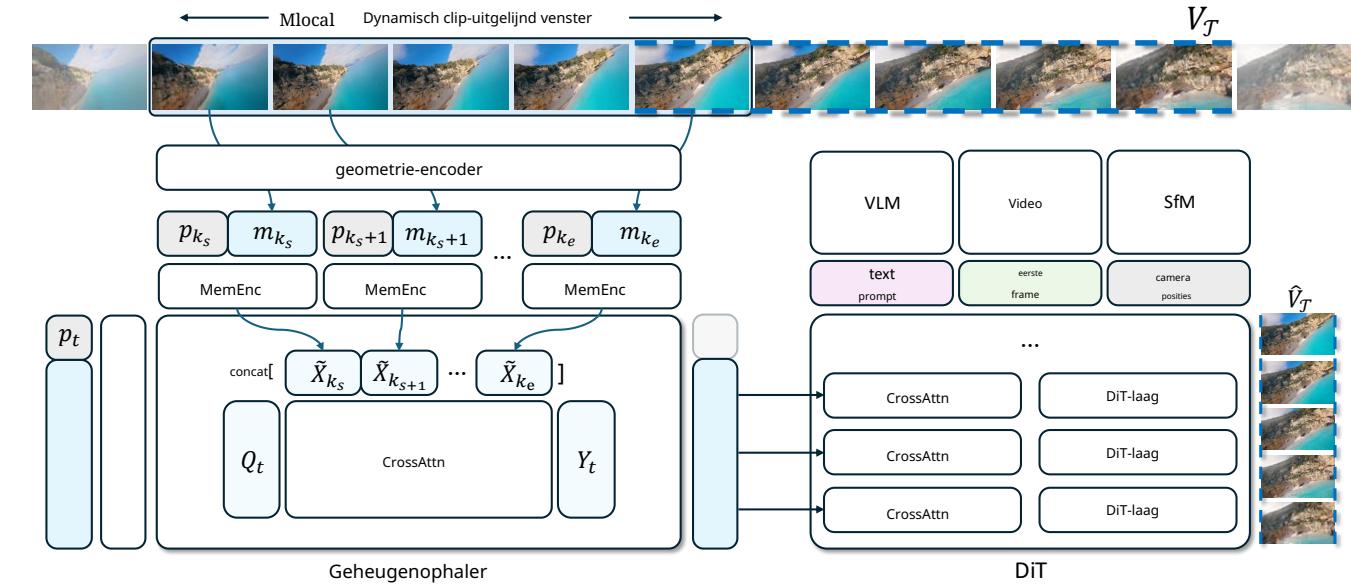


Figure 2. **Method overview.** *Captain Safari* bouwt een lokaal wereldgeheugen op en haalt, gegeven een query camerapositie, pose-uitgelijnde tokens op die de scène sa menvatten. Deze tokens conditioneren vervolgens de videogeneratie langs het door de gebruiker gespecificeerde traject, waarbij een stabiele 3D-indeling behouden blijft.

gebaseerd op fysieke aannames. Complementair aan deze generatieve reconstructies, Geometrie Dwang [44] en Geheugenforcing [14] koppelen trainingssignalen expliciet met geometrische supervisie en spatio-temporeel geheugen, wat consistentie tijdens lange uitrol garandeert. Ondertussen integreren open-wereldmodellen zoals Wonderland [20], Wonderworld [51], Wonder-Turbo [26], en EvoWorld [40] verder geometrie-geïndexeerde of adaptieve geheugens om blijvende wereldstaten te behouden tijdens interacties. Deze benaderingen gebruiken echter nog steeds impliciete, clip-gebonden geheugens, terwijl wij een expliciet pose-geïndexeerd wereldgeheugen introduceren dat op aanvraag wordt opgevraagd voor camera-gestuurde generatie.

2.2. Cameragestuurd Videogeneratie

Vroege T2V/I2V-modellen leerden camerabeweging impliciet en hebben moeite om expliciete trajecten betrouwbaar te herhalen [11, 42, 61]. Recent werk zoals CameraCtrl [10] behandelt cameraparameters als expliciete voorwaarden, waarbij camera-extrinsieken en trajecten worden gecodeerd of padbeperkingen worden afgedwongen—om de bestuurbaarheid en nauwkeurigheid te verbeteren [2, 25, 41, 47, 55]. Motion-Prompting [9] implementeert compositiecontrole door punt-track conditioning en MotionPro [56] gebruikt pad-uitlijningsverliezen die de rotatie- en translatiefout verminderen; controle zonder training wordt ook bereikt door een lichte puntenwolk aan te passen en een ruis-layout prior te gebruiken om het denoisen te sturen [11]. Scènebehoudende geometrische priors versterken verder de clip-niveau consistentie. Cami2V[57] behandelt camerapositie als een fysieke prior en benut epipolaire en multiview-beperkingen; RealCam-I2V [19] herstelt metrische diepte met DepthAnything v2 [48] om een schaalstabiele scène te reconstrueren; PoseTraj[16] maakt gebruik van pose-bewuste pretraining om

rotatie-uitgelijnde beweging te verkrijgen. vergeleken met alleen parameterconditionering verminderen deze priors de layout-drift binnen clips en behouden ze beter de lokale geometrie bij veranderingen in het zicht. Verder koppelt recent werk camerabesturing aan wereldmodellering. CVD [17], Cavia [46], en WoVoGen [22] synthetiseren gezamenlijk multi-view en multi-traject video's vanuit een gedeelde scène-representatie, waarbij cross-pad consistentie wordt afgedwongen. Ondertussen kunnen methoden die conditioneren vanuit expliciete renderbare 3D-representaties (bijv. 3D Gaussians) geometrie verankeren, cross-view 3D consistentie verbeteren en padnauwkeurigheid bevorderen [15, 29, 33, 52, 53]. Deze benaderingen bouwen echter meestal eenmalige 3D-scènes, terwijl wij lange-termijn camerabesturing verenigen met een persistent pose-geïndexeerd wereldgeheugen dat wordt gedeeld over trajecten.

3. Captain Safari

We introduceren *Captain Safari*, een geheugen-gestuurde videogeneratie framework. Sectie 3.1 presenteert een impliciet wereldgeheugen voor stabiele scènerepresentatie, terwijl Sectie 3.2 een pose-geconditioneerd ophaalsysteem beschrijft dat camerabeelden naar werelddaten mapt, wat een DiT-gesubstraatte generator leidt voor coherente outputs langs willekeurige trajecten.

3.1. Implicit Geheugen van Wereldgeometrie

Probleemopstelling. We representeren een video als $V = \{I_t\}_{t=0}^T$, waarbij I_t het frame is op tijdstip t . Op dezelfde tijds definieren we cameraposities $\mathcal{C} = \{(R_t, T_t)\}_{t=0}^T$ en verkrijgen we een 3D-bewust geheugenmerk m_t op elk tijdstip t met behulp van een voorgetrainde geometrie-encoder. Alle geheugenmerken vormen een globale geheugenbank $\mathcal{M} = \{m_t\}_{t=0}^T$.

Gegeven een teksprompt p , de cameraposities \mathcal{C} , en een doel

clip time step $\mathcal{T} = [t_0, t_1]$, together with its associated local world memory $\mathcal{M}_{\text{local}} \subset \mathcal{M}$, our goal is to synthesize a video segment $\hat{V}_{\mathcal{T}}$ that (i) aligns with p , (ii) respects the prescribed poses $\{(R_t, T_t)\}_{t \in \mathcal{T}}$, and (iii) maintains a coherent 3D world across viewpoints.

Local world memory. Directly conditioning on the full memory bank \mathcal{M} for every clip would be computationally expensive and dominated by temporally distant observations. Instead, for each target clip time step $\mathcal{T} = [t_0, t_1]$ we define a *local* memory $\mathcal{M}_{\text{local}} = \{m_{\tau} \mid \tau \in [k_s, k_e]\}$ whose endpoints are sampled under

$$\begin{aligned} t_0 - L &\leq k_s \leq t_0, \\ \max(k_s, t_0) + 1 &\leq k_e \leq \min(k_s + L, t_1), \end{aligned} \quad (1)$$

where L is a fixed bound and all time steps are integers. These constraints enforce that: (i) the memory window starts at most L seconds before the clip entrance t_0 , tying it to nearby observations; (ii) its duration is at most L , which keeps the conditioning set compact; and (iii) its end time k_e always touches or overlaps t_0 while remaining within $[t_0, t_1]$, ensuring that each clip is supported by a temporally compatible world prior. All $\mathcal{M}_{\text{local}}$ are constructed as such dynamic clip-aligned window of the shared bank \mathcal{M} , so neighboring clips naturally share overlapping memory entries, constraining computation while coupling their generations to a 3D-consistent underlying world.

Pose-retrieved memory. Within a given clip time step \mathcal{T} , we treat the local memory $\mathcal{M}_{\text{local}}$ as a static hypothesis of the surrounding world built from key frames. Each time step τ provides a pose token p_{τ} (derived from (R_{τ}, T_{τ})) and a set of 3D-aware memory tokens $m_{\tau,1}, \dots, m_{\tau,M}$. The collection $\{(p_{\tau}, m_{\tau,1:M})\}_{\tau}$ forms an implicit world table: pose token indicates *where* the camera has observed the scene, while memory tokens encode *what* the world looks like from those configurations. For any target time step $t \in \mathcal{T}$, we derive its camera pose to a query pose token p_t , embed it as $q_t = \phi_p(p_t)$, and use a dedicated retrieval module to read from this static table in a pose-dependent manner. Concretely, q_t is concatenated with a bank of learnable query tokens and processed into retrieval queries, which perform cross-attention over the encoded memory \tilde{X}^{mem} (defined in Sec. 3.2), yielding a set of world tokens

$$w_t = \text{Agg}\left(\text{CrossAttn}(Q_t, \tilde{X}^{\text{mem}})\right). \quad (2)$$

corresponding to the updated learnable queries. These pose-aligned world tokens w_t are directly used as the reconstructed memory at pose t . Thus, all frames in \mathcal{T} access local memory through pose-conditioned queries instead of raw time indices, encouraging multi-view observations to remain tied to a consistent static 3D world.

3.2. Memory Retrieval and Conditioning

Memory retriever design. As shown in Figure 2, given the local memory, we represent each time step τ by a pose

token p_{τ} and its associated memory tokens $m_{\tau,1:M}$. Our retriever is designed to (i) jointly encode pose–memory pairs into a coherent world representation, and (ii) extract, for any query pose, a compact set of pose-aligned tokens that summarize the most relevant parts of this local world.

We first embed pose and memory features into a shared space and form a joint sequence per time step:

$$\hat{X}_{\tau} = [\phi_p(p_{\tau}), \phi_m(m_{\tau,1}), \dots, \phi_m(m_{\tau,M})], \quad (3)$$

where ϕ_p and ϕ_m denote learnable embeddings for pose and memory tokens, respectively. A stack of transformer blocks (MemEnc) with 3D-aware positional encoding refines these sequences,

$$\tilde{X}_{\tau} = \text{MemEnc}(\hat{X}_{\tau}), \quad (4)$$

and we obtain the encoded local world memory by concatenation

$$\tilde{X}^{\text{mem}} = [\tilde{X}_{k_s}, \dots, \tilde{X}_{k_e}], \quad (5)$$

optionally masked to exclude padded or non-key entries.

