

# Generalized Gradient Flows for Stochastic Prediction, Filtering, Learning and Control

Abhishek Halder

Department of Applied Mathematics  
University of California, Santa Cruz  
Santa Cruz, CA 95064

Joint work with S. Haddad, K.F. Caluya (UC Santa Cruz)  
B. Singh (Ford), T.T. Georgiou (UC Irvine), W. Krichene (Google)

Optimal Transport and Mean Field Games Seminar  
University of South Carolina, SC

January 26, 2022

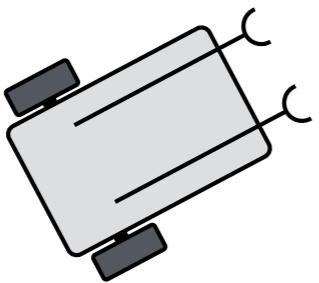


# Overarching Theme

**Systems-control theory and algorithms  
for densities**

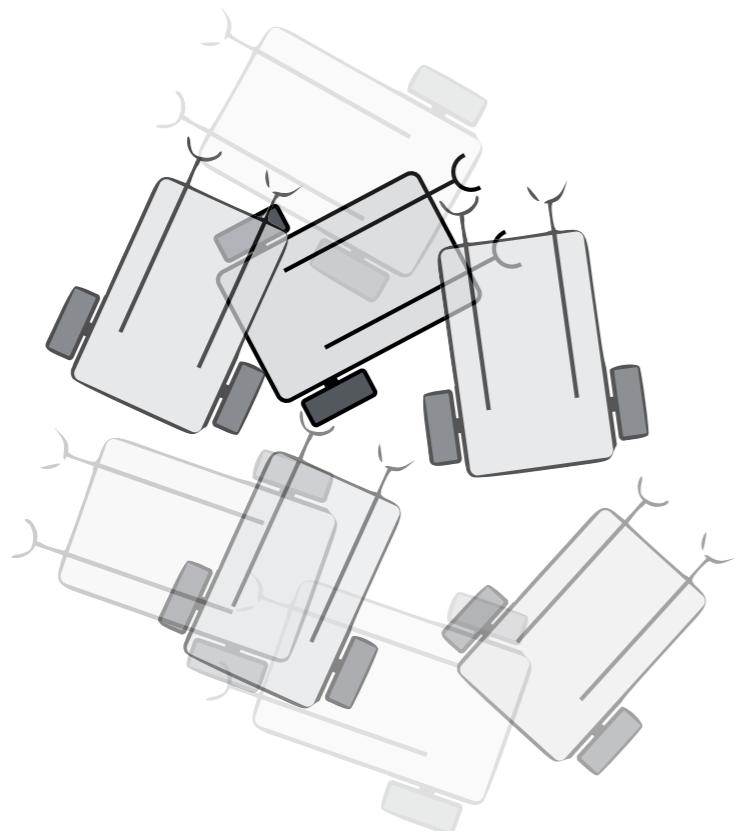
# What is density?

# Probability Density Fn.



$$x(t) \in \begin{pmatrix} x \\ y \\ \theta \end{pmatrix} \in \mathcal{X} \equiv \mathbb{R}^2 \times \mathbb{S}^1$$

# Probability Density Fn.

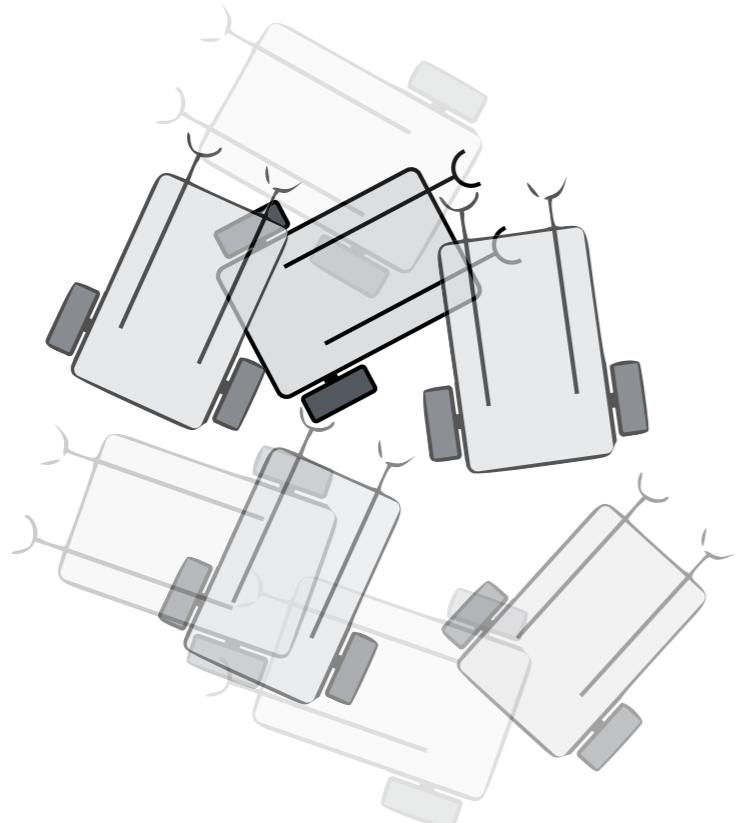


$$x(t) \in \begin{pmatrix} x \\ y \\ \theta \end{pmatrix} \in \mathcal{X} \equiv \mathbb{R}^2 \times \mathbb{S}^1$$

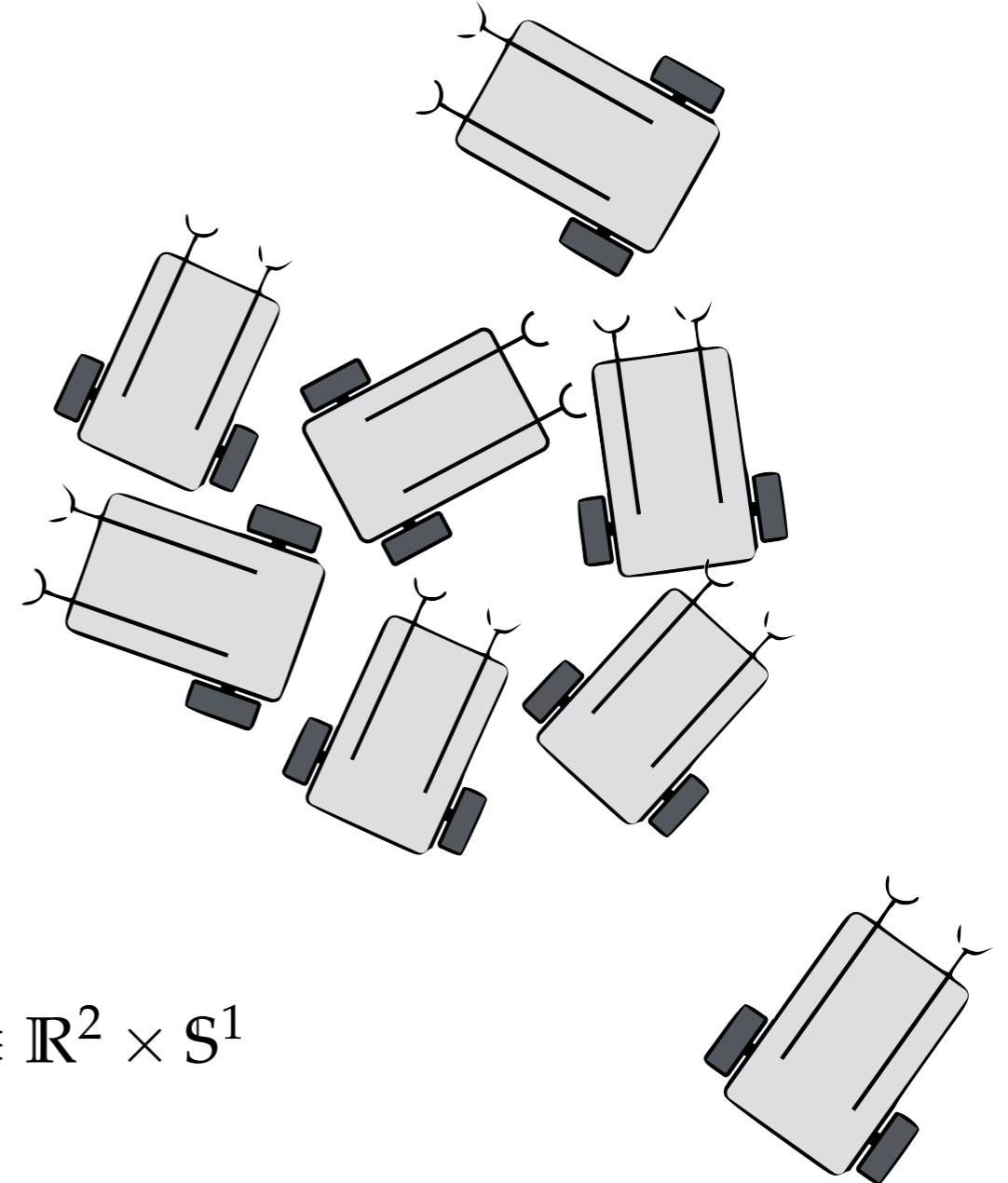
$$\rho(x, t) : \mathcal{X} \times [0, \infty) \mapsto \mathbb{R}_{\geq 0}$$

$$\int_{\mathcal{X}} \rho \, dx = 1 \quad \text{for all } t \in [0, \infty)$$

# Probability Density Fn.



# Population Density Fn.



$$x(t) \in \begin{pmatrix} x \\ y \\ \theta \end{pmatrix} \in \mathcal{X} \equiv \mathbb{R}^2 \times \mathbb{S}^1$$

