

WildSmoke: Ready-to-Use Dynamic 3D Smoke Assets from a Single Video in the Wild

Anonymous CVPR submission

Paper ID 20

Abstract

001 We propose a pipeline to extract and reconstruct **dynamic**
002 **3D smoke assets** from a **single in-the-wild video**, and fur-
003 **ther integrate interactive simulation for smoke design and**
004 **editing**. Recent developments in 3D vision have signif-
005 **icantly improved reconstructing and rendering fluid dy-**
006 **namics, supporting realistic and temporally consistent view**
007 **synthesis**. However, current fluid reconstructions rely
008 **heavily on carefully controlled clean lab environments,**
009 **whereas real-world videos captured in the wild are largely**
010 **underexplored**. We pinpoint three key challenges of recon-
011 **structing smoke in real-world videos and design targeted**
012 **techniques, including smoke extraction with background re-**
013 **moval, initialization of smoke particles and camera poses,**
014 **and inferring multi-view videos**. Our method not only out-
015 **performs previous reconstruction and generation methods**
016 **with high-quality smoke reconstructions (+2.22 average**
017 **PSNR on wild videos), but also enables diverse and realis-**
018 **tic editing of fluid dynamics by simulating our smoke assets.**
019 **We promise to release our models, data, and 4D smoke as-**
020 **sets**. We also include more video results in the supplement.
021

022 1. Introduction

023 Fluid phenomena, from eddies around high-speed trains
024 to smoke rings from jet engines, are ubiquitous. A key
025 challenge is reconstructing unobserved quantities such as
026 velocity and density over the full spatiotemporal domain
027 (3+1D) from visual inputs (2D photos or video). This
028 task, formally inferring **3D fluid fields** from imagery, i.e.,
029 reconstructing fluid dynamics from visual observations,
030 underpins applications from high-fidelity smoke render-
031 ing in visual effects [18, 37] to diagnostics in industrial
032 flows [3, 28, 33]. Recent 3D vision advances have accel-
033 erated progress: multi-view benchmarks [9] provide ac-
034 curately calibrated flow videos, and new methods infer
035 fluid fields by jointly optimizing differentiable physics (via

physics-based constraints) and neural scene representations
(via rendering losses) from video [6–8, 10, 11, 41].

Despite notable advances, most existing approaches are
designed and heavily evaluated on **deliberately controlled**
multi-view recordings for reconstructing fluid phenomena,
which **differ substantially from real-world smoke in the**
wild. Datasets such as ScalarFlow [9], TomoFluid [44],
and FluidNexus [10] rely on fixed, precisely calibrated
cameras, controlled smoke generation in laboratory set-
tings, and pre-captured simplistic backgrounds. Such con-
ditions are impractical for in-the-wild smoke, which is of-
ten recorded with a single moving camera (e.g., handheld
or drone footage), with cluttered backgrounds and unpre-
dictable camera motion.

At the same time, demand for ready-to-use 3D as-
sets has surged across graphics and simulation applica-
tions [27, 38, 46]. **Dynamic 3D smoke assets are partic-**
ularly valuable for downstream visual effects editing and
real-time simulations [2, 5, 47], owing to their intricate mo-
tion and evolving volumetric properties. Yet, despite their
importance, the generation of dynamic 3D smoke assets re-
mains largely underexplored. We are motivated to ask:

*How to bridge the gap between existing lab-
controlled smoke reconstruction methods and in-the-
wild smoke videos?*

In this work, we aim to reconstruct ready-to-use dy-
namic 3D smoke assets from a single in-the-wild video. We
identify three key challenges of reconstructing smoke fields
from the wild (Fig. 1, Section 2.2): noisy backgrounds, un-
known camera poses, and coupled camera viewpoints and
timesteps. Correspondingly, we introduce three key tech-
niques: 1) In-context smoke segmentation, plus dehazing
for light smoke background removal; 2) Initializations of
smoke particles and pose estimation; 3) Decoupling spatial
and temporal camera trajectories by multi-view generation
and local pose perturbation. More importantly, we demon-
strate practical applications of our reconstructed smoke as-
sets through fluid simulation.

Our main contributions are summarized as follows:

- 073 1. We design a comprehensive pipeline that makes smoke
074 reconstruction practical for unconstrained in-the-wild
075 videos, extracting and reconstructing dynamic 3D smoke
076 assets with consistent improvements over prior work
077 (+2.22 dB PSNR on average on in-the-wild videos) (Sec-
078 tion 4.3).
- 079 2. We comprehensively study alternative options for each
080 component in our framework, and quantify the impact of
081 each alternative on the final reconstruction performance
082 (Section 3 and Supplement 8).
- 083 3. Our smoke assets are ready-to-use: we support realis-
084 tic editing of smoke through interactive fluid simulations
085 (Section 3.5 and 4.4).

086 2. Background

087 2.1. Smoke Reconstruction from Videos

088 Given a single-view video capturing upward-rising smoke,
089 our goal is to reconstruct a dynamic 3D representation of
090 the smoke field over time ($t = 1, \dots, T$) that is coherent,
091 view-consistent, and editable. Following FluidNexus [10],
092 we represent fluid with two types of particles.

- 093 1. 3D physical particles for positions and velocities:
094 $\mathbf{p}_t^{\text{phy}}, \mathbf{u}_t \in \mathbb{R}^{N_t^{\text{phy}} \times 3}$, where N_t^{phy} is the number of phys-
095 ical particles at t . The density field $\rho_t : \mathbb{R}^3 \mapsto \mathbb{R}$ and
096 velocity field $\mathbf{V}_t : \mathbb{R}^3 \mapsto \mathbb{R}^3$ can be further mapped
097 from particles to grids via kernel-weighted interpola-
098 tion [10, 14–16, 25].
- 099 2. N_t^{vis} visual particles (grayscale) at t with attributes:
100 $\{\mathbf{p}_t^{\text{vis}} \in \mathbb{R}^{N_t^{\text{vis}} \times 3}, \mathbf{c}_t \in \mathbb{R}^{N_t^{\text{vis}}}, \mathbf{s}_t \in \mathbb{R}^{N_t^{\text{vis}} \times 3}, \mathbf{o}_t \in$
101 $\mathbb{R}^{N_t^{\text{vis}}}, \mathbf{r}_t \in \mathbb{R}^{N_t^{\text{vis}} \times 4}\}$, representing position, color, scale,
102 opacity, and rotation, respectively.

103 2.2. Smoke Reconstruction in the Wild: Importance 104 and Challenges

105 While smoke can be plausibly simulated with strong artis-
106 tic control, simulation alone is often underconstrained to
107 match a specific real shot and its scene interactions. Re-
108 constructing smoke from in-the-wild videos recovers scene-
109 conditioned density/velocity that can initialize and steer
110 downstream simulation-based editing.