For a target time step t , we derive the query pose token p_t , embed it as $q_t = \phi_p(p_t)$, and concatenate it with M learnable query tokens r_1, \dots, r_M ,

$$\hat{Q}_t = [q_t, r_1, \dots, r_M]. \quad (6)$$

This sequence is refined by transformer blocks sharing the same architecture as MemEnc, denoted as QryEnc, yielding pose-aware retrieval queries

$$Q_t = \text{QryEnc}(\hat{Q}_t). \quad (7)$$

We then perform cross-attention from Q_t to the encoded memory \tilde{X}^{mem} ,

$$Y_t = Q_t + \text{CrossAttn}(Q_t, \tilde{X}^{\text{mem}}), \quad (8)$$

and take the subset of tokens in Y_t corresponding to the learnable queries as the retrieved world tokens

$$w_t = [w_{t,1}, \dots, w_{t,M}], \quad (9)$$

which form a pose-aligned world feature for time t . During training, a linear head maps w_t back to the original memory space to reconstruct the target memory tokens at the query pose. Stacking multiple retrieval blocks iteratively refines both the queries and the retrieved tokens, enabling the model to softly route each query pose to the most relevant subset of past observations, instead of relying on a rigid temporal neighborhood or a single nearest frame.

Memory-conditioned DiT. For a given target clip time step \mathcal{T} , the retriever consumes $\mathcal{M}_{\text{local}}$ and the query pose p_t and outputs a pose-aligned set of world tokens $w_t \in \mathbb{R}^{M \times d_m}$, which summarize the static local world relevant to this segment. These tokens are mapped into the DiT hidden space by the memory embedding MLP

$$W_{\mathcal{T}} = \phi_w(w_t) \in \mathbb{R}^{M \times D}. \quad (10)$$

clip tijdstap $\mathcal{T} = [t_0, t_1]$, samen met het bijbehorende lokale wereldgeheugen $\mathcal{M}_{\text{local}} \subset \mathcal{M}$, is ons doel om een videosegment $V_{\mathcal{T}}$ te synthetiseren dat (i) overeenkomt met p , (ii) de voorgeschreven posities respecteert $\{(R_t, T_t)\}_{t \in \mathcal{T}}$, en (iii) een samenhangende 3D-wereld behoudt over verschillende gezichtspunten.

Lokale wereldgeheugen. Direct conditioneren op de volledige geheugenbank \mathcal{M} voor elke clip zou computationally duur zijn en gedomineerd worden door temporele verre observaties. In plaats daarvan definiëren we voor elke doel-clip tijdstap $\mathcal{T} = [t_0, t_1]$ een *lokale* geheugen $\mathcal{M}_{\text{local}} = \{m_{\tau} \mid \tau \in [k_s, k_e]\}$ waarvan de eindpunten worden bemonsterd onder

$$\begin{aligned} t_0 - L &\leq k_s \leq t_0, \\ \max(k_s, t_0) + 1 &\leq k_e \leq \min(k_s + L, t_1), \end{aligned} \quad (1)$$

waar L een vaste grens is en alle tijdstappen gehele getallen zijn. Deze beperkingen zorgen ervoor dat: (i) het geheugenvenster begint maximaal L seconden voor de clipingang t_0 , waardoor het aan nabije observaties wordt gekoppeld; (ii) de duur is maximaal L , wat de conditionerende set compact houdt; en (iii) de eindtijd k_e altijd samen met of overlapt t_0 terwijl het binnen $[t_0, t_1]$ blijft, wat ervoor zorgt dat elke clip wordt ondersteund door een temporele compatibele wereldvooraan. Alle $\mathcal{M}_{\text{local}}$ worden geconstrueerd als een dergelijk dynamisch clip-uitgelijnd venster van de gedeelde bank \mathcal{M} , zodat naburige clips van nature overlappende geheugenvervelingen delen, wat de berekening beperkt terwijl hun generaties worden gekoppeld aan een 3D-consistente onderliggende wereld.

Pose-opgehaald geheugen. Binnen een gegeven clip tijdstap \mathcal{T} beschouwen we het lokale geheugen $\mathcal{M}_{\text{local}}$ als een statische hypothese van de omringende wereld, opgebouwd uit sleutelbeelden. Elke tijdstap τ levert een pose-token p_{τ} (afgeleid van (R_{τ}, T_{τ})) en een set van 3D-bewuste geheugentokens $m_{\tau,1}, \dots, m_{\tau,M}$. De verzameling $\{(p_{\tau}, m_{\tau,1:M})\}_{\tau}$ vormt een impliciete werelddatabase: het pose-token geeft aan waar de camera de scène heeft waargenomen, terwijl geheugentokens coderen hoe de wereld eruitziet vanuit die configuraties. Voor elke doel-tijdstap $t \in \mathcal{T}$ leiden we de camerapositie af naar een query pose-token p_t , embedden het als $q_t = \phi_p(p_t)$, en gebruiken een speciale ophaalmodule om op een pose-afhankelijke manier uit deze statische tabel te lezen. Concreet wordt q_t geconcateneerd met een bank van leerbare querytokens en verwerkt tot ophaalqueries, die cross-attentie uitvoeren over het gecodeerde geheugen X^{mem} (gedefinieerd in Sectie 3.2), wat resulteert in een set van wereldtokens.

$$w_t = \text{Agg}\left(\text{CrossAttn}(Q_t, \tilde{X}^{\text{mem}})\right). \quad (2)$$

overeenkomend met de bijgewerkte leerbare queries. Deze pose-uitgelijnde wereldtokens w_t worden direct gebruikt als het gereconstrueerde geheugen op pose t . Zo krijgen alle frames in \mathcal{T} toegang tot lokaal geheugen via pose-geconditioneerde queries in plaats van ruwe tijdsindices, wat aanmoedigt dat multi-view observaties verbonden blijven met een consistent statische 3D-wereld.

3.2. Geheugenophaling en Conditionering

Geheugenophaalontwerp. Zoals getoond in Figuur 2, gegeven het lokale geheugen, representeren we elke tijdstap τ door een pose

token p_{τ} en de bijbehorende geheugentokens $m_{\tau,1:M}$. Onze ophaler is ontworpen om (i) pose-geheugenparen gezamenlijk te coderen in een samenhangende wereldrepresentatie, en (ii) voor elke querypose een compacte set pose-uitgelijnde tokens te extraheren die de meest relevante delen van deze lokale wereld samenvatten.

We embedden eerst pose- en geheugenmerken in een gedeelde ruimte en vormen een gezamenlijke reeks per tijdstap:

$$\hat{X}_{\tau} = [\phi_p(p_{\tau}), \phi_m(m_{\tau,1}), \dots, \phi_m(m_{\tau,M})], \quad (3)$$

waarbij ϕ_p en ϕ_m respectievelijk leerbare embeddings voor pose- en geheugentokens aanduiden. Een stapel transformerblokken (MemEnc) met 3D-bewuste positionele codering verfijnt deze reeksen,

$$\tilde{X}_{\tau} = \text{MemEnc}(\hat{X}_{\tau}), \quad (4)$$

en we verkrijgen het gecodeerde lokale wereldgeheugen door concatenatie

$$\tilde{X}^{\text{mem}} = [\tilde{X}_{k_s}, \dots, \tilde{X}_{k_e}], \quad (5)$$

optioneel gemaskeerd om opgevulde of niet-sleutelitems uit te sluiten.

Voor een doel-tijdstap t leiden we de query pose-token p_t af, embedden deze als $q_t = \phi_p(p_t)$, en voegen deze samen met M leerbare query-tokens r_1, \dots, r_M ,

$$\hat{Q}_t = [q_t, r_1, \dots, r_M]. \quad (6)$$

Deze reeks wordt verfijnd door transformerblokken die dezelfde architectuur delen als MemEnc, aangeduid als QryEnc, wat leidt tot pose-bewuste ophaalqueries

$$Q_t = \text{QryEnc}(\hat{Q}_t). \quad (7)$$

We voeren vervolgens cross-attentie uit van Q_t naar het gecodeerde geheugen X^{mem} ,

$$Y_t = Q_t + \text{CrossAttn}(Q_t, \tilde{X}^{\text{mem}}), \quad (8)$$

en nemen de subset van tokens in Y_t die overeenkomen met de leerbare queries als de opgehaalde wereldtokens

$$w_t = [w_{t,1}, \dots, w_{t,M}], \quad (9)$$

die een pose-uitgelijnde wereldfunctie vormen voor tijd t . Tijdens de training map een lineaire kop w_t terug naar de oorspronkelijke geheugenruimte om de doelgeheugentokens bij de query pose te reconstrueren. Het stapelen van meerdere ophaalblokken verfijnt iteratief zowel de queries als de opgehaalde tokens, waardoor het model in staat is om elke query pose zachtjes te routeren naar de meest relevante subset van eerdere observaties, in plaats van te vertrouwen op een rigide temporele buurt of een enkel dichtstbijzijnd frame.

Geheugen-geconditioneerde DiT. Voor een gegeven doel-clip tijdstap \mathcal{T} , verwerkt de retriever $\mathcal{M}_{\text{local}}$ en de queryholding p_t en levert een pose-uitgelijnde set van wereldtokens $w_t \in \mathbb{R}^{M \times d_m}$, die de statische lokale wereld samenvatten die relevant is voor dit segment. Deze tokens worden in de verborgen ruimte van DiT gemapt door de geheugen-embedding MLP.

$$W_{\mathcal{T}} = \phi_w(w_t) \in \mathbb{R}^{M \times D}. \quad (10)$$

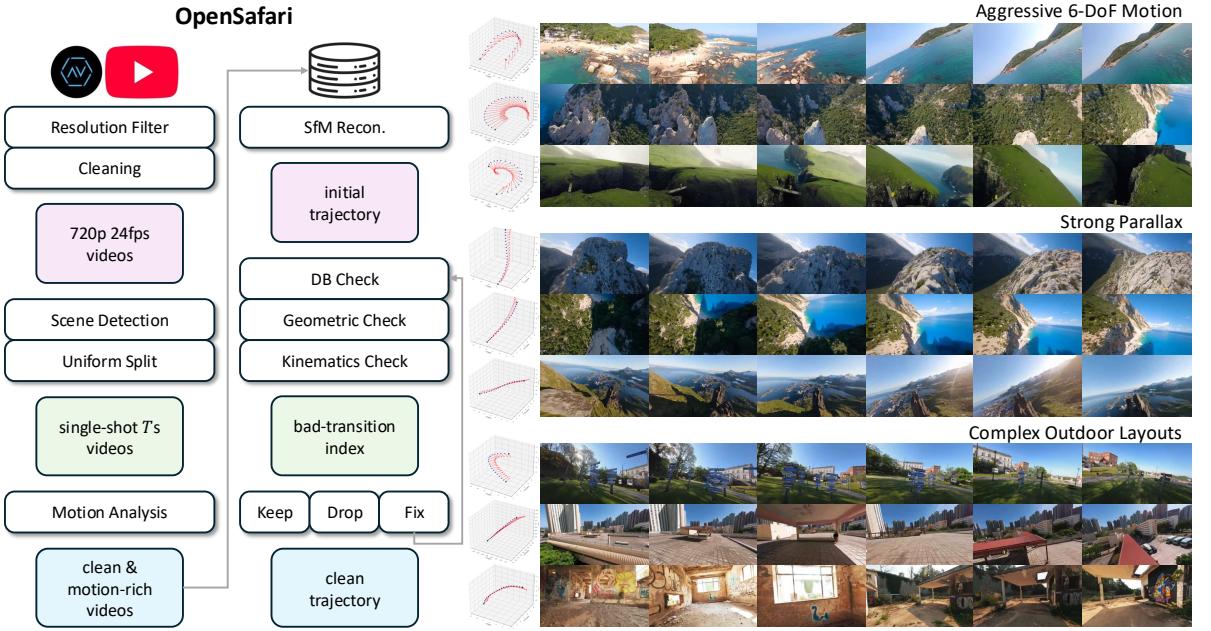


Figure 3. *OpenSafari*. A new in-the-wild FPV dataset with rigorously verified camera trajectories, designed to stress-test geometry-consistent, camera-controllable video generation. We curate clips through a compact, multi-stage pipeline that filters, reconstructs, and verifies trajectories, yielding clean, motion-rich videos with reliable camera paths.