$$\rho(x, t) : \mathcal{X} \times [0, \infty) \mapsto \mathbb{R}_{\geq 0}$$

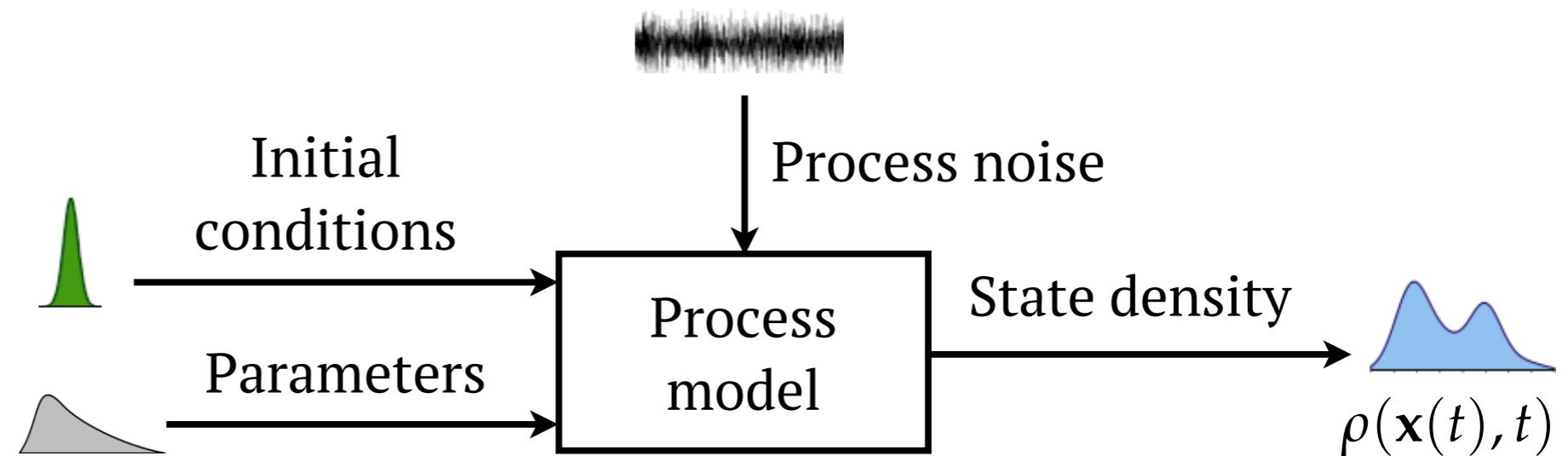
$$\int_{\mathcal{X}} \rho \, dx = 1 \quad \text{for all } t \in [0, \infty)$$

# Why care about densities?

# Prediction Problem

Compute  
joint state PDF

$$\rho(x, t)$$



Trajectory flow:

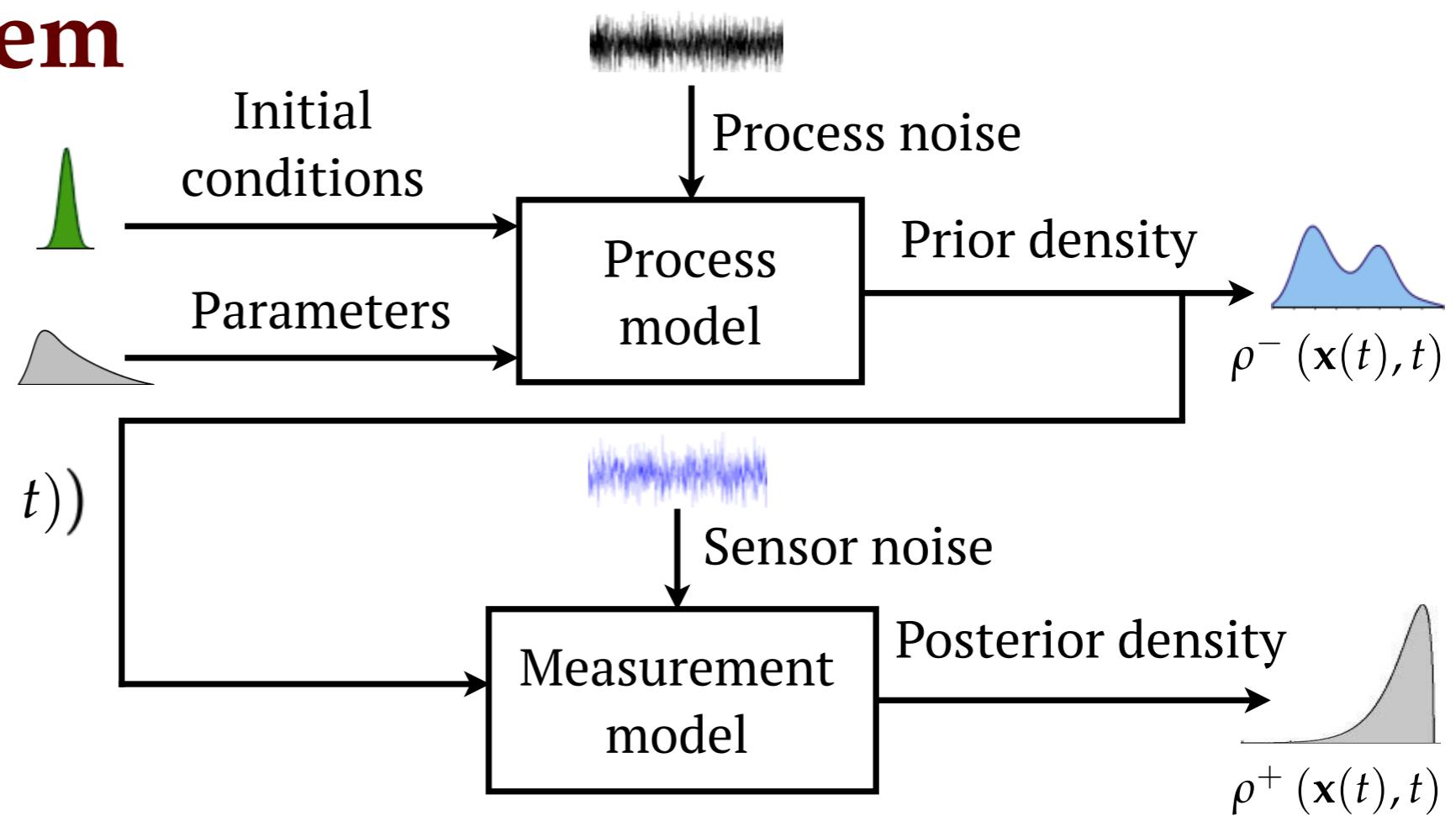
$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{x}, t) dt + \mathbf{g}(\mathbf{x}, t) dw(t), \quad dw(t) \sim \mathcal{N}(0, Qdt)$$

Density flow:

$$\frac{\partial \rho}{\partial t} = \mathcal{L}_{\text{FP}}(\rho) := -\nabla \cdot (\rho \mathbf{f}) + \frac{1}{2} \sum_{i,j=1}^n \frac{\partial^2}{\partial x_i \partial x_j} \left( \left( \mathbf{g} \mathbf{Q} \mathbf{g}^\top \right)_{ij} \rho \right)$$

# Filtering Problem

Compute conditional joint state PDF



$$\rho^+ := \rho (\mathbf{x}, t \mid \mathbf{z}(s), 0 \leq s \leq t))$$

Trajectory flow:

$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{x}, t) dt + \mathbf{g}(\mathbf{x}, t) dw(t), \quad dw(t) \sim \mathcal{N}(0, \mathbf{Q} dt)$$

