# Compositional Conditioning Consistency Model

賴柏宏 Oral Defense
College of Artificial Intelligence,
National Yang Ming Chiao Tung University
2025  07.02

# Outline

# 1 Introduction

## Our Contributions

1. We proposed **CCCM --** the first consistency model capable of compositional zero-shot generation, effectively transferring CCDM's unseen image generation ability into 2–4 step

2. **Modified consistency distillation,** combining teacher-predicted and forward-process-formulated supervision. Three fusion strategies: Switch, Step Fuse, and Loss Fuse

3. CCCM achieves superior FID scores and maintains zero-shot accuracy despite requiring only a fraction of CCDM's sampling steps.
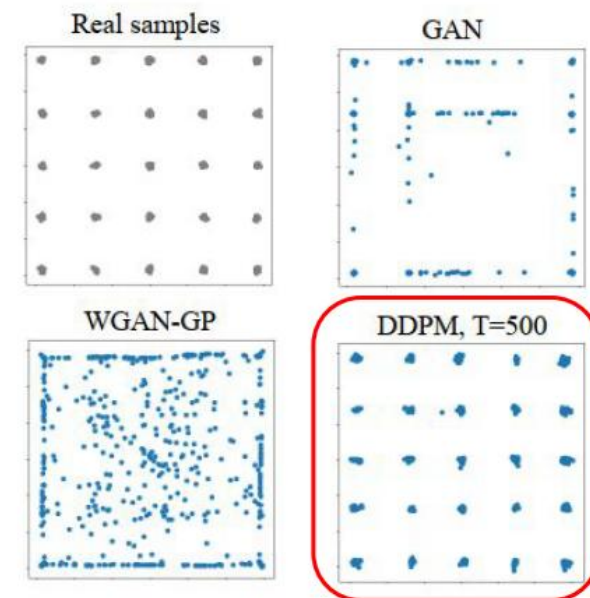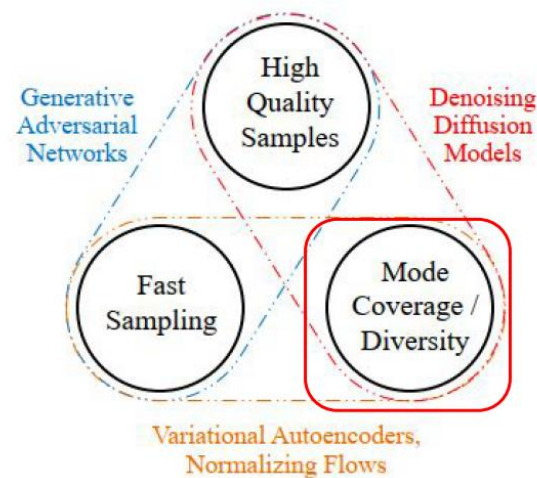
# **Backgrounds**

- Are there neural networks capable of generating unseen classes?



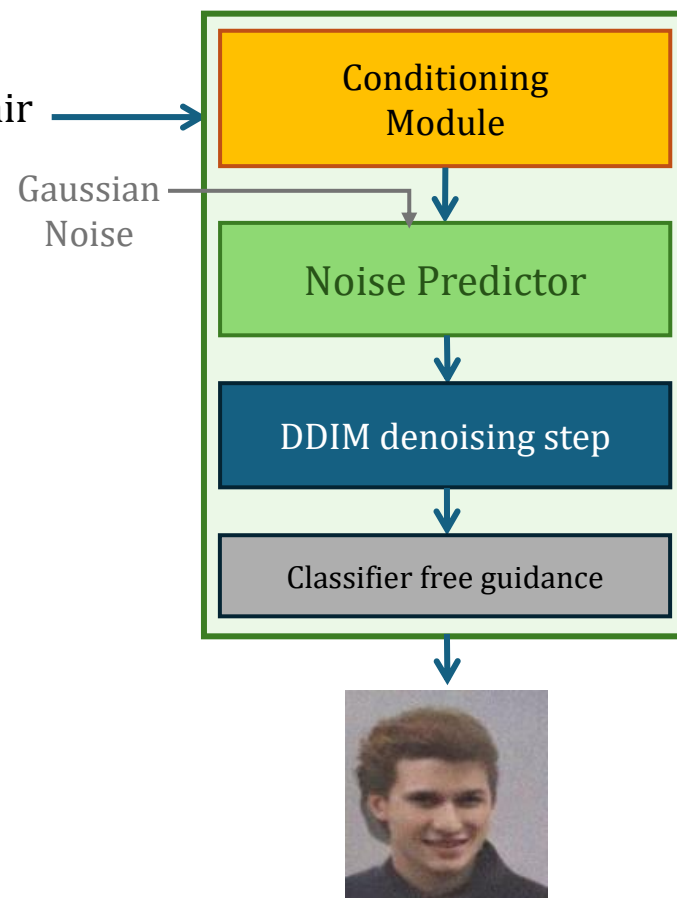Dataset divided into 8 classes with compositional labels

# Backgrounds

- **Compositional Conditional Diffusion Models**
  - Capable of generating unseen class images



$c_1$: Brown hair
$c_2$: Male

**Compositional Conditional Diffusion Model**



$c_1 = \{Gray, Black, Blonde, Brown\},$
$c_2 = \{Male, Female\}$

# Motivation

- **Can we make it faster (than DDIM) for sampling?**

**DDIM Reverse Process**

**20~50 denoising steps!**

$$X_0 \dashleftarrow X_t \dashleftarrow X_{t'} \dashleftarrow X_T$$

**CFG formulation:**

$$(1 + \omega) * x_0^{Cond} - \omega * x_0^{Uncond}$$

**CCDM**

Conditioning Module

Noise Predictor

Repeat 20~50 times (inference steps)

DDIM denoising step

40~100 steps for Conditioned and Unconditioned outputs

Classifier free guidance

# 2

# Related Works

and

# Preliminaries

# Diffusion Denoising Probabilistic Models

- Forward diffusion process and reparameterization

**Forward Diffusion Process**



$$X_0 \dashrightarrow X_{t-1} \rightarrow X_t \dashrightarrow X_T$$

$$q(x_t \mid x_{t-1}) = \mathcal{N}(x_t \,; \sqrt{1 - \beta_t} \cdot x_{t-1}, \beta_t \mathbf{I})$$
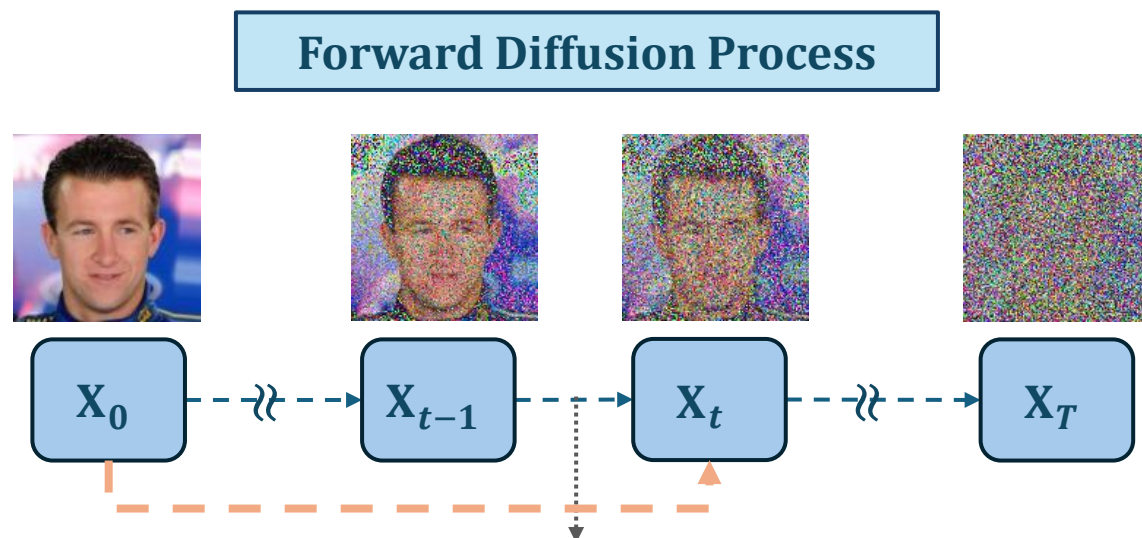
$$q(x_{1:T}|x_0) = \prod_{t=1}^{T} q(x_t|x_{t-1})$$

$$\bar{\alpha}_t = \prod_{s=1}^{t} (1 - \beta_s)$$
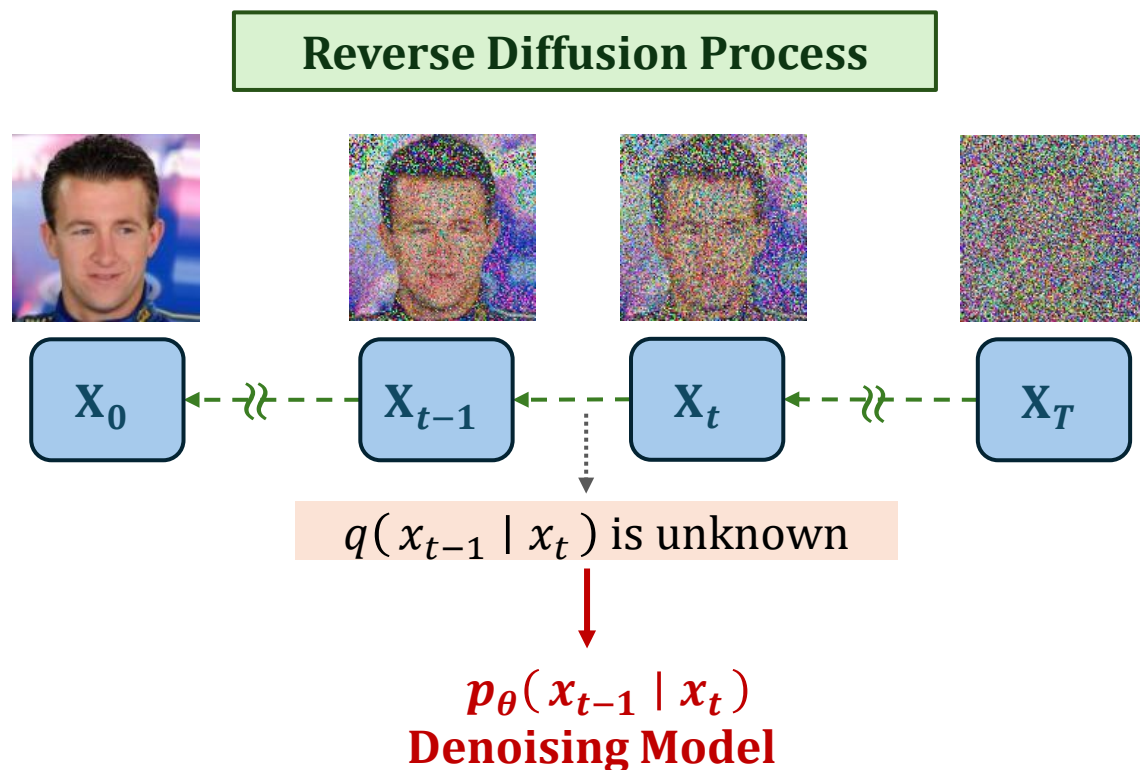
**forward process formula**

$$x_t = \sqrt{\bar{\alpha}_t} \cdot x_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon, \qquad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

Reparameterization

$$q(x_t \mid x_0) = \mathcal{N}(x_t \,; \sqrt{\bar{\alpha}_t} \cdot x_0, (1 - \bar{\alpha}_t) \cdot \mathbf{I})$$

# Diffusion Denoising Probabilistic Models

- Reverse diffusion process, training and sampling



Reverse Diffusion Process

$$X_0 \dashleftarrow X_{t-1} \leftarrow X_t \dashleftarrow X_T$$

$q(x_{t-1} \mid x_t)$ is unknown

$p_\theta(x_{t-1} \mid x_t)$
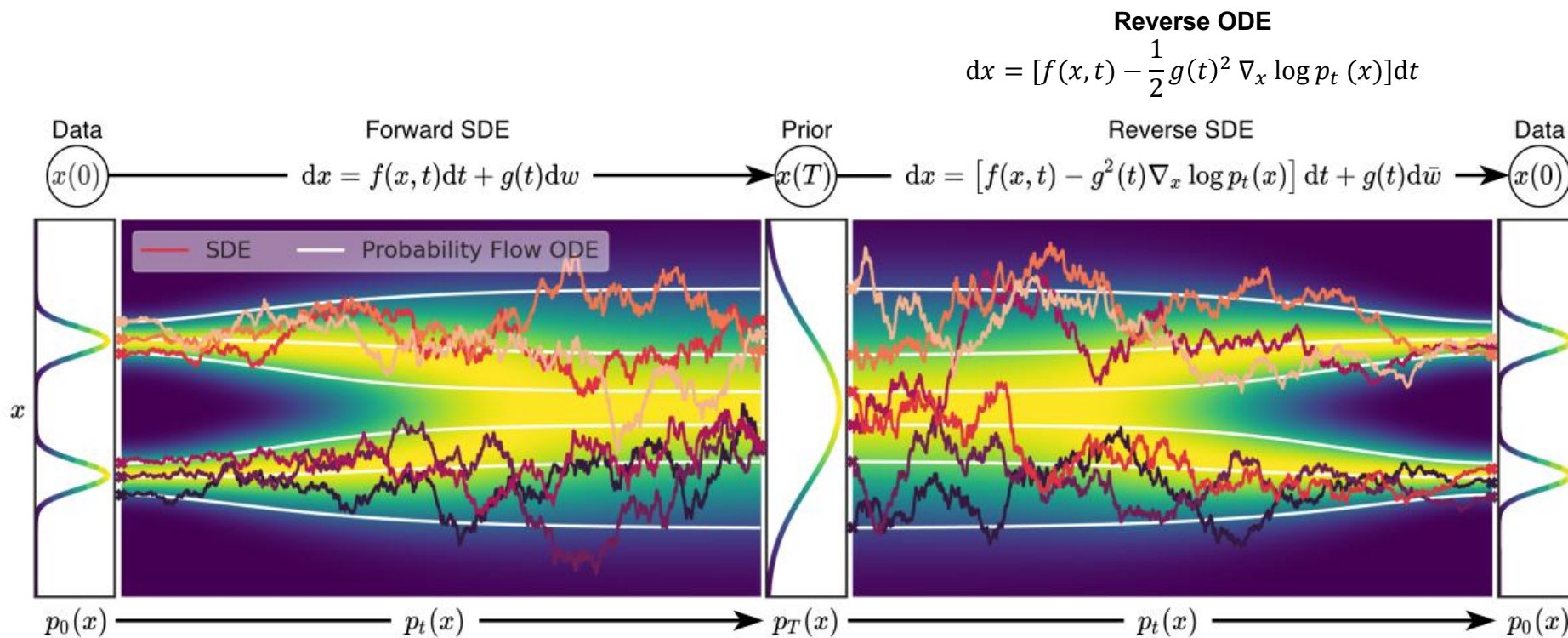**Denoising Model**

**Algorithm 1** Training
1: **repeat**
2:   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:   $t \sim \text{Uniform}(\{1, \ldots, T\})$
4:   $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:   Take gradient descent step on
**MSE loss** $\nabla_\theta \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2$
6: **until** converged

**Algorithm 2** Sampling
1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:   $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4:   $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)\right) + \sigma_t \mathbf{z}$
5: **end for**
6: **return** $\mathbf{x}_0$

**stochasticity**

# Score-Based Generative Modeling Through SDE

- Stochastic differential equation & Ordinary differential equation

**Reverse ODE**

$$dx = [f(x,t) - \frac{1}{2}g(t)^2 \nabla_x \log p_t(x)]dt$$



"For all diffusion processes, there exists a corresponding *deterministic process*, whose trajectories share the same marginal probability densities as the SDE."

