



Carnegie Mellon University

# Lecture 6: The Design Space of Diffusion Models & Solvers for Fast Sampling

---

Yutong (Kelly) He

*10-799 Diffusion & Flow Matching, Jan 27<sup>th</sup>, 2026*



Modal



# Housekeeping Announcements

- Homework 2 is out! <https://kellyyutonghe.github.io/10799S26/homework/>
  - Due date: 2/3 Tue, Late Due date: 2/5 Thur
  - Training models takes time! Start early!
  - Maybe change in due date? Vote in Discord!
- Quiz 3 next class!

# Now that we have learned the basics

Is DDPM perfect now? What can we improve?



# Us defending our DDPM model from HW1 be like...



**Carnegie  
Mellon  
University**

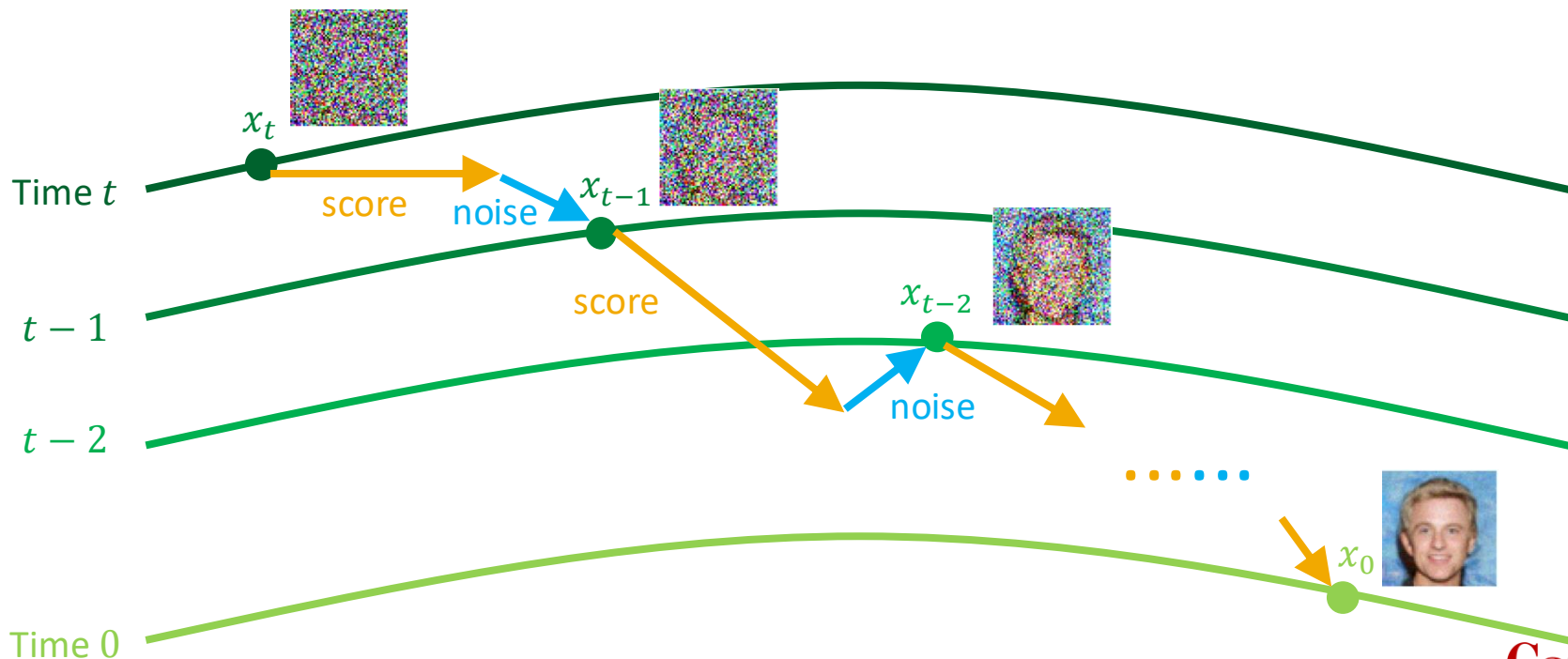
# Vanilla DDPM is sloooooooooow

- Take 1000 steps to generate an image
- The image quality severely degrades if you use fewer steps
- Scaling to higher resolution/larger model size becomes nightmare

=>

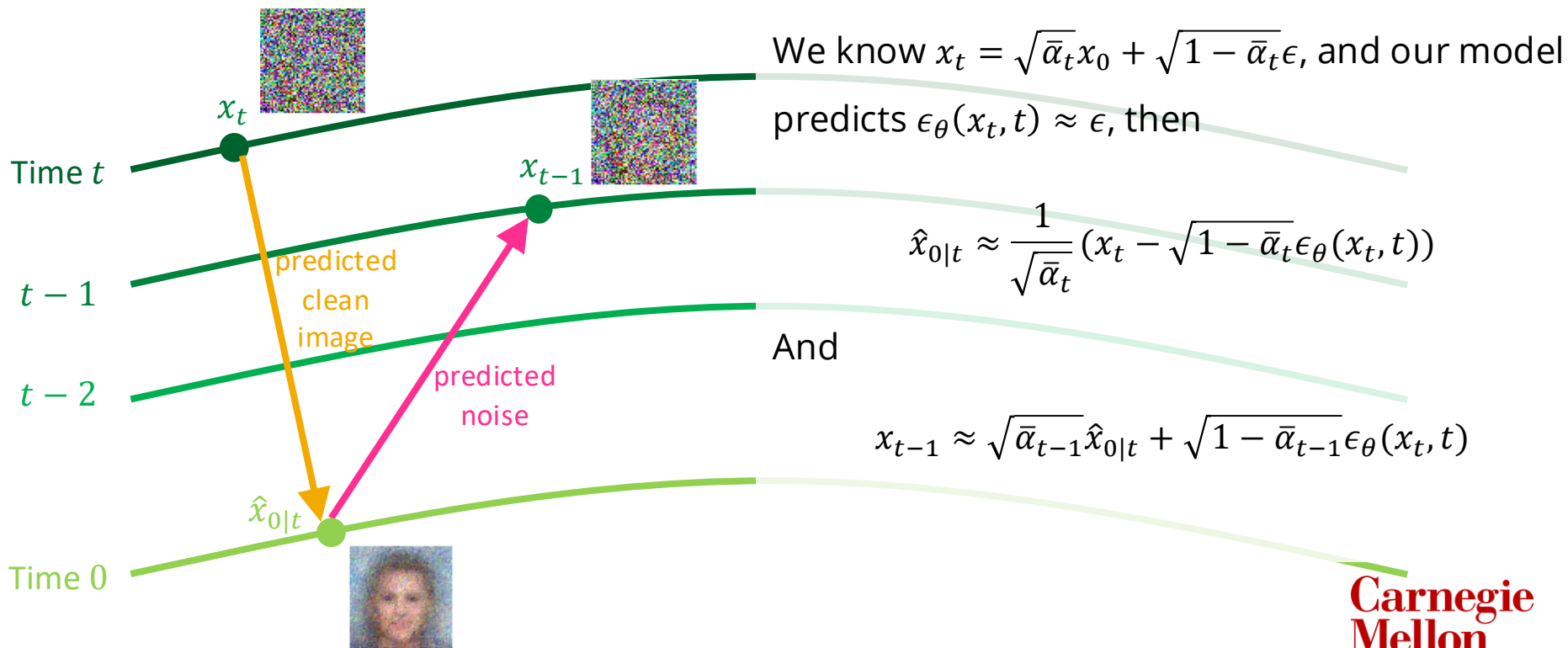
- Any application that needs real time generation ❌
- Video is going to be a big pain ☹️
- GPU cost 💰

