

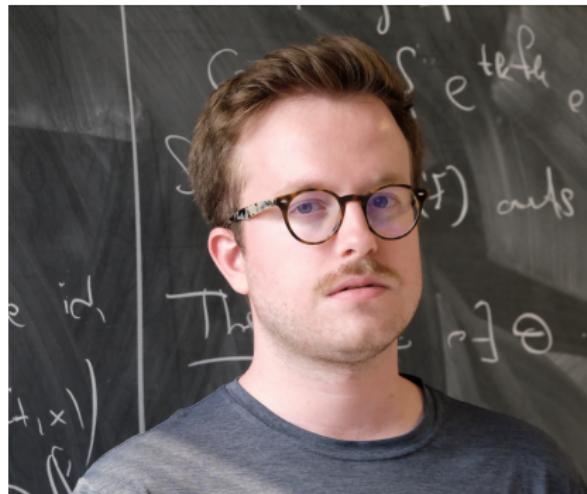
Entropy-driven control of the continuity equation for normalizing flows

Antonio Álvarez-López

FAU Erlangen–Nürnberg
Universidad Autónoma de Madrid

September 29, 2025

Joint work with...



Borjan Geshkovski
Inria Paris — Sorbonne Université (LJLL).



Domènec Ruiz-Balet
Universidad de Barcelona.

Outline

- 1 Motivation
- 2 Problem
- 3 Main result
- 4 Proof (sketch)
- 5 Extensions

Generative modeling

What is generative modeling?

Generative modeling

What is generative modeling?

Generate samples that follow a law... (**complex and high dimension!**)
by transforming samples from a base law via a learned map T :

Generative models

Flow-based generative models:

T = Flow map of an ODE/SDE.

- **Continuous Normalizing Flows (CNFs).** Invertible ODE flows trained by maximum likelihood. *Chen et al., 2018; Grathwohl et al., 2019.*
- **Diffusion Models.** Forward noising and reverse SDE/ODE denoiser using scores. *Ho et al., 2020; Song et al., 2021.*
- **Flow Matching.** Regress a time-dependent field to prescribed probability paths for fast few-step sampling. *Lipman et al., 2023.*

Common thread:



Generative models

Flow-based generative models:

T = Flow map of an ODE/SDE.

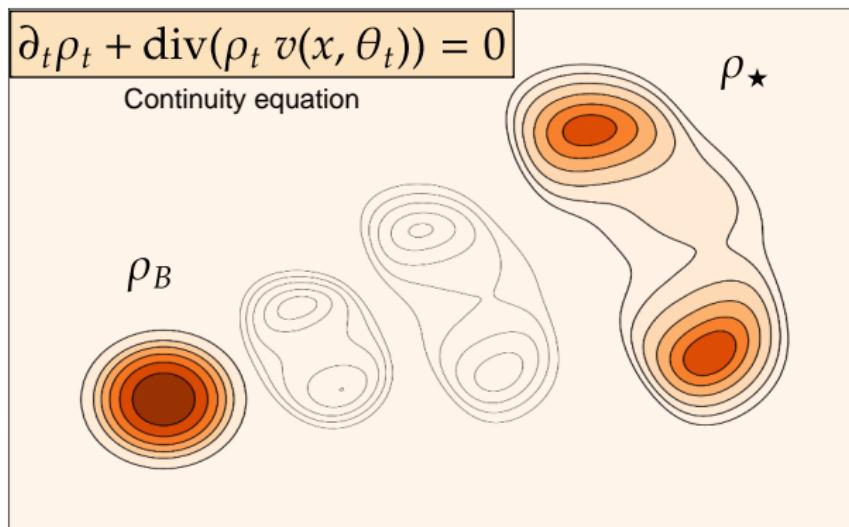
- **Continuous Normalizing Flows (CNFs).** Invertible ODE flows trained by maximum likelihood. *Chen et al., 2018; Grathwohl et al., 2019.*
- **Diffusion Models.** Forward noising and reverse SDE/ODE denoiser using scores. *Ho et al., 2020; Song et al., 2021.*
- **Flow Matching.** Regress a time-dependent field to prescribed probability paths for fast few-step sampling. *Lipman et al., 2023.*

Common thread:



Continuous Normalizing Flows

- Base $\rho_B \equiv \mathcal{N}(0, \text{Id}) \mapsto$ Target ρ_\star known or estimated from data $x_1, \dots, x_N \sim \rho_\star$.



- Fix Lipschitz control family $\{v(\cdot, \theta)\}_\theta$. Find θ and $T > 0$ such that $\rho_T^\theta \approx \rho_\star$.
- Sample $x \sim \rho_B$ follows the flow:
$$\begin{cases} \frac{d}{dt} \Phi_t^\theta(x) &= v(\Phi_t^\theta(x), \theta_t) \\ \Phi_0^\theta(x) &= x \end{cases}$$

Continuous Normalizing Flows

- **Optimization:**

$$\min_{\theta} \text{KL}(\mu_* \| \mu_T^\theta)$$

$$\begin{aligned}\mu_T^\theta &= (\Phi_T^\theta)_\# \mu_B \\ d\mu_* &= \rho_* dx \\ d\mu_T^\theta &= \rho_T^\theta dx\end{aligned}$$

where the Kullback-Leibler divergence/relative entropy is

$$\text{KL}(\mu_* \| \mu_T^\theta) := \begin{cases} \mathbb{E}_{\mu_*} \left[\log \left(\frac{d\mu_*}{d\mu_T^\theta} \right) \right] = \int_{\mathbb{R}^d} \log \left(\frac{\rho_*}{\rho_T^\theta} \right) \rho_*, & \mu_* \ll \mu_T^\theta, \\ +\infty, & \text{otherwise.} \end{cases}$$