# Diffusion Denoising Implicit Models

- Non-Markovian, deterministic denoising

- Speeds up sampling

- No **retraining** required to utilize DDIM sampling

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

# Classifier free guidance

- **Doesn't** need **another classifier** to guide the Diffusion model.

- When training, mostly learns **conditioned** output $x^c$ with input $c$.

- Occasionally **drops** condition for **unconditioned** output $x^\emptyset$.

- Weighted combination of conditioned output and unconditioned output.

$x_t \rightarrow$

$t \rightarrow$ **Diffusion Model** $\rightarrow x^c$

$c \rightarrow$

$x_t \rightarrow$

$t \rightarrow$ **Diffusion Model** $\rightarrow x^\emptyset$

$\emptyset \rightarrow$

$$(1 + \omega) * x^c - \omega * x^\emptyset \rightarrow x^{c,\omega}$$

➢ Problem: Needs **twice** (cond+uncond) as **many inference steps.**

# Distillation on Diffusion models

Student from previous round as the new Teacher

$t = 1$

$z_{3/4} = f(z_1; \eta)$

$z_{1/2} = f(z_{3/4}; \eta)$

Distillation

Distillation

$z_{1/4} = f(z_{1/2}; \eta)$

Distillation

$x = f(z_{1/4}; \eta)$

$x = f(z_1; \theta)$

$t = 0$

X X X

Progressive Distillation [6]

Stage 1-
eliminate CFG

$\omega$

Guidance scale embedding layer

$t$ $x_t$ $c$

$t$ $x_t$ $c/\emptyset$

student

teacher

$\hat{\epsilon}^{c,\omega}$

$\hat{\epsilon}^c$ $\hat{\epsilon}^{\emptyset}$

CFG ← $\omega$

Loss

Back Propagation

Stage 2-
Progressive Distillation

$t = 1$

$z_{3/4} = f(z_1; \eta)$

Distillation

$z_{1/2} = f(z_{3/4}; \eta)$

Distillation

$x = f(z_1; \theta)$

$z_{1/4} = f(z_{1/2}; \eta)$

Distillation

$x = f(z_{1/4}; \eta)$

$t = 0$
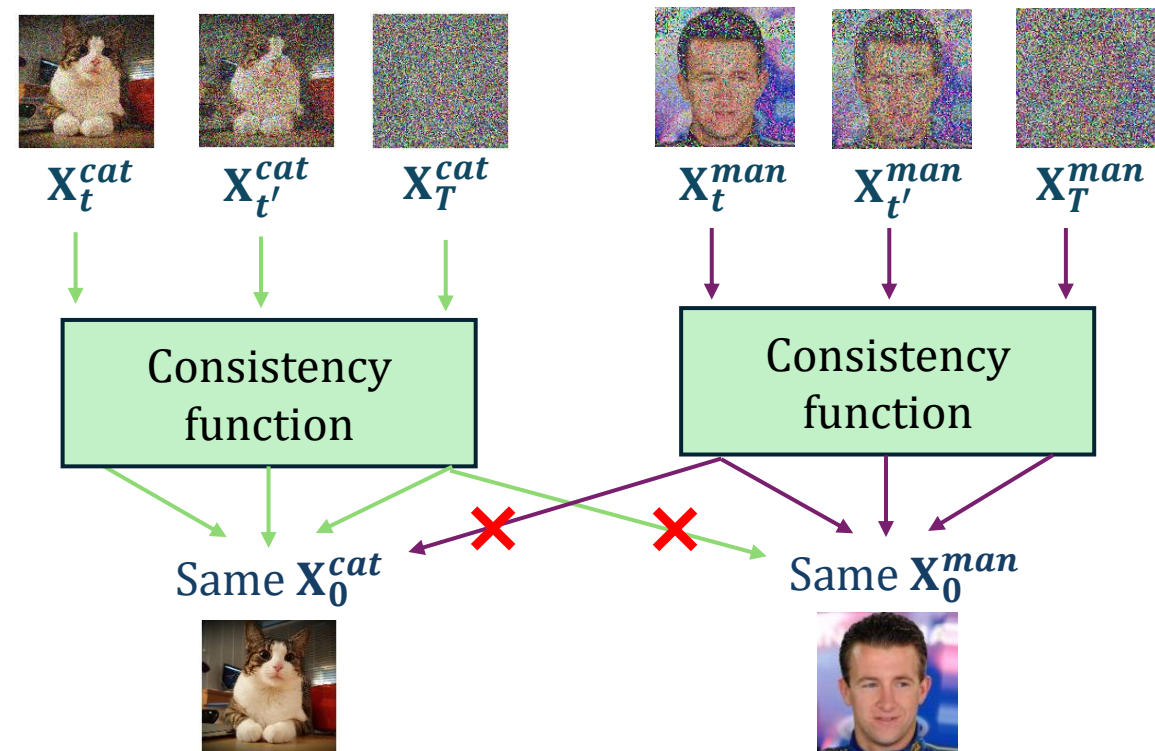
X X X

2-stage Distillation on Guided Diffusion models[7]

# Consistency Models

- Consistency function

  - Given any $t \in [0, ..T]$, exists a function: $f(x_t) \rightarrow x_0$

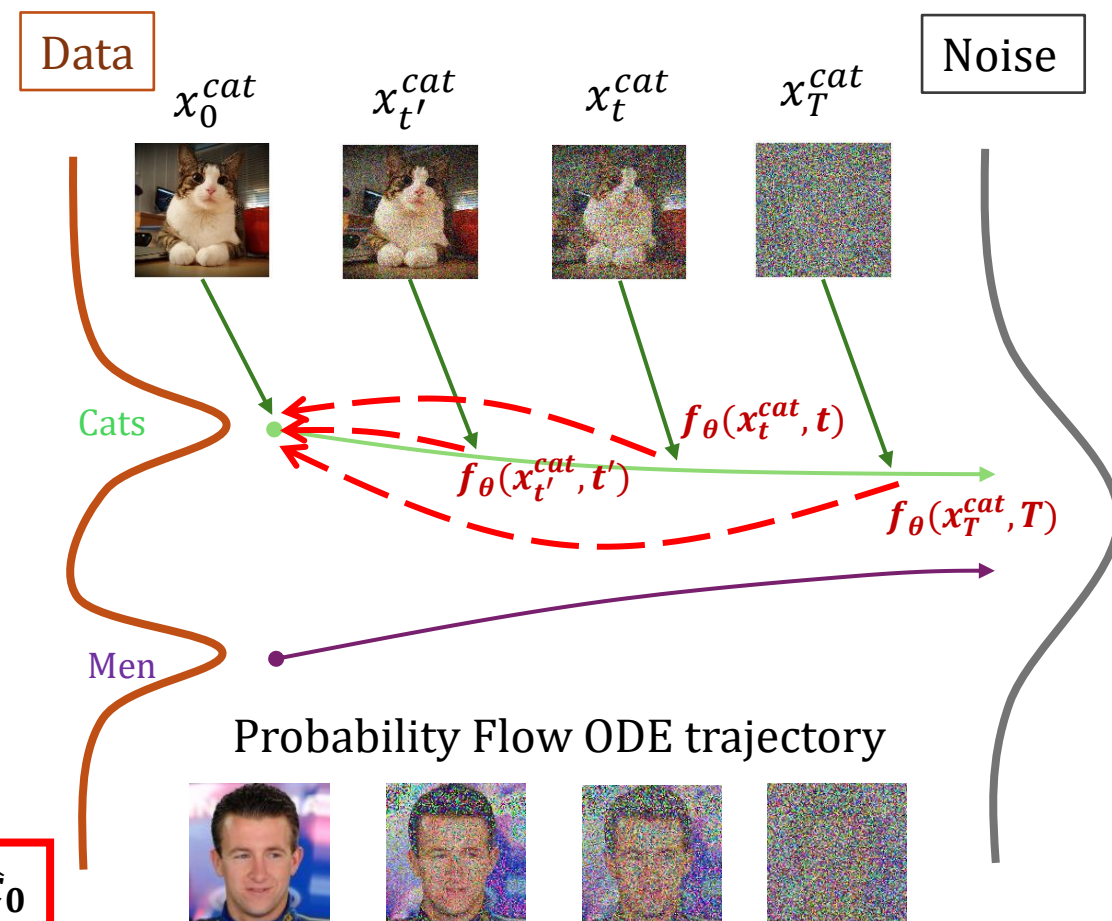  - Same $x_0$ for $x_t, x_{t'}, x_T$ on the same diffusion trajectory

# Consistency Models

- Loss function for training

  - Train a **model $f_\theta$** to approximate this function

  - Due to the **deterministic** property of the **PF-ODE**, we may train a consistency model by enforcing the **Consistency objective**

  - Loss function for training:

  $$\min_\theta [\, d\left(f_\theta(x_{t_n}, t_n), f_\theta\left(x_{t_{n-1}}, t_{n-1}\right)\right)],$$

  So that :

  $$f_\theta(x_{t_n}, t_n) = f_\theta(x_{t_{n-1}}, t_{n-1}) = f_\theta(x_{t_1}, t_1) = f_\theta(x_0, t_0) = \boxed{\hat{x}_0}$$

Data    $x_0^{cat}$    $x_{t'}^{cat}$    $x_t^{cat}$    $x_T^{cat}$    Noise

Cats

$f_\theta(x_t^{cat}, t)$

$f_\theta(x_{t'}^{cat}, t')$

$f_\theta(x_T^{cat}, T)$

Men

Probability Flow ODE trajectory

# Consistency Models

- Consistency Distillation
  - How to obtain adjacent point $x_{t_{n-1}}$?
    - ➤ May use forward process formula
    - ➤ **Or Better –** Teacher model $f_{\boldsymbol{\phi}}$ as the ODE solver $\boldsymbol{\Phi}$
  - **Consistency Distillation**

