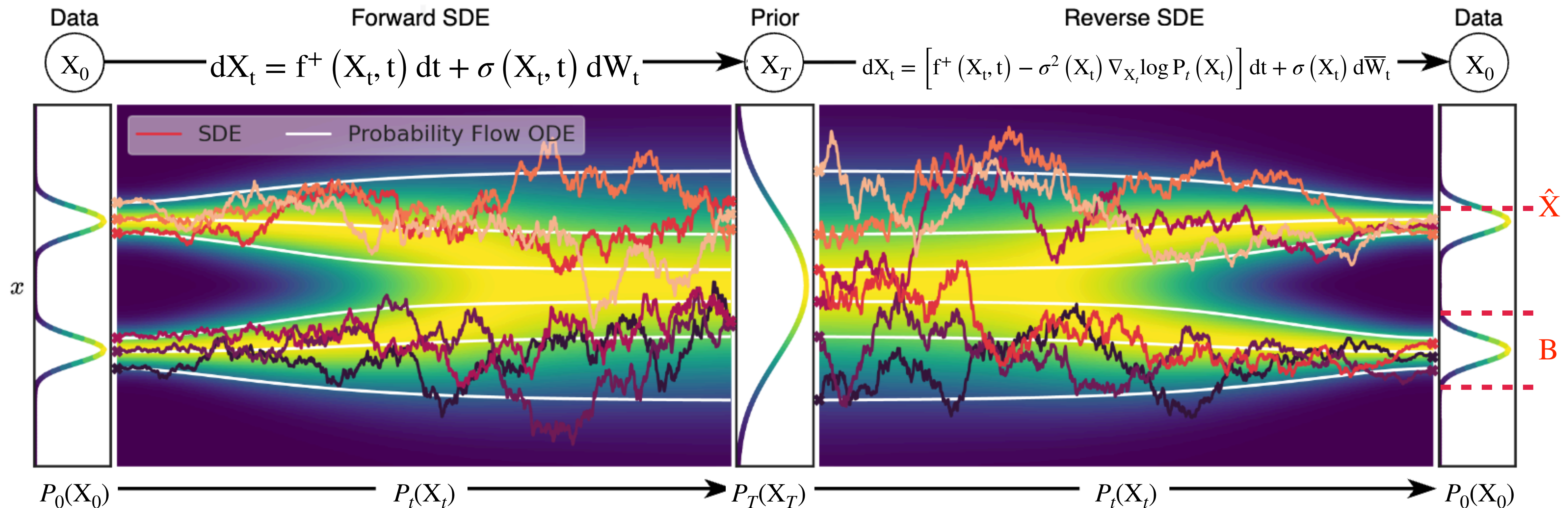# Conditioning in SDEs

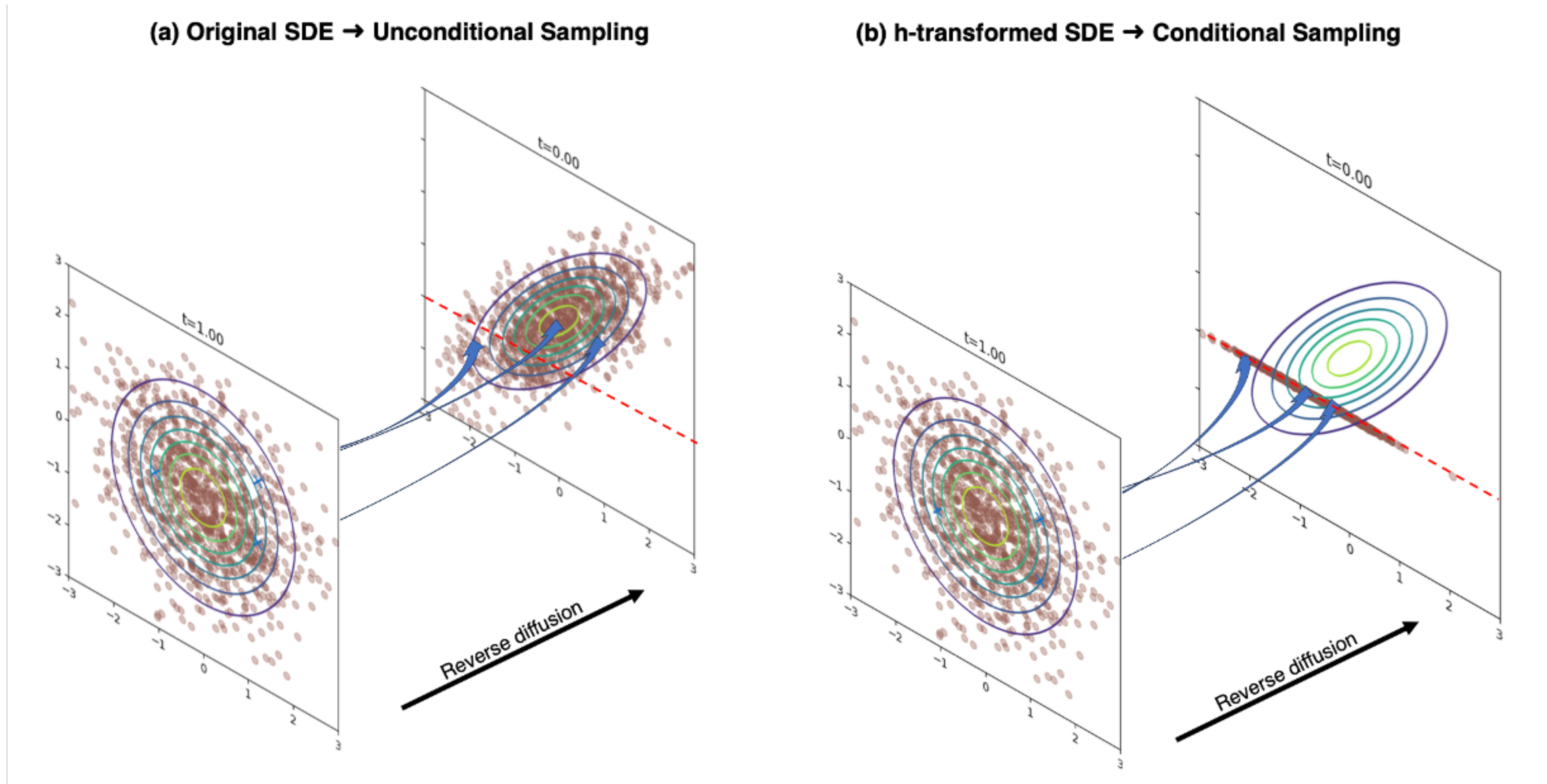**9 February 2024**

**R255 Lecture**

Shreyas Padhy

# So far - Reversal of SDEs

- We have a forward SDE $\qquad dX_t = f^+(X_t, t)\, dt + \sigma(X_t, t)\, dW_t$

- We can reverse the SDE $\qquad dX_t = f^-(X_t, t)\, dt + \sigma(X_t, t)\, d\overline{W}_t$

# Conditioning of SDEs

- We want the reverse SDE to hit a certain constraint set $X_0 \in B$ at time $T = 0$



