



Carnegie Mellon University

Lecture 7: Guidance & Controllable Generation

Yutong (Kelly) He

10-799 Diffusion & Flow Matching, Jan 29th, 2026



Modal



Quiz time!

10 minutes

Closed-book

On Gradescope!!!

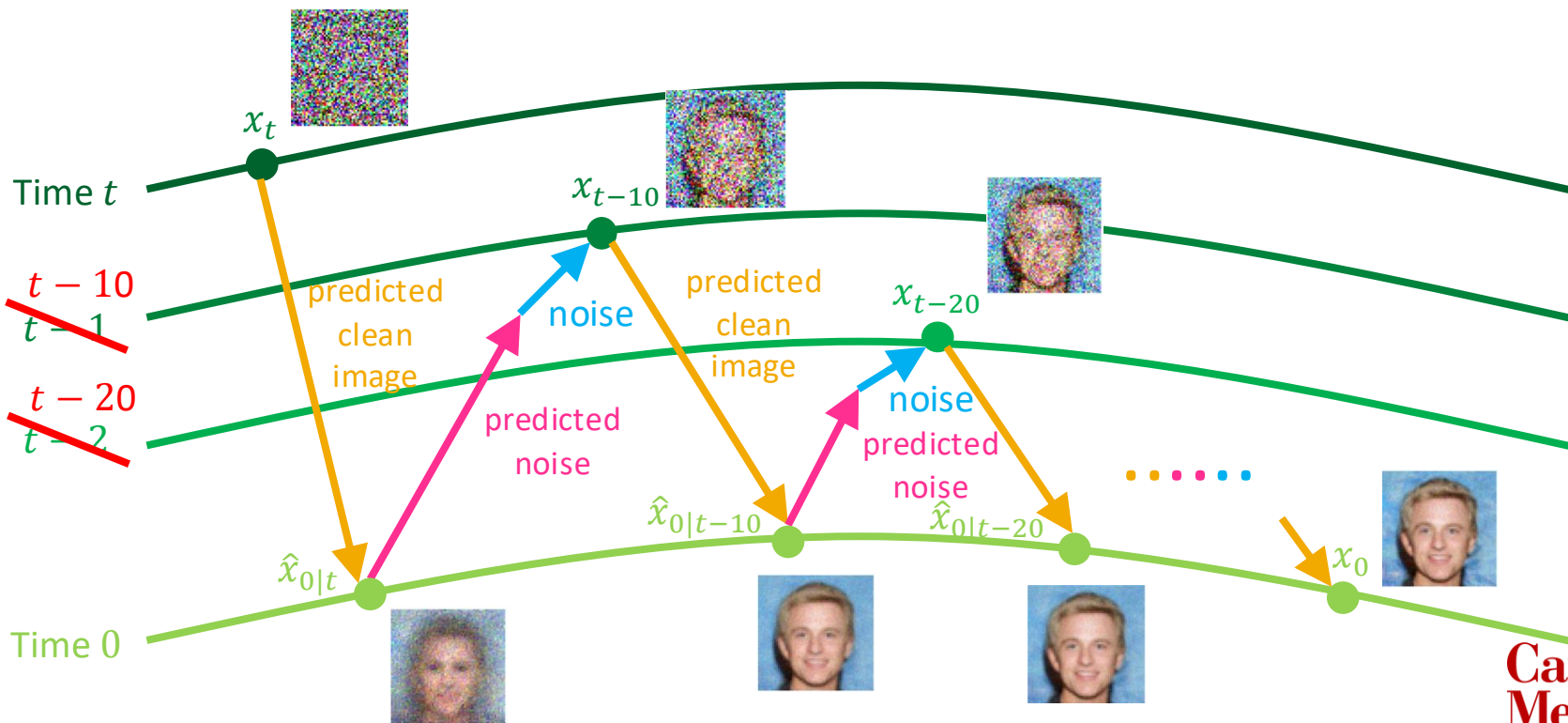


If you don't want to stay for the lecture, feel free to leave after submitting your quiz!

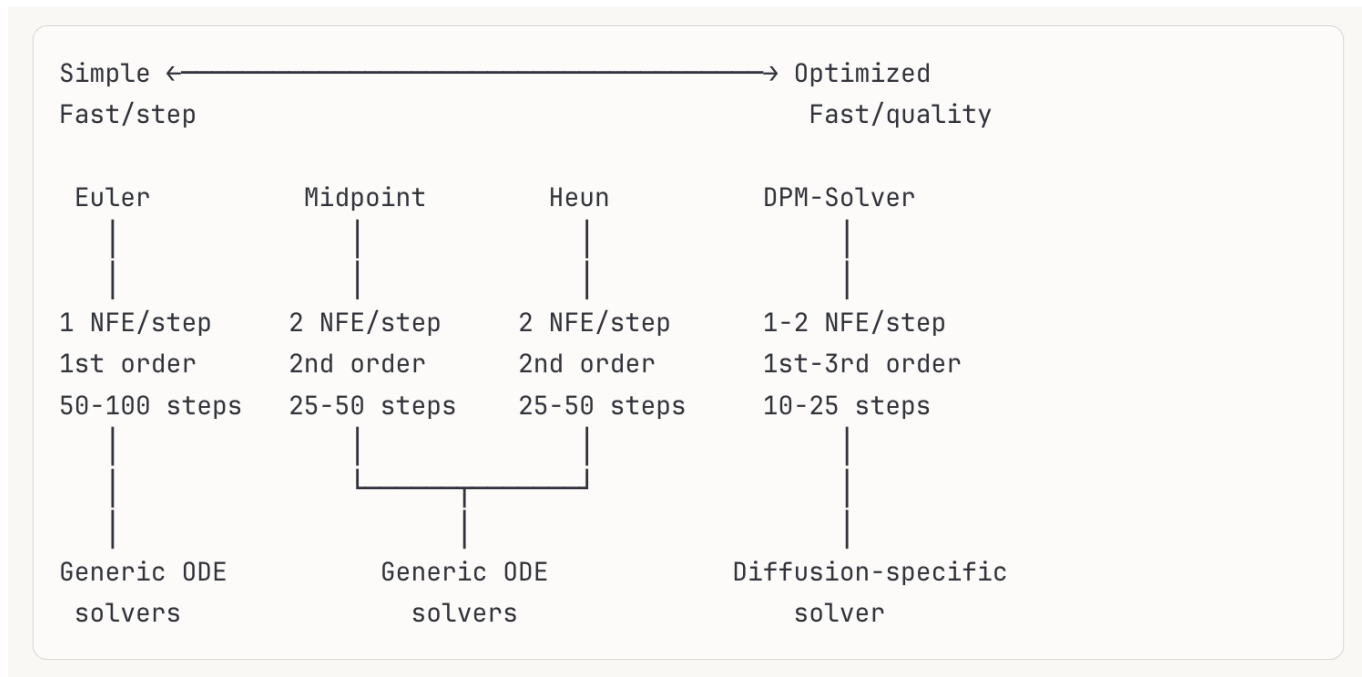
Housekeeping Announcements

- 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!
- **Next class** we will have **Max Simchowitz** come for a guest lecture/panel! Come to class **in person** if you want to ask him questions about diffusion/flow for robotics, control & decision making

Last time we learned DDIM



Solver comparison (from Claude)



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

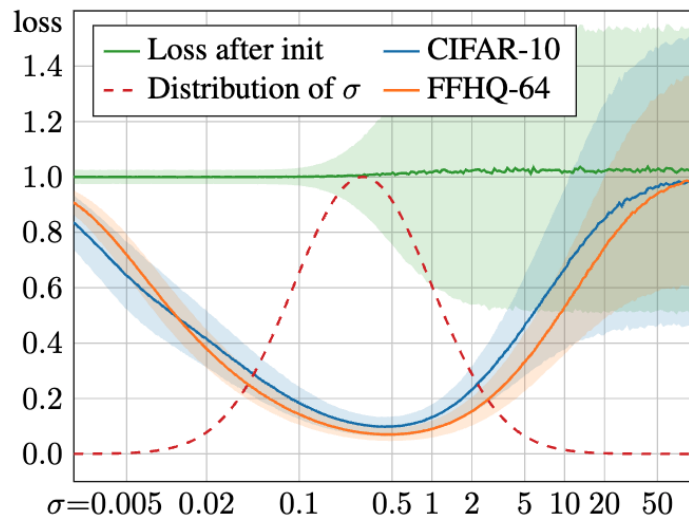
- Reparameterization
- Input/Output scaling
- How to do time conditioning

Sampling

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

Clarification on training noise distribution

Just sampling the **time steps that are more effective** during training time!

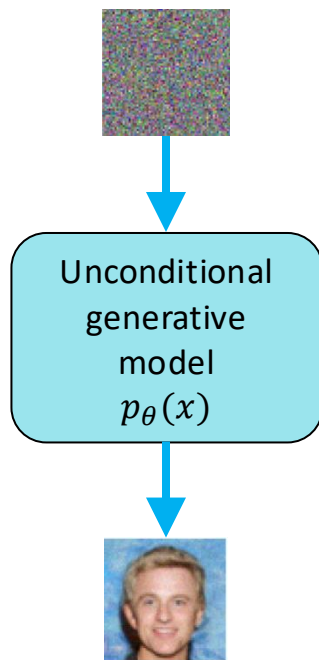


Are diffusion models perfect now?



How do normal people do conditional generation

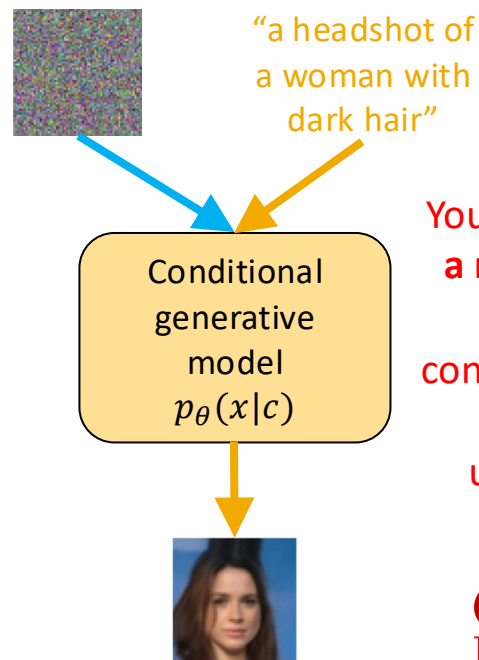
Unconditional generation:



But the unconditional model already has a lot of good knowledge on how to generate human faces in general...



Conditional generation:



"a headshot of a woman with dark hair"

You'll have to **train a new model** for each type of condition and can't reuse the unconditional model!

Remember the No.1 rule in probability theory

We know that diffusion models model the score function

$$\nabla_{x_t} \log p_\theta(x_t, t)$$

By Bayes' Theorem, we know that

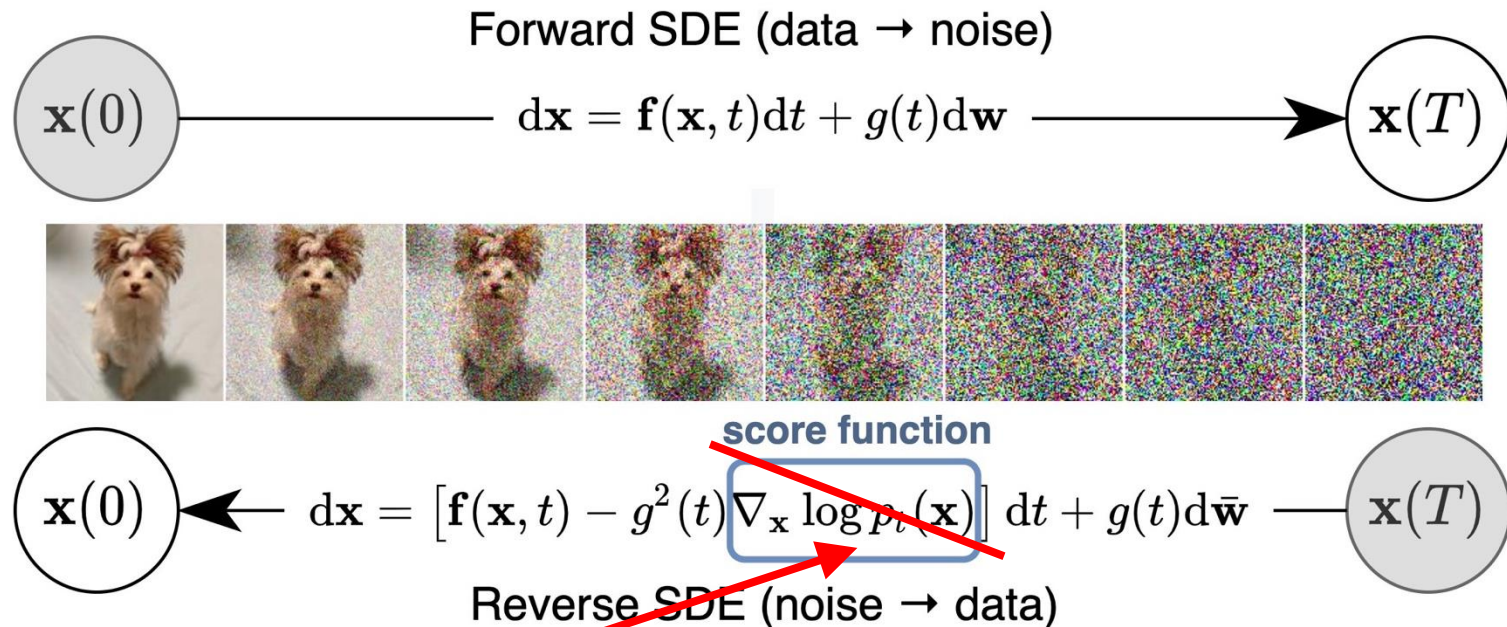
$$p(x|c) = \frac{p(x,c)}{p(c)} = \frac{p(c|x)p(x)}{p(c)} \propto p(c|x)p(x)$$

$$\nabla_{x_t} \log p_\theta(x_t, t|c) = \nabla_{x_t} \log p_\theta(c|x_t, t) + \nabla_{x_t} \log p_\theta(x_t, t)$$