111 Our primary technical contribution is, *for the first time*,
112 to identify and summarize three key challenges in recon-
113 structing 4D smoke assets from a single in-the-wild video
114 (Fig. 1), and to propose a dedicated pipeline that explicitly
115 addresses these challenges and outperforms existing smoke
116 reconstruction methods [10, 41] (Section 3).

117 **Noisy Backgrounds and Boundaries.** High-quality
118 videos of real smoke like ScalarFlow [9], TomoFluid [44],
119 and FluidNexus [10] were collected with carefully con-
120 trolled backgrounds and further post-processed to remove
121 environmental noise. However, videos in the wild are

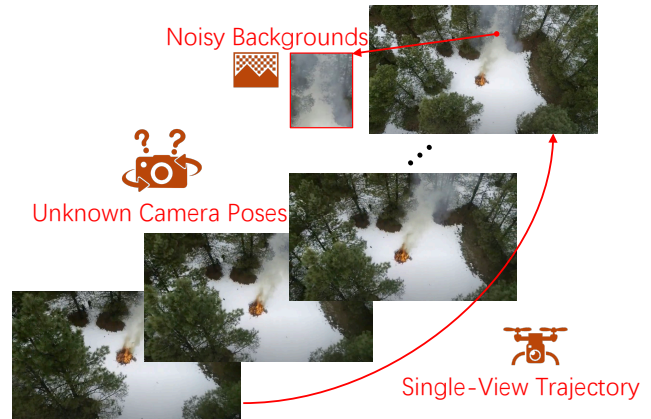


Figure 1. Challenges of smoke reconstruction from a single in-the-wild video: 1) noisy backgrounds and boundaries; 2) unknown camera poses; 3) single video with coupled camera viewpoints and timesteps.

inevitably contaminated by complex backgrounds visible
through semi-transparent smoke.

Unknown Camera Poses. Although videos of smoke in the wild are captured with moving cameras, their camera poses are typically neither measured nor released with the footage.

Single-Camera Trajectory. As a video of smoke in the wild typically provides only one camera trajectory over time, spatial viewpoints and temporal frames become inherently entangled: there is an (approximately) one-to-one correspondence between camera viewpoints and timesteps.

133 3. Methods

To address the above challenges, we propose a unified pipeline to reconstruct clean, background-free smoke from unconstrained real-world videos, as shown in Fig. 2.

137 3.1. Smoke Extraction

In real-world scenarios, smoke often exhibits complex, irregular boundaries and is intertwined with significant noise from backgrounds. To address this issue, we propose a pipeline that extracts clean smoke regions from a single in-the-wild video.

(a) Smoke Mask Extraction. We first segment the smoke and obtain a binary mask sequence $\mathbf{M} = \{M_t\}$ over all timesteps t in the video. Owing to the *widely observed limitation* of off-the-shelf semantic segmentation models, which generally struggle to generalize to in-the-wild smoke, we adopt a one-shot learning strategy, using SAM [19] for annotation and SegGPT [36] for propagation¹. For each

¹Although it is possible to fine-tune a segmentation model for general smoke extraction, the required data collection and annotation effort is beyond the scope of this work.

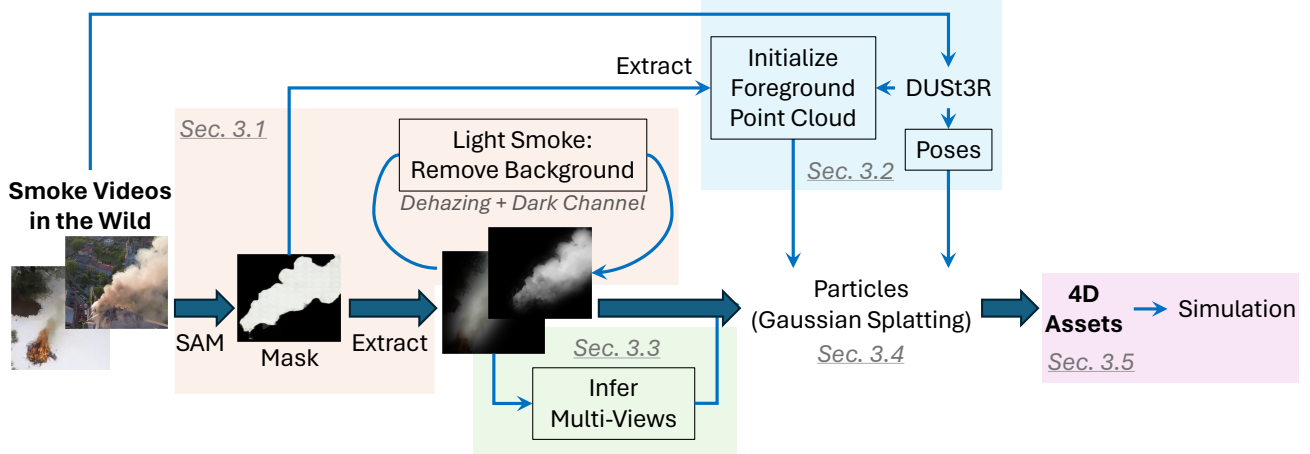


Figure 2. To reconstruct smoke from a single video in the wild, our pipeline includes five steps: 1) Smoke extraction, with background removed for light smoke due to translucence (Section 3.1); 2) Pose estimation and coarse initialization for 3D point cloud (Section 3.2); 3) Inferring multi-views for smoke (Section 3.3); 4) Training smoke particles (Section 3.4); 5) composition and simulation of 4D smoke assets (Section 3.5). “SAM”: segment anything [19]. “DUST3R” [35].

video, we select an early reference frame I_{ref} , interactively annotate its smoke with SAM² to obtain M_{ref} , and then feed the pair $(I_{\text{ref}}, M_{\text{ref}})$ to SegGPT [36] to perform one-shot inference and produce masks M_t for the remaining frames I_t . This combines SAM’s high-quality per-frame annotation with SegGPT’s frame-wise generalization.

The boundary of segmented smoke could still be noisy. We apply a Gaussian filter f_g to smooth the boundary of inferred masks, yielding a more natural smoke contour. The smoke areas are extracted via $\tilde{S}_t = f_g(M_t) \cdot I_t$. As shown in Fig. 3(a), this in-context segmentation can reliably segment smoke boundaries. See Supplement9.2 in the supplement for more smoke segmentation examples.

At this stage, it is necessary to further distinguish dense from light smoke. We categorize them based on the visibility of the background. For dense smoke, the background is completely occluded due to high optical density (low transmittance). In this case, background contamination is negligible, and the smoke can be directly treated as the masked region. However, for the light smoke, the background remains partially visible (nonzero transmittance).

