# Towards Efficient Inference-Time Scaling without Distillation
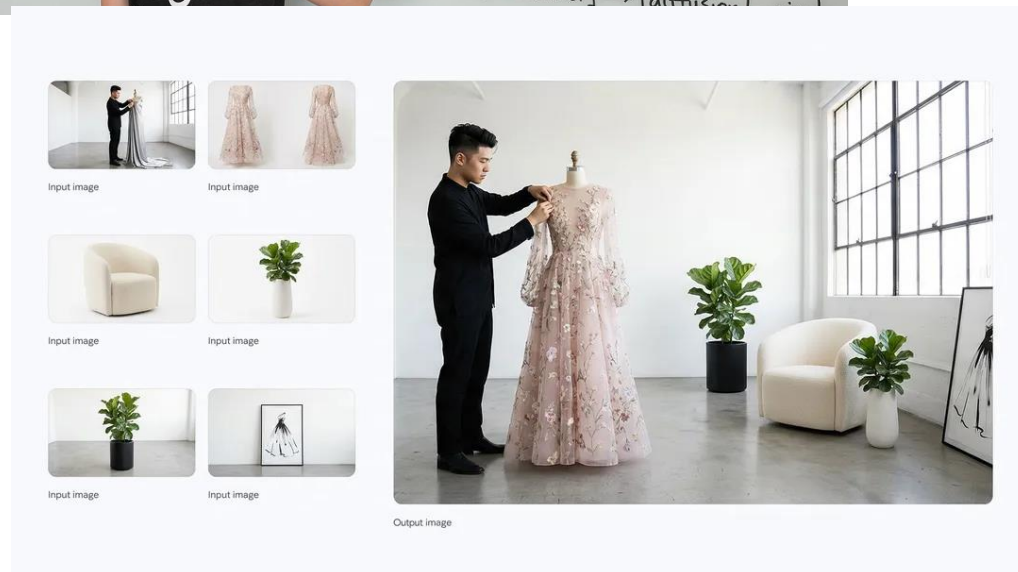
Linqi (Alex) Zhou

Luma

# Rise of Diffusion Models

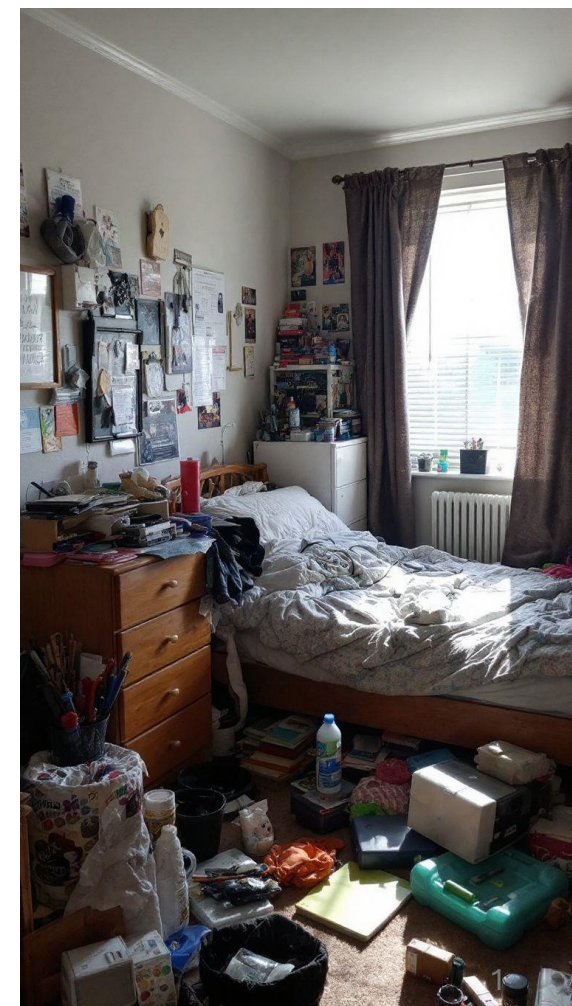

GPT-4o Image

Midjourney v7

FLUX

Nano Banana

# Text-to-Video Models at Luma



Luma
Ray 3

# Why are Diffusion and Flow Matching so Good?



Data

Prior

Not known in closed form!

stic process

• Simple L2 loss → stable and scalable!

# Problems of Diffusion and Flow Matching

$$\mathrm{d}\mathbf{x}_t = \mathbf{u}_t \mathrm{d}t \qquad \text{(probability-flow ODE)}$$
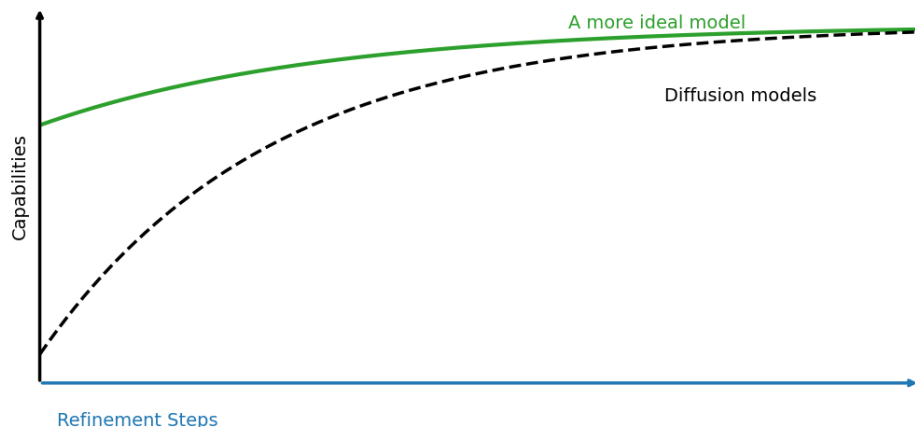
NOT optimal in utilizing network capacity.

ODE simulation error       Slow inference

Ideal case: one- or few-step mapping from prior to data.
(**efficient** inference-time scaling)