Continuous Normalizing Flows

- Optimization:

$$\min_{\theta} \text{KL}(\mu_{\star} \| \mu_T^{\theta})$$

$$\begin{aligned}\mu_T^{\theta} &= (\Phi_T^{\theta})_{\#} \mu_B \\ d\mu_{\star} &= \rho_{\star} dx \\ d\mu_T^{\theta} &= \rho_T^{\theta} dx\end{aligned}$$

where the Kullback-Leibler divergence/relative entropy is

$$\text{KL}(\mu_{\star} \| \mu_T^{\theta}) := \begin{cases} \mathbb{E}_{\mu_{\star}} \left[\log \left(\frac{d\mu_{\star}}{d\mu_T^{\theta}} \right) \right] = \int_{\mathbb{R}^d} \log \left(\frac{\rho_{\star}}{\rho_T^{\theta}} \right) \rho_{\star}, & \mu_{\star} \ll \mu_T^{\theta}, \\ +\infty, & \text{otherwise.} \end{cases}$$

- Why KL? Max. likelihood \iff Min. KL divergence:

Empirical MLE: $\max_{\theta} \mathbb{P}_{\theta}(x_1, \dots, x_n) \stackrel{\text{i.i.d.}}{=} \max_{\theta} \prod_{i=1}^n \rho_T^{\theta}(x_i) \iff \max_{\theta} \frac{1}{n} \sum_{i=1}^n \log \rho_T^{\theta}(x_i)$

Population objective: $\text{KL}(\mu_{\star} \| \mu_T^{\theta}) = \underbrace{\mathbb{E}_{\mu_{\star}} [\log \rho_{\star}]}_{\text{indep. of } \theta} - \mathbb{E}_{\mu_{\star}} [\log \rho_T^{\theta}]$.

Continuous Normalizing Flows

We can view it as an **approximate controllability problem**:

Given $\varepsilon > 0$, find $\theta = (\theta_t)_{t \in [0, T]}$ such that

$$\text{KL}(\mu_\star \| \mu_T^\theta) < \varepsilon, \quad \text{s.t.} \quad \begin{cases} \partial_t \rho_t^\theta(x) + \text{div}(\rho_t^\theta(x) \nu(x; \theta_t)) = 0, \\ \rho_0(x) = \rho_B(x). \end{cases}$$

Continuous Normalizing Flows

We can view it as an **approximate controllability problem**:

Given $\varepsilon > 0$, find $\theta = (\theta_t)_{t \in [0, T]}$ such that

$$\text{KL}(\mu_\star \| \mu_T^\theta) < \varepsilon, \quad \text{s.t.} \quad \begin{cases} \partial_t \rho_t^\theta(x) + \operatorname{div}(\rho_t^\theta(x) \nu(x; \theta_t)) = 0, \\ \rho_0(x) = \rho_B(x). \end{cases}$$

- **Difficulty 1: Not a distance**

KL is non-symmetric + no triangle inequality!

Continuous Normalizing Flows

We can view it as an **approximate controllability problem**:

Given $\varepsilon > 0$, find $\theta = (\theta_t)_{t \in [0, T]}$ such that

$$\text{KL}(\mu_\star \| \mu_T^\theta) < \varepsilon, \quad \text{s.t.} \quad \begin{cases} \partial_t \rho_t^\theta(x) + \operatorname{div}(\rho_t^\theta(x) \nu(x; \theta_t)) = 0, \\ \rho_0(x) = \rho_B(x). \end{cases}$$

- **Difficulty 1: Not a distance**

KL is non-symmetric + no triangle inequality!

- **Difficulty 2: Blow-up.** We have $\text{KL} \rightarrow +\infty$ whenever

- **Smaller support.** $\exists A \subset \mathbb{R}^d$ with $\mu_\star(A) > 0$ and $\mu_T^\theta(A) = 0$.
- **Relative decay.** $\mathbb{E}_{\mu_\star} [\log \rho_T^\theta(X)] \rightarrow +\infty$.

(Example) $\rho_\star(x) = \mathbf{1}_{(0,1)}(x)$, $\rho_T^\theta(x) = \frac{e^{-1/x}}{Z} \mathbf{1}_{(0,1)}(x)$, $\text{KL} = \log Z + \int_0^1 \frac{dx}{x} = +\infty$

Continuous Normalizing Flows

We can view it as an **approximate controllability problem**:

Given $\varepsilon > 0$, find $\theta = (\theta_t)_{t \in [0, T]}$ such that

$$\text{KL}(\mu_\star \| \mu_T^\theta) < \varepsilon \quad \text{s.t.} \quad \begin{cases} \partial_t \rho_t^\theta(x) + \operatorname{div}(\rho_t^\theta(x) v_t(x; \theta_t)) = 0, \\ \rho_0(x) = \rho_B(x). \end{cases}$$

⇒ We need control over the full support and over the relative tail decay!

- Our controllability result is valid for densities supported on \mathbb{R}^d .
- But μ_\star is an empirical measure built with n samples...

In practice: build new target $\hat{\mu}_{\star,n} \approx \mu_\star$ via (Gaussian) KDE.

Controllability of continuity/Liouville equation:

Brief review

- **Moser (1965); Dacorogna–Moser (1990)** — Mass-preserving diffeomorphisms; baseline density-to-density flows.
- **Duprez–Morancey–Rossi (2019)** — Continuity eq. with localized control; approx. controllability (crossing); exact with low regularity.
- **Duprez–Morancey–Rossi (2020)** — Minimal time for localized control; sharp mass-through-control-region condition.
- **Raginsky (2024)** — Liouville controllability; implements broad classes of diffeomorphisms for linear systems; OT links.
- **Brockett (notes)** — Ensemble controllability framework; state-dependent feedback viewpoint.
- **OT/OC link** — Benamou–Brenier (2000) dynamic OT; Elamvazhuthi (2023) dynamical OT for control-affine systems.