Pre-trained
unconditional
diffusion model

A discriminative model!
(e.g. a classifier, a CLIP model, etc)

Classifier guidance diffusion

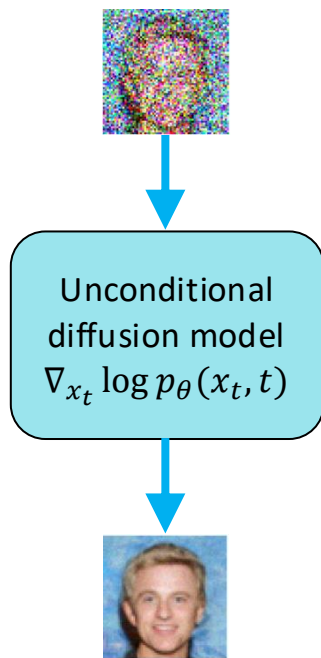


Conditional score:

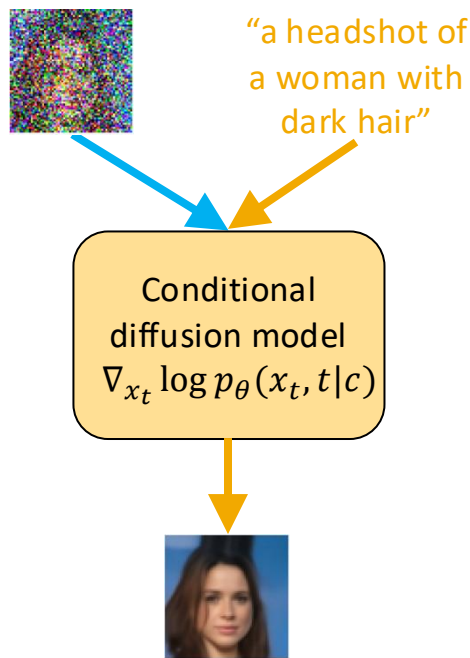
$$\nabla_{x_t} \log p_{\theta}(x_t, t|c) = \nabla_{x_t} \log p_{\theta}(c|x_t, t) + \nabla_{x_t} \log p_{\theta}(x_t, t)$$

Classifier guidance diffusion v.s. normies

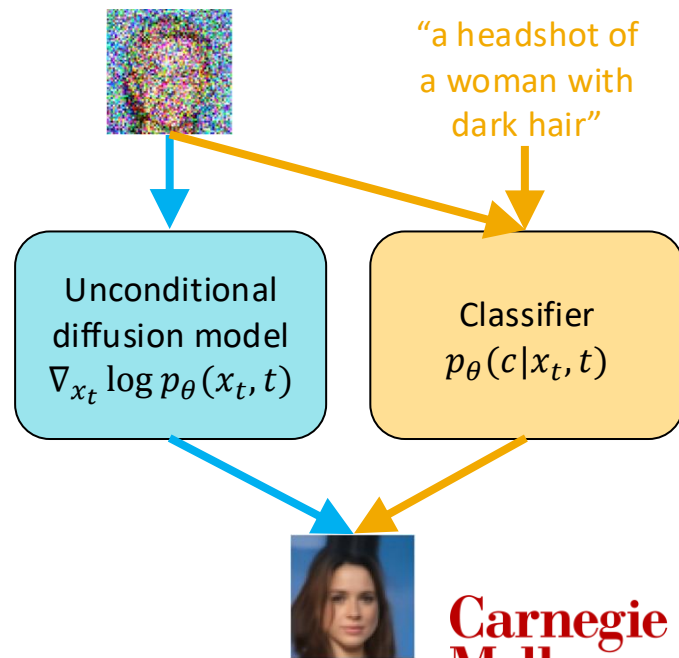
Unconditional diffusion:



Conditional diffusion:



Classifier guidance diffusion:



Classifier guidance diffusion/Conditional score SDE

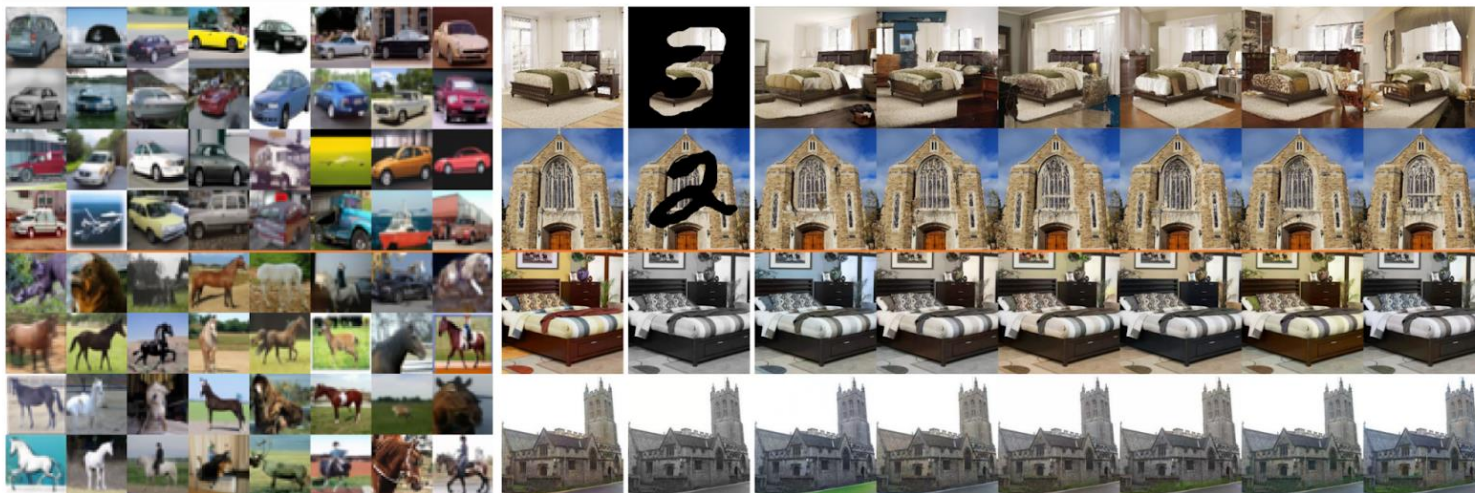


Figure 4: *Left*: Class-conditional samples on 32×32 CIFAR-10. Top four rows are automobiles and bottom four rows are horses. *Right*: Inpainting (top two rows) and colorization (bottom two rows) results on 256×256 LSUN. First column is the original image, second column is the masked/gray-scale image, remaining columns are sampled image completions or colorizations.

Remember the No.1 rule in probability theory

We know that diffusion models model the score function

$$\nabla_{x_t} \log p_\theta(x_t, t)$$

By Bayes' Theorem, we know that

$$p(x|c) = \frac{p(x,c)}{p(c)} = \frac{p(c|x)p(x)}{p(c)} \propto p(c|x)p(x)$$

$$\nabla_{x_t} \log p_\theta(x_t, t|c) = \nabla_{x_t} \log p_\theta(c|x_t, t) + \nabla_{x_t} \log p_\theta(x_t, t)$$

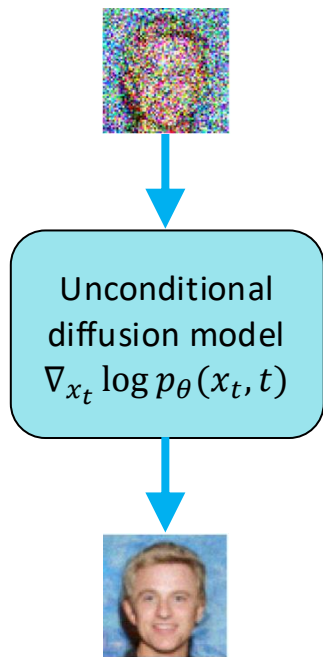
Pre-trained
unconditional
diffusion model

A discriminative model!
(e.g. a classifier, a CLIP model, etc)

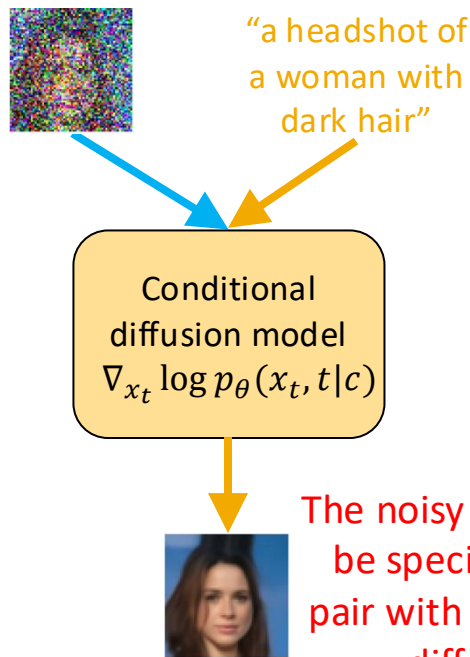
The classifier needs
to be able to process
noisy images!

We need to train a noisy classifier for classifier guidance diffusion!

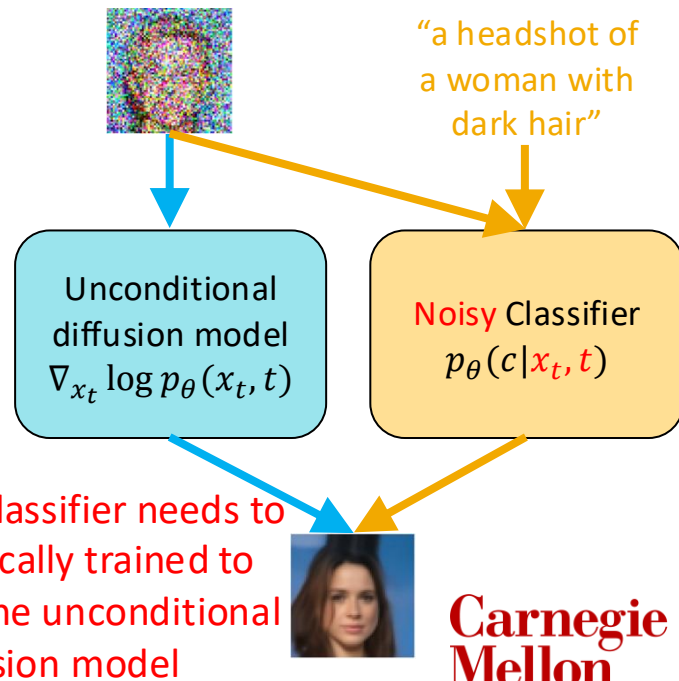
Unconditional diffusion:



Conditional diffusion:

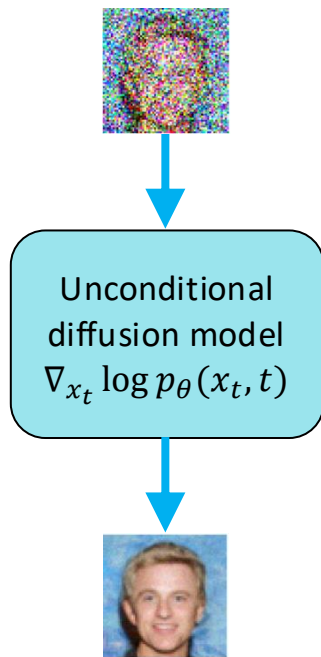


Classifier guidance diffusion:

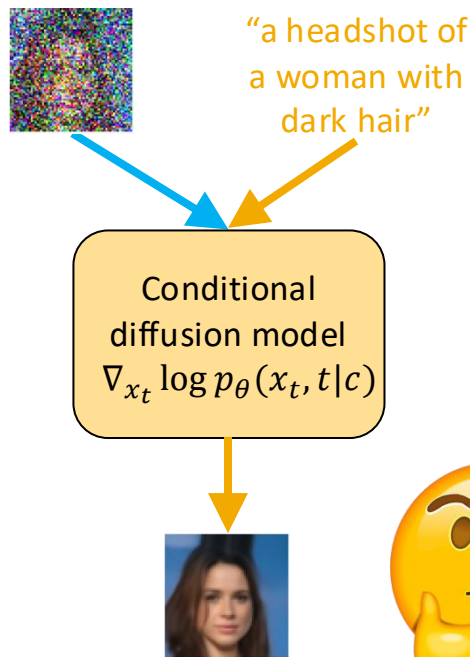


Is it possible to use an off-the-shelf classifier?

