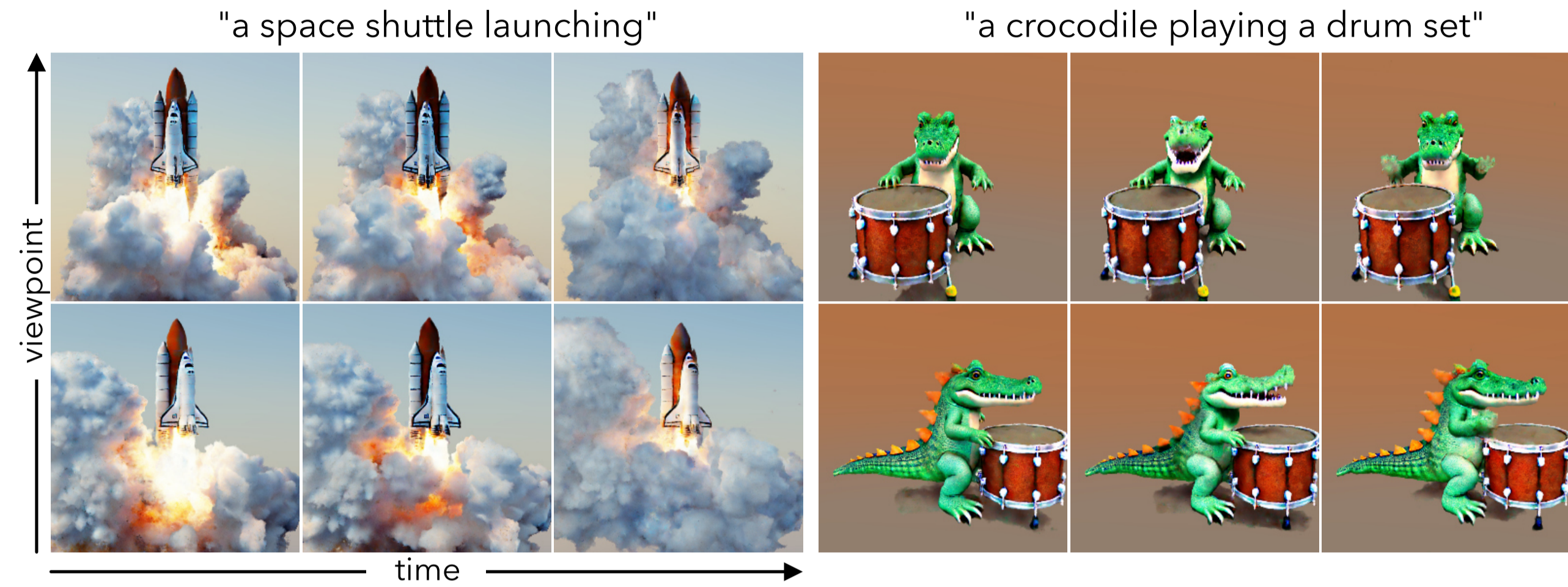
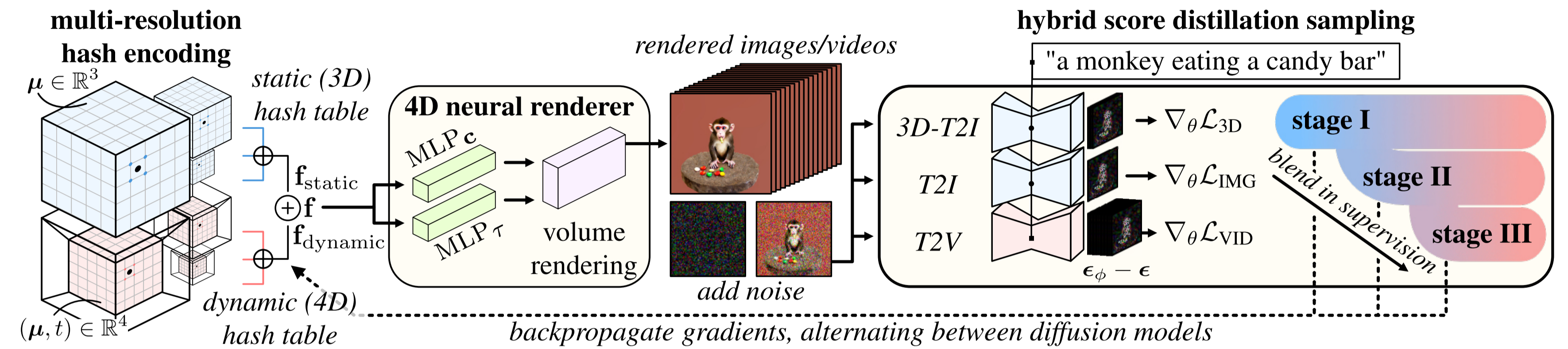