Control family

Neural ODEs (ReLU, one-neuron)

$$v(x, \theta_t) = w(t)(\langle a(t), x \rangle + b(t))_+, \quad \theta \in L^\infty(0, T; \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R})$$

- **Simple** control but **high** expressivity!¹

¹ Ruiz-Balet--Zuazua, 2023., Á-L.--Orive--Illera--Zuazua, 2024.

Control family

Neural ODEs (ReLU, one-neuron)

$$v(x, \theta_t) = w(t)(\langle a(t), x \rangle + b(t))_+, \quad \theta \in L^\infty(0, T; \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R})$$

- Simple control but high expressivity!¹
- The solution of

$$\partial_t \rho_t^\theta(x) + \operatorname{div} \left(\rho_t^\theta(x) w(\langle a, x \rangle + b)_+ \right) = 0, \quad \rho_0 = \rho_B$$

has closed-form (by the method of characteristics):

$$\rho_t^\theta(x) = \begin{cases} \rho_B(x), & x \in H^- := \{\langle a, x \rangle + b \leq 0\}, \\ e^{-t \langle w, a \rangle} \rho_B(A_t(x)), & x \in H^+ := \{\langle a, x \rangle + b > 0\}, \end{cases}$$

with

$$A_t(x) = e^{-t w a^\top} x - \frac{b}{\langle w, a \rangle} (1 - e^{-t \langle w, a \rangle}) w \rightarrow \text{Affine in } x.$$

¹ Ruiz-Balet-Zuazua, 2023., Á.-L.-Orive--Illera--Zuazua, 2024.

Control family

Neural ODEs (ReLU, one-neuron)

$$v(x, \theta_t) = w(t)(\langle a(t), x \rangle + b(t))_+, \quad \theta \in L^\infty(0, T; \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R})$$

- Simple control but high expressivity!¹
- The solution of

$$\partial_t \rho_t^\theta(x) + \operatorname{div} \left(\rho_t^\theta(x) w(\langle a, x \rangle + b)_+ \right) = 0, \quad \rho_0 = \rho_B$$

has closed-form (by the method of characteristics):

$$\rho_t^\theta(x) = \begin{cases} \rho_B(x), & x \in H^- := \{\langle a, x \rangle + b \leq 0\}, \\ e^{-t \langle w, a \rangle} \rho_B(A_t(x)), & x \in H^+ := \{\langle a, x \rangle + b > 0\}, \end{cases}$$

with

$$A_t(x) = e^{-t w a^\top} x - \frac{b}{\langle w, a \rangle} (1 - e^{-t \langle w, a \rangle}) w \rightarrow \text{Affine in } x.$$

Degenerate case. If $\langle w, a \rangle = 0$, then $A_t(x) = (I - t w a^\top)x - t b w$.

¹ Ruiz-Balet--Zuazua, 2023., Á-L.--Orive--Illera--Zuazua, 2024.

Main result

Theorem

Let $\mu_B = \mathcal{N}(m_B, \Sigma_B)$ and a probability density ρ_* s.t. for some $M, \sigma_* > 0$

$$\rho_*(x) \leq \rho_*(x) := \frac{1}{(2\pi\sigma_*^2)^{d/2}} e^{-\|x\|^2/(2\sigma_*^2)} \quad (\|x\| \geq M).$$

Then for all $T > 0, \varepsilon > 0$ there exists a piecewise-constant θ such that

$$\text{KL}(\mu_* \| \mu_T^\theta) = \int \rho_* \log \frac{\rho_*}{\rho_T^\theta} < \varepsilon.$$

Main result

Theorem

Let $\mu_B = \mathcal{N}(m_B, \Sigma_B)$ and a probability density ρ_\star s.t. for some $M, \sigma_\bullet > 0$

$$\rho_\star(x) \leq \rho_\bullet(x) := \frac{1}{(2\pi\sigma_\bullet^2)^{d/2}} e^{-\|x\|^2/(2\sigma_\bullet^2)} \quad (\|x\| \geq M).$$

Then for all $T > 0, \varepsilon > 0$ there exists a piecewise-constant θ such that

$$\text{KL}(\mu_\star \| \mu_T^\theta) = \int \rho_\star \log \frac{\rho_\star}{\rho_T^\theta} < \varepsilon.$$

Number of discontinuities: We can choose

$$N_d(\varepsilon) \leq (d + 10) \left\lceil \frac{2R_\varepsilon}{h_\varepsilon} \right\rceil^d + 2d$$

where $[-R_\varepsilon, R_\varepsilon]^d$ contains $\geq 1 - \varepsilon$ mass of both μ_B and μ_\star ; and h_ε mesh size of uniform grid on which piecewise-constant interpolants have L^1 error $\leq \varepsilon$.

Main result

Theorem

Let $\mu_B = \mathcal{N}(m_B, \Sigma_B)$ and a probability density ρ_* s.t. for some $M, \sigma_* > 0$

$$\rho_*(x) \leq \rho_*(x) := \frac{1}{(2\pi\sigma_*^2)^{d/2}} e^{-\|x\|^2/(2\sigma_*^2)} \quad (\|x\| \geq M).$$

Then for all $T > 0, \varepsilon > 0$ there exists a piecewise-constant θ such that

$$\text{KL}(\mu_* \| \mu_T^\theta) = \int \rho_* \log \frac{\rho_*}{\rho_T^\theta} < \varepsilon.$$

Number of discontinuities: We can choose

$$N_d(\varepsilon) \leq (d+10) \underbrace{\left[\frac{2R_\varepsilon}{h_\varepsilon} \right]^d}_{\text{Number of cells}} + 2d$$

where $[-R_\varepsilon, R_\varepsilon]^d$ contains $\geq 1 - \varepsilon$ mass of both μ_B and μ_* ; $h_\varepsilon \equiv$ mesh size of a uniform grid on which Riemann approximation has L^1 -error $\leq \varepsilon$.