The latent clip is encoded as a single spatio-temporal token sequence $Z \in \mathbb{R}^{L_z \times D}$, obtained by patchifying all frames in V_T . At each DiT layer l , we first apply self-attention over the full sequence and then inject the world tokens through a dedicated memory cross-attention:

$$Z^{(l+1)} = Z^l + \text{CrossAttn}(Z^l, W_T, W_T). \quad (11)$$

The clip-level world tokens W_T are reused as keys and values across all layers, providing a stable, 3D-consistent prior that shapes the denoising of every spatio-temporal token.

4. OpenSafari

4.1. Video Data Curation

Existing camera-conditioned datasets do not match our target regime. RealEstate10K [59] focuses on slow, mostly indoor real-estate walkthroughs with gentle motion and clean, quasi-static scenes, while Minecraft [4] is a synthetic voxel world with simplified geometry and engine-constrained dynamics. Neither captures aggressive, in-the-wild 6-DoF drone flight with strong parallax, large elevation changes, and complex outdoor layouts that truly stress long-horizon 3D consistency. We therefore propose *OpenSafari*, a new dataset of real-world FPV-style drone videos with verified camera trajectories tailored to this challenging setting.

We construct Safari-FPV from FPV-style drone videos

collected on AirVuz¹ and YouTube², and retain only clips that pass a strict multi-stage preprocessing pipeline. As shown in Figure 3, we: (i) download the highest available resolution for each URL and discard sources below the target resolution; (ii) normalize all videos to 720p, 24 fps, and a fixed 16:9 center crop, removing letterboxing and black borders so that subsequent camera estimation operates on a clean field of view; (iii) run scene detection to obtain single-shot segments; (iv) split segments into fixed-length T videos via uniform temporal slicing.

We then filter videos with a single diagnostic based on motion. Specifically, we run RAFT [36] to estimate optical-flow magnitude; videos with too little motion are removed, while videos with stable, coherent motion are kept to emphasize informative, parallax-rich trajectories rather than static views. Only videos satisfying the motion constraint enter the final dataset. This yields a large-scale, in-the-wild drone corpus explicitly tailored to stress-test geometry-aware, trajectory-following video generation.

4.2. Camera Trajectory Reconstruction

For each curated video, we estimate camera intrinsics and extrinsics at 4 fps using Hierarchical Localization [30, 31]. We extract local features, build exhaustive image pairs within each video, run feature matching, and reconstruct a COLMAP-style SfM model; from this model we export

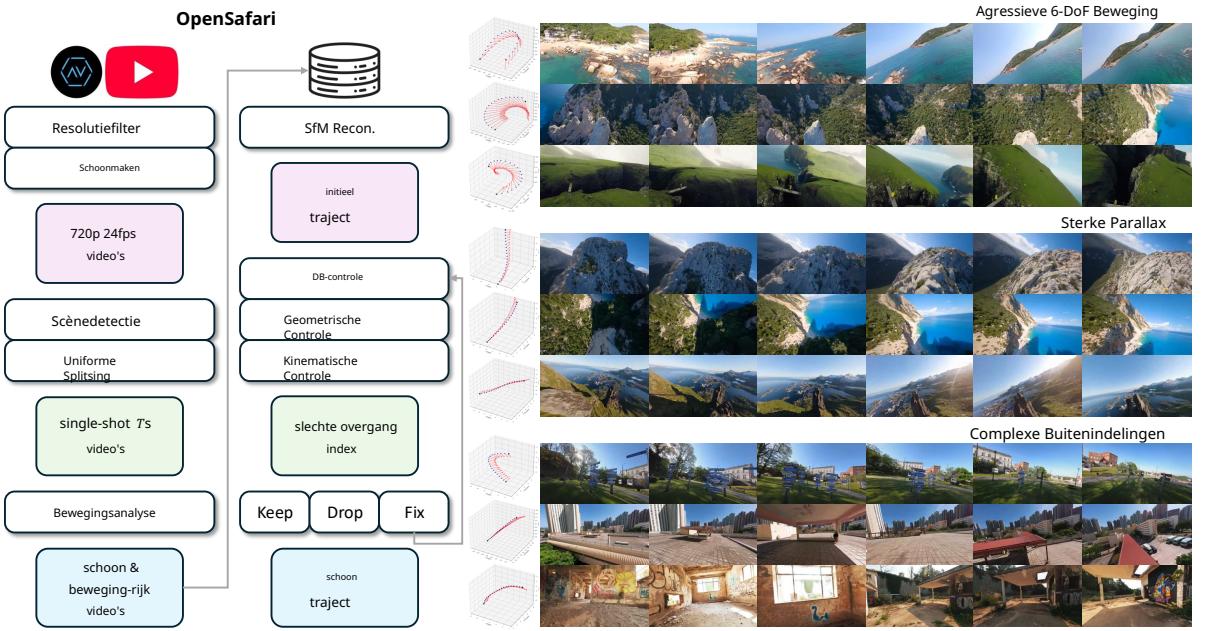


Figure 3. *OpenSafari*. Een nieuwe FPV-dataset in de natuur met zorgvuldig geverifieerde cameratrajecten, ontworpen om geometrie-consistente, camera-controleerbare videogeneratie te testen. We selecteren clips via een compacte, meerfasige pijplijn die trajecten filtert, reconstrueert en verifieert, resulterend in schone, bewegingrijke video's met betrouwbare camerapaden.

De latente clip wordt gecodeerd als een enkele spatio-temporele tokenreeks $Z \in \mathbb{R}^{L_z \times D}$, verkregen door alle frames in V_T te patchen. Bij elke DiT-laag l passen we eerst zelf-attentie toe over de volledige reeks en injecteren vervolgens de wereldtokens via een speciale geheugen-cross-attentie:

$$Z^{(l+1)} = Z^l + \text{CrossAttn}(Z^l, W_T, W_T). \quad (11)$$

De wereldtokens op clipniveau W_T worden hergebruikt als sleutels en waarden over alle lagen, wat een stabiele, 3D-consistente basis biedt die de ruisonderdrukking van elke spatio-temporele token vormgeeft.

4. OpenSafari

4.1. Video Data Curation

Bestaande camera-geconditioneerde datasets komen niet overeen met ons doelregime. RealEstate10K [59] richt zich op langzame, meestal binnenhuis-wandelingen met zachte bewegingen en schone, quasi-statische scènes, terwijl Minecraft [4] een synthetische voxelwereld is met vereenvoudigde geometrie en door de engine beperkte dynamiek. Geen van beide legt agressieve, in-the-wild 6-DoF dronevluchten vast met sterke parallax, grote hoogteverschillen en complexe buitenindelingen die echt de lange-termijn 3D consistentie testen. Daarom stellen we *OpenSafari* voor, een nieuwe dataset van real-world FPV-stijl dronevideo's met geverifieerde cameratrajecten die zijn afgestemd op deze uitdagende omgeving.

We construeren Safari-FPV uit FPV-stijl dronevideo's

verzameld op AirVuz¹ en YouTube², en behouden alleen clips die een strikte meerfasige voorverwerkingspijplijn doorstaan. Zoals getoond in Figuur 3, doen we het volgende: (i) we downloaden de hoogste beschikbare resolutie voor elke URL en verwijderen bronnen onder de doelresolutie; (ii) we normaliseren alle video's naar 720p, 24fps, en een vaste 16:9 centrale uitsnede, waarbij we letterboxen en zwarte randen verwijderen zodat de daaropvolgende camera-estimatie op een schoon gezichtsveld werkt; (iii) we voeren scènedetectie uit om enkelvoudige segmenten te verkrijgen; (iv) we splitsen segmenten in video's van vaste lengte via uniforme temporele slicing.

We filteren vervolgens video's met een enkele diagnose op basis van beweging. Specifiek draaien we RAFT [36] om de optische-stroomgrootte te schatten; video's met te weinig beweging worden verwijderd, terwijl video's met stabiele, coherente beweging worden behouden om informatie, parallax-rijke trajecten te benadrukken in plaats van statische beelden. Alleen video's die voldoen aan de bewegingsbeperkingen komen in de uiteindelijke dataset. Dit levert een grootschalige, in-the-wild dronecorpus op die expliciet is afgestemd om geometriebewuste, trajectvolgende videogeneratie te stress-testen.

4.2. Camera Trajectoireconstructie

Voor elke samengestelde video schatten we de interne en externe camerakenmerken bij 4fps met behulp van Hiërarchische Lokalisatie [30, 31]. We extraheren lokale kenmerken, bouwen uitputtende beeldparen binnen elke video, voeren kenmerkmatching uit en reconstrueren een C-OLMAP-stijl SfM-model; uit dit model exporteren we

¹<https://www.airvuz.com/>

²<https://www.youtube.com/>

Table 1. **Benchmark camera-controlled video generation.** *Captain Safari* ranks first in 3D consistency and trajectory following with competitive video quality. Compared to the ablated variant without memory, *Captain Safari* substantially improves 3D consistency and trajectory following, with only a slight trade-off in video quality. (Recon. = reconstruction rate. CosSim = cosine similarity.)

Model	Video Quality		3D consistency		Trajectory Following		
	FVD ↓	LPIPS ↓	MEt3R ↓	Recon. ↑	AUC@30 ↑	AUC@15 ↑	CosSim ↑
Geometry Forcing [44]	2662.75	0.667	0.4834	0.877	0.168	0.056	0.429
Real-CamI2V [19]	1585.61	0.513	<u>0.3703</u>	<u>0.923</u>	0.174	0.051	0.296
Wan2.2-5B-Control-Camera [38]	1387.75	0.545	0.3932	0.767	0.181	0.054	0.420
Captain Safari w/o Mem.	998.47	0.504	0.3720	0.912	<u>0.193</u>	0.068	<u>0.508</u>
Captain Safari	1023.46	0.512	0.3690	0.968	0.200	0.068	0.563

Table 2. **Human preference.** Users overwhelmingly prefer *Captain Safari* across all criteria, capturing **67%** of total votes. The memory-removed variant ranks a distant second, while baselines competitors receive single-digit preference.

Model	Video Quality	3D consistency	Trajectory Following	Average
Geometry Forcing [44]	0.20%	0.00%	0.20%	0.13%
Real-CamI2V [19]	4.20%	6.40%	4.40%	5.00%
Wan2.2-5B-Control-Camera [38]	3.20%	3.80%	6.40%	4.47%
Captain Safari w/o Mem.	25.00%	24.20%	20.00%	23.07%
Captain Safari	67.40%	65.60%	69.00%	67.33%

per-frame camera parameters as initial trajectories.

