

702 6. Supplementary Materials

703 6.1. The proof of Proposition 1

704 We observe the forward spatio-temporal noising process as
705 follows:

$$706 z_t^i = \sqrt{\alpha_t}(\sqrt{\rho_t} * z_{t-1}^i + \sqrt{1-\rho_t} * z_{t-1}^{i-1}) + \sqrt{1-\alpha_t}\varepsilon_t^i. \quad (24)$$

707 Then, with condition $\sqrt{\rho_t} = c, \sqrt{1-\rho_t} = d$, we can
708 calculate the form $z_1^i, z_2^i, \dots, z_t^i$:

$$709 z_1^i = \sqrt{\alpha_1}(c * z_0^i + d * z_0^{i-1}) + \sqrt{1-\alpha_1}\varepsilon_1^i, \quad (25)$$

$$710 z_2^i = \sqrt{\alpha_1\alpha_2}(c^2 * z_0^i + 2cd * z_0^{i-1} + d^2 * z_0^{i-2}) \\ 711 + \sqrt{\alpha_2}\sqrt{1-\alpha_1}(c\varepsilon_1^i + d\varepsilon_1^{i-1}) + \sqrt{1-\alpha_2}\varepsilon_2^i \\ 712 = \sqrt{\alpha_1\alpha_2} \sum_{m=0}^2 \binom{m}{2} c^{t-m} d^m z_0^{i-m} + \sqrt{1-\alpha_1\alpha_2}\hat{\varepsilon}_2^i \quad (26)$$

713 ...

$$714 z_t^i = \sqrt{\bar{\alpha}_t} \sum_{m=0}^t \binom{m}{t} c^{t-m} d^m z_0^{i-m} + \sqrt{1-\bar{\alpha}_t}\hat{\varepsilon}_t^i, \quad (27)$$

715 where, $\hat{\varepsilon}_2^i \sim \mathcal{N}(0, 1)$. Therefore, Then, we can generalize
716 the mathematical form of z_t^i :

$$717 z_t^i = \sqrt{\bar{\alpha}_t} \sum_{m=0}^t \binom{m}{t} c^{t-m} d^m z_0^{i-m} + \sqrt{1-\bar{\alpha}_t}\hat{\varepsilon}_t^i$$

$$718 = \sqrt{\bar{\alpha}_t}\rho^{t/2}z_0^i + \sqrt{\bar{\alpha}_t}g(t, i) + \sqrt{1-\bar{\alpha}_t}\varepsilon_t^i, \quad (28)$$

$$719 \text{s.t., } g(t, i) = \sum_{m=1}^t \binom{t}{m} \rho^{(t-m)/2} (1-\rho)^{m/2} z_0^{i-m}. \quad (29)$$

720 6.2. The proof of Eqn.(16)

721 In order to solve the posterior distribution of STDD, we
722 need to rethink our forward diffusion model in Eqn. (28).
723 Based on the special operator F in Eqn. (9), the forward
724 diffusion process can be written as

$$725 Z_t = \sqrt{\bar{\alpha}_t}F_t(Z_0) + \sqrt{1-\bar{\alpha}_t}\hat{\varepsilon}_t, \quad (30)$$

726 Since $\sum_{m=0}^t \binom{t}{m} \rho^{(t-m)/2} (1-\rho)^{m/2} \hat{\varepsilon}_t^m = (\sqrt{\rho^2 + \sqrt{1-\rho^2}})^{t/2} * \varepsilon = \varepsilon$, based on the properties of the Gaussian distribution, we can get:

$$729 Z_t = \sqrt{\bar{\alpha}_t}F_t(Z_0) + \sqrt{1-\bar{\alpha}_t}F_t(\bar{\varepsilon}_t), \quad (31)$$

730 where, Gaussian noise ε_t is not independent, and can be obtained by the independent noise $\bar{\varepsilon}_t$ diffusion of each frame

732 image (i.e., $\varepsilon = F_t(\bar{\varepsilon}_t)$). Therefore, the probabilistic model
733 of forward diffusion can be written as:

$$\begin{cases} q(Z_t|Z_{t-1}, Z_0) = \mathcal{N}(Z_t|\sqrt{\alpha_t}F_t(F_{t-1}^{-1}(Z_{t-1})), \beta_t I) \\ q(Z_t|Z_0) = \mathcal{N}(Z_t|\sqrt{\bar{\alpha}_t}F_t(Z_0), \sqrt{1-\bar{\alpha}_t}I) \end{cases} \quad (32)$$