Proof idea: (i) + (ii) + (iii)

(i) **TV-controllability** (Prior work, *Ruiz-Balet--Zuazua, 2024*):

$$[0, T/2] : \quad |\mu_{T/2}^\theta - \mu_\star|_{\text{TV}} = \frac{1}{2} \|\rho_{T/2}^\theta - \rho_\star\|_{L^1} < \varepsilon \quad (\text{constructive } \theta).$$

Proof idea: (i) + (ii) + (iii)

(i) **TV-controllability** (Prior work, *Ruiz-Balet--Zuazua, 2024*):

$$[0, T/2] : \quad |\mu_{T/2}^\theta - \mu_\star|_{\text{TV}} = \frac{1}{2} \|\rho_{T/2}^\theta - \rho_\star\|_{L^1} < \varepsilon \quad (\text{constructive } \theta).$$

(ii) **Positivity/Tail control:**

$$[T/2, T] : \quad \rho_T^\theta > 0 \quad \text{and} \quad \lim_{\|x\| \rightarrow \infty} \frac{\rho_\star(x)}{\rho_T^\theta(x)} < \infty \quad (\text{bounded Radon--Nikodym derivative})$$

Proof idea: (i) + (ii) + (iii)

(i) **TV-controllability** (Prior work, *Ruiz-Balet--Zuazua, 2024*):

$$[0, T/2] : \quad |\mu_{T/2}^\theta - \mu_\star|_{\text{TV}} = \frac{1}{2} \|\rho_{T/2}^\theta - \rho_\star\|_{L^1} < \varepsilon \quad (\text{constructive } \theta).$$

(ii) **Positivity/Tail control:**

$$[T/2, T] : \quad \rho_T^\theta > 0 \quad \text{and} \quad \lim_{\|x\| \rightarrow \infty} \frac{\rho_\star(x)}{\rho_T^\theta(x)} < \infty \quad (\text{bounded Radon--Nikodym derivative})$$

(iii) Pinsker inequality: $|\mu_\star - \mu_T^\theta|_{\text{TV}} \leq \sqrt{2 \cdot \text{KL}(\mu_\star \| \mu_T^\theta)}$ useless here.

Proof idea: (i) + (ii) + (iii)

(i) **TV-controllability** (Prior work, *Ruiz-Balet--Zuazua, 2024*):

$$[0, T/2] : \quad |\mu_{T/2}^\theta - \mu_\star|_{\text{TV}} = \frac{1}{2} \|\rho_{T/2}^\theta - \rho_\star\|_{L^1} < \varepsilon \quad (\text{constructive } \theta).$$

(ii) **Positivity/Tail control:**

$$[T/2, T] : \quad \rho_T^\theta > 0 \quad \text{and} \quad \lim_{\|x\| \rightarrow \infty} \frac{\rho_\star(x)}{\rho_T^\theta(x)} < \infty \quad (\text{bounded Radon--Nikodym derivative})$$

(iii) Pinsker inequality: $|\mu_\star - \mu_T^\theta|_{\text{TV}} \leq \sqrt{2 \cdot \text{KL}(\mu_\star \| \mu_T^\theta)}$ useless here.

Reverse Pinsker inequality lifts TV \rightarrow KL:

$$\text{KL}(\mu_\star \| \mu_T^\theta) \leq \frac{1}{2} \cdot \frac{\log \left\| \frac{d\mu_\star}{d\mu_T^\theta} \right\|_\infty}{1 - \left\| \frac{d\mu_\star}{d\mu_T^\theta} \right\|_\infty^{-1}} \cdot |\mu_\star - \mu_T^\theta|_{\text{TV}} = \underbrace{\frac{1}{2} \cdot \frac{\log \left\| \frac{\rho_\star}{\rho_T^\theta} \right\|_\infty}{1 - \left\| \frac{\rho_\star}{\rho_T^\theta} \right\|_\infty^{-1}}}_{\text{Finite thanks to (ii)}} \cdot \|\rho_\star - \rho_T^\theta\|_{L^1}$$

Tool I: Contraction in L^1

- (**Lemma**). Let $\varepsilon > 0$ and consider any perturbation

$$\rho_B(\cdot) + \varepsilon(\cdot) \quad \text{with} \quad \|\varepsilon(\cdot)\|_{L^1(\mathbb{R}^d)} < \varepsilon.$$

Then, for any $v \in L^\infty(0, T; \text{Lip}(\mathbb{R}^d; \mathbb{R}^d))$,

$$\partial_t \rho_t^\varepsilon(x) + \text{div}(\rho_t^\varepsilon(x)v(t, x)) = 0, \quad \rho_0^\varepsilon(\cdot) = \rho_B(\cdot) + \varepsilon(\cdot)$$

satisfies

$$\|\rho_t^0 - \rho_t^\varepsilon\|_{L^1(\mathbb{R}^d)} < \varepsilon \quad \text{for all } t > 0.$$

Step 1: Control in L^1 (Basic transformations)

$$\partial_t \rho + \partial_x (\mathbf{w} \cdot \sigma(x - b) \rho) = 0, \quad d\rho_B = \eta \mathbb{1}_{(x_0, x_0 + 1/\eta)}$$

