# MOTIONCRAFT:
# Physics-based Zero-Shot Video Generation

**Luca Savant Aira**[*]　　　　　**Antonio Montanaro**[*]

**Emanuele Aiello,**　　　　**Diego Valsesia**　　　　**Enrico Magli**

Politecnico di Torino
`{name.surname}@polito.it`

## Abstract

Generating videos with realistic and physically plausible motion is one of the main recent challenges in computer vision. While diffusion models are achieving compelling results in image generation, video diffusion models are limited by heavy training and huge models, resulting in videos that are still biased to the training dataset. In this work we propose MotionCraft, a new zero-shot video generator to craft physics-based and realistic videos. MotionCraft is able to warp the noise latent space of an image diffusion model, such as Stable Diffusion, by applying an optical flow derived from a physics simulation. We show that warping the noise latent space results in coherent application of the desired motion while allowing the model to generate missing elements consistent with the scene evolution, which would otherwise result in artefacts or missing content if the flow was applied in the pixel space. We compare our method with the state-of-the-art Text2Video-Zero reporting qualitative and quantitative improvements, demonstrating the effectiveness of our approach to generate videos with finely-prescribed complex motion dynamics. Project page: `https://mezzelfo.github.io/MotionCraft/`.

## 1 Introduction

As human beings, we have always exploited our creativity to generate art, in different forms such as visual art, music or poetry. In vision, we are often inspired by the natural world since our visual system continuously acquire images perceived as a video sequence. Indeed, videos or movies are one of the best visual stimuli since they contain images, motion and audio.

Recent generative models for still images based on diffusion models [24, 27, 28] achieved remarkable results with quality almost indistinguishable from real images. It is therefore clear that the next big goal is video generation. However, it seems that including the dimension of time remains challenging. Some works such as Sora [5] achieve astonishing temporal consistency and photorealism at the expense of enormous computational and data requirements. Moreover, we argue that fine-grained control over the motion dynamics is impossible with a simple text prompt. If one wants to synthesize a video according to some precise physical dynamics, they would not be able to do it with current models. Interestingly, explicitly controlling the motion dynamics also allows to decouple temporal evolution from content generation. Indeed, explicitly injecting the physics of the real world as motion dynamics allows to develop more parsimonious models, that do not need to brute-force learn them from data.

---

[*]indicates equal contribution.

Figure 1: Melting man simulation. Top: MOTIONCRAFT; Bottom: T2V0 [20]. MOTIONCRAFT uses a fluid dynamics simulation to warp noise latents and synthetize video frames. T2V0 is unable to simulate the evolution of the melting statue and simply moves the object towards the bottom of the frame.

For this reason, in this paper, we investigate the possibility to create a zero-shot video generation model that only requires a pretrained still image generator and knowledge of physical laws regarding motion. Indeed, since videos are temporal sequences of images correlated by physical laws, we only need to devise a way to include physical laws in the diffusion prior to animate a starting image. We thus advocate for physics simulators as appropriate sources of motion, output as a sequence of optical flows, while also being completely user-controllable, plausible, and explainable.

We propose MOTIONCRAFT, a physics-based zero-shot video generator that uses optical flow extracted from a physical simulation to warp the noise latent space of a pretrained image diffusion model to generate videos with complex dynamics without the need to train anything. While using a projection of motion onto the camera plane as a pixelwise displacement field (optical flow) may seem limiting due to the fact that, if applied in the pixel space, it would not be able to synthesise novel coherent content but only displace pixels, the trick lies in its application in the noise latent domain. Backed by evidence that motion vectors correlate between pixel and noise space, warping of the latter by means that MOTIONCRAFT allows to simultaneously apply the desired motion and exploit the powerful image prior of the generative model. This is capable of adapting the scene to the prescribed motion without significant artefacts, generate novel content and shows impressive global consistency (reflections, illumination, etc., consistent with the desired evolution).

We present quantitative and qualitative experimental results where we show that our zero-shot MOTIONCRAFT is capable of synthesising realistic videos with finely controlled temporal evolution governed by fluid-dynamics equations, rigid body physics, and multi-agent interaction models, while zero-shot state-of-art techniques cannot.

## 2   Related work

**Diffusion Based Video Generation**    Video Generation [2] is a longstanding problem in computer vision aiming to learn the distribution of and synthesise realistic videos. Recently, text-based Denoising Diffusion Probabilistic Models (DDPM) [27, 30] have been studied to tackle this challenge delivering impressive results. These approaches include Sora [5], Video Diffusion models [17], Imagen-video [16] and Align your Latents [3]. They require sophisticated spatio-temporal denoising architectures at the expense of huge computational requirements and large amounts of paired text-video data for training. To reduce the data requirements, different approaches investigate few-shot and unsupervised learning techniques. Make-a-Video [26] proposes an unsupervised training with only videos, coupled with a retrieval strategy to sample using text. On the other hand, Ni et al. [22] train a diffusion-based optical flow generator that outputs a flow conditioned on a reference image and a textual prompt, that reduces the computational burden of generating videos by training the diffusion process on small flow fields. Differently from them, our approach is zero-shot and we do not train anything.

To the best of our knowledge, Text-to-video-Zero [20] and Generative Rendering [6] are the only zero-shot video generators. However, Generative Rendering (concurrent work, with no code available)

Figure 2: A qualitative example of the image and latent flows correlation. This figure shows, from left to right, (a) the first RGB frame, (b) the second RGB frame superimposed with the estimated flow in the RGB domain, (c) the first latent frame, (d) the second latent frame superimposed with the estimated flow in the latent domain and (e) the correlation map of the two non-zero flows.

has significant extra requirements beyond Stable Diffusion (SD) as image generator, in the form of a depth-conditioned ControlNet [35], and a 3D mesh manually animated, leveraging UV maps to render the scene. Moreover, Generative Rendering cannot render fluids, since they are difficult to represent as 3D meshes.

In this paper, we compare our method to Text-to-video-Zero (T2V0), as zero-shot video generator baseline. T2V0 applies a constant shift (with a fixed direction) to the initial latent noise of SD, sampling each frame sequentially by means of DDPM. As shown in our work, since the motion in the noise latent space directly translates into the motion of the pixel space, the generated videos result in a overall shift in the same fixed direction. The largest part of the motion is caused by the stochastic fluctuations of the DDPM sampling strategy leading to unnatural motion and inconsistency of the objects in the different frames. On the contrary, in this work, we avoid the use of a constant warping operation derived from physics simulation flows in the latent space in order to incorporate complex motion dynamics.

**Diffusion Based Video and Image Editing**    Recently, different methods exploit the prior of text-to-image diffusion models for video editing. In particular, Tune-A-Video [34] finetunes a text-to-image diffusion model to edit a video. They start from the inverted frames in the latent space and use the text prompt as an editing tool. Pix2Video [7] employs a self-attention injection mechanism to edit videos using a pretrained image diffusion model.