# Towards Efficient Inference-Time Scaling

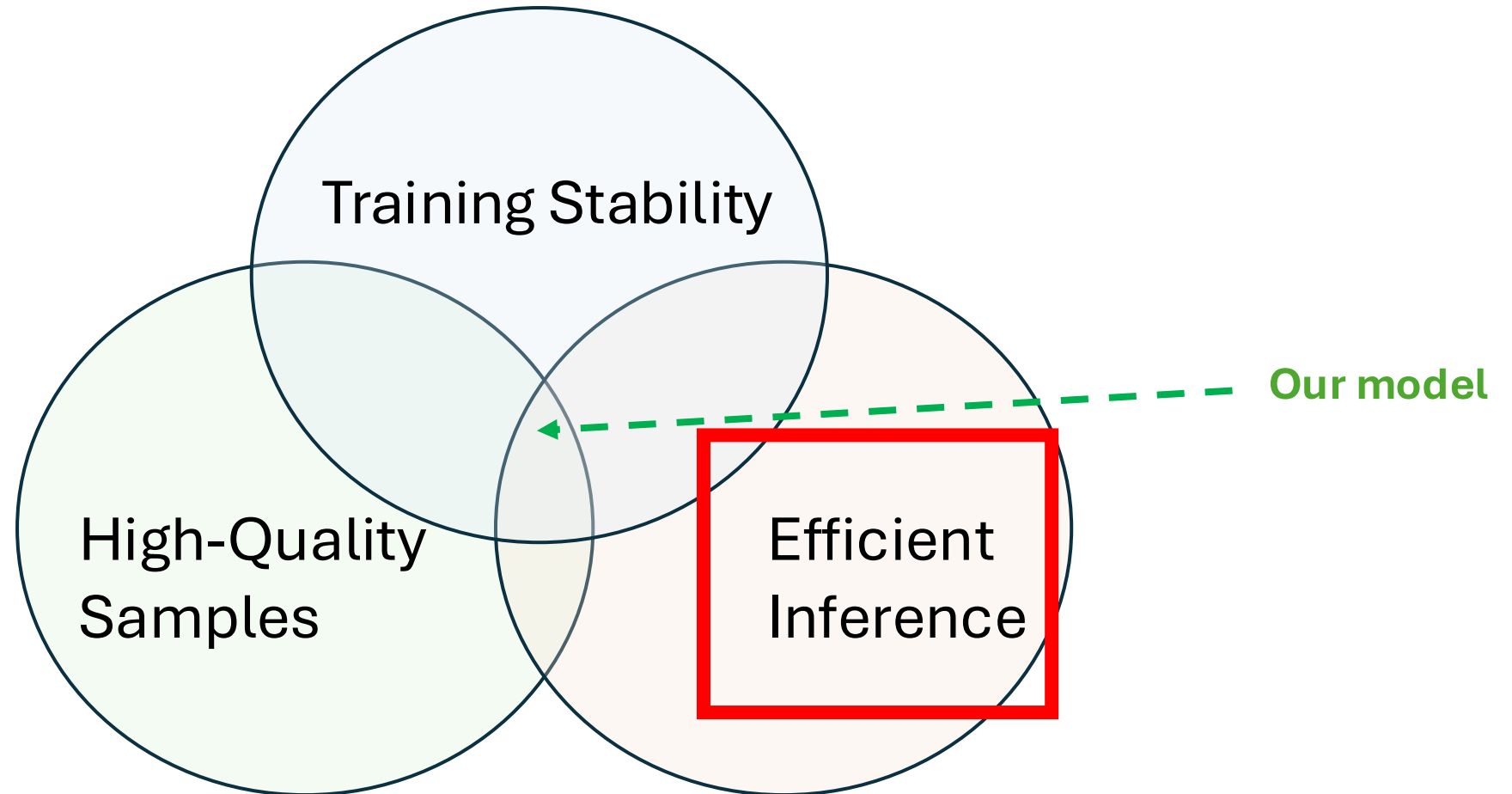**Diffusion Pretrain + Step Distillation**

**Problems:**

- Can be unstable

- Too many models and additional engineering complexities

**Single-Stage Pretraining**

- Consistency Training

- MeanFlow

- Inductive Moment Matching

- Terminal Velocity Matching

This talk

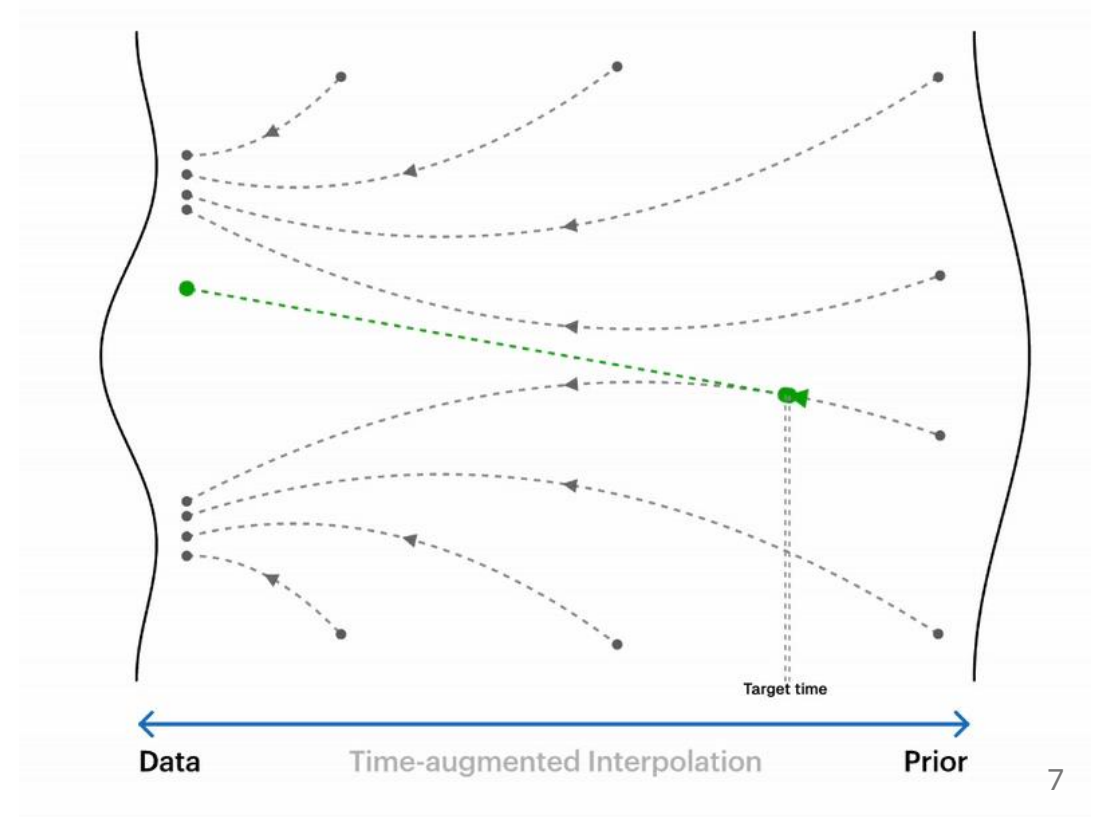# Desiderata of Efficient Inference-Time Scaling

# Problems with Diffusion Inference

- Want: **Large** jump in timesteps (**NOT** infinitesimally tiny flow ODE)
- Denoising Diffusion Implicit Models (DDIM)
  - Euler under FM schedule

$$\mathbf{x}_s = \mathbf{x}_t + \hat{\mathbf{u}} \cdot (\boxed{s} - t)$$
$$\hat{\mathbf{u}} = \hat{\mathbf{u}}(\mathbf{x}_t; t)$$

- <span style="color:red">Linear w.r.t. $s$</span>



Data     Time-augmented Interpolation     Prior

Target time

# Fixing the Capacity Issue

- Inject $s$ into the network
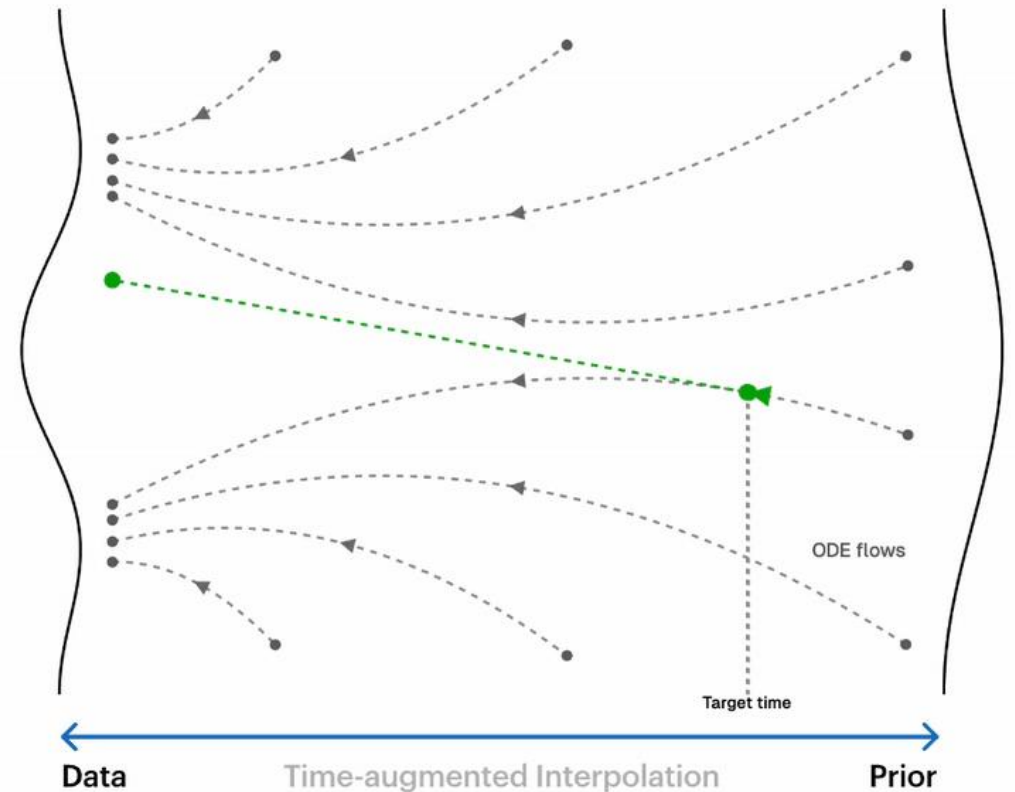
$$\mathbf{x}_s = \mathbf{x}_t + \hat{\mathbf{u}} \cdot (s - t)$$