To obtain deployment-ready data, we apply a three-stage verification-and-fix pipeline to every reconstructed trajectory. First, *database check* consumes SfM statistics (inlier counts and ratios) to flag potentially unreliable transitions. Next, *geometric check* revisits suspicious pairs using stored keypoints and matches, recomputes essential matrices, and thresholds symmetric epipolar errors. Last, *kinematics check* analyzes the pose sequence for translation spikes, rotation jumps, forward-direction flips, and higher-order smoothness violations, using robust MAD-based scores to detect implausible motion.

The per-transition decisions are fused into a binary bad-index, which drives a strict policy. If bad transitions are sparse and localized, we invoke a targeted fix: we linearly interpolate camera centers and apply SLERP to rotations with a capped interpolation angle, optionally extrapolating at video boundaries. The fixed segments are then re-validated by the same database/geometric/kinematics criteria. If post-fix validation succeeds, the trajectory is exported into the final dataset. If the bad-index is too dense, violations are too severe, or fixed trajectories still fail verification, the entire video is discarded.

The resulting *OpenSafari* couples high-dynamic, in-the-wild FPV drone video with rigorously verified camera trajectories. It departs from existing benchmarks by emphasizing aggressive 6-DoF motion, strong parallax, and complex outdoor layouts, while enforcing strict geometric and kinematic validation. This makes *OpenSafari* a challenging testbed for camera-controllable video generation.

5. Experiments

5.1. Implementation Details

Training recipe. We adopt a two-stage recipe. We first warm up the pose-conditioned memory retriever using pose-aligned memory tokens m_t . We then jointly train the retriever and DiT end-to-end, updating the DiT via LoRA [12]. Memory cross-attention is initialized from the corresponding context cross-attention weights, and other new layers use standard initialization.

Dataset. We extract overlapping clips with 1 s stride, yielding 51,997 training candidates. A diversity-based trajectory filter removes clips with near-static motion, resulting in 11,481 final training clips. We additionally construct a non-overlapping test set of 787 clips for evaluation. For each clip, we generate a single descriptive caption using Qwen2.5-VL-7B [3] and use it as the text condition.

Configuration and notation. We generate $\mathcal{T} = 5$ s clips at 24 fps from $T = 15$ s videos. Camera poses and memory features are sampled at 4 fps. For a target 5 s clip with interval $[t_0, t_1]$, we use the terminal pose p_{t_1} as the query. The memory window is limited to $L = 5$ s. We use Wan2.2-Fun-5B-Control-Camera [38] as our base DiT with a hidden dimension $D = 3072$. Retriever and DiT are trained with 1 and 5 epochs, respectively. For each video, we extract 3D-aware memory feature from a pretrained StreamVGG [62]. We select four layers $\{4, 11, 17, 23\}$; at each layer, the feature contains 782 tokens. Concatenating across the four layers yields $M = 4 \times 782$ and $d_m = 1024$ memory tokens per frame.

Tabel 1. **Benchmark camera-gestuurde videogeneratie.** *Captain Safari* staat op de eerste plaats in 3D consistentie en trajectvolging met concurrerende videokwaliteit. vergeleken met de geablateerde variant zonder geheugen, verbetert *Captain Safari* de 3D consistentie en trajectvolging aanzienlijk, met slechts een kleine concessie in videokwaliteit. (Recon. = reconstructiesnelheid. CosSim = cosinusgelijkenis.)

Model	Videokwaliteit		3D consistentie		Trajectvolging		CosSim ↑
	FVD ↓	LPIPS ↓	MEt3R ↓	Recon. ↑	AUC@30 ↑	AUC@15 ↑	
Geometrie Dwang [44]	2662.75	0.667	0.4834	0.877	0.168	0.056	0.429
Real-CamI2V [19]	1585.61	0.513	<u>0.3703</u>	<u>0.923</u>	0.174	0.051	0.296
Wan2.2-5B-Control-Camera [38]	1387.75	0.545	0.3932	0.767	0.181	0.054	0.420
Captain Safari z/oGeheugen.	998.47	0.504	0.3720	0.912	<u>0.193</u>	0.068	<u>0.508</u>
Captain Safari	1023.46	0.512	0.3690	0.968	0.200	0.068	0.563

Tabel 2. **Menselijke voorkeur.** Gebruikers geven overweldigend de voorkeur aan *Captain Safari* op alle criteria, met **67%** van de totale stemmen. De variant zonder geheugen staat op een verre tweede plaats, terwijl de basisconcurrenten een voorkeur in enkelcijferige percentages ontvangen.

Model	Videokwaliteit		3D Consistentie		Trajectvolging		Gemiddeld
	FVD ↓	LPIPS ↓	MEt3R ↓	Recon. ↑	AUC@30 ↑	AUC@15 ↑	
Geometrie Dwang [44]	0,20%	0,00%	0,20%	0,20%	0,00%	0,00%	0,13%
Real-CamI2V [19]	4,20%	6,40%	4,40%	4,40%	3,20%	3,80%	5,00%
Wan2.2-5B-Control-Camera [38]	3,20%	3,80%	6,40%	6,40%	24,20%	24,20%	4,47%
Captain Safari z/oGeheugen.	25,00%	24,20%	20,00%	23,07%	67,40%	65,60%	23,07%
Captain Safari	67,40%	65,60%	69,00%	67,33%	67,40%	65,60%	67,33%

per-frame cameragegevens als initiële trajecten.

Om gegevens gereed voor implementatie te verkrijgen, passen we een verificatie- en herstelpipeline in drie fasen toe op elke gereconstrueerde trajectorie. Eerst gebruikt *databasecontrole* SfM-statistieken (inlier-aantallen en -verhoudingen) om potentieel onbetrouwbare overgangen te markeren. Vervolgens herbekijkt *geometrische controle* verdachte paren met behulp van opgeslagen sleutelpunten en overeenkomsten, herberekent essentiële matrices en stelt drempels voor symmetrische epipolaire fouten. Ten slotte analyseert *kinematische controle* de posevolgorde op vertaalpielen, rotatiesprongen, omkeringen van de voorwaartse richting en schendingen van hogere-orde gladheid, waarbij robuuste MAD-gebaseerde scores worden gebruikt om onwaarschijnlijke bewegingen te detecteren.

De beslissingen per overgang worden samengevoegd tot een binaire slechte-index, die een strikt beleid aanstuurt. Als slechte overgangen schaars en lokaal zijn, passen we een gerichte correctie toe: we interpoleren lineair de cameracentra en passen SLERP toe op rotaties met een begrenste interpolatiehoek, waarbij we desgewenst extrapoleren aan de videoranden. De gecorrigeerde segmenten worden vervolgens opnieuw gevalideerd volgens dezelfde database-/geometrische-/kinematische criteria. Als de validatie na correctie slaagt, wordt de trajectorie geëxporteerd naar het definitieve dataset. Als de slechte-index te dicht is, de schendingen te ernstig zijn, of de gecorrigeerde trajectorieën nog steeds niet slagen voor verificatie, wordt de gehele video verworpen.

Het resulterende *OpenSafari* combineert dynamische, in-het-wild FPV drone video met rigoureus geverifieerde camera-trajecten. Het wijkt af van bestaande benchmarks door de nadruk te leggen op agressieve 6-DoF beweging, sterke parallax en complexe buitenlayouts, terwijl strikte geometrische en kinematische validatie wordt afgeworpen. Dit maakt *OpenSafari* een uitdagende testomgeving voor camera-controleerbare videoproductie.

5. Experimenten

5.1. Implementatiедетали

Trainingsреcept. We hanteren een tweestapsреcept. Eerst warmen we de pose-geconditioneerde geheugenophaler op met behulp van pose-uitgelijnde geheugentokens m_t . Vervolgens trainen we de ophaler en DiT gezamenlijk end-to-end, waarbij we de DiT bijwerken via LoRA [12]. Geheugen cross-attentie wordt geïnitialiseerd vanuit de overeenkomstige context cross-attentie gewichten, en andere nieuwe lagen gebruiken standaardinitialisatie.

Dataset. We extraheren overlappende clips met een stap van 1s, resulterend in 51,997 trainingskandidaten. Een diversiteitsgebaseerd trajectfilter verwijdert clips met bijna-statistische beweging, wat resulteert in 11,481 definitieve trainingsclips. We construeren daarnaast een niet-overlappende testset van 787 clips voor evaluatie. Voor elke clip genereren we een enkele beschrijvende caption met behulp van Qwen2.5-VL-7B [3] en gebruiken deze als tekstvoorwaarde.

Configuratie en notatie. We genereren $\mathcal{T} = 5$ s clips bij 24 fps van $T = 15$ s video's. Cameraposities en geheugenmerken worden gesampled bij 4 fps. Voor een doelclip van 5 s met interval $[t_0, t_1]$ gebruiken we de eindpositie p_{t_1} als de query. Het geheugenvenster is beperkt tot $L = 5$ s. We gebruiken Wan2.2-Fun-5B-Control-Camera [38] als onze basis DiT met een verborgen dimensie $D = 3072$. Retriever en DiT worden respectievelijk getraind met 1 en 5 epochs. Voor elke video extraheren we 3D-bewuste geheugenmerken van een voorgetrainde StreamVGGT [62]. We selecteren vier lagen $\{4, 11, 17, 23\}$; op elke laag bevat het kenmerk 782 tokens. Het samenvoegen over de vier lagen levert $M = 4 \times 782$ en $d_m = 1024$ geheugentokens per frame op.

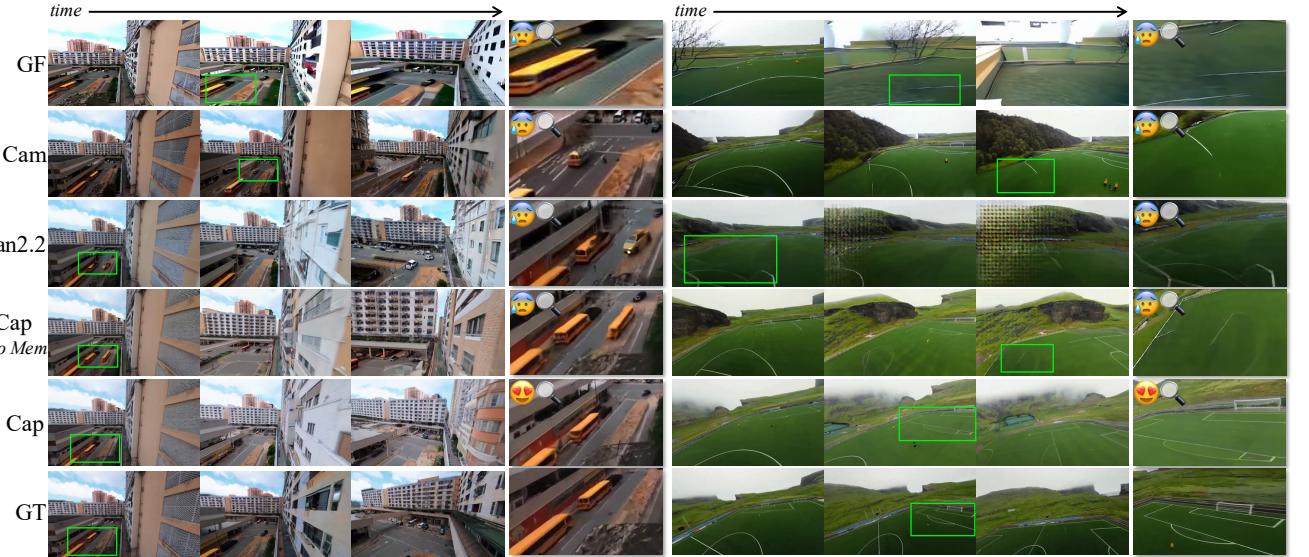


Figure 4. **Qualitative comparisons.** **Left:** Baselines—including the memory-removed variant—exhibit abrupt popping/vanishing of the school bus, and GF is low-quality. **Captain Safari** alone renders the bus smoothly exiting the frame. **Right:** Baselines distort or lose field marking, with Wan2.2 collapsing under large camera motion, affirming the challenge of 3D consistency under rapid trajectories. **Captain Safari** preserves crisp markings and coherent layout while following the fast 6-DoF path.

