

# From Entropy to Artistry: on Thermodynamics and Generative AI

Stephan Mandt  
Department of Computer Science  
University of California, Irvine

# A physicist's perspective on diffusion models

## Part 1:

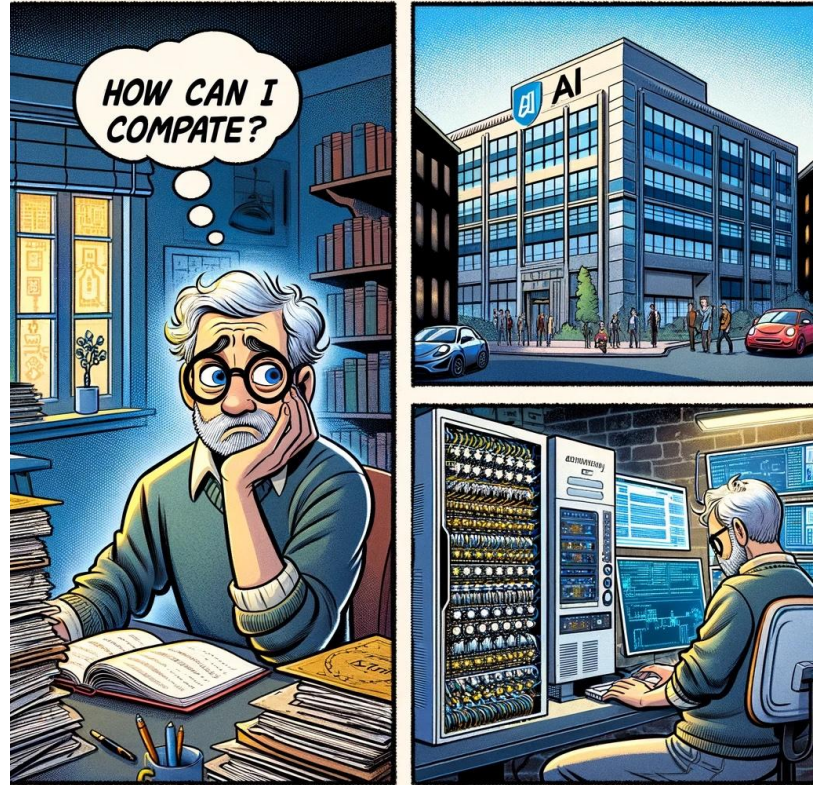
- Understanding the design space of diffusion models
- Efficient samplers for accelerating diffusion models
- Data communication with diffusion models

## Part 2:

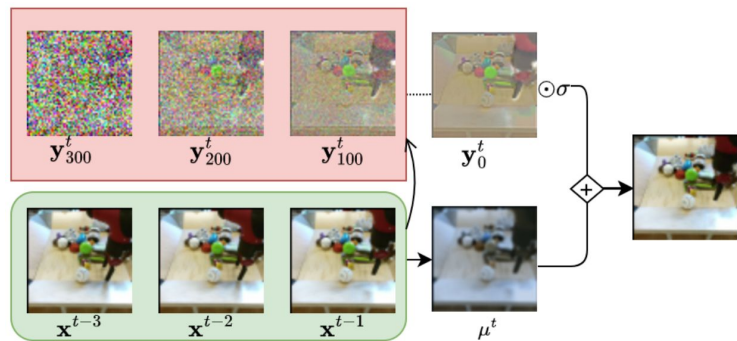
- Super-resolving atmospheric convection with diffusion models

# Diffusion Models for Image Generation

Prompt:  
“I’m an AI professor desperately worried to compete with industry research on computing resources. Depict my situation in a comic.”



# Diffusion Models for Video Generation



R. Yang, P. Srivastava, S. Mandt. arXiv 2022  
J. Ho et al., NeurIPS 2022  
Sora, OpenAI, 2024

# Diffusion Models for Precipitation Super-Resolution

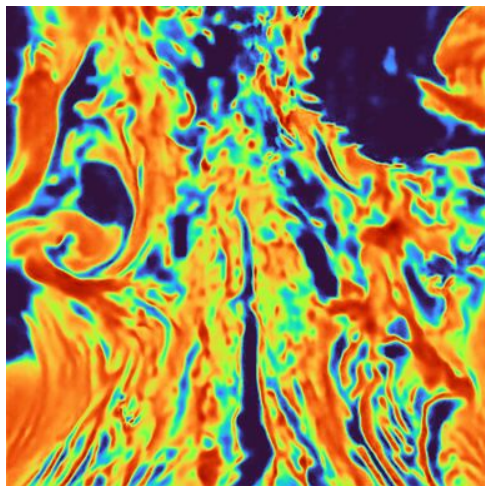
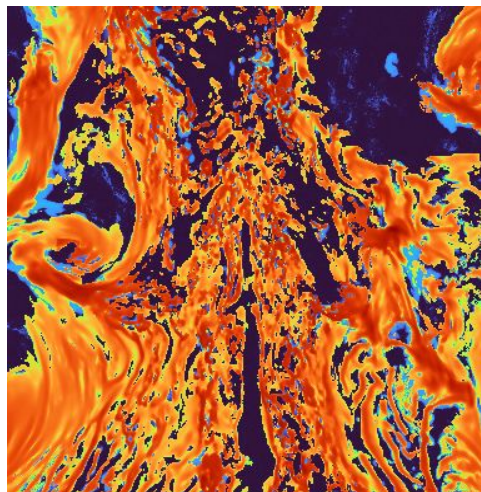
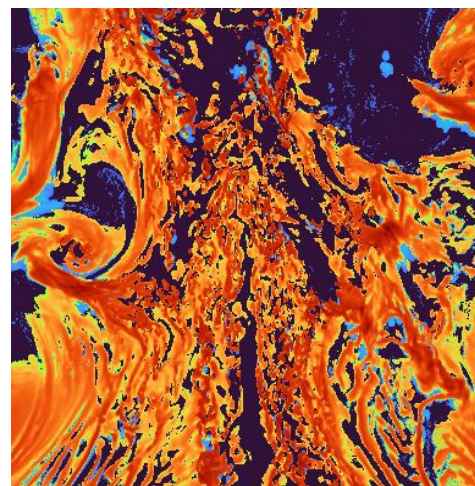


Image SR



Video SR  
(ours)



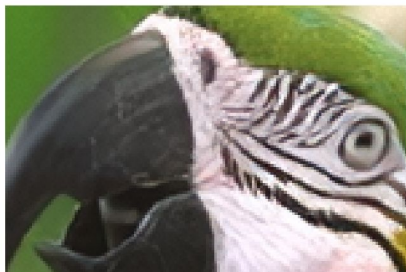
Truth



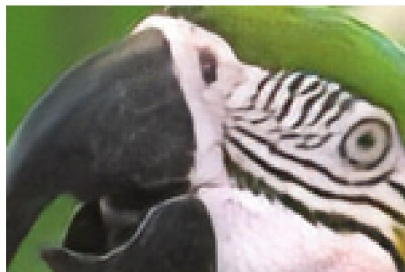
# Diffusion Models for Data Compression



(a) Ground Truth



(b) CDC( $\rho = 0.9$ ) (bpp=0.205)



(c) HiFiC (bpp=0.207)



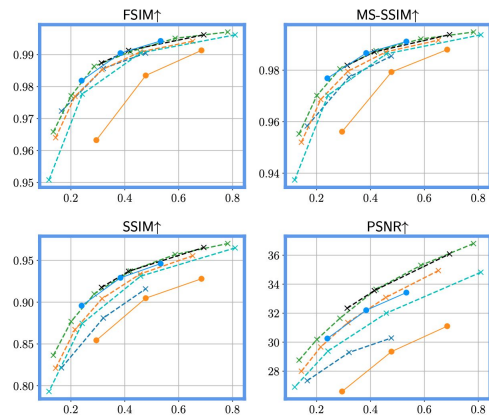
(d) Ground Truth



(e) CDC( $\rho = 0.9$ ) (bpp=0.398)



(f) HiFiC (bpp=0.456)



# This talk: physics, information, and generative modeling

---

## Deep Unsupervised Learning using Nonequilibrium Thermodynamics

---

**Jascha Sohl-Dickstein**  
Stanford University

JASCHA@STANFORD.EDU

**Eric A. Weiss**  
University of California, Berkeley

EWEISS@BERKELEY.EDU

**Niru Maheswaranathan**  
Stanford University

NIRUM@STANFORD.EDU

**Surya Ganguli**  
Stanford University

SGANGULI@STANFORD.EDU

### Abstract

A central problem in machine learning involves modeling complex data-sets using highly flexible families of probability distributions in which learning, sampling, inference, and evaluation are computationally intractable.

these models are unable to aptly describe structure in rich datasets. On the other hand, models that are *flexible* can be molded to fit structure in arbitrary data. For example, we can define models in terms of any (non-negative) function  $\phi(\mathbf{x})$  yielding the flexible distribution  $p(\mathbf{x}) = \frac{\phi(\mathbf{x})}{Z}$ , where  $Z$  is a normalization constant. However, computing this

## Nonequilibrium Measurements of Free Energy Differences for Microscopically Reversible Markovian Systems

Gavin E. Crooks<sup>1</sup>

*Received October 20, 1997; final December 16, 1997*

---

An equality has recently been shown relating the free energy difference between two equilibrium ensembles of a system and an ensemble average of the work required to switch between these two configurations. In the present paper it is shown that this result can be derived under the assumption that the system's dynamics is Markovian and microscopically reversible.