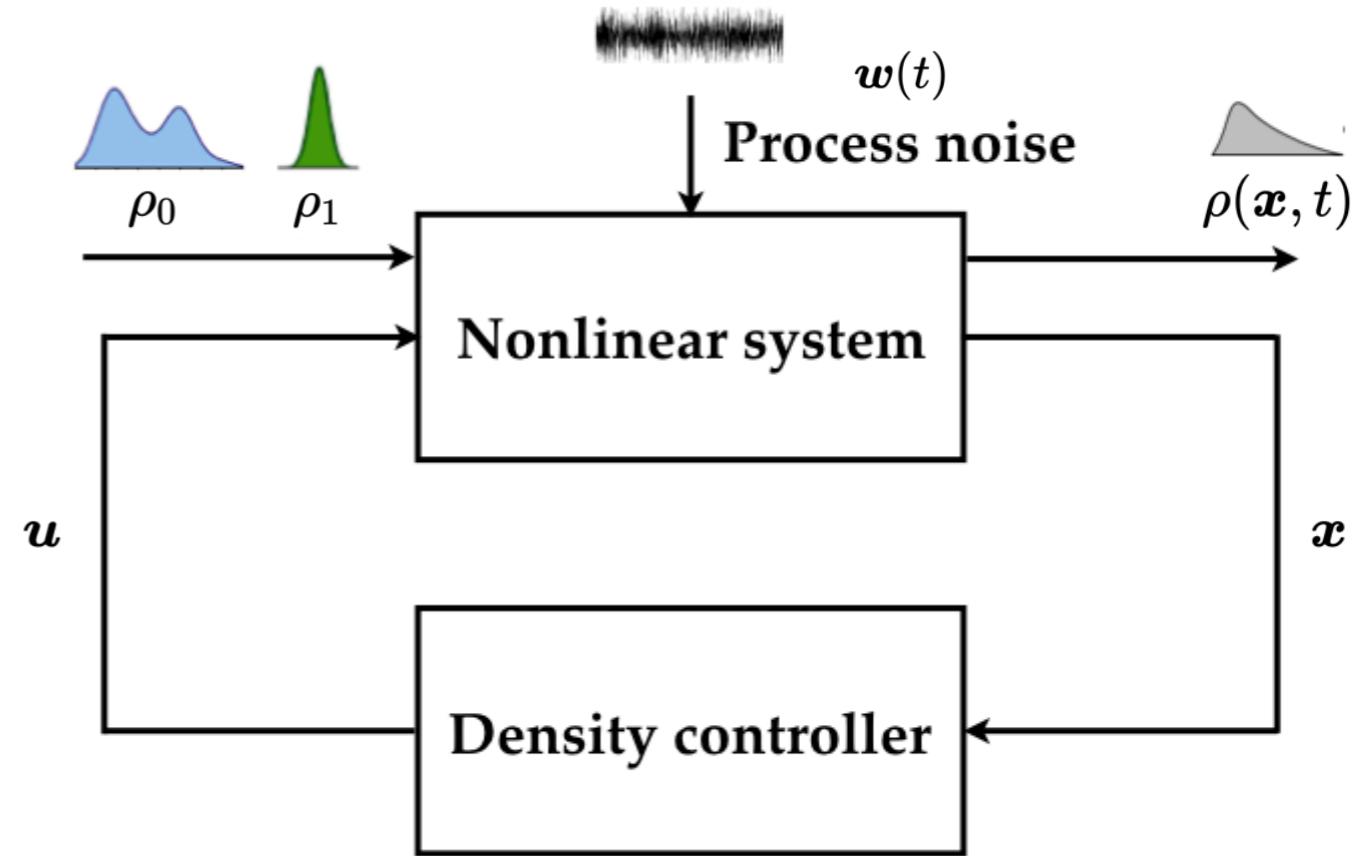
$$d\mathbf{z}(t) = \mathbf{h}(\mathbf{x}, t) dt + dv(t), \quad dv(t) \sim \mathcal{N}(0, \mathbf{R} dt)$$

Density flow:

$$d\rho^+ = \left[ \mathcal{L}_{FP} dt + (\mathbf{h}(\mathbf{x}, t) - \mathbb{E}_{\rho^+}\{\mathbf{h}(\mathbf{x}, t)\})^\top \mathbf{R}^{-1} (d\mathbf{z}(t) - \mathbb{E}_{\rho^+}\{\mathbf{h}(\mathbf{x}, t)\} dt) \right] \rho^+$$

# Control Problem

Steer joint state PDF via feedback control over finite time horizon



$$\underset{u \in \mathcal{U}}{\text{minimize}} \quad \mathbb{E} \left[ \int_0^1 \|u\|_2^2 \, dt \right]$$

subject to

$$dx = f(x, u, t) \, dt + g(x, t) \, dw,$$

$$x(t=0) \sim \rho_0, \quad x(t=1) \sim \rho_1$$

# Neural Network Learning Problem

Consider fully connected NN

Think “layers” as interacting population of neurons

Mean field learning problem:

$$\inf_{\rho \in \mathcal{P}_2(\mathbb{R}^p)} R\left(\int \Phi(x, \theta) \rho(\theta) d\theta\right)$$

PDF dynamics:

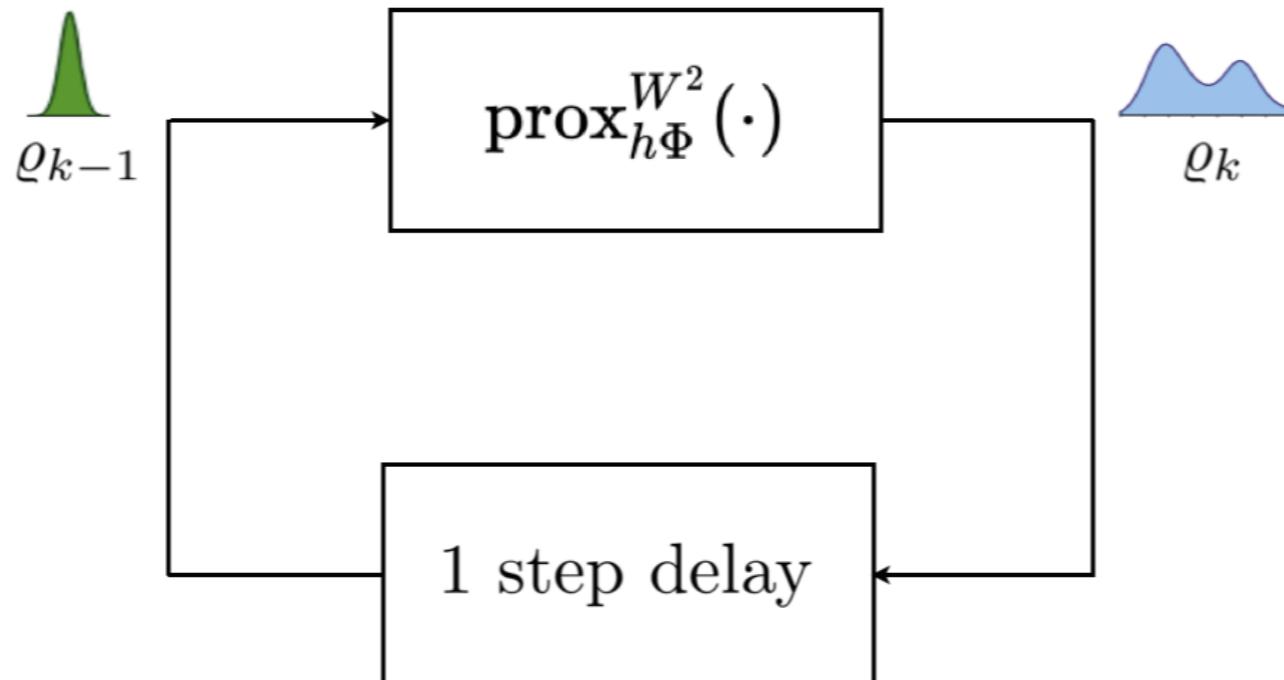
$$\frac{\partial \rho}{\partial t} = -\nabla^W R\left(\int \Phi \rho\right) = \nabla \cdot \left( \rho \nabla \frac{\delta}{\delta \rho} R\left(\int \Phi \rho\right) \right)$$

# Solving prediction problem as generalized gradient flow

# What's New?

Main idea: Solve  $\frac{\partial \rho}{\partial t} = \mathcal{L}_{\text{FP}} \rho$ ,  $\rho(x, t=0) = \rho_0$  as gradient flow in  $\mathcal{P}_2(\mathcal{X})$

Infinite dimensional variational recursion:



Proximal operator:  $\rho_k = \text{prox}_{h\Phi}^{W^2}(\rho_{k-1}) := \arg \inf_{\rho \in \mathcal{P}_2(\mathcal{X})} \left\{ \frac{1}{2} W^2(\rho, \rho_{k-1}) + h\Phi(\rho) \right\}$

Optimal transport cost:  $W^2(\rho, \rho_{k-1}) := \inf_{\pi \in \Pi(\rho, \rho_{k-1})} \int_{\mathcal{X} \times \mathcal{X}} c(x, y) d\pi(x, y)$

Free energy functional:  $\Phi(\rho) := \int_{\mathcal{X}} \psi \rho dx + \beta^{-1} \int_{\mathcal{X}} \rho \log \rho dx$

# Geometric Meaning of Gradient Flow

## Gradient Flow in $\mathcal{X}$

$$\frac{d\mathbf{x}}{dt} = -\nabla \varphi(\mathbf{x}), \quad \mathbf{x}(0) = \mathbf{x}_0$$

## Gradient Flow in $\mathcal{P}_2(\mathcal{X})$

$$\frac{\partial \rho}{\partial t} = -\nabla^W \Phi(\rho), \quad \rho(\mathbf{x}, 0) = \rho_0$$

### Recursion:

$$\begin{aligned}\mathbf{x}_k &= \mathbf{x}_{k-1} - h \nabla \varphi(\mathbf{x}_k) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{x}_{k-1}\|_2^2 + h \varphi(\mathbf{x}) \right\} \\ &=: \text{prox}_{h\varphi}^{\|\cdot\|_2}(\mathbf{x}_{k-1})\end{aligned}$$

### Recursion:

$$\begin{aligned}\rho_k &= \rho(\cdot, t = kh) \\ &= \arg \min_{\rho \in \mathcal{P}_2(\mathcal{X})} \left\{ \frac{1}{2} W^2(\rho, \rho_{k-1}) + h \Phi(\rho) \right\} \\ &=: \text{prox}_{h\Phi}^{W^2}(\rho_{k-1})\end{aligned}$$

### Convergence:

$$\mathbf{x}_k \rightarrow \mathbf{x}(t = kh) \quad \text{as} \quad h \downarrow 0$$