  Distill a pretrained Diffusion model into a CM

**Algorithm 2** Consistency Distillation (CD)

**Input:** dataset $\mathcal{D}$, initial model parameter $\boldsymbol{\theta}$, learning rate $\eta$, ODE solver $\Phi(\cdot, \cdot; \phi)$, $d(\cdot, \cdot)$, $\lambda(\cdot)$, and $\mu$
$\boldsymbol{\theta}^- \leftarrow \boldsymbol{\theta}$
**repeat**
   Sample $\mathbf{x} \sim \mathcal{D}$ and $n \sim \mathcal{U}[\![1, N-1]\!]$
   Sample $\mathbf{x}_{t_{n+1}} \sim \mathcal{N}(\mathbf{x}; t_{n+1}^2 \boldsymbol{I})$
   $\hat{\mathbf{x}}_{t_n}^{\phi} \leftarrow \mathbf{x}_{t_{n+1}} + (t_n - t_{n+1}) \Phi(\mathbf{x}_{t_{n+1}}, t_{n+1}; \phi)$
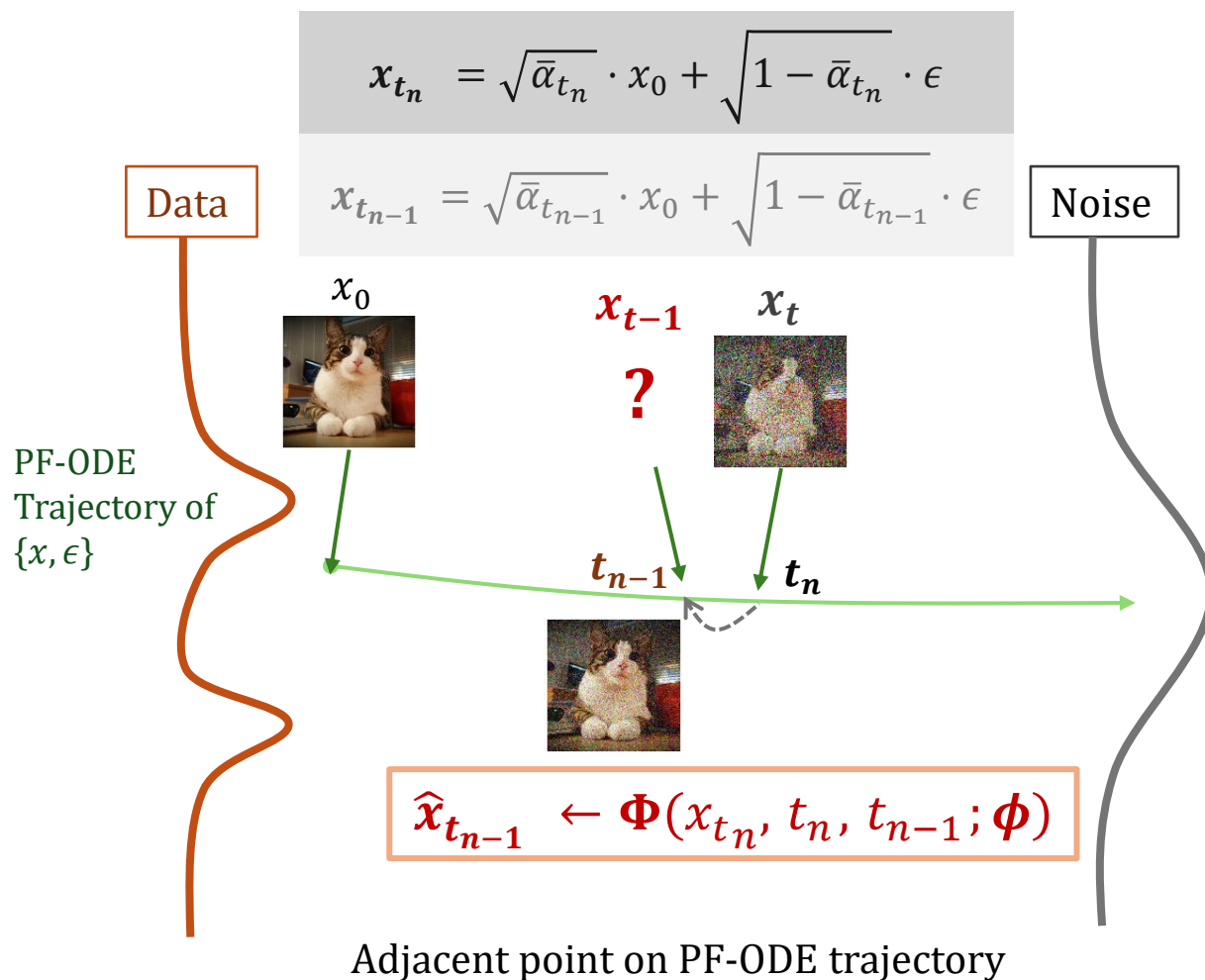   $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^-; \phi) \leftarrow$
      $\lambda(t_n) d(f_{\boldsymbol{\theta}}(\mathbf{x}_{t_{n+1}}, t_{n+1}), f_{\boldsymbol{\theta}^-}(\hat{\mathbf{x}}_{t_n}^{\phi}, t_n))$
   $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^-; \phi)$
   $\boldsymbol{\theta}^- \leftarrow \text{stopgrad}(\mu \boldsymbol{\theta}^- + (1-\mu)\boldsymbol{\theta})$
**until** convergence

$$x_{t_n} = \sqrt{\bar{\alpha}_{t_n}} \cdot x_0 + \sqrt{1 - \bar{\alpha}_{t_n}} \cdot \epsilon$$

$$x_{t_{n-1}} = \sqrt{\bar{\alpha}_{t_{n-1}}} \cdot x_0 + \sqrt{1 - \bar{\alpha}_{t_{n-1}}} \cdot \epsilon$$

Data

Noise

$x_0$

$x_{t-1}$    $x_t$

**?**

PF-ODE Trajectory of $\{x, \epsilon\}$

$t_{n-1}$    $t_n$

$$\hat{x}_{t_{n-1}} \leftarrow \Phi(x_{t_n}, t_n, t_{n-1}; \boldsymbol{\phi})$$

Adjacent point on PF-ODE trajectory

# Consistency Models

- Sampling

  Single step or Multistep Consistency Sampling



**Algorithm 1** Multistep Consistency Sampling

**Input:** Consistency model $f_\theta(\cdot, \cdot)$, sequence of time points $\tau_1 > \tau_2 > \cdots > \tau_{N-1}$, initial noise $\hat{x}_T$
$x \leftarrow f_\theta(\hat{x}_T, T)$
**for** $n = 1$ **to** $N - 1$ **do**
$\quad$ Sample $z \sim \mathcal{N}(0, I)$
$\quad \hat{x}_{\tau_n} \leftarrow x + \sqrt{\tau_n^2 - \epsilon^2} z$
$\quad x \leftarrow f_\theta(\hat{x}_{\tau_n}, \tau_n)$
**end for**
**Output:** $x$

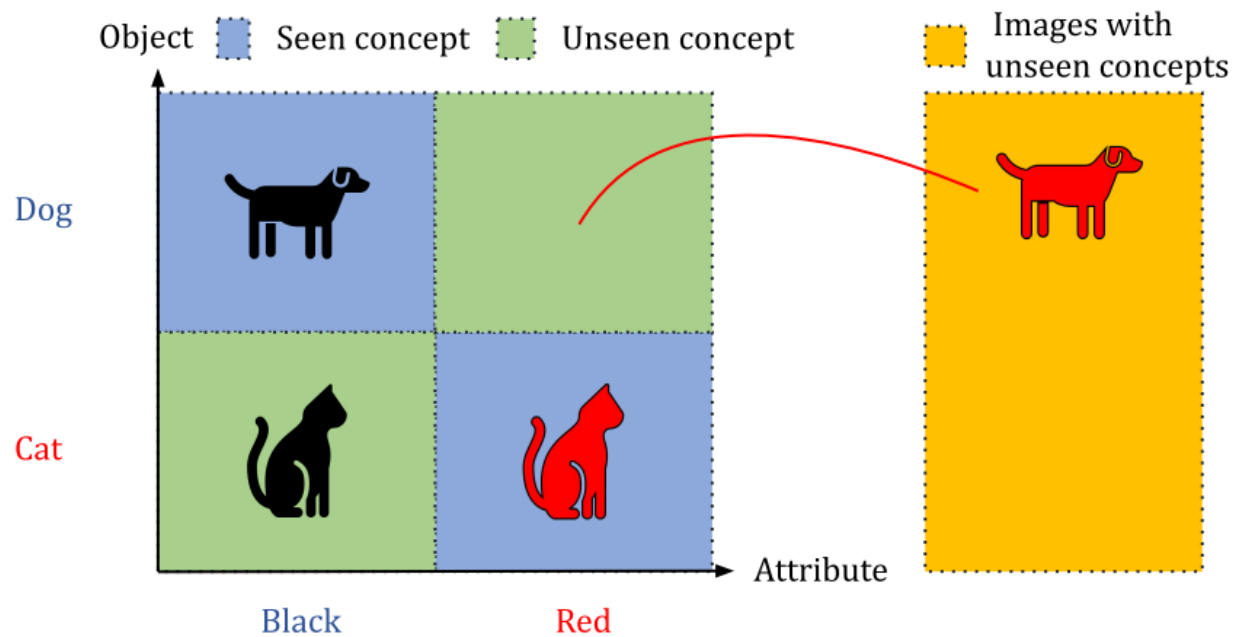Example diagram of 3 step sampling

$\widehat{\mathbf{X}}_T$ $f_\theta$ → predicted $\widehat{\mathbf{X}}_0^T$

Add noise &
Scale to $\boldsymbol{\tau_1}$

$\widehat{\mathbf{X}}_{\tau_1}$ $f_\theta$ → predicted $\widehat{\mathbf{X}}_0^{\tau_1}$

Add noise &
Scale to $\boldsymbol{\tau_2}$

$\widehat{\mathbf{X}}_{\tau_2}$ $f_\theta$ → predicted final $\widehat{\mathbf{X}}_0$
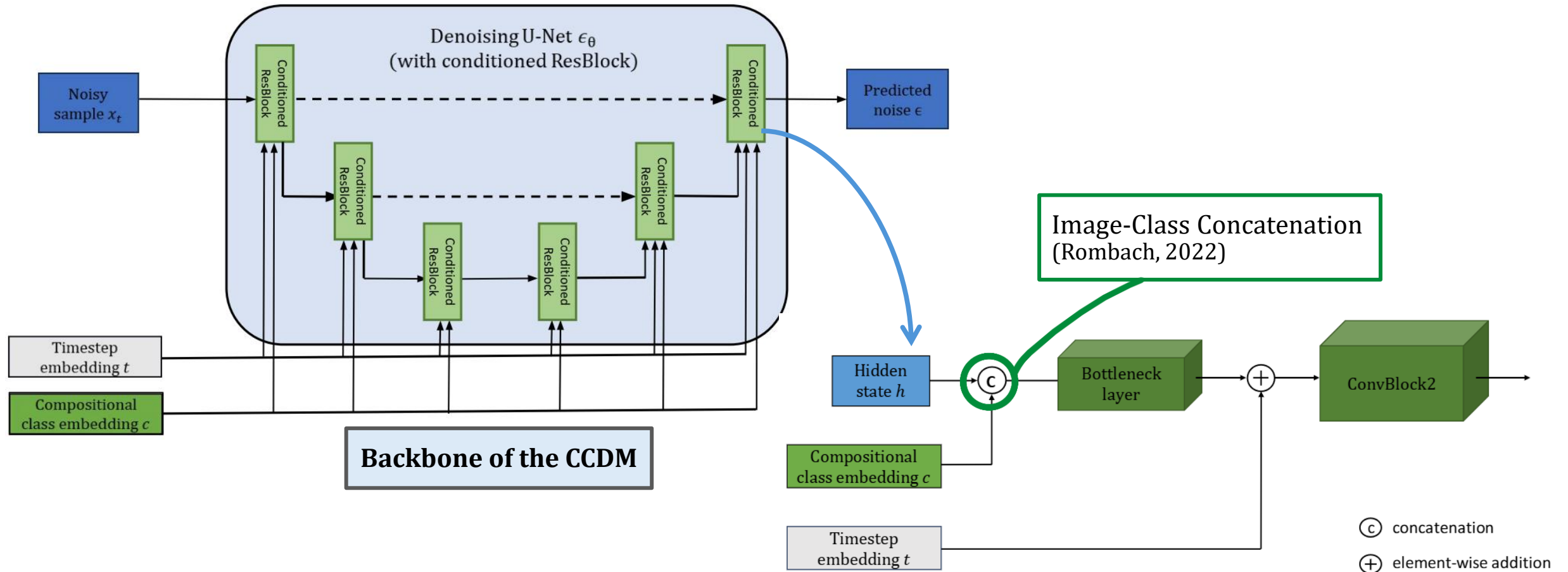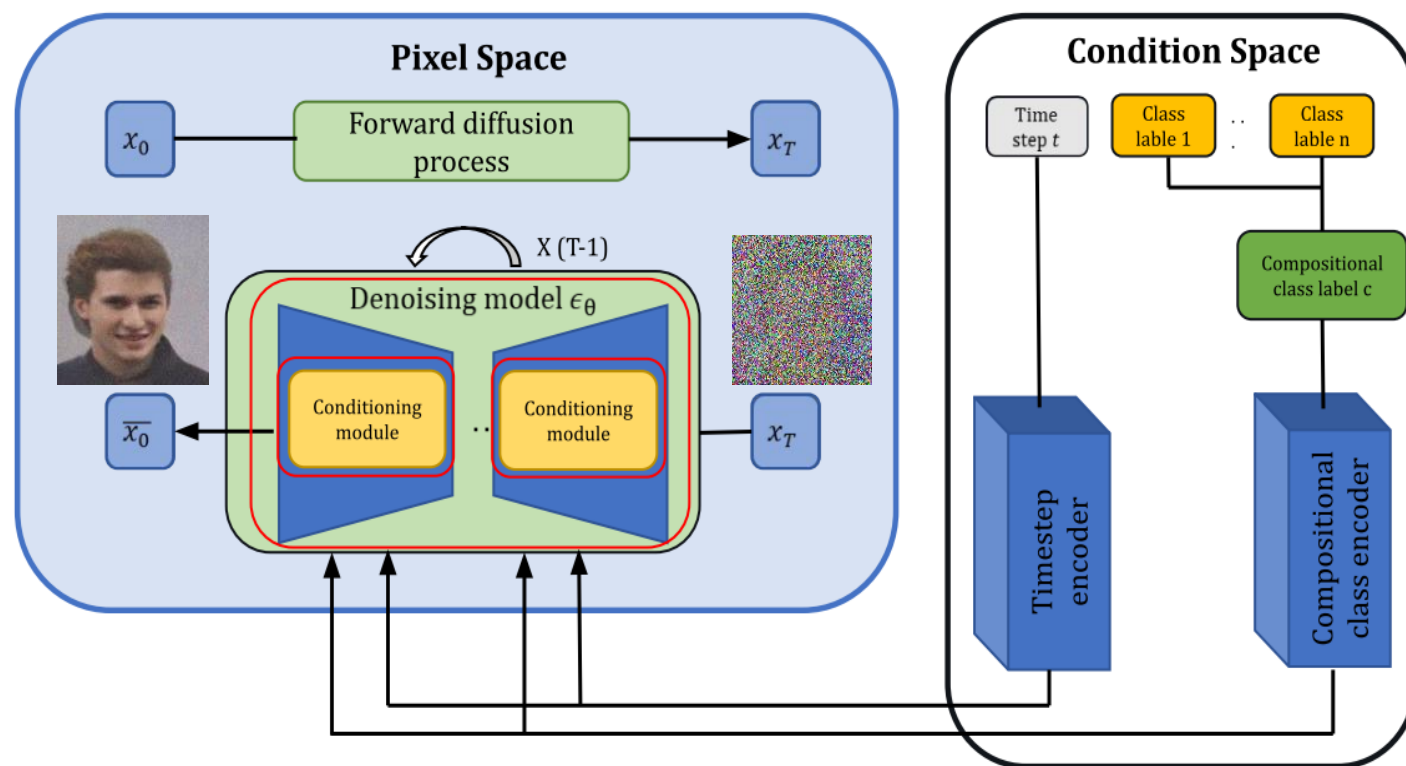
## Compositional Zero-Shot Learning

■ CZSL allows models to predict unseen classes by leveraging a combination of zero-shot learning and compositional understanding.