---

**KEY WORDS:** Nonequilibrium statistical mechanics; free energy; work; thermodynamic integration; thermodynamic perturbation.

### 1. INTRODUCTION

Consider a classical system in contact with a constant temperature heat bath where some degree of freedom of the system can be controlled.

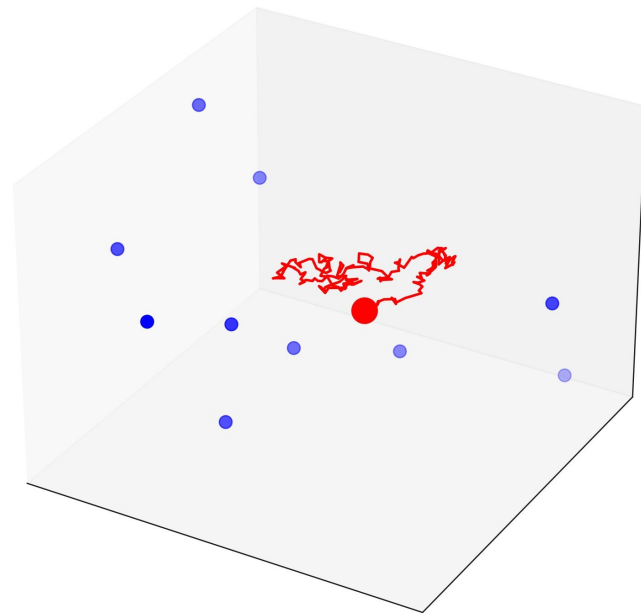
- Diffusion models are rooted non-equilibrium thermodynamics
  - Theory of irreversible processes; entropy production
  - Also connects to information theory and efficient data communication

# Diffusion Models



# Background: Brownian Motion

- Heavy particle (red) in a “bath” of particles (blue), frequent collisions
- Whole system is Newtonian/deterministic, but subsystem *appears* stochastic
- Stochastic process perspective:
  - deterministic “drift”  $f$  (external forces)
  - stochastic “diffusion”  $dW$  (due to collisions)

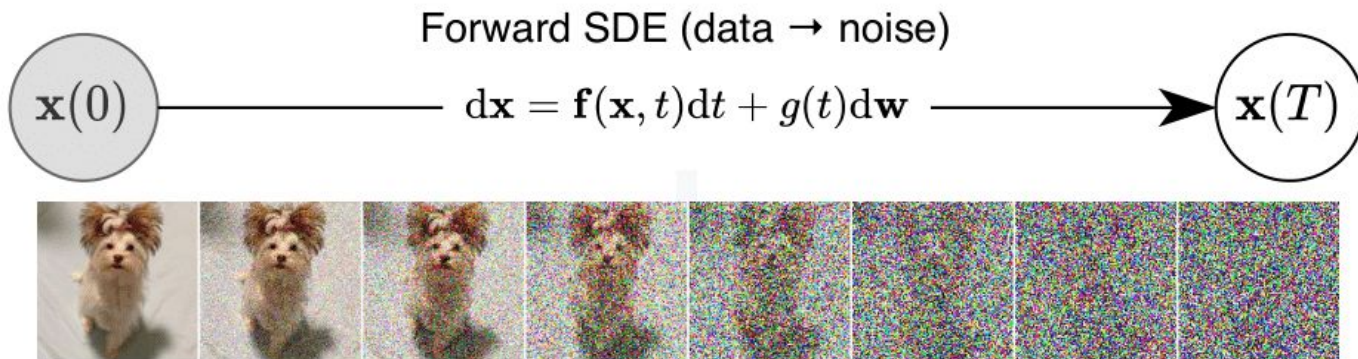


$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}$$

drift  $\swarrow$   $\nwarrow$  diffusion

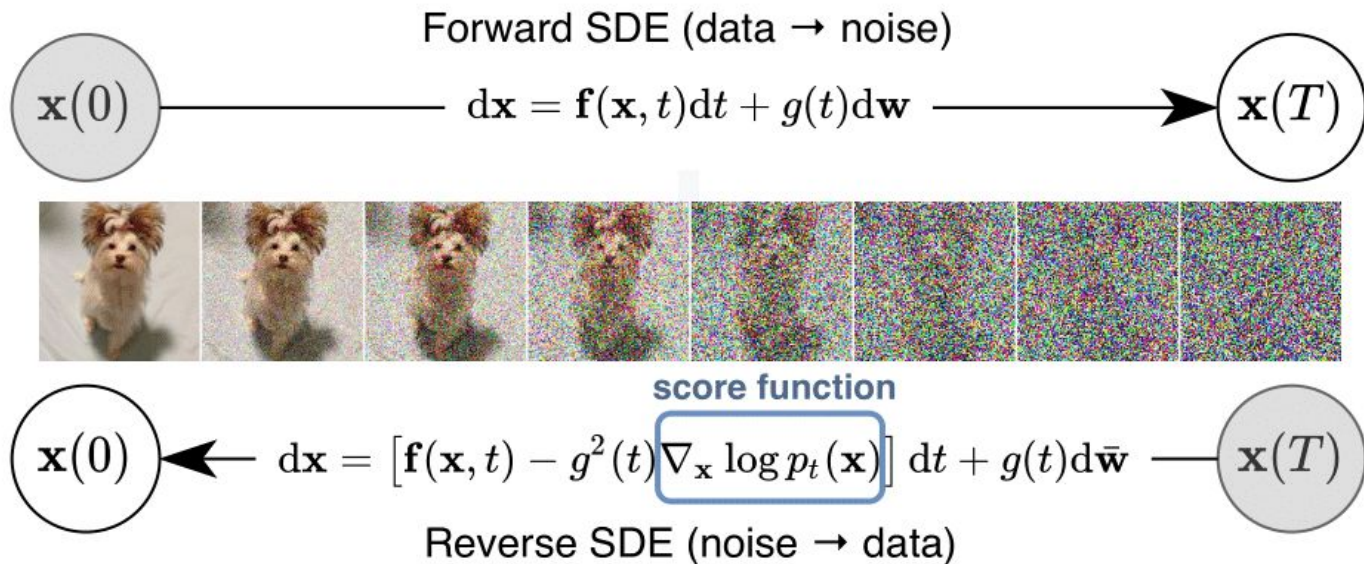
# Diffusion in Generative AI

“Creating noise from data is easy...”



# Diffusion in Generative AI

“Creating noise from data is easy... **creating data from noise is generative modeling**”



Forward Process



$$dz_t = F_t z_t dt + G_t dw_t, \quad t \in [0, T]$$

- Let  $z = (x, m)$  be a more general coordinate (e.g., pixel space + “momenta”)
- Increases **entropy** by transforming **data** to **noise**
- Usually no trainable parameters, no neural networks
- Desired: convergence to a Gaussian (so that we can sample from the inverse process)

## Reverse Process



$$dz_t = [F_t z_t - G_t G_t^\top \underbrace{\nabla \log p(z_t)}_{\text{Score}}] dt + G_t dw_t$$

- Reduces **entropy** by transforming **noise** to **data**
- Notably, the **score**  $\nabla \log p(z_t)$  enters the process (flow to high density regions)
- Unfortunately, we don't know the score!



## Reverse Process



$$dz_t = [F_t z_t - G_t G_t^\top \boxed{s_\theta(z_t, t)}] dt + G_t dw_t$$

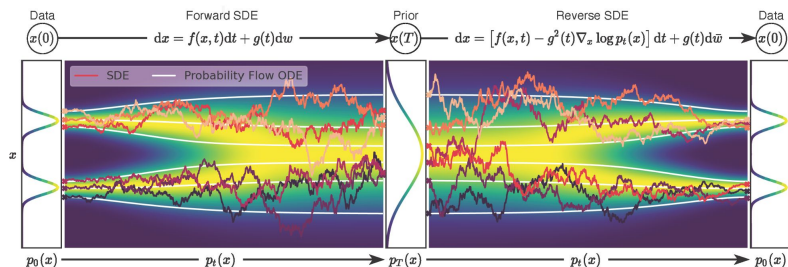
Approximated score

- Reduces **entropy** by transforming **noise** to **data**
- The score  $\nabla \log p(z_t)$  is approximated using a neural network:  $s_\theta(z_t, t)$
- Once the score is learned, we can sample from the model by solving the stochastic differential equation numerically
- The score is learned using **score matching** (regression); skip details.