### Convergence:

$$\rho_k \rightarrow \rho(\cdot, t = kh) \quad \text{as} \quad h \downarrow 0$$

### $\varphi$ as Lyapunov function:

$$\frac{d}{dt} \varphi = - \|\nabla \varphi\|_2^2 \leq 0$$

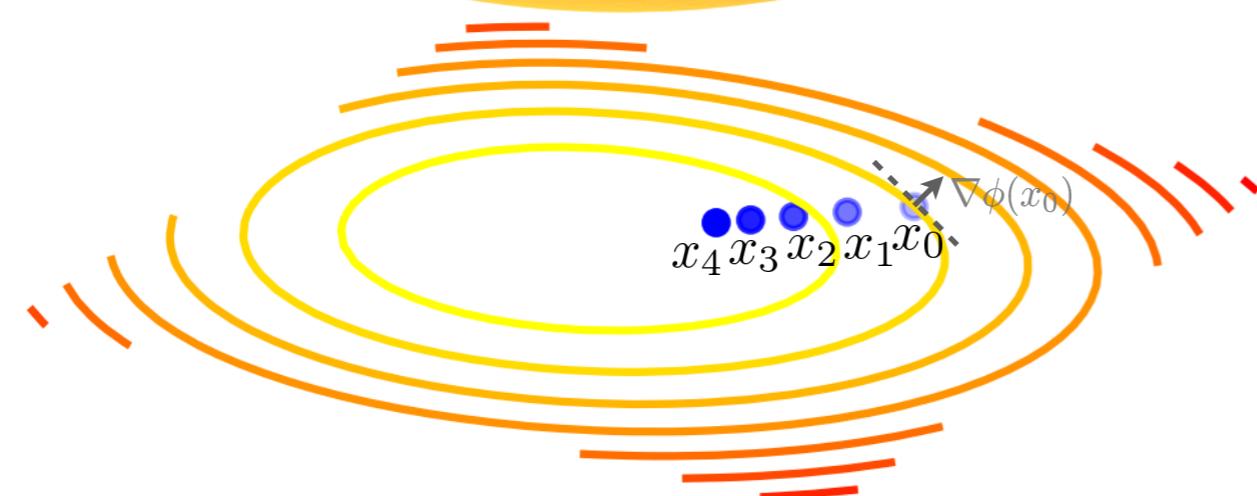
### $\Phi$ as Lyapunov functional:

$$\frac{d}{dt} \Phi = -\mathbb{E}_\rho \left[ \left\| \nabla \frac{\delta \Phi}{\delta \rho} \right\|_2^2 \right] \leq 0$$

# Geometric Meaning of Gradient Flow

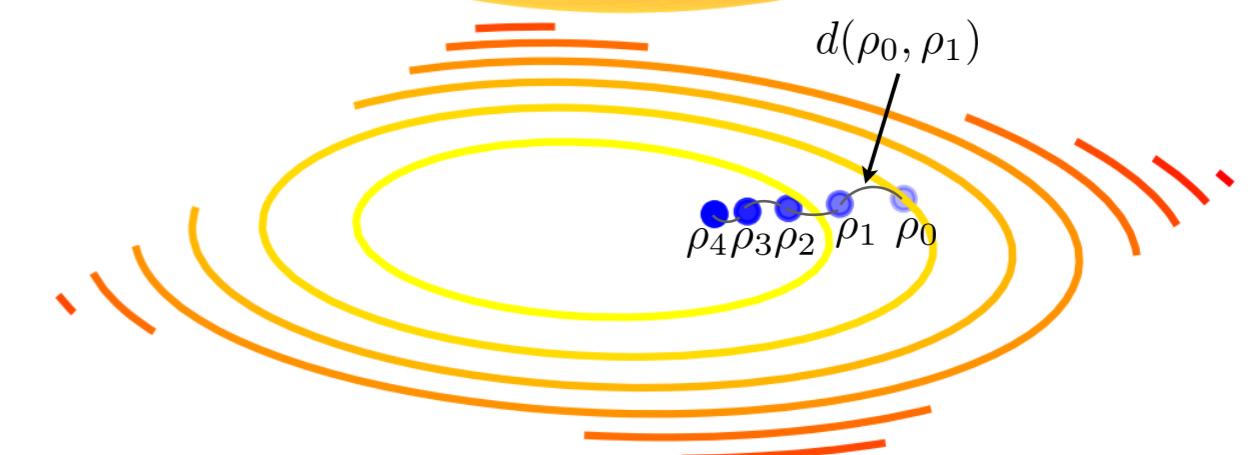
## Gradient Flow in $\mathcal{X}$

$$z = \phi(x), \quad x \in \mathbb{R}^2$$



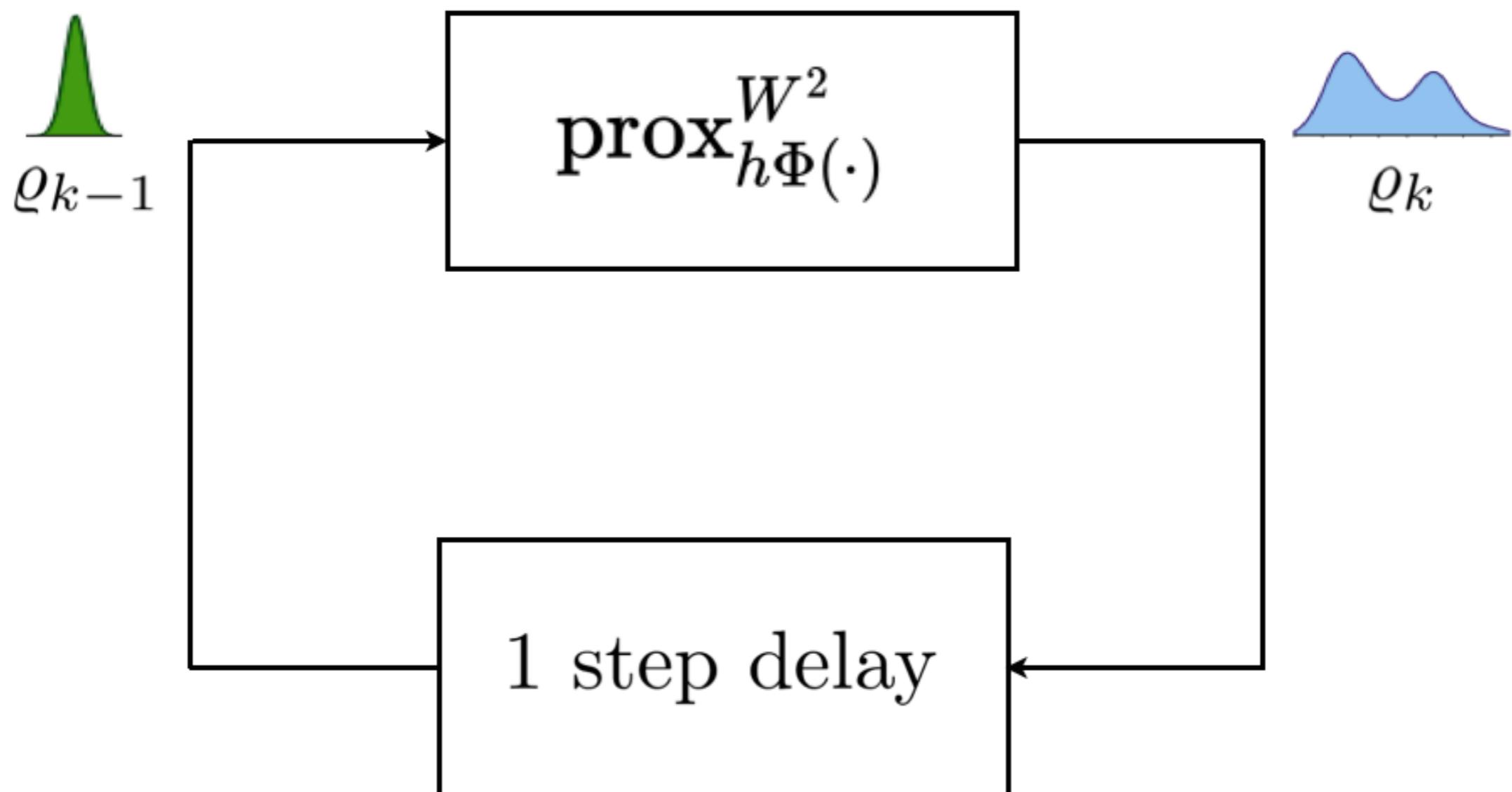
## Gradient Flow in $\mathcal{P}_2(\mathcal{X})$

$$z = \Phi(\rho), \quad \rho \in \mathcal{P}_2(\mathcal{X})$$



# Algorithm: Gradient Ascent on the Dual Space

Uncertainty propagation via point clouds



No spatial discretization or function approximation

# Algorithm: Gradient Ascent on the Dual Space

$$\frac{\partial \rho}{\partial t} = \nabla \cdot (\nabla \psi \rho) + \beta^{-1} \Delta \rho$$

⇓

**Proximal Recursion**

$$\rho_k = \rho(\mathbf{x}, t = kh) = \arg \inf_{\rho \in \mathcal{P}_2(\mathbb{R}^n)} \left\{ \frac{1}{2} W^2(\rho, \rho_{k-1}) + h \Phi(\rho) \right\}$$

⇓

**Discrete Primal Formulation**

$$\varrho_k = \arg \min_{\varrho} \left\{ \min_{\mathbf{M} \in \Pi(\varrho_{k-1}, \varrho)} \frac{1}{2} \langle \mathbf{C}_k, \mathbf{M} \rangle + h \langle \psi_{k-1} + \beta^{-1} \log \varrho, \varrho \rangle \right\}$$

⇓

**Entropic Regularization**

$$\varrho_k = \arg \min_{\varrho} \left\{ \min_{\mathbf{M} \in \Pi(\varrho_{k-1}, \varrho)} \frac{1}{2} \langle \mathbf{C}_k, \mathbf{M} \rangle + \epsilon H(\mathbf{M}) + h \langle \psi_{k-1} + \beta^{-1} \log \varrho, \varrho \rangle \right\}$$