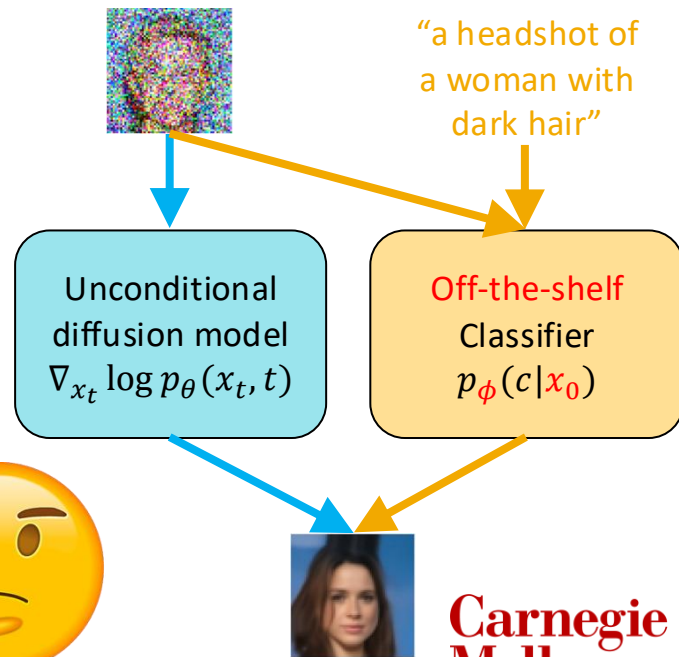
Unconditional diffusion:



Conditional diffusion:

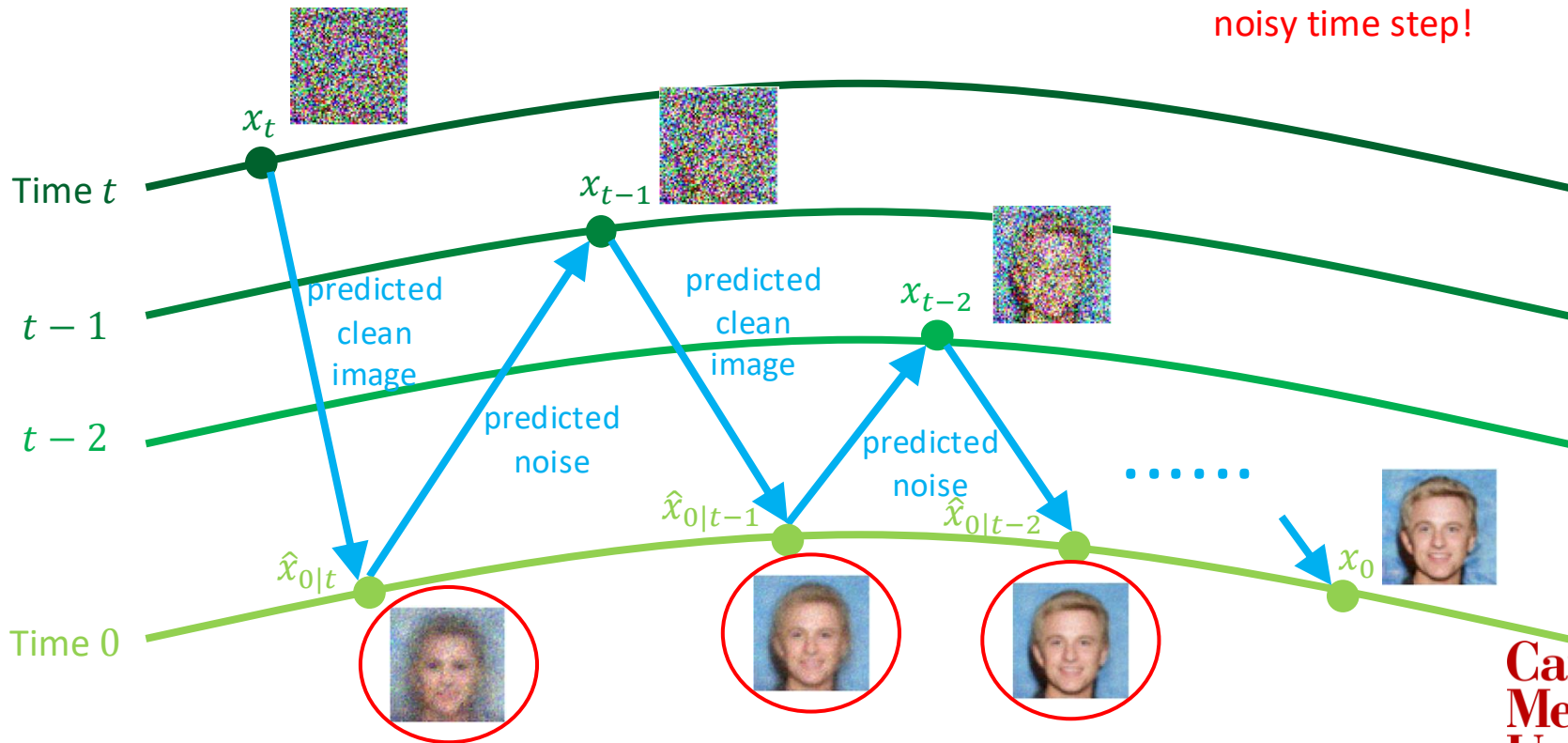


Classifier guidance diffusion:

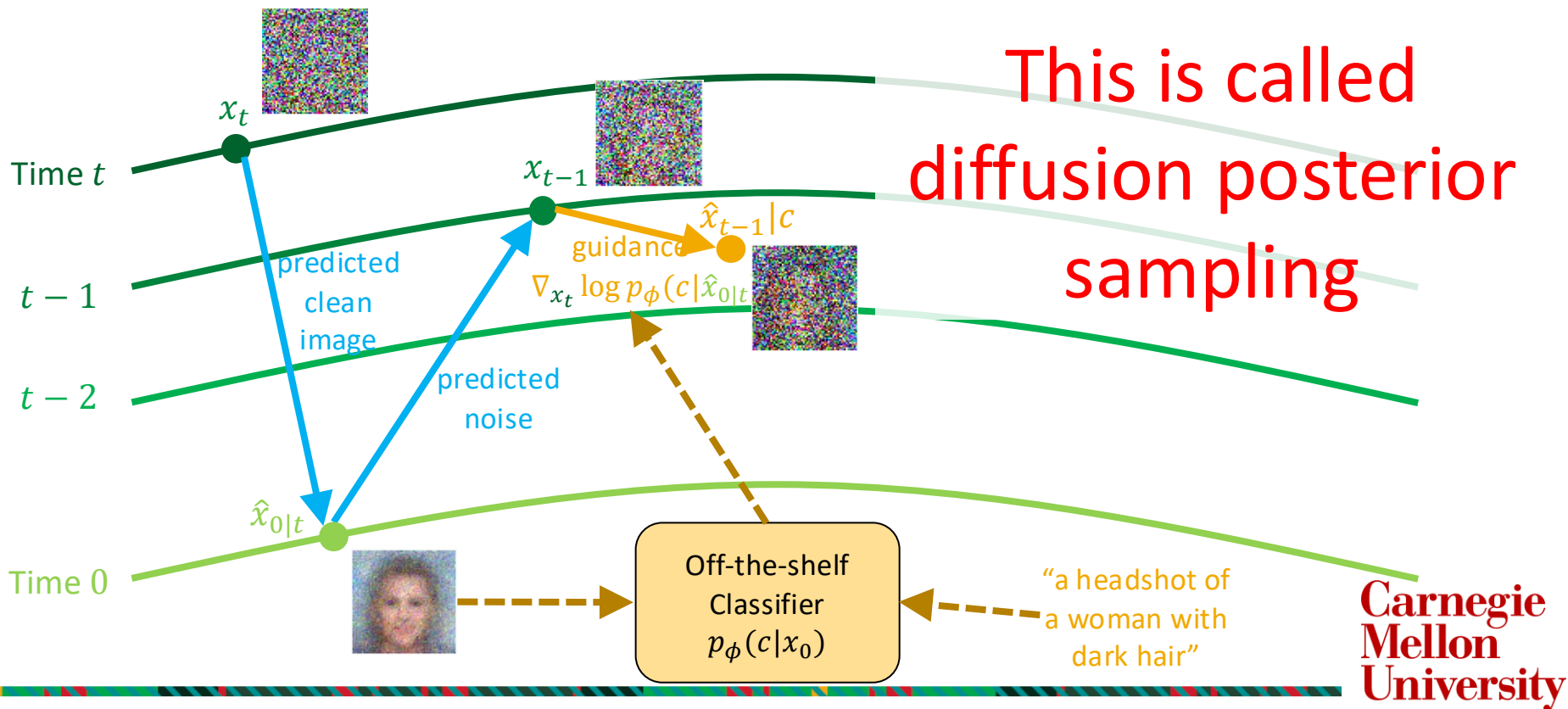


Remember how we do DDIM

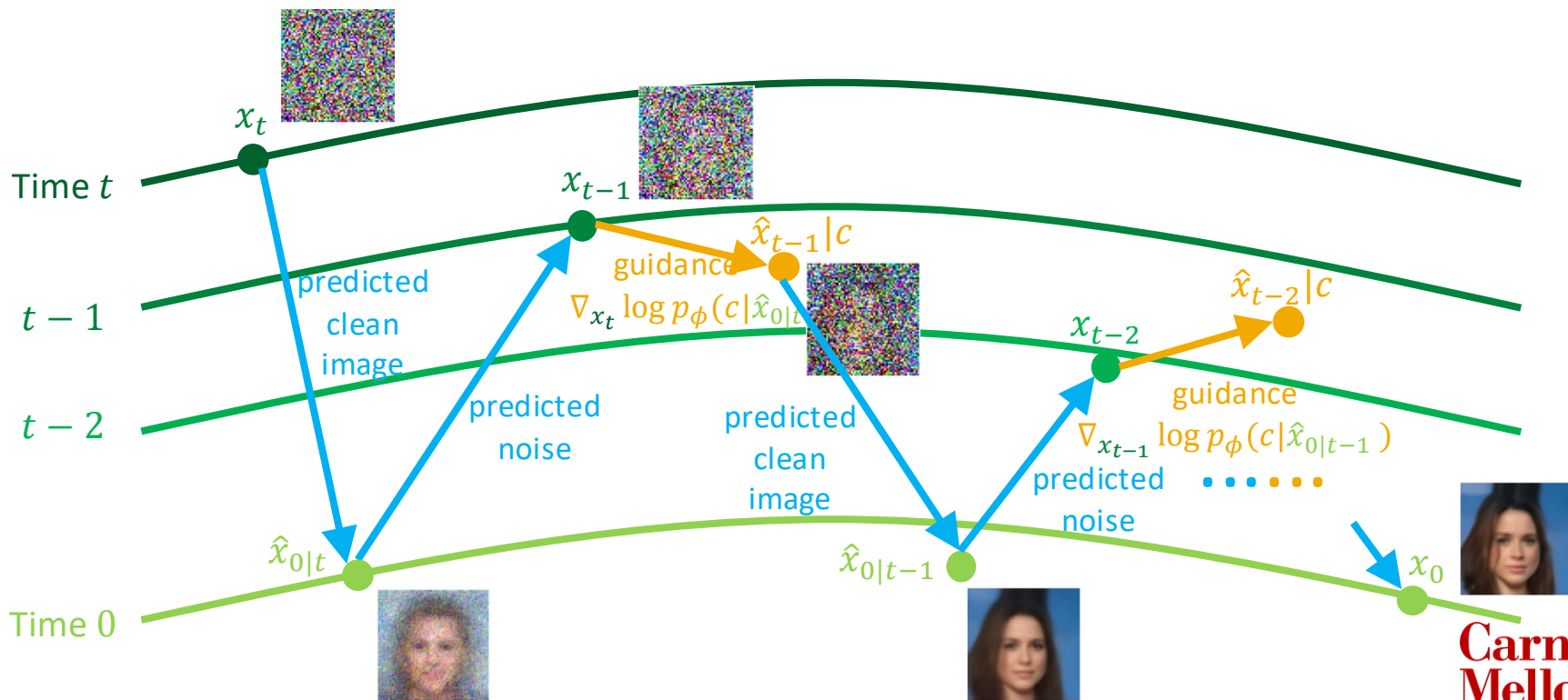
We do have some clean image estimates from each noisy time step!



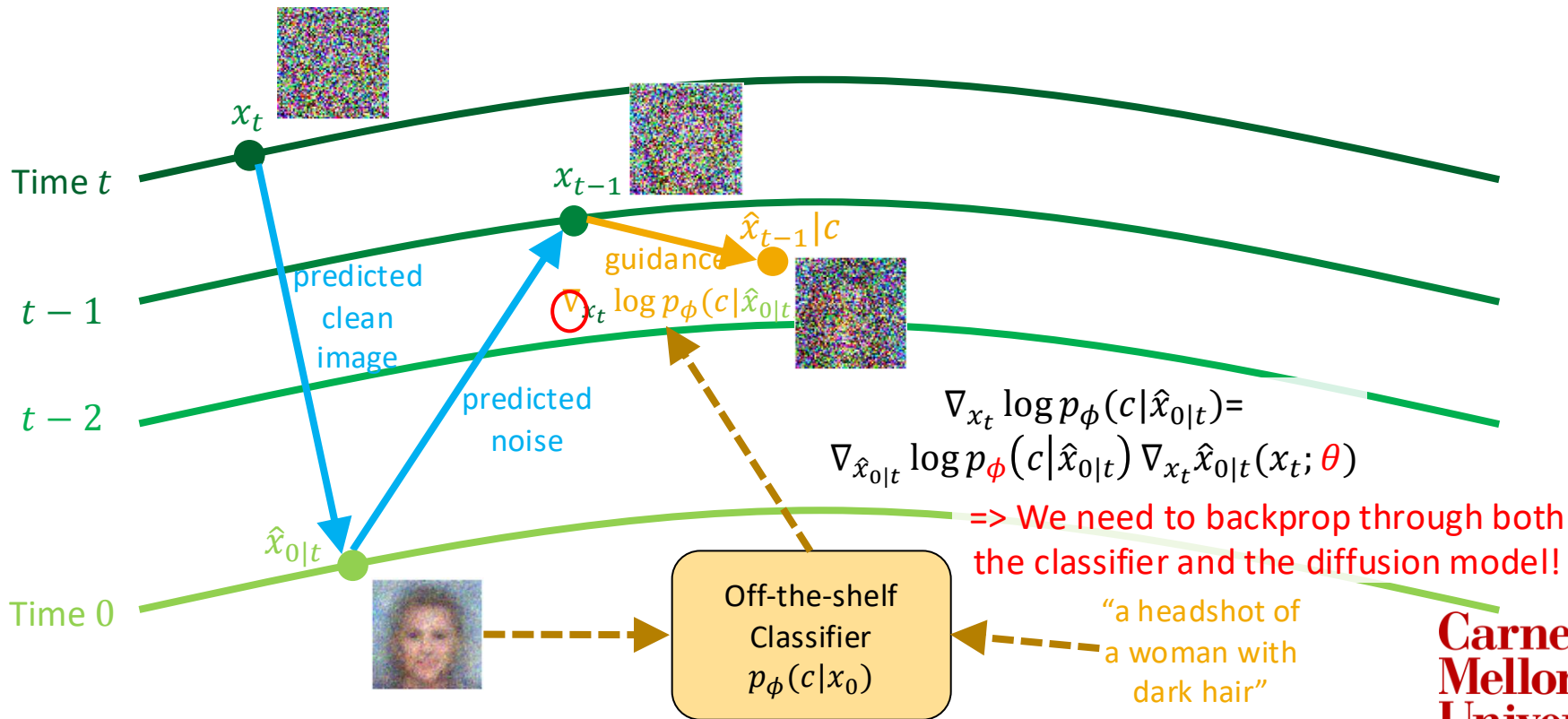
Using off-the-shelf classifiers for diffusion guidance