(b) Extract Clean Smoke. Similar to fog or haze, background objects remain visible through light smoke due to light transmission, introducing artifacts during reconstruction. For light smoke, we fine-tune a pretrained DehazeFormer [31] to isolate the smoke and remove the background. We construct a fine-tuning dataset of synthetic smoky inputs \tilde{I} (pixels normalized to $[0, 1]$) by blending coarsely extracted smoke \tilde{S} with video frames of clean background without smoke I^{clean} , which serve as supervision targets:

$$\tilde{I} = I^{\text{clean}} \cdot \tilde{T} + A \cdot \tilde{S}, \quad (1)$$

where $\tilde{T} = 1 - \tilde{S}$ is the coarse transmission map, representing the amount of background light that reaches the camera through the smoke. We estimate the atmospheric light A using the dark channel prior [12]. As shown in Fig. 3 (c), the fine-tuned DehazeFormer can remove the foreground smoke and recover the background \tilde{I}^{clean} .

For dense smoke where background contamination is negligible, the extracted smoke \tilde{S} from masks can be used directly without this dehazing step.

For light smoke, we extract the clean foreground smoke based on Equation (1):

$$S = 1 - \frac{I - A}{\tilde{I}^{\text{clean}} - A}, \quad (2)$$

where I is the smoky frame in the original video, and \tilde{I}^{clean} is the background recovered by DehazeFormer. The resulting foreground-smoke images are then used for smoke-field reconstruction.

Alternative Approach: Mask-to-Matte Refinement. We also provide an alternative for dehazing. We refine the coarse binary masks into pixel-accurate alpha mattes using VideoMaMa [21], a mask-guided video matting model that converts masks to alpha.

We apply VideoMaMa to the coarse smoke mask sequence $\{M_t\}$ and obtain an alpha matte sequence $\{\alpha_t\}$, where $\alpha_t \in [0, 1]$ encodes smoke opacity. We then use $S = \alpha_t * I_t$ as our grayscale smoke observation for all subsequent steps (i.e., background-free smoke input to reconstruction).

3.2. Pose Estimation and Coarse Geometry

To obtain camera intrinsics and extrinsics for frames, we initialize camera poses and focal lengths with estimates

²We use an online SAM-based annotator <https://roboflow.com/>.

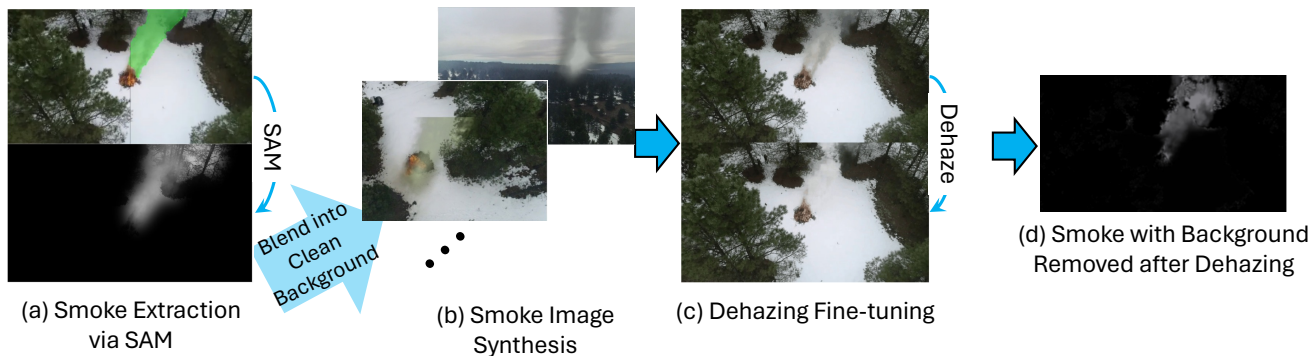


Figure 3. Smoke extraction, with background removal for light smoke.

212 from the pretrained DUST3R [35]. We further initialize both
 213 physical and visual particles from a sparse 3D point cloud
 214 $\{\mathbf{p}_i\}_{i=1}^N$ estimated by DUST3R, which provides a rough ge-
 215 ometry of the scene, including both smoke and background
 216 regions. We filter the DUST3R-generated 3D point cloud
 217 using the extracted smoke masks, retaining only the smoke-
 218 relevant points to initialize our particles. Although DUST3R
 219 produces coarse results, this initialization provides mean-
 220 ingful spatial priors for optimizing our particles; without it,
 221 training may fail to converge. See more details in Supple-
 222 ment 7.2 in the supplement.

223 **Alternative Approach: VGGT [34].** This initialization
 224 can be replaced by any pose estimation work pretrained on
 225 outdoor datasets. We also test the results on VGGT [34]
 226 (See Supplement 8 for experiments and more details).

227 3.3. Inferring Multi-View Videos

228 In wild smoke videos, cameras typically follow a single
 229 trajectory. For example, a drone may fly from the ground
 230 upwards along with the smoke plume to capture the video.
 231 Although the camera may span a wide spatiotemporal ext-
 232 ent, camera viewpoints and timesteps are highly coupled,
 233 i.e., each camera viewpoint is associated with a unique
 234 timestep, and vice versa. This coupling can cause over-
 235 fitting and degrade novel-view synthesis for unseen view-
 236 points or timesteps.

237 To decouple the spatiotemporal camera trajectory, in our
 238 work we employ generative multi-view synthesis. To obtain
 239 multi-view supervision from single-view input, we use pre-
 240 trained SV4D 2.0 [40] to generate smoke videos from novel
 241 views. We choose azimuth angles of $[-10^\circ, 10^\circ, 20^\circ, 30^\circ]$
 242 relative to the pose of the current frame as our novel view-
 243 points. Since SV4D 2.0 does not support long temporally
 244 consistent video generation, we split our video sequences
 245 into short clips, and overlap one frame between neighbor-
 246 ing clips, ensuring valid conditioning for every novel-view
 247 segment. Per-frame camera-to-world poses are obtained by
 248 applying the angle set to the DUST3R-initialized poses.

249 Moreover, to address the unreliable frames produced by
 250 SV4D at later timesteps, we apply an exponentially decay-
 251 ing weight over the frame index. This strategy progressively

252 down-weights the influence of later frames generated by
 253 SV4D 2.0 during Gaussian-particle training. See Supple-
 254 ment 7.4 and 7.5 for more details.

255 **Alternative Approach:** The multi-view generator can be
 256 replaced by other novel-view video synthesis models, e.g.,
 257 TrajectoryCrafter [42]. However, for outdoor videos with
 258 cluttered backgrounds, their quality degrades severely un-
 259 der large viewpoint changes and long sequences. In our
 260 experiments, SV4D 2.0 remains the most reliable drop-in
 261 choice after background removal. See examples in Supple-
 262 ment 8.

263 3.4. Training Gaussian Particles

264 Following the training strategy of FluidNexus [10], we sep-
 265 arately train visual and physical particle representations.
 266 Particles are optimized by minimizing both visual supervi-
 267 sion (photometric errors between input frames and rendered
 268 views using 3D Gaussian Splatting [17]) and also physi-
 269 cal regularization based on fluid simulation. To overcome
 270 the “numerical diffusion” during fluid simulation that leads
 271 to artificial smoothing [1], we introduce a frequency-aware
 272 regularization during training. Moreover, to determine the
 273 *unknown buoyancy in the wild*, we make its strength learn-
 274 able. For more training details, we refer readers to Supple-
 275 ment 7.6.