Part 1:  
Augmented Diffusions  
and Efficient Integrators

# Current diffusion models are inspired by physics

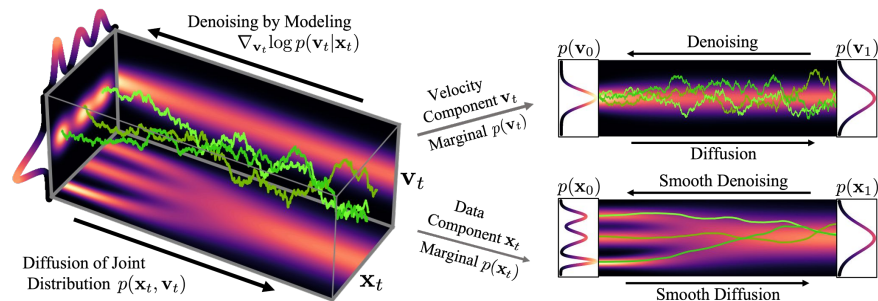
## Diffusion in “position” space



- Perform diffusion only in the data space,  $z = x$
- Follow an Ornstein-Uhlenbeck process

$$F_t = -\frac{1}{2}\beta_t I_d \quad G_t = \sqrt{\beta_t} I_d$$

## Diffusion in “phase space”



- Perform diffusion in an **augmented** space i.e.  $z_t = [x_t, m_t]$
- Inspired from Molecular Dynamics

$$F_t = \left( \beta \begin{pmatrix} 0 & M^{-1} \\ -1 & -\Gamma M^{-1} \end{pmatrix} \otimes I_d \right), \quad G_t = \left( \begin{pmatrix} 0 & 0 \\ 0 & \sqrt{2\Gamma\beta} \end{pmatrix} \otimes I_d \right)$$

# How to design new diffusion models beyond physics intuition?

- **Augmented dimensions** that aren't necessarily momenta
  - Current models require  $\dim(x) = \dim(m)$  to match
- **Noise sources** that aren't necessarily thermal noise
  - Current models couple thermal noise with momentum, if available
- **Drift forces** that aren't necessarily conservative
  - Forces do not necessarily have to be gradients of scalar potentials

Ma, Chen, and Fox. A Complete Recipe for Stochastic Gradient MCMC. NeurIPS 2015

Singhal, Goldstein, Ranganath. Where to diffuse, how to diffuse, and how to get back. ICLR 2023

Pandey and Mandt. A Complete Recipe for Diffusion Generative Models. ICCV 2023.

# A Complete Recipe for Diffusion Models

- **A1:** Consider position and auxiliary variables:  $z = [x, m]^\top$
- **A2:** Consider continuous-time, first-order Markov process:  $dz = f(z)dt + \sqrt{2D(z)}dw_t$
- **A3:** Demand converge to a **simple**, pre-specified prior:  $p_s(z) \propto \exp(-H(z))$   
e.g.,  $H(z) = |x|^2 + |m|^2$

- **Result:** The following parameterization is **complete (always exists & unique)**:

$$f(z) = -(D(z) + Q(z))\nabla H + \tau(z)$$

$$\tau_i(z) = \sum_{j=1}^d \frac{\partial}{\partial z_j} (D_{ij}(z) + Q_{ij}(z))$$

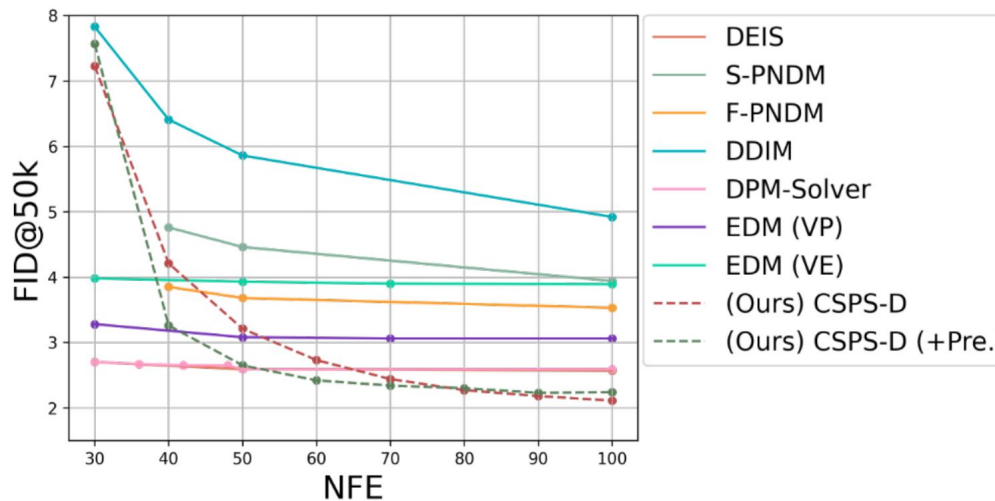


# Phase-Space Langevin Diffusion: A New Sampler

$$d \begin{pmatrix} x \\ m \end{pmatrix} = \frac{\beta}{2} \begin{pmatrix} -\Gamma & M^{-1} \\ -1 & -\nu \end{pmatrix} \begin{pmatrix} x \\ m \end{pmatrix} dt + \begin{pmatrix} \sqrt{\Gamma\beta} & 0 \\ 0 & \sqrt{M\nu\beta} \end{pmatrix} dw_t,$$

CIFAR-10 (ODE)

- The noise parameter  $\Gamma$  is “unphysical”, but improves convergence
- Noise sources in both position and momentum
- Works even better with splitting integrators



# Phase-Space Langevin Diffusion: A New Sampler



# Diffusion models may live in poorly-conditioned geometries

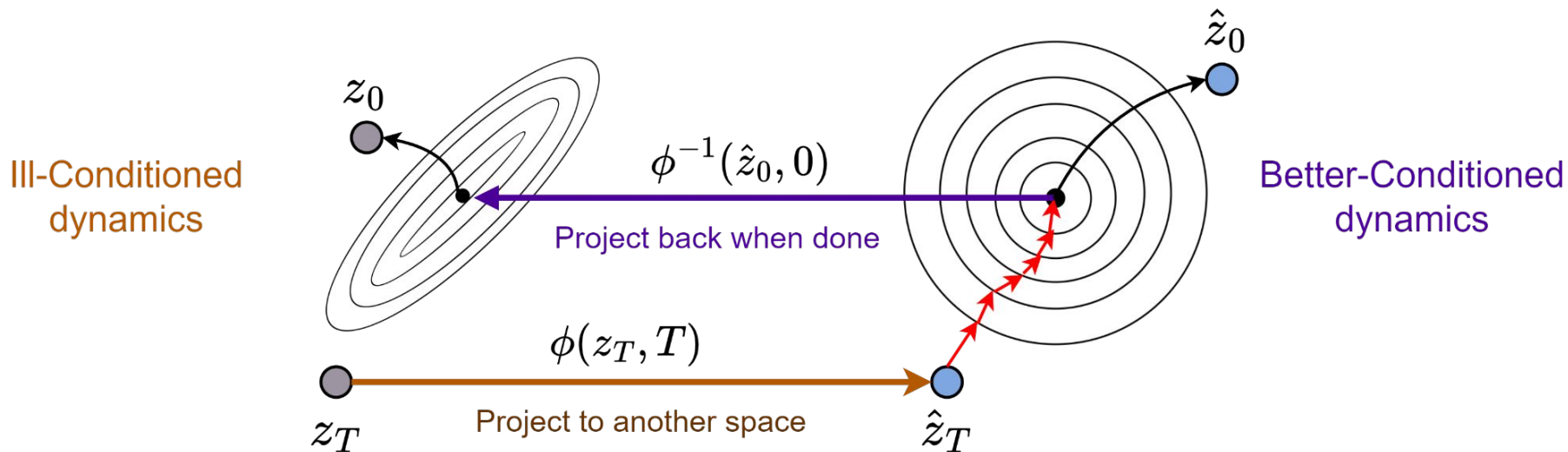
$$dz_t = \left[ F_t z_t - \frac{1}{2} G_t G_t^\top s_\theta(z_t, t) \right] dt$$

Linear drift; could be solved analytically in isolation

Diffusion coefficient matrix; can complicate sampling geometry

- Claim: the geometry is “ill-conditioned” (as in optimization)
- Intuition: change of coordinates before simulating the equations
  - Simpler equations can be solved in fewer discretization steps
  - Eliminate the state-independent drift