Other methods use the optical flow to edit reference images or videos. Motion Guidance [10] leverages a user defined optical flow that allow zero-shot image editing. It works by guiding the diffusion sampling process with the gradient from a pretrained optical flow network via a guidance loss. LatentWarp [1] and TokenFlow [11], use an optical flow estimated from a reference video to warp the latent space of the diffusion model to achieve consistent editing. These methods leverage both diffusion models priors and other components such as ControlNet for structural control, and trained flow estimators such as RAFT [31]. Alternatively, we propose a zero-shot video generation method, using only vanilla SD. This means that MOTIONCRAFT does not require a reference video but it can animate an image, generated by the SD model or obtained by inverting a real one. Moreover, the physics simulations allow to generate different videos from the same starting image.

## 3   Method

This section describes MOTIONCRAFT, a zero-shot video generation method, where the meaning of "zero-shot" is twofold: we do not train or finetune any component of the text-to-image diffusion model, nor we do not use reference video or optical flow estimators as starting point. In the following, we used used Stable Diffusion as pretrained text-to-image model.

### 3.1   Optical Flow is preserved in the Latent Space of Stable Diffusion

Our proposed method stems from a key observation: the optical flow estimated between two frames in the pixel space is correlated with the flow estimated between the corresponding noise latent representations of SD. We conjecture that this is related to the specific design of the SD variational auto-encoder and denoiser architectures. In fact, by largely using convolution operations, they enforce a locality prior which preserves spatial information to some extent.
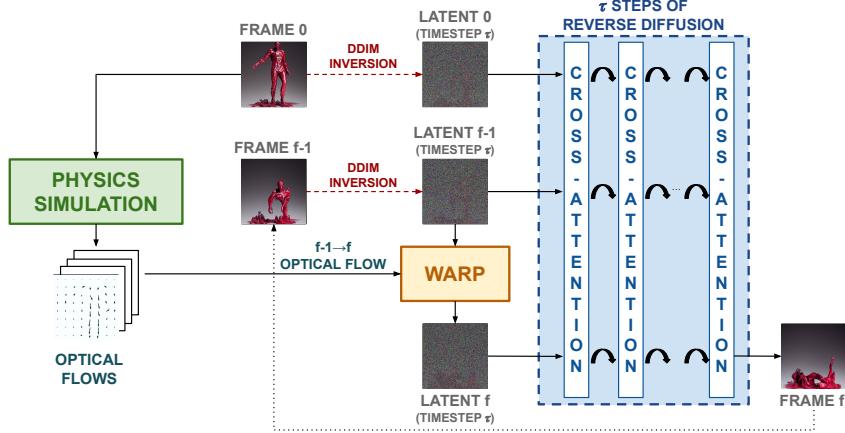
Figure 3: MOTIONCRAFT overview. A video is generated from a starting image using a pretrained still image generative model by warping noise latents according to an optical flow description of the motion to be synthesised.

In order to empirically investigate this phenomenon, we conducted a quantitative experiment using the MSU Video Frame Interpolation Benchmark dataset [12], considering only real videos. For each pair of consecutive video frames, the following steps have been taken. We first estimate the optical flow in the RGB space by using a well-established method, based on the Gunnar Farneback's algorithm, provided by OpenCV [19]. Then, we compute the noise latent representations of the two frames, first encoding the image in the variational autoencoder (VAE) of SD at timestep $\tau = 0$, followed by DDIM inversion [28] up to timestep $\tau = 400$ (same value for all experiments in this work, empirically determined). Finally, a correlation coefficient based on cosine similarity is computed between the optical flows estimated in the RGB and noise latent spaces. The resulting correlations are then averaged across all pairs of consecutive frames in the dataset, obtaining an average value of 0.727, which indicates a strong correlation between the optical field in the RGB and noise latent domains. An example of this experiment is presented in Fig. 2, showcasing the two estimated flows in the image and latent space and their correlation.

## 3.2 Physics-based zero shot video generation

Based on the analysis presented in the previous section, we propose a novel zero-shot video generation method, named MOTIONCRAFT, where an image (real or generated), serving as a starting frame $I^0$, is animated according to a physical simulation, by means of a (possibly time-varying) optical-flow generator $\mathcal{W}$ in the noise latent space. The outcome is a video made of $N$ frames $I^0, \ldots, I^{N-1}$ that follows the motion prescribed by the physical simulation and evolves the content of the first frame coherently. Inspired by the previous observation, this animation is obtained by warping the noisy latent representation of an image in the latent diffusion space. Regarding the physics simulation for the optical flow generation, we use different libraries to simulate different physics, as explained in the experimental section, such as fluid dynamics, rigid motion and multi-agent systems. It is also possible, albeit not shown in this paper, to use animation software to generate the required optical flows.

Fig. 3 illustrates an overview of MOTIONCRAFT highlighting the autoregressive generation of the video. At each iteration $f \geq 1$, the frame $I^f$ is generated using only the information contained in the first frame $I^0$ and the previous frame $I^{f-1}$. Given this Markovian structure, MOTIONCRAFT is characterized by $\mathcal{O}(1)$ space complexity and $\mathcal{O}(N)$ time complexity with respect to the total number $N$ of frames to be generated. More in detail, first, the two RGB frames $I^0, I^{f-1}$ are encoded into the latent space and they are independently inverted with the reversed DDIM sampling scheme up to a fixed diffusion timestep $\tau$, obtaining $z_\tau^0$, and $z_\tau^{f-1}$, respectively. Then, the optical flow warping operator $\mathcal{W}^{f-1 \to f}$ prescribed by the physical simulation is applied to $z_\tau^{f-1}$, obtaining $\zeta_\tau^f$. Finally, the next RGB frame $I^f$ is generated by performing $\tau$ steps of reverse diffusion using the DDIM sampling scheme with a novel cross-frame attention mechanism and a novel spatial noise map $\eta^f$ weighting technique, explained below. Furthermore, we exploit the classifier-free guidance (CFG)

**Algorithm 1:** Pseudocode of MOTIONCRAFT

**Input**   : $I^0, \mathcal{W}, \eta, \mathcal{P}, \mathcal{P}_\emptyset$
**Output** : $I^0, \ldots, I^{N-1}$

1  **for** $f = 1$ **to** $N - 1$ **do**
2      $z_0^{f-1} = \mathcal{E}(I^{f-1})$ ;                                              // Encode the frame
3      **for** $t = 0$ **to** $\tau - 1$ **do**                                       // Inversion loop
4         $\hat{\epsilon} \leftarrow \epsilon_t(z_t^{f-1}, \mathcal{P}; \{z_t^{f-1}\})$ ;                   // Self-Attention, No MCFA
5         $z_{t+1}^{f-1} \leftarrow \text{DDIMInversion}_{t \to t+1}(z_t^{f-1}, \hat{\epsilon}, 0)$ ;         // $\eta = 0 \iff$ DDIM
6      **end**
7      $\zeta_\tau^f = \mathcal{W}^{f-1 \to f}(z_\tau^{f-1})$ ;                          // Warp the latent
8      **for** $t = \tau - 1$ **to** $0$ **do**                                       // Generation loop
9         $\hat{\epsilon}_\mathcal{P} \leftarrow \epsilon_t(\zeta_{t+1}^f, \mathcal{P}; \{z_t^0, z_t^{f-1}\})$ ;           // MCFA with $I^0$ and $I^{f-1}$
10        $\hat{\epsilon}_\emptyset \leftarrow \epsilon_t(\zeta_{t+1}^f, \mathcal{P}_\emptyset; \{z_t^0, z_t^{f-1}\})$ ;           // MCFA with $I^0$ and $I^{f-1}$
11        $\hat{\epsilon} \leftarrow \hat{\epsilon}_\emptyset + \gamma(\hat{\epsilon}_\mathcal{P} - \hat{\epsilon}_\emptyset)$ ;           // Classifier-free guidance
12        $\zeta_t^f \leftarrow \text{DDIM}_{t+1 \to t}(\zeta_{t+1}^f, \hat{\epsilon}, \eta^f)$ ;           // Perform Spatial-$\eta$
13     **end**
14     $I^f \leftarrow \mathcal{D}(\zeta_0^f)$ ;                                              // Decode the latent
15 **end**
16 **return** $I^0, \ldots, I^{N-1}$

technique for generation proposed in [14], with $\mathcal{P}$ and $\mathcal{P}_\emptyset$ being the positive and negative prompt, respectively, and $\gamma > 1$ being the strength of the CFG. More details can be found in the Appendix A.