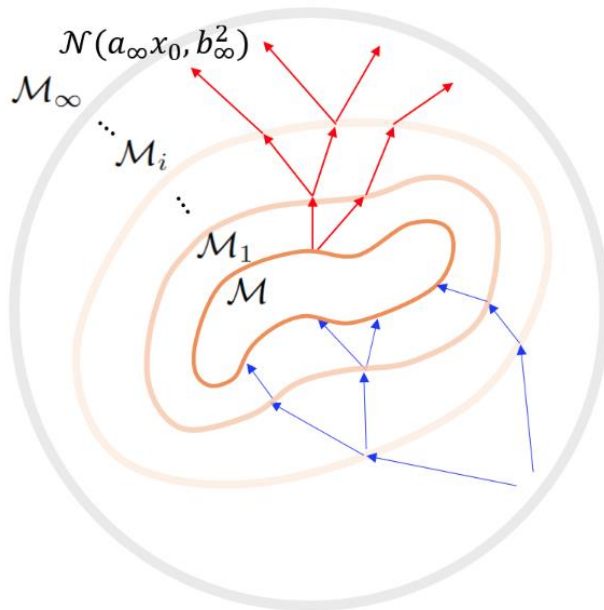
Diffusion posterior sampling (DPS)



DPS is great, but



Diffusion manifolds are layered bubble shells



(a) Geometry of diffusion model

Why?

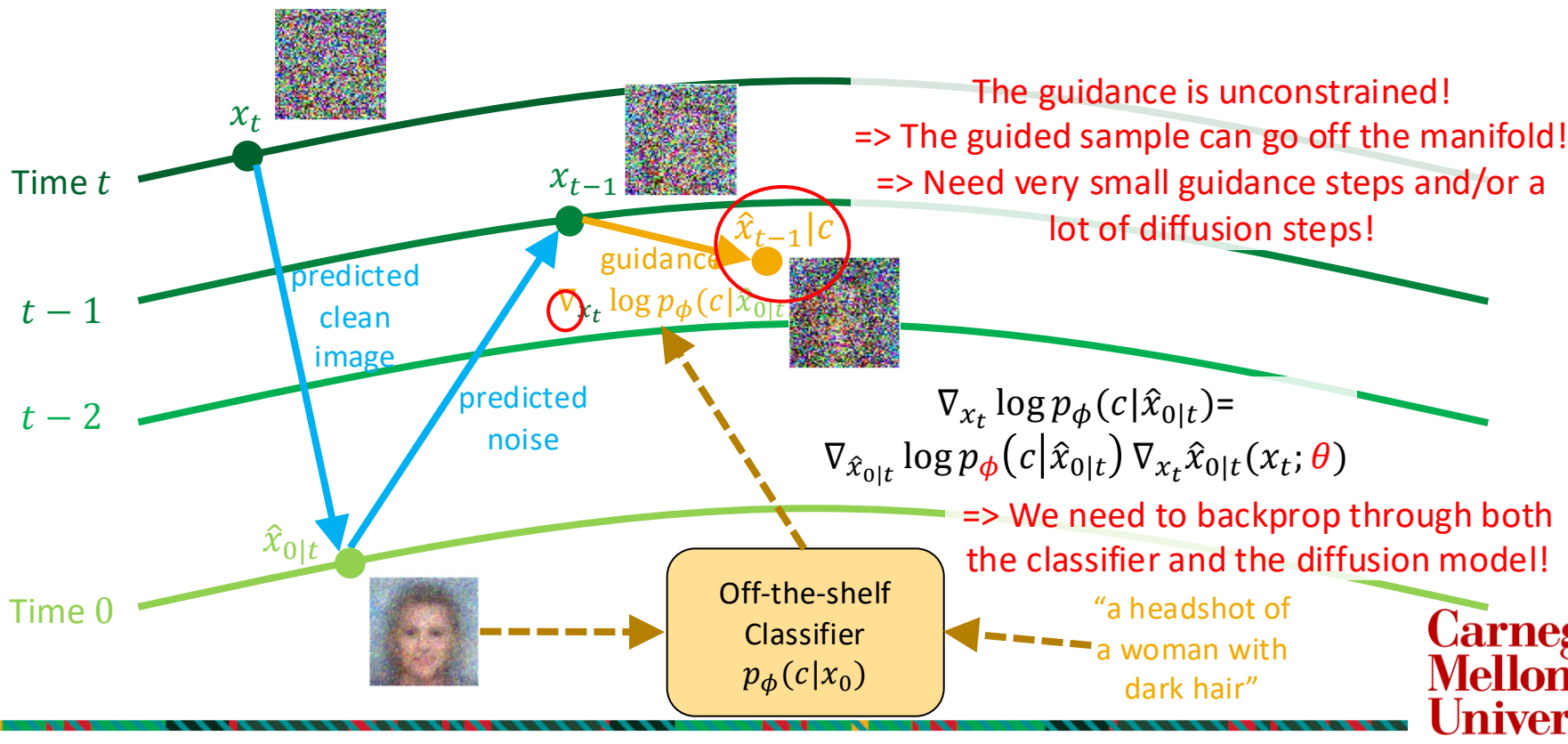
Because high dimensional Gaussians are soap bubbles!

Proof hint:

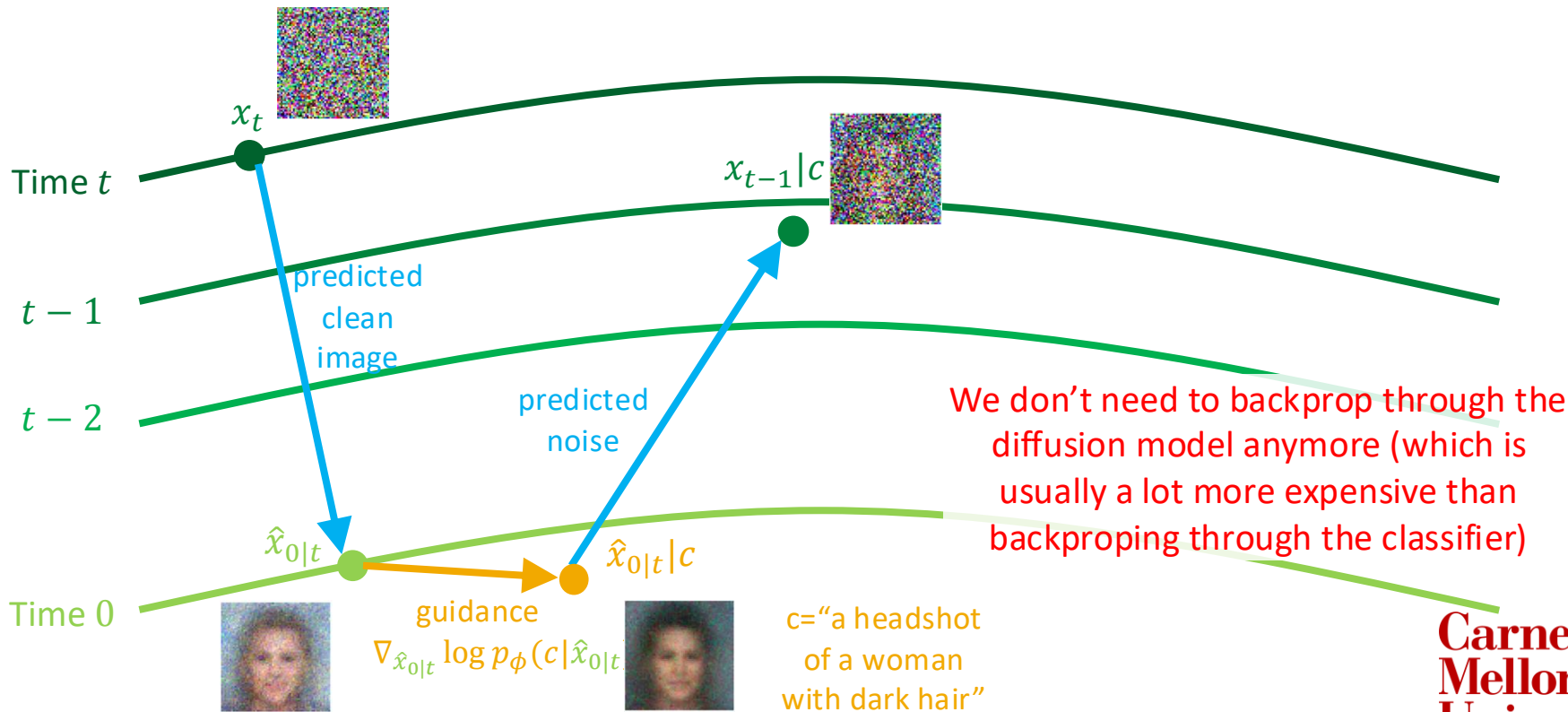
Try to see where is the majority of the density concentrated at in high dimensional Gaussians

DPS is great, but

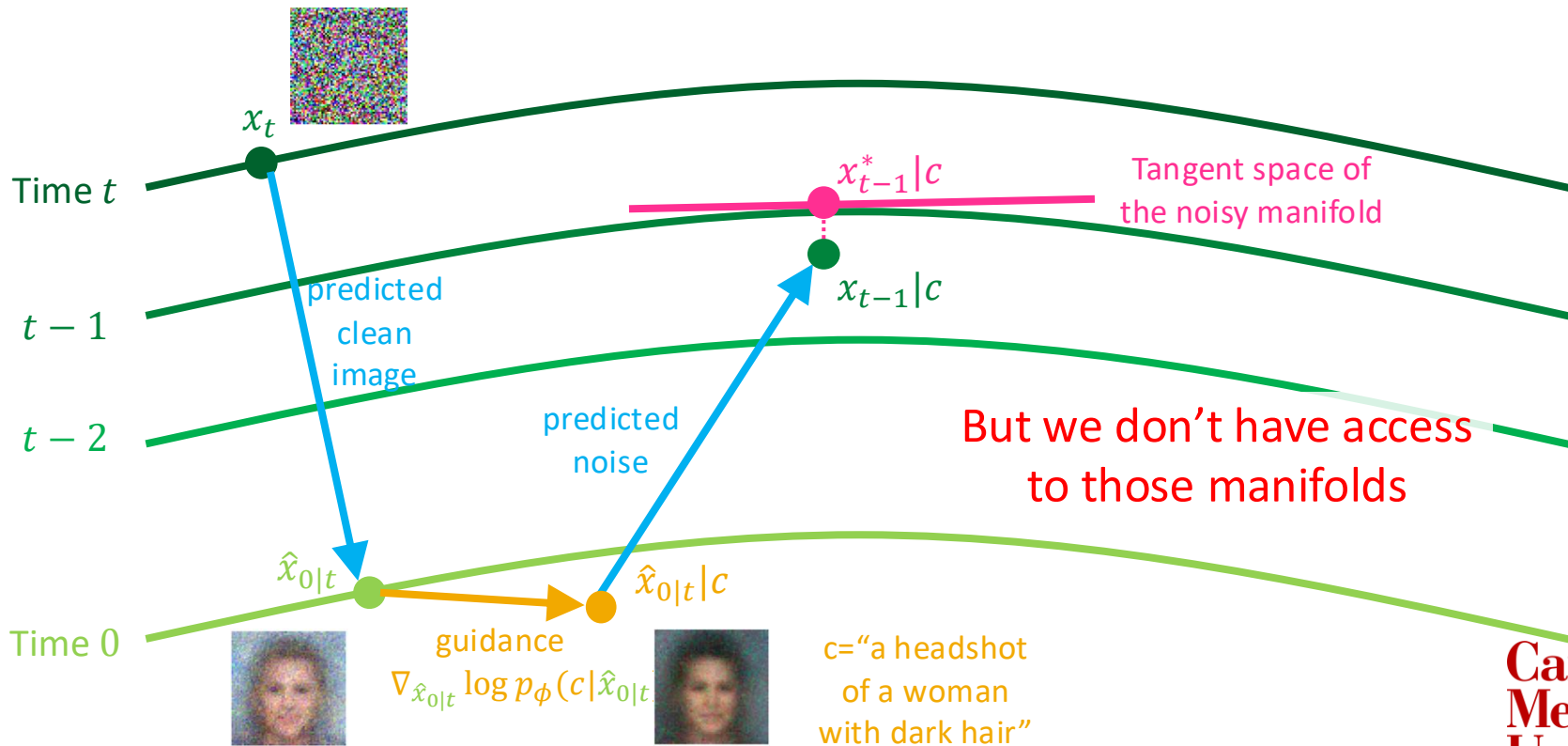
DPS is very very slow!