- $b = x_0$: $\rho_t^\theta = \eta e^{-wt} \mathbb{1}_{(x_0, b + (x_0 + \frac{1}{\eta} - b)e^{wt})}$ dilation/compression.
- $b < x_0$: $\rho_t^\theta = \eta e^{-wt} \mathbb{1}_{(b + (x_0 - b)e^{wt}, b + (x_0 + \frac{1}{\eta} - b)e^{wt})}$ translation + dilation.
- $x_0 < b < x_0 + \frac{1}{\eta}$: $\rho_t^\theta = \eta \mathbb{1}_{(x_0, b)} + \eta e^{-wt} \mathbb{1}_{(b, b + (x_0 + \frac{1}{\eta} - b)e^{wt})}$ split

STEP 1: Translation and dilation

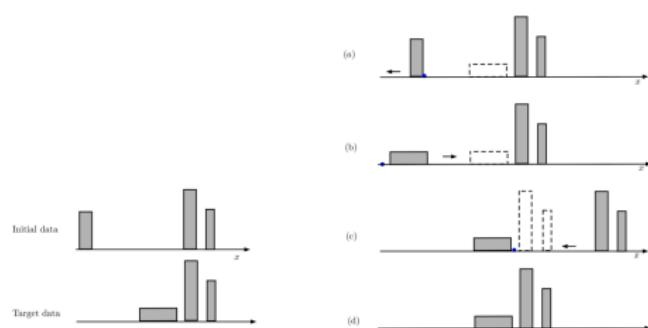


STEP 2: Compression/dilation



(a) Translation as composition of two steps.

(b) Localized translation.



Step 1: Control in L^1 ($d = 1$)

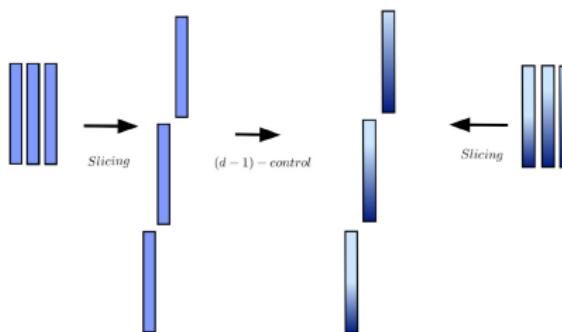
Mechanism: Approximation by K rectangles + inductive exact control

- ① Set $\rho_B = \mathbf{1}_{[x_0-1, x_0]}$ with $x_0 < \inf \text{supp } \rho_\star$.
- ② For $\varepsilon > 0$ choose $h = h(\varepsilon)$ and ε -Riemann approximations ρ_\star^h and ρ_B^h in L^1 consisting of $K(\varepsilon)$ disjoint blocks. Each ordered pair must contain the same mass.
- ③ Inductively (i) compress/dilate and (ii) translate locally the rectangles, then $\rho_B^h \rightarrow \rho_\star^h$ exactly. By L^1 -contraction, it follows that $\|\rho_\star - \rho_{T/2}^\theta\|_{L^1} < \varepsilon$.

Step 1: Control in L^1 ($d > 1$)

- ➊ **Discretization.** Let $[-R_\varepsilon, R_\varepsilon]^d$ contain $1 - \varepsilon$ mass of ρ_B and ρ_* . Grid it with step h , remove tiny strips, and take ε -Riemann approximations for ρ_B^h, ρ_*^h in L^1 .
- ➋ **Simplified goal.** Transform any Riemann approximation on a grid into a constant function on the same grid.

Step 1: Control in L^1 ($d > 1$)



Iterative process to reduce control to $d = 1$.

- ➊ **Discretization.** Let $[-R_\varepsilon, R_\varepsilon]^d$ contain $1 - \varepsilon$ mass of ρ_B and ρ_* . Grid it with step h , remove tiny strips, and take ε -Riemann approximations for ρ_B, ρ_*^h in L^1 .
- ➋ **Simplified goal.** Transform any Riemann approximation on a grid into a constant density with the same mass and defined on the same grid.
- ➌ **Dimension reduction.** Translate “mass columns” until they are ordered in coordinate d . Project onto the first $d - 1$ coords.
- ➍ **Iterate to $d = 1$ and apply Step 1.**

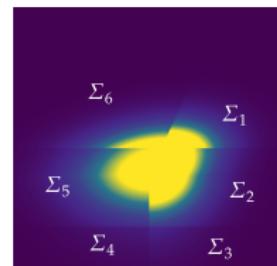
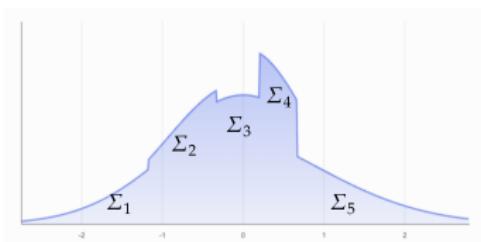
Step 2: Tail control

- Recall: Under a constant control (w, a, b) if $\langle a, x \rangle + b > 0$ then

$$x \mapsto \rho_t^\theta(x) = e^{-t\langle w, a \rangle} \rho_B \left(e^{-t w a^\top} x + \text{shift} \right).$$

- \implies Step I transforms the initial Gaussian tail into a **piecewise Gaussian tail** with covariance matrices depending on each $(\langle a, x \rangle + b)_+$:

$$\Sigma \longmapsto \Sigma_i(t) := \left(e^{-t w a^\top} \right) \Sigma \left(e^{-t w a^\top} \right)^\top \quad \text{or} \quad \Sigma \longmapsto \Sigma_i(t) := \Sigma.$$



Evolution of Gaussian after piecewise constant controls.