# Conjugate Integrators



Works in already trained models! E.g., OpenAI diffusion models trained on ImageNet

# Conjugate Integrators

Consider linear transformation,  $\hat{z}_t = A_t z_t$

Score parameterization:  $s_\theta(z_t, t) = C_{\text{skip}}(t)z_t + C_{\text{out}}(t)\epsilon(z_t, t)$

Define  $A_t = \exp\left(\int_0^t B_s - F_s + \frac{1}{2}G_s G_s^\top C_{\text{skip}}(s)ds\right)$ ,

$$\Phi_t = -\int_0^t \frac{1}{2}A_s G_s G_s^\top C_{\text{out}}(s)ds$$

After some straightforward math:

$$d\hat{z}_t = A_t B_t A_t^{-1} \hat{z}_t dt + d\Phi_t \epsilon_\theta(A_t^{-1} \hat{z}_t, t)$$

Optional damping term;  
can set  $B=0$  or  $B = -\lambda I$

Linear preconditioner, can be precomputed



# Conjugate Integrators - Theory

- **Connections to ODE stability criteria**

- Let  $\mathcal{G}$  be the flow map of an ODE integrator. *Stability* implies that  $\forall \Delta \exists \epsilon$ :

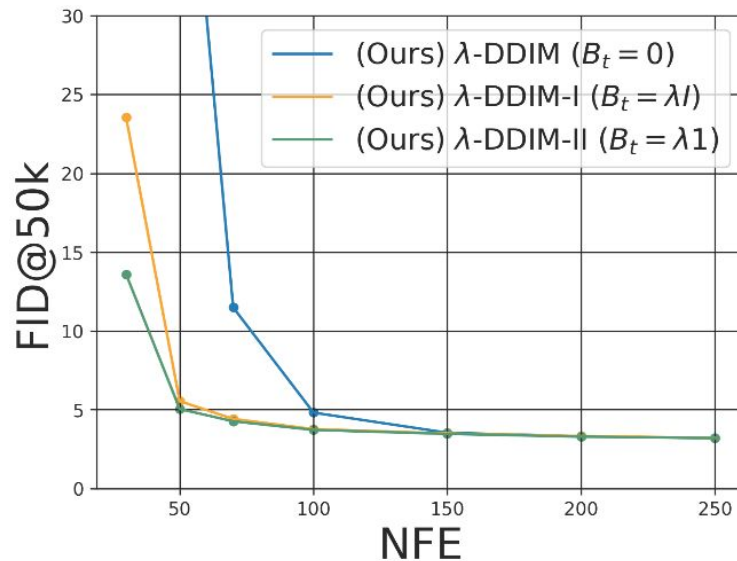
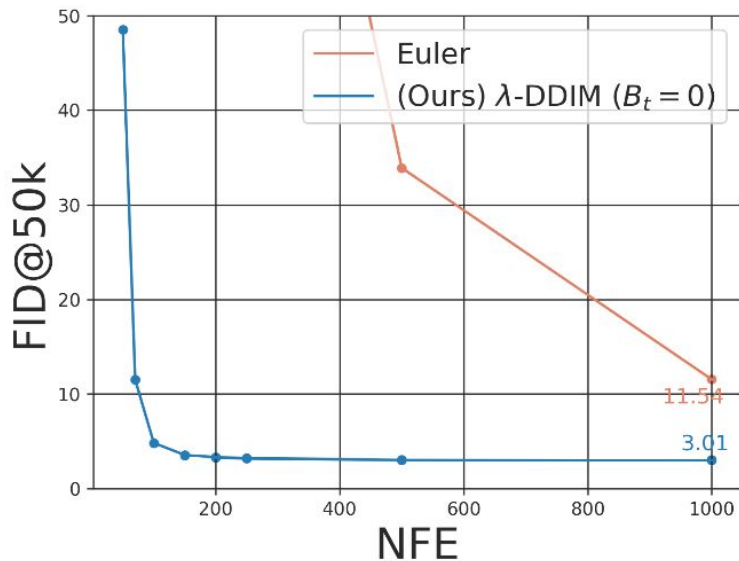
$$\|\mathcal{G}_h(\hat{z}(t)) - \mathcal{G}_h(\hat{z}_t)\| \leq \Delta, \quad \text{s.t. } \|\hat{z}(t) - \hat{z}_t\| < \epsilon, \epsilon > 0, \Delta > 0$$

Stability analysis considers the Eigenvalues of the linearized flow operator  
Different choices of B can enhance stability, improving over existing samplers

- **Connections to prior methods:**

- Interestingly,  $B_t = 0$  corresponds to DDIM for diffusion models.
- Also connections with fast samplers based on exponential integrators

# Speed vs. Sample Quality Tradeoffs

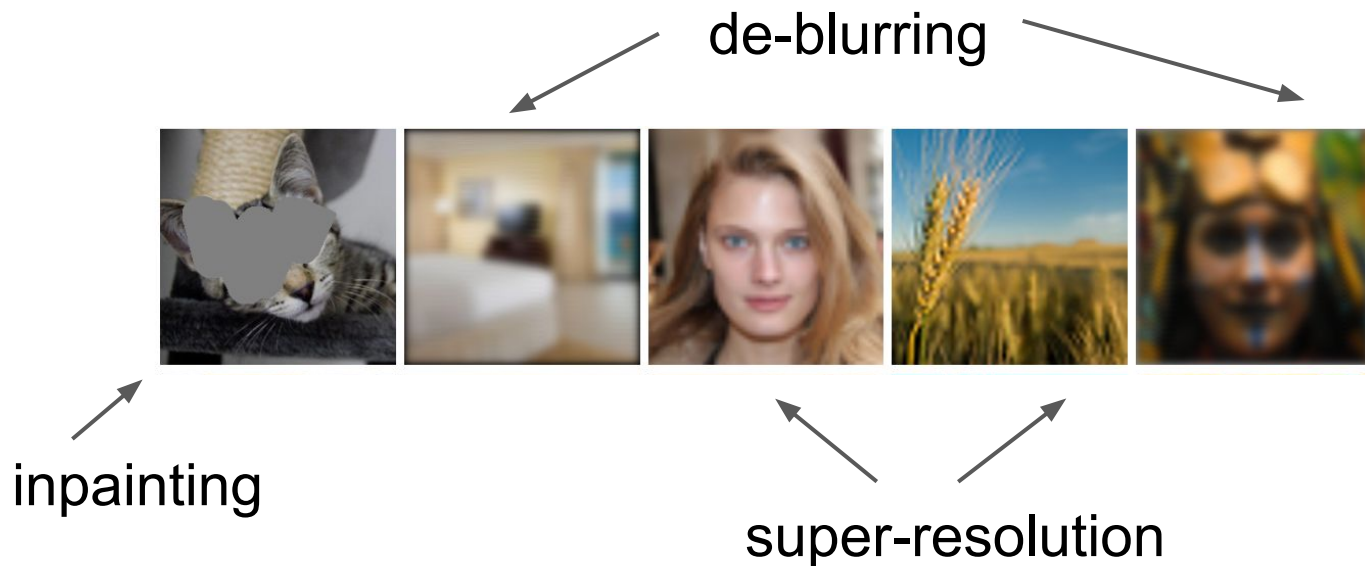


- Dramatic speed improvements; in particular for phase-space diffusions
- Non-zero damping term  $\lambda$  can further improve performance

## Now, consider inverse problems

Degradation operator  $H$ :

$$\mathbf{y} = \mathbf{H}\mathbf{x}_0 + \sigma_y \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}), \quad \mathbf{x}_0 \sim p_{\text{data}}$$



## Now, consider inverse problems

Degradation  
operator  $H$ :

$$\mathbf{y} = \mathbf{H}\mathbf{x}_0 + \sigma_y \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}), \quad \mathbf{x}_0 \sim p_{\text{data}}$$

We will consider two types of iterative refinement models:

- Diffusion models
- Flow matching models

Diffusion: 
$$d\mathbf{x}_t = \left[ \mathbf{F}_t \mathbf{x}_t - \frac{1}{2} \mathbf{G}_t \mathbf{G}_t^\top \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) \right] dt, \quad (3)$$

Flows: 
$$d\mathbf{x}_t = \mathbf{b}(\mathbf{x}_t, \mathbf{y}, t) dt,$$

# Fast Samplers for Inverse Problems in Iterative Refinement Models

Degradation  
operator  $\mathbf{H}$ :

$$\mathbf{y} = \mathbf{H}\mathbf{x}_0 + \sigma_y \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}), \quad \mathbf{x}_0 \sim p_{\text{data}}$$