5.2. Benchmark

Metrics. We evaluate video generation along three complementary axes: video quality, 3D consistency, and trajectory following. For video quality, we report FVD [37] and LPIPS [54]. For 3D consistency, we use MEt3R [1], computed between GT and generated videos at matched time steps and a reconstruction rate that measures the percentage of frames successfully registered in the recovered 3D model [30, 31]. For trajectory following, we report camera relocation accuracy (AUC [39]) and the cosine similarity between the flattened camera pose, capturing how the model adheres to the desired camera parameters over time.

Baselines. We compare against representative camera-controllable video generation models, including Geometry Forcing [44], Real-CamI2V [19, 57], and Wan2.2-5B-Control-Camera [38], which cover geometry-constrained, reconstruction-driven, and large-scale diffusion-based approaches to trajectory-conditioned video synthesis.

Human Study. We conduct a human study with 50 participants. Each participant is presented with 10 cases, where each case contains the GT video and five anonymized model-generated videos (three baselines, our model, and its ablated variant). For every case, participants are asked to select the best video under three criteria: Video Quality, 3D Consistency, and Trajectory Following. In total, the study collects $50 \times 10 \times 3 = 1,500$ human preference votes.

5.3. Generation Quality

As shown in Table 1, our *Captain Safari* attains a substantially lower FVD (1023.46 vs. 1387.75) and a slightly

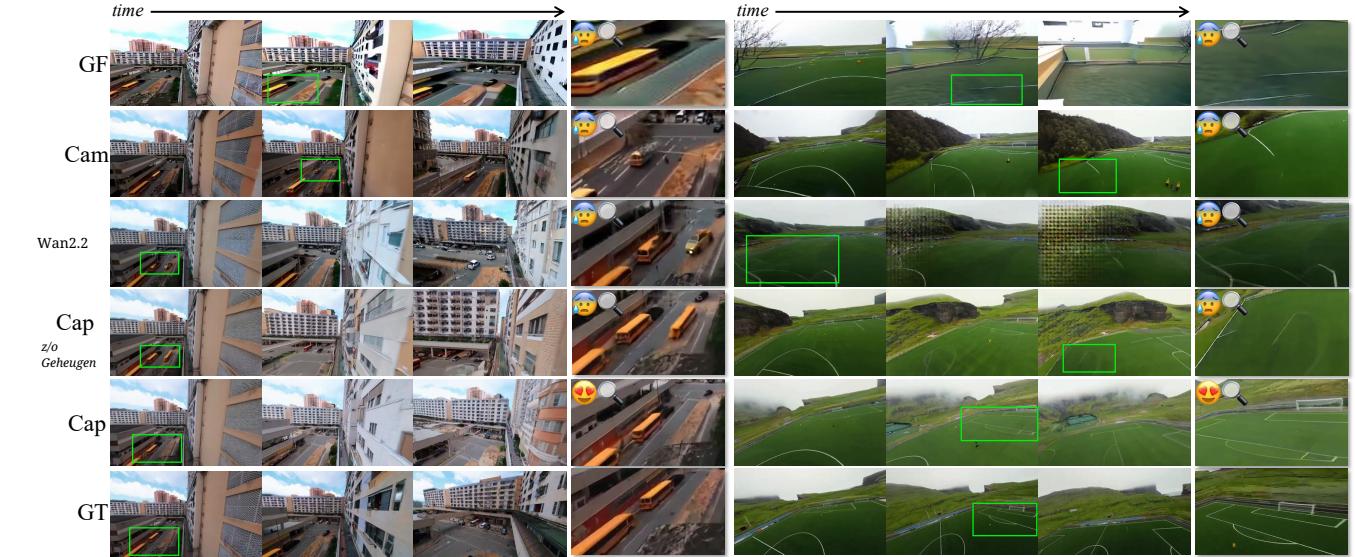
improved LPIPS score (0.512 vs. 0.513) compared to the SOTA baseline, demonstrating more stable temporal dynamics and sharper spatial details. Moreover, the human study in Table 2 indicates that **67.40%** of participants prefer our videos over all competing methods, highlighting the perceptual realism and overall fidelity of our generations.

Qualitative comparisons in Figure 4 further reveal that *Captain Safari* produces visually compelling, realistic, and highly authentic scene dynamics. These findings are also consistent with the samples shown in Figure 1, where our method delivers vivid, coherent, and natural-looking drone videos that closely resemble real-world captures.

5.4. 3D Consistency

Captain Safari achieves state-of-the-art 3D consistency. As shown in Table 1, our method lowers MEt3R by 0.0013 (0.3690 vs. 0.3703) and raises the reconstruction rate by 0.045 (0.968 vs. 0.923) compared to the strongest baseline. Consistently, the human study in Table 2 shows that **65.60%** of participants prefer *Captain Safari* for 3D consistency, substantially surpassing all competing approaches.

Qualitative visualizations further confirm these quantitative gains. In Figure 1, structures such as the Greek-style columns remain geometrically stable across large viewpoint changes. In Figure 4, our model produces (*left*) a school bus that smoothly moves out of the frame, and (*right*) preserves crisp, globally consistent field markings on the soccer pitch, whereas baselines exhibit distortions and disappearance. Figure 5 and Figure 1 further show that our reconstructions yield sharper façades and well-formed windows



Figuur 4. **Kwalitatieve vergelijkingen.** **Links:** Baselines—waaronder de variant zonder geheugen—vertonen abrupt opduiken/verdwijnen van de schoolbus, en GF is van lage kwaliteit. **Captain Safari** alleen laat de bus soepel uit het frame verdwijnen. **Rechts:** Baselines vervormen of verliezen veldmarkeringen, waarbij Wan2.2 instort bij grote camerabewegingen, wat de uitdaging van 3D consistentie onder snelle trajecten bevestigt. **Captain Safari** behoudt scherpe markeringen en een samenhangende lay-out terwijl het het snelle 6-DoF pad volgt.

5.2. Benchmark

Metrieken. We evalueren videogeneratie langs drie complementaire assen: videokwaliteit, 3D consistentie en trajectvolg. Voor videokwaliteit rapporteren we FVD [37] en LPIPS [54]. Voor 3D consistentie gebruiken we MEt3R [1], berekend tussen GT en gegenereerde video's op overeenkomende tijdstappen en een reconstructiesnelheid die het percentage frames meet dat succesvol is geregistreerd in het herstelde 3D-model [30, 31]. Voor trajectvolg rapporteren we de nauwkeurigheid van cameraverplaatsing (AUC [39]) en de cosinusgelijkenis tussen de afgevlakte camerapositie, die vastlegt hoe het model zich houdt aan de gewenste cameraparameters in de loop van de tijd.

Baselines. We vergelijken met representatieve camera-controleerbare videogeneratiemodellen, waaronder Geometry Forcing [44], Real-CamI2V [19, 57], en Wan2.2-5B-Control-Camera [38], die geometrie gebonden, reconstructie gedreven en grootschalige diffusie gebaseerde benaderingen voor traject geconditioneerde videosynthese omvatten.

Menselijke Studie. We voeren een menselijke studie uit met 50 deelnemers. Elke deelnemer krijgt 10 gevallen gepresenteerd, waarbij elk geval de GT-video en vijf geanonimiseerde door modellen gegenereerde video's bevat (drie baselines, ons model en de geablateerde variant). Voor elk geval wordt de deelnemers gevraagd de beste video te selecteren op basis van drie criteria: Videokwaliteit, 3D Consistentie en Trajectvolg. In totaal verzamelt de studie $50 \times 10 \times 3 = 1,500$ menselijke voorkeurstemmen.

5.3. Generatiekwaliteit

Zoals getoond in Tabel 1, behaalt onze *Captain Safari* een aanzienlijk lagere FVD (1023,46 vs. 1387,75) en een iets

verbeterde LPIPS-score (0,512 vs. 0,513) vergeleken met de SOTA-basislijn, wat stabielere temporele dynamiek en scherpere ruimtelijke details aantoon. Bovendien geeft de menselijke studie in Tabel 2 aan dat **67,40%** van de deelnemers onze video's verkiezen boven concurrerende methoden, wat het perceptuele realisme en de algehele getrouwheid van onze generaties benadrukt.

Kwalitatieve vergelijkingen in Figuur 4 laten verder zien dat *Captain Safari* visueel aantrekkelijke, realistische en zeer authentieke scènes-dynamiek produceert. Deze bevindingen komen ook overeen met de voorbeelden getoond in Figuur 1, waar onze methode levendige, coherente en natuurlijk ogende dronevideo's levert die sterk lijken op opnames uit de echte wereld.

5.4. 3D Consistentie

Captain Safari bereikt state-of-the-art 3D consistentie. Zoals getoond in Tabel 1, verlaagt onze methode MEt3R met 0,0013 (0,3690 vs. 0,3703) en verhoogt het reconstructiepercentage met 0,045 (0,968 vs. 0,923) vergeleken met de sterkste basislijn. Consistent hiermee toont de menselijke studie in Tabel 2 aan dat **65,60%** van de deelnemers de voorkeur geeft aan *Captain Safari* voor 3D consistentie, wat alle concurrerende benaderingen aanzienlijk overtreft.

Kwalitatieve visualisaties bevestigen verder deze kwantitatieve verbeteringen. In Figuur 1 blijven structuren zoals de Griekse zuilen geometrisch stabiel bij grote veranderingen in het gezichtspunt. In Figuur 4 produceert ons model (*links*) een schoolbus die soepel uit het frame beweegt, en (*rechts*) behoudt scherpe, wereldwijd consistentie veldmarkeringen op het voetbalveld, terwijl basislijnen vervormingen en verdwijningen vertonen. Figuur 5 en Figuur 1 laten verder zien dat onze reconstructies scherpere gevallen en goed gevormde ramen opleveren.