Tool II: Variance/tail domination

- (**Lemma**). Suppose

$$\langle (\Sigma_2 - \Sigma_1) e_k, e_k \rangle > 0, \quad \text{for some } k \in [d].$$

Then, for any $M > 0$, there exists $M(k) > 0$ such that

$$\rho_1(x) < \rho_2(x) \quad \text{in } \left\{ x \in \mathbb{R}^d : |x_k| > M(k) \text{ and } |x_j| < M \text{ for } j \neq k \right\}.$$

- (**Lemma**). Suppose $\Sigma_1 \prec \Sigma_2$, then there exists some $\bar{M} > 0$ such that

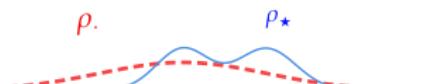
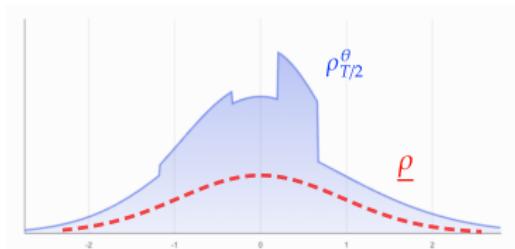
$$\rho_1(x) < \rho_2(x) \quad \text{for all } \|x\| > \bar{M}.$$

Step 2: Tail control

- ➊ New base $\underline{\rho} = \alpha \exp(-\|x\|^2/2\underline{\sigma}^2)$ with $\underline{\sigma}^2 I_d \prec \Sigma_i \forall i$, and $\alpha > 0$ so that $\underline{\rho} < \rho_{T/2}^\theta$.
- ➋ New target $\rho_\bullet := (2\pi\sigma_\bullet^2)^{-d/2} e^{-\|x\|^2/(2\sigma_\bullet^2)}$ (recall $\rho_\star(x) \leq \rho_\bullet(x)$ for $\|x\| \geq M$).
- ➌ New goal:

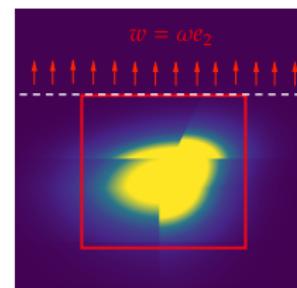
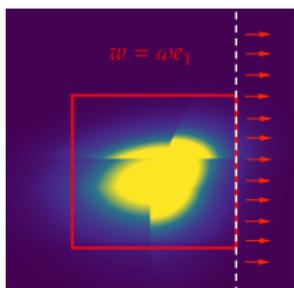
Obtain $\begin{cases} \langle (\Sigma_i(T) - \sigma_\bullet^2 I_d) e_k, e_k \rangle > 0 & \text{on the } k\text{-th face of } [-R_\varepsilon, R_\varepsilon]^d \\ \sigma_\bullet^2 I_d \prec \Sigma_i(T) & \text{on each corner.} \end{cases}$

$\left(\text{This induces full tail control since } \rho_0 < \tilde{\rho}_0 \implies \rho_t < \tilde{\rho}_t \text{ for all } t \in [0, T] \right).$



Step 2: Tail control

- ➊ Choose (a, b) so that $H := \{a^\top x + b = 0\} \equiv e_k$ -face of $[-R_\varepsilon, R_\varepsilon]^d$ and $[-R_\varepsilon, R_\varepsilon]^d \subset H^-$. Then the density will remain fixed inside $[-R_\varepsilon, R_\varepsilon]^d$.
- ➋ Take $w = \omega e_k$ so in H^+ . For $\omega > \frac{2d}{T} \log \left(\frac{\sigma_\bullet^2}{\sigma^2} \right)$, the goal is achieved at face k .
- ➌ Repeat for each of the $2d$ faces of the cube. Note that in the corner regions d controls intervene so we obtain $\Sigma_* \prec \Sigma_i(T)$.



Step 3: Combine both ingredients

- (Step I) We achieve

$$\|\rho_{T/2} - \rho_\star\|_{L^1} < \varepsilon_0.$$

- (Step II) We achieve:

- 1 $\rho_T^\theta \geq \rho_\star$ on $\mathbb{R}^d \setminus [-S_\varepsilon, S_\varepsilon]^d$ for some $S_\varepsilon > R_\varepsilon$.
- 2 $\rho_T = \rho_{T/2}$ on $[-R_\varepsilon, R_\varepsilon]^d$ containing $1 - \varepsilon$ mass of ρ_\star and ρ_T^θ .

- Positivity of ρ_B is preserved to ρ_T^θ by the continuity equation.

- Then $\sup_{\mathbb{R}^d} \frac{\rho_\star}{\rho_T^\theta} < +\infty$, and moreover

$$\left\| \rho_T^\theta - \rho_\star \right\|_{L^1} \leq \int_{[-S_\varepsilon, S_\varepsilon]^d} |\rho_{T/2} - \rho_\star| + \int_{\mathbb{R}^d \setminus [-S_\varepsilon, S_\varepsilon]^d} \rho_T^\theta + \int_{\mathbb{R}^d \setminus [-S_\varepsilon, S_\varepsilon]^d} \rho_\star < 3\varepsilon_0.$$

By Reverse Pinsker Inequality, choosing ε_0 small enough we conclude

$$\text{KL}\left(\mu_\star \parallel \mu_T^\theta\right) < \varepsilon. \quad \square$$

Reverse KL (variational inference)