# Currently, the DDPM sampling is like

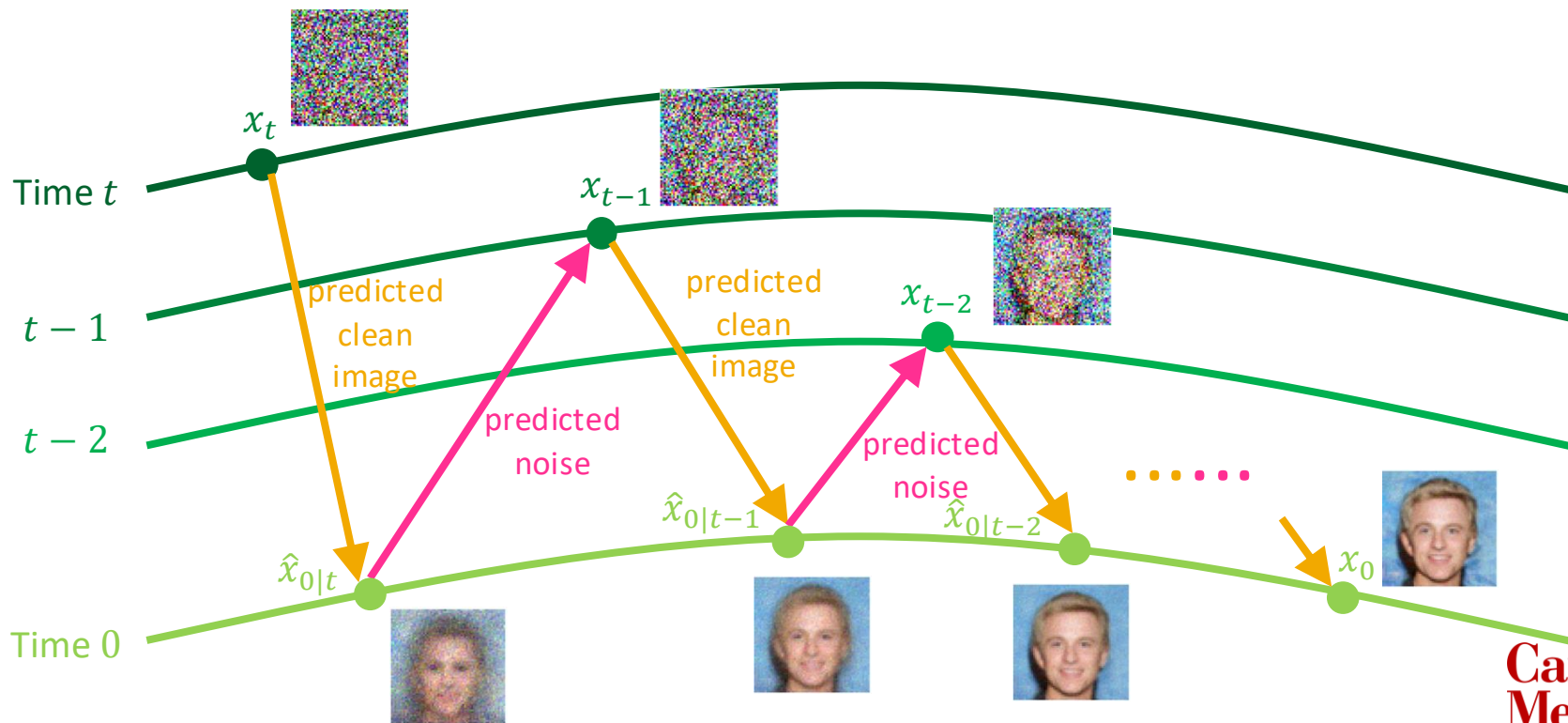


# But there is another way!

No, you can't! There must be another way.



# So we can also do sampling like this



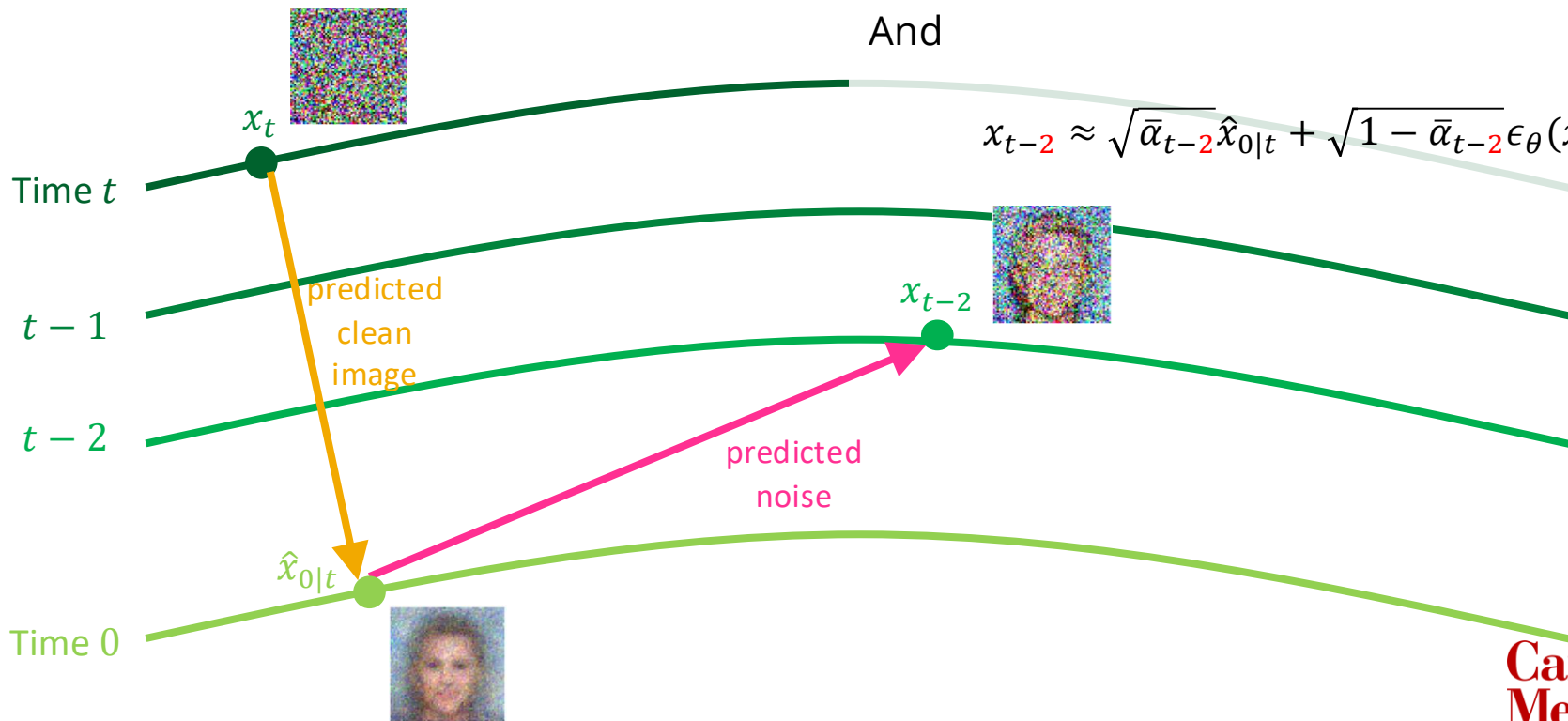


# Now what if we skip t-1

$$\hat{x}_{0|t} \approx \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta}(x_t, t))$$