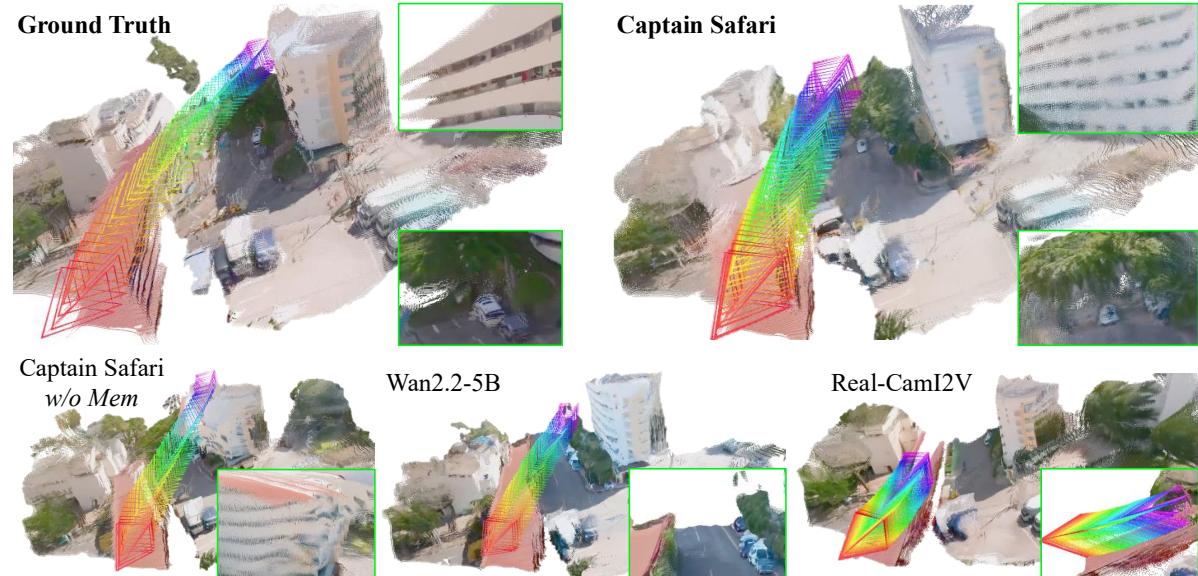


Figure 5. Scene reconstruction and camera trajectory. With pose-aligned memory, *Captain Safari* reconstructs a well-structured building façade (the memory-removed variant blurs/warps it), demonstrating the benefit of memory. It also preserves fine details—parked cars and the tree on their roofs—that Wan2.2-5B fails to retain. Meanwhile, Real-CamI2V follows only a short path, whereas *Captain Safari* covers the full trajectory with stable 3D structure, highlighting the challenge of maintaining 3D consistency under fast motion.

without collapsing geometry. Together, these results validate that the implicit world memory and pose-conditioned retrieval of *Captain Safari* effectively stabilize the underlying 3D world under aggressive camera motion.

5.5. Trajectory Following

Captain Safari delivers the most accurate trajectory following among all competing models. As shown in Table 1, our method achieves the highest AUC@30 (0.200) and AUC@15 (0.068), along with the best cosine similarity (0.563), outperforming the strongest baseline by clear margins. The human study in Table 2 further reinforces this observation, with **69.00%** of participants identifying our model as the most faithful to the target camera path.

Figure 5 provides a clear visualization of these improvements. *Captain Safari*'s predicted trajectory closely aligns with the ground-truth path, while the ablated variant deviates and flies over the rooftop, and RealCam-I2V fails to follow the intended forward motion, advancing only slightly rather than committing to the prescribed trajectory. Furthermore, our method demonstrates stable and coherent generation under challenging viewpoint changes with complex camera maneuvers in Figure 1. These results highlight the effectiveness of our memory-augmented, pose-conditioned design for precise trajectory adherence.

5.6. Ablation Study

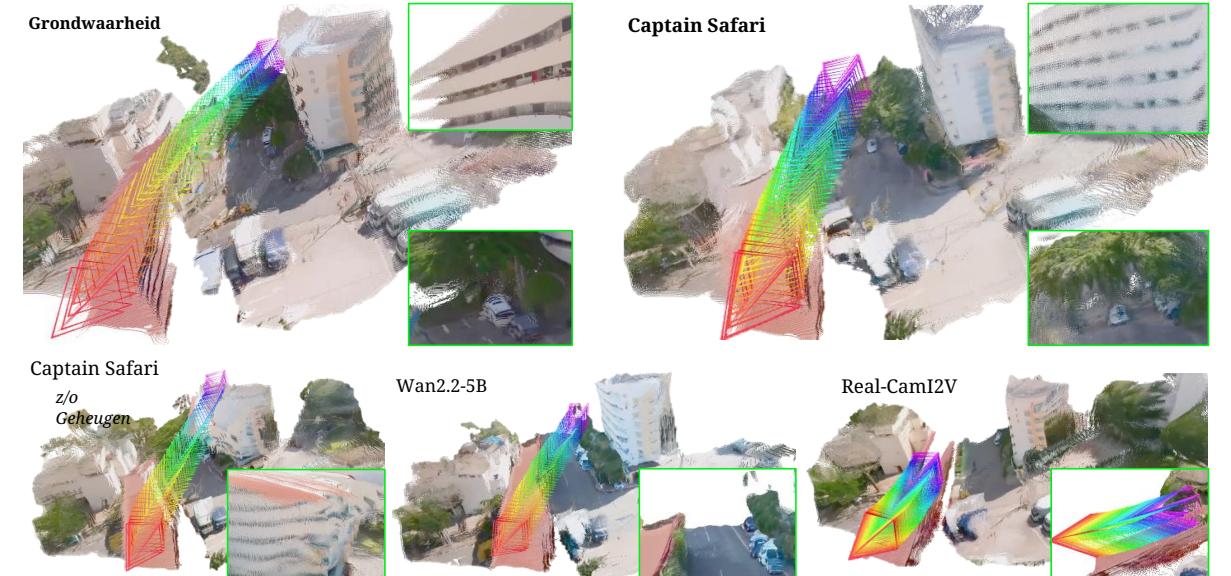
Our results highlight the importance of the proposed pose-conditioned world memory. As shown in Table 1, adding memory yields substantial improvements in both 3D con-

sistency and trajectory following. These gains confirm that retrieving pose-aligned world features at the target frame provides the model with an explicit understanding of what the scene *should* look like, enabling stable geometry and accurate motion alignment.

Qualitative comparisons in Figure 4 and Figure 5 further illustrate these effects. With memory, the generated scenes preserve global structure, maintain consistent geometry across viewpoints. In contrast, the ablated variant often drifts and exhibits geometric inconsistencies. Together, these results validate the effectiveness of our memory-augmented design in stabilizing the underlying 3D world and guiding precise camera motion.

6. Conclusion

We introduced *Captain Safari*, a pose-conditioned world engine built on a world memory that enables long-range, 3D-consistent video generation under complex FPV trajectories. Together with *OpenSafari*, our curated dataset of in-the-wild drone videos with verified camera poses, this establishes a rigorous benchmark for controllable video generation. *Captain Safari* markedly improves 3D consistency and trajectory accuracy over prior methods while maintaining strong visual fidelity. Although the system incurs non-trivial inference overhead, future work will explore real-time world engines with lightweight memory and faster generative backbones. We hope *Captain Safari* and *OpenSafari* encourage further research in persistent world models and long-horizon controllable video generation.



Figuur 5. Scenereconstructie en cameratraject. Met pose-uitgelijnd geheugen reconstrueert *Captain Safari* een goed gestructureerde gebouwgevel (de variant zonder geheugen vervaagt/vervormt het), wat het voordeel van geheugen aantoont. Het behoudt ook fijne details—geparkeerde auto's en de boom op hun daken—die Wan2.2-5B niet kan behouden. Ondertussen volgt Real-CamI2V slechts een kort pad, terwijl *Captain Safari* het volledige traject met stabiele 3D-structuur bestrijkt, wat de uitdaging benadrukt om 3D consistentie te behouden bij snelle bewegingen.

zonder dat de geometrie instort. Samen bevestigen deze resultaten dat het impliciete wereldgeheugen en pose-geconditioneerde retrievel van *CaptainSafari* effectief de onderliggende 3D-wereld stabiliseren bij agressieve camerabewegingen.

5.5. Trajectvolging

CaptainSafari levert de meest nauwkeurige trajectvolging van alle concurrenente modellen. Zoals te zien in Tabel 1, behaalt onze methode de hoogste AUC@30 (0.200) en AUC@15 (0.068), samen met de beste cosinusgelijkenis (0.563), waarmee het de sterkste basislijn met duidelijke marges overtreft. De menselijke studie in Tabel 2 versterkt deze observatie verder, met **69,00%** van de deelnemers die ons model als het meest trouw aan het doelcamerapad identificeren.

Figuur 5 biedt een duidelijke visualisatie van deze verbeteringen. *Captain Safari*'s voorspelde traject komt nauw overeen met het werkelijke pad, terwijl de geablateerde variant afwijkt en over het dak vliegt, en RealCam-I2V er niet in slaagt de bedoelde voorwaartse beweging te volgen, slechts licht vooruitgaand in plaats van zich aan het voorgeschreven traject te houden. Bovendien toont onze methode stabiele en coherente generatie onder uitdagende veranderingen in het gezichtspunt met complexe camerabewegingen in Figuur 1. Deze resultaten benadrukken de effectiviteit van ons geheugen-verrijkte, pose-geconditioneerde ontwerp voor nauwkeurige trajectrouw.

5.6. Ablatie Studie

Onze resultaten benadrukken het belang van het voorgestelde pose-geconditioneerde wereldgeheugen. Zoals getoond in Tabel 1, het toevoegen van geheugen levert aanzienlijke verbeteringen op in zowel 3D-con-

sistentie als trajectvolging. Deze verbeteringen bevestigen dat het ophalen van pose-uitgelijnde wereldkenmerken bij het doelkader het model een expliciet begrip geeft van hoe de scène zou moeten zijn, wat stabiele geometrie en nauwkeurige bewegingsuitlijning mogelijk maakt.

Kwalitatieve vergelijkingen in Figuur 4 en Figuur 5 illustreren deze effecten verder. Met geheugen behouden de gegenereerde scènes de globale structuur en handhaven ze consistente geometrie over verschillende gezichtspunten. Daarentegen vertoont de geablateerde variant vaak afwijkingen en geometrische inconsistenties. Samen bevestigen deze resultaten de effectiviteit van ons geheugen-verrijkte ontwerp in het stabiliseren van de onderliggende 3D-wereld en het begeleiden van nauwkeurige camerabewegingen.

6. Conclusie

We hebben *Captain Safari* geïntroduceerd, een pose-geconditioneerde wereldmotor gebouwd op een wereldgeheugen dat langeafstands, 3D-consistente videoproduktie mogelijk maakt onder complexe FPV-trajecten. Samen met *OpenSafari*, ons samengestelde dataset van dronevideo's in het wild met geverifieerde cameraposities, vormt dit een rigoureuze benchmark voor controleerbare videoproduktie. *CaptainSafari* verbetert de 3D consistentie en trajectnauwkeurigheid aanzienlijk ten opzichte van eerdere methoden, terwijl het een sterke visuele kwaliteit behoudt. Hoewel het systeem aanzienlijke inferentie-overhead met zich meebrengt, zal toekomstig werk zich richten op real-time wereldmotoren met lichtgewicht geheugen en snellere generatieve basissen. We hopen dat *Captain Safari* en *OpenSafari* verder onderzoek aanmoedigen naar persistente wereldmodellen en lange-termijn controleerbare videoproduktie.