Motivation



Method

4D-fy uses a three-stage training process that smoothly blends supervision from different pre-trained diffusion models using alternating optimization.

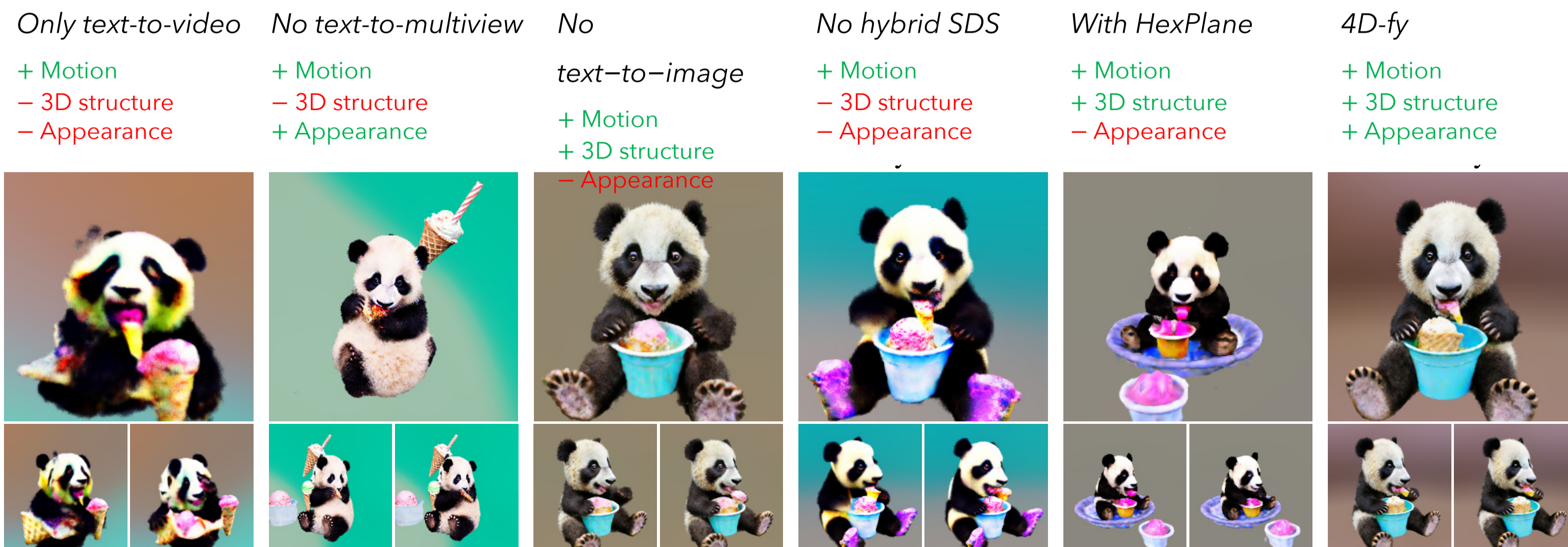


We generate text-to-4D scenes using alternating optimization (hybrid SDS) between pre-trained (i) text-to-video, (ii) text-to-multiview, and (iii) text-to-image models.