**Proposition 1.** *For the conditional diffusion dynamics defined in Eqn. 3, introducing a diffeomorphism,  $\bar{\mathbf{x}}_t = \mathbf{A}_t \mathbf{x}_t$ , where,*

$$\mathbf{A}_t = \exp\left(\int_0^t \mathbf{B}_s - \mathbf{F}_s ds\right), \quad \Phi_t = -\int_0^t \frac{1}{2} \mathbf{A}_s \mathbf{G}_s \mathbf{G}_s^\top \mathbf{C}_{out}(s) ds, \quad (6)$$

induces the following projected diffusion dynamics,

$$d\hat{\mathbf{x}}_t = \mathbf{A}_t \mathbf{B}_t \mathbf{A}_t^{-1} \hat{\mathbf{x}}_t dt + d\Phi_t \epsilon_\theta(\mathbf{x}_t, t) - \frac{w_t r_t^{-2}}{2} \mathbf{G}_t \mathbf{G}_t^\top \frac{\partial \hat{\mathbf{x}}_0^\top}{\partial \mathbf{x}_t} (\mathbf{H}^\dagger \mathbf{y} - \mathbf{P} \hat{\mathbf{x}}_0) dt, \quad (7)$$

where  $\mathbf{H}^\dagger = \mathbf{H}^\top (\mathbf{H}\mathbf{H}^\top)^{-1}$  and  $\mathbf{P} = \mathbf{H}^\top (\mathbf{H}\mathbf{H}^\top)^{-1} \mathbf{H}$  represent the pseudoinverse and the orthogonal projector operators for the degradation operator  $\mathbf{H}$ . (Proof in Appendix A.2)

Main idea: assign different dynamics to degradation operator's null space and its orthogonal complement.



# Conjugate Integrators: high-quality generation in only five iterations



Pandey, Yang,  
Mandt.  
<https://arxiv.org/pdf/2405.17673>

## Conjugate Integrators: high-quality generation in only five iterations



Ground Truth



Regular  
integrator  
(NFE=5)



Conjugate  
Integrator  
(NFE=5)

# Fast Samplers for Inverse Problems in Iterative Refinement Models

Flow Results	NFE	LPIPS↓		KID×10 <sup>-3</sup> ↓		FID↓	
		C-PIGFM	PIGFM	C-PIGFM	PIGFM	C-PIGFM	PIGFM
Inpainting	5	<b>0.125</b>	0.240	<b>17.6</b>	167.0	<b>26.95</b>	161.49
	10	<b>0.074</b>	0.188	<b>8.0</b>	86.6	<b>14.64</b>	94.91
	20	<b>0.065</b>	0.144	<b>4.6</b>	54.4	<b>10.93</b>	65.39
Super-Resolution	5	<b>0.063</b>	0.091	<b>5.5</b>	17.5	<b>13.08</b>	21.84
	10	<b>0.058</b>	0.076	<b>3.6</b>	12.2	<b>10.65</b>	16.73
	20	<b>0.064</b>	0.069	3.9	<b>3.5</b>	11.07	<b>10.23</b>
Deblurring	5	<b>0.083</b>	0.114	<b>3.7</b>	10.9	<b>12.86</b>	18.97
	10	<b>0.077</b>	0.088	<b>5.0</b>	7.0	<b>14.41</b>	15.09
	20	0.080	<b>0.073</b>	7.9	<b>3.1</b>	17.10	<b>11.35</b>

Diffusion Results		C-PIGDM	PIGDM	DPS	DDRM	C-PIGDM	PIGDM	DPS	DDRM	C-PIGDM	PIGDM	DPS	DDRM
		Super-Resolution	5	<b>0.220</b>	0.306			<b>2.7</b>	6.3			<b>37.31</b>	49.06
	10	<b>0.206</b>	0.252	0.252	0.318	<b>1.6</b>	4.8	5.8	14.1	<b>34.22</b>	44.30	38.18	51.64
	20	<b>0.207</b>	0.222			<b>1.7</b>	2.5			<b>34.28</b>	37.36		
Deblurring	5	<b>0.272</b>	0.349			<b>3.89</b>	14.1			<b>44.42</b>	63.94		
	10	<b>0.272</b>	0.294	0.619	0.336	<b>3.6</b>	5.3	59.5	12.3	<b>43.37</b>	47.80	139.58	62.53
	20	0.268	<b>0.259</b>			<b>3.5</b>	4.2			<b>43.70</b>	44.20		

# Reverse diffusion as progressive decompression

A diffusion model can be understood as:

- Denoising autoencoder at multiple noise levels [Vincent 2011, Song & Ermon, 2019]
- Learning to reverse an SDE (mostly this talk) [Song et al., 2021]
- **Deep hierarchical VAE** [Sohl-Dickstein et al., 2015, Ho et al., 2020, Kingma et al., 2021]

$$-\text{VLB}(\mathbf{x}) = \mathbb{E} [D_{\text{KL}}[q(\mathbf{z}_T | \mathbf{X}) \parallel p_T(\mathbf{z}_T)]] + \sum_{s=1}^{T-1} \mathbb{E} [D_{\text{KL}}[q(\mathbf{z}_s | \mathbf{Z}_{s+1}, \mathbf{X}) \parallel p(\mathbf{z}_s | \mathbf{Z}_{s+1})]]$$

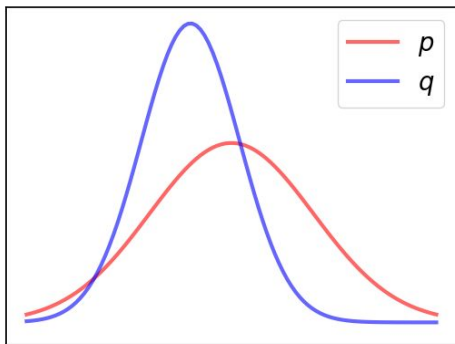
Entropy reduced in every reverse diffusion step

- Can this information be efficiently transmitted between a sender and receiver?

## Background: Relative entropy coding

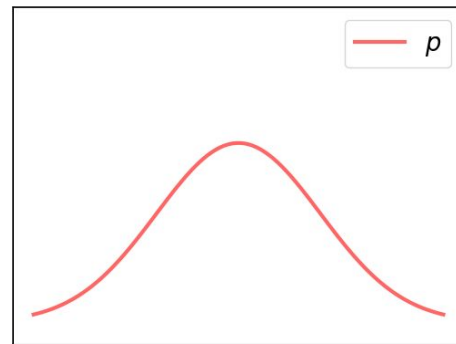
Problem: Alice wants to transmit a sample from  $q$  to Bob, under shared prior  $p$ , using  $KL(q||p)$  bits.

Alice



01011010111...

Bob

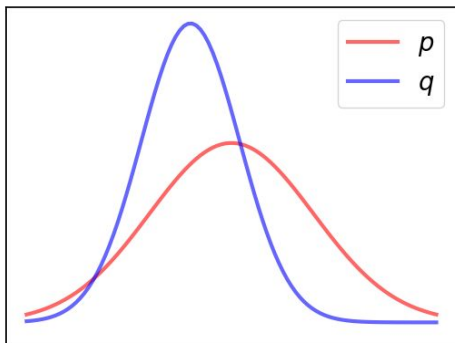


# Background: Relative entropy coding

Idea (sketch):

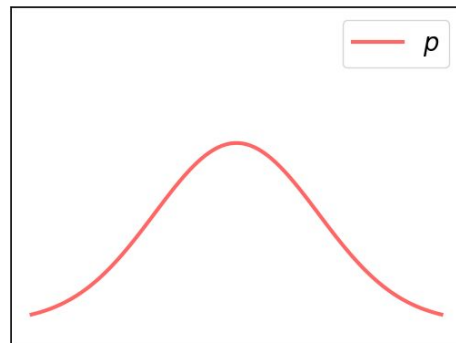
- Let Alice and Bob share a random seed.
- Sample from  $p$  many times until we hit a “good” (high likelihood under  $q$ ) sample
- Transmit the index  $K$  in binary.

Alice



01011010111...