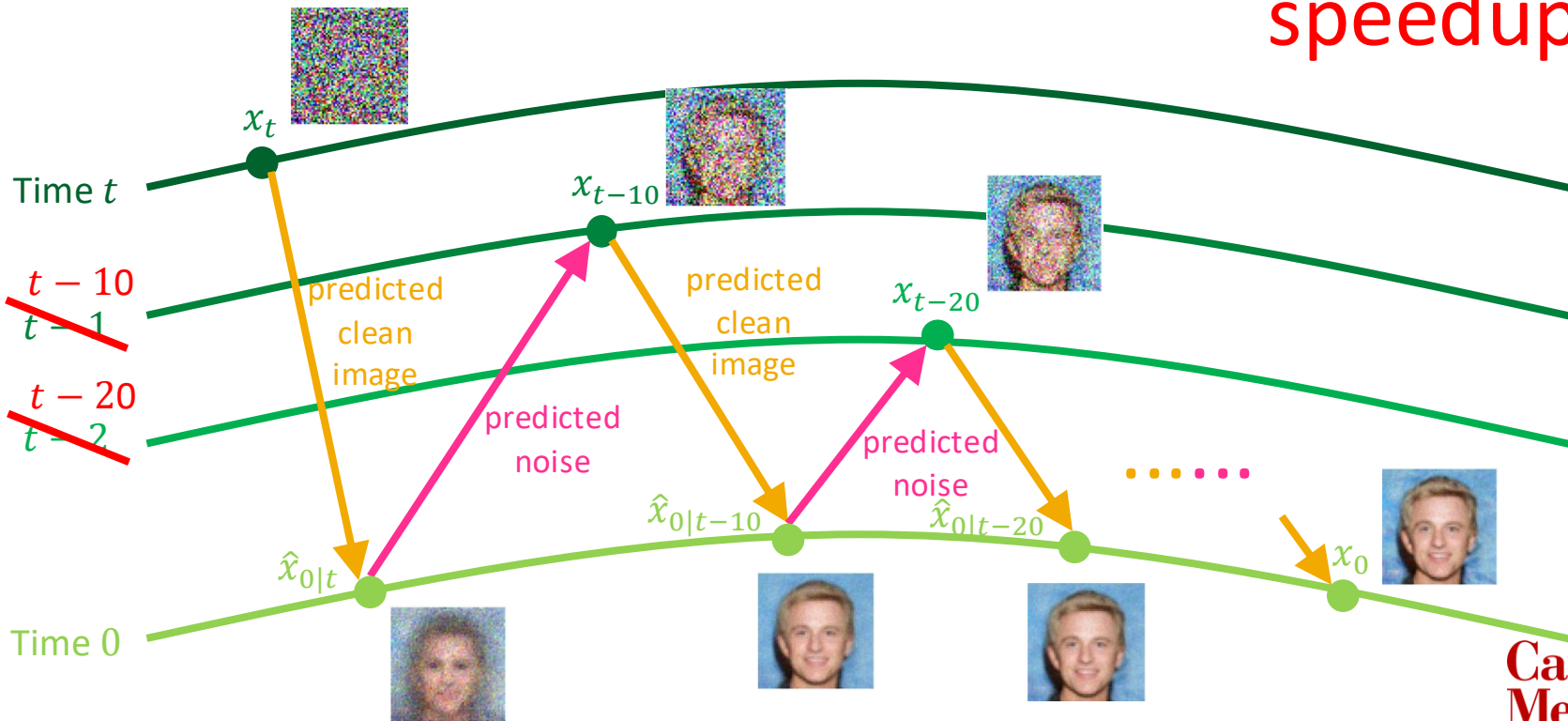
And

$$x_{t-2} \approx \sqrt{\bar{\alpha}_{t-2}} \hat{x}_{0|t} + \sqrt{1 - \bar{\alpha}_{t-2}} \epsilon_{\theta}(x_t, t)$$

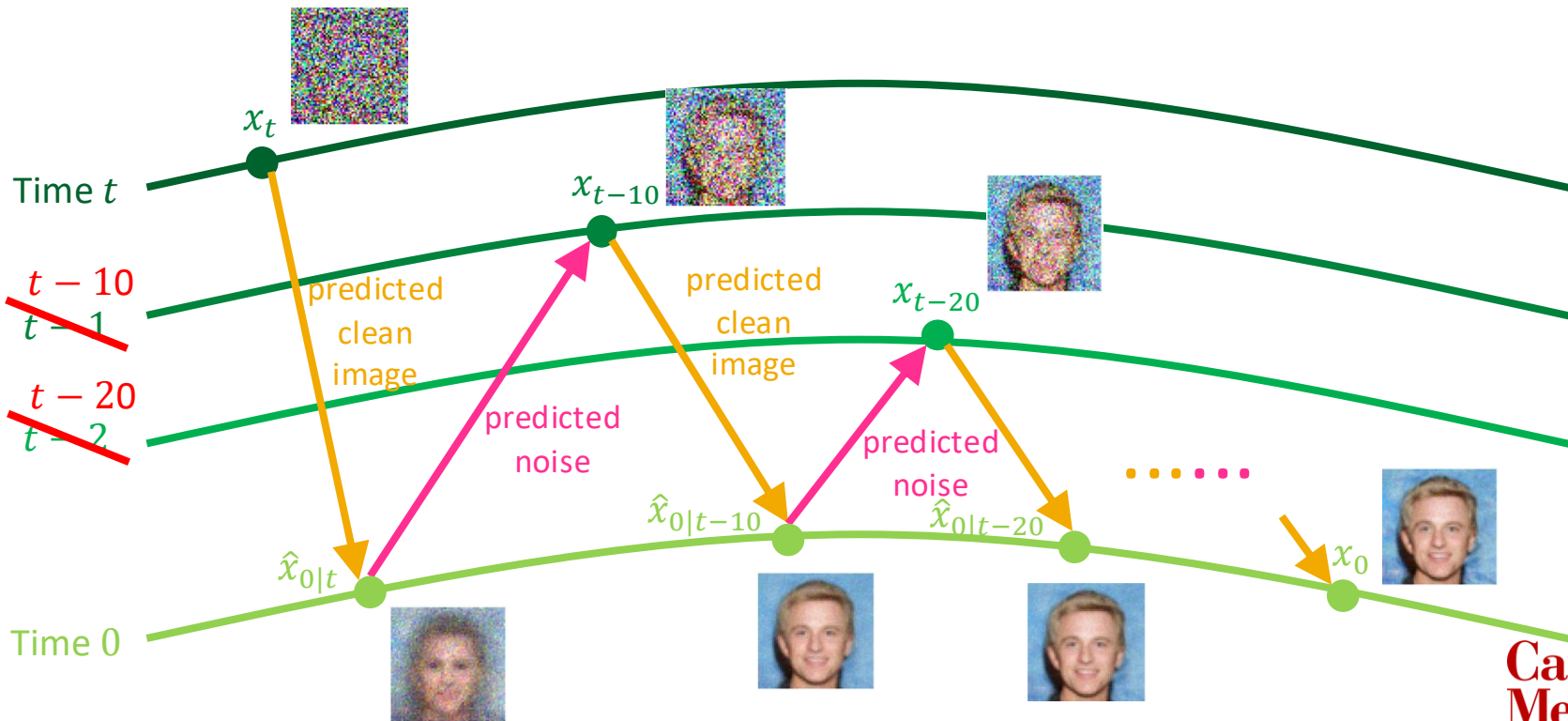


# We can skip more than 1 step!

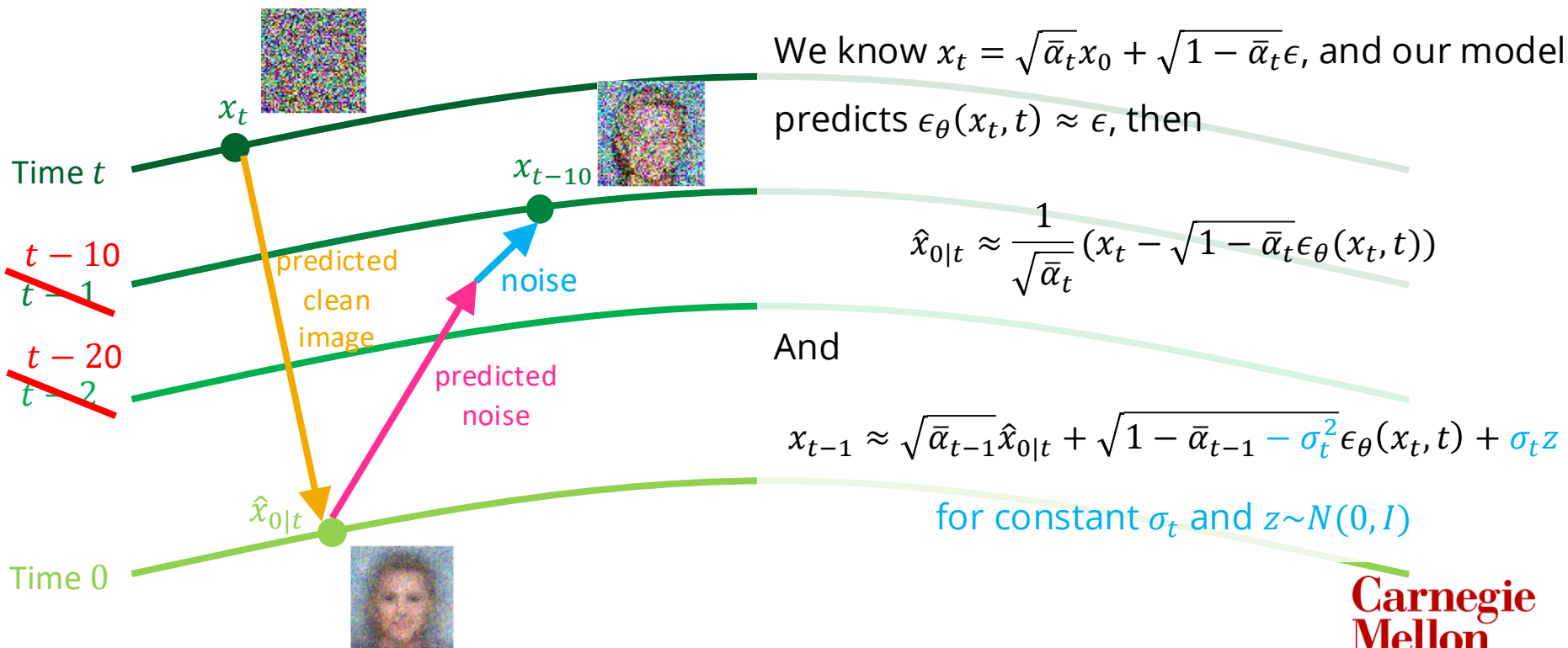
# 10X speedup!