Algorithm 1 reports the pseudocode of MOTIONCRAFT. Lines $2 - 6$ include the DDIM inversion up to timestep $\tau$. Starting current frame $I^{f-1}$ that was previously generated, in line 2 we embed it with the VAE encoder $\mathcal{E}$, obtaining $z_0^{f-1}$. Then we apply DDIM inversion on $z_0^{f-1}$ for $\tau$ timesteps (line $3 - 6$). This involves the UNet with the standard self-attention (note the repetition of the noisy latent $z_t^{f-1}$) and the positive prompt $\mathcal{P}$. As briefly reported in [21], we have also experienced that DDIM inversion is not compatible with CFG; hence, during the inversion, we do not use the negative prompt $\mathcal{P}_\emptyset$. The resulting estimated noise is used in line 5 for applying the DDIM inversion step (note that the $\eta = 0$, so pure DDIM is performed). Upon completion of the DDIM inversion process, we obtain $z_\tau^{f-1}$, the noisy latent corresponding to the frame $I^{f-1}$.

In line 7, the optical flow warping operator $\mathcal{W}^{f-1 \to f}$ is applied to the noisy latent of the current frame $z_\tau^{f-1}$ to obtain a new noisy latent $\zeta_\tau^f$ that will generate the successive frame. Finally, in lines $8 - 14$, the frame is generated. During this generation phase we use CFG to increase the quality of the generated images, hence also the negative prompt $\mathcal{P}_\emptyset$ is used. To create new content while preserving the original image, we propose two direct generalization of two known techniques: the multiple cross-frame attention (MCFA) mechanism and a spatial noise map weighting (Spatial-$\eta$).

The MCFA technique generalizes the Cross Frame Attention (CFA) [20], as it enables the to-be-generated frame to attend to an arbitrary number of frames. We choose to attend to the first frame and the previous frame (as shown in lines $9 - 10$ of Alg. 1 and Fig. 3) to ensure long-range and short-range temporal consistency, respectively. MCFA intervenes in all the self-attention blocks of the SD UNet, by replacing the keys and values, that are originally computed from projections of the generating frame features, with the ones computed from the attended frames.

We also propose Spatial-$\eta$ (line 12), that is a novel technique that enables to choose, on a pixel-by-pixel basis, whether to use DDIM or DDPM as a sampling scheme. This enables the usage of DDPM in regions of the images where novel content should be created (for example, when a new part of an object is entering the scene), while using DDIM in the other regions to ensure consistency and determinism where the already-present content is just moving. Note that this spatial map $\eta^f$ can be obtained in multiple ways from the physical simulation. For example, $\eta^f$ can be set to 1 in regions of the image where the flow is not well-defined (pointing outside of the image boundaries) or in regions where the optical flow field has discontinuities.

Table 1: Quantitative results.

| | | Frame Consistency | | Motion Consistency | |
|---|---|---|---|---|---|
| | | T2V0 [20] | **MOTIONCRAFT** | T2V0 [20] | **MOTIONCRAFT** |
| Fluid | Dragons | 0.9664 | **0.9991** | 0.6846 | **0.9637** |
| | Melting Man | 0.9463 | **0.9566** | 0.7817 | **0.8252** |
| Rigid Body | Satellite Scan | 0.9588 | **0.9875** | 0.2852 | **0.9219** |
| | Revolving Earth | **0.9812** | 0.9696 | **0.7213** | 0.6783 |
| Agents | Birds | 0.9765 | **0.9968** | 0.8973 | **0.9385** |
| Average | | $0.9658 \pm 0.01$ | **0.9819** $\pm 0.02$ | $0.6740 \pm 0.23$ | **0.8655** $\pm 0.12$ |



Figure 4: Rigid motion simulation: satellite orbit. Top: MOTIONCRAFT; Bottom: T2V0 [20].

## 4 Experimental results

### 4.1 Experimental setting

In this section, we show examples of video generation based on different physics simulations: rigid body motion, fluid dynamics and multi-agent systems. Given an optical flow, we apply it on the SD latent space using MOTIONCRAFT. Then, we compare our method to Text2Video-Zero [20] that, to the best of our knowledge, is the only diffusion-based zero-shot method for video generation.

We show qualitative results in Figs. 1, 4, 5, 6, 7, which we separately describe in the following sections. Table 1 reports two metrics to evaluate the quality of the generated videos. As done in previous works, we use the *Frame Consistency* metric, defined as the average cosine similarity of the CLIP embeddings of consecutive frames. However, this metric presents some limitations, as CLIP focuses on high-level semantic features and not on low-level details, resulting in high correlations even if the content changes but its semantics do not (as an example, see the video generated by T2V0 of the dragon in Fig. 6, which has a *Frame Consistency* of 0.97 even if the dragons are not the same dragons in each frame). To overcome some of these limitations, we propose a novel metric, named *Motion Consistency*, that measures how similar two frames are while accounting for the motion between them. We start from the observation that, if an object moves through the scene, its textures should remain almost the same, and, if we know its flow, we can bring back that object to overlap with its starting position. Then we can apply a similarity distance between the initial image and the next frame brought back by the reversed flow. Given two consecutive frames, we use a high-quality flow estimator (RAFT [31]) to estimate the optical flow between them and apply it on the second frame to reverse the motion. Then we compute the SSIM metric [33] on the first frame and the registered one.