⇓

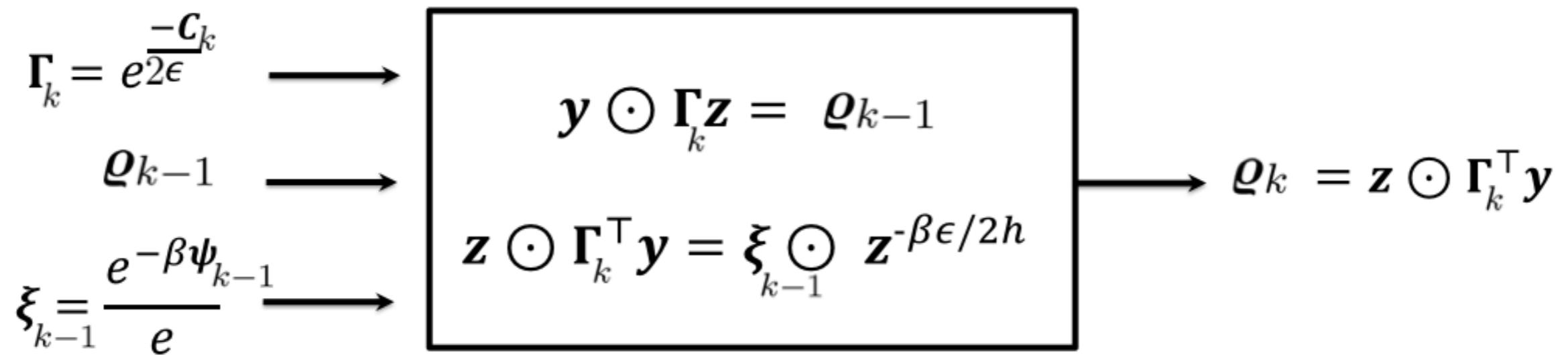
**Dualization**

$$\begin{aligned} \lambda_0^{\text{opt}}, \lambda_1^{\text{opt}} &= \arg \max_{\lambda_0, \lambda_1 \geq 0} \left\{ \langle \lambda_0, \varrho_{k-1} \rangle - F^*(-\lambda_1) \right. \\ &\quad \left. - \frac{\epsilon}{h} \left( \exp(\lambda_0^\top h/\epsilon) \exp(-\mathbf{C}_k/2\epsilon) \exp(\lambda_1 h/\epsilon) \right) \right\} \end{aligned}$$

# Recursion on the Cone

$$y = e^{\frac{\lambda_0^*}{\epsilon} h} \quad z = e^{\frac{\lambda_1^*}{\epsilon} h}$$

Coupled Transcendental Equations in  $y$  and  $z$

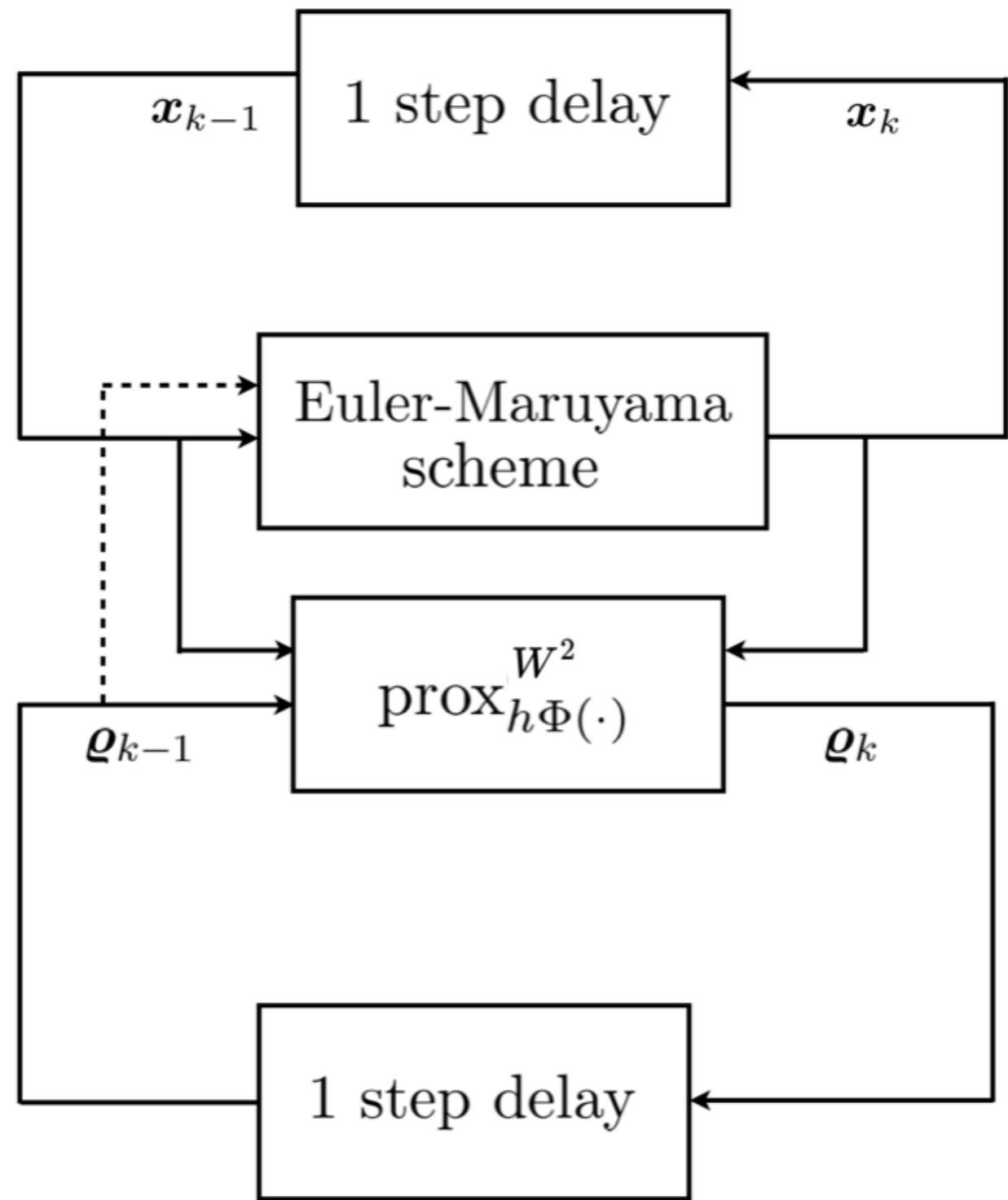


**Theorem:** Consider the recursion on the cone  $\mathbb{R}_{\geq 0}^n \times \mathbb{R}_{\geq 0}^n$

$$y \odot (\Gamma_k z) = Q_{k-1}, \quad z \odot (\Gamma_k^T y) = \xi_{k-1} \odot z^{-\frac{\beta\epsilon}{h}},$$