Before: $\hat{\mathbf{u}} = \hat{\mathbf{u}}(\hat{\mathbf{x}}_t; t)$    After: $\hat{\mathbf{u}} = \hat{\mathbf{u}}(\hat{\mathbf{x}}_t; t, s)$

- Covers complex solutions
  - ODE integration

- Can perform large time jump (speeds up sampling)



ODE flows

Target time

Data          Time-augmented Interpolation          Prior

# Desiderata of Efficient Inference-Time Scaling

# Inductive Moment Matching

A Distribution-Matching Few-Step Method

# Inductive Moment Matching

- Key components

**Sample-based Distribution Matching**.
- Maximum Mean Discrepancy (MMD)

**Inductive Learning**.
- Mathematical induction to match the model's own distribution

# 1. Sample-based Distribution Matching

Data ←——————————————————————————————→ Prior



$$q_s(\mathbf{x}_s)$$

$$q_t(\mathbf{x}_t)$$

**A naïve approach:**

Match in distribution

$$p_{s,t}^{\theta}(\mathbf{x}_s)$$

DDIM (w/ 2 timesteps)
$$\mathbf{x}_s = \mathbf{x}_t + \hat{\mathbf{u}}_\theta(s - t)$$

We learn this DDIM mapping

- Use MMD due to its stability

# Maximum Mean Discrepancy

**Advantages:**

- No adversarial training:
  - GAN-like, **optimal discriminator** chosen in RKHS

- Standard kernel functions:
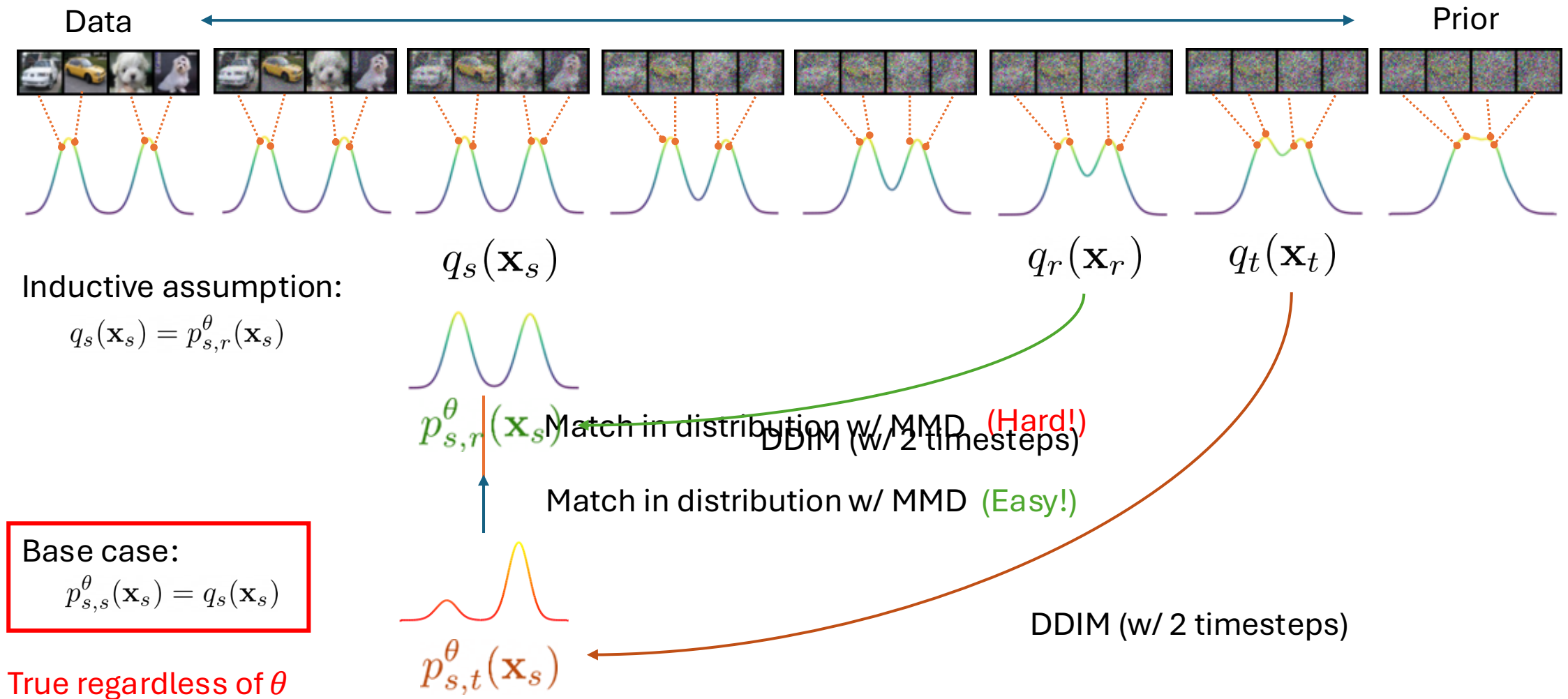  - RBF, Laplace, etc. match **all moments**

**Empirical implementation:**

- Multiple particles to estimate expectation

# 2. Inductive Learning



Data ← — — — — — — — — — — — — → Prior

$q_s(\mathbf{x}_s)$          $q_r(\mathbf{x}_r)$    $q_t(\mathbf{x}_t)$

Inductive assumption:

$q_s(\mathbf{x}_s) = p_{s,r}^{\theta}(\mathbf{x}_s)$

$p_{s,r}^{\theta}(\mathbf{x}_s)$ Match in distribution w/ MMD  (Hard!)

DDIM (w/ 2 timesteps)

Match in distribution w/ MMD  (Easy!)

Base case:
$p_{s,s}^{\theta}(\mathbf{x}_s) = q_s(\mathbf{x}_s)$

$p_{s,t}^{\theta}(\mathbf{x}_s)$

DDIM (w/ 2 timesteps)