### 4.2 Rigid Body flows

Fig. 4 shows a pivotal example where MOTIONCRAFT can be directly compared to the state-of-the-art T2V0, as in this case we use an optical flow equivalent to a their proposed shift along the vertical axis. This example shows a video generated starting from a satellite view of a city, and, by simulating the rectilinear motion of the satellite, new portions of the city appear from the top of the image. While
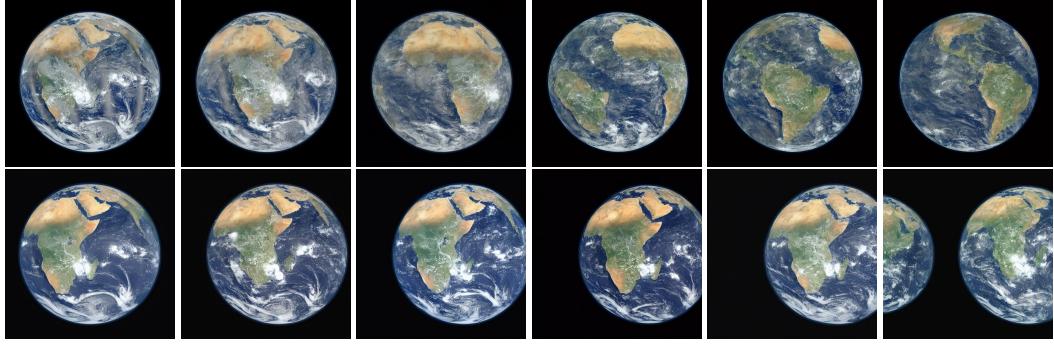
Figure 5: Rigid motion simulation: revolving Earth. Top: MOTIONCRAFT; Bottom: T2V0 [20].



Figure 6: Fluid simulation: dragon fire. Top: MOTIONCRAFT; Bottom: T2V0 [20].

T2V0 struggles with keeping temporal consistency, even with large structural elements (e.g., the river), MOTIONCRAFT is able to coherently scroll down the already present part of the city, while also generating new plausible content in the top part of the frames.

A similar case study is the Earth rotation in Fig. 5. Here, the optical flow is obtained by simulating a rotating sphere that was fitted to the first frame while keeping track of the starting and ending position of each point. As the Earth rotates, a slice disappears from one side and a new one needs to be generated on the opposite side. Thanks to the powerful natural image prior of SD, MOTIONCRAFT is able to autonomously generate other continents in the correct position, even if the text prompt contains no reference about them (see Appendix D for all the text prompts used in this paper). On the contrary, T2V0 is not able to rotate the Earth consistently while creating new content, as visible in the same Fig. 5.

### 4.3 Fluids dynamics

In this set of experiments, we use the $\Phi$-flow [18] library to simulate fluid dynamics (by numerically solving Navier-Stokes equations) with the shape and position provided by the first frame $I^0$. Moreover, we can set up the simulation in different ways, depending on the numerical solvers, i.e. *Eulerian* (particle-based) or *Lagrangian* (grid-based), we can add rigid obstacles to the fluid or we can define a initial velocity and force fields. All these different options result in videos that can have the same starting frame but differ in their evolution according to the simulation constraints. We extract the velocity field of the simulation as a proxy for the optical flow. Examples of the velocity field can be seen in Appendix C.

Fig. 6 shows a fluid simulation of two dragons breathing fire. We can approximate the two initial fire breaths with two centered smoke balls, obtaining a binary mask that will be fed to the simulation. At this point, we run the simulation, solving the Navier-Stokes equations by sequentially evaluating advection, diffusion and pressure. The vorticity and the expansion of the smoke is due to the buoyancy force set in the desired direction, that in this case is such that the two balls cross near the middle of the image.

7

Figure 7: Multi-agent system simulation: bird flock. Top: MOTIONCRAFT; Bottom: T2V0 [20].

The figure shows that MOTIONCRAFT produces a consistent scene with a realistic animation of the fire breaths. Moreover, the global scene illumination seems to change accordingly, and a realistic occlusion of a dragon due to smoke gradually appears. This is mainly due to the MCFA mechanism, as we ablate in Sec. 4.5. In T2V0, the scene is not temporally consistent and shows increasingly more artefacts, such as color shifts or the fact that the right dragon changes with time, while the left one even disappears.

A similar analysis can be conducted for Fig. 1, where a simulation of a melting statue is shown. We can see that the generated video includes bouncing of parts on the ground before the fluid settles.

## 4.4 Multi-agent systems

Multi-agent systems are another interesting family of simulated dynamics. A simple agents model is the *Boids* model [23], consisting of a set of point-like agents (named boids) that move according to three steering behaviour rules: separation, as boids avoid collisions with nearby agents by steering away from them, alignment, as boids align their direction with that of nearby agents, and cohesion, as boids move towards the average position of nearby agents to stay together as a group. To simulate this system we used the agentpy [9] library, in which the number of agents, the simulation time-steps and different physical parameters related to the steering rules can be chosen.

An example is shown in Fig. 7, generating the temporal evolution of a flock of birds. As SD is not able to generate images with a controllable number of agents in specified positions, we start from an image where there is a single agent (a bird in the example). Then, we extract the corresponding latent vector patch with the attention map [8] related to the CLIP token containing the word "bird", and clone it to the simulated positions of the other agents. At this point, we evolve the frames according to the optical flow derived from the simulation velocity field.

While MOTIONCRAFT produces a realistic flock motion, T2V0 motion is not consistent and the number of birds changes in each frame.

## 4.5 Ablations

In this section, we ablate the contribution of the most important components/hyperparameters in the proposed pipeline. First, we start from investigating the impact of the cross-attention mechanism by comparing four different variants: i) each frame attends to itself (no MCFA); ii) each frame attends to the previous frame; iii) each frame attends to the first frame; iv) each frame attends to both the previous frame and the first frame (proposed MCFA). Visual results are shown in Fig. 8. As can be seen, the MCFA mechanism is necessary to generate plausible frames; moreover, attending only to the first frame reduces the overall motion, (e.g., always showing Africa as in the first frame), while only attending to the previous frame reduces color consistency. Overall, we demonstrate that the proposed MCFA, attending to both the first and the previous frame, represents the optimal solution to keep global consistency with the initial image and local consistency with the preceding frame.