# Notice how right now everything is deterministic

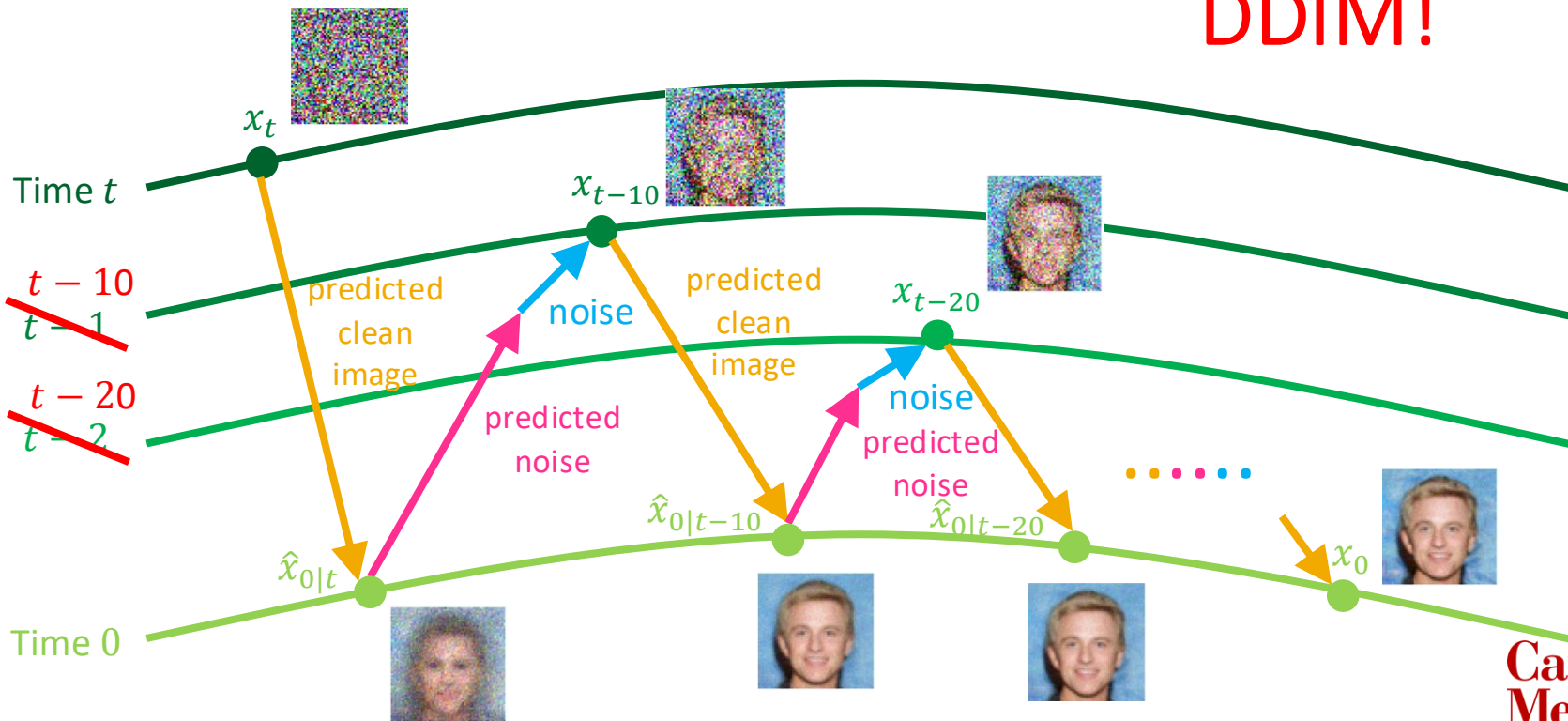


# But we can add the stochasticity back!



# Putting everything together

Now you have  
DDIM!



# DDIM: The OG diffusion fast sampling algorithm

$$\mathbf{x}_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \left( \frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right)}_{\text{"predicted } \mathbf{x}_0"} + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}_{\text{"direction pointing to } \mathbf{x}_t"} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

The pseudocode of the deterministic version of DDIM (also in HW2)

---

## Algorithm 1 DDIM Sampling

---

**Require:** trained noise predictor  $\epsilon_{\theta}$ , number of steps  $S$ , noise schedules  $\bar{\alpha}$

- 1: Sample  $x_T \sim \mathcal{N}(0, I)$
  - 2: Create timestep subsequence  $[\tau_S, \tau_{S-1}, \dots, \tau_1]$  from  $[T, \dots, 1]$  ▷ e.g., [1000, 900, 800, ...]
  - 3: **for**  $i = S, S-1, \dots, 1$  **do**
  - 4:    $t \leftarrow \tau_i$
  - 5:    $t_{\text{prev}} \leftarrow \tau_{i-1}$  (or 0 if  $i = 1$ )
  - 6:    $\epsilon \leftarrow \epsilon_{\theta}(x_t, t)$  ▷ Predict noise using your trained DDPM
  - 7:    $\hat{x}_0 \leftarrow \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon}{\sqrt{\bar{\alpha}_t}}$  ▷ Predict clean image
  - 8:    $x_{t_{\text{prev}}} \leftarrow \sqrt{\bar{\alpha}_{t_{\text{prev}}}} \cdot \hat{x}_0 + \sqrt{1 - \bar{\alpha}_{t_{\text{prev}}}} \cdot \epsilon$  ▷ DDIM step
  - 9: **end for**
  - 10: **return**  $x_0$
-