(a) Original SDE → Unconditional Sampling
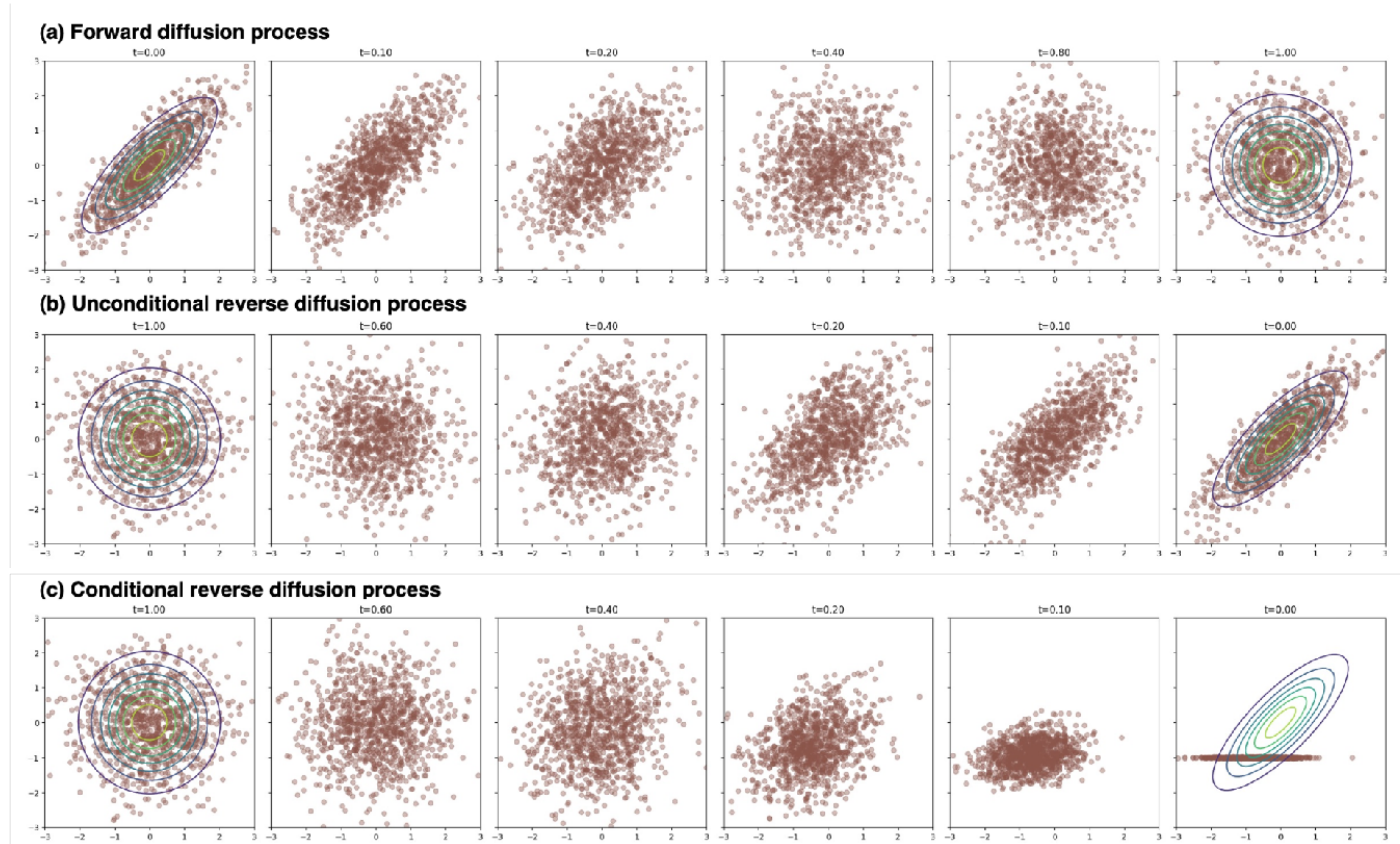
(b) h-transformed SDE → Conditional Sampling

# Conditioning of SDEs

- We want the reverse SDE to hit a certain constraint set $X_0 \in B$ at time $T = 0$
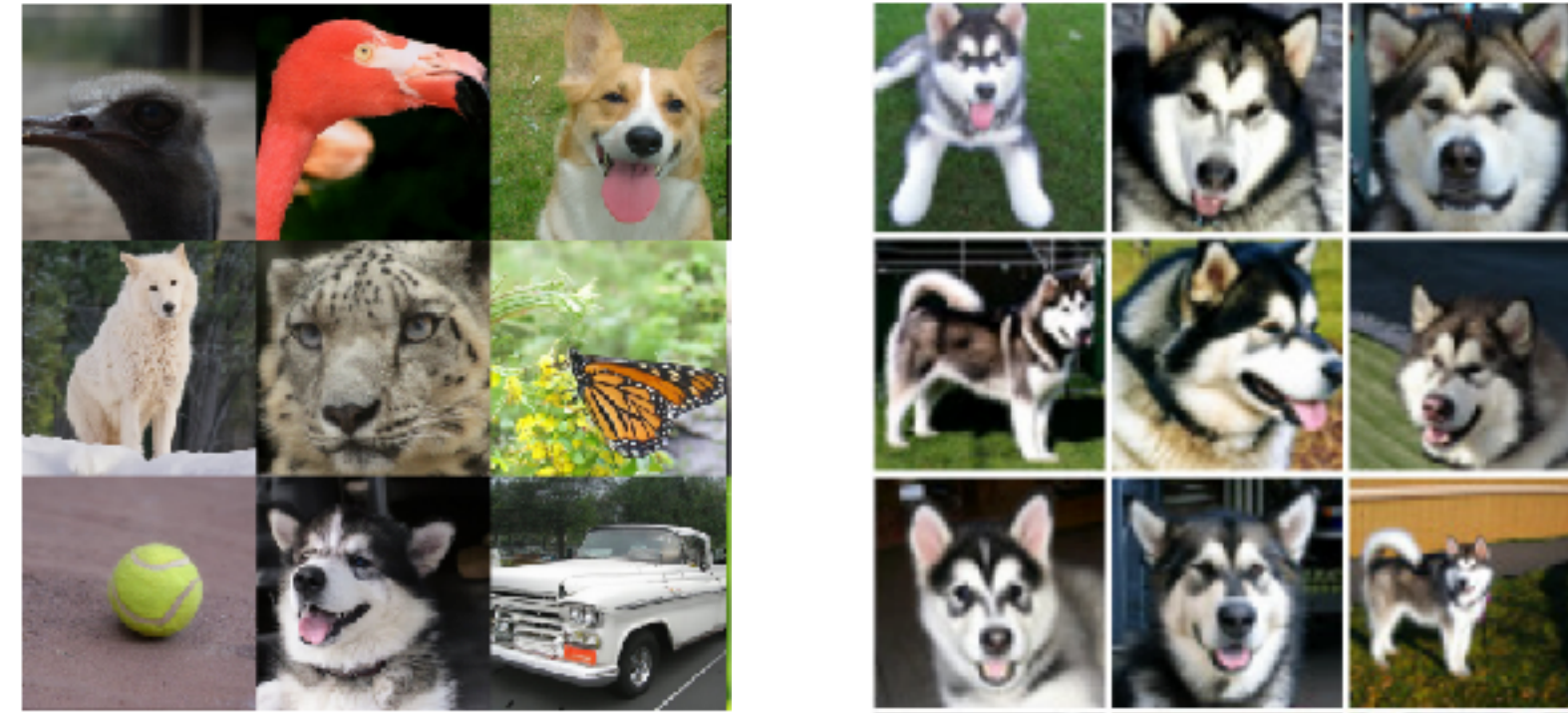
# Conditioning of SDEs

- We want the reverse SDE to hit a certain constraint set $X_0 \in B$ at time $T = 0$