Bob



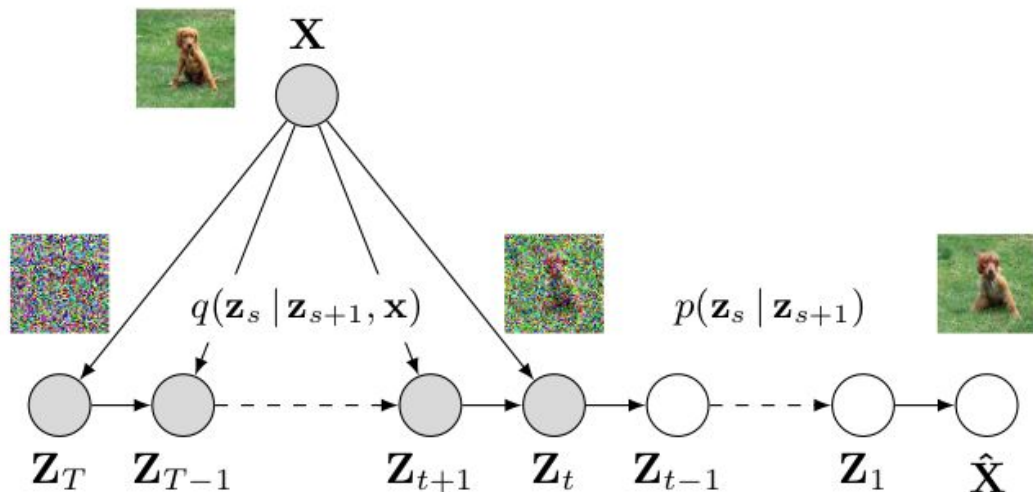
# Diffusion models for transmitting information

Problem setup: Alice wants to transmit data  $\mathbf{x}$  to Bob using  $-VLB(\mathbf{x})$  many bits.

$$\mathbb{E} [D_{\text{KL}}[q(\mathbf{z}_T | \mathbf{X}) \parallel p_T(\mathbf{z}_T)]] + \sum_{s=1}^{T-1} \mathbb{E} [D_{\text{KL}}[q(\mathbf{z}_s | \mathbf{z}_{s+1}, \mathbf{X}) \parallel p(\mathbf{z}_s | \mathbf{z}_{s+1})]]$$

posterior

shared prior



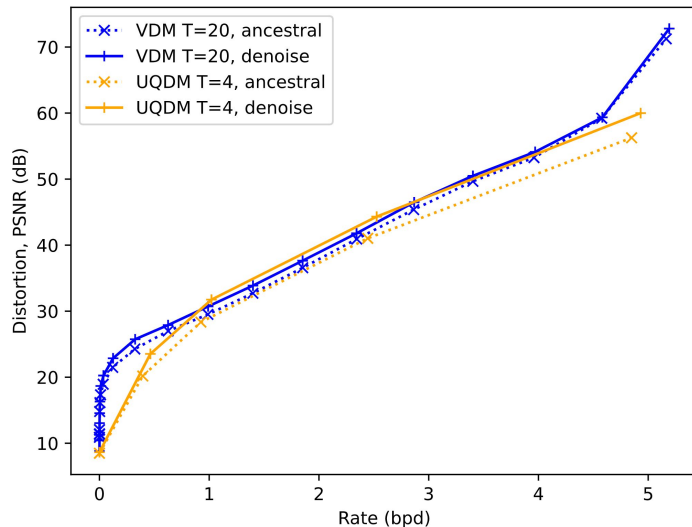
- Works in theory; doesn't scale with dimension
- Practical solution: *universal quantization*
- Requires re-design of the diffusion process, Gaussian  $\rightarrow$  Uniform



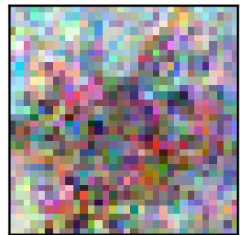
# Preliminary results

## CIFAR data

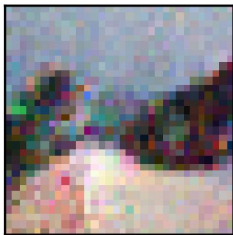
- Promising results when compared to Gaussian diffusion at  $T=20$  discretization steps
- Still work in progress; not competitive at larger  $T$  & compared to SOTA models



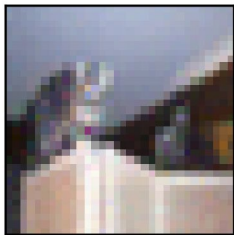
bpd=0.00,  
psnr=11.49



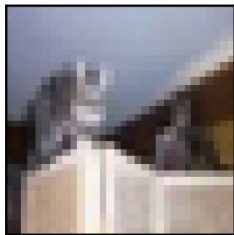
bpd=0.20,  
psnr=19.86



bpd=0.78,  
psnr=29.36



bpd=2.27,  
psnr=40.95



bpd=4.80,  
psnr=56.07

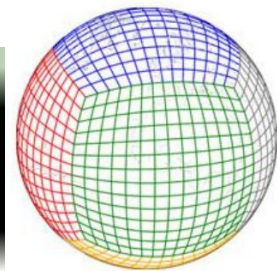


Yang, Mandt, and Theis.  
An introduction to neural data  
compression. Foundations  
and Trends in Computer  
Vision, 2023

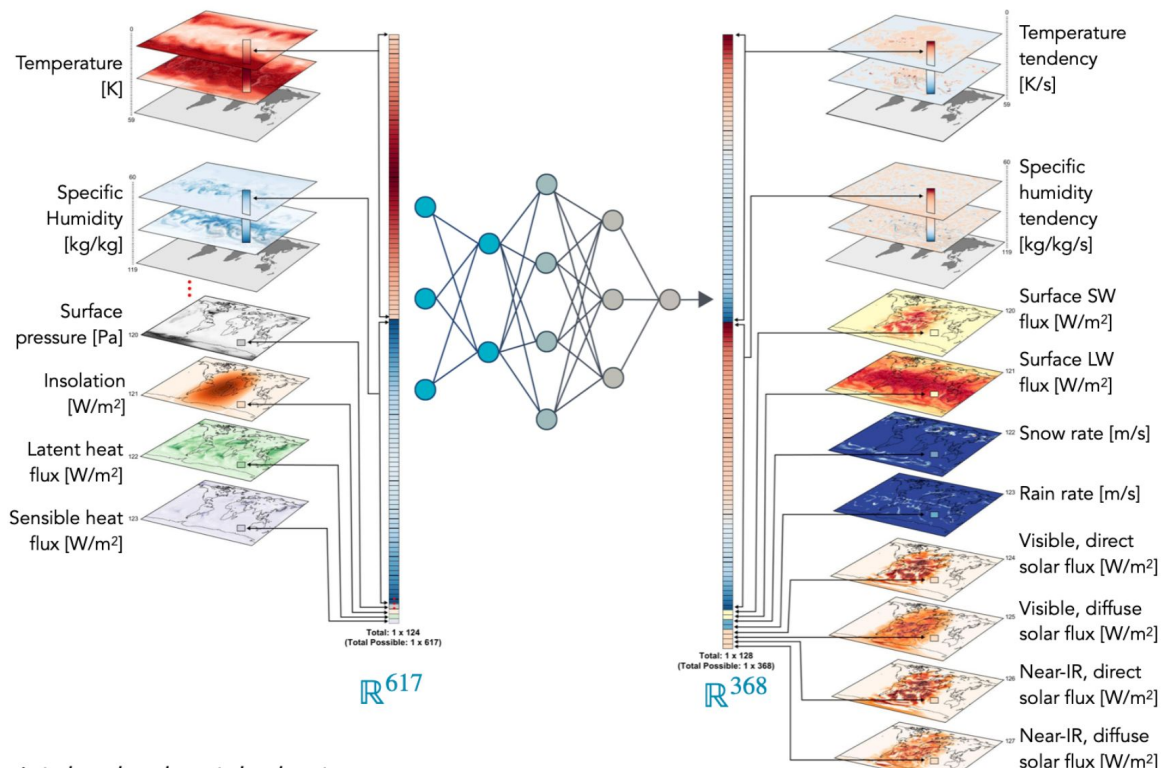
# Part 2: Emulating Thermodynamic Processes with Diffusion Models

# Generative Modeling for Atmospheric Convection

- The climate modeling dilemma:
  - **Either** simulate the climate at sufficiently high spatial resolution (e.g., a few km) to capture, e.g., cloud-related processes
  - **Or**, simulate the climate for a long-enough time (several decades) to make accurate predictions on global warming
- Unresolved processes are huge drivers of uncertainty and introduce randomness and bias
- Can we use generative modeling to stochastically downscale a low-resolution simulation?