True regardless of $\theta$

# An Intuition on Inductive Learning



Data ←——————————————————————————————→ Prior

$q_s(\mathbf{x}_s)$   $q_t(\mathbf{x}_t)$

$s \le r < t,$
$s = r$ when $t$ is close to $s$

$p_{s,r}^{\theta}(\mathbf{x}_s) = p_{s,s}^{\theta}(\mathbf{x}_s)$

Base case:
$p_{s,s}^{\theta}(\mathbf{x}_s) = q_s(\mathbf{x}_s)$

True regardless of $\theta$

$p_{s,t}^{\theta}(\mathbf{x}_s)$

# An Intuition on Inductive Learning

Prior



$q_s(\mathbf{x}_s)$  $q_r(\mathbf{x}_r)$  $q_t(\mathbf{x}_t)$

Inductive assumption:

$$q_s(\mathbf{x}_s) = p_{s,r}^{\theta}(\mathbf{x}_s)$$

$$p_{s,r}^{\theta}(\mathbf{x}_s)$$

Base case:

$$p_{s,s}^{\theta}(\mathbf{x}_s) = q_s(\mathbf{x}_s)$$

True regardless of $\theta$

$$p_{s,t}^{\theta}(\mathbf{x}_s)$$

# An Intuition on Inductive Learning



Data ⟷ Prior

$q_s(\mathbf{x}_s)$　　　$q_r(\mathbf{x}_r)$　　　$q_t(\mathbf{x}_t)$

Inductive assumption:
$$q_s(\mathbf{x}_s) = p_{s,r}^\theta(\mathbf{x}_s)$$

$p_{s,r}^\theta(\mathbf{x}_s)$

Base case:
$$p_{s,s}^\theta(\mathbf{x}_s) = q_s(\mathbf{x}_s)$$

True regardless of $\theta$

$p_{s,t}^\theta(\mathbf{x}_s)$

# An Intuition on Inductive Learning



Data ⟵⟶ Prior

$q_s(\mathbf{x}_s)$  $q_r(\mathbf{x}_r)$  $q_t(\mathbf{x}_t)$

Inductive assumption:
$$q_s(\mathbf{x}_s) = p_{s,r}^\theta(\mathbf{x}_s)$$

$p_{s,r}^\theta(\mathbf{x}_s)$

Base case:
$$p_{s,s}^\theta(\mathbf{x}_s) = q_s(\mathbf{x}_s)$$

True regardless of $\theta$

$p_{s,t}^\theta(\mathbf{x}_s)$

# Stable Training

- Stable training as long as >=4 particles

- Consistency training is a 1-particle special case



Inductive Moment Matching

Consistency Model

100K

50K

Training Iterations

Training Iterations
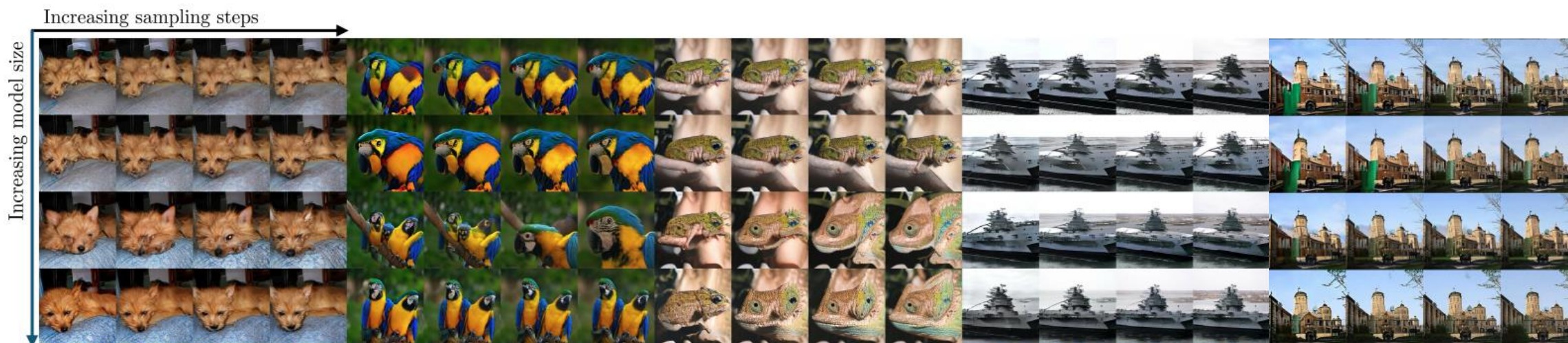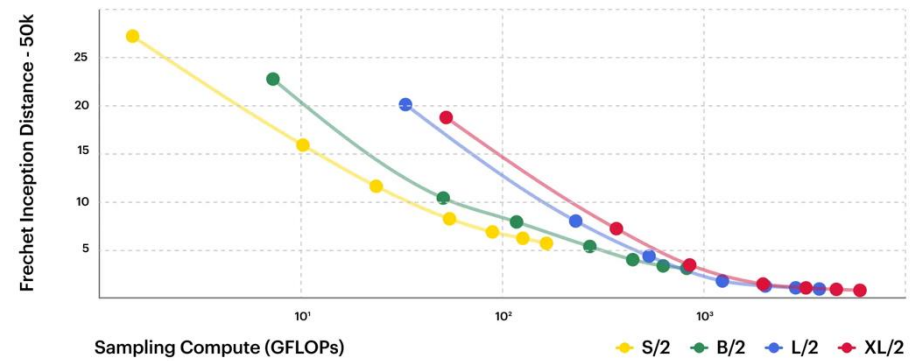
19

# Image Generation

- ImageNet-256x256

|  | FID |
| --- | --- |
| DiT         (250-step) | 2.27 |
| SiT         (250-step) | 2.15 |
| VAR-d20 | 2.57 |
| VAR-d30 | 1.92 |
| **IMM**      **(8-step)** | **1.99** |
|              **(16-step)** | **1.90** |