- Applications -

  - Class-conditional sampling



  - Inverse modeling, where we rewrite $X_0 \in B$ to equivalently $\mathscr{A}(X_0) = y$
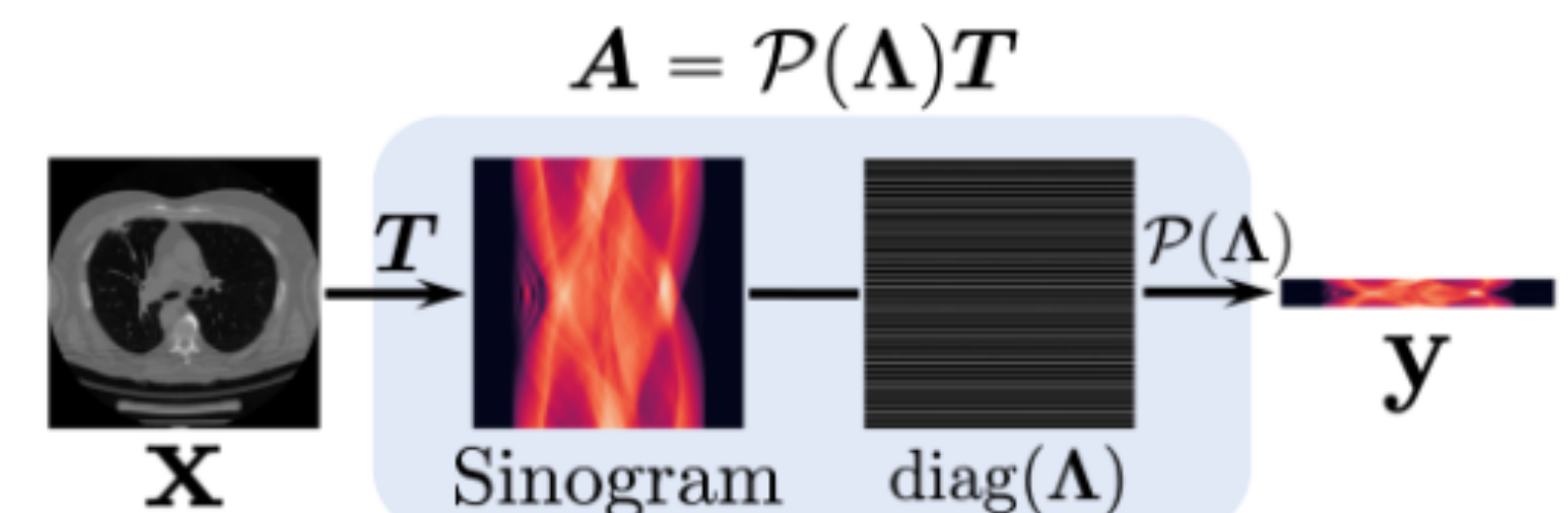


In-painting

Deblurring

Colorisation

$A = \mathcal{P}(\Lambda)T$

Medical Imaging

# Conditioning of SDEs - Doob's $h$-transform

[Rogers and Williams 2000] Consider the forward SDE:

$$dX_t = f^+\left(X_t, t\right) dt + \sigma_t dW_t$$

where time flows forwards and with transition densities $p_{t|s}$. It then follows that the conditioned process $X_t \mid X_T \in B$ is a solution of

$$dH_t = \left(f_t^+\left(H_t, t\right) + \sigma_t^2 \nabla_{H_t} \ln P_{T|t}\left(X_T \in B \mid H_t\right)\right) dt + \sigma_t dW_t,$$

such that $\text{Law}\left(H_t \mid H_s\right) = \vec{p}_{t|s,T}\left(h_t \mid h_s, x_T \in B\right)$ and $\mathbb{P}\left(H_T \in B\right) = 1$.

- By conditioning a diffusion process to hit a particular event $X_T \in B$, we get a resulting conditional process that is also an SDE with an additional drift term.

- $h\left(t, H_t\right) \triangleq P_{T|t}\left(X_T \in B \mid H_t\right)$ is known as the $h$-**transform**

# Aside - FPK and Infinitesimal Generators

- The FPK, a PDE that describes the evolution of the probability density of an SDE

  - $dX_t = \mu\left(X_t, t\right) dt + \sigma\left(X_t, t\right) dW_t$

  - $\dfrac{\partial}{\partial t} p(x, t) = -\dfrac{\partial}{\partial x}[\mu(x, t)p(x, t)] + \dfrac{\partial^2}{\partial x^2}[D(x, t)p(x, t)]$

- Let us define an operator called the (infinitesimal) generator of a stochastic process

  - $\mathscr{A}\phi(x_t, t) = \lim\limits_{s \to 0} \dfrac{\mathbb{E}\left[\phi(x_{t+s}, t + s)\right] - \phi(x_t, t)}{s}$

- Why is it called the generator?

  - $\mathbb{E}\left[\phi(x_{t+s}, t + s)\right] \simeq \phi(x_t, t) + s\mathscr{A}\phi(x_t, t)$

# Aside - FPK and Infinitesimal Generators

$$\mathscr{A}\phi(x_t, t) = \lim_{s \to 0} \frac{\mathbb{E}\left[\phi(x_{t+s}, t+s)\right] - \phi(x_t, t)}{s}$$

- For the Ito process $dX_t = \mu\left(X_t, t\right) dt + \sigma\left(X_t, t\right) dW_t$, the generator is given by

$$\mathscr{A}(\bullet) = \sum_i \frac{\partial(\bullet)}{\partial x_i} f_i(x_t, t) + \frac{1}{2} \sum_{i,j} \left(\frac{\partial^2(\bullet)}{\partial x_i \partial x_j}\right) \left[\sigma^2(X_t, t)\right]_{ij}$$

- We can show that the transition densities of the SDE satisfy the backwards Kolmogorov equation defined by the generator acting on the transition density

$$-\frac{\partial p(x_{t+s}, t+s \mid x_t, t)}{\partial t} = \mathscr{A}p(x_{t+s}, t+s \mid x_t, t)$$

# Doob's *h*-transform (Quick Proof)

- Given the SDE with transition density $p_{t|s}(x_t \mid x_s)$, where $t > s$,

- We wish to obtain the process after conditioning the SDE to hit a deterministic end point $z$ at time $T$,

$$p_{t|s,T}\left(x_t \mid x_s, x_T = z\right) = \frac{p_{t|s,T}\left(x_T = z \mid x_t\right) p_{t|s}\left(x_t \mid x_s\right)}{p\left(x_T = z \mid x_s\right)}$$

- Let us try to find the SDE that would have this transition density

# Doob's *h*-transform (Quick Proof)

- Apply Bayes' rule

$$p_{t|s,T}\left(x_t \mid x_s, x_T = z\right) = \frac{p_{T|s,t}\left(x_T = z \mid x_s, x_t\right) p_{t|s}\left(x_t \mid x_s\right)}{p_{T|s}\left(x_T = z \mid x_s\right)}$$

- Now, apply the Markov property $p_{T|s,t}\left(x_T = z \mid x_s, x_t\right) \rightarrow p_{T|s,t}\left(x_T = z \mid x_t\right)$,

$$p_{t|s,T}\left(x_t \mid x_s, x_T = z\right) = \frac{p_{T|t}\left(x_T = z \mid x_t\right) p_{t|s}\left(x_t \mid x_s\right)}{p_{T|s}\left(x_T = z \mid x_s\right)}$$