# Compositional Conditional Diffusion Models

# Compositional Conditional Diffusion Models

**3**

# Proposed Method

# Compositional Conditional Consistency Model

- Using CCDM as the teacher model for consistency distillation

Follows Latent Consistency Model's (Luo, 2023) implementation of :

- skipping step $k$ =20

    Reduces training time

- guidance scale embedding layer

Notations:
$n \in [1,2,..T]$: step index
$t_n \in [0,1]$: time at step $n$
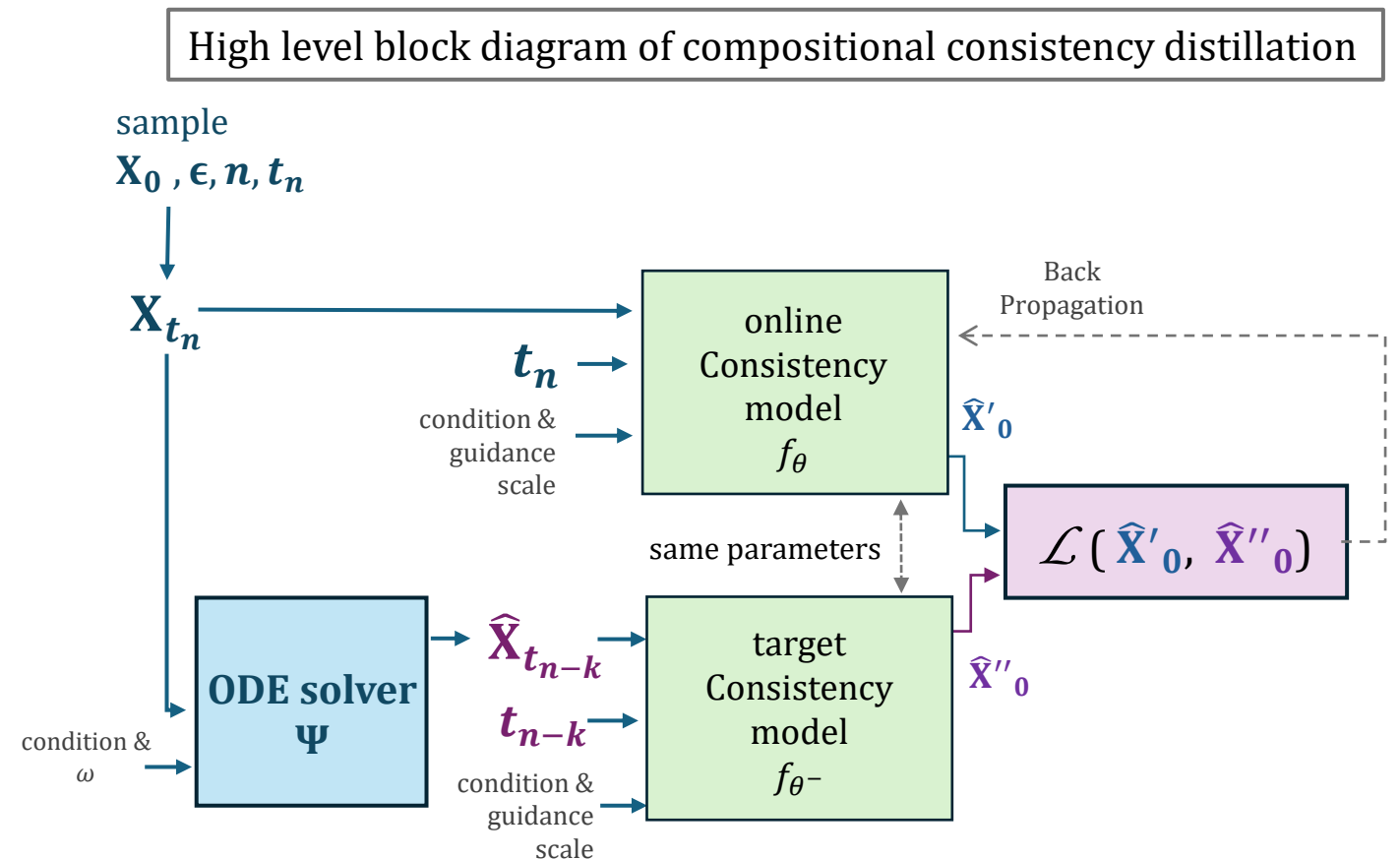$\epsilon \sim \mathcal{N}(0, \mathbf{I})$: Gaussian Noise
$X_{t_n}$: noisy image at $t_n$
$\widehat{X}_{t_n}$: predicted image at $t_n$
$X_0$: clean image from dataset
$\widehat{X}_0$: predicted clean image
$\omega$ : guidance scale

High level block diagram of compositional consistency distillation
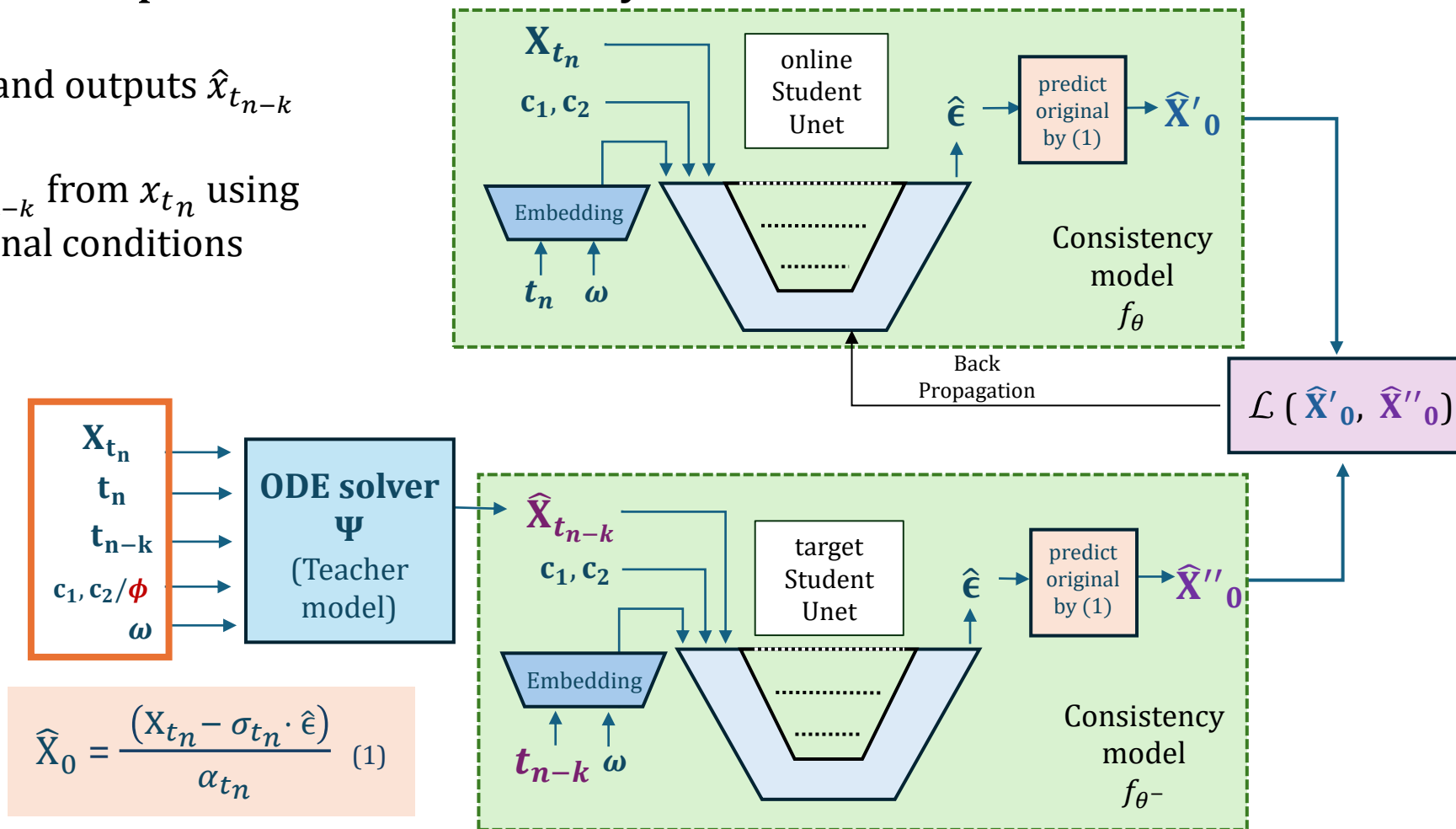
# 3. Proposed Method

## Compositional Conditional Consistency Model

- Detailed block diagram of compositional consistency distillation

  ▪ ODE solver takes input $x_{t_n}$ and outputs $\hat{x}_{t_{n-k}}$

  ▪ Deterministically estimates $\hat{x}_{t_{n-k}}$ from $x_{t_n}$ using CCDM + DDIM under compositional conditions

Notations:
$n \in [1,2,..T]$: step index
$t_n \in [0,1]$: time at step $n$
$X_{t_n}$: noisy image at $t_n$
$\hat{X}_{t_n}$: predicted image at $t_n$
$X_0$ : clean image from dataset
$\hat{X}_0$ : predicted clean image
$\alpha_{t_n}$ : signal rate at $t_n$
$\sigma_{t_n}$ : noise rate at $t_n$
$\hat{\epsilon}$ : predicted noise
$\omega$ : guidance scale

$$\hat{X}_0 = \frac{(X_{t_n} - \sigma_{t_n} \cdot \hat{\epsilon})}{\alpha_{t_n}} \quad (1)$$
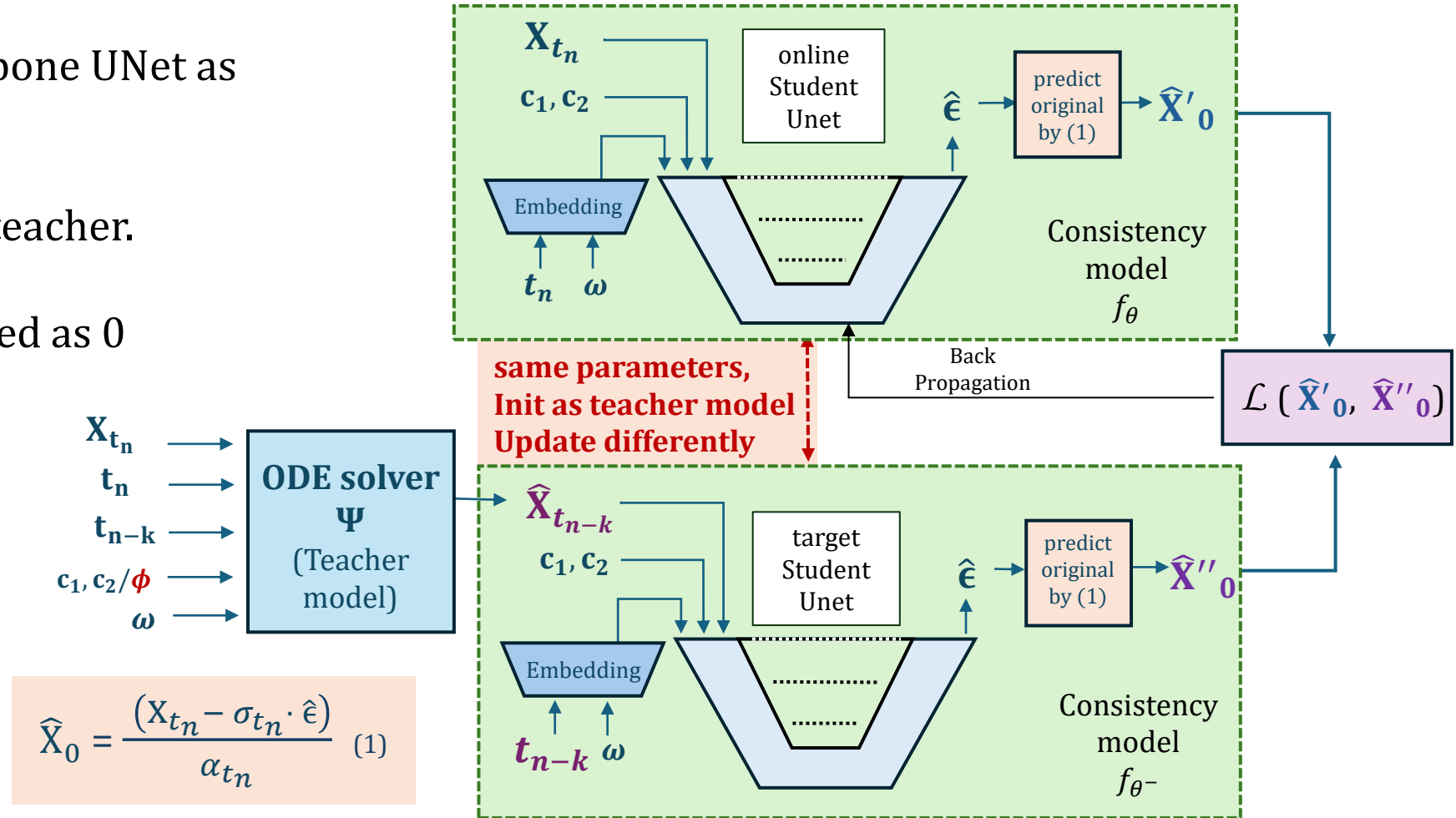
## Compositional Conditional Consistency Model

- Detailed block diagram of compositional consistency distillation

  - $f_\theta$ and $f_{\theta^-}$ using same backbone UNet as teacher (CCDM)

  - both weights initialized as teacher.