276 To further disentangle spatial and temporal viewpoints
 277 in single-camera trajectories from in-the-wild videos, we
 278 not only train with our generated multi-view trajectories
 279 (Section 3.3), but also progressively enrich the view tra-
 280 jectory by **perturbing local camera poses**. Specifically,
 281 after the particles at time t converge, we introduce per-
 282 turbed viewpoints by shifting the camera pose forward by
 283 Δt along the trajectory (modulo the sequence length T).
 284 Concretely, together with the original pose (\mathbf{R}_t, t) , we also
 285 include $(\mathbf{R}_{(t+\Delta t) \bmod T}, t)$ as an input and its correspond-
 286 ing rendering result as the target. Here, \mathbf{R} is the rotation
 287 matrix of the camera’s extrinsics; $\Delta t \ll T$, i.e., we perturb
 288 within a very short period relative to the whole temporal
 289 domain. Essentially, we effectively decouple viewpoint and
 290 timestep by associating pose $\mathbf{R}_{(t+\Delta t) \bmod T}$ with timestep t
 291 via “local pose perturbation.” Since Δt is small, perturbed

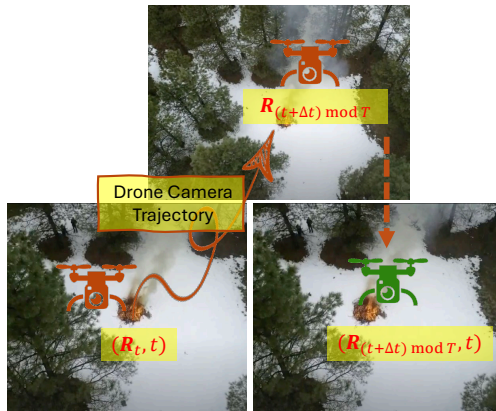


Figure 4. Local pose perturbation. Given the original camera pose (\mathbf{R}_t, t) on a single-video trajectory (“red drone”), we additionally introduce a perturbed viewpoint by shifting the camera pose by Δt along the trajectory with wrap-around $(t + \Delta t) \bmod T$ (“green drone”). For this perturbed viewpoint, the rendered image serves as supervision. Both original and perturbed samples are used during following training.

292 novel viewpoints still reside in neighborhoods of the origi-
293 nal video trajectory. In practice, we progressively increase
294 Δt from 2 to 4 during training.

295 **Alternative Approach.** We also consider HyFluid [41]
296 and FluidNexus [10], but both are less suitable for in-the-
297 wild videos in our setting. FluidNexus assumes a fixed
298 buoyancy term and lacks our spatiotemporal augmenta-
299 tion (local pose perturbation), which can lead to incorrect
300 force balance and poorer reconstruction, which motivates
301 our learnable buoyancy and augmentation design. Detailed
302 comparisons are reported in Supplement 8.

303 3.5. Ready-to-Use 4D Smoke Assets

304 A key use case of our reconstructed smoke assets is visual-
305 effects (VFX) editing of novel smoke scenes. In our
306 work, we consider simulating smoke interactions using Phi-
307 Flow [13]. The simulation is initialized with the recon-
308 structed density field ρ from visual particles, and velocity
309 field \mathbf{V} from physical particles:

310 **Density Field from Visual Particles.** For each visual par-
311 ticle i with center $\mathbf{p}_{t,i}^{\text{vis}}$, scale $\mathbf{s}_{t,i}$, rotation $\mathbf{r}_{t,i}$, and opacity
312 $\mathbf{o}_{t,i}$, we define a Gaussian kernel

$$313 \phi_{t,i}(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{p}_{t,i}^{\text{vis}})^\top \Sigma_{t,i}^{-1}(\mathbf{x} - \mathbf{p}_{t,i}^{\text{vis}})\right), \quad (3)$$

314 where $\Sigma_{t,i} = R(\mathbf{r}_{t,i}) \text{diag}(\mathbf{s}_{t,i}^2) R(\mathbf{r}_{t,i})^\top$, $\mathbf{r}_{t,i} \in \mathbb{R}^4$ is
315 the unit quaternion representing the particle rotation, and
316 $R(\mathbf{r}_{t,i}) \in \mathbb{R}^{3 \times 3}$ is the corresponding rotation matrix ob-
317 tained from $\mathbf{r}_{t,i}$. The density field is the weighted sum of
318 all particle kernels:

$$319 \rho_t(\mathbf{x}) = \sum_{i=1}^{N_t^{\text{vis}}} \mathbf{o}_{t,i} \phi_{t,i}(\mathbf{x}). \quad (4)$$

Velocity Field from Physical Particles. Given a physical
320 particle with position $\mathbf{p}_{t,i}^{\text{phy}}$ and velocity $\mathbf{u}_{t,i}$, we map them
321 to the grid. The Gaussian kernel $\phi'_{t,i}$ for the velocity field
322 has the same form as $\phi_{t,i}$ but is centered at $\mathbf{p}_{t,i}^{\text{phy}}$:
323

$$324 \mathbf{V}_t(\mathbf{x}) = \frac{\sum_{i=1}^{N_t^{\text{phy}}} \phi'_{t,i}(\mathbf{x}) \mathbf{u}_{t,i}}{\sum_{i=1}^{N_t^{\text{phy}}} \phi'_{t,i}(\mathbf{x}) + \varepsilon}, \quad (5)$$

325 here we use $\varepsilon = 10^{-8}$. To ensure spatial consistency be-
326 tween the density and velocity fields, both visual and phys-
327 ical splatting are restricted to the same bounding region de-
328 fined by the visual particles.

329 After the initialization of density and velocity fields, we
330 continue the fluid dynamics in PhiFlow [13] with the stan-
331 dard MacCormack semi-Lagrangian method [29] for advec-
332 tion and the projection-based pressure-Poisson solver for in-
333 compressible flow. We consider either external wind forces
334 or an inserted obstacle to interact with the smoke.

335 4. Experiments

336 4.1. Settings

337 **Datasets.** We evaluate our pipeline on both synthetic and
338 real-world collections of smoke, assessing the reconstruc-
339 tion quality using view rendering. Note that, in the **ab-**
340 **sence** of any existing large-scale benchmark dedicated to
341 real-world smoke videos, we undertake a *substantial data*
342 *collection effort* and construct our dataset from scratch.

- 343 • We **synthesize** two high-resolution and photorealistic
344 smoke videos by combining a 4D smoke VDB sequence
345 with a real 3D scene. These paired videos, rendered
346 with distinct camera trajectories, allow direct comparison
347 between our reconstructed results and the ground-truth.
348 Specifically, we download a synthetic 4D smoke VDB
349 sequence and a 3D scene from CGTrader³, combine them
350 and render two videos in Blender⁴ with two camera trajec-
351 tories. One trajectory is used for training and another for
352 evaluation. See Supplement 7.1 for more rendered sam-
353 ples.
- 354 • We also evaluate on two **real-world** testbeds: 1) FLAME
355 dataset [30] is a fire-imaging dataset collected by drones
356 during a prescribed burning of piled detritus in an Ari-
357 zona pine forest. The dataset includes video recordings
358 and thermal heatmaps captured by infrared cameras. 2)
359 We further collect three videos of smoke from Pixabay⁵
360 , covering diverse scenarios.