- Now, define $h(x_t, t) = p_{T|t}(x_T = z \mid x_t)$, and replace above to get

$$p_{t|s,T}\left(x_t \mid x_s, x_T = z\right) = \frac{h(x_t, t) p_{t|s}\left(x_t \mid x_s\right)}{h(x_s, s)}$$

# Doob's $h$-transform (Quick Proof)

- In order for this to be a valid Markov kernel, we require

$$\int_{x_t} p_{t|s,T}\left(x_t \mid x_s, x_T = z\right) dx_t = 1$$

$$\int_{x_t} \frac{h(x_t, t)p_{t|s}\left(x_t \mid x_s\right)}{h(x_s, s)} dx_t = 1$$

- Taking $h(x_s, s)$ to the RHS (independent of $x_t$), we get the following property that $h(x_s, s)$ satisfies

$$\bullet \quad h(x_s, s) = \int_{x_t} h(x_t, t)p_{t|s}\left(x_t \mid x_s\right) = \mathbb{E}\left[h(x_t, t)\right]$$

- From the definition for infinitesimal generator $\mathscr{A}$, we have (assuming $(t, t+s)$ instead of $(s, t)$)

$$\bullet \quad \mathscr{A}\phi(\mathbf{x}, t) = \lim_{s\downarrow 0} \frac{\mathrm{E}[\phi(\mathbf{x}(t+s), t+s)] - \phi(\mathbf{x}(t), t)}{s}$$

- Therefore, $\mathscr{A}h(x_t, t) = 0$

# Doob's $h$-transform (Quick Proof)

- $\mathscr{A}h(x_t, t) = 0$, and the new transition probability $p^h = p\dfrac{h(x_{t+s}, t+s)}{h(x_t, t)}$

- Writing down the infinitesimal generator for some arbitrary $\phi$, $\mathbb{E}$ w.r.t $p^h$, we get

$$\mathscr{A}\phi = \lim_{s\to 0} \frac{\mathbb{E}^h[\phi(x_{t+s}, t+s)] - \phi(x_t, t)}{s}$$

- Noting that $\mathbb{E}^h[\phi] = \mathbb{E}\left[\phi\dfrac{h(x_{t+s}, t+s)}{h(x_t, t)}\right]$, we have

$$\mathscr{A}^h\phi = \lim_{s\downarrow 0} \frac{\mathrm{E}[\phi(x_{t+s})h(x_{t+s}, t+s)] - \phi(x_t)h(x_t, t)}{sh(x_t, t)}$$

$$= \frac{1}{h(x_t, t)}\mathscr{A}\{h(x_t, t))\phi(x_t)\}$$

- Now, remember the form of $\mathscr{A}$ below, and apply the product rule to $h(\,.\,)\phi(\,.\,)$

$$\mathscr{A}(\bullet) = \sum_i \frac{\partial(\bullet)}{\partial x_i}f_i(x_t, t) + \frac{1}{2}\sum_{i,j}\left(\frac{\partial^2(\bullet)}{\partial x_i\partial x_j}\right)\left[\sigma^2(X_t, t)\right]_{ij}$$

# Doob's *h*-transform (Quick Proof)

- Now, remember the form of $\mathcal{A}$ below, and apply the product rule to $h(\,.\,)\phi(\,.\,)$

$$\mathcal{A}(\,\bullet\,) = \sum_i \frac{\partial(\bullet)}{\partial x_i} f_i(x_t, t) + \frac{1}{2}\sum_{i,j}\left(\frac{\partial^2(\bullet)}{\partial x_i \partial x_j}\right)\left[\sigma^2(X_t, t)\right]_{ij}$$

- We get this behemoth

$$\mathcal{A}^h\phi = \frac{1}{h(t, x_t)}\left\{ \frac{\partial h(t, x_t)}{\partial t}\phi + \sum_i\left[\frac{\partial h(t, x_t)}{\partial x_i}\phi(x_t) + h(t, x_t)\frac{\partial \phi(x_t)}{\partial x_i}\right]f_i(x_t, t) \right.$$
$$\left. + \frac{1}{2}\sum_{i,j}\frac{\partial^2[h(t, x_t)\phi(x_t)]}{\partial x_i \partial x_j}\left[\sigma^2(x_t, t)\right]_{ij} \right\}$$

- This gives us an even bigger behemoth $\qquad \color{red}{\mathcal{A}h = 0}$

$$= \frac{1}{h(t, x_t)}\left\{ \color{red}{\left[\frac{\partial h(t, x_t)}{\partial t} + \sum_i \frac{\partial h(t, x_t)}{\partial x_i}f_i(x_t, t) + \frac{1}{2}\sum_{i,j}\frac{\partial^2 h(t, x_t)}{\partial x_i \partial x_j}\left[\sigma^2(x_t, t)\right]_{ij}\right]}\phi(x_t) + \sum_i h(t, x_t)\frac{\partial \phi(x_t)}{\partial x_i}f_i(x_t, t) + \right.$$

- $$\left. \frac{1}{2}\sum_{i,j}\left[\frac{\partial h(t, x_t)}{\partial x_j}\frac{\partial \phi(x_t)}{\partial x_i} + \frac{\partial h(t, x_t)}{\partial x_i}\frac{\partial \phi(x_t)}{\partial x_j} + h(t, x_t)\frac{\partial^2 \phi(x_t)}{\partial x_i \partial x_j}\right]\left[\sigma^2(x_t, t)\right]_{ij} \right\}$$

# Doob's $h$-transform (Quick Proof)

- Finally, arranging terms we get

$$= \sum_i \left[ f_i(x_t, t) + \sigma^2(x_t, t) \frac{\nabla h(t, x_t)}{h(t, x_t)} \right] \frac{\partial \phi(x_t)}{\partial x_i} + \frac{1}{2} \sum_{i,j} \frac{\partial^2 \phi(x_t)}{\partial x_i \partial x_j} \left[ \sigma^2(x_t, t) \right]_{ij}$$

- Comparing to $\mathscr{A}(\bullet) = \sum_i \frac{\partial(\bullet)}{\partial x_i} f_i(x_t, t) + \frac{1}{2} \sum_{i,j} \left( \frac{\partial^2(\bullet)}{\partial x_i \partial x_j} \right) \left[ \sigma^2(X_t, t) \right]_{ij}$

- The new drift is $f_i(x_t, t) + \sigma^2(x_t, t) \frac{\nabla h(t, x_t)}{h(t, x_t)} = f_i(x_t, t) + \sigma^2(x_i, t) \nabla \log h(x_t, t)$

# Example: Pinned Brownian Motion

- Consider a Brownian motion starting from arbitrary $X_0 \sim p_{\text{data}}$, and $\mathrm{d}X_t = \sigma \mathrm{d}W_t$

- Let us condition this Brownian motion to hit $X_T = 0$ at time T,

  - Therefore $h(x_t, t) = P(X_T = 0 \mid X_t) = \mathcal{N}\left(X_t, \sigma(T-t)\right)$

  - Which gives us $\nabla_{X_t} \log h(X_t, t) = \nabla_{X_t}\left( C - \dfrac{(X_t - 0)^2}{2\sigma^2(T-t)} \right) = -\dfrac{X_t}{\sigma^2(T-t)}$

- Plugging into the conditioned process formula, we get

  - $\mathrm{d}\mathrm{H}_t = \left( f_t^+\left(\mathrm{H}_t, t\right) + \sigma_t^2 \nabla_{\mathrm{H}_t} \ln P_{T|t}\left(X_T \in B \mid \mathrm{H}_t\right) \right) \mathrm{d}t + \sigma_t \mathrm{d}\mathrm{W}_t,$