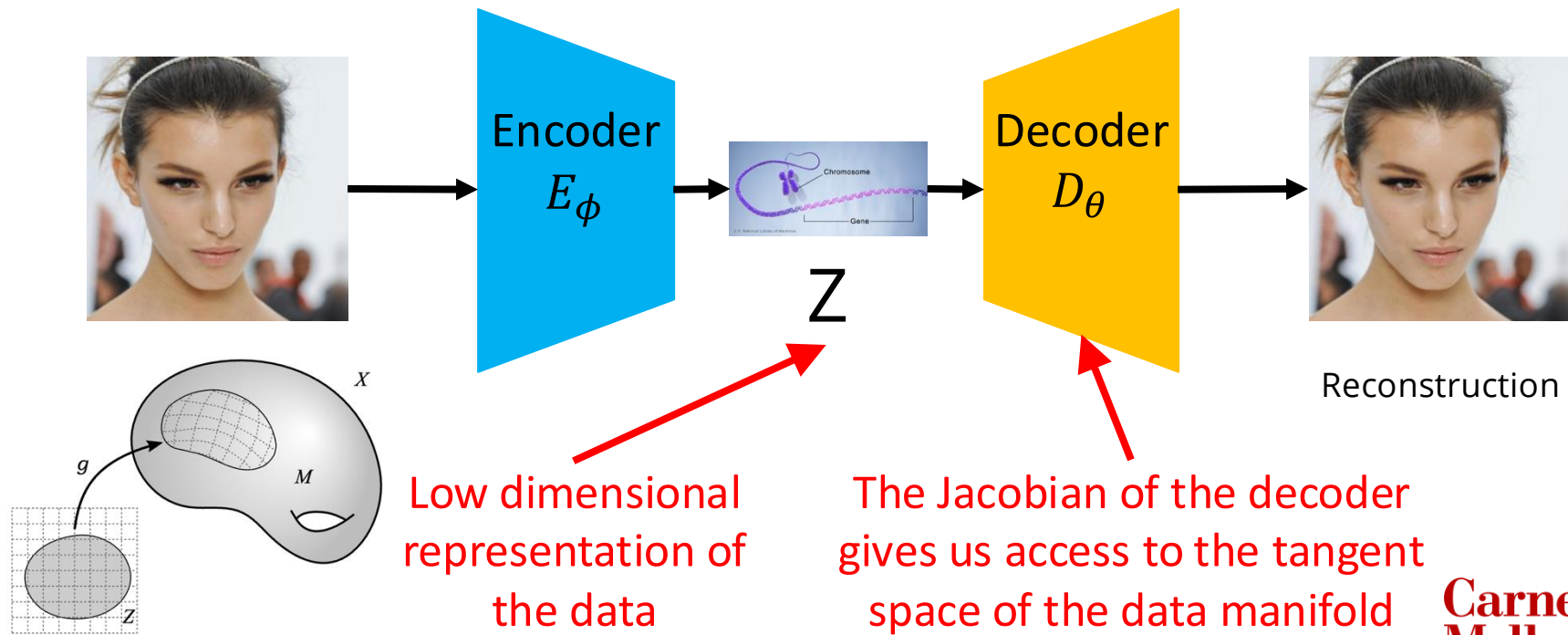
Remedy 1: Take gradient w.r.t. the clean image instead



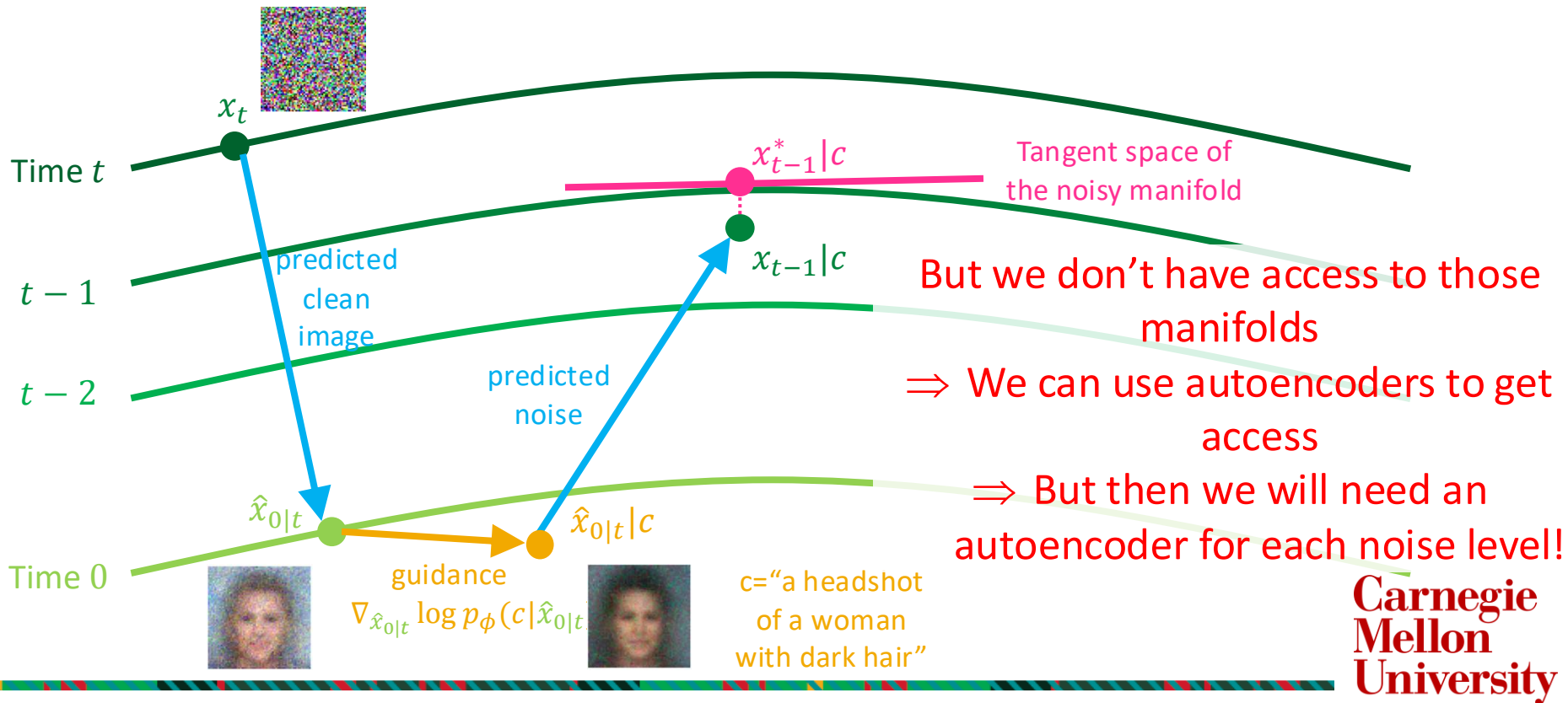
Remedy 2: Project the guided samples to the manifold



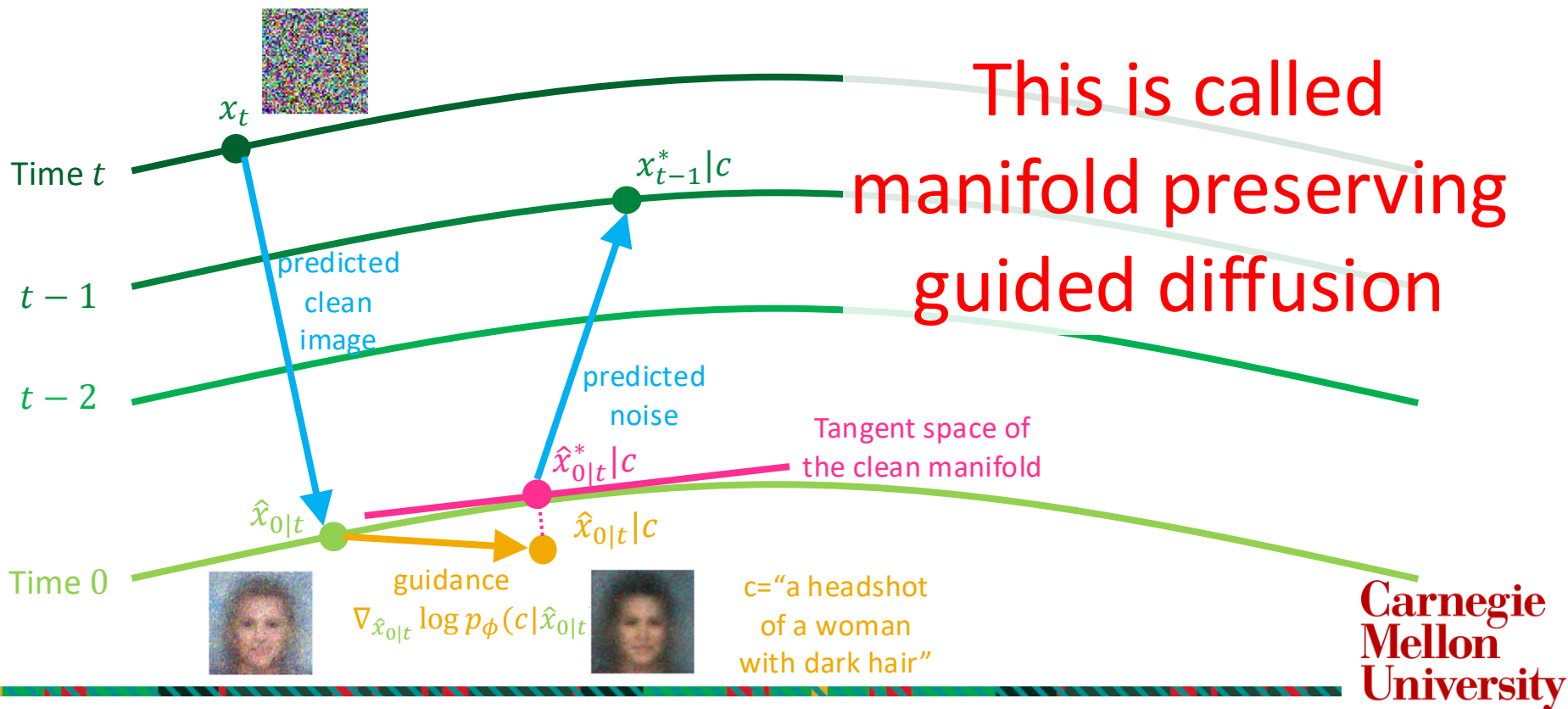
Getting access to data manifolds via the decoder of an autoencoder



Remedy 2: Project the guided samples to the manifold



Remedy 2.5: Project the clean data estimation to the clean data manifold



Manifold preserving guided diffusion (MPGD) is much faster (and better) than DPS

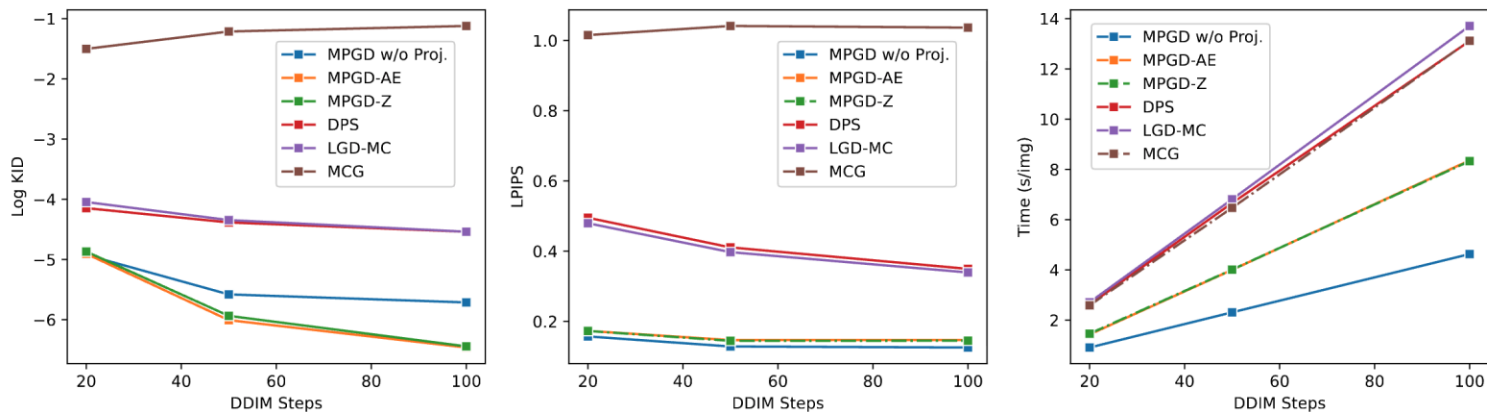
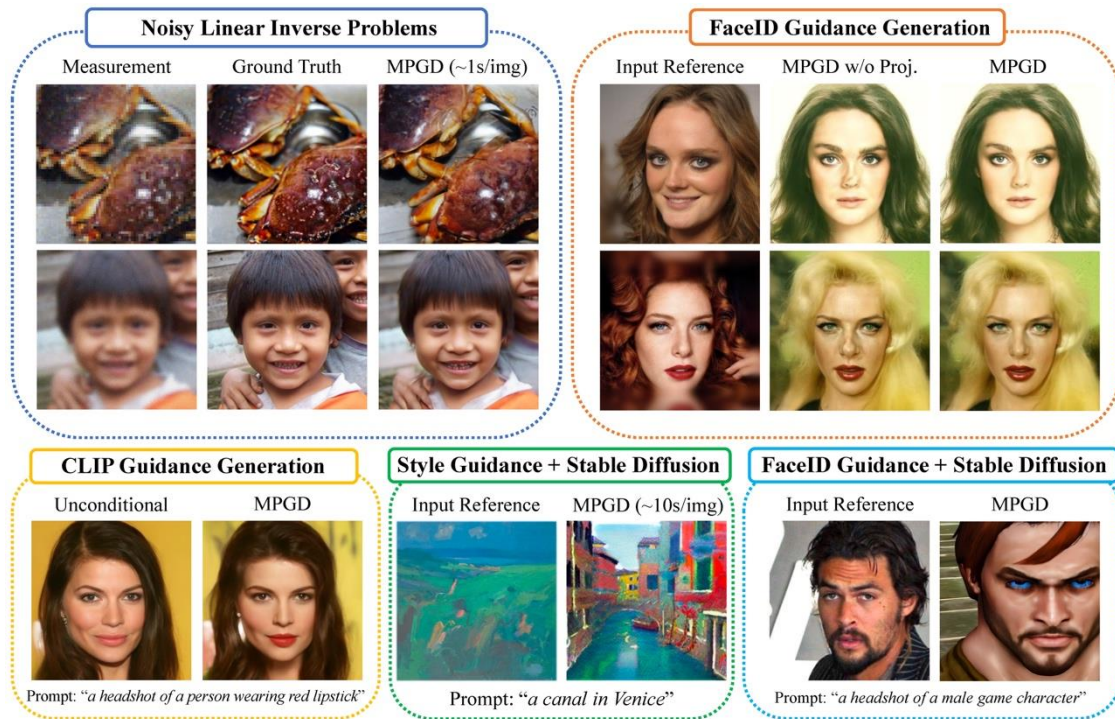


Figure 5: Quantitative results of FFHQ super-resolution experiment that compares fidelity (log KID), guidance quality (LPIPS) and inference time across different numbers of DDIM steps.