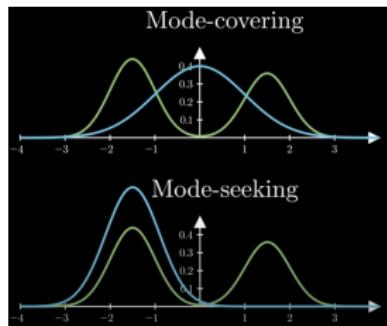
Theorem 2

Assume $\rho_\star > 0$, $\rho_\star \log \rho_\star \in L^1$ and, for some $M, \sigma_\bullet > 0$,

$$\rho_\star(x) \geq \rho_\bullet(x) = \frac{1}{(2\pi\sigma_\bullet^2)^{d/2}} e^{-\|x\|^2/(2\sigma_\bullet^2)} \quad (\|x\| \geq M).$$

Then $\forall T, \varepsilon > 0$ there is a piecewise-constant θ with finitely many switches s.t.

$$\text{KL}(\mu_T^\theta \| \mu_\star) \leq \varepsilon.$$



Beyond Gaussians

- Write

$$\rho_B(x) = e^{-U_B(x)}, \quad \rho_*(x) = e^{-U_*(x)}.$$

- Assumptions:

- 1 Linear growth:

$$U_*(x) \gtrsim \|x\| \quad \text{for } \|x\| \gg 1.$$

- 2 Convex envelope condition:

Given $A \in M_{d \times d}(\mathbb{R})$ with $\text{spec}(A) \subset (0, +\infty)$, there exists $\lambda_A > 0$ s.t.

$$\frac{\text{conv } U_*(\lambda_A x)}{U_B(Ax)} \rightarrow +\infty.$$

- **Idea.** The convex envelope condition guarantees that, after affine transport, the tails of $\mu_T^\theta > \rho_*$ eventually dominate some λ_A rescale of ρ_* 's tails.
- Then $\forall T, \varepsilon > 0$ there is piecewise-constant θ with finitely many switches s.t.

$$\text{KL}\left(\mu_T^\theta \parallel \mu_*\right) \leq \varepsilon.$$

Beyond KL

Csiszár f -divergences. For a convex $f : (0, \infty) \rightarrow \mathbb{R}$ with $f(1) = 0$,

$$D_f(\mu_1, \mu_2) := \int f\left(\frac{d\mu_1}{d\mu_2}\right) d\mu_2.$$

- Examples.
 - $f(t) = t \ln t \Rightarrow$ KL divergence.
 - $f(t) = |t - 1| \Rightarrow$ Total Variation.
 - $f(t) = (\sqrt{t} - 1)^2 \Rightarrow$ Squared Hellinger distance.
 - $f(t) = (t - 1)^2 \Rightarrow \chi^2$ -divergence.
- Our main result extends as:

$$\sup_{t>0} \frac{f(t)}{t \ln t} < \infty \quad \implies \quad D_f(\mu_\star \| \mu_T^\theta) < \varepsilon.$$

Beyond KL

Rényi divergences. For $\lambda > 0$, $\lambda \neq 1$ ($\lambda \rightarrow 1$: recovers KL)

$$D_\lambda(\mu_1, \mu_2) := \frac{1}{\lambda-1} \log \int \left(\frac{d\mu_1}{d\mu_2} \right)^\lambda d\mu_2.$$

- Our main result extends as:

$$0 \leq \lambda \leq 1 \implies D_\lambda(\mu_\star \| \mu_T^\theta) < \varepsilon.$$

Beyond Lipschitz

Motivation. Can we match densities with **different tail regimes**?

It is possible by breaking global Lipschitz constraint of the vector field!

Beyond Lipschitz

- Take the locally Lipschitz activation $\sigma(x) = x \log x \mathbf{1}_{\{x > 1\}}$:

$$\partial_t \rho + \partial_x (-\sigma(x)\rho) = 0, \quad \rho_B(x) \propto e^{-|x|^p}.$$

- Explicit solution:

$$\Phi_t^\theta(x) = \begin{cases} x, & x \leq 1, \\ x^{e^{-t}}, & x > 1 \end{cases} \quad \rho_t^\theta(x) = \begin{cases} \rho_B(x), & x \leq 1, \\ e^t \rho_B(x^{e^t}) x^{e^t - 1}, & x > 1. \end{cases}$$

- Result.** For each $q > 0$, there exists² $T = T(p, q) \geq 0$ such that

$$\lim_{x \rightarrow \infty} \frac{\rho_T(x)}{e^{-|x|^q}} = \lim_{x \rightarrow \infty} e^{|x|^q - |x|^{pe^t}} = 0.$$

- Idea.** Non-Lipschitz $x \log x$ accelerates decay in the tails while preserving mass.

Characteristics increase exponentially but **don't blow up!**

²Take any $T > \log(q/p)$ if $q > p$ and $T = 0$ if $q \leq p$.

Thank you for your attention!



The article.

Takeaways

- ➊ One-neuron ReLU perceptrons suffice to train a CNF by minimizing KL (under tail-compatibility).
- ➋ Explicit construction \Rightarrow blueprint for architectures.
- ➌ Extensions to variational inference and other metrics.
- ➍ Locally-Lipschitz activation $\sigma(x) = x \log x 1_{\{x>1\}}$ enables matching between different tail decays.