361 To unify our training and evaluation settings, we stan-
362 dardize all our videos to 270 frames; training timesteps are

³<https://www.cgtrader.com/3d-models>

⁴<https://www.blender.org>

⁵Pixabay content is released under the Pixabay Content License; see
<https://pixabay.com/service/terms/>.

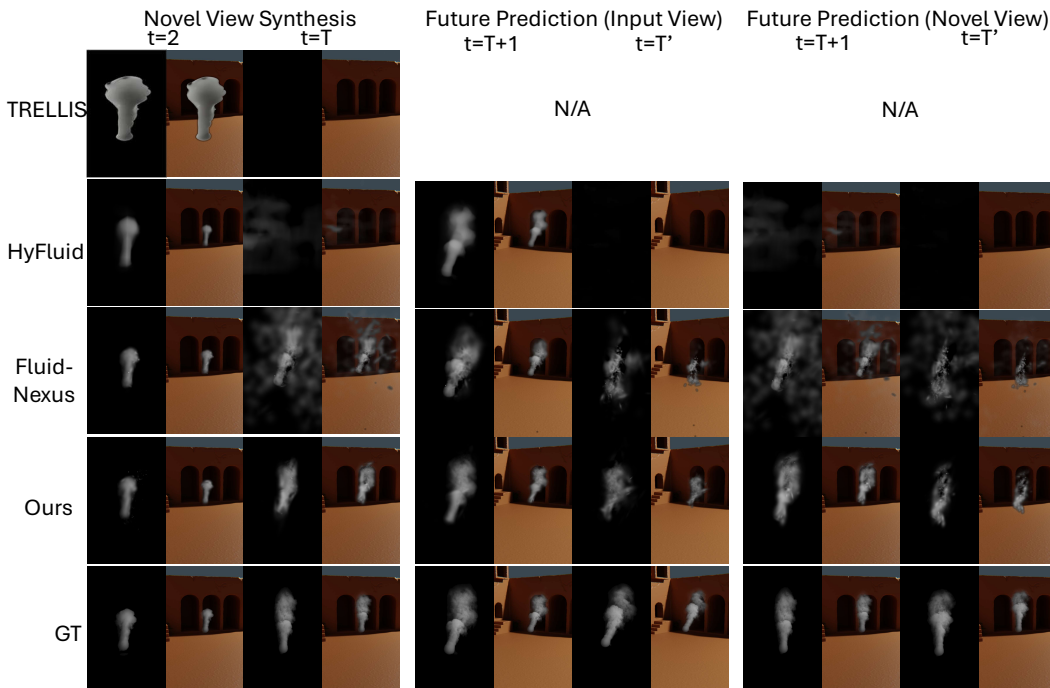


Figure 5. Visualization of novel view synthesis and future predictions on *synthetic smoke videos*. *Novel view* uses the camera pose at $t = 1$, and *input view* means the camera poses along the source video. Training uses $T = 240$ frames, and the unseen future extends to $T' = 270$ frames. *GT*: Ground Truth.

$t = 1, \dots, T$ with $T = 240$, and unseen future timesteps are $t = T + 1, \dots, T'$ with $T' = 270$.

We measure view changes along the video trajectory via the camera’s relative rotation angle $\Delta\theta_{t_1, t_2} = \arccos\left(\frac{\text{tr}(\Delta\mathbf{R}_{t_1, t_2}) - 1}{2}\right)$, where $\Delta\mathbf{R}_{t_1, t_2} = \mathbf{R}_{t_2}^\top \mathbf{R}_{t_1}$ (\mathbf{R} is the rotation matrix of the camera’s extrinsics). In our synthetic data setting, the camera undergoes a rotation of 53° from $t = 1$ to $t = T$, which is larger than the 7° rotation during the future steps from $t = T + 1$ to $t = T'$. Smoke in our videos is always centered in the scene.

Tasks. As explained by $\Delta\theta_{1, T}$ and $\Delta\theta_{T+1, T'}$ above, larger pose changes occur within the training timesteps than in the future steps. Synthesizing frames from the fixed viewpoint at the beginning of video ($t = 1$) is more challenging than following the camera pose trajectory.

Therefore, in our evaluation, the **novel view** is defined with camera pose at $t = 1$, and the **input view** follows the ground-truth camera trajectory.

Following [41], we consider the following two tasks:

- **Novel view synthesis:** We render smoke views from the fixed novel viewpoint over the training timesteps. Specifically, we fix the camera pose at $t = 1$ and synthesize novel views through $t = 2, \dots, T$. We study novel view synthesis only on our synthetic dataset, due to the lack of ground-truth multi-view videos on real-world videos (FLAME and Pixabay).
- **Future prediction:** We extrapolate the fluid dynamics into the future steps $t = T + 1, \dots, T'$. No model is

ever trained with ground-truth future frames from videos. On synthetic videos, we study both quantitative and visual results, and predict futures for both input view (i.e. follow the ground-truth camera pose trajectory during $t = T + 1, \dots, T'$) and novel view (fixed camera pose at $t = 1$). Due to the lack of ground-truth multi-view videos, on real-world videos (FLAME and Pixabay), we can only study future predictions based on the input view.

During inference, the learned velocity field is used to *ad-* *evolve* (evolve) the visual particles for future prediction. We refer the reader to [41] for more details about these tasks.

Evaluation Metrics. We report the peak signal-to-noise ratio (PSNR) averaged over frames, and defer the *structural similarity index measure* (SSIM) and the *perceptual metric LPIPS* [45] in [Supplement 9.4](#). These metrics are *widely adopted* in prior smoke reconstruction [10, 41] and deblurring works [20, 26].

While perfectly faithful recovery of turbulent smoke from a single in-the-wild video is arguably ill-posed, and novel-view PSNR from a single input trajectory cannot fully resolve the resulting ambiguities, we design our method and evaluation protocol to provide, to the best of our knowledge, the most faithful reconstructions currently attainable under these challenging conditions.

Baselines. We compare with both reconstruction and generation methods, including HyFluid [41], FluidNexus [10], and Trellis [38]. We follow FluidNexus to train with the default grayscale input setting.