References

- Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 6
- [1] Mohammad Asim, Christopher Wewer, Thomas Wimmer, Bernt Schiele, and Jan Eric Lenssen. Met3r: Measuring multi-view consistency in generated images. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6034–6044, 2025. 7
- [2] Sherwin Bahmani, Xian Liu, Wang Yifan, Ivan Skorokhodov, Victor Rong, Ziwei Liu, Xihui Liu, Jeong Joon Park, Sergey Tulyakov, Gordon Wetzstein, et al. Tc4d: Trajectory-conditioned text-to-4d generation. In *European Conference on Computer Vision*, pages 53–72. Springer, 2024. 3
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 6
- [4] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Eoeffet, Brandon Houghton, Raul Sampaio, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022. 5
- [5] Yuanhao Cai, He Zhang, Kai Zhang, Yixun Liang, Mengwei Ren, Fujun Luan, Qing Liu, Soo Ye Kim, Jianming Zhang, Zhifei Zhang, et al. Baking gaussian splatting into diffusion denoiser for fast and scalable single-stage image-to-3d generation and reconstruction. *arXiv preprint arXiv:2411.14384*, 2024. 2
- [6] Yaru Cao, Zhijian He, Lujia Wang, Wenguan Wang, Yixuan Yuan, Dingwen Zhang, Jinglin Zhang, Pengfei Zhu, Luc Van Gool, Junwei Han, et al. Visdrone-det2021: The vision meets drone object detection challenge results. In *Proceedings of the IEEE/CVF International conference on computer vision*, pages 2847–2854, 2021. 2
- [7] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017. 1
- [8] Shenyuan Gao, Siyuan Zhou, Yilun Du, Jun Zhang, and Chuang Gan. Adaworld: Learning adaptable world models with latent actions. *arXiv preprint arXiv:2503.18938*, 2025. 1, 2
- [9] Daniel Geng, Charles Herrmann, Junhwa Hur, Forrester Cole, Serena Zhang, Tobias Pfaff, Tatiana Lopez-Guevara, Yusuf Aytar, Michael Rubinstein, Chen Sun, et al. Motion prompting: Controlling video generation with motion trajectories. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1–12, 2025. 3
- [10] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for text-to-video generation. *arXiv preprint arXiv:2404.02101*, 2024. 1, 3
- [11] Chen Hou and Zhibo Chen. Training-free camera control for video generation. *arXiv preprint arXiv:2406.10126*, 2024. 3
- [12] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2
- Lora: Low-rank aanpassing van grote taalmodellen. *ICLR*, 1(2):3, 2022. 6
- [13] Ronghang Hu, Nikhila Ravi, Alexander C Berg en Deepak Pathak. Worldsheets: De wereld in een 3D-vel wikkelen voor weergavesynthese vanuit een enkele afbeelding. In *Proceedings of the IEEE/CVF Internationale Conferentie over Computer Vision*, pagina's 6034–6044, 2025. 7
- [2] Sherwin Bahmani, Xian Liu, Wang Yifan, Ivan Skorokhodov, Victor Rong, Shaoshuai Shi, Zhiota Tian, Tianyu He en Li Jiang. Geheugenforcing: Spatio-temporeel geheugen voor consistente scènegenereatie in Minecraft. *arXiv preprint arXiv:2510.03198*, 2025. 2, 3
- [15] Tianyu Huang, Wangguandong Zheng, Tengfei Wang, Yuhao Liu, Zhenwei Wang, Junta Wu, Jie Jiang, Hui Li, Rynson WH Lau, Wangmeng Zuo, et al. Voyager: Langeafstands- en wereldconsistente videodiffusie voor verkenbare 3D-scènegenereatie. *arXiv preprint arXiv:2506.04225*, 2025. 3
- [16] Longbin Ji, Lei Zhong, Pengfei Wei, and Changjian Li. Pose-traj: Pose-aware trajectory control in video diffusion. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 22776–22785, 2025. 1, 2, 3
- [17] Zhengfei Kuang, Shengqu Cai, Hao He, Yinghao Xu, Hongsheng Li, Leonidas J Guibas, and Gordon Wetzstein. Collaborative video diffusion: Consistent multi-video generation with camera control. *Advances in Neural Information Processing Systems*, 37:16240–16271, 2024. 3
- [18] Zhengfei Kuang, Shengqu Cai, Hao He, Yinghao Xu, Hongsheng Li, Leonidas J Guibas, and Gordon Wetzstein. Collaborative video diffusion: Consistent multi-video generation with camera control. *Advances in Neural Information Processing Systems*, 37:16240–16271, 2024. 2
- [19] Teng Li, Guangcong Zheng, Rui Jiang, Shuigen Zhan, Tao Wu, Yehao Lu, Yining Lin, Chuanyun Deng, Yepan Xiong, Min Chen, et al. Realcam-i2v: Real-world image-to-video generation with interactive complex camera control. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 28785–28796, 2025. 1, 3, 6, 7
- [20] Hanwen Liang, Junli Cao, Vedit Goel, Guocheng Qian, Sergei Korolev, Demetri Terzopoulos, Konstantinos N Plataniotis, Sergey Tulyakov, and Jian Ren. Wonderland: Navigating 3d scenes from a single image. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 798–810, 2025. 1, 3
- [21] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9298–9309, 2023. 2
- [22] Jiachen Lu, Ze Huang, Zeyu Yang, Jiahui Zhang, and Li Zhang. Wovogen: World volume-aware diffusion for controllable multi-camera driving scene generation. In *European Conference on Computer Vision*, pages 329–345. Springer, 2024. 2, 3
- [23] Taiming Lu, Tianmin Shu, Junfei Xiao, Luoxin Ye, Jiahao Wang, Cheng Peng, Chen Wei, Daniel Khashabi, Rama Chellappa, Alan Yuille, et al. Genex: Generating an explorable world. *arXiv preprint arXiv:2412.09624*, 2024. 1, 2
- [13] Ronghang Hu, Nikhila Ravi, Alexander C Berg en Deepak Pathak. Worldsheets: De wereld in een 3D-vel wikkelen voor weergavesynthese vanuit een enkele afbeelding. In *Proceedings of the IEEE/CVF Internationale Conferentie over Computer Vision*, pagina's 6034–6044, 2025. 7
- [2] Junchao Huang, Xinting Hu, Boyao Han, Shaoshuai Shi, Zhiota Tian, Tianyu He, and Li Jiang. Geheugenforcing: Spatio-temporeel geheugen voor consistente scènegenereatie in Minecraft. *arXiv preprint arXiv:2510.03198*, 2025. 2, 3
- [15] Tianyu Huang, Wangguandong Zheng, Tengfei Wang, Yuhao Liu, Zhenwei Wang, Junta Wu, Jie Jiang, Hui Li, Rynson WH Lau, Wangmeng Zuo, et al. Voyager: Langeafstands- en wereldconsistente videodiffusie voor verkenbare 3D-scènegenereatie. *arXiv preprint arXiv:2506.04225*, 2025. 3
- [16] Longbin Ji, Lei Zhong, Pengfei Wei, and Changjian Li. Pose-traj: Pose-bewuste trajectcontrole in videodiffusie. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pagina's 22776–22785, 2025. 1, 2, 3
- [17] Zhengfei Kuang, Shengqu Cai, Hao He, Yinghao Xu, Hongsheng Li, Leonidas J Guibas en Gordon Wetzstein. Collaboratieve videodiffusie: Consistente multi-video generatie met camerabesturing. Voortgang in Neurale Informatie Verwerkende Systemen, 37:16240–16271, 2024. 3
- [18] Zhengfei Kuang, Shengqu Cai, Hao He, Yinghao Xu, Hongsheng Li, Leonidas J Guibas en Gordon Wetzstein. Collaboratieve videodiffusie: Consistente multi-video generatie met camerabesturing. Voortgang in Neurale Informatie Verwerkende Systemen, 37:16240–16271, 2024. 2
- [19] Teng Li, Guangcong Zheng, Rui Jiang, Shuigen Zhan, Tao Wu, Yehao Lu, Yining Lin, Chuanyun Deng, Yepan Xiong, Min Chen, et al. Realcam-i2v: Real-world afbeelding-naar-video generatie met interactieve complexe camerabesturing. In *Proceedings of the IEEE/CVF Internationale Conferentie over Computer Vision*, pagina's 28785–28796, 2025. 1, 3, 6, 7
- [20] Hanwen Liang, Junli Cao, Vedit Goel, Guocheng Qian, Sergei Korolev, Demetri Terzopoulos, Konstantinos N Plataniotis, Sergey Tulyakov en Jian Ren. Wonderland: Navigeren door 3D-scènes vanuit een enkele afbeelding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pagina's 798–810, 2025. 1, 3
- [21] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, en Carl Vondrick. Zero-1-to-3: Zero-shot één afbeelding naar 3D-object. In *Proceedings of the IEEE/CVF internationale conferentie over computer vision*, pagina's 9298–9309, 2023. 2
- [22] Jiachen Lu, Ze Huang, Zeyu Yang, Jiahui Zhang, en Li Zhang. Wovogen: Wereldvolume-bewuste diffusie voor controleerbare multi-camera rijscènegenereatie. In *Europese Conferentie over Computer Vision*, pagina's 329–345. Springer, 2024. 2, 3
- [23] Taiming Lu, Tianmin Shu, Junfei Xiao, Luoxin Ye, Jiahao Wang, Cheng Peng, Chen Wei, Daniel Khashabi, Rama Chellappa, Alan Yuille, et al. Genex: Genereren van een verkenbare wereld. *arXiv preprint arXiv:2412.09624*, 2024. 1, 2