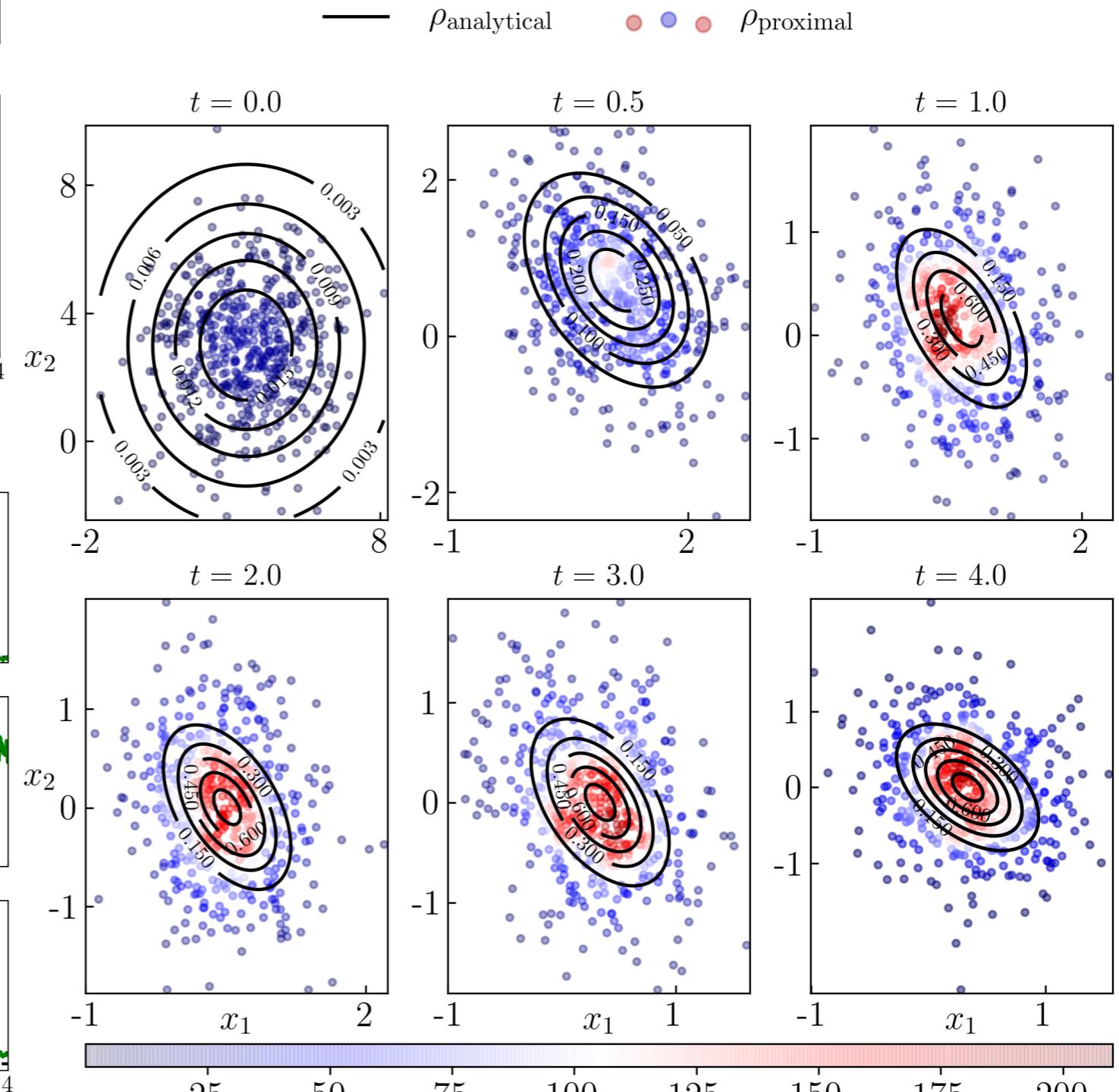
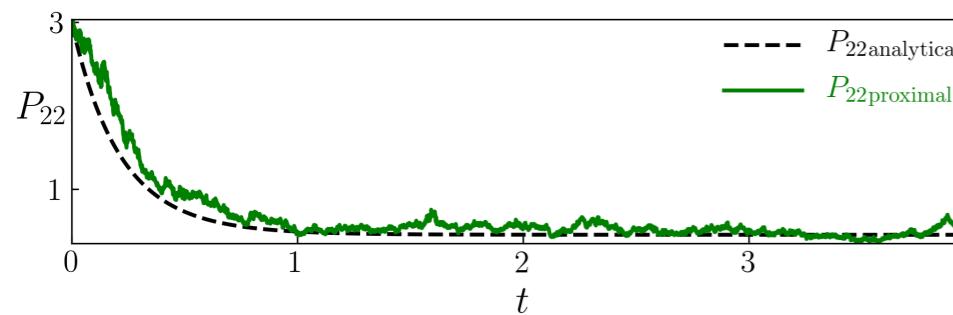
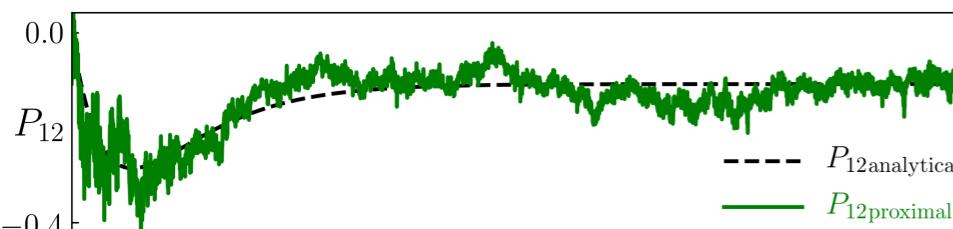
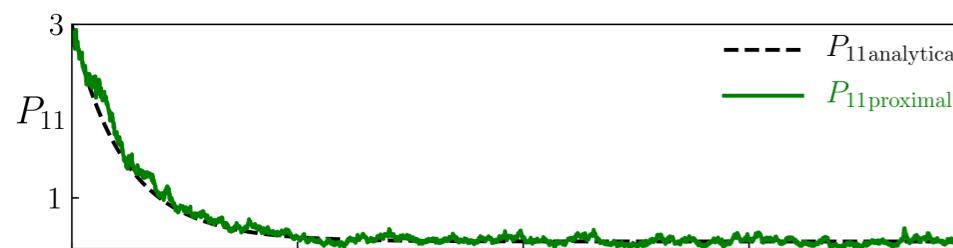
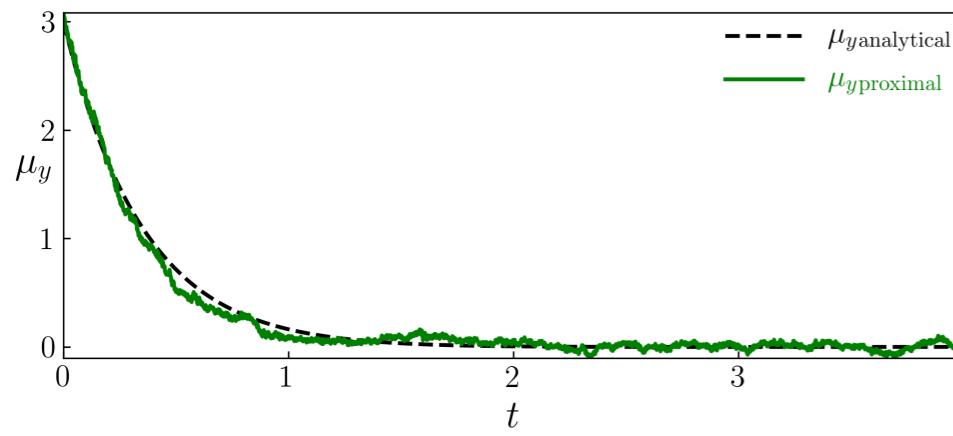
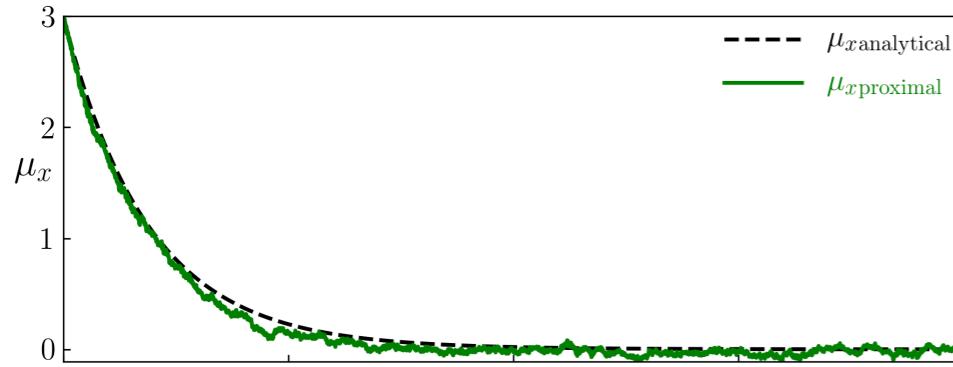
Then the solution  $(y^*, z^*)$  gives the proximal update  $Q_k = z^* \odot (\Gamma_k^T y^*)$

# Algorithmic Setup

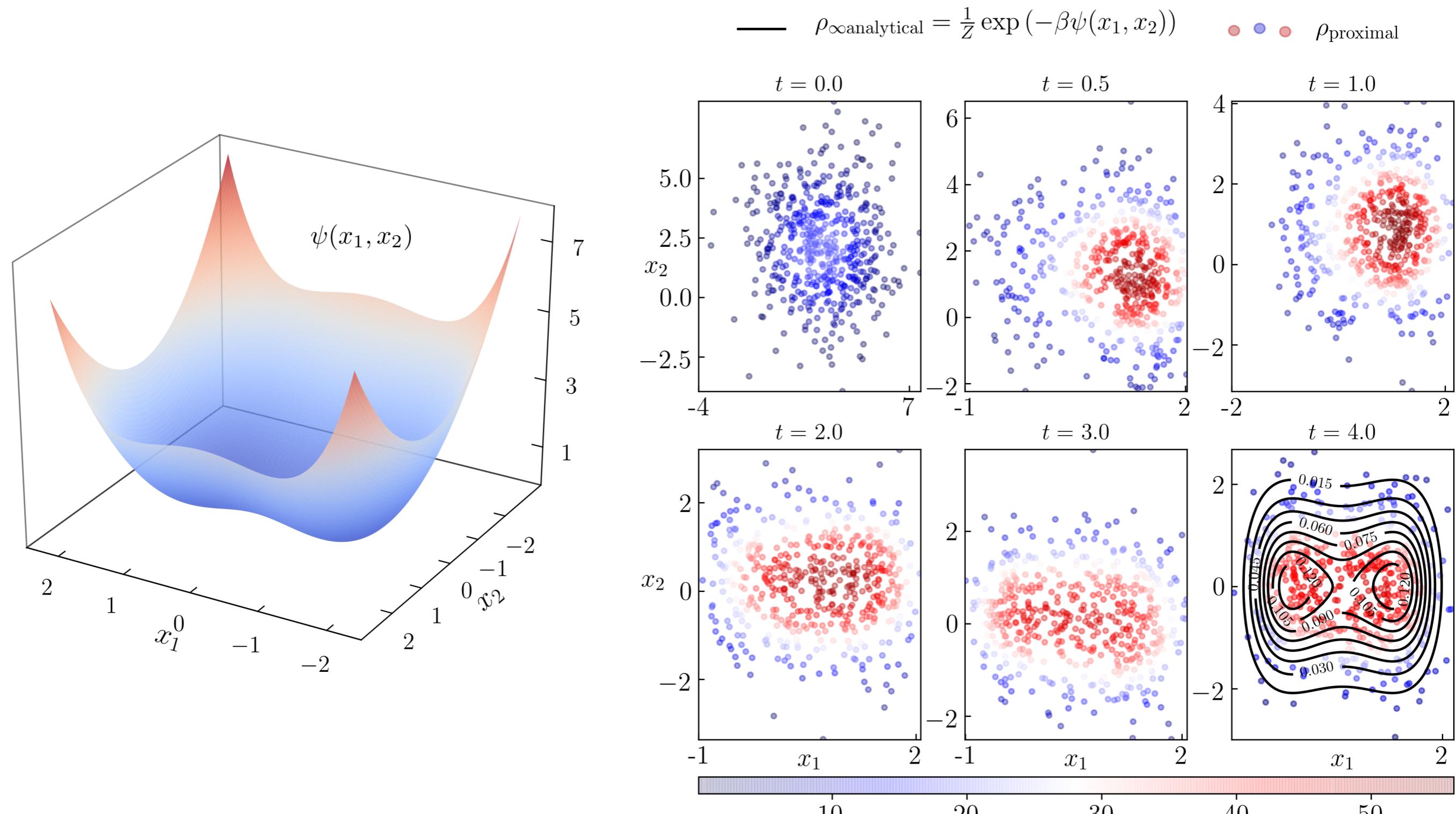


**Theorem:** Block co-ordinate iteration of  $(y, z)$  recursion is contractive on  $\mathbb{R}_{>0}^n \times \mathbb{R}_{>0}^n$ .