  - $\omega$ embedding layer initialized as 0

Notations:
$n \in [1, 2, \ldots T]$: step index
$t_n \in [0,1]$: time at step $n$
$X_{t_n}$: noisy image at $t_n$
$\widehat{X}_{t_n}$: predicted image at $t_n$
$X_0$: clean image from dataset
$\widehat{X}_0$: predicted clean image
$\alpha_{t_n}$: signal rate at $t_n$
$\sigma_{t_n}$: noise rate at $t_n$
$\hat{\epsilon}$: predicted noise
$\omega$: guidance scale

$$\widehat{X}_0 = \frac{(X_{t_n} - \sigma_{t_n} \cdot \hat{\epsilon})}{\alpha_{t_n}} \quad (1)$$
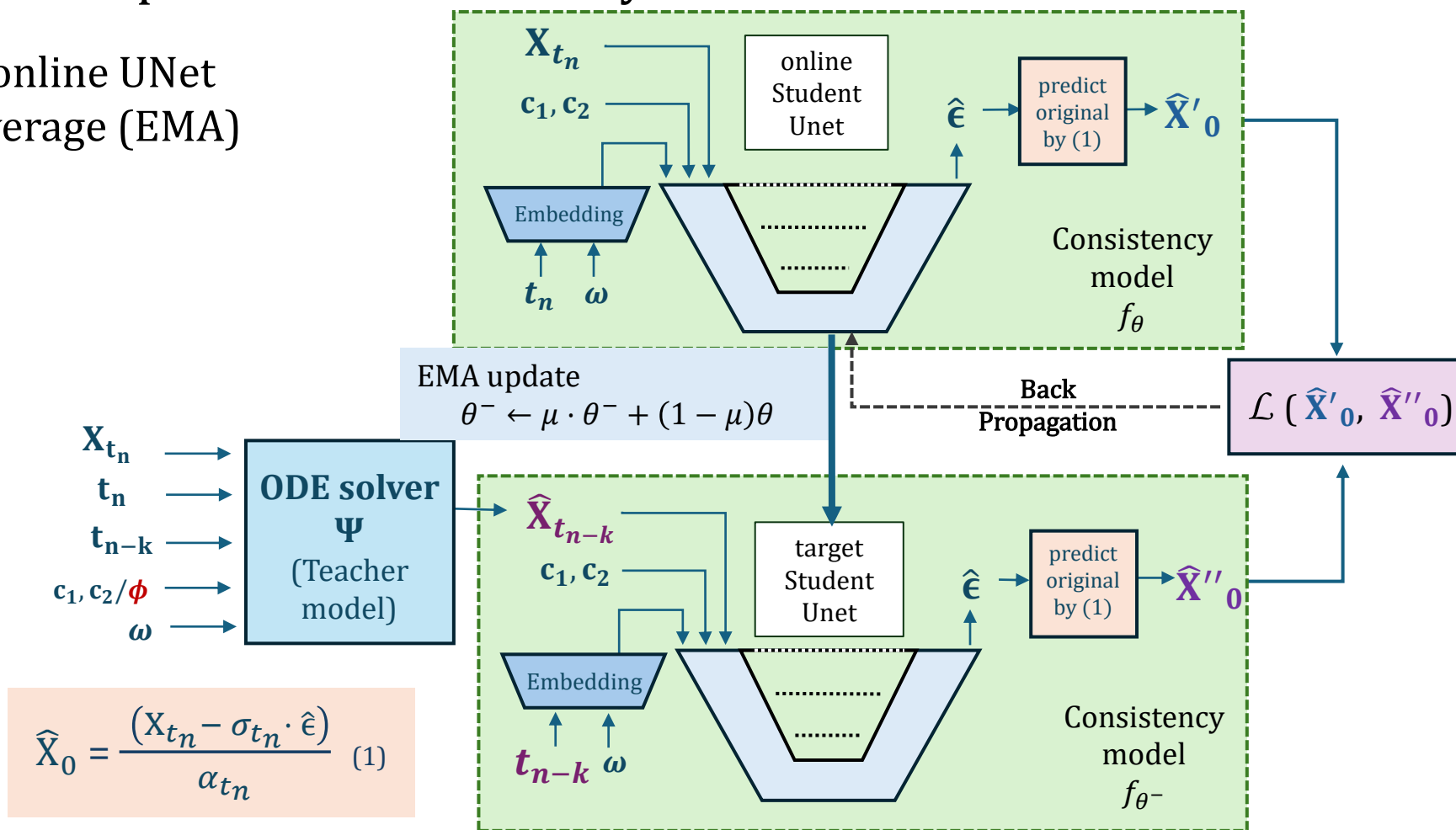
## Compositional Conditional Consistency Model

- Detailed block diagram of compositional consistency distillation

■ Target UNet updated from online UNet using Exponential Moving Average (EMA)

$X_{t_n}$

$c_1, c_2$

online Student Unet

$\hat{\epsilon}$

predict original by (1)

$\hat{X}'_0$

Embedding

$t_n$   $\omega$

Consistency model $f_\theta$

EMA update
$$\theta^- \leftarrow \mu \cdot \theta^- + (1-\mu)\theta$$

Back Propagation

$\mathcal{L}(\hat{X}'_0, \hat{X}''_0)$

Notations:
$n \in [1,2,..T]$: step index
$t_n \in [0,1]$: time at step $n$
$X_{t_n}$: noisy image at $t_n$
$\hat{X}_{t_n}$: predicted image at $t_n$
$X_0$ : clean image from dataset
$\hat{X}_0$ : predicted clean image
$\alpha_{t_n}$ :signal rate at $t_n$
$\sigma_{t_n}$ :noise rate at $t_n$
$\hat{\epsilon}$ : predicted noise
$\omega$ : guidance scale

$X_{t_n}$

$t_n$

$t_{n-k}$

$c_1, c_2/\phi$

$\omega$

**ODE solver $\Psi$** (Teacher model)

$\hat{X}_{t_{n-k}}$

$c_1, c_2$

target Student Unet

$\hat{\epsilon}$

predict original by (1)

$\hat{X}''_0$

Embedding

$t_{n-k}$   $\omega$

Consistency model $f_{\theta^-}$

$$\hat{X}_0 = \frac{(X_{t_n} - \sigma_{t_n} \cdot \hat{\epsilon})}{\alpha_{t_n}} \quad (1)$$
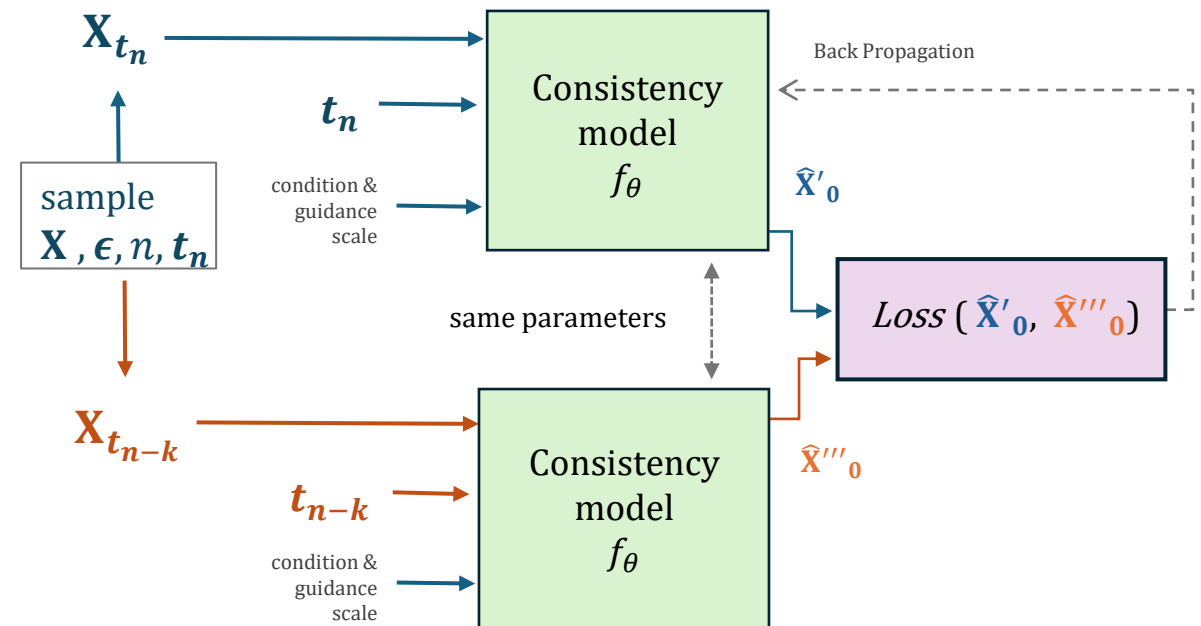
# Ways of Computing $X_{t_{n-k}}$

- Teacher predicted or Forward-process formulation?

    - CCCM with teacher's supervision offers high-quality samples.

    - Could formulated $x_{t_{n-k}}$ do better?
        - ➤ Not quite if solely rely on formulated.

    $$x_{t_{n-k}} = \sqrt{\bar{\alpha}_{t_{n-k}}} \cdot x_0 + \sqrt{1 - \bar{\alpha}_{t_{n-k}}} \cdot \epsilon, \epsilon \text{ same as in } x_{t_n}$$

    - Gradually shifting from teacher to formulated $x_{t_{n-k}}$ might help?

## Epoch function $\lambda(e)$ - Fusion implementation

- To control the strength of 2 branches: $\hat{x}_{t_{n-k}}$ and $x_{t_{n-k}}$

  - Epoch function $\lambda(e): \{1, 2, \ldots, e_{max}\} \rightarrow [0,1]$.
    1 = fully teacher signal, 0 = formulated signal

  - Switch, Step Fuse, and Loss Fuse implemented via $\lambda(e)$
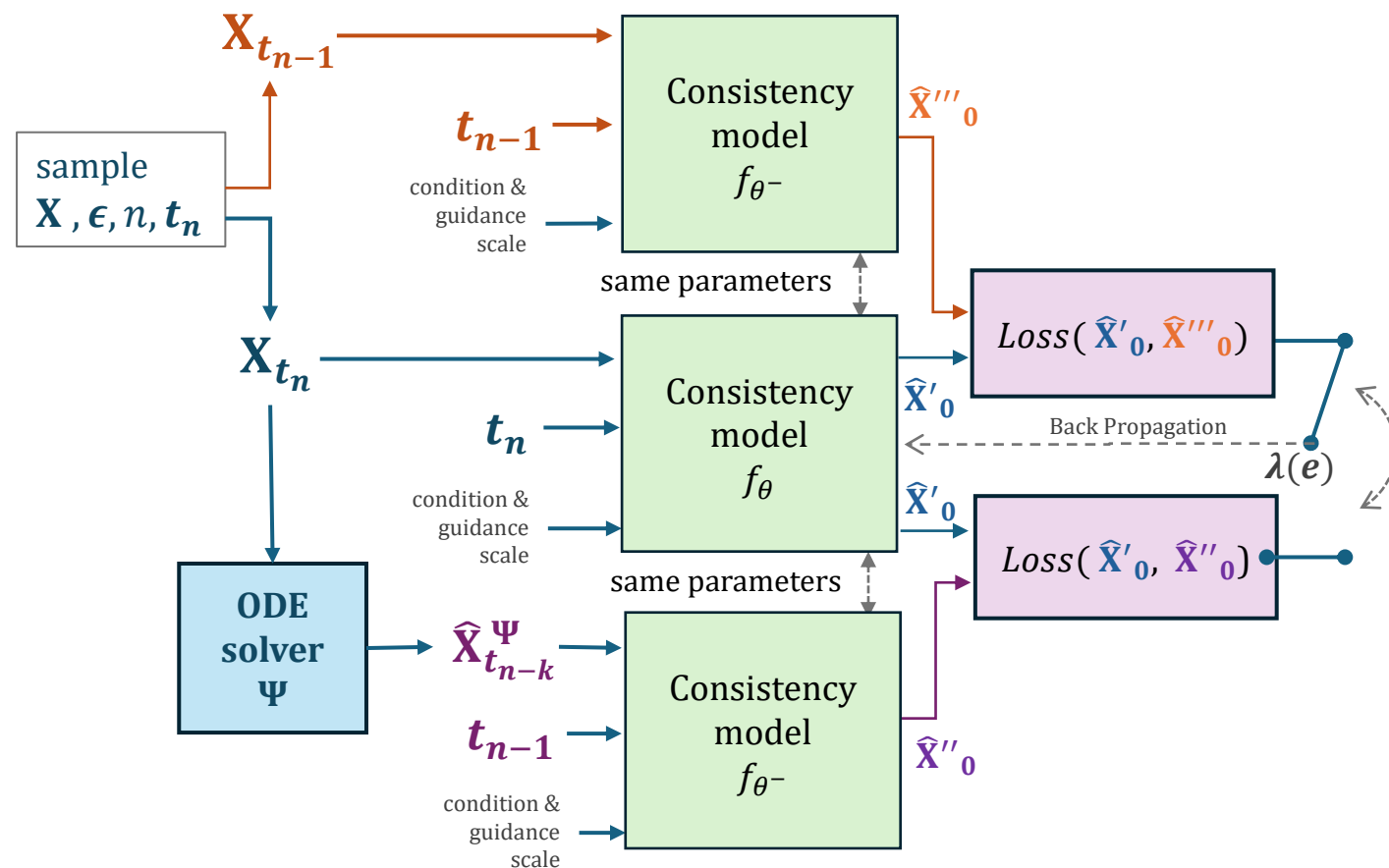
  - Controlled by Fuse Scheduler

# Modified Consistency Distillation

- Switch Strategy

  - Switches source of $x_{t_{n-k}}$ based on an epoch threshold :

  $$\lambda(e) = \begin{cases} 1, & e < threshold \\ 0, & e \geq threshold \end{cases}$$