  - $\mathrm{d}\mathrm{H}_t = -\dfrac{X_t}{T - t} + \sigma \mathrm{d}\mathrm{W}_t$

# Doob's *h*-transform on the reverse SDE

[Rogers and Williams 2000] Consider the reverse SDE:

$$\mathrm{d}X_t = f^- \left(X_t, t\right) \mathrm{d}t + \sigma_t \mathrm{d}\overline{W}_t$$

where time flows backwards and with transition densities $\bar{p}_{t|s}$. It then follows that the conditioned process $X_t \mid X_0 \in B$ is a solution of

$$\mathrm{d}H_t = \left(f_t^- \left(H_t, t\right) - \sigma_t^2 \nabla_{H_t} \ln \bar{P}_{0|t} \left(X_0 \in B \mid H_t\right)\right) \mathrm{d}t + \sigma_t \bar{\mathrm{d}}_t \bar{W}_t,$$

such that $\mathrm{Law}\left(H_s \mid H_t\right) = \vec{p}_{s|t,0}\left(h_s \mid h_t, x_0 \in B\right)$ and $\mathbb{P}\left(H_0 \in B\right) = 1$.

# Specifying hard constraints

- We consider events of the form $X_0 \in B$ that can be described by

    - An equality constraint $\mathscr{A}(X_0) = y$

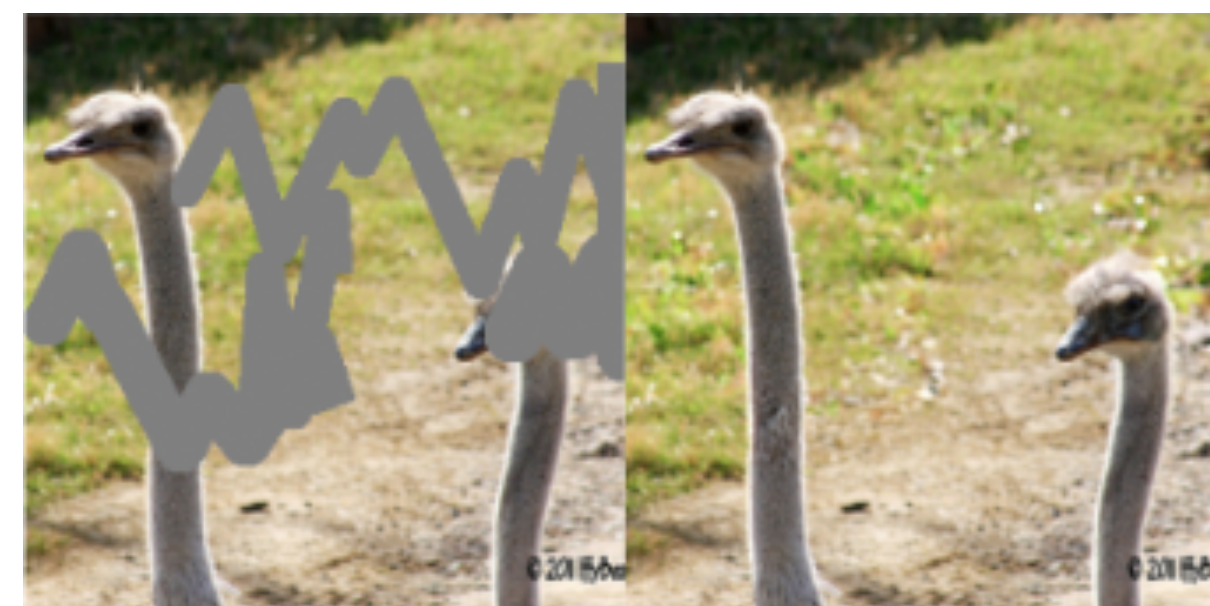- Consider Doob's $h$-transform with hard constraints

    - $\mathrm{d}\mathrm{H}_t = \left( f_t^- \left( \mathrm{H}_t, t \right) - \sigma_t^2 \, \nabla_{\mathrm{H}_t} \ln \bar{P}_{0|t} \left( \mathscr{A}(X_0) = y \mid \mathrm{H}_t \right) \right) \mathrm{d}t + \sigma_t \bar{\mathrm{d}}_t \bar{\mathrm{W}}_t,$

- Replacing $f_t^-(H_t, t) = f^+(H_t, t) - \sigma^2 \nabla_{H_t} \ln P_t(H_t)$, we get

    - $\mathrm{d}\mathrm{H}_t = f_t^+ \left( \mathrm{H}_t, t \right) - \sigma_t^2 \left( \boxed{\nabla_{H_t} \ln P_t(H_t)} + \boxed{\nabla_{\mathrm{H}_t} \ln \bar{P}_{0|t} \left( \mathscr{A}(X_0) = y \mid \mathrm{H}_t \right)} \right) \mathrm{d}t + \sigma_t \bar{\mathrm{d}}_t \bar{\mathrm{W}}_t,$