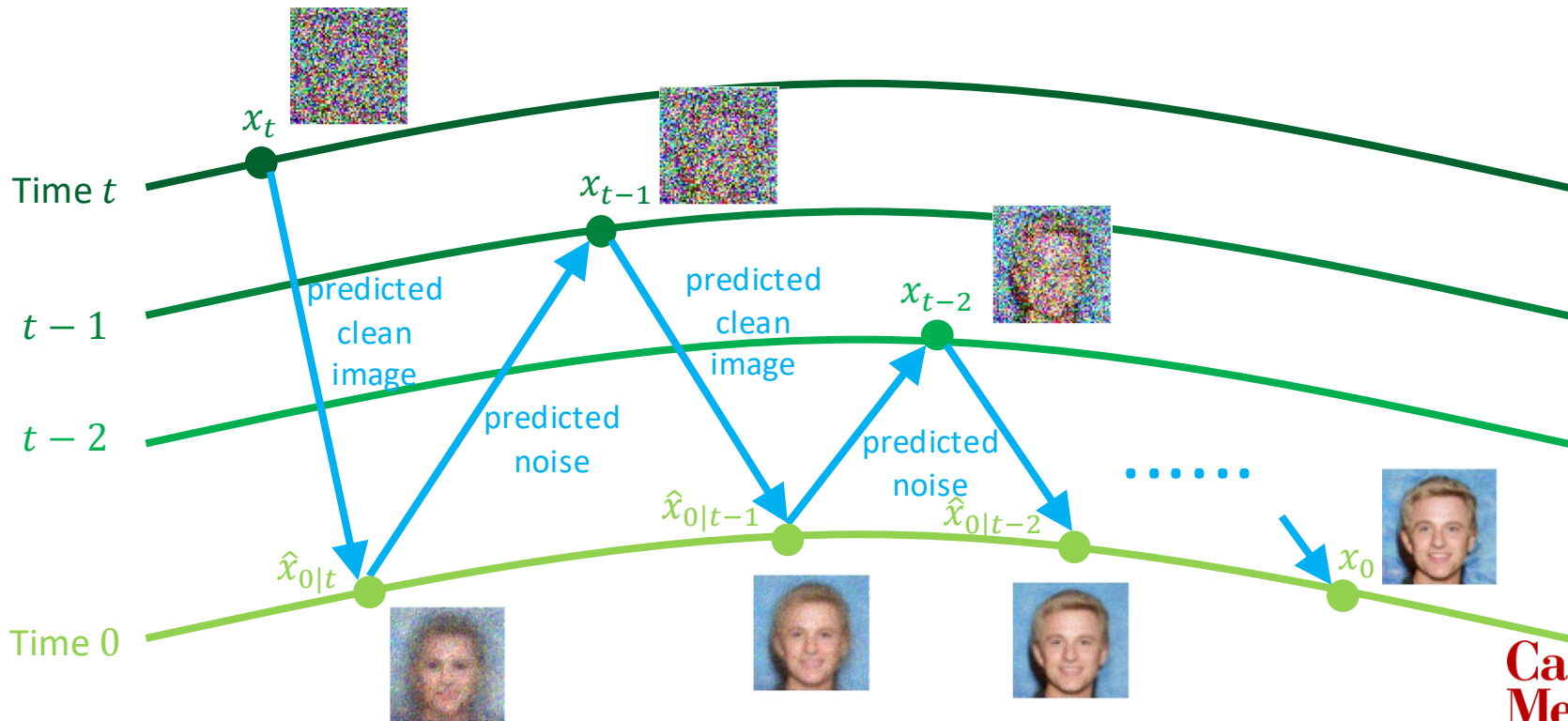
MPGD is training-free and applicable to all differentiable guidance/classifier/reward



Any other training-free ways to inject conditions into diffusion models?

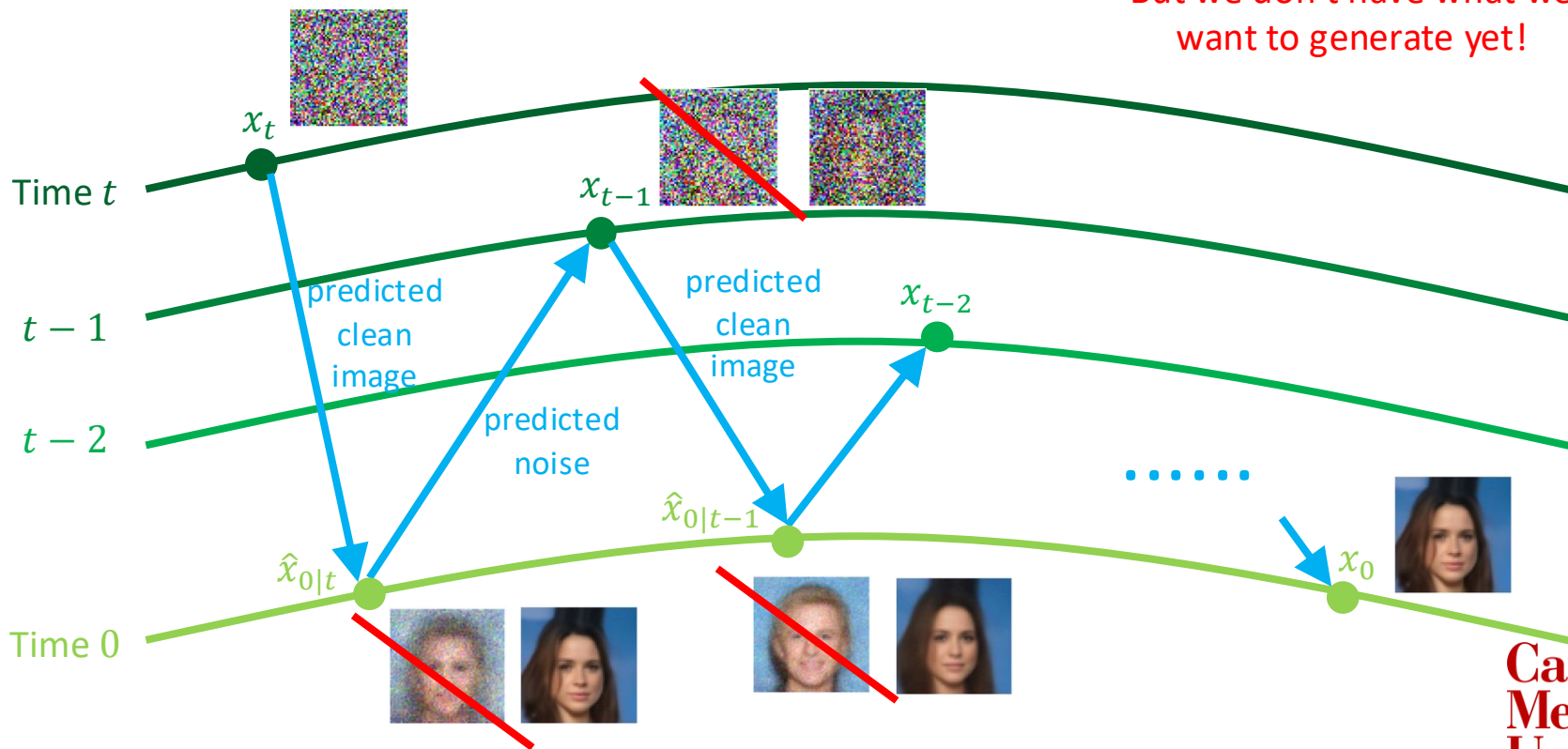


Remember how we do DDIM

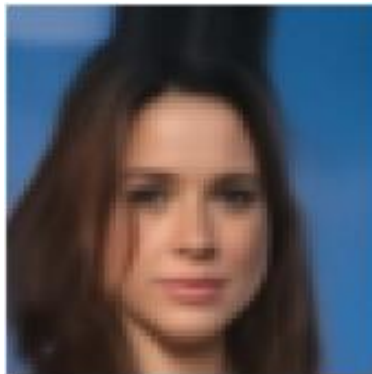


If we already know what we want to generate

But we don't have what we want to generate yet!



But do we really need to know EXACTLY what we need to generate?



We only need a rough draft/preliminary version of what we want to generate!



But we also don't have a draft!

We only need a rough draft/preliminary version of what we want to generate!



Or do we?

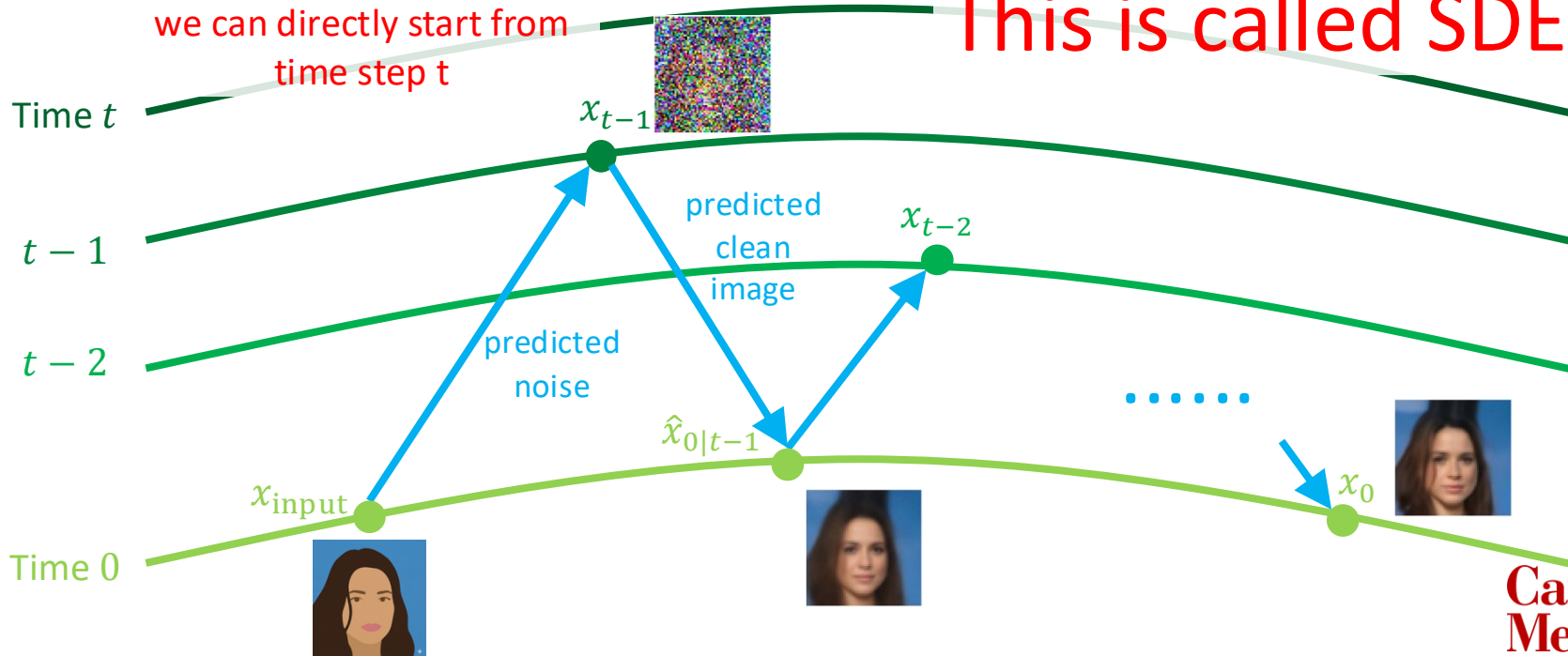
How about just draw one lol



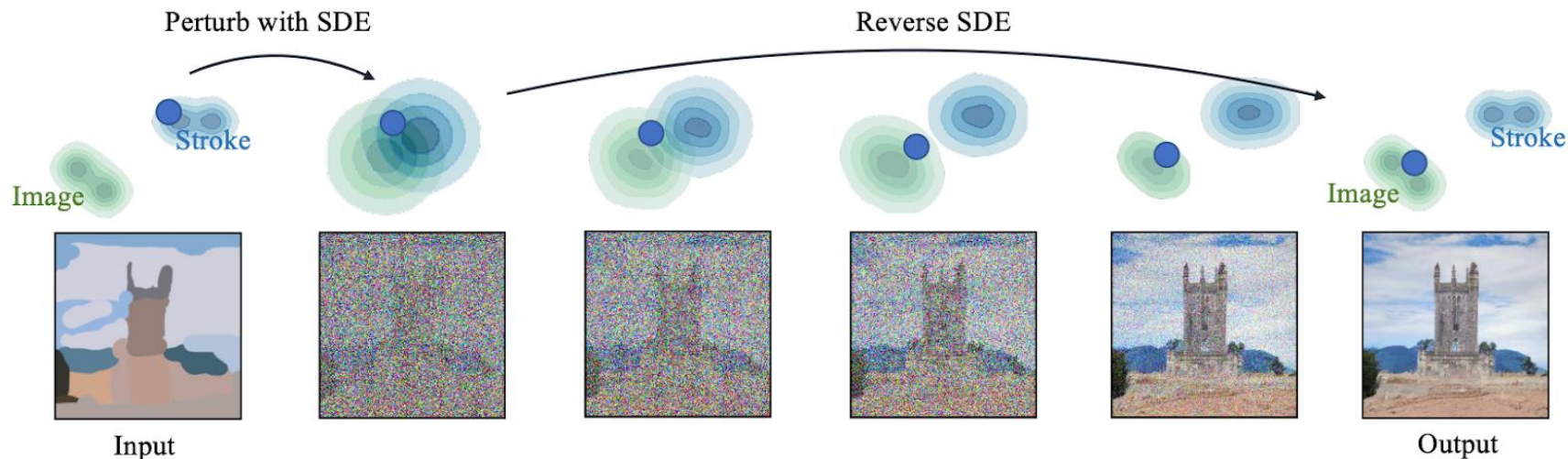
With this rough painting, we can do

Notice that we don't need to
do any steps before t now,
we can directly start from
time step t