# Scaling Property

# Relation to Other Works

- Consistency Training (2023)
  - Particle number = 1, L2 kernel

- Generative Moment Matching Networks (2015)
  - t = 1, s=r=0

- Generative Modeling via Drifting (2026)
  - Drift field parameterized via MMD attraction + repulsion
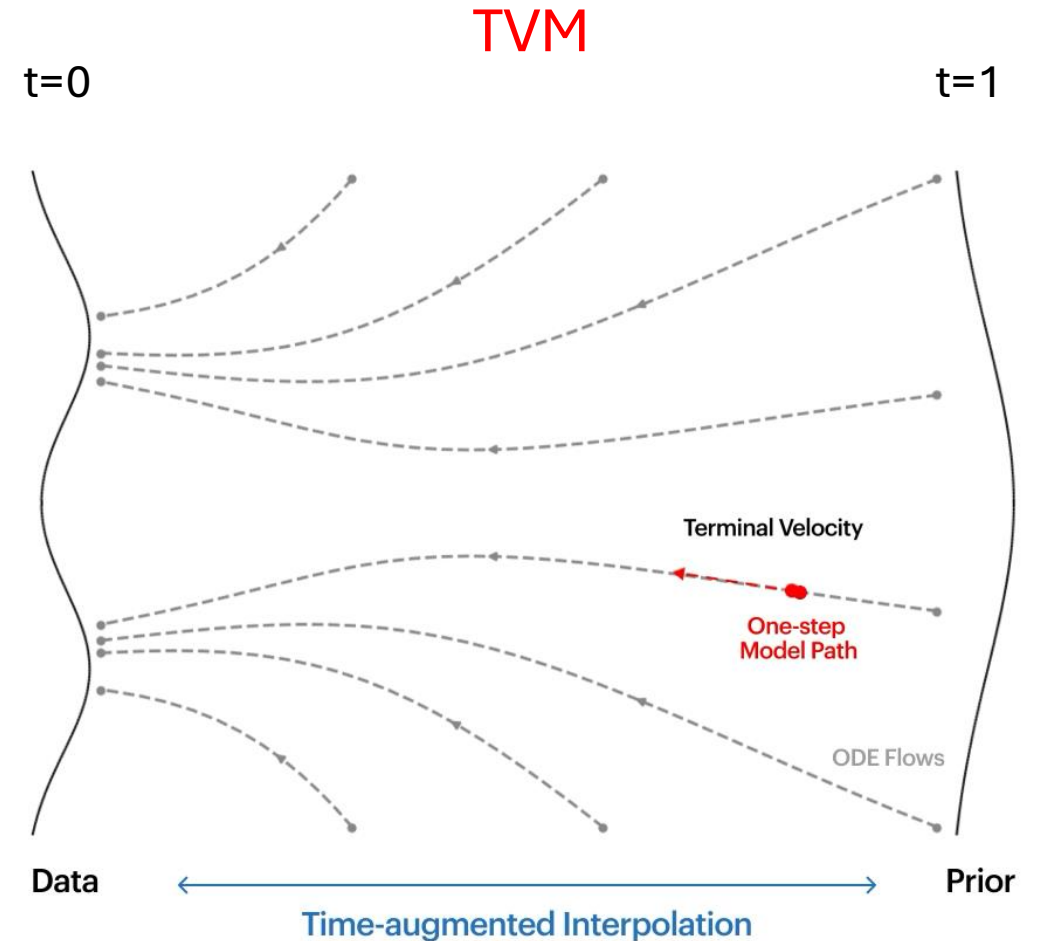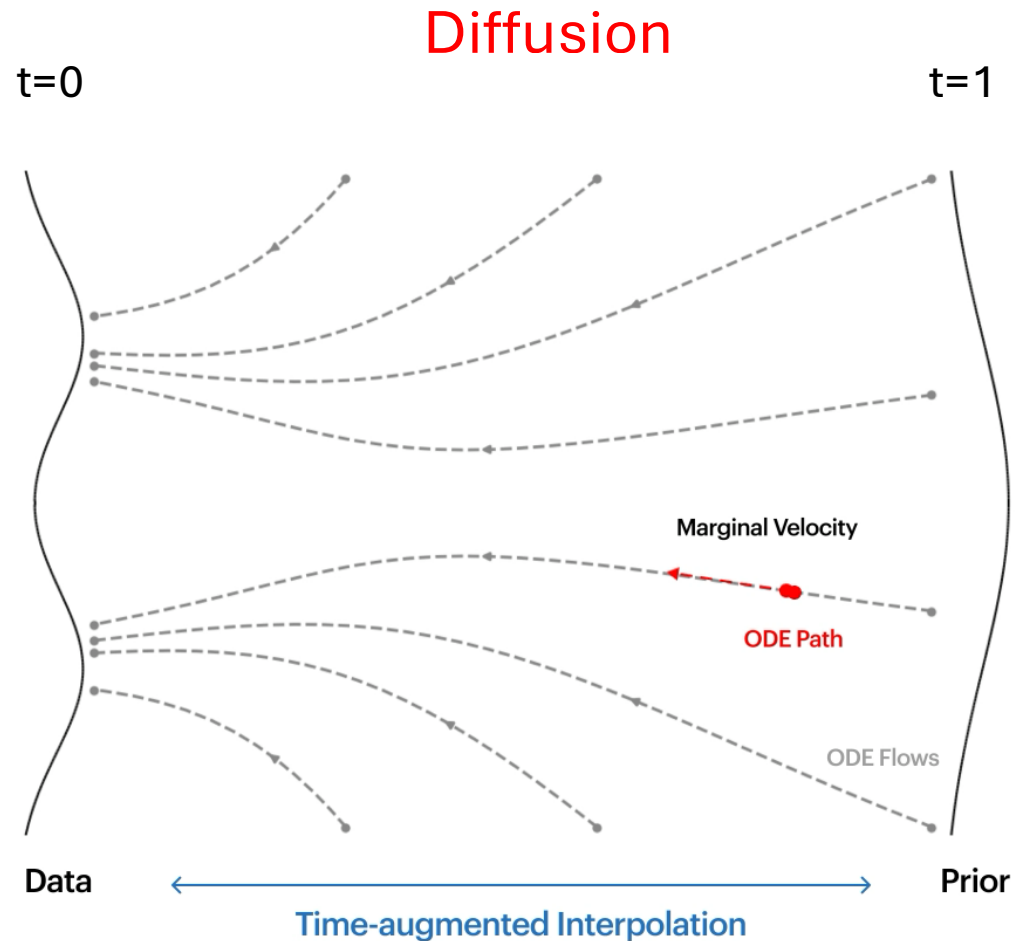
# Limitations of IMM

- Multi-sample objective
  - => difficult to scale to high-dimension

- Mapping function $r(s, t)$ requires high precision (e.g. FP16)
  - => Large scale models require BF16 or lower

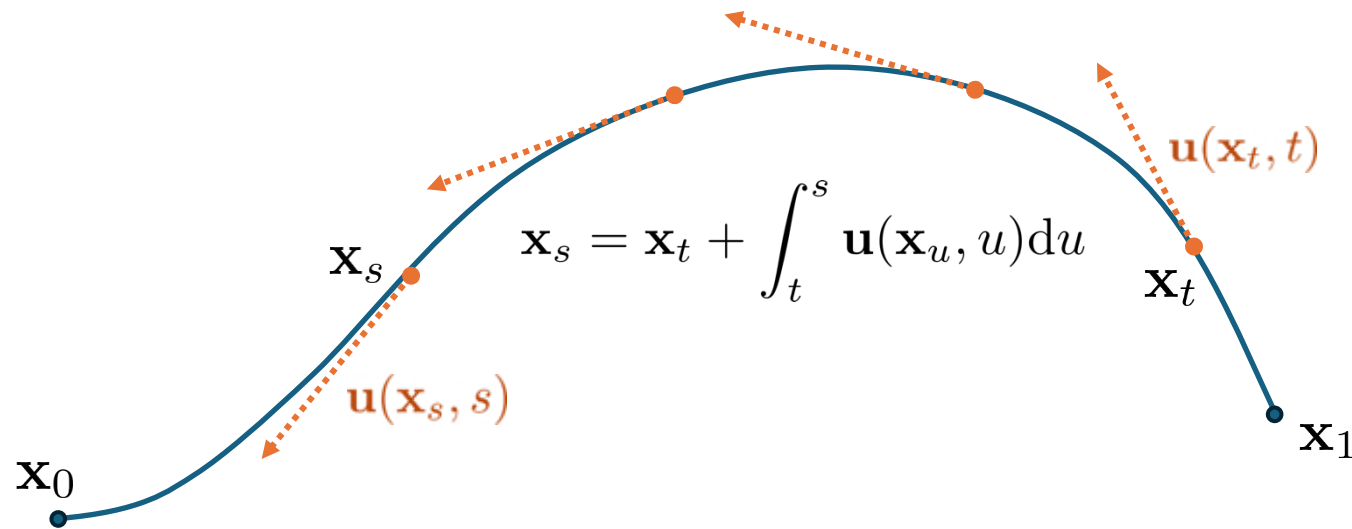# Terminal Velocity Matching

Going Back to Flows

# Intuition for Terminal Velocity Matching



**Diffusion**

t=0                                                                    t=1

Marginal Velocity

ODE Path

ODE Flows

Data  ←————————————————→  Prior

Time-augmented Interpolation

**TVM**

t=0                                                                    t=1

Terminal Velocity

One-step
Model Path

ODE Flows

Data  ←————————————————→  Prior

Time-augmented Interpolation

# Learning One-Step Trajectory Mapping

$$\mathbf{f}(\mathbf{x}_t, t, s) = \int_t^s \mathbf{u}(\mathbf{x}_r, r)\mathrm{d}r$$

Diffusion/FM

$$\mathbf{x}_s = \mathbf{x}_t + \int_t^s \mathbf{u}(\mathbf{x}_u, u)\mathrm{d}u$$