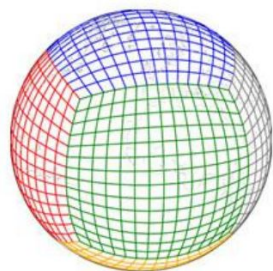
# Aside: physics-ML hybrid models may solve the resolution dilemma



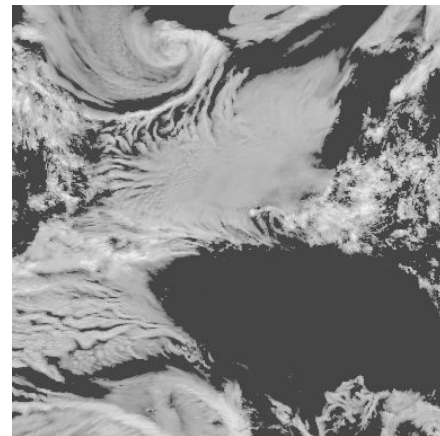
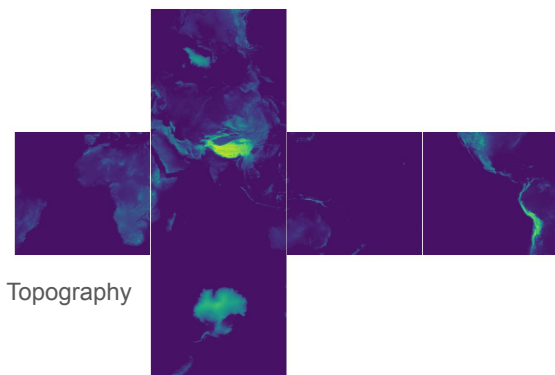
# The data: high-resolution atmospheric simulation

Focus on precipitation channel

- FV3GFS global atmosphere simulation dataset (Allen Institute for AI)
- Captures all relevant physical fields, including precipitation (rain, snow)
- 25 km resolution, 3-hourly



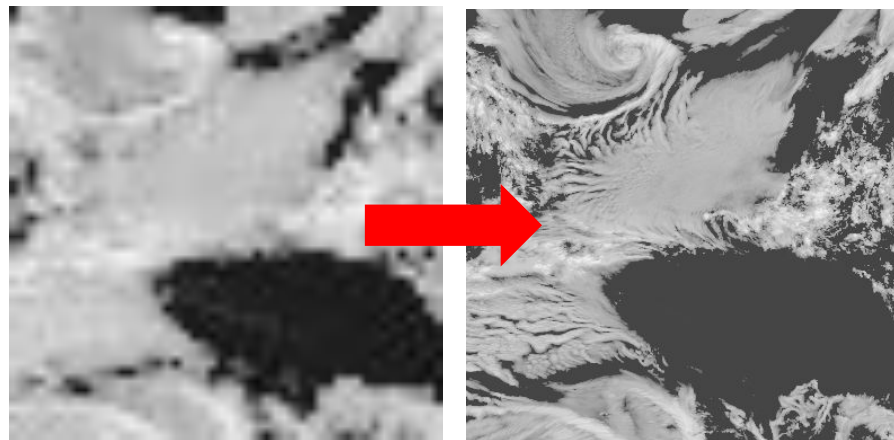
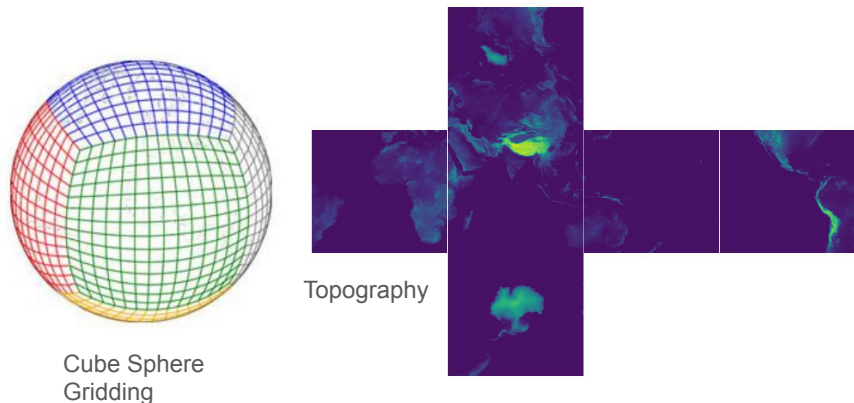
Cube Sphere  
Gridding



# The data: high-resolution atmospheric simulation

Focus on precipitation channel

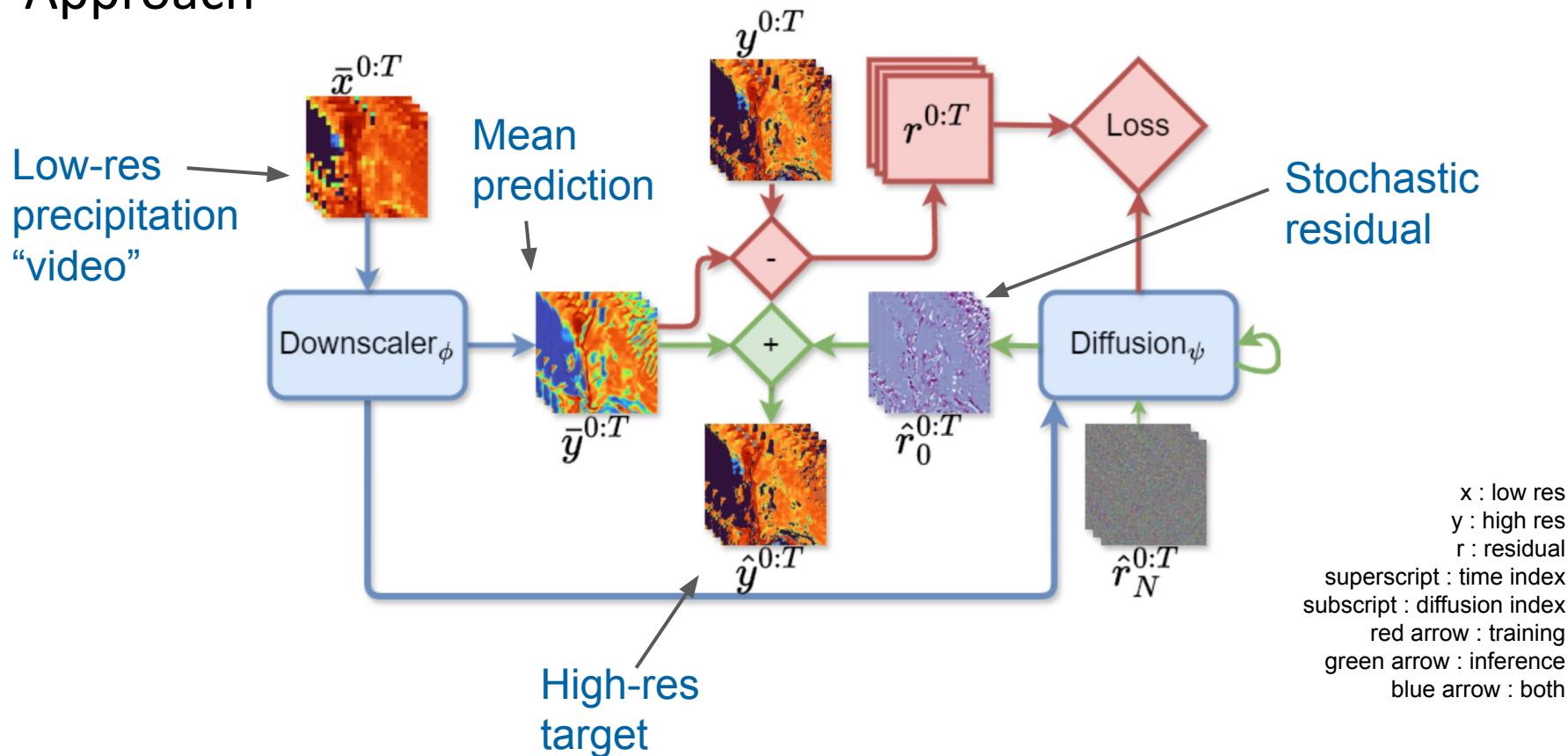
- FV3GFS global atmosphere simulation dataset (Allen Institute for AI)
- Captures all relevant physical fields, including precipitation (rain, snow)
- 25 km resolution, 3-hourly



**Goal: super-resolve the precipitation channel**

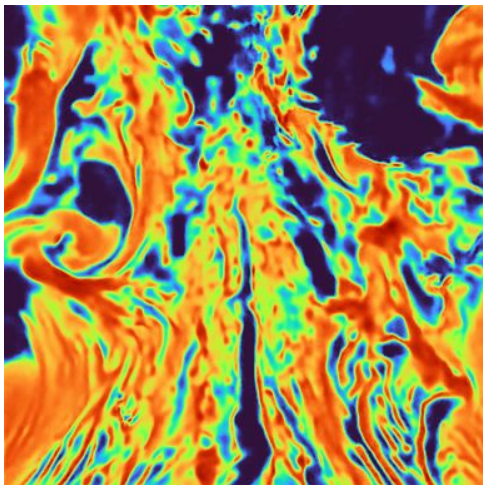
- Allows to run a cheap model to simulate many years
- Use a super-resolution model to convert the data to high spatial resolution