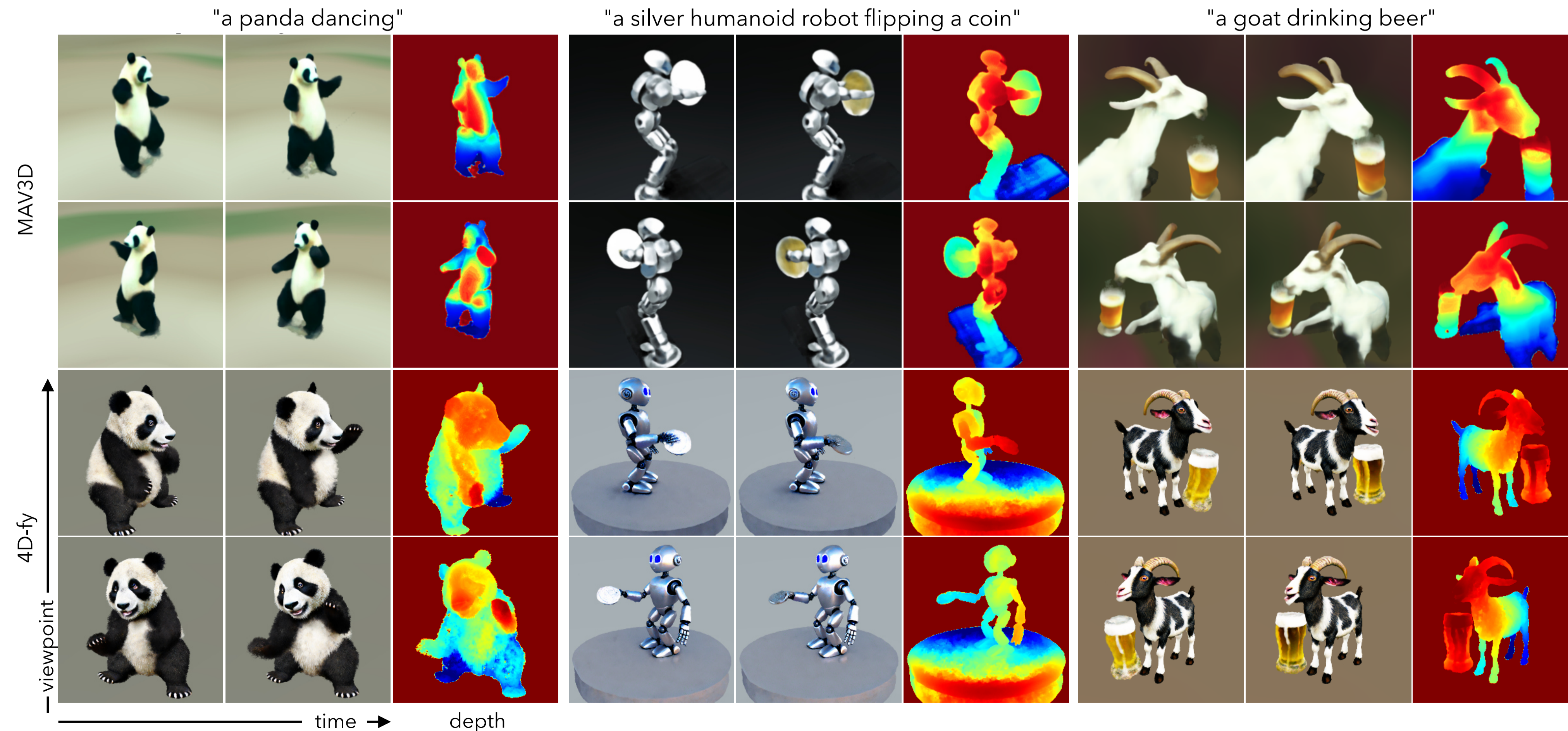
Challenges:

- **Motion.** The scene is animated using a text-to-video model (ZeroScope).
- **3D structure.** The text-to-multiview model addresses the multi-face problem (MVDream).
- **Appearance.** Single-view text-to-image models offer the best visual quality (StableDiffusion).

Ablation Study



Text-to-4D Generation



"a baby panda eating ice cream"