# Any other ways to sample?



# Remember how everything can be ODE now

## Cond-OT flow matching:

$$p_0 = N(0, I), p_1 = \delta(x_1)$$

$$x_t = tx_1 + (1 - t)x_0, \quad x_0 \sim p_0$$

$$p_t(x_t|x_1) = N(tx_1, (1 - t)^2 I)$$

$$\frac{dx_t}{dt} = u(x_t|x_1) = x_1 - x_0$$

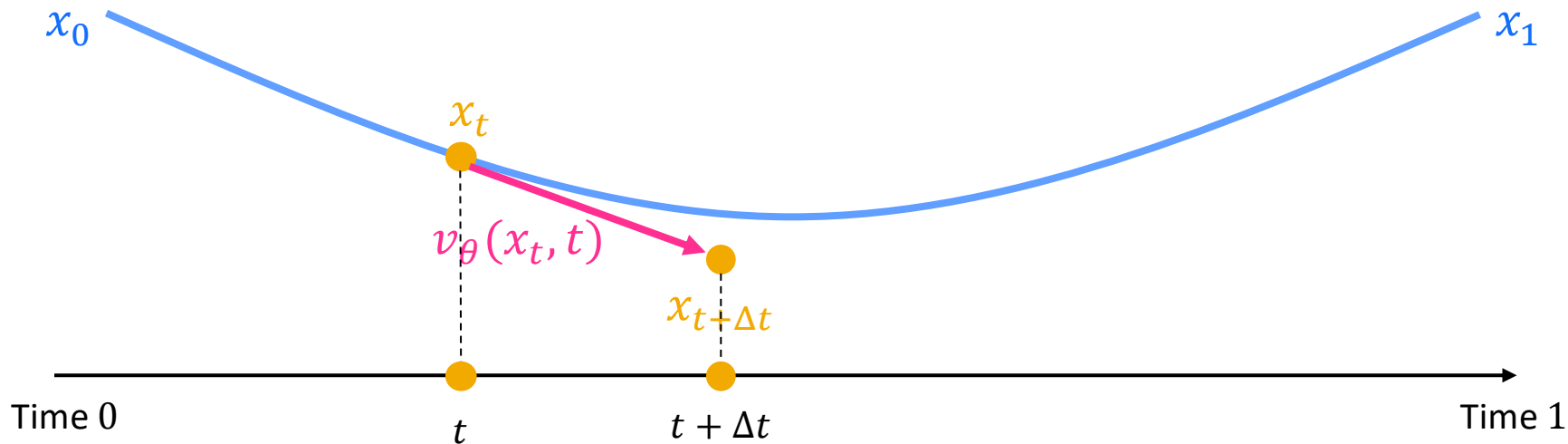
## Probability flow ODE:

$$d\mathbf{x} = \left[ \mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt,$$

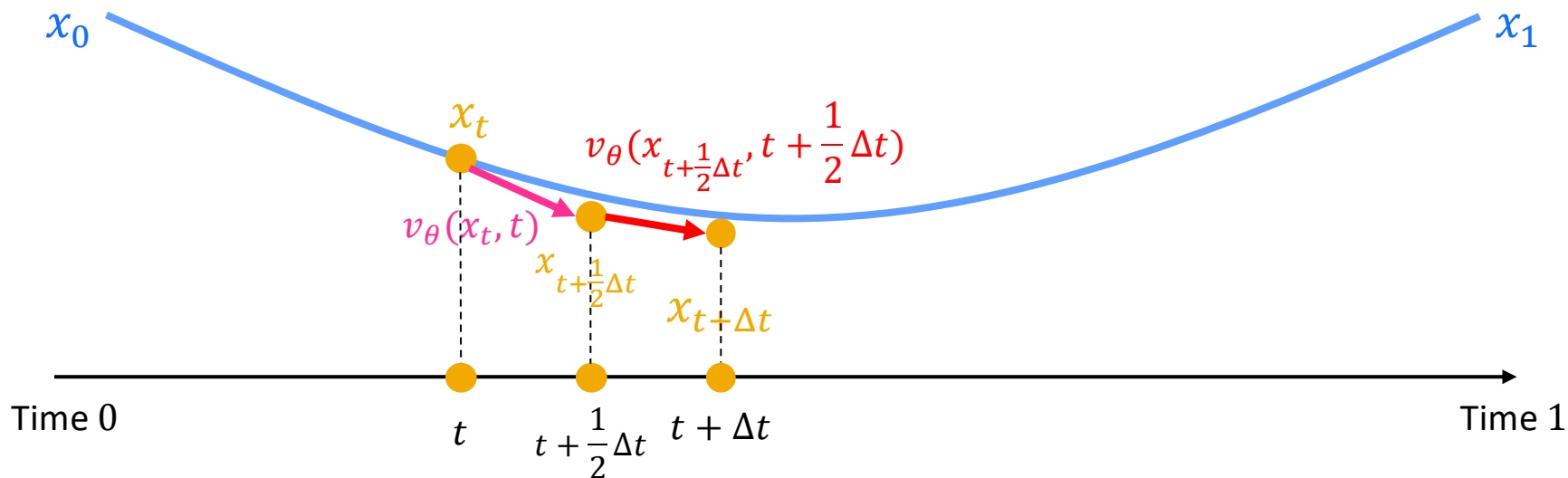
We can use different solvers to solve the ODEs!



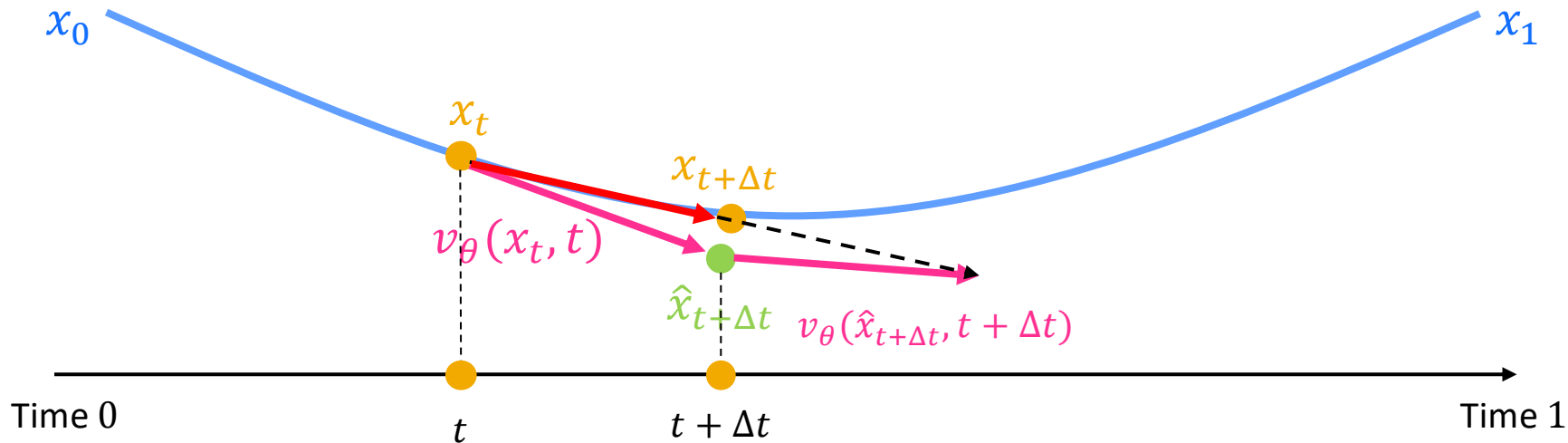
# Classic ODE solver 1 – Euler solver



# Classic ODE solver 2 – Midpoint solver



# Classic ODE solver 3 – 2<sup>nd</sup> order Heun solver



# Diffusion specific solver – DPM solver

The diffusion ODE is

$$\frac{d\mathbf{x}_t}{dt} = \underbrace{f(t)\mathbf{x}_t}_{\text{linear}} - \underbrace{\frac{1}{2}g^2(t)\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)}_{\text{model learned}},$$

$f(t) = d \log \alpha_t / dt$ 
 $g^2(t) = \frac{d\sigma_t^2}{dt} - 2 \frac{d \log \alpha_t}{dt} \sigma_t^2 = 2\sigma_t^2 \left( \frac{d \log \sigma_t}{dt} - \frac{d \log \alpha_t}{dt} \right) = -2\sigma_t^2 \frac{d\lambda_t}{dt}.$

$\Downarrow$

$$\mathbf{x}_t = e^{\int_s^t f(\tau) d\tau} \mathbf{x}_s + \int_s^t \left( e^{\int_\tau^t f(r) dr} \frac{g^2(\tau)}{2\sigma_\tau} \epsilon_\theta(\mathbf{x}_\tau, \tau) \right) d\tau.$$

# Diffusion specific solver – DPM solver

From here 
$$\mathbf{x}_{t_{i-1} \rightarrow t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \hat{\epsilon}_{\theta}(\hat{\mathbf{x}}_{\lambda}, \lambda) d\lambda.$$