$\mathbf{u}(\mathbf{x}_t, t)$

$\mathbf{x}_s$

$\mathbf{u}(\mathbf{x}_s, s)$

$\mathbf{x}_t$

$\mathbf{x}_1$

$\mathbf{x}_0$

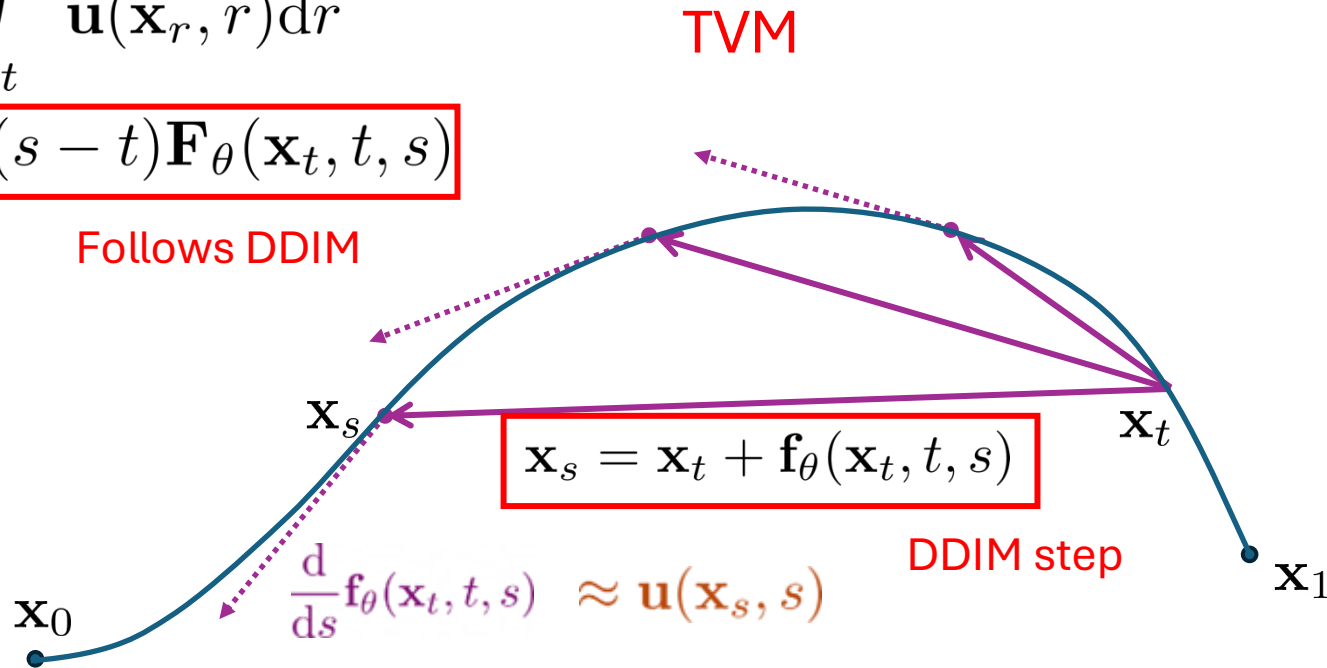# Learning One-Step Trajectory Mapping

$$\mathbf{f}(\mathbf{x}_t, t, s) = \int_t^s \mathbf{u}(\mathbf{x}_r, r)\mathrm{d}r$$

$$\mathbf{f}_\theta(\mathbf{x}_t, t, s) = \boxed{(s-t)\mathbf{F}_\theta(\mathbf{x}_t, t, s)}$$

Follows DDIM

TVM

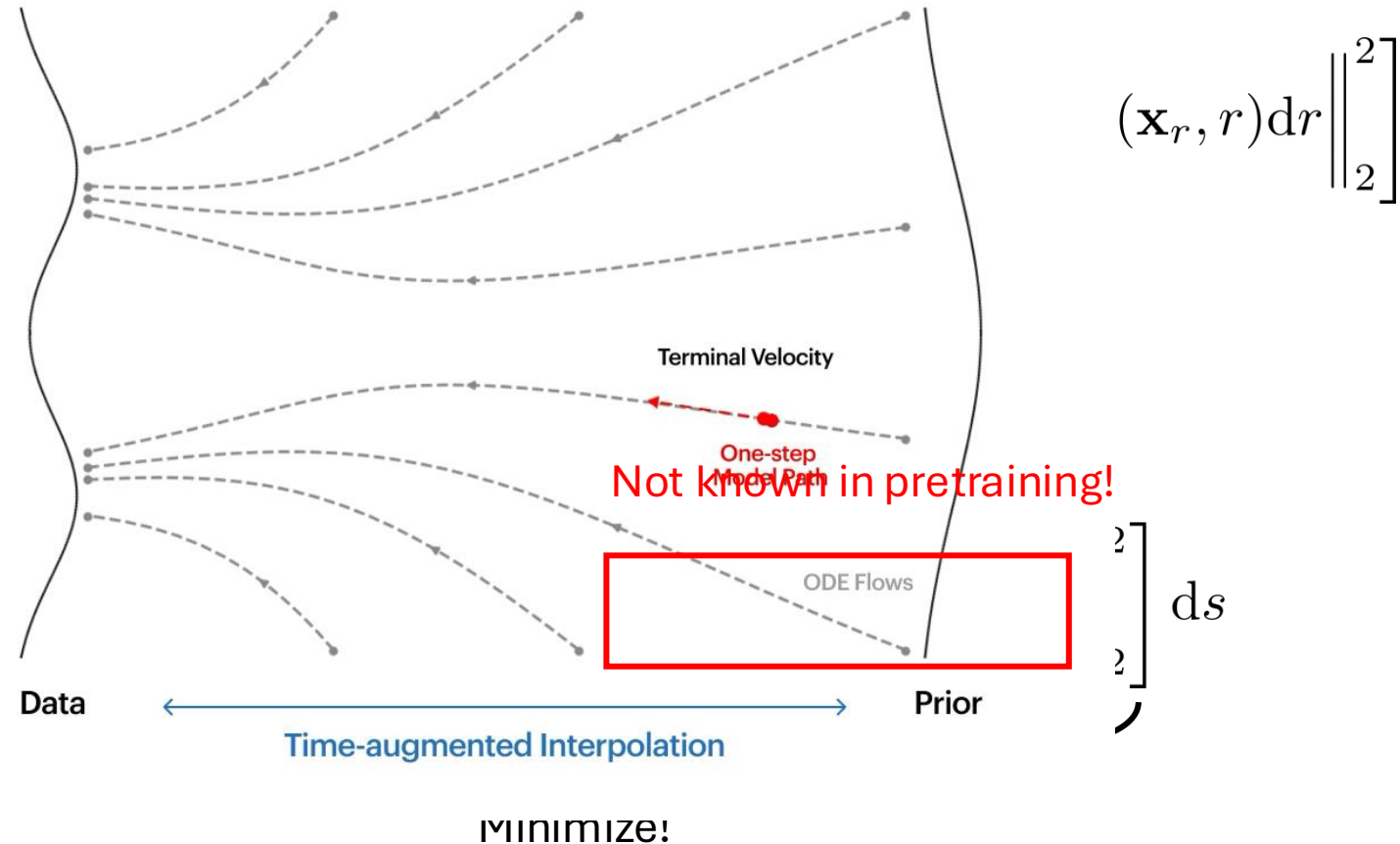$$\boxed{\mathbf{x}_s = \mathbf{x}_t + \mathbf{f}_\theta(\mathbf{x}_t, t, s)}$$

DDIM step

$$\frac{\mathrm{d}}{\mathrm{d}s}\mathbf{f}_\theta(\mathbf{x}_t, t, s) \approx \mathbf{u}(\mathbf{x}_s, s)$$

$\mathbf{x}_s$

$\mathbf{x}_t$

$\mathbf{x}_1$

$\mathbf{x}_0$

Naively: Match Displacement!

Terminal Velocity Condition

$$\mathcal{L}_{\mathrm{disp}}^t \frac{\mathrm{d}}{\mathrm{d}s}(\mathbf{f})(\mathbf{x}_t, \mathbf{F}t, \mathbf{x}_t s)\left[\left\|\mathbf{f}_\theta(\mathbf{u}(\mathbf{x}_t, 0) - \mathbf{f}(\int_t^0 \mathbf{x}_t, \mathbf{u}(\mathbf{x}s), r, s)\mathrm{d}r\right\|_2^2\right]$$

27

# Learning One-Step Trajectory Mapping

- Displacemen

- Terminal Velc

$$(\mathbf{x}_r, r)\mathrm{d}r \Big\|_2^2 \Big]$$

**Terminal Velocity**

One-step
Model Path

<span style="color:red">Not known in pretraining!</span>

ODE Flows

$$\mathcal{L}_{\mathrm{displ}}^t(\theta) \leq \qquad \qquad \Big] \mathrm{d}s$$

**Data**

**Prior**

Time-augmented Interpolation

Minimize!