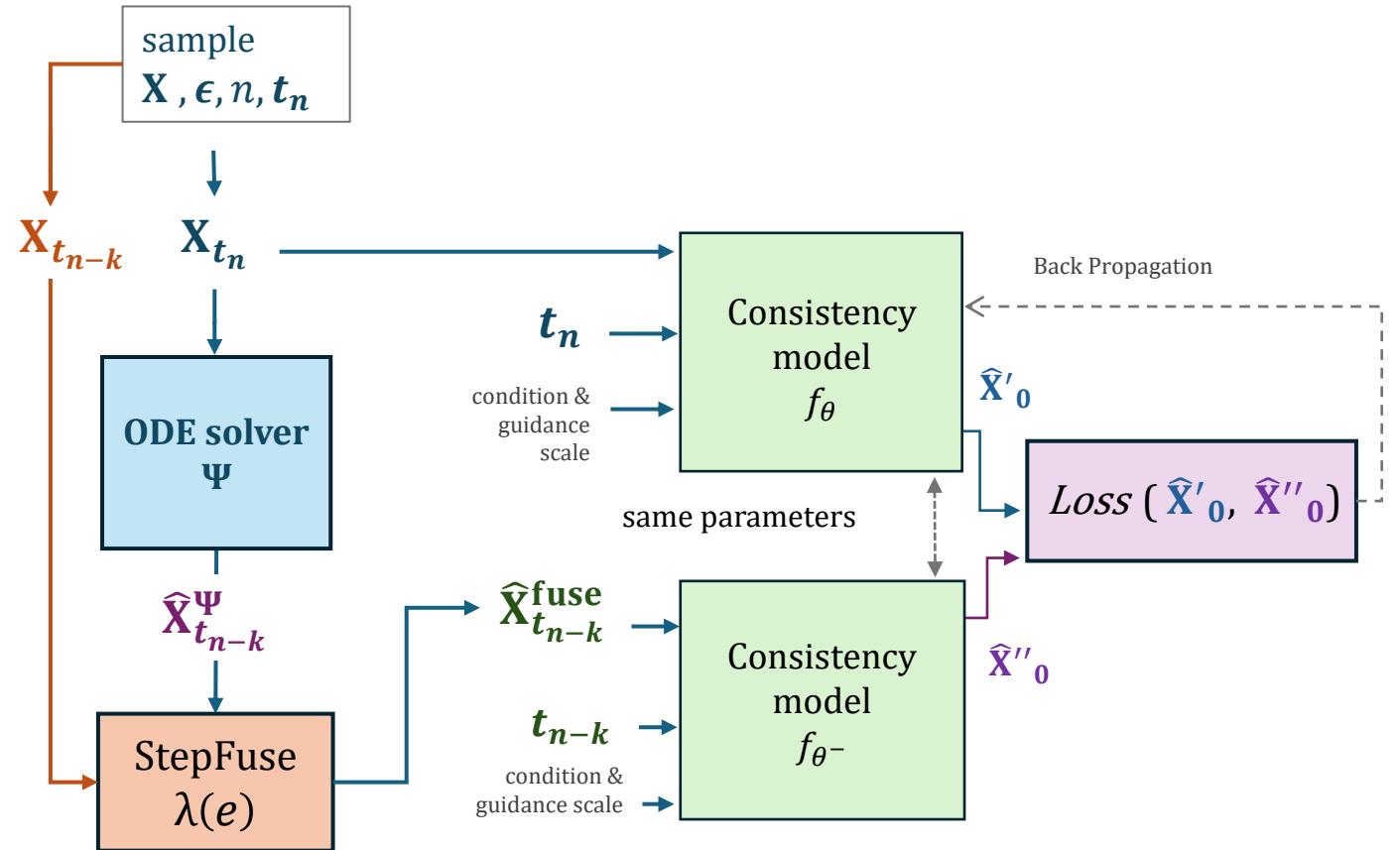# Modified Consistency Distillation

- Step-Fuse Strategy

  ▪ Pixel-wise weighted sum of two $x_{t-k}$ sources:

$$\hat{x}_{t_{n-k}}^{\text{fuse}} = \lambda(e) \cdot \hat{x}_{t_{n-k}}^{\Psi} + (1 - \lambda(e)) \cdot x_{t-k}$$

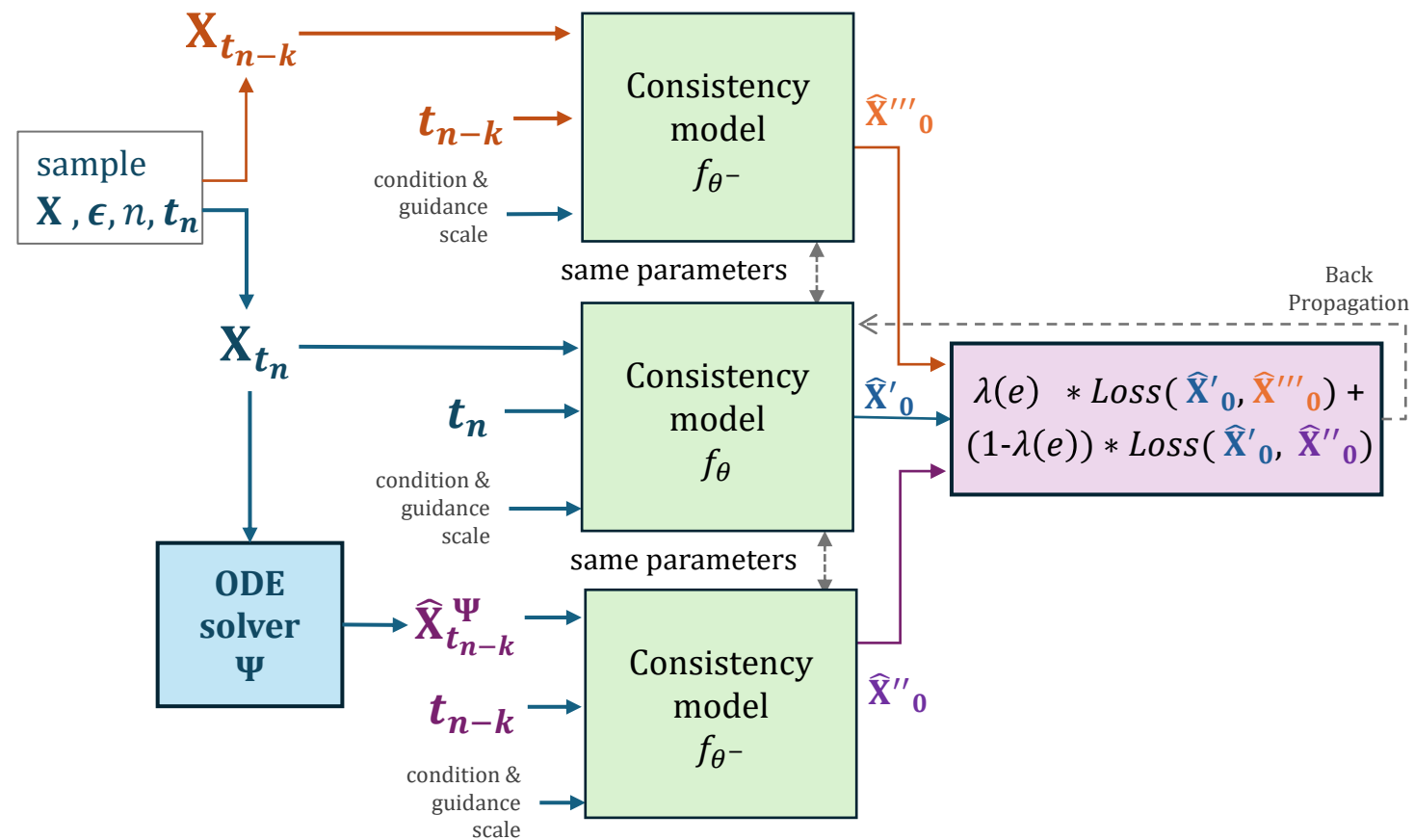# Modified Consistency Distillation

- Loss-Fuse Strategy

  ▪ Weighted sum of two loss terms:

$$\mathcal{L}_{teacher} =$$
$$\mathbb{E}\left[d\left(f_\theta(x_{t_n}, t_n, c_1, c_2, \omega), f_\theta(\hat{x}_{t_{n-k}}^\Psi, t_{n-k}, c_1, c_2, \omega)\right)\right]$$

$$\mathcal{L}_{formulated} =$$
$$\mathbb{E}\left[d\left(f_\theta(x_{t_n}, t_n, c_1, c_2, \omega), f_\theta(x_{t_{n-k}}, t_{n-k}, c_1, c_2, \omega)\right)\right]$$

$$loss = \lambda(e) \cdot \mathcal{L}_{teacher} + \left(1 - \lambda(e)\right) \cdot \mathcal{L}_{formulated}$$

$X_{t_{n-k}}$ → Consistency model $f_{\theta^-}$ → $\hat{X}'''_0$

sample $X, \epsilon, n, t_n$

condition & guidance scale

same parameters

$X_{t_n}$ → Consistency model $f_\theta$ → $\hat{X}'_0$

$t_n$, condition & guidance scale

same parameters

ODE solver $\Psi$ → $\hat{X}_{t_{n-k}}^\Psi$ → Consistency model $f_{\theta^-}$ → $\hat{X}''_0$

$t_{n-k}$, condition & guidance scale

$\lambda(e) * Loss(\hat{X}'_0, \hat{X}'''_0) + (1-\lambda(e)) * Loss(\hat{X}'_0, \hat{X}''_0)$

Back Propagation

## Fuse Scheduler design

- Constant:

  $$\lambda(epoch) = \lambda_0 \, ,$$

  fixed blending weight throughout training.

- $epoch \in [1, epoch_{max}]$

- $prog = \dfrac{epoch}{epoch_{Max}} \in [0,1]$

- Exponential Decay:

  $$\lambda(epoch) = e^{-\gamma \cdot prog} \cdot (1 - prog),$$

  $\gamma$ = decay rate

- Piecewise Linear:

  $$\lambda(epoch) = \lambda_i + \frac{\lambda_{i+1} - \lambda_i}{p_{i+1} - p_i}(prog - p_i)$$

  Defined by control points :$\{(p_i, \lambda_i)\}, \quad p_i, \lambda_i \in [0,1]$.

  Linearly interpolated over epochs

## Summary of proposed Strategies

| Method | Source of $x_{t-k}$ | Description |
|---|---|---|
| Fully Teacher | Teacher Model Prediction (via ODE Solver) | Uses noise predicted by teacher model to estimate $x_{t-k}$ through an ODE solver. |
| Fully diffusion formula | Reused Forward Process Noise | $x_{t-k}$ formulated using the same noise as $x_t$. |
| Step Fusion | Mixed (Teacher & Diffusion Formula) | Pixel-wise weighted combination of both teacher and formulated $x_{t-k}$. |
| Loss Fusion | Both (Teacher & Diffusion Formula) | Computes separate losses, combines with weighting. |
| Switch | Alternating | Uses teacher before an epoch threshold, switches to formulated $x_{t-k}$ after. |

**4**

# **Experiments**

# Experiment setup

- Dataset preparation

  - $c_1$ : Hair color, 4 classes
  - $c_2$ : Gender,    2 classes
  - Total of 4*2=8 composition classes,
  - Unseen: (Brown hair, Male)
  - Dataset image counts: 1k~20k
  - Image size: 128*128



Training set with  compositional class labels $c_1, c_2$

# Experiment setup

- CCCM and Baseline Training Configuration

### Baseline CCDM

| Training Epochs | 120 |
|---|---|
| Learning rate, scheduler | 5e-5 linear with warmup |
| Loss type | L2-norm |

### Our CCCM

| Training Epochs | 80 |
|---|---|
| Learning rate, scheduler | 5e-6 constant with warmup |
| Loss type | Huber loss |
| Guidance scale interval | 2.6, 3.0 |
| Fuse Method | Fully teacher |

# Experiment setup

- Modified Consistency Distillation Configuration

**Fully Teacher**

| Fuse Method | Fully Teacher |
|---|---|
| Fuse Scheduler | Constant = 1 |

**Fully Formulated $x_{t-k}$**

| Fuse Method | Fully Formulated |
|---|---|
| Fuse Scheduler | Constant = 0 |

**Switch 32**

| Fuse Method | Switch |
|---|---|
| Fuse Scheduler | Threshold = 32 |

**Switch 48**

| Fuse Method | Switch |
|---|---|
| Fuse Scheduler | Threshold = 48 |

**Unchanged Hyperparameters**

| Training Epochs | 80 |
|---|---|
| Learning rate, scheduler | 5e-6 constant with warmup |
| Loss type | Huber loss |
| Guidance scale interval | 2.6, 3.0 |

# Experiment setup

- Modified Consistency Distillation Configuration

**Loss Fuse Constant**

| Fuse Method | Loss Fuse |
|---|---|
| Fuse Scheduler | Constant = 0.8 |

**Loss Fuse piecewise**

| Fuse Method | Loss Fuse |
|---|---|
| Fuse Scheduler | Piecewise Linear, $(40, 0.5)$ |

$\lambda(epoch)$ drops to 0.5 at epoch 40 linearly, holds at 0.5 till finished.

**Loss Fuse exponential**

| Fuse Method | Loss Fuse |
|---|---|
| Fuse Scheduler | Exponential Decay, $\gamma = 2.0$ |

**Step Fuse exponential**

| Fuse Method | Step Fuse |
|---|---|
| Fuse Scheduler | Exponential Decay, $\gamma = 2.0$ |

**Unchanged Hyperparameters**

| Training Epochs | 80 |
|---|---|
| Learning rate, scheduler | 5e-6 constant with warmup |
| Loss type | Huber loss |
| Guidance scale interval | 2.6, 3.0 |

# Experiment setup

- Evaluation metrics

  - **FID** score:

    measures the similarity between generated and real image distributions.

  - Each model generates ~7k images (500 or 1000 per class) for evaluation

  - **Unseen Class Accuracy** by human evaluation:

    Compositional class $(Brown\_hair, Male)$ is generated, to test the compositional zero-shot image generation ability.

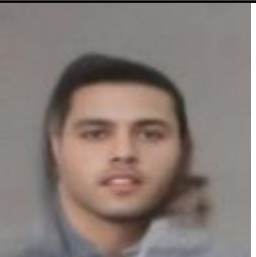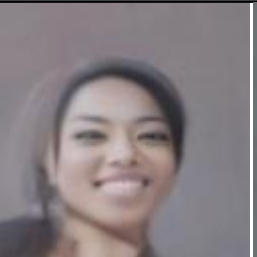  - Each model generates 300 $(Brown\_hair, Male)$ images.

**Standards:**

- **Unseen class judged by both attributes.**
  Masculine features required (e.g., short hair, no makeup). Feminine traits → failure, even if hair color is correct.