734 The posterior distribution can be obtained from the prior
735 distribution via Bayes' formula:
736

$$\begin{aligned} q(Z_{t-1}|Z_t) &= q(Z_{t-1}|Z_t, Z_0) \\ &= q(Z_t|Z_{t-1}, Z_0) \frac{q(Z_{t-1}|Z_0)}{q(Z_t|Z_0)} \quad (33) \\ &= N(\tilde{\mu}_t, \tilde{\beta}_t). \end{aligned}$$

737 We obtain the mean and variance of the posterior distribution as:
738

$$\begin{cases} \tilde{\mu}_t = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}F_{t-1}(Z_0) + \frac{\sqrt{\bar{\alpha}_t}(1-\alpha_{t-1})}{1-\bar{\alpha}_t}F_{t-1}(F_t^{-1}(Z_t)) \\ \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t \end{cases} \quad (34)$$

740 By substituting Z_o for Z_t with Eqn. (31), the mean and variance
741 of the STDD posterior distribution are finally obtained:
742

$$\begin{cases} \tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(F_{t-1}(F_t^{-1}(Z_t)) - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_t \right) \\ \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t, \end{cases} \quad (35)$$

743 6.3. Visualization of generated videos

744 We also provide some sampled videos on UCF101 and Sky
745 Time-Lapse, shown in Fig. 5 and Fig. 6. Each line is 16
746 frames of a sampled video,
747



Figure 5. Sample videos on UCF101.

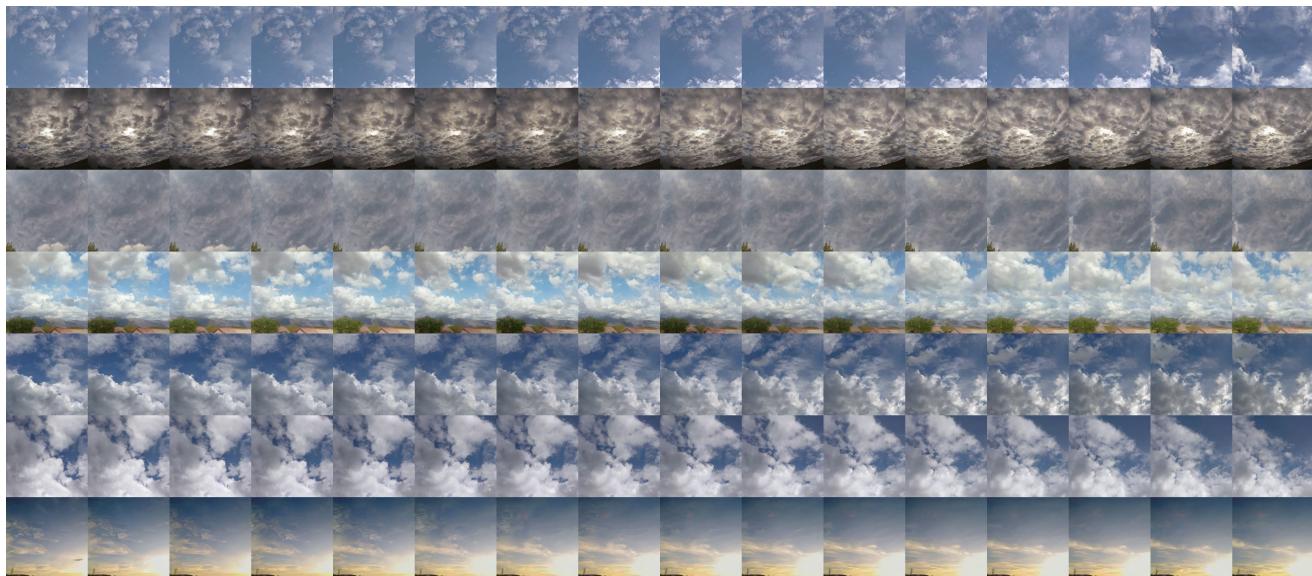


Figure 6. Sample videos on Sky Time-Lapse.