- [24] Andrew Melnik, Michal Ljubljjanac, Cong Lu, Qi Yan, Weiming Ren, and Helge Ritter. Video diffusion models: A survey. *arXiv preprint arXiv:2405.03150*, 2024. 2
- [25] Chong Mou, Mingdeng Cao, Xintao Wang, Zhaoyang Zhang, Ying Shan, and Jian Zhang. Revideo: Remake a video with motion and content control. *Advances in Neural Information Processing Systems*, 37:18481–18505, 2024. 3
- [26] Chaojun Ni, Xiaofeng Wang, Zheng Zhu, Weijie Wang, Haoyun Li, Guosheng Zhao, Jie Li, Wenkang Qin, Guan Huang, and Wenjun Mei. Wonderturbo: Generating interactive 3d world in 0.72 seconds. *arXiv preprint arXiv:2504.02261*, 2025. 1, 3
- [27] Ava Pun, Gary Sun, Jingkang Wang, Yun Chen, Ze Yang, Sivabalan Manivasagam, Wei-Chiu Ma, and Raquel Urtasun. Neural lighting simulation for urban scenes. *Advances in Neural Information Processing Systems*, 36:19291–19326, 2023. 2
- [28] Santhosh K Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alex Clegg, John Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, et al. Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai. *arXiv preprint arXiv:2109.08238*, 2021. 1
- [29] Xuanchi Ren, Tianchang Shen, Jiahui Huang, Huan Ling, Yifan Lu, Merlin Nimier-David, Thomas Müller, Alexander Keller, Sanja Fidler, and Jun Gao. Gen3c: 3d-informed world-consistent video generation with precise camera control. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6121–6132, 2025. 1, 3
- [30] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *CVPR*, 2019. 5, 7
- [31] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. In *CVPR*, 2020. 5, 7
- [32] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9339–9347, 2019. 1
- [33] Manuel-Andreas Schneider, Lukas Hölein, and Matthias Nießner. Worldexplorer: Towards generating fully navigable 3d scenes. *arXiv preprint arXiv:2506.01799*, 2025. 3
- [34] Xincheng Shuai, Henghui Ding, Zhenyuan Qin, Hao Luo, Xingjun Ma, and Dacheng Tao. Free-form motion control: Controlling the 6d poses of camera and objects in video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12449–12458, 2025. 2
- [35] Shuhan Tan, Kelvin Wong, Shenlong Wang, Sivabalan Manivasagam, Mengye Ren, and Raquel Urtasun. Scenegen: Learning to generate realistic traffic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 892–901, 2021. 2
- [36] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European conference on computer vision*, pages 402–419. Springer, 2020. 5
- [37] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018. 7
- [38] 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, Jingen 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, Tiany Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenying Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yang 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. 6, 7
- [39] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5294–5306, 2025. 7
- [40] Jiahuo Wang, Luoxin Ye, TaiMing Lu, Junfei Xiao, Jiahua Zhang, Yuxiang Guo, Xijun Liu, Rama Chellappa, Cheng Peng, Alan Yuille, et al. Evoworld: Evolving panoramic world generation with explicit 3d memory. *arXiv preprint arXiv:2510.01183*, 2025. 1, 2, 3
- [41] Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingen Zhou. Videocomposer: Compositional video synthesis with motion controllability. *Advances in Neural Information Processing Systems*, 36:7594–7611, 2023. 3
- [42] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 3
- [43] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7467–7477, 2020. 1, 2
- [44] Haoyu Wu, Diankun Wu, Tianyu He, Junliang Guo, Yang Ye, Yueqi Duan, and Jiang Bian. Geometry forcing: Marrying video diffusion and 3d representation for consistent world modeling. *arXiv preprint arXiv:2507.07982*, 2025. 3, 6, 7
- [45] Jianzong Wu, Liang Hou, Haotian Yang, Xin Tao, Ye Tian, Pengfei Wan, Di Zhang, and Yunhai Tong. Vmoba: Mixture-of-block attention for video diffusion models. *arXiv preprint arXiv:2506.23858*, 2025. 1, 2
- [46] Dejia Xu, Yifan Jiang, Chen Huang, Liangchen Song, Thorsten Gernoth, Liangliang Cao, Zhangyang Wang, and Hao Tang. Cavia: Camera-controllable multi-view video diffusion with view-integrated attention. *arXiv preprint arXiv:2410.10774*, 2024. 3
- [47] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, en Sylvain Gelly. Naar nauwkeurige generatieve modellen van video: Een nieuwe metriek & uitdagingen. *arXiv preprint arXiv:1812.01717*, 2018. 7
- [48] 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, Jingen 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, Tiany Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenying Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yang Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yi-tong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, en Ziyu Liu. Wan: Open en geavanceerde grootschalige generatieve videomodellen. *arXiv preprint arXiv:2503.20314*, 2025. 6, 7
- [49] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, en David Novotny. Vggt: Visuele geometrie gegronde transformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pagina's 5294–5306, 2025. 7
- [50] Jiahuo Wang, Luoxin Ye, TaiMing Lu, Junfei Xiao, Jiahua Zhang, Yuxiang Guo, Xijun Liu, Rama Chellappa, Cheng Peng, Alan Yuille, et al. EvoWorld: Evoluerende panoramische wereldgeneratie met expliciete 3D-geheugen. *arXiv preprint arXiv:2510.01183*, 2025. 1, 2, 3
- [51] Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiniu Wang, Yingya Zhang, Dayou Chen, Jiniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, en Jingen Zhou. Videocomposer: Compositorische videosynthesen met bewegingscontroleerbaarheid. Voortgang in Neurale Informatie Verwerkende Systemen, 36:7594–7611, 2023. 3
- [52] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, en Ying Shan. Motionctr: Een verenigde en flexibele bewegingscontroller voor videogeneratie. In *ACM SIGGRAPH 2024 Conference Papers*, pagina's 1–11, 2024. 3
- [53] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, en Justin Johnson. Synsin: End-to-end weergavesynthese vanuit een enkele afbeelding. In *Proceedings van de IEEE/CVF-conferentie over computervisie en patroonherkenning*, pagina's 7467–7477, 2020. 1, 2 – .
- [54] Haoyu Wu, Diankun Wu, Tianyu He, Junliang Guo, Yang Ye, Yueqi Duan, en Jiang Bian. Geometrie Dwang: Het combineren van videodiffusie en 3D-representatie voor consistent wereldmodellering. *arXiv preprint arXiv:2507.07982*, 2025. 3, 6, 7
- [55] Jianzong Wu, Liang Hou, Haotian Yang, Xin Tao, Ye Tian, Pengfei Wan, Di Zhang, en Yunhai Tong. Vmoba: Mengsel-van-blok aandacht voor videodiffusiemodellen. *arXiv preprint arXiv:2506.23858*, 2025. 1, 2
- [56] Dejia Xu, Yifan Jiang, Chen Huang, Liangchen Song, Thorsten Gernoth, Liangliang Cao, Zhangyang Wang, en Hao Tang. Cavia: Camera-controleerbare multi-view videodiffusie met weergave-geïntegreerde aandacht. *arXiv preprint arXiv:2410.10774*, 2024. 3

- [47] Dejia Xu, Weili Nie, Chao Liu, Sifei Liu, Jan Kautz, Zhangyang Wang, and Arash Vahdat. Camco: Camera-controllable 3d-consistent image-to-video generation. *arXiv preprint arXiv:2406.02509*, 2024. 3
- [48] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2024. 3
- [49] Xiaoda Yang, Jiayang Xu, Kaixuan Luan, Xinyu Zhan, Hongshun Qiu, Shijun Shi, Hao Li, Shuai Yang, Li Zhang, Checheng Yu, et al. Omnicam: Unified multimodal video generation via camera control. *arXiv preprint arXiv:2504.02312*, 2025. 2
- [50] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4578–4587, 2021. 1, 2
- [51] Hong-Xing Yu, Haoyi Duan, Charles Herrmann, William T Freeman, and Jiajun Wu. Wonderworld: Interactive 3d scene generation from a single image. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5916–5926, 2025. 1, 3
- [52] Mark YU, Wenbo Hu, Jinbo Xing, and Ying Shan. Trajectorycrafter: Redirecting camera trajectory for monocular videos via diffusion models. *arXiv preprint arXiv:2503.05638*, 2025. 3
- [53] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan, and Yonghong Tian. Viewcrafter: Taming video diffusion models for high-fidelity novel view synthesis. *arXiv preprint arXiv:2409.02048*, 2024. 3
- [54] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 7
- [55] Zhenghao Zhang, Junchao Liao, Menghao Li, Zuozhuo Dai, Bingxue Qiu, Siyu Zhu, Long Qin, and Weizhi Wang. Tora: Trajectory-oriented diffusion transformer for video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 2063–2073, 2025. 3
- [56] Zhongwei Zhang, Fuchen Long, Zhaofan Qiu, Yingwei Pan, Wu Liu, Ting Yao, and Tao Mei. Motionpro: A precise motion controller for image-to-video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 27957–27967, 2025. 3
- [57] Guangcong Zheng, Teng Li, Rui Jiang, Yehao Lu, Tao Wu, and Xi Li. Cami2v: Camera-controlled image-to-video diffusion model. *arXiv preprint arXiv:2410.15957*, 2024. 1, 3, 7
- [58] Mengqi Zhou, Yuxi Wang, Jun Hou, Shougao Zhang, Yawei Li, Chuachen Luo, Junran Peng, and Zhaoxiang Zhang. Scenex: Procedural controllable large-scale scene generation. *arXiv preprint arXiv:2403.15698*, 2024. 2
- [59] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 2, 5
- [60] Yunsong Zhou, Michael Simon, Zhenghao Peng, Sicheng Mo, Hongzi Zhu, Minyi Guo, and Bolei Zhou. Simgen: Simulator-conditioned driving scene generation. *Advances in Neural Information Processing Systems*, 37:48838–48874, 2024. 2
- [61] Zhenghong Zhou, Jie An, and Jiebo Luo. Latent-reframe: Enabling camera control for video diffusion models without training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12779–12789, 2025. 3
- [62] Dong Zhuo, Wenzhao Zheng, Jiahe Guo, Yuqi Wu, Jie Zhou, and Jiwen Lu. Streaming 4d visual geometry transformer. *arXiv preprint arXiv:2507.11539*, 2025. 6
- [63] Dejia Xu, Weili Nie, Chao Liu, Sifei Liu, Jan Kautz, Zhangyang Wang, en Arash Vahdat. Camco: Camera-gestuurde 3D-consistente afbeelding-naar-video generatie. *arXiv preprint arXiv:2406.02509*, 2024. 3
- [64] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jias hi Feng, en Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2024. 3
- [65] Xiaoda Yang, Jiayang Xu, Kaixuan Luan, Xinyu Zhan, Hongshun Qiu, Shijun Shi, Hao Li, Shuai Yang, Li Zhang, Checheng Yu, et al. Omnicam: Geünificeerde multimodale video generatie via camera controle. *arXiv preprint arXiv:2504.02312*, 2025. 2
- [66] Alex Yu, Vickie Ye, Matthew Tancik, en Angjoo Kanazawa. pixelnerf: Neurale stralingsvelden van één of enkele afbeeldingen. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pagina's 4578–4587, 2021. 1, 2
- [67] Hong-Xing Yu, Haoyi Duan, Charles Herrmann, William T Freeman, en Jiajun Wu. Wonderworld: Interactieve 3D scène generatie vanuit een enkele afbeelding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pagina's 5916–5926, 2025. 1, 3
- [68] Mark YU, Wenbo Hu, Jinbo Xing, en Ying Shan. Trajectorycrafter: Het herleiden van camera trajecten voor monoculaire video's via diffusiemodellen. *arXiv preprint arXiv:2503.05638*, 2025. 3
- [69] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan, en Yonghong Tian. Viewcrafter: Het temmen van videodiffusiemodellen voor hoogwaardig nieuwe weergave synthese. *arXiv preprint arXiv:2409.02048*, 2024. 3
- [70] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, en Oliver Wang. De onredelijke effectiviteit van diepe kenmerken als een perceptuele maatstaf. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pagina's 586–595, 2018. 7
- [71] Zhenghao Zhang, Junchao Liao, Menghao Li, Zuozhuo Dai, Bingxue Qiu, Siyu Zhu, Long Qin, en Weizhi Wang. Tora: Trajectorie-georiënteerde diffusie transformator voor video generatie. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pagina's 2063–2073, 2025. 3
- [72] Zhongwei Zhang, Fuchen Long, Zhaofan Qiu, Yingwei Pan, Wu Liu, Ting Yao, en Tao Mei. Motionpro: Een precieze bewegingscontroller voor afbeelding-naar-video generatie. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pagina's 27957–27967, 2025. 3
- [73] Guangcong Zheng, Teng Li, Rui Jiang, Yehao Lu, Tao Wu, en Xi Li. Cami2v: Camera-gestuurde afbeelding-naar-video diffusiemodel. *arXiv preprint arXiv:2410.15957*, 2024. 1, 3
- [74] Mengqi Zhou, Yuxi Wang, Jun Hou, Shougao Zhang, Yawei Li, Chuachen Luo, Junran Peng, en Zhaoxiang Zhang. Scenex: Procedureel controleerbare grootschalige scène generatie. *arXiv preprint arXiv:2403.15698*, 2024. 2
- [75] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, en Noah Snavely. Stereo vergroting: Leren van weergave synthese met behulp van multiplane afbeeldingen. *arXiv preprint arXiv:1805.09817*, 2018. 2, 5
- [76] Yunsong Zhou, Michael Simon, Zhenghao Peng, Sicheng Mo, Hongzi Zhu, Minyi Guo, en Bolei Zhou. Simgen: Simulator-geconditioneerde generatie van rijscènes. *Voortgang in Neurale Informatie Verwerkende Systemen*, 3748874, 202448838:2–.
- [77] Zhenghong Zhou, Jie An, en Jiebo Luo. Latent-reframe: Camera controle mogelijk maken voor videodiffusiemodellen zonder training. In *Proceedings of the IEEE/CVF Internationale Conferentie over Computer Vision*, pagina's 12779–12789, 2025. 3
- [78] Dong Zhuo, Wenzhao Zheng, Jiahe Guo, Yuqi Wu, Jie Zhou, en Jiwen Lu. Streaming 4d visuele geometrie transformator. *arXiv preprint arXiv:2507.11539*, 2025. 6