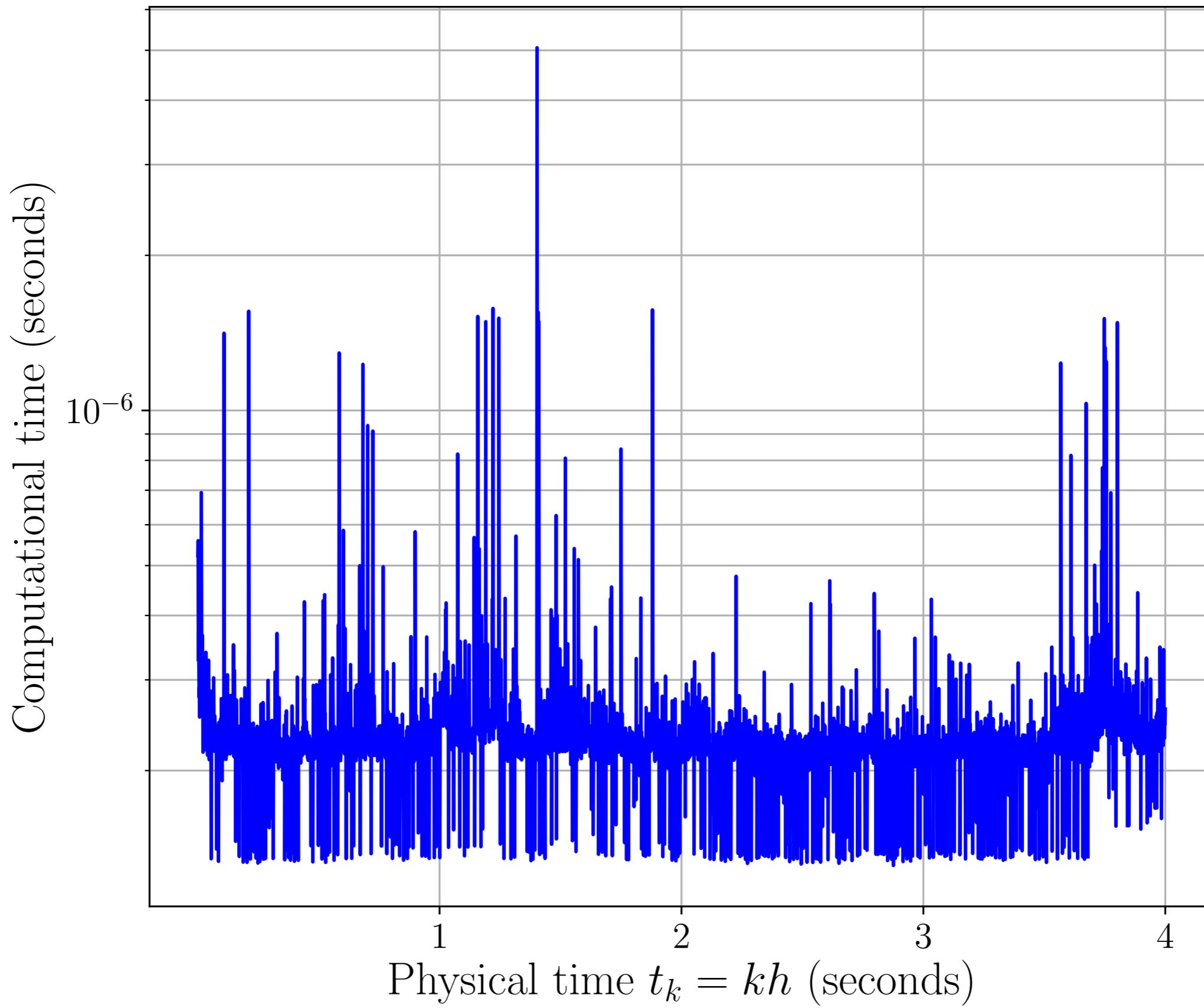
# Proximal Prediction: 2D Linear Gaussian



# Proximal Prediction: Nonlinear Non-Gaussian



# Computational Time: Nonlinear Non-Gaussian



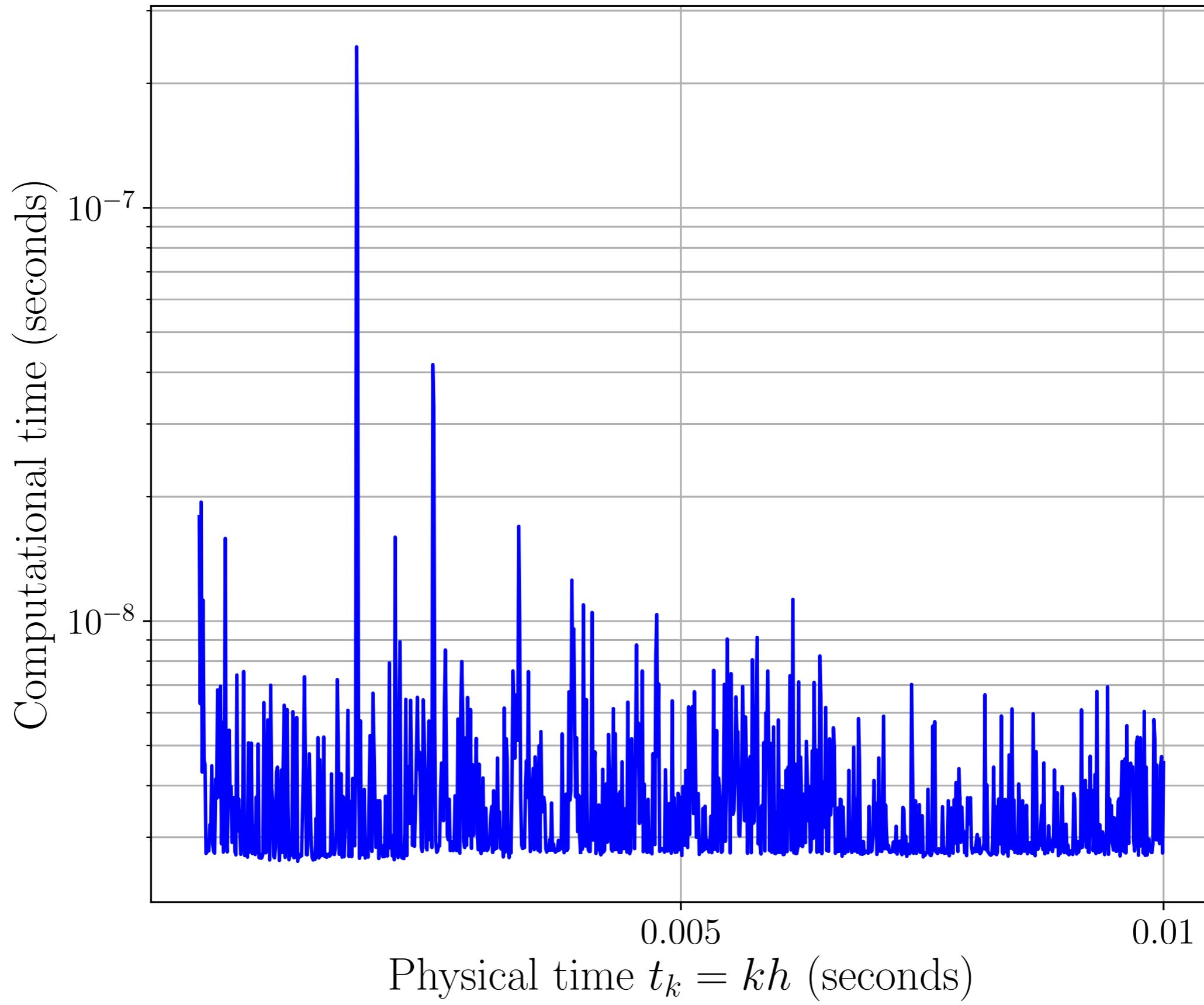
# Proximal Prediction: Satellite in Geocentric Orbit