Fig. 9 shows the ablation of the Spatial-$\eta$ weighting technique. As shown, being able to sample with DDPM in some parts of the image is crucial in order to generate novel plausible content. Indeed, DDPM adds, during each reverse diffusion step, random white noise to the latent. We suppose that
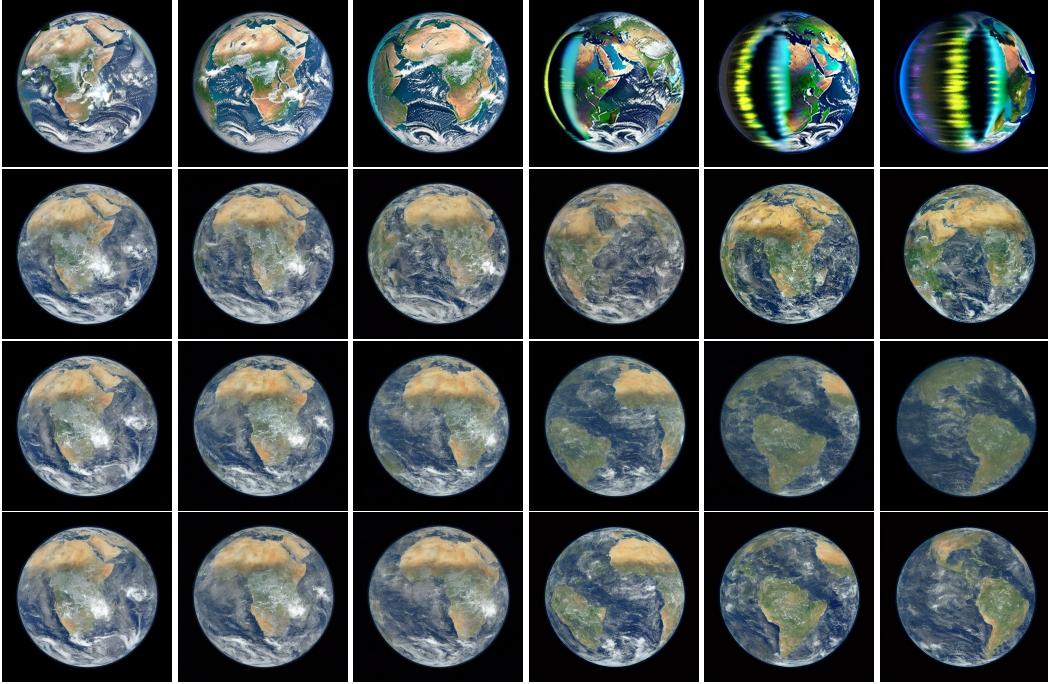
Figure 8: Ablation - Cross-Frame attention. First row: no cross frame attention; Second Row: Attend only to the initial frame; Third Row: Attend only to the previous frame; Fourth Row: Attend to the initial and preceding frame (ours).
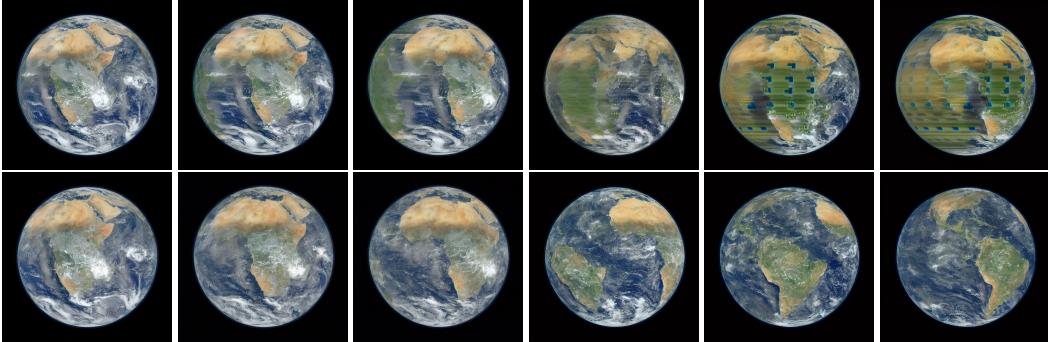


Figure 9: Ablation - Spatial-$\eta$. First Row: $\eta = 0$; Second Row: Spatial-$\eta$ on.

this allows to better sample from the real distribution, avoiding artefacts other components of the method, such as the warping operator or the MCFA, would otherwise introduce.

Finally, we ablated the partial inversion process, i.e., lines 2-6 in Alg. 1. Without the DDIM inversion, textures and details generated by SD cannot be brought into the next frame, resulting in corrupted videos. Visual results can be found in the Appendix B.3.

## 5 Conclusions

In this work, we have presented MOTIONCRAFT, a novel zero-shot approach for video generation. Our method allows to generate realistic videos with the image prior of Stable Diffusion and a physically-derived optical flow, without any additional training. MOTIONCRAFT warps the noise latent space according to the prescribed flow, and with a modified sampling process exploiting multi-frame cross-attention and the spatial-$\eta$ variable sampling scheme generates novel plausible contents following the prescribed motion and temporally consistent. For the evaluations of the results, we relied on a standard metric and a proposed one, showing that our method is not only qualitatively but also quantitatively superior to the state-of-the-art of zero-shot video generation.

# References

[1] Yuxiang Bao, Di Qiu, Guoliang Kang, Baochang Zhang, Bo Jin, Kaiye Wang, and Pengfei Yan. Latentwarp: Consistent diffusion latents for zero-shot video-to-video translation. *arXiv preprint arXiv:2311.00353*, 2023.

[2] Rishika Bhagwatkar, Saketh Bachu, Khurshed Fitter, Akshay Kulkarni, and Shital Chiddarwar. A review of video generation approaches. In *2020 International Conference on Power, Instrumentation, Control and Computing (PICC)*, pages 1–5, 2020. doi: 10.1109/PICC51425.2020.9362485.

[3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023.

[4] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.

[5] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. URL `https://openai.com/research/video-generation-models-as-world-simulators`.

[6] Shengqu Cai, Duygu Ceylan, Matheus Gadelha, Chun-Hao Paul Huang, Tuanfeng Yang Wang, and Gordon Wetzstein. Generative rendering: Controllable 4d-guided video generation with 2d diffusion models. *arXiv preprint arXiv:2312.01409*, 2023.

[7] Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23206–23217, 2023.

[8] Dave Epstein, Allan Jabri, Ben Poole, Alexei Efros, and Aleksander Holynski. Diffusion self-guidance for controllable image generation. *Advances in Neural Information Processing Systems*, 36:16222–16239, 2023.

[9] Joël Foramitti. Agentpy: A package for agent-based modeling in python. *Journal of Open Source Software*, 6(62):3065, 2021.

[10] Daniel Geng and Andrew Owens. Motion guidance: Diffusion-based image editing with differentiable motion estimators. In *The Twelfth International Conference on Learning Representations*, 2024. URL `https://openreview.net/forum?id=WIAO4vbnNV`.

[11] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. In *The Twelfth International Conference on Learning Representations*, 2024. URL `https://openreview.net/forum?id=lKK50q2MtV`.

[12] MSU Graphics and Media Lab. Msu video frame interpolation benchmark dataset, 2022. URL `https://videoprocessing.ai/benchmarks/video-frame-interpolation-dataset.html`.

[13] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=_CDixzkzeyb`.

[14] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. URL `https://openreview.net/forum?id=qw8AKxfYbI`.

[15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

[16] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

[17] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.

[18] Philipp Holl, Nils Thuerey, and Vladlen Koltun. Learning to control pdes with differentiable physics. In *International Conference on Learning Representations*, 2020. URL `https://openreview.net/forum?id=HyeSin4FPB`.

[19] Itseez. Open source computer vision library. `https://github.com/itseez/opencv`, 2015.

[20] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15954–15964, 2023.

[21] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6038–6047, 2023.

[22] Haomiao Ni, Changhao Shi, Kai Li, Sharon X Huang, and Martin Renqiang Min. Conditional image-to-video generation with latent flow diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18444–18455, 2023.

[23] Craig W Reynolds. Flocks, herds and schools: A distributed behavioral model. In *Proceedings of the 14th annual conference on Computer graphics and interactive techniques*, pages 25–34, 1987.

[24] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.

[25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015.

[26] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=nJfylDvgzlq`.

[27] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.

[28] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. In *International Conference on Learning Representations*, 2020.

[29] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.

[30] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.

[31] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020.

[32] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

[33] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4): 600–612, 2004.

[34] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7623–7633, 2023.

[35] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.