We can channel Taylor expansion

$$\hat{\epsilon}_{\theta}(\hat{\mathbf{x}}_{\lambda}, \lambda) = \sum_{n=0}^{k-1} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} \hat{\epsilon}_{\theta}^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}}) + \mathcal{O}((\lambda - \lambda_{t_{i-1}})^k),$$



$$\mathbf{x}_{t_{i-1} \rightarrow t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \sum_{n=0}^{k-1} \hat{\epsilon}_{\theta}^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}}) \underbrace{\int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda}_{\text{Can be calculated analytically}} + \mathcal{O}(h_i^{k+1}),$$

Can be calculated analytically

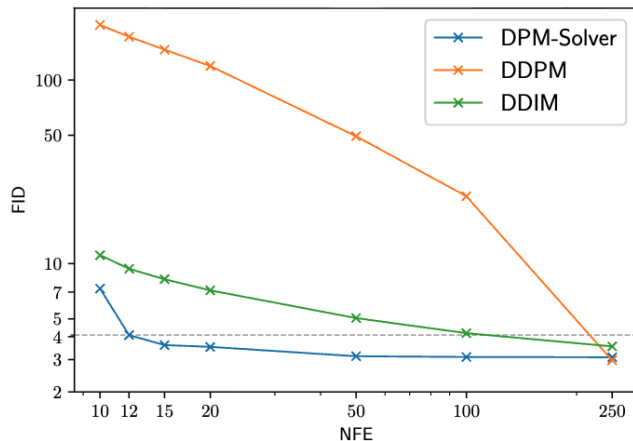
# Diffusion specific solver – DPM solver

$$\mathbf{x}_{t_{i-1} \rightarrow t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \sum_{n=0}^{k-1} \underbrace{\hat{\epsilon}_{\theta}^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}})}_{\text{Can be estimated}} \underbrace{\int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda}_{\text{Can be calculated analytically}} + \underbrace{\mathcal{O}(h_i^{k+1})}_{\text{Ignored}},$$

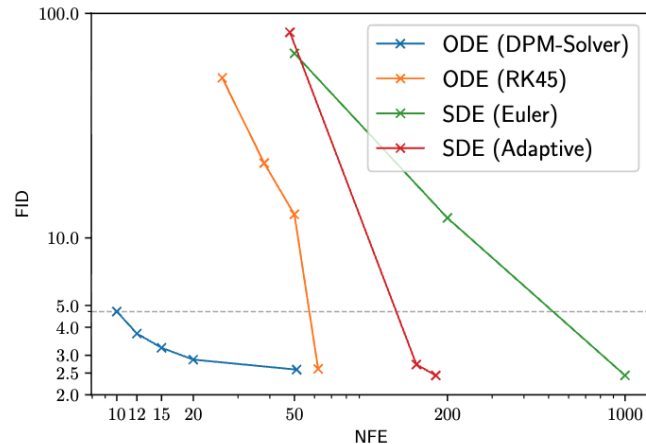
Finally we get for  $k=1$

$$\tilde{\mathbf{x}}_{t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \sigma_{t_i} (e^{h_i} - 1) \epsilon_{\theta}(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1})$$

# Solver comparison

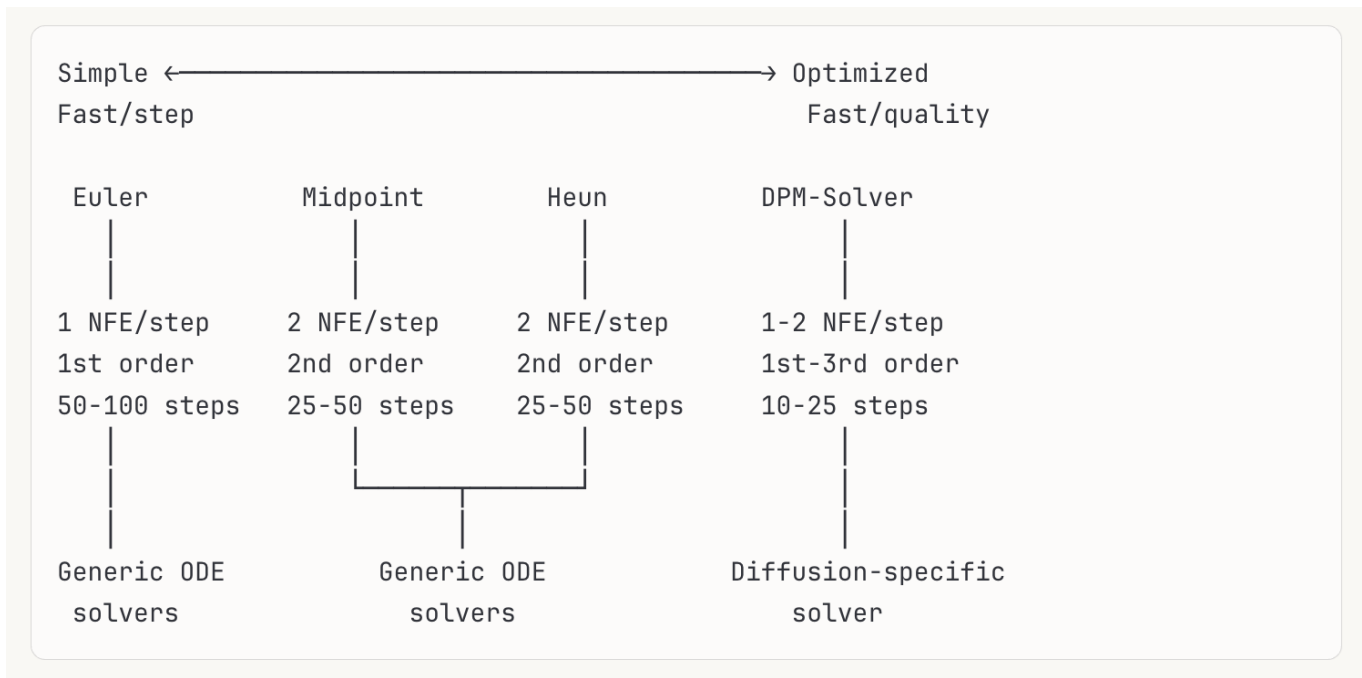


(e) ImageNet 128x128 (discrete)



(a) CIFAR-10 (continuous)

# Solver comparison (from Claude)





**Besides sampling, any other ways to improve DDPM?**



# Actually there is a paper that studied them all

**EDM:** What are the actual independent design choices of diffusion models and let's disentangle them as study them one by one

---

## Elucidating the Design Space of Diffusion-Based Generative Models

---

**Tero Karras**  
NVIDIA

**Miika Aittala**  
NVIDIA

**Timo Aila**  
NVIDIA

**Samuli Laine**  
NVIDIA

# The design space of diffusion models

## The design space of diffusion models



### Training

- Prefixed noise schedule
- Training noise sampling schedule
- Loss weighting w.r.t. time