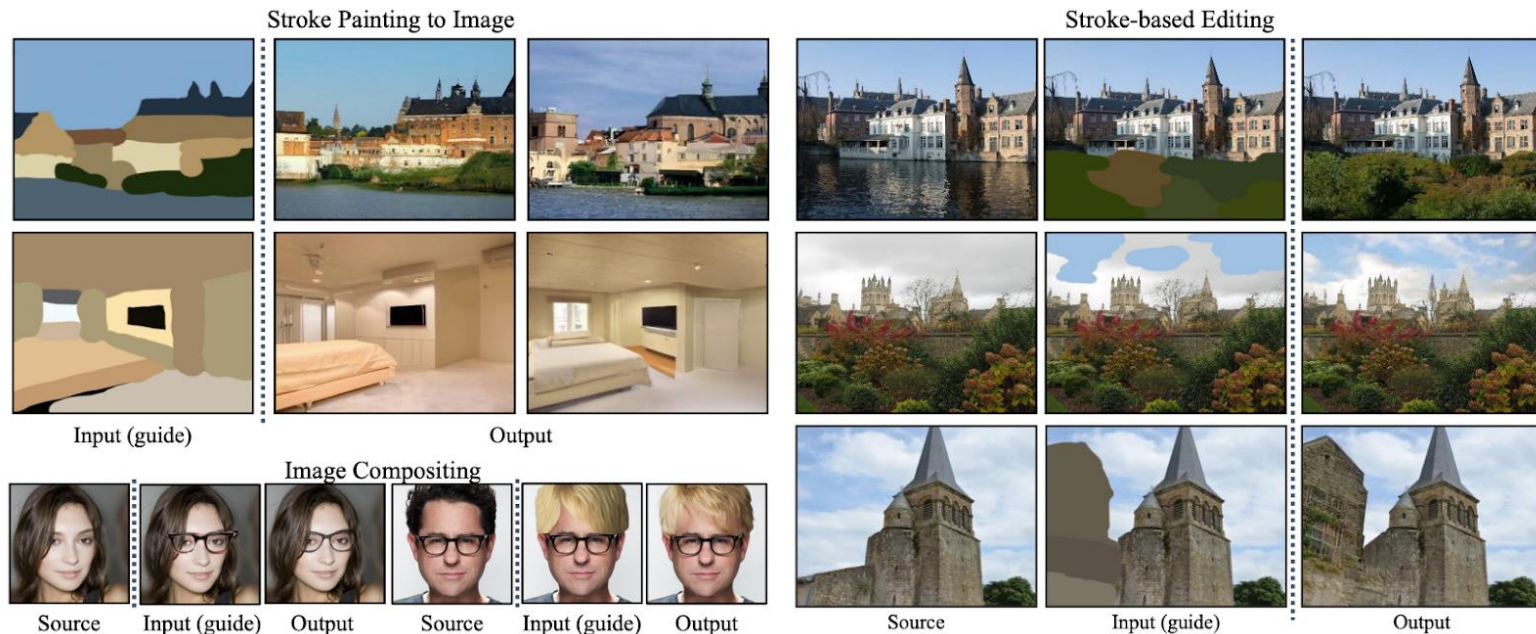
This is called SDEdit



SDEdit hijacks the diffusion process with user inputs



SDEdit can turn your preliminary paintings/edits into realistic looking images



SDEdit can turn your preliminary paintings/edits into realistic looking images

Original Image



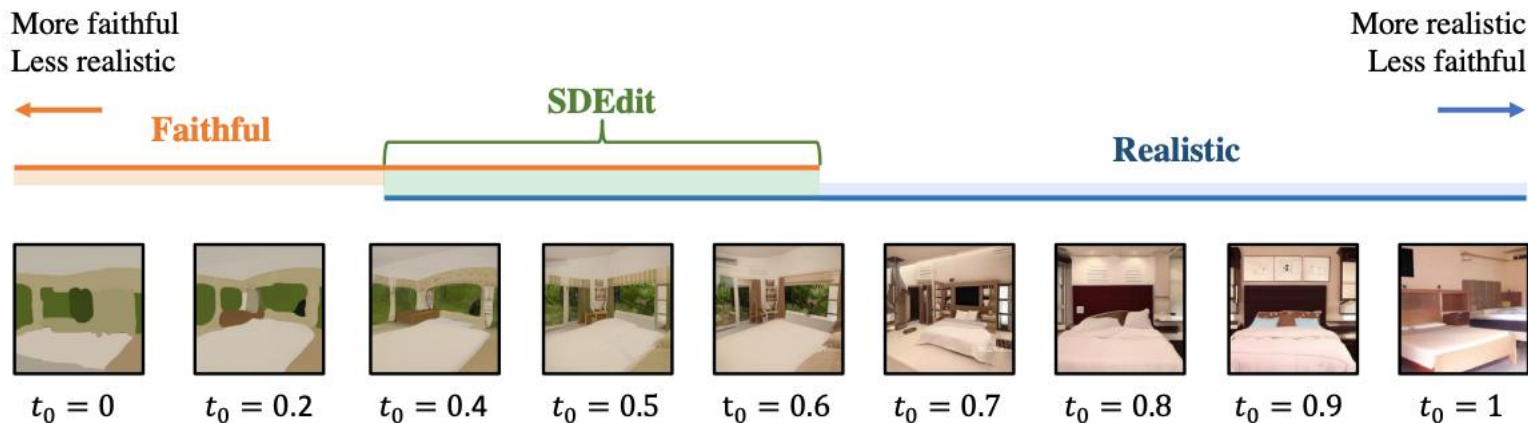
User Input



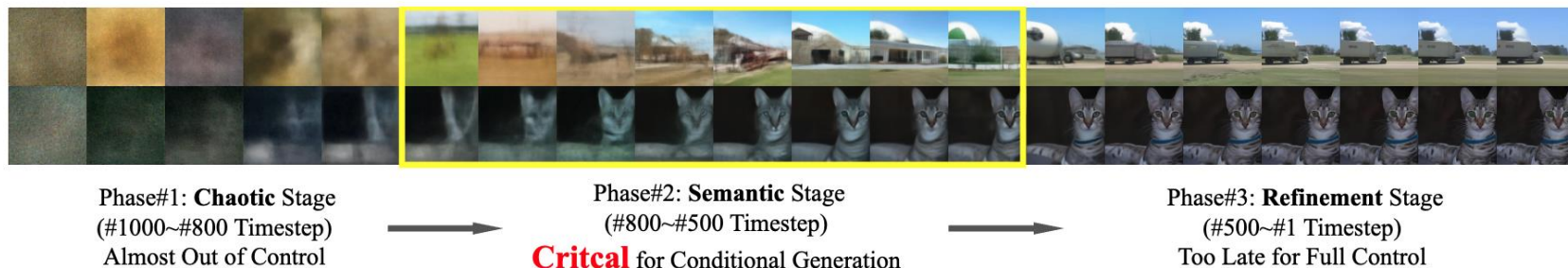
SDEdit Output



The controllability-fidelity tradeoff

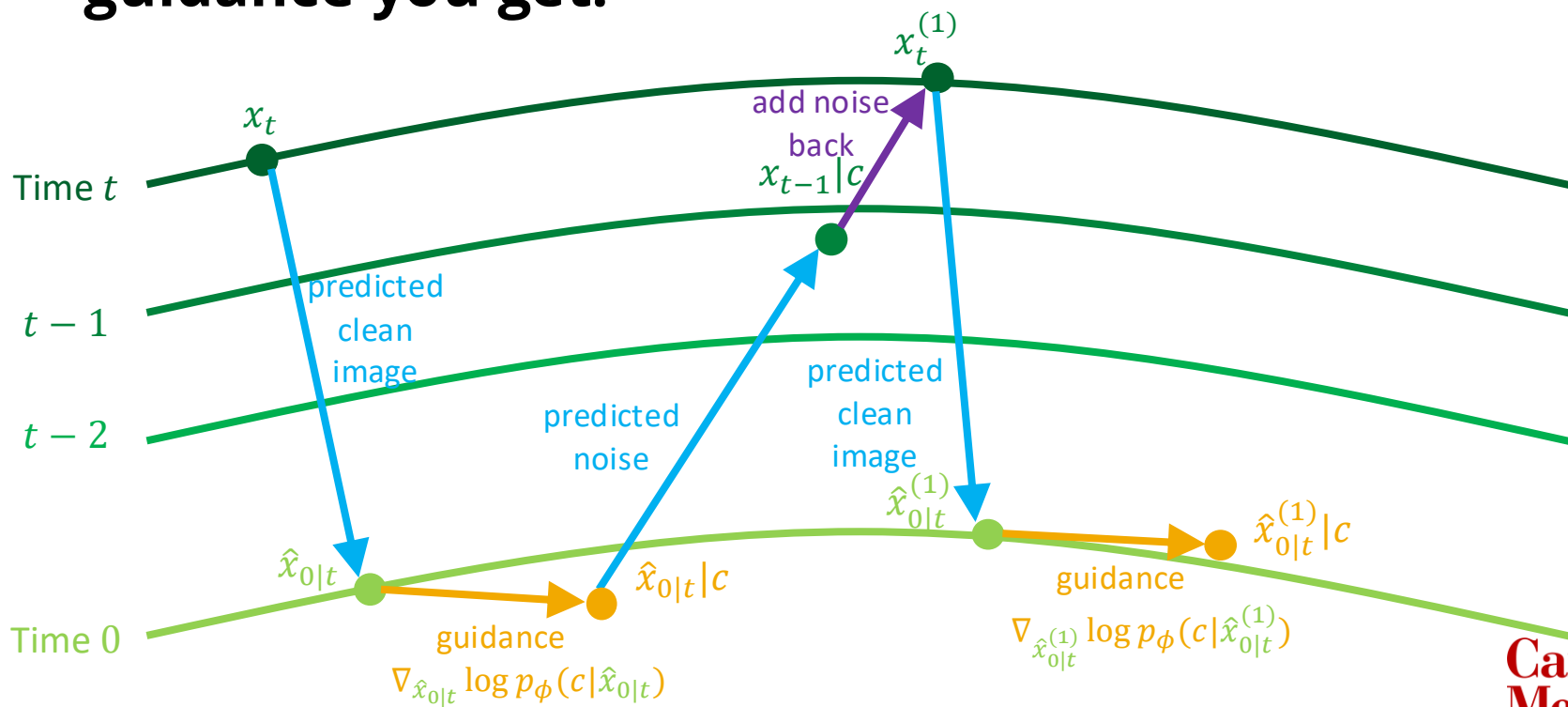


Guidance pro tips 1: You should guide more during the middle time steps (wisdom from EDM still holds)

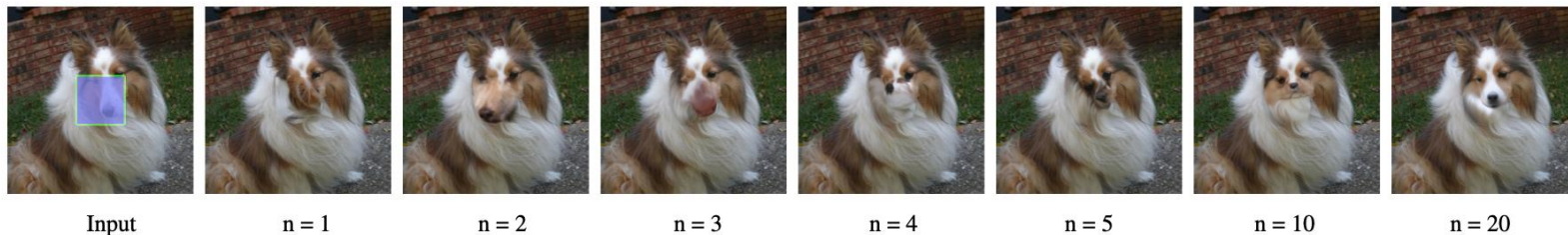


Guidance pro tips 2: You can even time travel back to the previous time step if you want to refine the guidance you get!