KL Divergence

- **Interpretation:** Expected excess surprise when using ν as a model for the real distribution μ .
- **(Basic) properties:**
 - **Non-symmetric:** In general, $\text{KL}(\mu||\nu) \neq \text{KL}(\nu||\mu)$. This implies KL is not a distance.
 - **Non-negative:** $\text{KL}(\mu||\nu) \geq 0$ for all μ, ν and
 $\text{KL}(\mu||\nu) = 0 \iff \mu = \nu$.³
 - **Jointly convex:** For any $(\mu_1, \nu_1), (\mu_2, \nu_2) \in \mathcal{P}(\mathbb{R}^d)$, $\lambda \in [0, 1]$
$$\text{KL}(\lambda\mu_1 + (1 - \lambda)\mu_2 || \lambda\nu_1 + (1 - \lambda)\nu_2) \leq \lambda\text{KL}(\mu_1 || \nu_1) + (1 - \lambda)\text{KL}(\mu_2 || \nu_2).$$

³Proof:

$$\text{KL}(\mu||\nu) = \mathbb{E}_\mu \underbrace{\left[-\log \left(\frac{d\nu}{d\mu} \right) \right]}_{\text{convex}} \stackrel{\text{(Jensen ineq.)}}{\geq} -\log \left(\mathbb{E}_\mu \left[\frac{d\nu}{d\mu} \right] \right) = \log 1 = 0$$

Basic properties of divergences

- **More (basic) properties:**

- **Bounds.** For any convex $f : (0, \infty) \rightarrow \mathbb{R}$ with $f(1) = 0$,

$$D_f(\mu_1 \| \mu_2) := \int f\left(\frac{d\mu_1}{d\mu_2}\right) d\mu_2.$$

If there exists $C > 0$ such that

$$f(u) \leq C(u \log u - u + 1) \quad \forall u > 0,$$

then (whenever $\mu_1 \ll \mu_2$)

$$D_f(\mu_1 \| \mu_2) \leq C \text{KL}(\mu_1 \| \mu_2).$$

- **Rényi divergence (order $\lambda \in (0, 1)$).**

$$D_\lambda(\mu_1 \| \mu_2) := \frac{1}{\lambda - 1} \log \int \left(\frac{d\mu_1}{d\mu_2} \right)^\lambda d\mu_2 \leq \text{KL}(\mu_1 \| \mu_2),$$

using the monotonicity of D_α in α and $D_1 = \text{KL}$.

- **Unbalanced mass.** If $\mu(\mathbb{R}^d) = \alpha \neq \beta = \nu(\mathbb{R}^d)$ and $\mu \ll \nu$, write $\bar{\mu} := \mu/\alpha$, $\bar{\nu} := \nu/\beta$ (probability measures). Then

$$\text{KL}(\mu \| \nu) = \alpha \log\left(\frac{\alpha}{\beta}\right) + \alpha \text{KL}(\bar{\mu} \| \bar{\nu}).$$

Proof sketch: 1D simultaneous control

Step A (approximate target): Riemann (or P_0 FEM) approximation

$$\rho_T^h = \sum_{k=1}^K \rho_{T,k} \mathbb{1}_{(x_{k-1}, x_k - \delta)}, \quad \|\rho_T - \rho_T^h\|_{L^1} < \varepsilon.$$

Step B (approximate source): Split ρ_B into K disjoint intervals with the *same* masses as the K target blocks, introducing tiny gaps of total mass $\tau \leq \varepsilon$:

$$\rho_B^h = \sum_{k=1}^{K-1} \mathbb{1}_{(y_{k-1}, y_k - \tau/(K-1))} + \mathbb{1}_{(y_{K-1}, y_K)}, \quad \|\rho_B - \rho_B^h\|_{L^1} \leq \varepsilon.$$

Step C (simultaneous control of blocks):

$$\partial_t \rho^{(k)} + \partial_x (w(t) \sigma(a(t)x + b(t)) \rho^{(k)}) = 0,$$

$$\rho_0^{(k)} = \mathbb{1}_{I_k} \longrightarrow \rho_T^{(k)} = \rho_{T,k} \mathbb{1}_{J_k}, \quad k = 1, \dots, K.$$

Mechanism: Induction over k using (i) local compression/dilation and (ii) local translation , with choices of (a, b) that make the field vanish on protected regions.

Conclusion: $\rho_B^h \rightarrow \rho_T^h$ exactly; by L^1 -contraction $\|\rho_\star - \rho_T\|_{L^1} \leq 2\varepsilon$.

Proof sketch: parallel translation in $d \geq 2$

Half-space transport with ReLU: Choose $a \perp w$ so $\operatorname{div}(w \sigma(\langle a, x \rangle + b)) = 0$ and the field vanishes on one side of the hyperplane $\{\langle a, x \rangle + b = 0\}$.

Two-step composition: Pick b_1, b_2 and apply

$$\begin{cases} (a, w) = ((1, 0, \dots, 0), -(0, 1, 0, \dots, 0)), & b = -b_1 \\ (a, w) = ((1, 0, \dots, 0), +(0, 1, 0, \dots, 0)), & b = -b_2 \end{cases}$$

to obtain a net *parallel translation* of selected mass along e_2 , leaving the complementary half-space invariant.

Use: Enables independent movement of separated parts of the support in \mathbb{R}^d .

Proof sketch: Dimension descent

- (1) **Truncate & grid:** pick cube $\mathcal{R} = [-R, R]^d$, mesh size h , hyperplanes $H_{k\ell} : x^{(k)} = c_{k,\ell}$.

$$\rho_{0,T} \approx \sum_j m_{j,0/T} \mathbb{1}_{\square_{j+h(1,\dots,1)/2, h-\delta}},$$

removing strips Ω_δ around hyperplanes.

- (2) **Align along $x^{(d)}$:** By parallel translations slab-by-slab so projected boxes in $(x^{(1)}, \dots, x^{(d-1)})$ are *disjoint* and independent.

- (3) **Recurse:** the new densities are constant in $x^{(d)}$ on their support \Rightarrow reduce to a $(d-1)$ -D problem. Iterate down to 1-D.

- (4) **Stability:** L^1 -contraction propagates the ε -approximations from truncation and meshing.