### Model

- Reparameterization
- Input/Output scaling
- How to do time conditioning

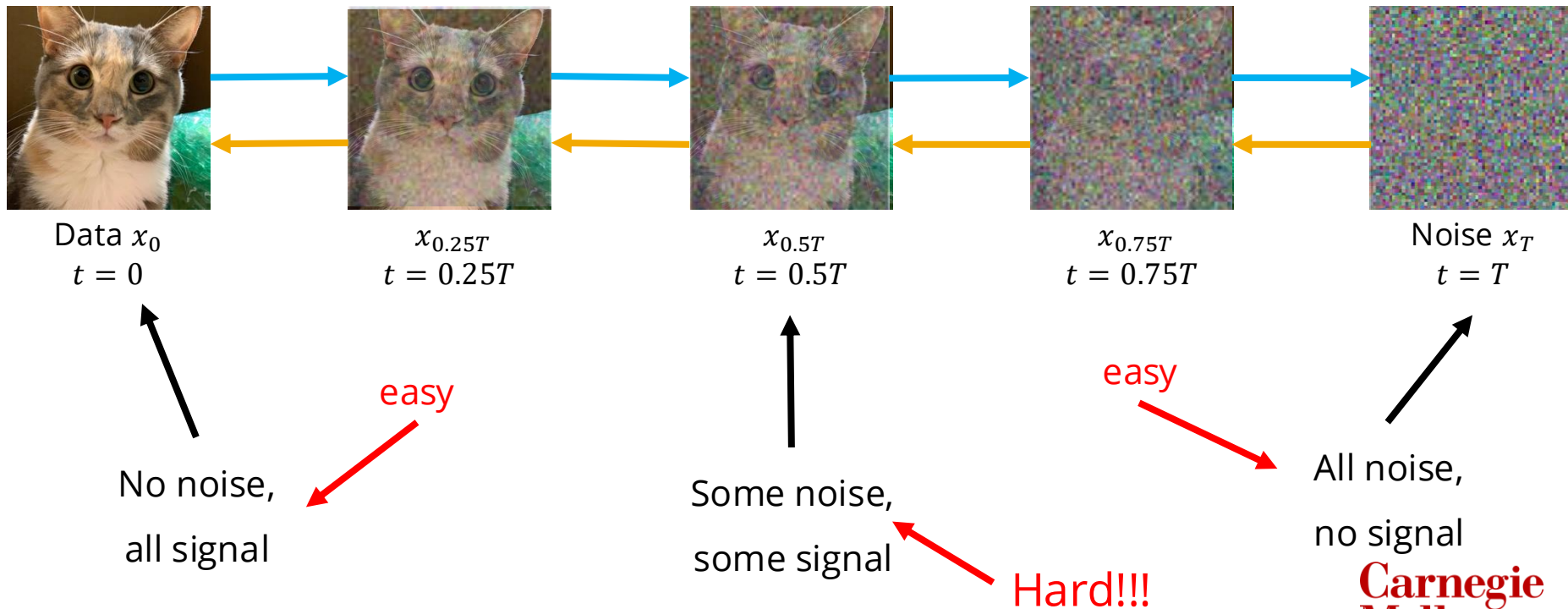
### Sampling

- Solver 
- Sampling time noise schedule
- Number of time steps 

SNR = Signal/Noise

$$\text{SNR} = \frac{\alpha_t^2}{\sigma_t^2}$$

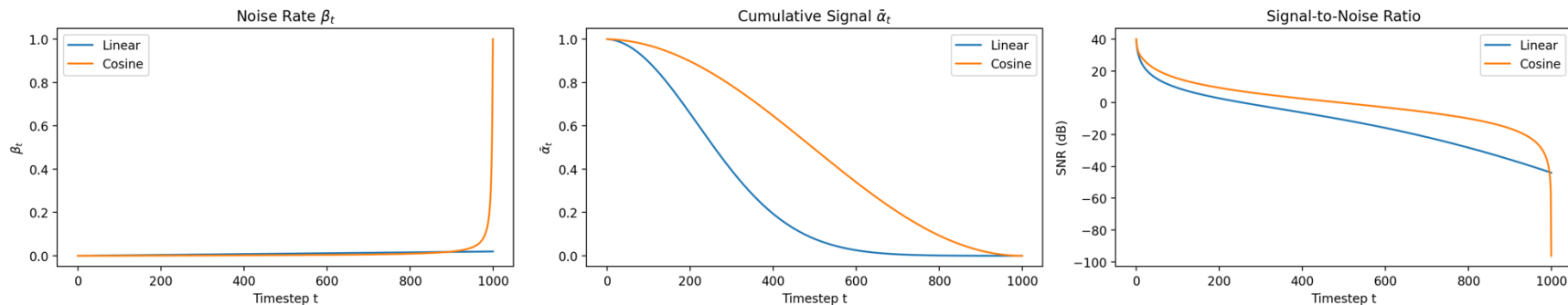
# Signal-to-noise ratio (SNR)



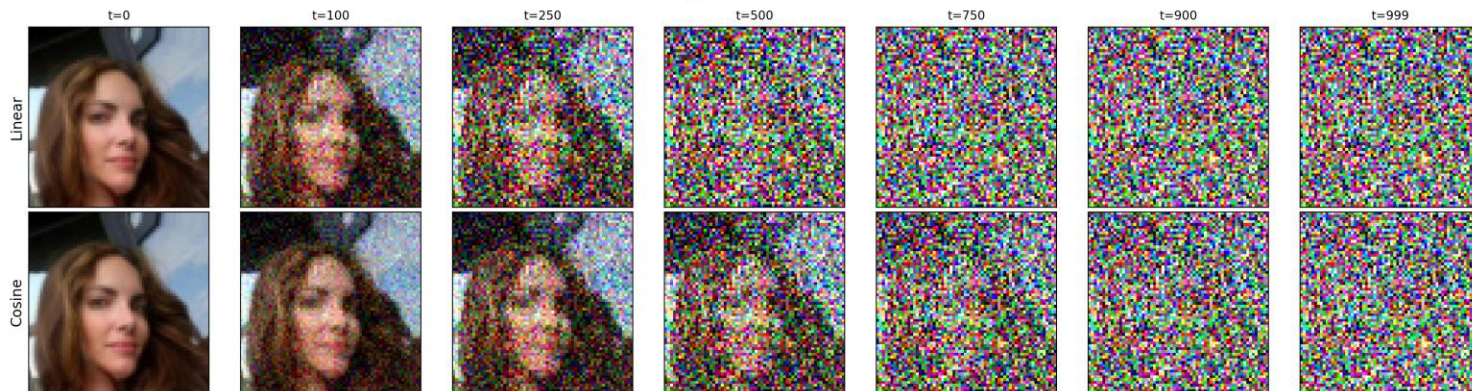
**Carnegie  
Mellon  
University**

# SNR Attempt 1: Linear scheduler -> Cosine scheduler

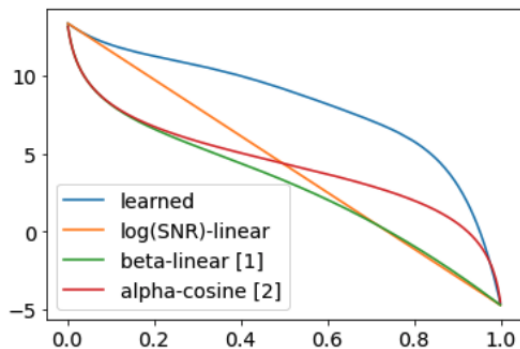
Linear vs Cosine Noise Schedule (T=1000)



Forward Noising: Linear vs Cosine Schedule



# SNR Attempt 1.5: You can even learn your schedule



(a) log SNR vs time  $t$

SNR( $t$ ) schedule	Var(BPD)
<b>Learned (ours)</b>	<b>0.53</b>
log SNR-linear	6.35
$\beta$ -Linear [1]	31.6
$\alpha$ -Cosine [2]	31.1

(b) Variance of VLB estimate

**SNR Attempt 2: At sampling time, you can also spend more time on the more difficult noise levels!**

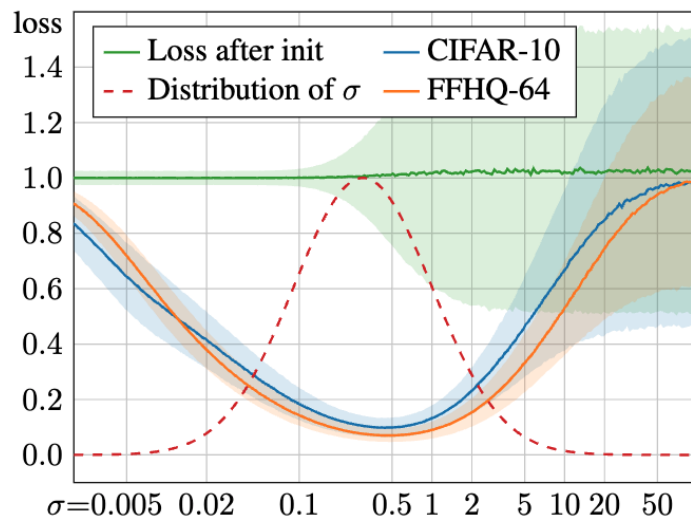
$$\sigma_{i < N} = \left( \sigma_{\max}^{\frac{1}{\rho}} + \frac{i}{N-1} (\sigma_{\min}^{\frac{1}{\rho}} - \sigma_{\max}^{\frac{1}{\rho}}) \right)^{\rho} \quad \text{and} \quad \sigma_N = 0.$$

**SNR Attempt 3: You should parameterize your model to take the actual noise level ( $\sigma$  or  $\log \sigma$ ) instead of timesteps!**