43

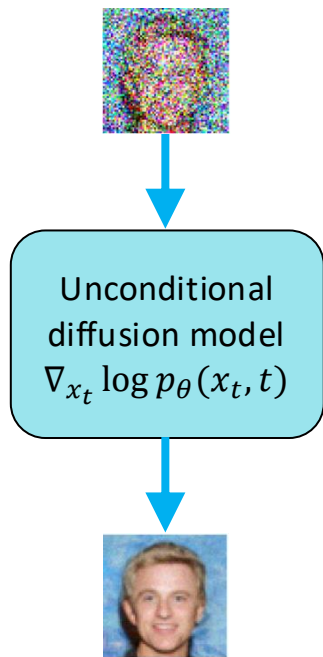


Guidance pro tips 2: You can even time travel back to the previous time step if you want to refine the guidance you get!

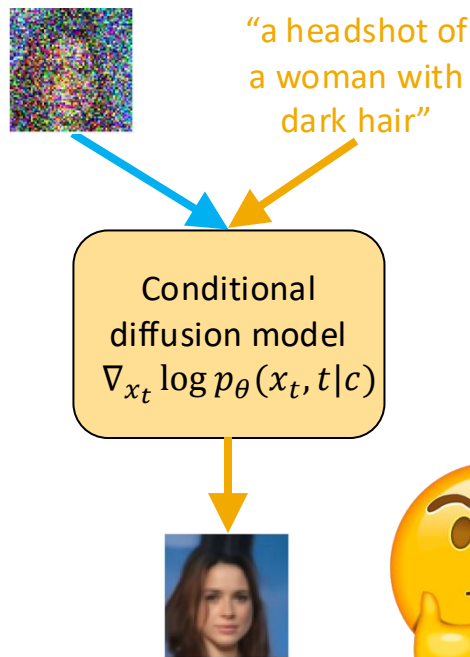


What if we are normies

Unconditional diffusion:



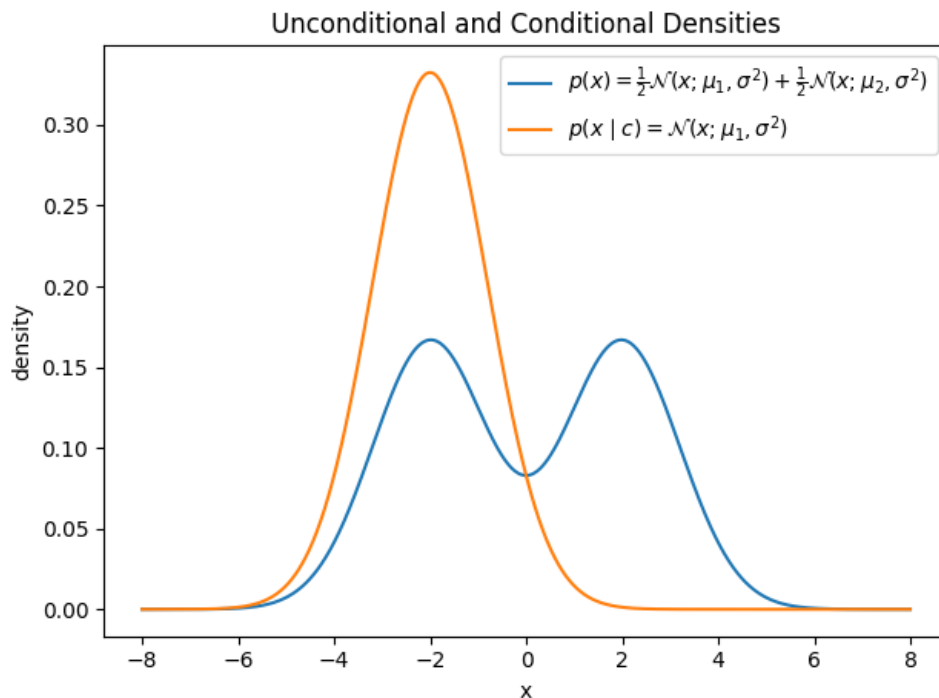
Conditional diffusion:



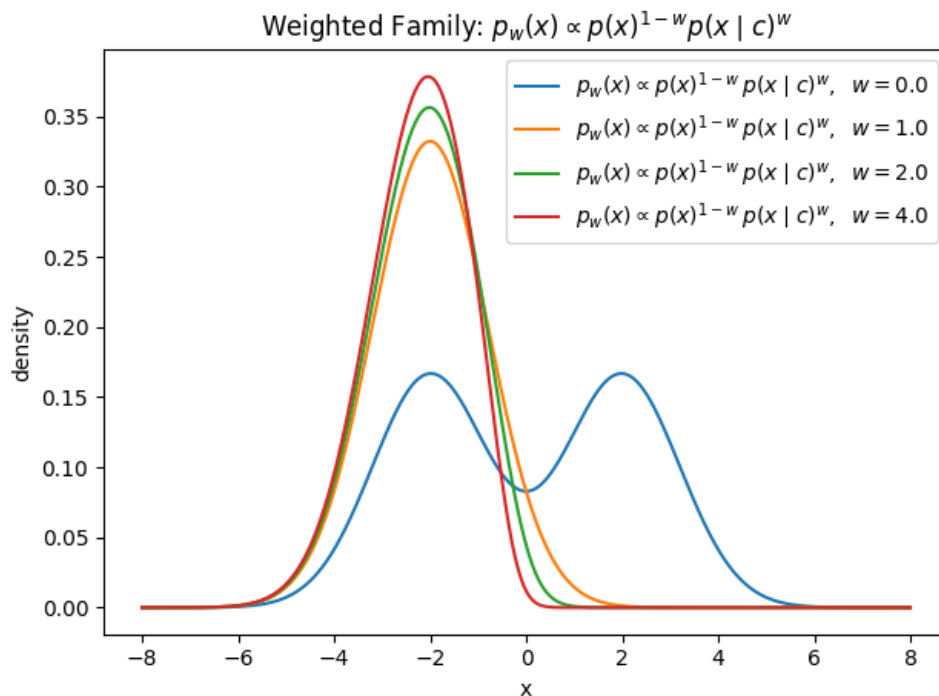
If we have already trained an unconditional diffusion model and a separate conditional diffusion model, can we somehow **use both** to make our conditional generation better?



Unconditional distribution v.s. Conditional distribution



Conditional mixed with unconditional distributions



Classifier-free guidance: Sample from the weighted mixture of unconditional + conditional score

$$\nabla_x \log p_\gamma(x | y) = (1 - \gamma) \nabla_x \log p(x) + \gamma \nabla_x \log p(x | y).$$

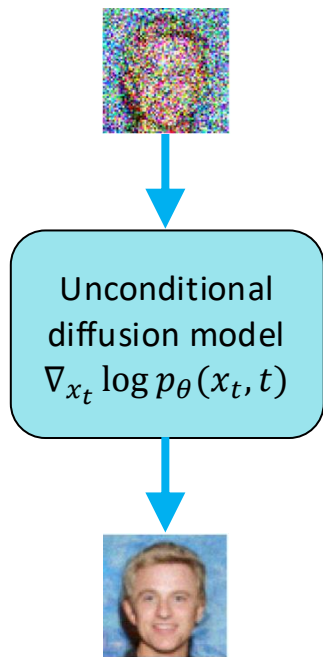
Algorithm 1 Joint training a diffusion model with classifier-free guidance

Require: p_{uncond} : probability of unconditional training

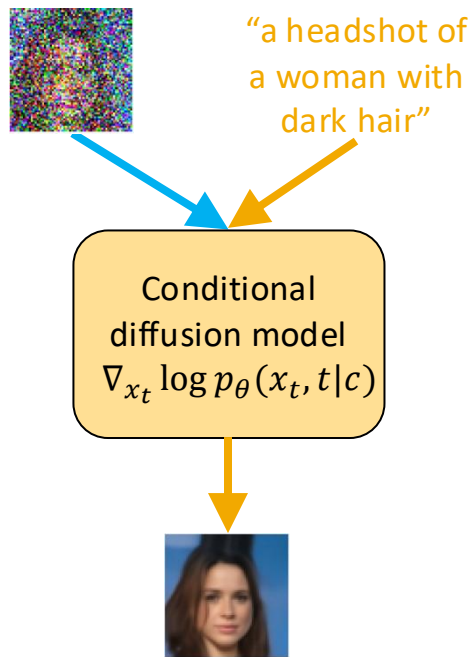
- 1: **repeat**
 - 2: $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$ ▷ Sample data with conditioning from the dataset
 - 3: $\mathbf{c} \leftarrow \emptyset$ with probability p_{uncond} ▷ Randomly discard conditioning to train unconditionally
 - 4: $\lambda \sim p(\lambda)$ ▷ Sample log SNR value
 - 5: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 6: $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \epsilon$ ▷ Corrupt data to the sampled log SNR value
 - 7: Take gradient step on $\nabla_\theta \|\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \epsilon\|^2$ ▷ Optimization of denoising model
 - 8: **until** converged
-

Do we even need to train two separate models?

Unconditional diffusion:

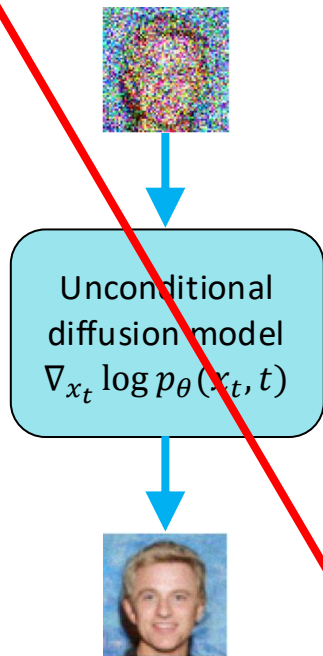


Conditional diffusion:

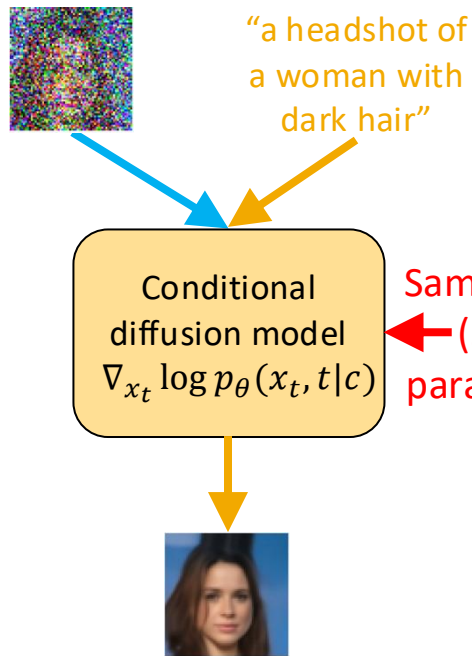


Do we even need to train two separate models?

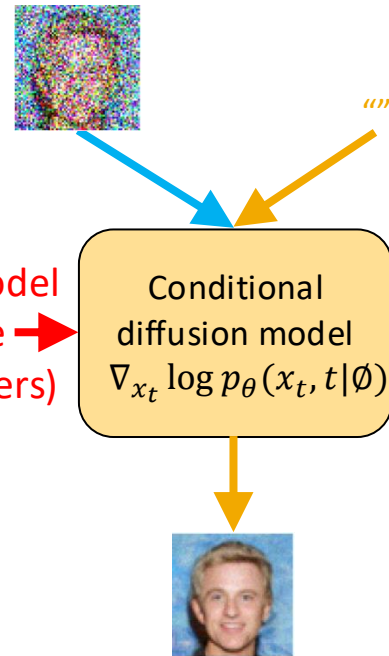
Unconditional diffusion:



Conditional diffusion:



Pseudo unconditional diffusion:



Classifier-free guidance: Sample from the weighted mixture of unconditional + conditional score

$$\nabla_x \log p_\gamma(x | y) = (1 - \gamma) \nabla_x \log p(x) + \gamma \nabla_x \log p(x | y).$$



Two sets of samples from OpenAI's GLIDE model, for the prompt 'A stained glass window of a panda eating bamboo.', taken from their paper. Guidance scale 1 (no guidance) on the left, guidance scale 3 on the right.

Now you know how to turn an unconditional diffusion into a conditional one!

Training-based methods:

- Classifier guidance diffusion
- Classifier free guidance (if you don't already have a conditional model)

Training-free methods:

- DPS
- MPGD
- SDEdit
- Classifier free guidance (if you already have a conditional model)

In the next two weeks, we are going into the realm of the state-of-the-arts

- **2/3 (Tue):** How to use diffusion/flow models for robotics, control & decision making (**Max Simchowitz**, MLD Prof.)
- **2/5 (Thur):** Text-to-image models and SOTA techniques
- **2/10 (Tue):** How to sample from diffusion/flow models with a single step
 - Distillation
 - Consistency models
 - Flow maps
- **2/12 (Thur):** Real-time generation techniques for video (**Linqi (Alex) Zhou**, Luma AI)

In the next two weeks, we are going into the realm of the state-of-the-arts

- **2/3 (Tue):** How to use diffusion/flow models for robotics, control & decision making (**Max Simchowitz**, MLD Prof.) (Will be chill, please come in person!)
- **2/5 (Thur):** Text-to-image models and SOTA techniques
- **2/10 (Tue):** How to sample from diffusion/flow models with a single step
 - Distillation
 - Consistency models
 - Flow maps
- **2/12 (Thur):** Real-time generation techniques for video (**Linqi (Alex) Zhou**, Luma AI)