# A  Background

DDMs (Denoising Diffusion Models) [15, 27, 29] represents a generative modeling approach that leverage a noise diffusion process to model a data distribution starting from random noise. These models are based on a predefined Markovian forward noising chain that progressively adds Gaussian noise to the data $\boldsymbol{x}_0$ in an iterative procedure of $T$ steps. The reverse diffusion process traverses back the Markov Chain and can be written as:

$$p_\theta(\boldsymbol{x}_{0:T}) = p(\boldsymbol{x}_T) \prod_{t=1}^{T} p_\theta(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t) \qquad p_\theta(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t) = \mathcal{N}(\boldsymbol{x}_{t-1} \mid \mu_\theta(\boldsymbol{x}_t, t), \sigma_t^2 \boldsymbol{I}) . \quad (1)$$

The training phase optimizes the parameters of the reverse process $p_\theta$ maximising an evidence lower bound (ELBO) over the target data. The work of [28] shows that is possible to construct a non-Markovian process defining a faster sampler (DDIM) that is compatible with the pretrained model. So starting from $p_\theta(\boldsymbol{x}_{0:T})$, it is possible to sample $\boldsymbol{x}_{t-1}$ using:

$$x_{t-1} = \sqrt{\alpha_{t-1}} \left( \frac{x_t - \sqrt{1-\alpha_t}\hat{\epsilon}_t}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t(\eta)^2} \cdot \hat{\epsilon}_t + \sigma_t(\eta)\varepsilon_t \qquad (2)$$

where $\sigma_t(\eta) = \eta\sqrt{\frac{1-\alpha_{t-1}}{1-\alpha_t}}\sqrt{\frac{1-\alpha_t}{\alpha_{t-1}}}$ and $\eta \in (0,1)$ is a parameter controlling the forward process, when $\eta = 0$, the sampling becomes deterministic, when $\eta = 1$, the process result in DDPM sampling. $\hat{\epsilon}_t$ is the estimated noise present in $x_t$, typically estimated with a UNet architecture [25]: $\epsilon_t(\cdot)$. Finally, $\varepsilon_t$ is an independent normal stochastic variable. In this work we employ a Latent Diffusion Model [24] that perform the diffusion process over a compressed latent space, reducing the computational burden of training in pixel space, while keeping high perceptual quality. Before the diffusion process, a VQ-VAE [32] is trained; the input image is then encoded by the VQ-VAE Encoder $\mathcal{E}$ that reduces the spatial dimension. The generated features are decoded back to the image space when generating images by means of th VQ-VAE Decoder $\mathcal{D}$. The UNet architecture is tipically composed by convolutional layers followed by spatial self-attention layers and cross-attention conditioning layers. Recent works [4, 13, 20] propose to reprogram this mechanism to enhance consistency between frames by letting the currently generated frame to attend to the first frame by swapping the original attention keys (K) and values (V) with the keys and values of the first frame, leading to the Cross-Frame Attention (CFA) mechanism:

$$\text{Cross-Frame-Attn(Q,K,V)} = \text{Softmax}\left( \frac{Q^f \cdot K^1}{\sqrt{d_k}} \right) V^1 \qquad (3)$$

where $V^1$ and $K^1$ represent the keys and values of the first frame, while $Q^f$ represents the queries of the current frame, and $d_k$ is the channel dimension of the keys. In this work we will use the notation $\epsilon_t(z, \mathcal{P}; \{a, b, c, \dots\})$, where $z$ is a latent, $\mathcal{P}$ is the prompt, and $\{a, b, c, \dots\}$ is a *list* of latents to attend to, as MCFA enables to attends to a list of latents and not only to a single one.

Classifier-Free Guidance (CFG) [14] is a widely used technique to guide conditional generation process using a linear combination of conditional and unconditonal estimated scores:

$$\hat{\epsilon} = \epsilon_t(z, \mathcal{P}_\emptyset, \{\dots\}) + \gamma\left[ \epsilon_t(z, \mathcal{P}, \{\dots\}) - \epsilon_t(z, \mathcal{P}_\emptyset, \{\dots\}) \right] \qquad (4)$$

where $\gamma$ is the scaling factor, $\mathcal{P}_\emptyset$ represents the null condition and $\mathcal{P}$ is the target text prompt.

# B  Extendend Ablation Study

In this section we show the remaining ablations for the scene *Earth*, and additional ablations on two new scenes: *Dragons* and *Satellite Scan*. The ablations for cross frame attention mechanism can be found in Figures 10 and 11. The ablations of the Spatial-$\eta$ are shown in Figures 12 and 13. Moreover, we also show the contribution of the inversion mechanism in Figures 14, 15, and 16.

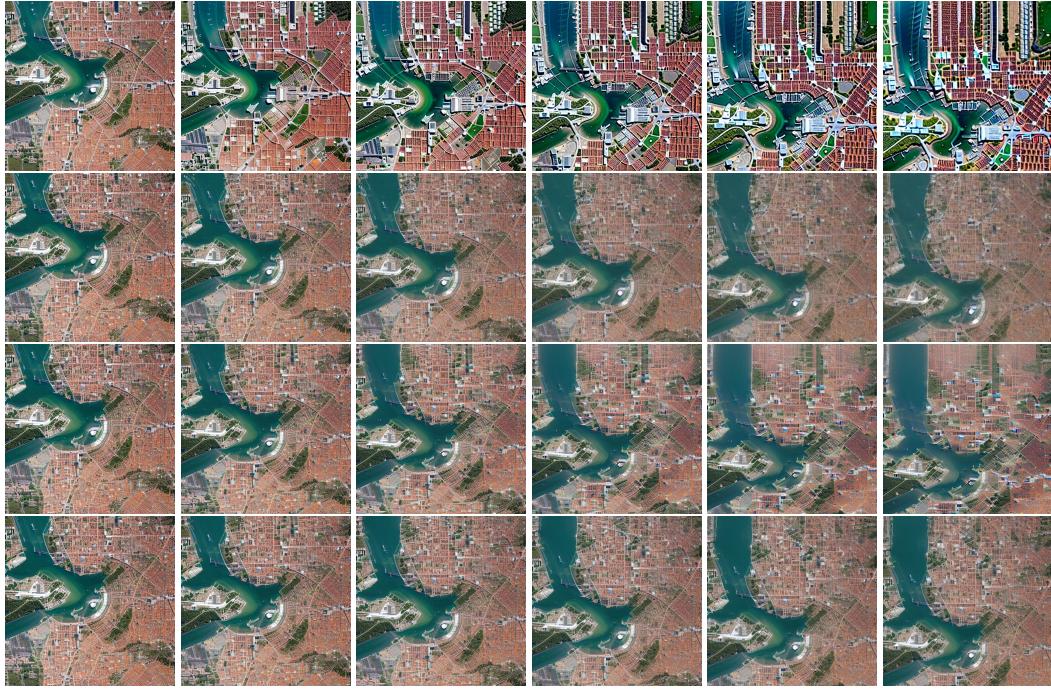## B.1 Multiple Cross-Frame Attention Mechanism Ablation