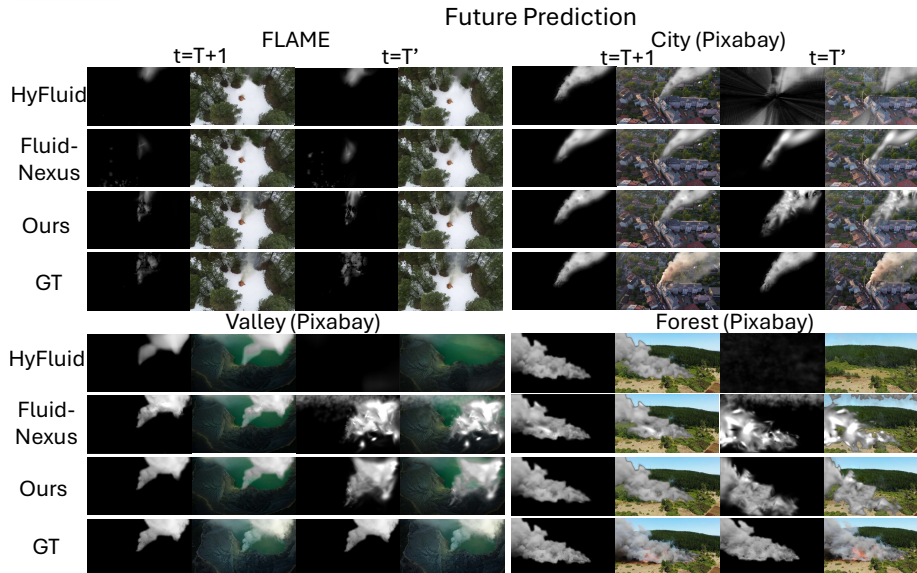


Figure 6. Visualization of future predictions (input view) on the FLAME dataset [30] and videos collected from Pixabay. Training uses $T = 240$ frames, and the unseen future extends to $T' = 270$ frames. *GT*: Ground Truth.

Table 1. Comparing PSNR (higher is better) of smoke reconstruction by different methods on the synthetic dataset. *Novel view* uses the camera pose at $t = 1$, and *input view* means the camera poses along the source video.

Methods	Novel View Synthesis	Future Prediction (Input View)	Future Prediction (Novel View)
Trellis [38]	19.98	-	-
HyFluid [41]	24.26	22.54	22.18
FluidNexus [10]	29.26	23.61	21.54
Ours	29.78	25.26	25.04

Table 2. Cumulative ablation study (PSNR) on the synthetic smoke video. *Novel view* uses the camera pose at $t = 1$, and *input view* means the camera poses along the source video.

	Novel View Synthesis	Future Prediction (Input View)	Future Prediction (Novel View)
Baseline	22.55	16.69	18.17
+ Smoke Extraction	29.26	23.61	21.54
+ DUS3R Init.	29.48	26.45	22.92
+ Local Perturbation	29.77	26.85	23.59
+ Multi-Views (Ours)	29.78	25.26	25.04

4.2. Synthetic Data

Results. We first show quantitative results in Table 1 and visualizations in Fig. 5. Trellis [38] fits a static 3D asset per frame (image-to-3D), estimating neither time-varying particles nor velocities; hence future prediction is ill-defined. Moreover, per-frame inconsistency in the novel-view synthesis setting causes the smoke to progressively fade, yielding a black frame at time T . Compared with HyFluid [41] and FluidNexus [10], our pipeline achieves higher PSNR and more stable visualizations.

Ablation Study. We further provide ablation studies in Table 2; see visualization results in Supplement9.6. Smoke extraction (segmentation), initializing poses/particles, and adding inferred multi-view supervision progressively im-

prove PSNR and visual quality. Note that, due to significant spatiotemporal deviations, the task of *future prediction at novel view* is substantially more challenging than either *novel view synthesis* or *future prediction at input view*. As a result, *future prediction at novel view* demands additional geometric cues, even coarse ones generated by SV4D, and therefore benefits most (+1.45 PSNR from 23.59 to 25.04) from incorporating “+ Multi-Views.”

Table 3. Comparing PSNR (higher is better) of smoke reconstruction by different methods on the FLAME dataset [30] and videos collected from Pixabay.

Dataset	FLAME	Pixabay		
		City	Valley	Forest
HyFluid [41]	21.67	23.24	12.43	13.70
FluidNexus [10]	21.78	24.18	16.44	14.63
Ours	22.88	24.68	20.42	17.91

4.3. Smoke Videos in the Wild

We further evaluate on in-the-wild smoke videos. Due to the lack of ground-truth novel views in real videos, here we evaluate only the task of *future prediction at input view*, meaning the camera poses along the source video.

FLAME [30] and Pixabay. The FLAME dataset [30] includes drone-captured videos in the wild forest. The smoke is light, and thus background removal is necessary. We also evaluate three thick-smoke videos from Pixabay. We show the results in Table 3 and Fig. 6.

Overall, across both light and dense smoke over diverse real-world videos, our method outperforms prior work, achieving an average PSNR improvement of +2.22 dB. Qualitatively, the reconstructions remain consistent and blend back into the original backgrounds. We also test

456 the physical plausibility via velocity divergence in Supple-
457 ment 9.1, and human validation in Supplement 9.8.

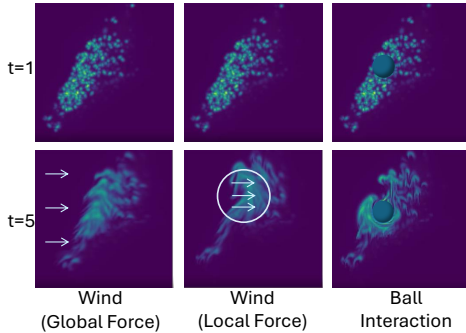


Figure 7. Visualizations of simulations with our smoke assets in different interaction scenarios. t means simulation time, different from the video frames we used in previous figures. The simulation starts with particles reconstructed at the last frame of the “city” video from Pixabay (Fig. 6 top).

458 4.4. Interactive Simulations

459 Finally, we demonstrate interactive simulations enabled by
460 our fluid reconstruction. We consider two scenarios: (i) ex-
461 ternal wind (both global and local forces), and (ii) a rigid
462 spherical obstacle. For the global wind, we apply a global
463 constant force $f_w = (0.005, 0, 0)$ uniformly across the do-
464 main. The local wind uses the same magnitude and direc-
465 tion but is confined to a sphere of radius 30 at the scene
466 center. For the obstacle case, we place a rigid ball (radius
467 10) centered at (50, 70) and simulate the resulting interac-
468 tions. All simulations use the reconstruction from the “city”
469 video from Pixabay (Fig. 6 top); results are shown in Fig. 7
470 (XY-plane from Z-axis). We can see that our smoke assets
471 support realistic editing via diverse simulation scenarios.

472 4.5. Training and Inference Cost

473 Table 4 in the Supplement 7.7 reports the per-stage GPU
474 hours. Across datasets and resolutions, our pipeline recon-
475 structs smoke from in-the-wild videos in 4.47 GPU-hours
476 on average among all 5 videos. Notably, with the DUS3R-
477 based pose initialization, smoke can be localized without
478 seeding a large number of particles, reducing the Gaussian-
479 particle training stage by 1.5 GPU-hours.