Known
unconditional
score

- Sampling from this SDE gives us $x \sim p_{\mathrm{data}}$ that satisfy $\mathscr{A}(x) = y$
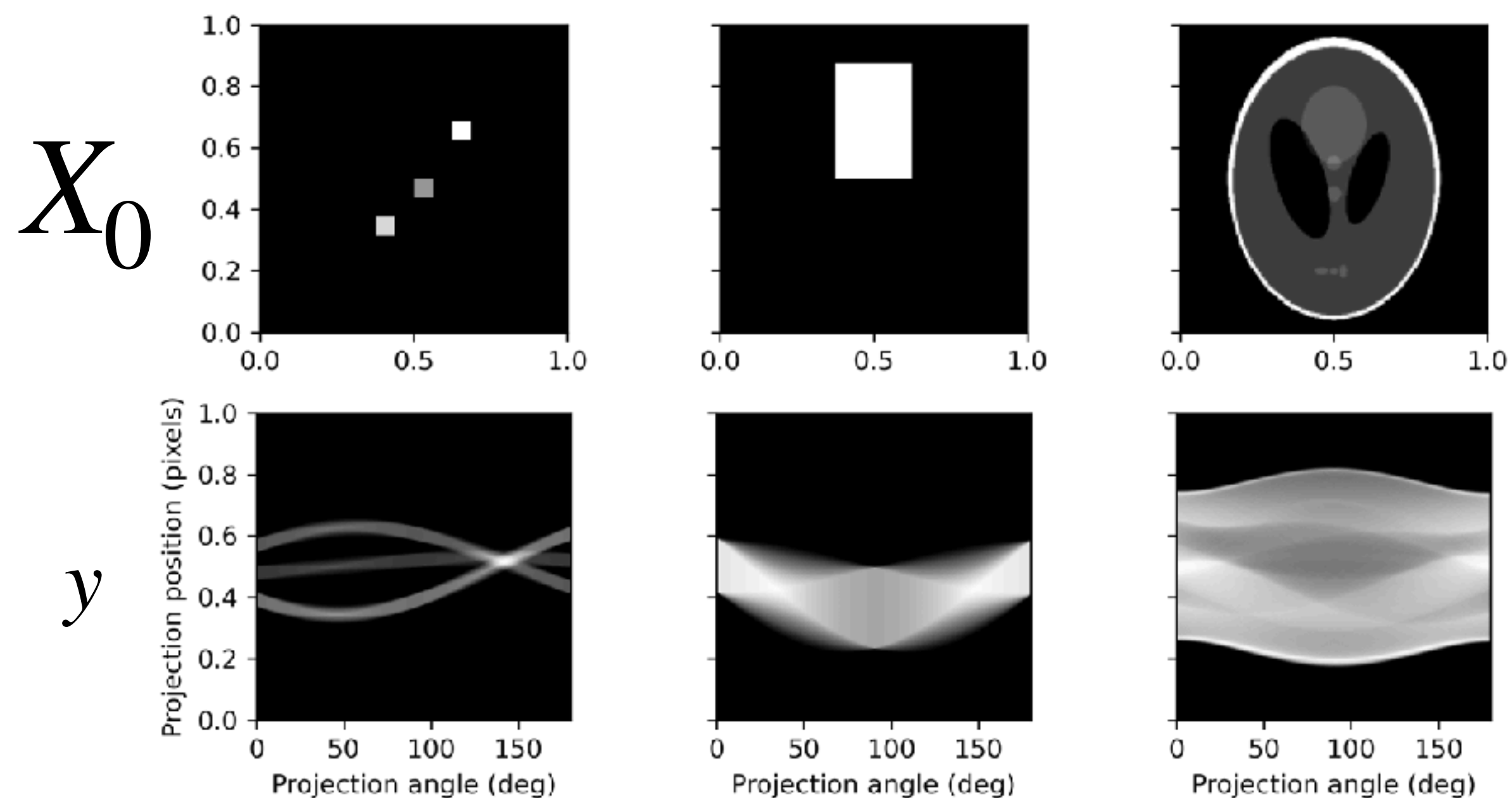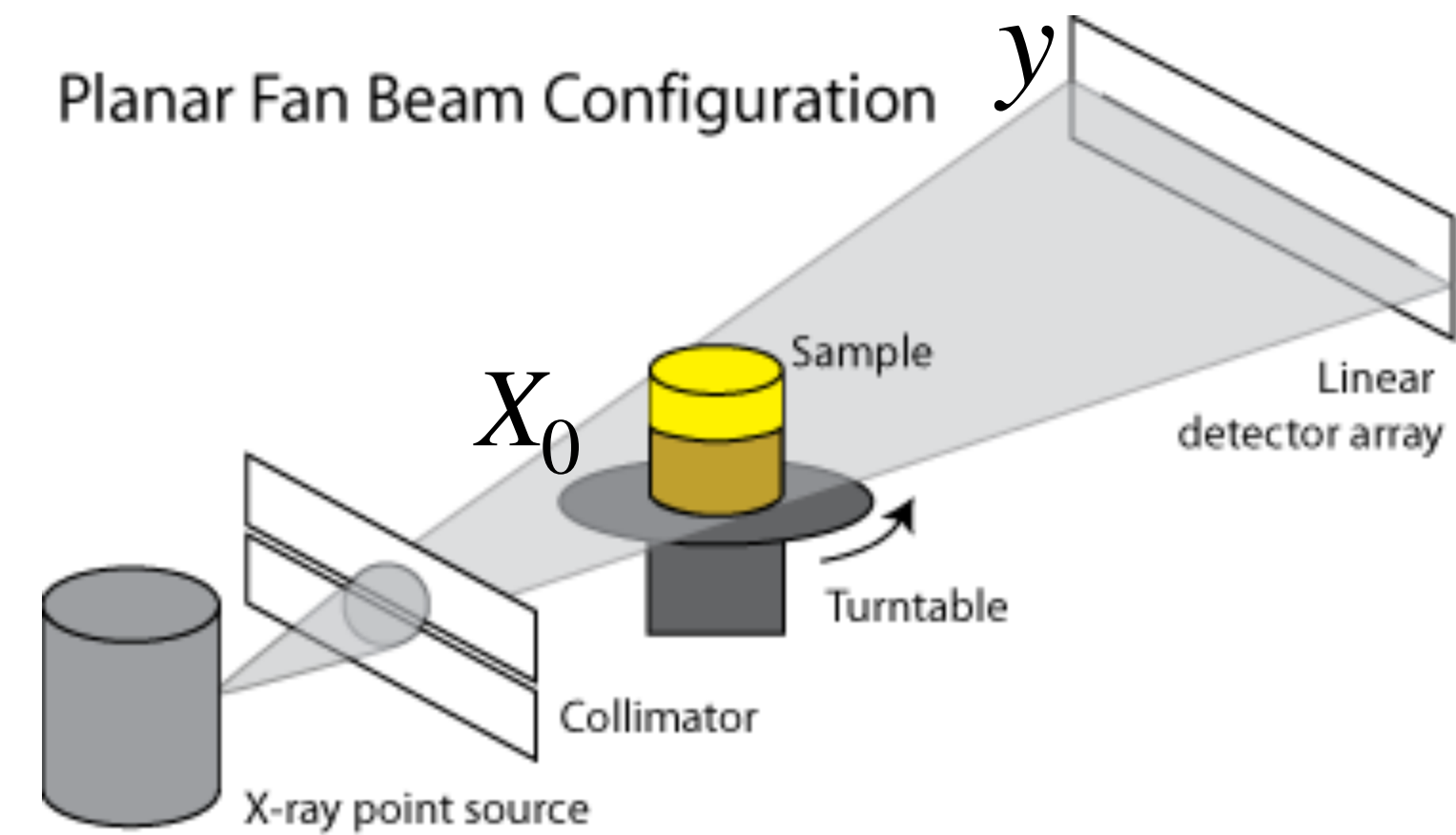
$$\mathscr{A} = \{0,1\}^{H \times W \times C}$$

Unknown
conditional
score

In-painting

# Specifying soft constraints

- We consider events of the form $X_0 \in B$ that can be described by

  - Noisy observations $y = \mathscr{A}(X_0) + \eta$, and a density $p(y \mid X_0)$

    - We wish to recover the posterior $p(X_0 \mid y)$

# Specifying soft constraints

- We consider events of the form $X_0 \in B$ that can be described by

  - Noisy observations $y = \mathcal{A}(X_0) + \eta$, and a density $p(y \mid X_0)$