Figure 10: Ablation - Cross-Frame attention. First row: no cross frame attention; Second Row: Attend only to the initial frame; Third Row: Attend only to the previous frame; Fourth Row: Attend to the initial and preceding frame (ours).
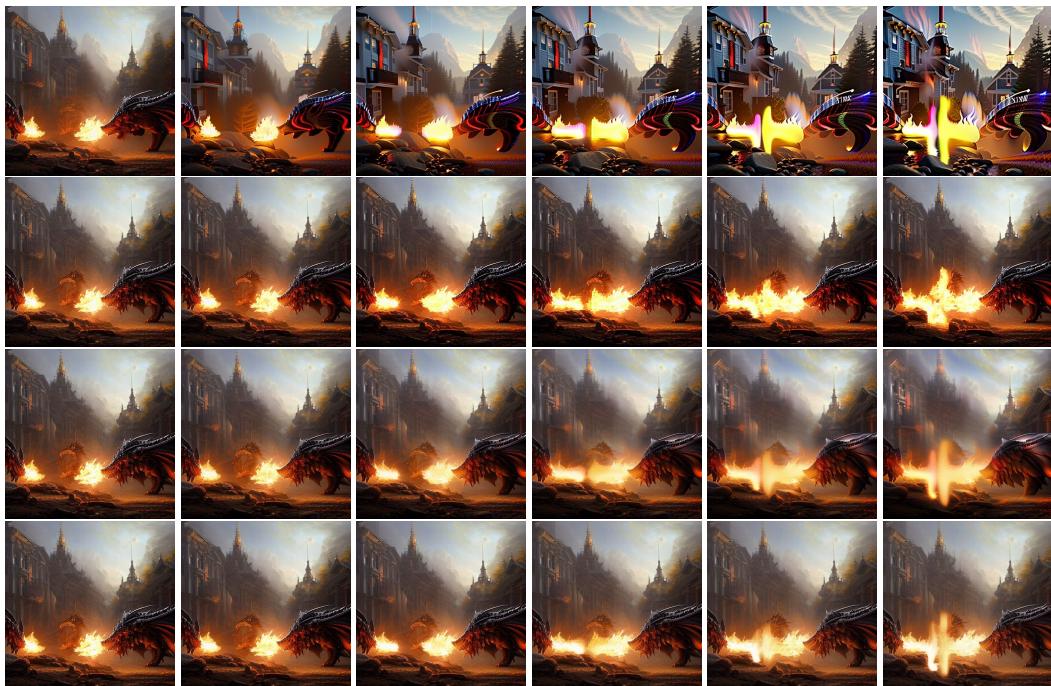


Figure 11: Ablation - Cross-Frame attention. First row: no cross frame attention; Second Row: Attend only to the initial frame; Third Row: Attend only to the previous frame; Fourth Row: Attend to the initial and preceding frame (ours).

## B.2  Spatial eta



Figure 12: Ablation - Spatial-$\eta$. First Row: Spatial-$\eta$ on; Second Row: $\eta = 0$.



Figure 13: Ablation - Spatial-$\eta$. First Row: Spatial-$\eta$ on; Second Row: $\eta = 0$.
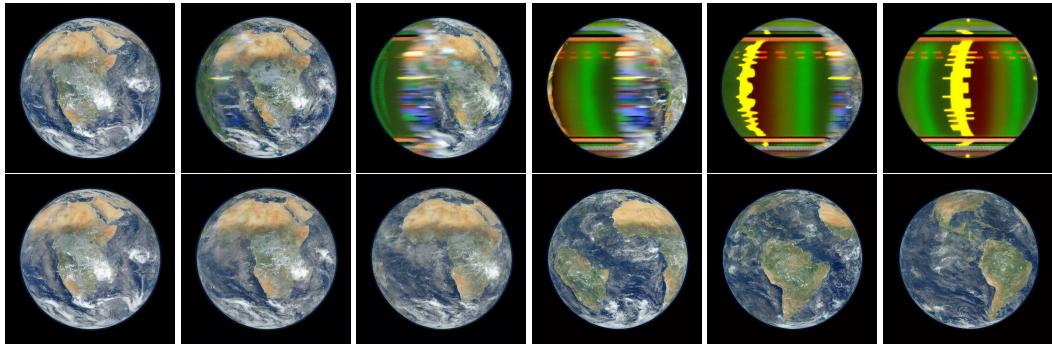
## B.3 Inversion



Figure 14: Ablation - Inversion Mechanism. First Row: Without Inversion; Second Row: With Inversion



Figure 15: Ablation - Inversion Mechanism. First Row: Without Inversion; Second Row: With Inversion



Figure 16: Ablation - Inversion Mechanism. First Row: Without Inversion; Second Row: With Inversion

# C  Obstacles, Different Physics and Additional Visual Results

In this section we showcase additional visual results of our method; all the generated videos can be found in the *Supplementary Material*. In Fig. 17 we show an example of a poured glass with MOTIONCRAFT (third row) and applying the same flow in image space (fourth row). In the first two rows of Fig. 17 we show the results of the Φ-flow physics simulator. Note that we simulate both the fluid as a set of particles (*Eulerian* simulations) in a specific position (blue balls in the first row of the figure) and two obstacles (orange objects) representing the glass and the jug. The corresponding optical flow that we used in MOTIONCRAFT is visualized in the second row of the figure.

As it can be seen, the optical flow applied to the image space produces some artefacts, such as deformations of the glass and the smoothness of the liquid due to the stretching of the pixels. On the other hand, when the same flow is applied to the noisy latent space through our method, the resulted video appears more realistic, avoiding such deformations.

Since Φ-flow is able to adopt both *Eulerian* and *Lagrangian* numerical solvers, we show the corresponding videos in Fig. 18 (second and fourth row). While the former decomposes the fluid in a set of particles, the latter models the fluid in the entire space as a fluid field. In both cases we extract from the simulation the (eventually extrapolated) velocity field (first and third row in Fig. 18) and we use it as the optical flow in the latent space, resulting in two different videos.
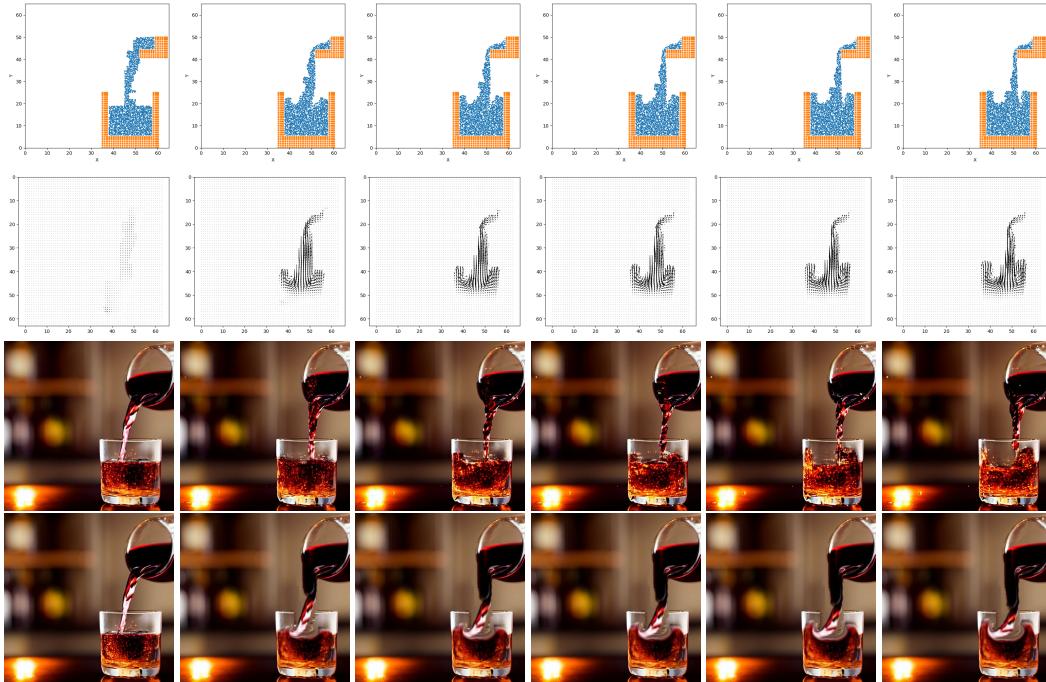