## Qualitative result

- Baseline CCDM vs Fully teacher CCCM

| 4-step sampling | Blonde Hair | | Black Hair | | Brown Hair | | Gray Hair | |
|---|---|---|---|---|---|---|---|---|
| | Male | Female | Male | Female | Male | Female | Male | Female |
| Baseline CCDM (with 8 inferences) | | | | | | | | |
| CCCM (with 4 inferences) | | | | | | | | |

# Qualitative result

- Baseline CCDM vs Fully teacher CCCM

| 2-step sampling | Blonde Hair | | Black Hair | | Brown Hair | | Gray Hair | |
|---|---|---|---|---|---|---|---|---|
| | Male | Female | Male | Female | Male | Female | Male | Female |
| Baseline CCDM (with 4 inferences) | | | | | | | | |
| CCCM (with 2 inferences) | | | | | | | | |

## Quantitative result

- FID scores under 2,3 and 4 steps sampling

**2 steps**

| Method | FID score ↓ |
|---|---|
| Baseline DDIM | 207.99 |
| Forward-process $x_{t-k}$ | 115.66 |
| Step Fuse (exponential) | 97.13 |
| Switch (threshold = 32) | 94.99 |
| Switch (threshold = 48) | 93.97 |
| Loss Fuse (exponential) | 91.27 |
| Loss Fuse (piecewise = 40:0.5) | 88.82 |
| Loss Fuse (constant = 0.8) | 88.47 |
| **Fully Teacher $x_{t-k}$** | **81.30** |

**3 steps**

| Method | FID score ↓ |
|---|---|
| Baseline DDIM | 136.68 |
| Forward-process $x_{t-1}$ | 92.15 |
| Step Fuse (exponential) | 83.73 |
| Switch (threshold = 32) | 85.66 |
| Switch (threshold = 48) | 84.23 |
| Loss Fuse (exponential) | 82.90 |
| Loss Fuse (piecewise = 40:0.5) | 80.16 |
| Loss Fuse (constant = 0.8) | 77.85 |
| **Fully Teacher $x_{t-k}$** | **73.27** |

**4 steps**

| Method | FID score ↓ |
|---|---|
| Baseline DDIM | 94.67 |
| Diffusion-formulated $x_{t-1}$ | 80.94 |
| Step Fuse (exponential) | 77.27 |
| Switch (threshold = 32) | 76.26 |
| Switch (threshold = 48) | 75.86 |
| Loss Fuse (exponential) | 75.34 |
| Loss Fuse (piecewise = 40:0.5) | 73.45 |
| Loss Fuse (constant = 0.8) | 70.15 |
| **Fully Teacher $x_{t-1}$** | **68.11** |

## Unseen Accuracy evaluation

- 300 images of Brown hair Male, compositional zero shot generation

4 steps sampling

| Method | Acc% ↑ |
|---|---|
| Baseline DDIM | 40.6% |
| Diffusion-formulated $x_{t_{n-k}}$ | 43.0% |
| Step Fuse (exponential) | 51.6% |
| Switch (threshold = 32) | 49.3% |
| Switch (threshold = 48) | 49.6% |
| Loss Fuse (exponential) | 51.6% |
| **Loss Fuse (piecewise = (40, 0.5))** | **52.0%** |
| Loss Fuse (constant = 0.8) | 50.6% |
| Fully Teacher $x_{t_{n-k}}$ | 47.0% |



Loss Fuse
piecewise (40,05)



Fully teacher $x_{t_{n-k}}$

# 5. Conclusions

# Conclusion

- We propose the **Compositional Conditional Consistency Model.**

- Achieving **faster sampling** speed than CCDM.

- **Preserves unseen** class generation.

- Observed that modified consistency distillation strategies yield **better unseen accuracy.**
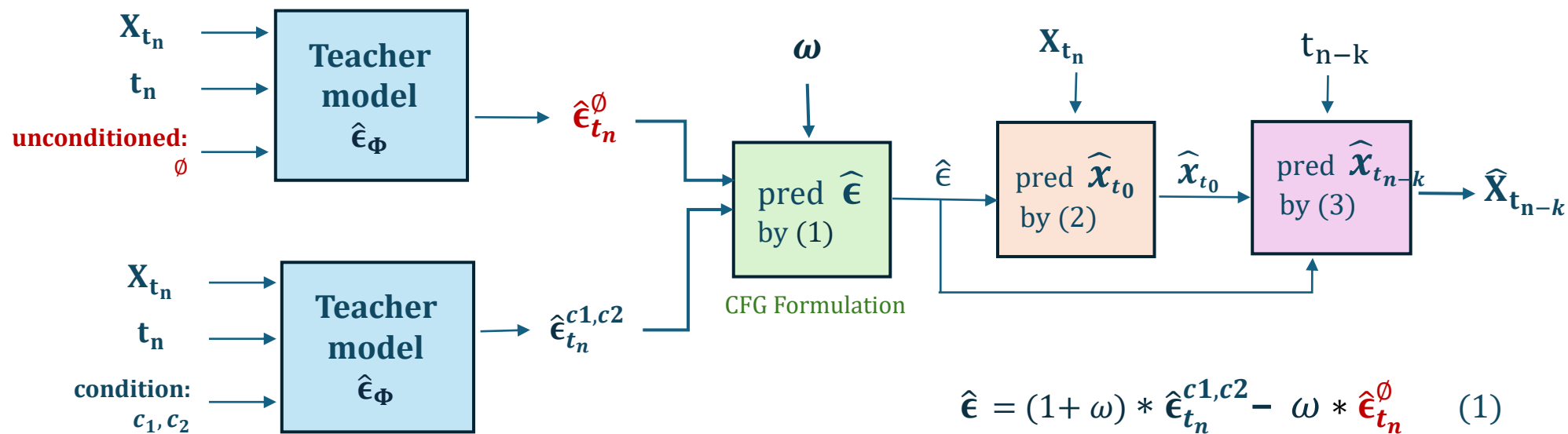
    Our guess is: the teacher model may introduce **bias**, encouraging the student to generate seen or

    **high-confidence images**, which **limits generalization**.

# References

- **"Compositional Conditional Diffusion Model",** S. -L. Lai, P. -C. Chen & C. -W. Ma, 2024.

- "**Consistency Models",** Y. Song, P. Dhariwal, M. Chen & I. Sutskever, 2023.

- **"Latent consistency models: Synthesizing high-resolution images with few-step inference"**, Luo, S., Tan, Y., Huang, L., Li, J., & Zhao, H. 2023.

- **"Improved techniques for training consistency models",** Y. Song and P. Dhariwal, 2023.

- **"Score-based generative modeling through stochastic differential equations",** Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon & B. Poole, 2020.

- **"Denoising diffusion implicit models",** J. Song, C. Meng & S. Ermon**,** 2020.

- **"Progressive distillation for fast sampling of diffusion models",** T. Salimans and J. Ho, 2022.

- **"On distillation of guided diffusion models",** C. Meng, R. Rombach, R. Gao, D. Kingma, S. Ermon, J. Ho, and T. Salimans, 2023.

- **"Classifier-free diffusion guidance",** J. Ho and T. Salimans, 2022.

- **"High-resolution image synthesis with latent diffusion models",** R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, 2022.

# ODE solver

- CCCM uses DDIM as ODE solver



Notations:
$t_n$ : time at step n, $t_n \in [0,1], n \in [1..T]$
$\alpha_{t_n}$ : signal rate $\qquad$ ; $\sigma_{t_n}$ : noise rate
$X_{t_n}$ : image at $t_n$ $\qquad$ ; $\widehat{X}_{t_n}$ : predicted image at $t_n$
$\hat{\epsilon}_{t_n}$ : predicted noise at $t_n$
$\omega$ : guidance scale

$$\hat{\epsilon} = (1+\omega) * \hat{\epsilon}_{t_n}^{c1,c2} - \omega * \hat{\epsilon}_{t_n}^{\emptyset} \qquad (1)$$

$$\widehat{X}_{t_0} = \frac{X_{t_n} - \sigma_{t_n} * \hat{\epsilon}}{\alpha_{t_n}} \qquad (2)$$
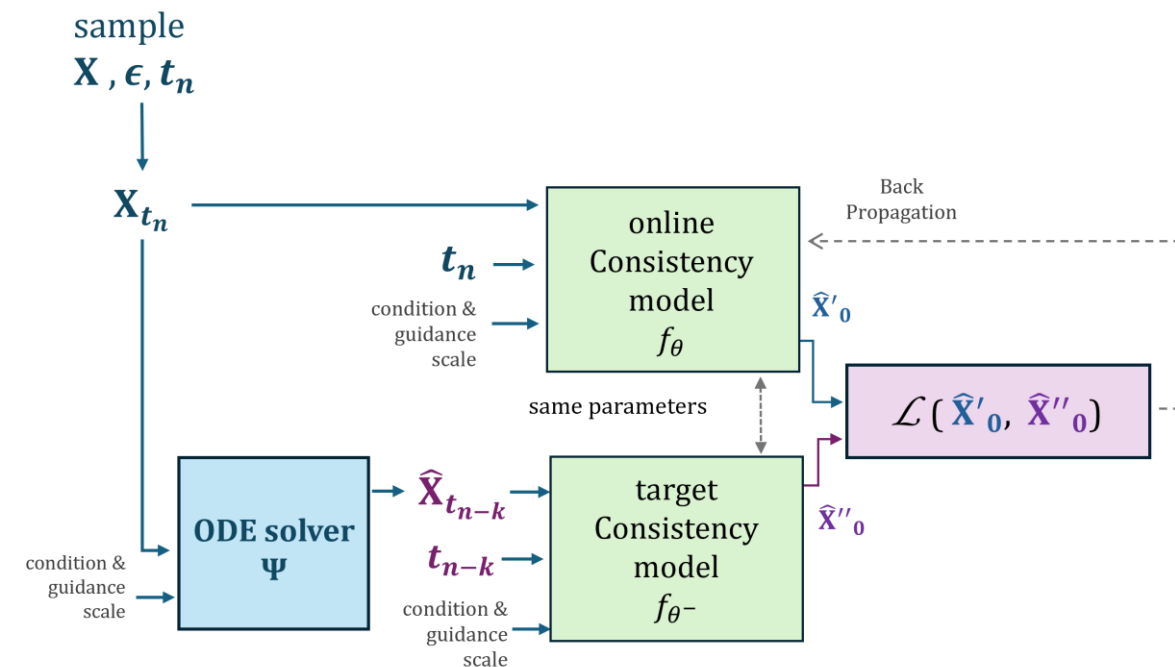
$$\widehat{X}_{t_{n-k}} = \alpha_{t_{n-k}} * \widehat{X}_0 + \sigma_{t_{n-k}} * \hat{\epsilon} \qquad (3)$$

# Guidance scale embedding

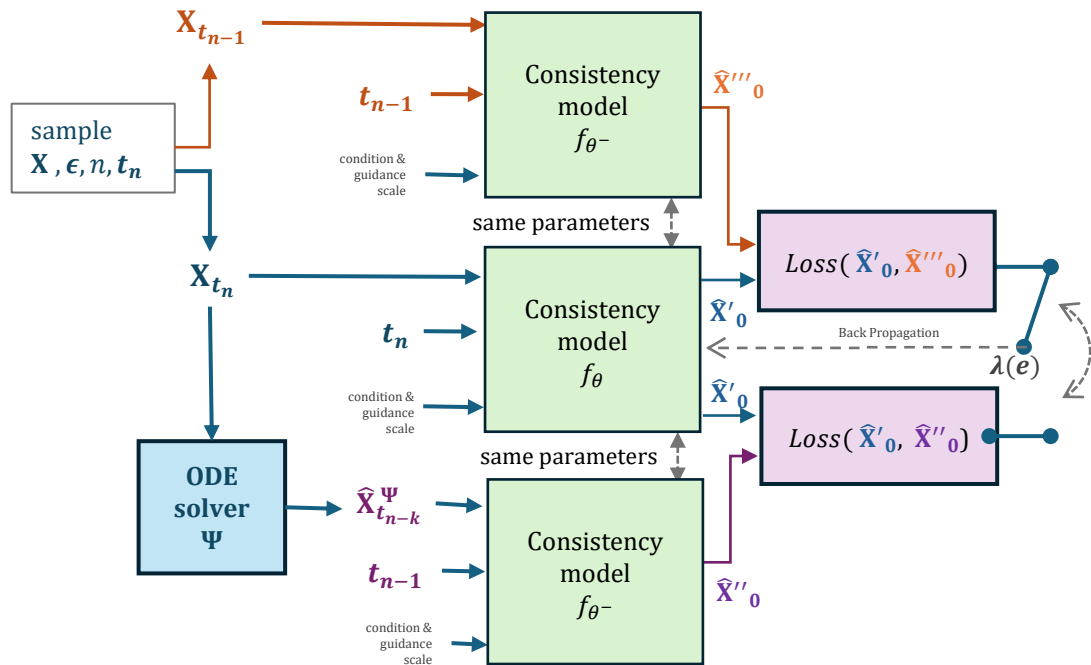# Compositional Consistency Distillation Algo

- Fully teacher



**Algorithm 4.1** Compositional label Consistency Distillation

**Input:** dataset $\mathcal{D}$, initial model parameter $\theta$, learning rate $\eta$, ODE solver $\Psi(\cdot, \cdot, \cdot, \cdot)$, distance metric $d(\cdot, \cdot)$, guidance scale $[\omega_{min}, \omega_{max}]$, skipping steps $k$