# Approach

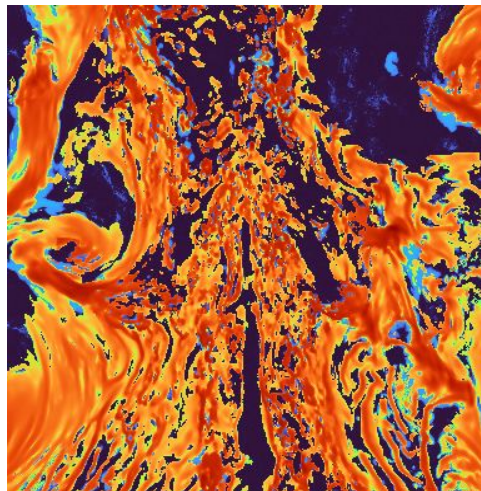




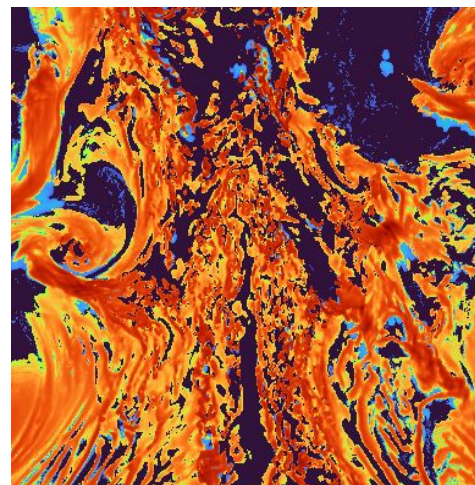
# Video Diffusion accurately captures temporal information



single-frame SR



multi-frame SR

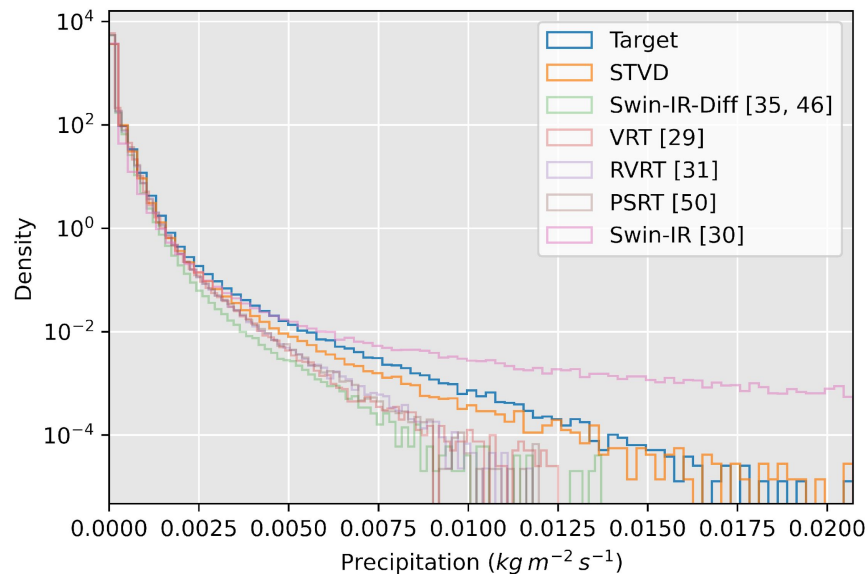


Truth



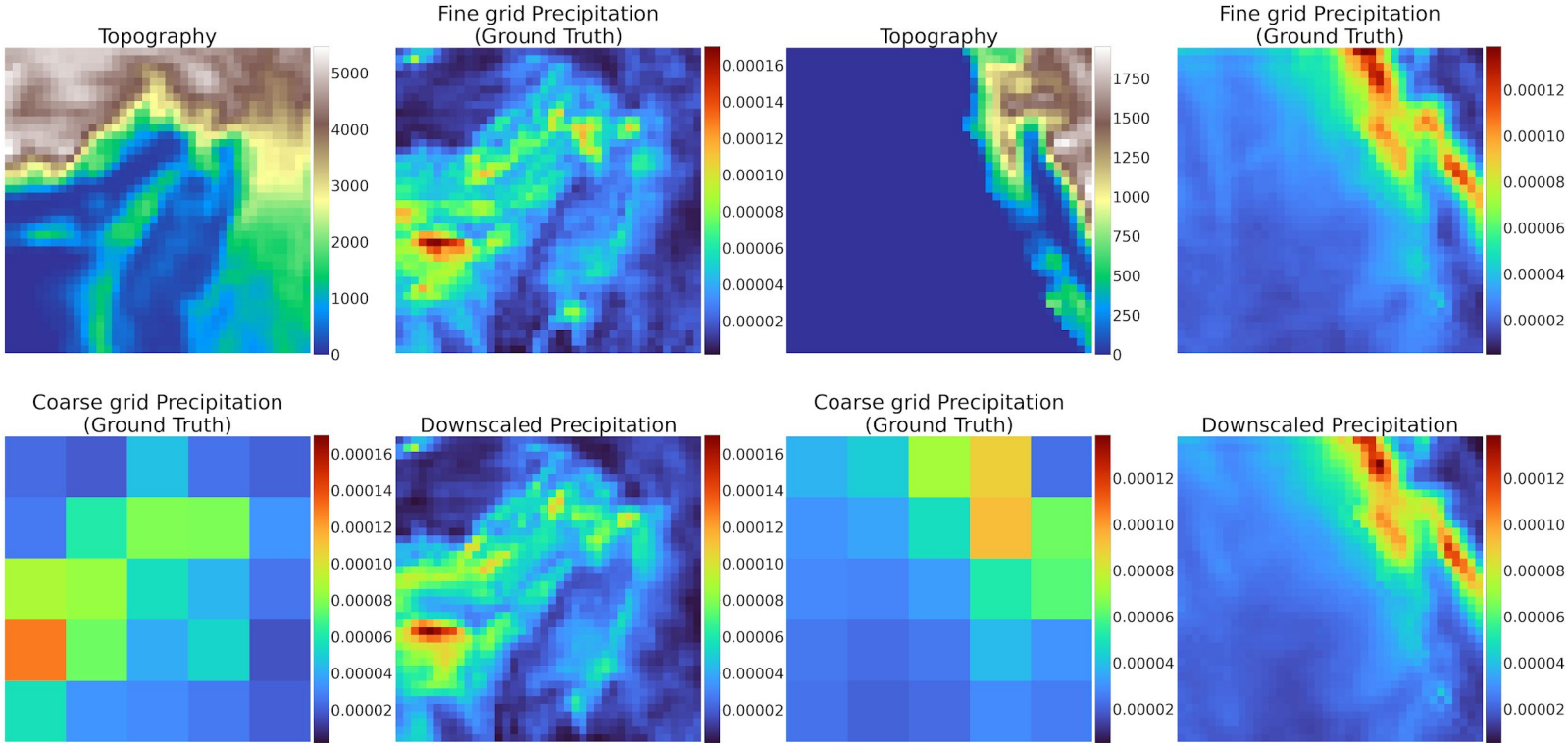
# Deterministic models underestimate extreme precipitation

- Downscaling especially challenging for precipitation because of rare extreme events
  - Few geographical regions with extreme rainfall
  - Cyclones/extreme weather events
- While most models can capture low precipitation regions quite well, they heavily downplay the extremes



Precipitation distribution over three-hour windows on all grid points around the globe

# Precipitation: Annual Averages



# Diffusion Modeling in Molecular Dynamics

A treasure trove for diffusion  
generative modeling

- Phase space methods
- Conjugate Integrators
- Splitting Integrators

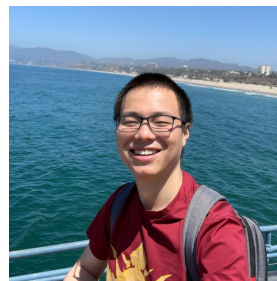
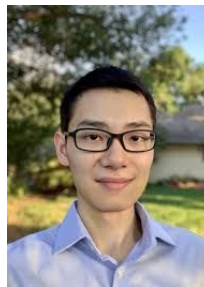
Interdisciplinary Applied Mathematics 39

Ben Leimkuhler  
Charles Matthews

# Molecular Dynamics

With Deterministic and Stochastic  
Numerical Methods

# Summary



Kushagra Pandey, Yibo Yang, Ruihan Yang, Prakhar Srivastava

- The thermodynamics of diffusion
  - Origins of diffusion models in thermodynamics
  - Can adopt ideas such as physics-inspired generative processes (and beyond)
  - Efficient sampler design
  - Diffusion models as efficient progressive coders
- Diffusion for thermodynamics/climate
  - Climate data as a playground for generative modeling
  - Requires stochasticity to capture distribution-level properties (e.g., annual precipitation)
- Open questions / future research:
  - How to incorporate physical constraints into modeling convection
  - Unpaired distribution-to-distribution translation between climate models
  - Theoretical analysis of diffusion models: dynamical phase transitions, critical slowing down