    - We wish to recover the posterior $p(X_0 \mid y)$
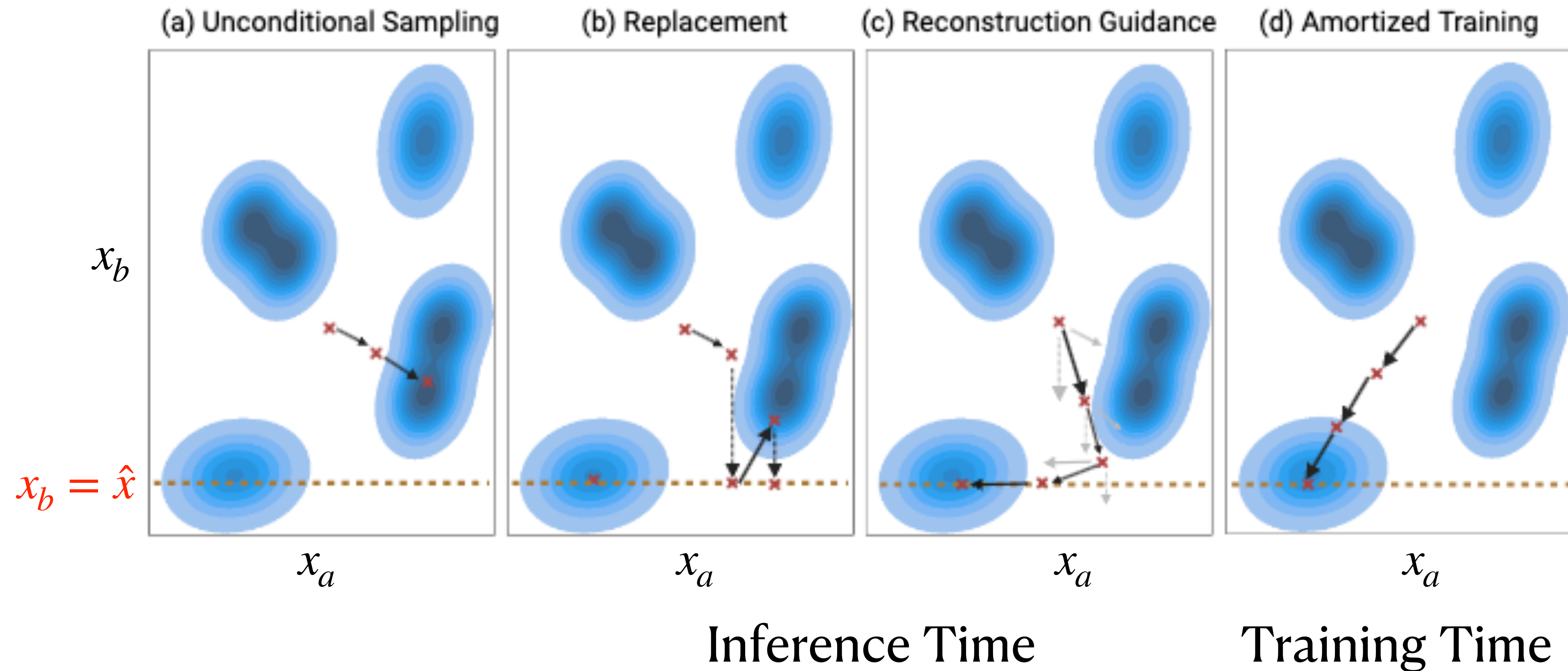
- Consider Doob's $h$-transform with soft constraints

  - $d\mathrm{H}_t = \left( f_t^- \left( \mathrm{H}_t, t \right) - \sigma_t^2 \nabla_{\mathrm{H}_t} \ln \bar{P}_{y|t} \left( Y = y \mid \mathrm{H}_t \right) \right) dt + \sigma_t \bar{d}_t \bar{\mathrm{W}}_t,$

- Again, replacing $f_t^-(H_t, t) = f_t^+(H_t, t) - \sigma^2 \nabla_{H_t} \ln P_t(H_t)$, we get

  - $d\bar{\mathrm{H}}_t = f_t^+ \left( \mathrm{H}_t, t \right) - \sigma_t^2 \left( \nabla_{H_t} \ln P_t(H_t) + \nabla_{\mathrm{H}_t} \ln \bar{P}_{y|t} \left( Y = y \mid \mathrm{H}_t \right) \right) dt + \sigma_t \bar{d}_t \bar{\mathrm{W}}_t,$

- Sampling from this SDE gives us $x \sim p(x_0 \mid Y = y) = \dfrac{p(y \mid x_0) p_{\mathrm{data}}(x_0)}{p(y)}$

# Different approaches for conditioning



(a) Unconditional Sampling  (b) Replacement  (c) Reconstruction Guidance  (d) Amortized Training

$x_b$

$x_b = \hat{x}$

$x_a$  $x_a$  $x_a$  $x_a$

Inference Time  Training Time

# How do we sample from conditioned SDE?

- We have
$$\mathrm{d}H_t = f_t^+ \left( H_t, t \right) - \sigma_t^2 \left( \nabla_{H_t} \ln P_t(H_t) + \textcolor{red}{\nabla_{H_t} \ln \bar{P}_{y|t} \left( Y = y \mid H_t \right)} \right) \mathrm{d}t + \sigma_t \bar{\mathrm{d}}_t \bar{W}_t,$$

- The challenge is $\bar{P}_{y|t} \left( Y = y \mid H_t \right) = \int p(y \mid x_0) \textcolor{red}{\bar{p}_{0|t} \left( x_0 \mid H_t \right)} \, \mathrm{d}x_0$

- We can sample from the reverse SDE to obtain samples from $p(x_0 \mid H_t)$

- We cannot estimate the integral without *many* samples, or evaluate at fixed $x_0$

- Can we approximate $p(x_0 \mid H_t)$?

# Reconstruction Guidance

- Approximate $p(x_0 \mid x_t)$ with a Gaussian approximation using **Tweedie's formula**

[**Tweedie's formula**]. Let $p(a \mid b) = p_0(a)\exp(b^\top T(a) - \psi(b))$. Then the posterior mean $\hat{b} = \mathbb{E}[b \mid a]$ should satisfy

$$\left( \nabla_a T(a) \right)^\top \hat{b} = \nabla_a \log p(a) - \nabla_a \log p_0(a)$$

- **Diffusion Posterior Sampling** [Chung et al 2023]: Forward process is VP-SDE or DDPM sampling, approximate posterior mean using Tweedie's formula:

$$\hat{x}_0 = \mathbb{E}[x_0 \mid x_t] = \frac{1}{\sqrt{\bar{\alpha}(t)}} \left( x_t + (1 - \bar{\alpha}(t)) \nabla_{x_t} \log p_t(x_t) \right)$$
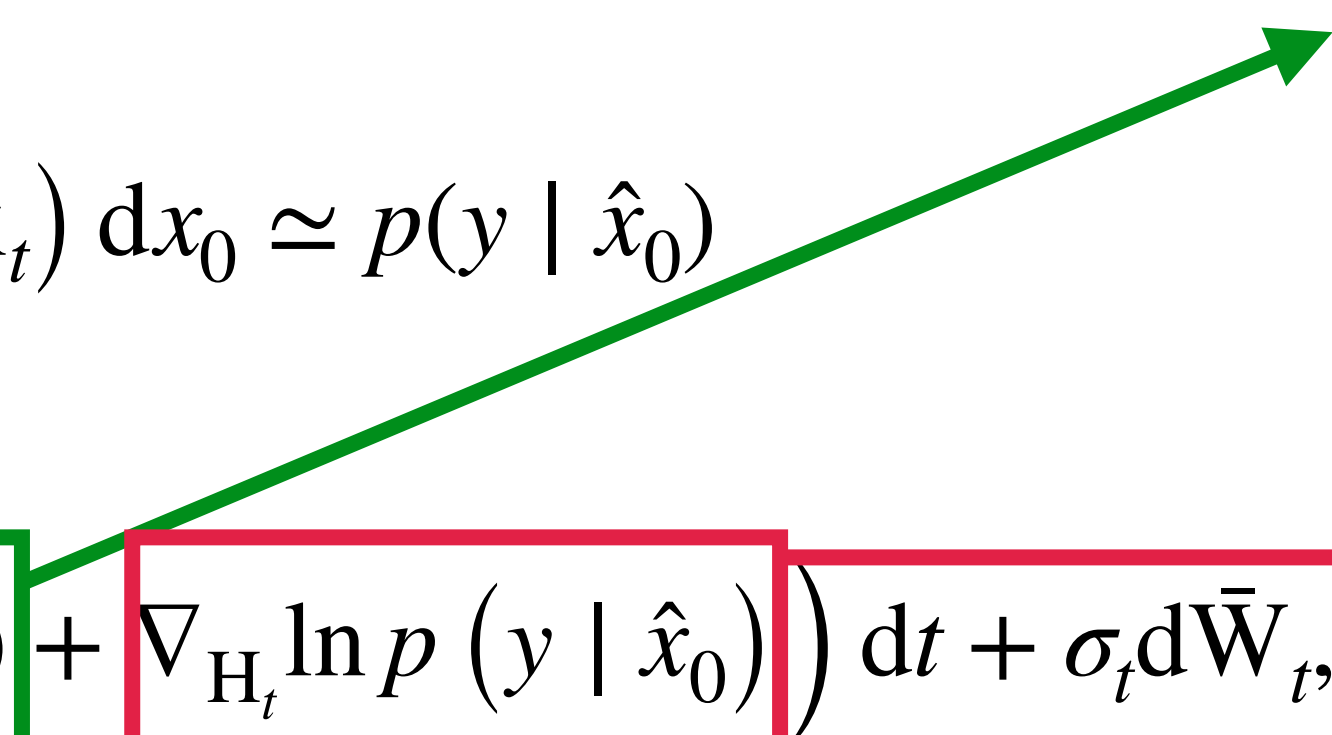
- Then approximate

