Vid Gen

Authors

Affiliations

August 16, 2024

#### Motivation 1

- A video with f frames is denoted as  $\mathbf{x} = [x^1, \cdots, x^f] \in \mathbb{R}^{f \times c \times h \times w}$ , where  $x^i \in \mathbb{R}^{c \times h \times w}$  denotes a frame (an image),
- we would like to build a diffusion model for the video distribution  $p(\mathbf{x})$ .

### 1.1 **Autoregressive**

We can autoregressively decompose the video distribution as:

$$p(\mathbf{x}) = \prod_{i=1}^{f} p(\mathbf{x}^i \mid \mathbf{x}^{< i}) \tag{1}$$

where  $x^{< i} = [x^1, \cdots, x^{i-1}]$  denotes the first i-1 frames. In this case, we would like to cast a pretrained T2I model to represent the conditional distribution. This can be simply achieved by modeling the autoregressively conditional score:

$$\nabla_{\boldsymbol{x}^i} \log p(\boldsymbol{x}^i \mid \boldsymbol{x}^{< i}, \boldsymbol{c}), \tag{2}$$

where c denotes the text prompt.

#### 1.2 Distributed Noisy Frame Conditioning (Deprecated) 13

- The score function  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$  is our modeling target, where  $\mathbf{x}_t$  denotes the noisy video at timestep t. The dimension 14
- of the score function is f times larger than an image diffusion model and is unaffordable.
- The key observation is that the score function of the noisy video  $\mathbf{x}_t=(m{x}_t^1,\cdots,m{x}_t^f)$  can be rewrote as:

$$\nabla_{\mathbf{x}_{t}} \log p_{t}(\mathbf{x}_{t}) = \begin{bmatrix} \nabla_{\mathbf{x}_{t}^{1}} \log p_{t}(\mathbf{x}_{t}) \\ \vdots \\ \nabla_{\mathbf{x}_{t}^{f}} \log p_{t}(\mathbf{x}_{t}) \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}_{t}^{1}} \left( \log p_{t}(\mathbf{x}_{t}^{1} \mid \mathbf{x}_{t}^{-1}) + \log p_{t}(\mathbf{x}_{t}^{-1}) \right) \\ \vdots \\ \nabla_{\mathbf{x}_{t}^{f}} \left( \log p_{t}(\mathbf{x}_{t}^{f} \mid \mathbf{x}_{t}^{-f}) + \log p_{t}(\mathbf{x}_{t}^{-f}) \right) \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}_{t}^{1}} \log p_{t}(\mathbf{x}_{t}^{1} \mid \mathbf{x}_{t}^{-1}) \\ \vdots \\ \nabla_{\mathbf{x}_{t}^{f}} \log p_{t}(\mathbf{x}_{t}^{f} \mid \mathbf{x}_{t}^{-f}) \end{bmatrix}.$$
(3

- Thus we can **distributedly** model the video score by modeling the conditional score  $\nabla_{m{x}_t^i} \log p(m{x}_t^i \mid m{x}_t^{-i})$ , and aggregate them into  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$  during sampling.
- And we would like to build  $\nabla_{{m{x}}_t^i} \log p({m{x}}_t^i \mid {m{x}}_t^{-i})$  with a pretrained image diffusion model.

## 2 Method

# 2.1 Distributed Noisy Frame Conditioning (Deprecated)

Most pretrained image diffusion models, e.g. Stable Diffusion, PixArt, are all trained as a DDPM. Consider a DDPM  $\epsilon_{\theta}(\boldsymbol{x}_t,t)$  that tries to predict the noise  $\epsilon$  added to its clean version  $\boldsymbol{x}_0$ , its link to the score function is:

$$\nabla_{\boldsymbol{x}_t} \log p_t(\boldsymbol{x}_t) = -\frac{\epsilon_{\theta}(\boldsymbol{x}_t, t)}{\sigma_{\star}^2},\tag{4}$$

24 and what we want is:

$$\nabla_{\boldsymbol{x}_{t}^{i}} \log p(\boldsymbol{x}_{t}^{i} \mid \boldsymbol{x}_{t}^{-i}) = -\frac{\epsilon_{\theta}(\boldsymbol{x}_{t}^{i}, t \mid \boldsymbol{x}_{t}^{-i})}{\sigma_{t}^{2}} = -\frac{\epsilon_{\theta}(\boldsymbol{x}_{t}^{i}, t \mid \boldsymbol{x}_{t}, \boldsymbol{x}_{t}^{-i}, i)}{\sigma_{t}^{2}} = -\frac{\epsilon_{\theta}(\boldsymbol{x}_{t}^{i}, t \mid \boldsymbol{x}_{t}, i)}{\sigma_{t}^{2}},$$
(5)

- so what we need to do is build a model to summarize  $(\mathbf{x}_t, i)$  and feed it into a pretrained image diffusion model  $\epsilon_{\theta}(\mathbf{x}_t, t)$  as additional conditions in addition to the text prompts.
- Note that  $\mathbf{x}_t$  is a noisy video, which can be very memory-costing and redundancy to be directly feeded into the image diffusion model, we want to first build a model that can summarize  $\mathbf{x}_t \in \mathbb{R}^{f \times c \times h \times w}$  to a fixed number of context tokens  $\mathbf{c} \in \mathbb{R}^{\ell \times d}$ , where  $\ell$  denotes the number of tokens and d denotes the dimension of each token.

### 30 References

- Z. Deng, J. Shi, and J. Zhu. Neuralef: Deconstructing kernels by deep neural networks. In <u>International Conference on Machine Learning</u>, pages 4976–4992. PMLR, 2022.
- 3 [Deng et al., 2022]