1: $\theta^- \leftarrow \theta$
2: **repeat**
3:      Sample $(x, c_1, c_2) \sim \mathcal{D}$, and $\omega \sim [\omega_{min}, \omega_{max}]$
4:      Sample $\epsilon \sim \mathcal{N}(0, I)$, $n \sim \mathcal{U}[1 + k, N]$
5:      $x_{t_n} \leftarrow \sqrt{\bar{\alpha}_{t_n}} \cdot x + \sqrt{1 - \bar{\alpha}_{t_n}} \cdot \epsilon$
6:      $x_{t_{n-k}}^{\Psi, \omega} \leftarrow x_{t_n} + (1 + \omega) \cdot \Psi(x_{t_n}, t_n, t_{n-k}, c_1, c_2) - \omega \cdot \Psi(x_{t_n}, t_n, t_{n-k}, \varnothing, \varnothing)$
7:      $\mathcal{L}(\theta, \theta^-; \Psi) \leftarrow d\left(f_\theta(x_{t_n}, t_n, c_1, c_2, \omega), \ f_{\theta^-}(\hat{x}_{t_{n-k}}, t_{n-k}, c_1, c_2, \omega)\right)$
8:      $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta, \theta^-)$
9:      $\theta^- \leftarrow \texttt{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$
10: **until** convergence

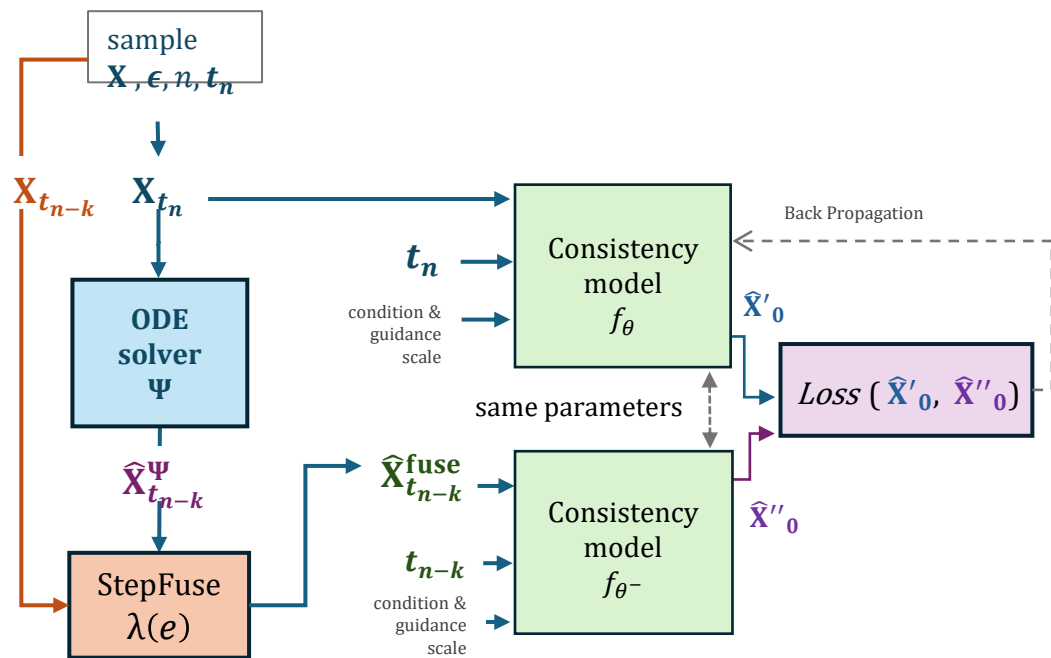# Compositional Consistency Distillation Algo

- Switch

**Algorithm 4.2** Modified Consistency Distillation–Switch

**Input:** dataset $\mathcal{D}$, initial model parameter $\theta$, learning rate $\eta$, ODE solver $\Psi(\cdot,\cdot,\cdot,\cdot)$, distance metric $d(\cdot,\cdot)$, guidance scale $[\omega_{min},\omega_{max}]$, skipping steps $k$, **switching threshold** $e_{\text{switch}}$

1: $\theta^- \leftarrow \theta$
2: **repeat**
3:      Sample $(x, c_1, c_2) \sim \mathcal{D},\ \omega \sim [\omega_{\min}, \omega_{\max}]$
4:      Sample $\epsilon \sim \mathcal{N}(0, I),\ n \sim \mathcal{U}[1+k, N]$
5:      $x_{t_n} \leftarrow \sqrt{\bar{\alpha}_{t_n}} \cdot x + \sqrt{1 - \bar{\alpha}_{t_n}} \cdot \epsilon$
6:      **if** $e < e_{\text{switch}}$ **then**
7:          $x_{t_{n-k}}^{\Psi,\omega} \leftarrow x_{t_n} + (1+\omega) \cdot \Psi(x_{t_n}, t_n, t_{n-k}, c_1, c_2) - \omega \cdot \Psi(x_{t_n}, t_n, t_{n-k}, \varnothing, \varnothing)$
8:          $\mathcal{L}_{\text{Switch}}(\theta, \theta^-; \Psi) \leftarrow d\left( f_\theta(x_{t_n}, t_n, c_1, c_2, \omega),\ f_{\theta^-}(x_{t_{n-k}}^{\Psi,\omega}, t_{n-k}, c_1, c_2, \omega) \right)$
9:      **else**
10:         $x_{t_{n-k}} \leftarrow \sqrt{\bar{\alpha}_{t_{n-k}}} \cdot x + \sqrt{1 - \bar{\alpha}_{t_{n-k}}} \cdot \epsilon$
11:         $\mathcal{L}_{\text{Switch}}(\theta, \theta^-) \leftarrow d\left( f_\theta(x_{t_n}, t_n, c_1, c_2, \omega),\ f_{\theta^-}(x_{t_{n-k}}, t_{n-k}, c_1, c_2, \omega) \right)$
12:      **end if**
13:      $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{\text{Switch}}(\theta, \theta^-)$
14:      $\theta^- \leftarrow \texttt{stopgrad}\left(\mu\theta^- + (1-\mu)\theta\right)$
15: **until** convergence

Diagram labels:
- $X_{t_{n-1}}$
- sample $X, \epsilon, n, t_n$
- Consistency model $f_{\theta^-}$, $t_{n-1}$, $\hat{X}'''_0$
- condition & guidance scale
- same parameters
- $X_{t_n}$, Consistency model $f_\theta$, $t_n$, $\hat{X}'_0$
- $Loss(\hat{X}'_0, \hat{X}'''_0)$
- $Loss(\hat{X}'_0, \hat{X}''_0)$
- Back Propagation, $\lambda(e)$
- ODE solver $\Psi$, $\hat{X}^{\Psi}_{t_{n-k}}$
- Consistency model $f_{\theta^-}$, $t_{n-1}$, $\hat{X}''_0$
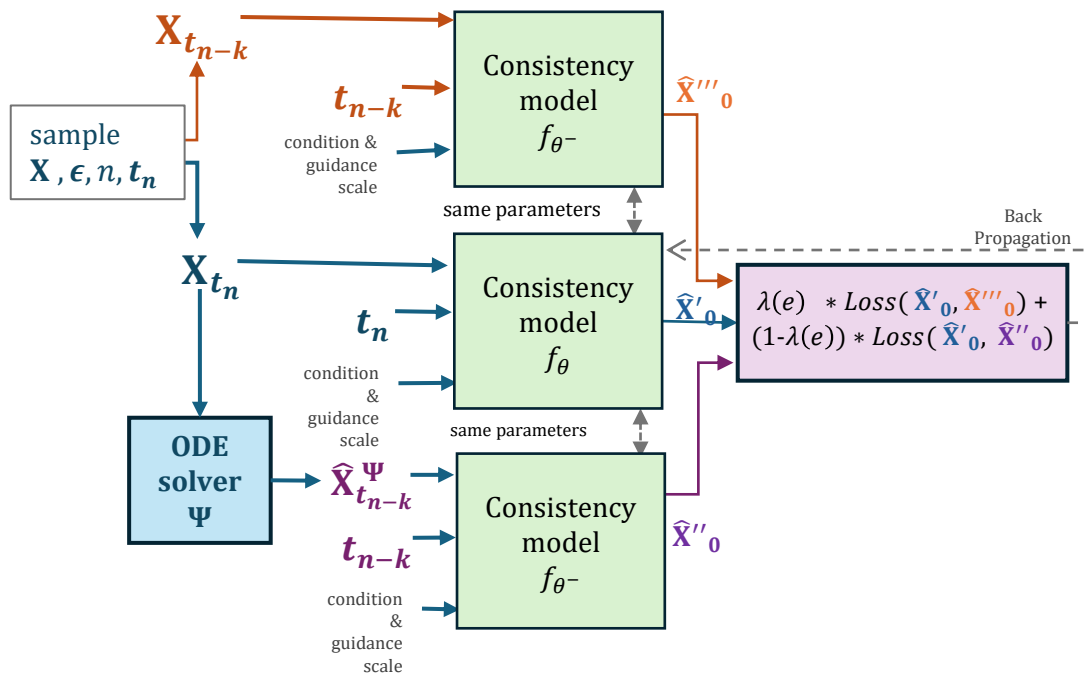
# Compositional Consistency Distillation Algo

- StepFuse



**Algorithm 4.3** Modified Compositional label Consistency Distillation–StepFuse

**Input**: dataset $\mathcal{D}$, initial model parameter $\theta$, learning rate $\eta$, ODE solver $\Psi(\cdot,\cdot,\cdot,\cdot)$, distance metric $d(\cdot,\cdot)$, guidance scale $[\omega_{min},\omega_{max}]$, skipping steps $k$, **fuse scheduler** $\lambda(e)$

1: $\theta^- \leftarrow \theta$
2: **repeat**
3:      Sample $(x, c_1, c_2) \sim \mathcal{D}$, and $\omega \sim [\omega_{min}, \omega_{max}]$
4:      Sample $\epsilon \sim \mathcal{N}(0, I)$, $n \sim \mathcal{U}[1 + k, N]$
5:      $x_{t_n} \leftarrow \sqrt{\bar{\alpha}_{t_n}} \cdot x + \sqrt{1 - \bar{\alpha}_{t_n}} \cdot \epsilon$
6:      $x_{t_{n-k}}^{\Psi,\omega} \leftarrow x_{t_n} + (1 + \omega) \cdot \Psi(x_{t_n}, t_n, t_{n-k}, c_1, c_2) - \omega \cdot \Psi(x_{t_n}, t_n, t_{n-k}, \varnothing, \varnothing)$
7:      $x_{t_{n-k}} \leftarrow \sqrt{\bar{\alpha}_{t_{n-k}}} \cdot x + \sqrt{1 - \bar{\alpha}_{t_{n-k}}} \cdot \epsilon$
8:      $\hat{x}_{t-k}^{\text{fuse}} \leftarrow \lambda(e) \cdot \hat{x}_{t_{n-k}}^{\Psi,\omega} + (1 - \lambda(e)) \cdot x_{t-k}$
9:      $\mathcal{L}_{Stepfuse}(\theta, \theta^-; \Psi) \leftarrow d\left( f_\theta(x_{t_n}, t_n, c_1, c_2, \omega), \ f_{\theta^-}(\hat{x}_{t-k}^{\text{fuse}}, t_{n-k}, c_1, c_2, \omega) \right)$
10:     $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{Stepfuse}(\theta, \theta^-)$
11:     $\theta^- \leftarrow \texttt{stopgrad}\left(\mu\theta^- + (1 - \mu)\theta\right)$
12: **until** convergence

# Compositional Consistency Distillation Algo

- LossFuse



**Algorithm 4.4** Modified Compositional label Consistency Distillation–LossFuse

**Input**: dataset $\mathcal{D}$, initial model parameter $\theta$, learning rate $\eta$, ODE solver $\Psi(\cdot, \cdot, \cdot, \cdot)$, distance metric $d(\cdot, \cdot)$, guidance scale $[\omega_{min}, \omega_{max}]$, skipping steps $k$, **fuse scheduler** $\lambda(e)$

1:   $\theta^- \leftarrow \theta$

2: **repeat**

3:      Sample $(x, c_1, c_2) \sim \mathcal{D}$, and $\omega \sim [\omega_{min}, \omega_{max}]$

4:      Sample $\epsilon \sim \mathcal{N}(0, I)$, $n \sim \mathcal{U}[1 + k, N]$

5:      $x_{t_n} \leftarrow \sqrt{\bar{\alpha}_{t_n}} \cdot x + \sqrt{1 - \bar{\alpha}_{t_n}} \cdot \epsilon$

6:      $x_{t_{n-k}}^{\Psi, \omega} \leftarrow x_{t_n} + (1 + \omega) \cdot \Psi(x_{t_n}, t_n, t_{n-k}, c_1, c_2) - \omega \cdot \Psi(x_{t_n}, t_n, t_{n-k}, \varnothing, \varnothing)$

7:      $x_{t_{n-k}} \leftarrow \sqrt{\bar{\alpha}_{t_{n-k}}} \cdot x + \sqrt{1 - \bar{\alpha}_{t_{n-k}}} \cdot \epsilon$

8:      $\mathcal{L}_{teacher}(\theta, \theta^-; \Psi) \leftarrow d\left( f_\theta(x_{t_n}, t_n, c_1, c_2, \omega), \ f_{\theta^-}(\hat{x}_{t-k}^{\Psi, \omega}, t_{n-k}, c_1, c_2, \omega) \right)$

9:      $\mathcal{L}_{fpf}(\theta, \theta^-) \leftarrow d\left( f_\theta(x_{t_n}, t_n, c_1, c_2, \omega), \ f_{\theta^-}(x_{t-k}, t_{n-k}, c_1, c_2, \omega) \right)$

10:     $\mathcal{L}_{lossfuse} \leftarrow \lambda(e) \cdot \mathcal{L}_{teacher} + (1 - \lambda(e)) \cdot \mathcal{L}_{fpf}$

11:     $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{lossfuse}(\theta, \theta^-)$

12:     $\theta^- \leftarrow \texttt{stopgrad}(\mu \theta^- + (1 - \mu)\theta)$

13: **until** convergence

# Experiment setup

- Training Configuration

Hyperparameters fixed across all experiments

| | |
|---|---|
| Training timesteps | 1000 |
| Batch Size | 24 |
| Noise scheduler $\beta$ | 0.0001-0.2 linear |
| Base channels dimension | 128 |
| # of Residual Blocks | 2 |
| Channel Multiplier | 1,2,4,4 |
| Optimizer | AdamW |
| $\omega$ embedding dimension | 128 |
| Huber loss param | $\delta = 0.001$ |
| EMA decay | $\mu = 0.995$ |