# Proxy for Ground-Truth

- Just use network itself as approximation

$$\mathbf{u}_\theta(\mathbf{x}_t + \mathbf{f}_\theta(\mathbf{x}_t, t, s), s) \approx \mathbf{u}(\mathbf{x}_t + \mathbf{f}(\mathbf{x}_t, t, s), s)$$

Terrible approximation at start of training!

- Solution: Train such that $\mathbf{u}_\theta(\mathbf{x}_t, t) \approx \mathbf{u}(\mathbf{x}_t, t)$

Flow Matching

- Final Objective:

$$\mathcal{L}_{\text{TVM}}^{t,s}(\theta) = \mathbb{E}_{\mathbf{x}_t, \mathbf{x}_s, \mathbf{v}_s}\left[\left\|\frac{\mathrm{d}}{\mathrm{d}s}\mathbf{f}_\theta(\mathbf{x}_t, t, s) - \mathbf{u}_\theta(\mathbf{x}_t + \mathbf{f}_\theta(\mathbf{x}_t, t, s), s)\right\|_2^2 + \left\|\mathbf{u}_\theta(\mathbf{x}_s, s) - \mathbf{v}_s\right\|_2^2\right]$$

# Relation to Distribution Divergence

- TVM loss bounds Wasserstein distance up to a constant

Network Lipschitzness

$$W_2^2(\mathbf{f}_{t\to 0}^\theta \# p_t, p_0) \leq \int_0^t \lambda[L](s)\mathcal{L}_{\text{TVM}}^{t,s}(\theta)\mathrm{d}s + C,$$

- **Implication:** diffusion transformers are **NOT Lipschitz-continuous**. Need "semi-Lipschitz control"!
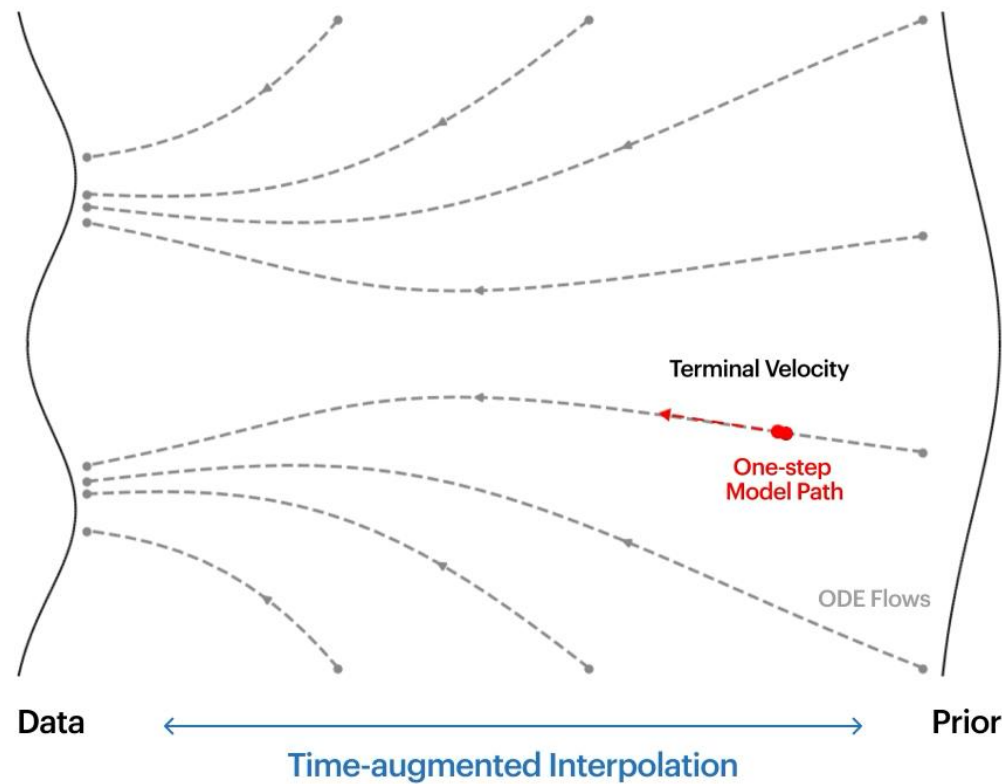
  - LayerNorm -> RMSNorm

  - QKNorm w/ RMSNorm

  - Add RMSNorm without parameters to time modulation

# Calculating Terminal Velocity

- How to calc

- Recall:

$$\frac{\mathrm{d}}{\mathrm{d}s}$$



Terminal Velocity

One-step
Model Path

ODE Flows

Data

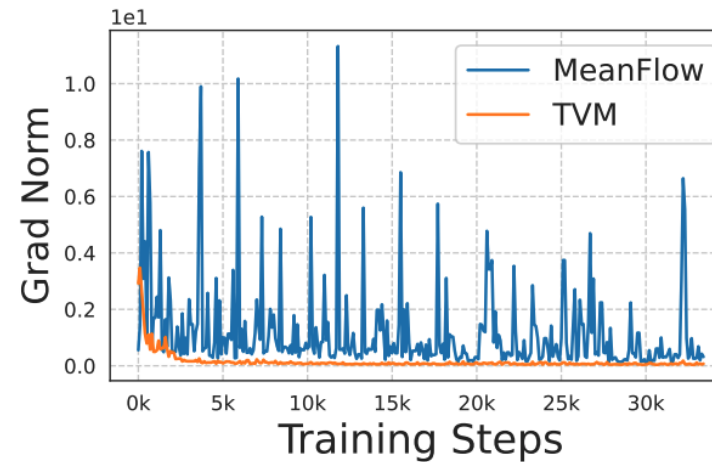Time-augmented Interpolation

Prior

(JVP)

:kwards

# Stable Training vs. MeanFlow

- Well-conditioned gradient profile



- Well-conditioned norm $\mathbf{u}_\theta(\mathbf{x}_t + \mathbf{f}_\theta(\mathbf{x}_t, t, s), s)$ w/ random CFG