480 5. Related Works

481 5.1. Fluid Field Extraction from Videos

482 A series of recent works tackled fluid reconstruction from
483 videos, but most relied on controlled, multi-view record-
484 ings that differed substantially from in-the-wild footage.
485 PINF [6] coupled Navier–Stokes PDEs with a contin-
486 uous spatiotemporal NeRF to recover flow; NeuroFluid [11]
487 introduced a particle-driven neural renderer that embed-
488 ded fluid properties and a particle transition model;

HyFluid [41] estimated hybrid neural fields for density and
velocity using physics-based losses; FluidNexus [10] recon-
structed smoke from a single video by leveraging generative
priors, yet its multi-view generator was fine-tuned on
laboratory smoke with fixed, calibrated cameras and clean
backgrounds. To our knowledge, reconstructing fluid fields
from a single in-the-wild video and turning them into ready-
to-use dynamic 3D smoke assets remains largely underex-
plored.

5.2. Physics-based Fluid Simulation and Editing

Classical fluid simulation in computer graphics has been
broadly categorized into Eulerian [32] and Lagrangian for-
mulations [4]. Production solvers edited flows by manipu-
lating external forces, boundary conditions, and solid–fluid
interactions (obstacles), which underpinned VFX and in-
teractive applications. Recent differentiable-physics sys-
tems such as DiffTaichi [15] and PhiFlow [13] made long-
horizon, gradient-based editing practical. These tools pro-
vided useful platforms for controlled physical interactions.
Building on these advances, our work reconstructs smoke
directly from in-the-wild videos and integrates simulation
for realistic editing, extending such physics-based editing
to unconstrained real-world scenarios.

5.3. 4D Generation

Recent advances in generative modeling have markedly
improved image-to-3D and image-to-novel-view-synthesis,
and video diffusion models with stabilized virtual cameras
produced short clips with a degree of multi-view consis-
tency [22, 43, 48]. SV4D [39] and SV4D 2.0 [40] extended
view-consistent generation to dynamic scenes, reinforcing
temporal continuity. In our pipeline, we build on this gen-
erative consistency to obtain auxiliary multi-view trajec-
tories, but go further by reconstructing physically plausible
dynamic smoke fields.

6. Conclusion

In this paper, we addressed the challenge of recovering
time-varying 3D smoke fields from a single in-the-wild
video and producing ready-to-use dynamic assets. The pro-
posed pipeline tackles three core challenges (noisy back-
grounds, unknown camera poses, and coupled spatiotem-
poral views) through one-shot segmentation and dehaz-
ing, pose estimation, and decoupled spatiotemporal trajec-
tory through multi-view generation and local pose perturba-
tion. We further showed that the reconstructed assets sup-
port physically consistent editing through fluid simulation,
aligning visual fidelity with physical plausibility and sup-
porting downstream simulation workflows.

536

References

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

- [1] Numerical diffusion. https://mooseframework.inl.gov/modules/porous_flow/numerical_diffusion.html. PorousFlow module documentation, Idaho National Laboratory, accessed 2025. 4
- [2] Kai Bai, Wei Li, Mathieu Desbrun, and Xiaopei Liu. Dynamic upsampling of smoke through dictionary-based learning. *ACM TOG*, 40(1):1–19, 2020. 1
- [3] Samuel J Baker, Michael A Hobley, Isabel Scherl, Xiaohang Fang, Felix CP Leach, and Martin H Davy. Enginebench: flow reconstruction in the transparent combustion chamber iii optical engine. *arXiv preprint arXiv:2406.03325*, 2024. 1
- [4] Andrew Bennett. *Lagrangian fluid dynamics*. Cambridge University Press, 2006. 8
- [5] Mengyu Chu, Nils Thuerey, Hans-Peter Seidel, Christian Theobalt, and Rhaleb Zayer. Learning meaningful controls for fluids. *ACM TOG*, 40(4):1–13, 2021. 1
- [6] Mengyu Chu, Lingjie Liu, Quan Zheng, Erik Franz, Hans-Peter Seidel, Christian Theobalt, and Rhaleb Zayer. Physics informed neural fields for smoke reconstruction with sparse data. *ACM Transactions on Graphics (ToG)*, 41(4):1–14, 2022. 1, 8
- [7] Yitong Deng, Hong-Xing Yu, Jiajun Wu, and Bo Zhu. Learning vortex dynamics for fluid inference and prediction. *arXiv preprint arXiv:2301.11494*, 2023.
- [8] Yitong Deng, Hong-Xing Yu, Diyang Zhang, Jiajun Wu, and Bo Zhu. Fluid simulation on neural flow maps. *ACM TOG*, 42(6):1–21, 2023. 1
- [9] Marie-Lena Eckert, Kiwon Um, and Nils Thuerey. Scalarflow: a large-scale volumetric data set of real-world scalar transport flows for computer animation and machine learning. *ACM TOG*, 38(6):1–16, 2019. 1, 2
- [10] Yue Gao, Hong-Xing Yu, Bo Zhu, and Jiajun Wu. Fluidnexus: 3d fluid reconstruction and prediction from a single video. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26091–26101, 2025. 1, 2, 4, 5, 6, 7, 8, 3
- [11] Shanyan Guan, Huayu Deng, Yunbo Wang, and Xiaokang Yang. Neurofluid: Fluid dynamics grounding with particle-driven neural radiance fields. In *International Conference on Machine Learning*, pages 7919–7929. PMLR, 2022. 1, 8
- [12] Kaiming He, Jian Sun, and Xiaoou Tang. Single image haze removal using dark channel prior. *IEEE TPAMI*, 33(12):2341–2353, 2010. 3
- [13] Philipp Holl and Nils Thuerey. Φ_{flow} (PhiFlow): Differentiable simulations for pytorch, tensorflow and jax. In *International Conference on Machine Learning*. PMLR, 2024. 5, 8
- [14] Yuanming Hu, Tzu-Mao Li, Luke Anderson, Jonathan Ragan-Kelley, and Frédo Durand. Taichi: a language for high-performance computation on spatially sparse data structures. *ACM TOG*, 38(6):201, 2019. 2
- [15] Yuanming Hu, Luke Anderson, Tzu-Mao Li, Qi Sun, Nathan Carr, Jonathan Ragan-Kelley, and Frédo Durand. DiffTaichi: Differentiable programming for physical simulation. *ICLR*, 2020. 8
- [16] Yuanming Hu, Jiafeng Liu, Xuanda Yang, Mingkuan Xu, Ye Kuang, Weiwei Xu, Qiang Dai, William T. Freeman, and Frédo Durand. Quantaichi: A compiler for quantized simulations. *ACM TOG*, 40(4), 2021. 2
- [17] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM TOG*, 42(4):139–1, 2023. 4
- [18] Byungsoo Kim, Vinicius C Azevedo, Markus Gross, and Barbara Solenthaler. Lagrangian neural style transfer for fluids. *ACM TOG*, 39(4):52–1, 2020. 1
- [19] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*, pages 4015–4026, 2023. 2, 3
- [20] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiří Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8183–8192, 2018. 6
- [21] Sangbeom Lim, Seoung Wug Oh, Jiahui Huang, Heeji Yoon, Seungryong Kim, and Joon-Young Lee. Videomama: Mask-guided video matting via generative prior. *arXiv preprint arXiv:2601.14255*, 2026. 3
- [22] 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 *ICCV*, pages 9298–9309, 2023. 8
- [23] Miles Macklin and Matthias Müller. Position based fluids. *ACM TOG*, 32(4):1–12, 2013. 2
- [24] Miles Macklin, Matthias Müller, Nuttapong Chentanez, and Tae-Yong Kim. Unified particle physics for real-time applications. *ACM TOG*, 33(4):1–12, 2014. 2
- [25] Matthias Müller, David Charypar, and Markus Gross. Particle-based fluid simulation for interactive applications. In *Proceedings of the 2003 ACM SIGGRAPH/Eurographics symposium on Computer animation*, pages 154–159, 2003. 2
- [26] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3883–3891, 2017. 6
- [27] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 1
- [28] Pankaj Saini, Christoph M Arndt, and Adam M Steinberg. Development and evaluation of gappy-pod as a data reconstruction technique for noisy piv measurements in gas turbine combustors. *Experiments in Fluids*, 57(7):122, 2016. 1
- [29] Andrew Selle, Ronald Fedkiw, Byungmoon Kim, Yingjie Liu, and Jarek Rossignac. An unconditionally stable mac-cormack method. *Journal of Scientific Computing*, 35(2):350–371, 2008. 5
- [30] Alireza Shamsoshoara, Fatemeh Afghah, Abolfazl Razi, Liming Zheng, Peter Z Fulé, and Erik Blasch. Aerial imagery pile burn detection using deep learning: The flame dataset. *Computer Networks*, 193:108001, 2021. 5, 7, 4