- Scheduler invariant parameterization
- Continuous => Easier to do make inference time changes (fewer/more steps)
- Better numerical scaling

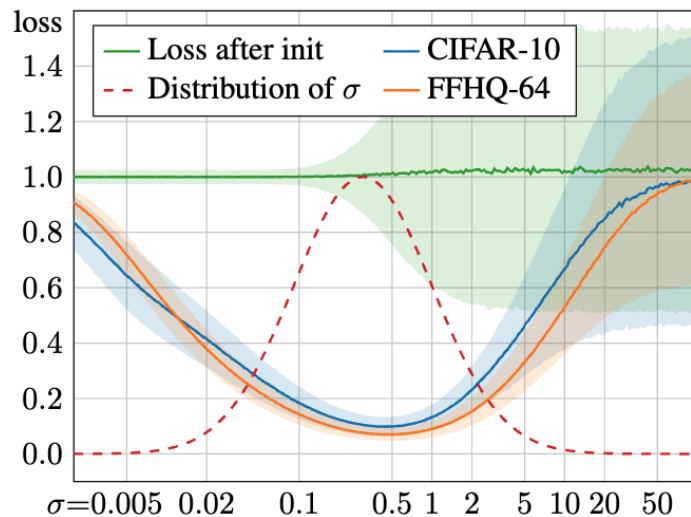
# SNR Attempt 4: Train on harder time steps more

Just sampling the more difficult time steps to more frequently during training time!





# SNR Attempt 5: Weight different timestep differently according to their SNR

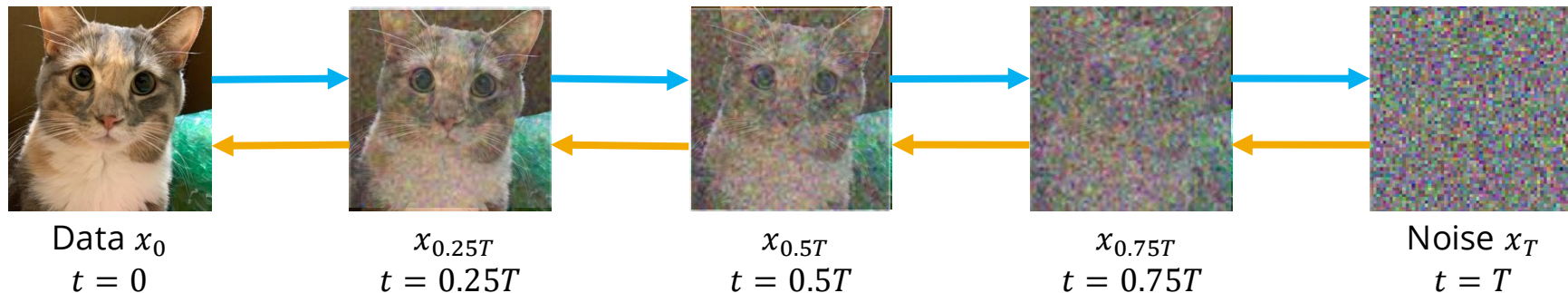


The gradient scale varies a lot depending on the SNR



Apply a weighting scalar to balance it out!

# SNR Attempt 6: Reparameterization



## Predict clean signal ( $x_{0,\theta}(x_t, t)$ ):

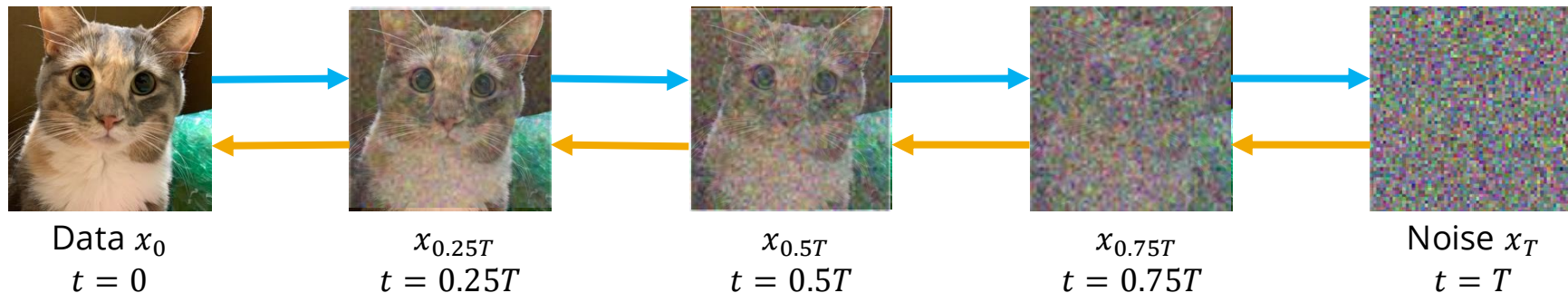
- Easy at low noise (high SNR)
- Hard at high noise (low SNR)

## Predict noise ( $\epsilon_\theta(x_t, t)$ ):

- Easy at high noise (low SNR)
- Hard at low noise (high SNR)

Something in between?

# SNR Attempt 6: Reparameterization



**Predict clean signal ( $x_{0,\theta}(x_t, t)$ ):**

- Easy at low noise (high SNR)
- Hard at high noise (low SNR)

**Predict interpolation**

**(a.k.a velocity  $v_\theta(x_t, t)$ ):**

- Good balance!
- $v = \alpha_t \epsilon - \sigma_t x_0$

**Predict noise ( $\epsilon_\theta(x_t, t)$ ):**

- Easy at high noise (low SNR)
- Hard at low noise (high SNR)

# Another way to do network parameterization?

Remember how previously we want to do v-prediction as our target

$$v = \alpha_t \epsilon - \sigma_t x_0$$

A mixture of noise  
and clean data

Now notice how  $x_t$  is also a mixture of noise and clean data

⇒ The network should reuse some information in  $x_t$

⇒ Skip connection!

# Input/Output scaling & Preconditioning

$$D_{\theta}(\mathbf{x}; \sigma) = c_{\text{skip}}(\sigma) \mathbf{x} + c_{\text{out}}(\sigma) F_{\theta}(c_{\text{in}}(\sigma) \mathbf{x}; c_{\text{noise}}(\sigma)),$$

Rescale input

Time/noise  
conditioning

Predicts clean data  $\mathbf{x}_0$

Skip connection weight

Network output weight

Trained network

- When  $t \rightarrow 0, \sigma \rightarrow 0 \Rightarrow$  input is mostly clean  $\Rightarrow$  can directly pass through more
- When  $t \rightarrow 1, \sigma \rightarrow \infty \Rightarrow$  input is mostly noise  $\Rightarrow$  should ignore most of the input and rely on the network prediction more

# Input/Output scaling & Preconditioning

$$D_{\theta}(\mathbf{x}; \sigma) = c_{\text{skip}}(\sigma) \mathbf{x} + c_{\text{out}}(\sigma) F_{\theta}(c_{\text{in}}(\sigma) \mathbf{x}; c_{\text{noise}}(\sigma)),$$

Rescale input

Time/noise  
conditioning

Predicts clean data  $\mathbf{x}_0$

Skip connection weight

Network output weight

Trained network

We also want:




- The network input to have unit variance
- Training target to have unit variance
- Reuse information in input as much as possible

$$\Rightarrow c_{\text{skip}}(\sigma) = \frac{\sigma_{\text{data}}^2}{\sigma^2 + \sigma_{\text{data}}^2}, c_{\text{out}}(\sigma) = \frac{\sigma \sigma_{\text{data}}}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}}, c_{\text{in}}(\sigma) = \frac{1}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}}, c_{\text{noise}}(\sigma) = \frac{1}{4} \log \sigma$$




# The design space of diffusion models

## The design space of diffusion models




### Training

- Prefixed noise schedule 
- Training noise sampling schedule 
- Loss weighting w.r.t. time 

### Model

- Reparameterization 
- Input/Output scaling 
- How to do time conditioning 

### Sampling

- Solver 
- Sampling time noise schedule 
- Number of time steps 

# Next class we will learn how to turn an unconditional diffusion model into a conditional one

Spoiler alert: You may or may not need to train for it!