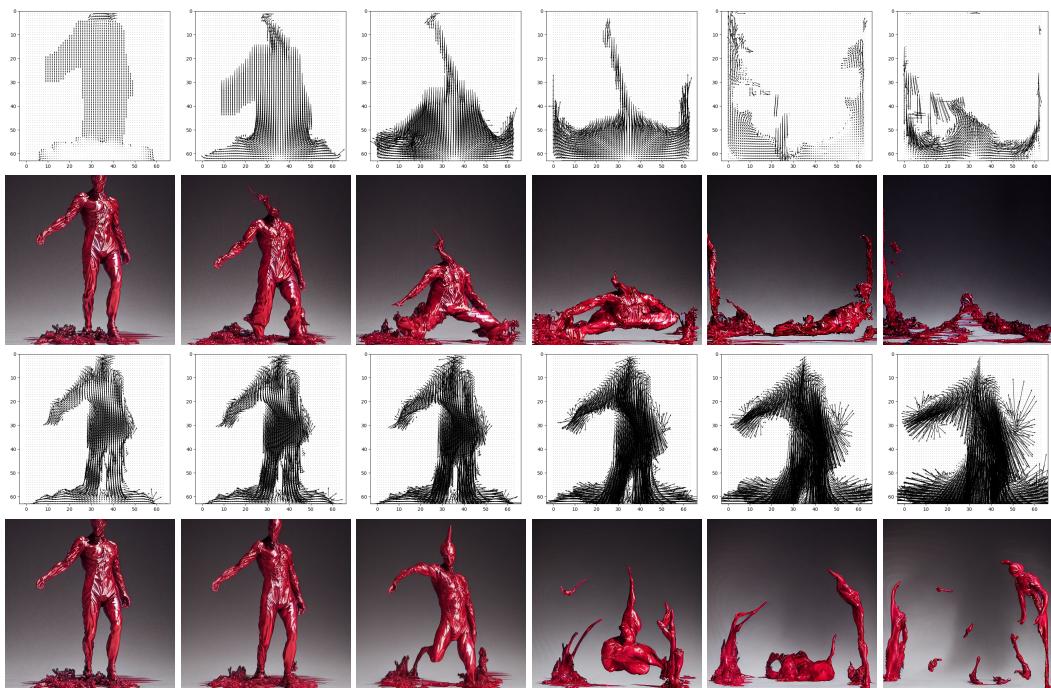Figure 17: Fluid simulation: pouring drink. Top: MOTIONCRAFT; Bottom: T2V0 [20].

Figure 18: Smoke simulation: Evaporating man. Top: MOTIONCRAFT; Bottom: T2V0.

## D   Text Prompts

In this section we state the text prompts used in the generated videos for our method and T2V0. Note that while MOTIONCRAFT is able to start from a real or generated image (with almost zero error for the real image reconstruction), T2V0 needs a hyper-parameters tuning due to a high guidance scale (not supporting direct inversion of real images).

- *Fighting Dragons*: "Two dragons fighting while breathing fires to each other. The flames are blazing and majestic light. Theatrical, character concept art by ruan jia, thomas kinkade, and trending on Artstation."
- *Melting Man* (both versions): "transparent man made by water and smoke, in style of Yoji Shinkawa and Hyung-tae Kim, trending on ArtStation, dark fantasy, great composition, concept art, highly human made of water and foam, in the style of Pierre Koenig, red pigment, pastel paint, pink color scheme"
- *Satellite Scan*: "a satellite image of a city"
- *Revolving Earth*: "a close up of a picture of the earth from space."
- *Flock of birds*: "a small flock bird flying in the sky at the sunset"
- *Pouring drink*: "wine falling on a empty glass"

For the text prompts of *Fighting Dragons* and *Melting Man* we leveraged MagicPrompt (for which we credit Gustavo Santana), a tool for rewriting simple text prompts to create more appealing starting images with Stable Diffusion.

For each example, the negative prompt $\mathcal{P}_\emptyset$ is equal to "poorly drawn, cartoon, 2d, disfigured, bad art, deformed, poorly drawn, extra limbs, close up, b&w, weird colors, blurry"

## E   Limitations and future work

In this section we discuss the limitations of the proposed approach. Being a zero-shot approach, MOTIONCRAFT relies on the pretrained text-to-image model, i.e. Stable Diffusion, and it can inherit some limitations from it, such as not exact DDIM inversion. Hence, by exploiting other diffusion models we could improve our method as well.
Experimentally, we observed a global color shift, getting stronger in the last frames of the generated videos. We noted that the proposed MCFA strategy partially solved this, but a better solution could be attending to all the previous generated frames (albeit resulting in a memory and run-time complexity increase). Moreover, MOTIONCRAFT depends on the optical flow derived from physics simulations but there are some dynamics that may be difficult to simulate (e.g. the motion of a dancer), thus limiting the generality of the generated videos. However, we speculate that it might be possible to devise a generative model of optical flows conditioned on a starting frames and a prompt, while also being constrained by a physics simulator. This could readily provide inputs to MOTIONCRAFT and have the advantage of disentangling learning of motion from learning of content. A future direction could also employ a better interaction between the image generator and the physics simulator, in order to have a closed feedback-loop framework leading to more physical fidelity in the generated frames. In this work we have shown videos generated by different physical simulations, but as future work we could also combine them to generate more complex scenes with different physics mixed together.

## F   Implementation Details and Licenses

We used the following hyperparameters throughout the work if not explicitly said otherwise. We set $\tau = 400$, the number of inference steps (both for DDIM inversion and for inverse diffusion) is set to 200 and the used model is `runwayml/stable-diffusion-v1-5` (license CreativeML Open RAIL-M). All our experiments are done on a single NVIDIA A6000 (48GB); video generation runs in minutes (1-5min) on a single GPU. Our provided code is available under MIT license. The *Earth* image is a composite of six separate orbits taken on January 23, 2012 by the Suomi National Polar-orbiting Partnership satellite (Credit: NASA/NOAA).

## G   Broader Impact

Synthetic video generation is a powerful technology that can be misused to create fake videos, hence it is important to limit and safely deploy these models. From a safety perspective, we emphasize that MOTIONCRAFT does not add any new restrictions nor does it relax any existing ones with respect to our base text-to-image model. Moreover MOTIONCRAFT, using existing text-to-image diffusion models, does not need extra training or adjustments. This means we avoid the large environmental costs associated with training new models. One possible broader impact of MOTIONCRAFT is its usage by scientists across various fields to visualize their simulations, thereby offering AI-based visualization of physical processes to a wider scientific audience.