- 649 [31] Yuda Song, Zhuqing He, Hui Qian, and Xin Du. Vision
650 transformers for single image dehazing. *IEEE TIP*, 32:1927–
651 1941, 2023. 3, 5
- 652 [32] M’hamed Souli and David J Benson. *Arbitrary Lagrangian
653 Eulerian and fluid-structure interaction: numerical simula-
654 tion*. John Wiley & Sons, 2013. 8
- 655 [33] Nils Thuerey, Konstantin Weissenow, Lukas Prantl, and Xi-
656 angyu Hu. Deep learning methods for reynolds-averaged
657 navier–stokes simulations of airfoil flows. *AIAA Journal*, 58
658 (1):25–36, 2020. 1
- 659 [34] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea
660 Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Vi-
661 sual geometry grounded transformer. In *Proceedings of the
662 Computer Vision and Pattern Recognition Conference*, pages
663 5294–5306, 2025. 4, 3
- 664 [35] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris
665 Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vi-
666 sion made easy. In *CVPR*, pages 20697–20709, 2024. 3,
667 4
- 668 [36] Xinlong Wang, Wen Wang, Yue Cao, Chunhua Shen, and
669 Tiejun Huang. Images speak in images: A generalist painter
670 for in-context visual learning. In *CVPR*, pages 6830–6839,
671 2023. 2, 3
- 672 [37] Xiaokun Wang, Yanrui Xu, Sinuo Liu, Bo Ren, Jiri Kosinka,
673 Alexandru C Telea, Jiamin Wang, Chongming Song, Jian
674 Chang, Chenfeng Li, et al. Physics-based fluid simulation in
675 computer graphics: Survey, research trends, and challenges.
676 *Computational Visual Media*, pages 1–56, 2024. 1
- 677 [38] Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng
678 Wang, Bowen Zhang, Dong Chen, Xin Tong, and Jiaolong
679 Yang. Structured 3d latents for scalable and versatile 3d gen-
680 eration. *arXiv preprint arXiv:2412.01506*, 2024. 1, 6, 7
- 681 [39] Yiming Xie, Chun-Han Yao, Vikram Voleti, Huaizu Jiang,
682 and Varun Jampani. Sv4d: Dynamic 3d content generation
683 with multi-frame and multi-view consistency. *arXiv preprint
684 arXiv:2407.17470*, 2024. 8
- 685 [40] Chun-Han Yao, Yiming Xie, Vikram Voleti, Huaizu Jiang,
686 and Varun Jampani. Sv4d 2.0: Enhancing spatio-temporal
687 consistency in multi-view video diffusion for high-quality 4d
688 generation. *arXiv preprint arXiv:2503.16396*, 2025. 4, 8
- 689 [41] Hong-Xing Yu, Yang Zheng, Yuan Gao, Yitong Deng, Bo
690 Zhu, and Jiajun Wu. Inferring hybrid neural fluid fields from
691 videos. *NeurIPS*, 36, 2024. 1, 2, 5, 6, 7, 8, 3, 4
- 692 [42] Mark Yu, Wenbo Hu, Jinbo Xing, and Ying Shan. Tra-
693 jectorycrafter: Redirecting camera trajectory for monocular
694 videos via diffusion models. In *ICCV*, pages 100–111, 2025.
695 4, 3
- 696 [43] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li,
697 Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan,
698 and Yonghong Tian. Viewcrafter: Taming video diffusion
699 models for high-fidelity novel view synthesis. *arXiv preprint
700 arXiv:2409.02048*, 2024. 8
- 701 [44] Guangming Zang, Ramzi Idoughi, Congli Wang, Anthony
702 Bennett, Jianguo Du, Scott Skeen, William L Roberts, Peter
703 Wonka, and Wolfgang Heidrich. Tomoffluid: Reconstruct-
704 ing dynamic fluid from sparse view videos. In *CVPR*, pages
705 1870–1879, 2020. 1, 2
- [45] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shecht-
man, and Oliver Wang. The unreasonable effectiveness of
deep features as a perceptual metric. In *Proceedings of the
IEEE conference on computer vision and pattern recogni-
tion*, pages 586–595, 2018. 6, 4
- [46] Zibo Zhao, Zeqiang Lai, Qingxiang Lin, Yunfei Zhao,
Haolin Liu, Shuhui Yang, Yifei Feng, Mingxin Yang, Sheng
Zhang, Xianghui Yang, et al. Hunyuan3d 2.0: Scaling diffu-
sion models for high resolution textured 3d assets generation.
arXiv preprint arXiv:2501.12202, 2025. 1
- [47] Zhiwei Zhao, Alan Zhao, Minchen Li, and Yixin Hu.
Vid2fluid: 3d dynamic fluid assets from single-view
videos with generative gaussian splatting. *arXiv preprint
arXiv:2503.00868*, 2025. 1
- [48] Jensen Zhou, Hang Gao, Vikram Voleti, Aaryaman Vasishta,
Chun-Han Yao, Mark Boss, Philip Torr, Christian Rup-
precht, and Varun Jampani. Stable virtual camera: Gener-
ative view synthesis with diffusion models. *arXiv preprint
arXiv:2503.14489*, 2025. 8