# ImageNet Generation

- SOTA 1-NFE
- Outperforms DiT, SiT in 4 NFE
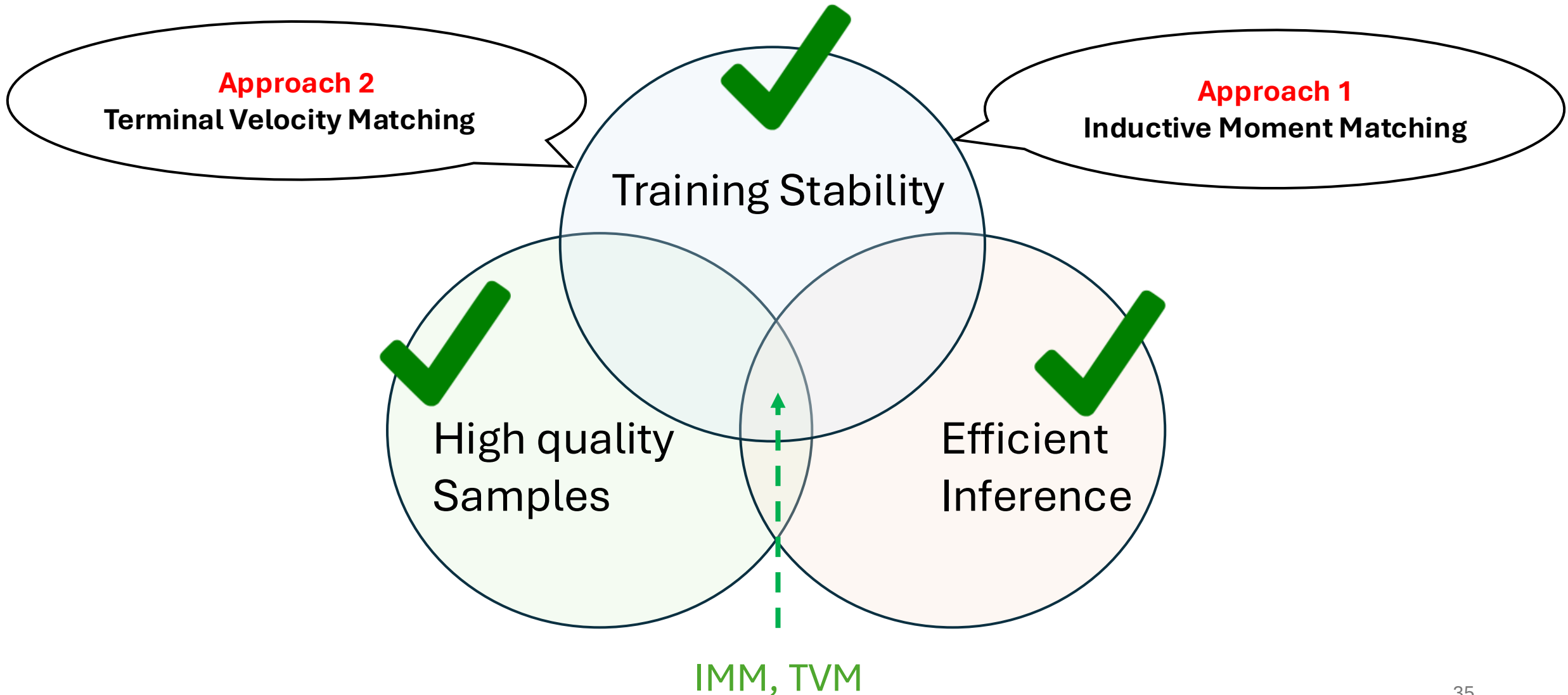


One-step samples

# TVM at 10B+ Scale
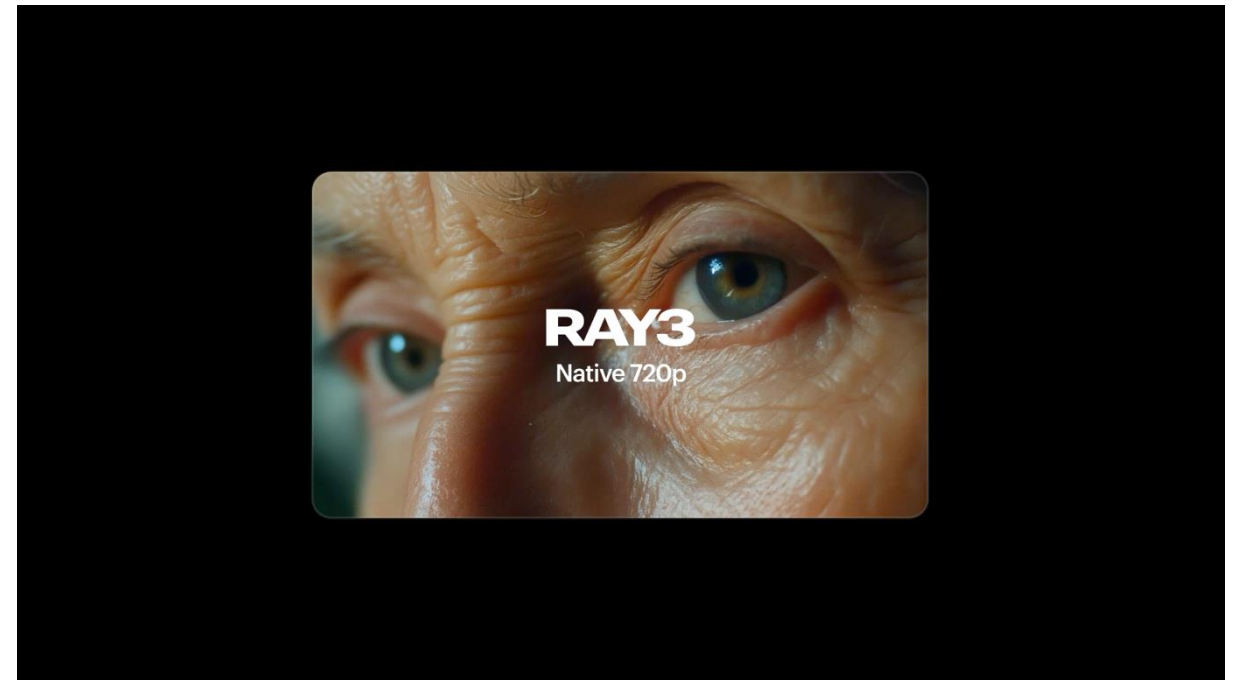
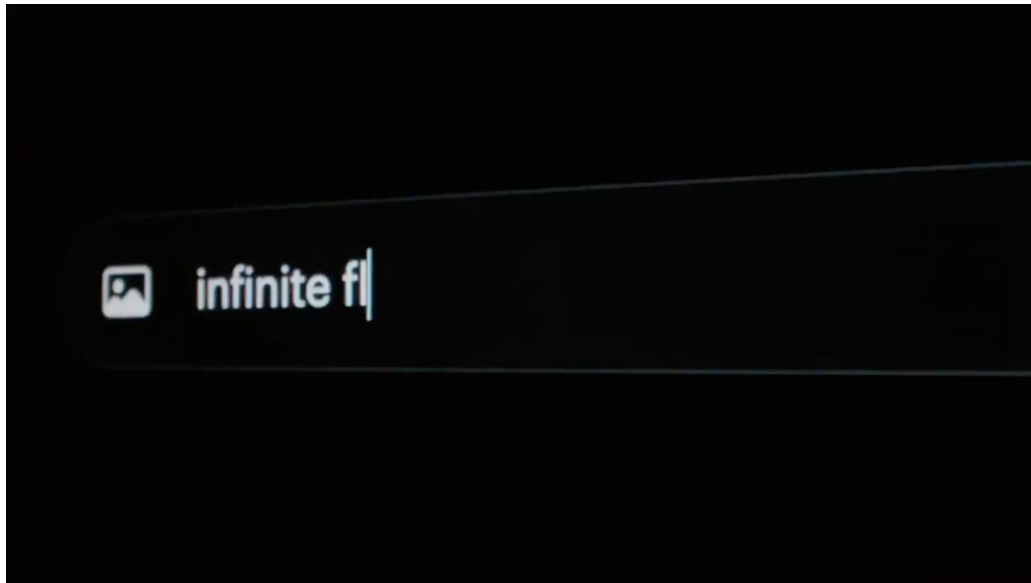4-NFE
TVM
on T2I

- Challenges at Scale:

  - FSDP with JVP – Solution: Wrap JVP inside each layer of FSDP.

  - Writing JVP kernel (with backwards support) for arbitrary sequence length

# Desiderata of Efficient Inference-Time Scaling



35

# Luma AI is a research and product lab aiming to build multimodal AGI.

- Recent Series-C: $900M.

- Average age in research team: ~27

- Inventors of **DDIM** & **NeRF** work here.

# If you are interested in advancing multimodal generative AI, join us!

https://lumalabs.ai/join

**Research team have multiple roles open around:**

1. Omni models (unifying understanding + generation)
2. Video / Audio models
3. Voice agents
4. World models
5. Multimodal agents
6. AI Infra

**("Internships / Residency" also available, but ideally >= 6 months and not driven by publications)**