$$\bar{P}_{y|t} \left( Y = y \mid x_t \right) = \int p(y \mid x_0) \bar{p}_{0|t} \left( x_0 \mid x_t \right) dx_0 \simeq p(y \mid \hat{x}_0)$$

Known unconditional score

- Simulate the following reverse SDE

  - $dH_t = f_t^+ \left( H_t, t \right) - \sigma_t^2 \left( \boxed{\nabla_{H_t} \ln P_t(H_t)} + \boxed{\nabla_{H_t} \ln p \left( y \mid \hat{x}_0 \right)} \right) dt + \sigma_t d\bar{W}_t,$

Extra backprop through score network

# Classifier Guidance

- When we have external labels to guide diffusion

- Consider $dX_t = f_t^+ \left( X_t, t \right) - \sigma_t^2 \left( \nabla_{X_t} \ln P_t(X_t) + \nabla_{X_t} \ln P \left( Y = y \mid X_t \right) \right) dt + \sigma_t d\bar{W}_t,$

  - We can train a time-dependent classifier by creating training data

    - Sample $(x_0, y)$ from dataset, then sample $x_t \sim$ SDE to obtain $(x_t, y)$

    - Using trained $p(y \mid x_t)$ classifier, we backprop to get $\nabla_{x_t} \ln p(y \mid x_t)$

  - $\nabla_{x_t} \ln p(x_t)$ is unconditional score model

  - We can overweight "guidance", $\nabla_{X_t} \ln P_t(X_t) + \gamma \nabla_{X_t} \ln P \left( y \mid X_t \right)$

# Classifier Free Guidance

- What if we don't want to train a separate classifier model?

- Consider $\mathrm{dX}_t = f_t^+\left(\mathrm{X}_t, t\right) - \sigma_t^2\left(\nabla_{X_t}\ln P_t(X_t) + \nabla_{X_t}\ln P\left(Y = y \mid X_t\right)\right)\mathrm{d}t + \sigma_t\mathrm{d}\bar{\mathrm{W}}_t,$

  - From Bayes' rule, $\nabla_{X_t}\ln p(X_t \mid y) = \nabla_{X_t}\ln p(X_t) + \nabla_{X_t}\ln p(y \mid X_t)$

  - Train a single model to approximate both $p(X_t \mid y)$ and $p(X_t)$ on $(X_0, y)$ inputs

  - Occasionally, 10-20% of the time, drop $y$ to learn $p(X_t)$, otherwise $p(X_t \mid y)$

  - At inference time, $\nabla_x \log p_\gamma(x \mid y) = \nabla_x \log p(x) + \gamma\left(\nabla_x \log p(x \mid y) - \nabla_x \log p(x)\right)$

# Amortised training of *h*-transform

- Can we learn a score model for the *h*-transform directly?

- Consider $d\mathrm{X}_t = f_t^+ \left( \mathrm{X}_t, t \right) - \sigma_t^2 \left( \nabla_{X_t} \ln P(X_t \mid y) \right) dt + \sigma_t d\bar{\mathrm{W}}_t,$

- Let us learn a score model $f(t, X_t, y, \mathscr{A})$ to approximate $\nabla_{X_t} \ln p(X_t \mid y)$

$$f* = \arg\min_{f} \mathbb{E}_{Y \sim p_{|\mathscr{A}, X_0}, \mathscr{A} \sim p, X_0 \sim p_{\mathrm{data}}} \left[ \int_0^T \left\| f\left(t, X_t, y, \mathscr{A}\right) - \nabla_{X_t} \ln \vec{p}_{t|0}\left(X_t \mid X_0\right) \right\|^2 dt \right]$$

- Sample $X_0 \sim p_{\mathrm{data}}$, sample an operator $\mathscr{A}$, sample $y \sim p(y \mid X_0, \mathscr{A})$

- We "amortise" the score model over $y$ and $\mathscr{A}$

- The minimiser is given by the conditional score

$$f* \left( t, \mathrm{X}_t, y, \mathrm{A} \right) = \nabla_{x_t} \ln p_{t|0} \left( \mathrm{X}_t \mid Y = y, \mathscr{A} = \mathrm{A} \right)$$

# Proof of minimiser of loss function

$$f\left(x_t, y, A\right) = \mathbb{E}_{X_0 | X_t = x_t, Y = y, \mathscr{A} = A}\left[\nabla_{X_t} \ln p_{t|0}\left(X_t \mid X_0\right)\right]$$

$$= \int \nabla_{x_t} \ln p_{t|0}\left(x_t \mid X_0\right) \, p_{0|t}\left(X_0 \mid x_t, y, A\right) dX_0 \qquad \text{(Write expectation as integral)}$$

$$= \int \nabla_{x_t} \ln p_{t|0}\left(x_t \mid X_0\right) \frac{p_{t|0}\left(x_t \mid X_0\right) p_0\left(X_0\right)}{p_t\left(x_t \mid y, A\right)} dX_0 \qquad \text{(Bayes' rule)}$$

$$= \int \frac{\nabla_{x_t} p_{t|0}\left(x_t \mid X_0\right)}{p_{t|0}\left(x_t \mid X_0\right)} \frac{p_{t|0}\left(x_t \mid X_0\right) p_0\left(X_0\right)}{p_t\left(x_t \mid y, A\right)} dX_0 \qquad \text{(Expand the gradient of log)}$$

$$= \frac{1}{p_t\left(x_t \mid y, A\right)} \int \frac{\nabla_{x_t} p_{t|0}\left(x_t \mid X_0\right)}{p_{t|0}\left(x_t \mid X_0\right)} p_{t|0}\left(x_t \mid X_0\right) p_0\left(X_0\right) dX_0 \qquad \text{(Cancel terms)}$$

$$= \frac{1}{p_t\left(x_t \mid y, A\right)} \int \nabla_{x_t}\left[p_{t|0}\left(x_t \mid X_0\right)\right] p_t\left(x_t \mid y, A\right) dX_0 \qquad \text{(Swap integral derivative - Dominated Convergence Theorem)}$$

$$= \frac{1}{p_t\left(x_t \mid y, A\right)} \nabla_{x_t} \int \left[p_{t|0}\left(x_t \mid x_0\right)\right] p_t\left(x_t \mid y, A\right) dX_0 \qquad \text{(Marginalise over } X_0\text{)}$$

$$= \frac{1}{p_t\left(x_t \mid y, A\right)} \nabla_{x_t} p_t\left(x_t \mid y, A\right)$$

$$= \nabla_{x_t} \ln p_t\left(x_t \mid y, A\right)$$