Here,  $\mathcal{X} \equiv \mathbb{R}^6$

$$\begin{pmatrix} dx \\ dy \\ dz \\ dv_x \\ dv_y \\ dv_z \end{pmatrix} = \begin{pmatrix} v_x \\ v_y \\ v_z \\ -\frac{\mu x}{r^3} + (f_x)_{\text{pert}} - \gamma v_x \\ -\frac{\mu y}{r^3} + (f_y)_{\text{pert}} - \gamma v_y \\ -\frac{\mu z}{r^3} + (f_z)_{\text{pert}} - \gamma v_z \end{pmatrix} dt + \sqrt{2\beta^{-1}\gamma} \begin{pmatrix} 0 \\ 0 \\ 0 \\ dw_1 \\ dw_2 \\ dw_3 \end{pmatrix},$$

$$\begin{pmatrix} f_x \\ f_y \\ f_z \end{pmatrix}_{\text{pert}} = \begin{pmatrix} s\theta \ c\phi & c\theta \ c\phi & -s\phi \\ s\theta \ s\phi & c\theta \ s\phi & c\phi \\ c\theta & -s\theta & 0 \end{pmatrix} \begin{pmatrix} \frac{k}{2r^4} (3(s\theta)^2 - 1) \\ -\frac{k}{r^5} s\theta \ c\theta \\ 0 \end{pmatrix}, \quad k := 3J_2 R_E^2, \mu = \text{constant}$$

# Computational Time: Satellite in Geocentric Orbit



# Extensions: Nonlocal Interactions

PDF dependent sample path dynamics:

$$dx = -(\nabla U(x) + \nabla \rho * V) dt + \sqrt{2\beta^{-1}} dw$$

McKean-Vlasov-Fokker-Planck-Kolmogorov integro PDE:

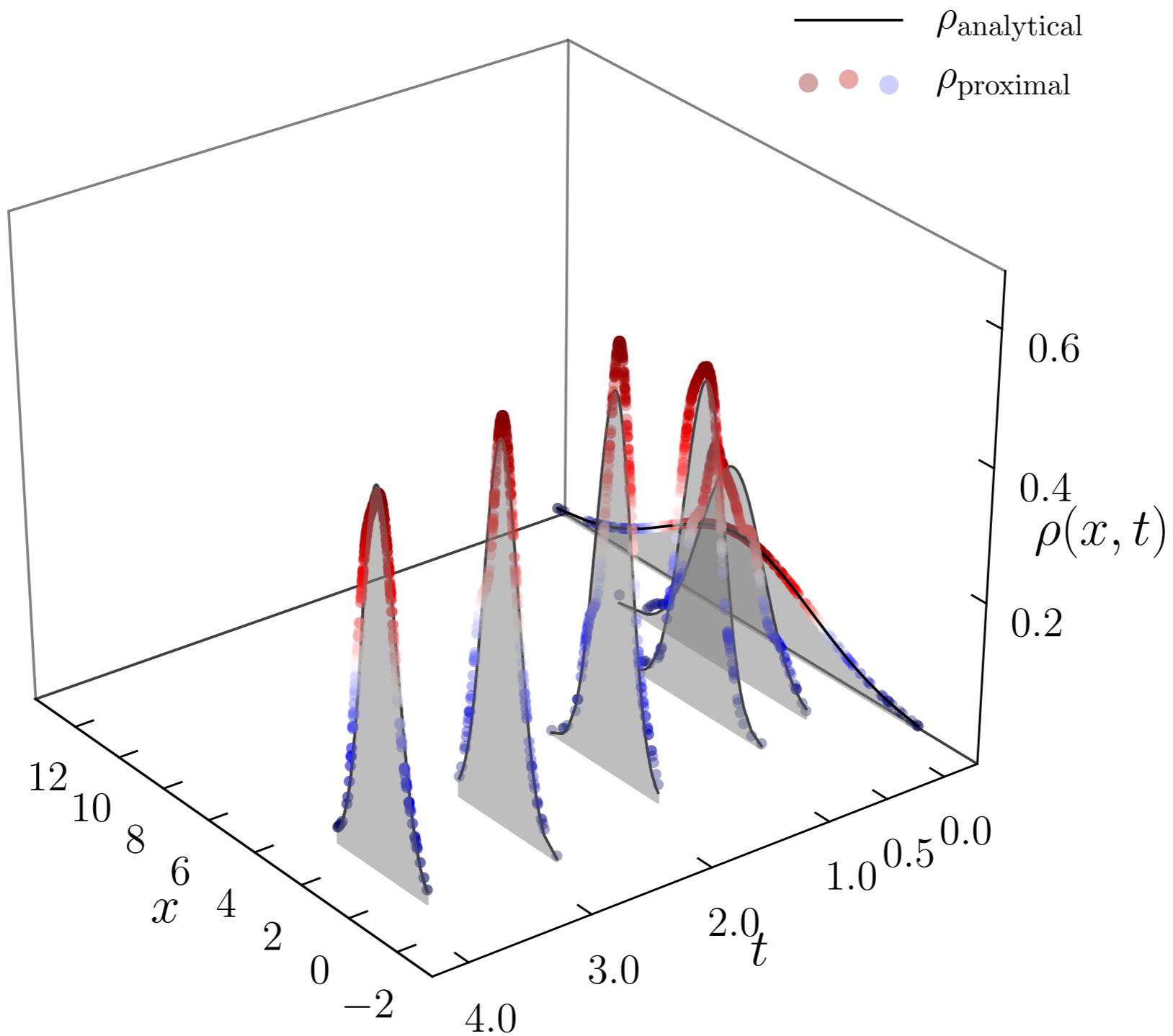
$$\frac{\partial \rho}{\partial t} = \nabla \cdot (\rho \nabla (U + \rho * V)) + \beta^{-1} \Delta \rho$$

Free energy:

$$F(\rho) := \mathbb{E}_\rho [U + \beta^{-1} \rho \log \rho + \rho * V]$$

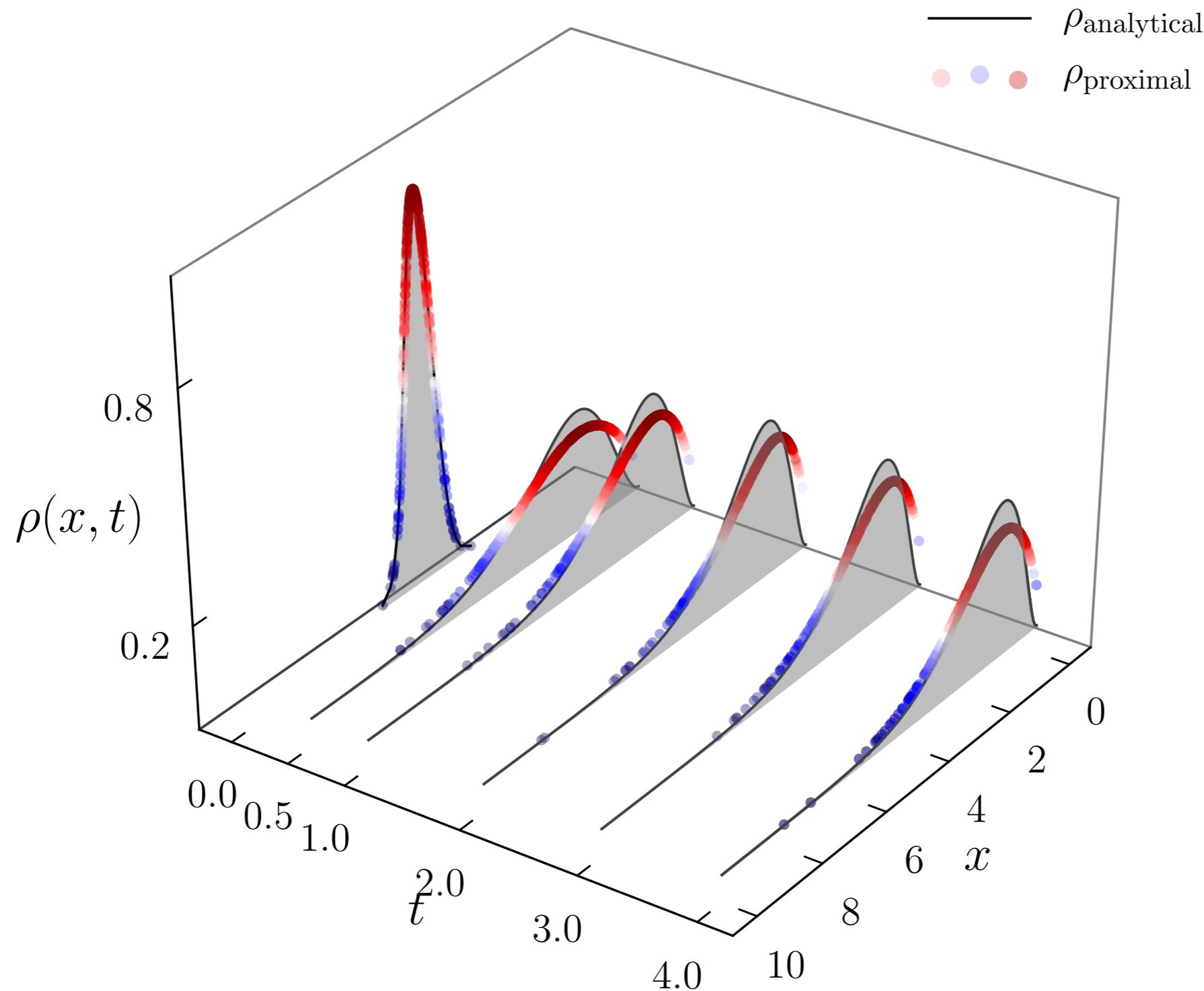
# Extensions: Nonlocal Interactions

$$U(\cdot) = V(\cdot) = \|\cdot\|_2^2$$



# Extensions: Multiplicative Noise

Cox-Ingersoll-Ross:  $dx = a(\theta - x) dt + b\sqrt{x} dw, 2a > b^2, \theta > 0$



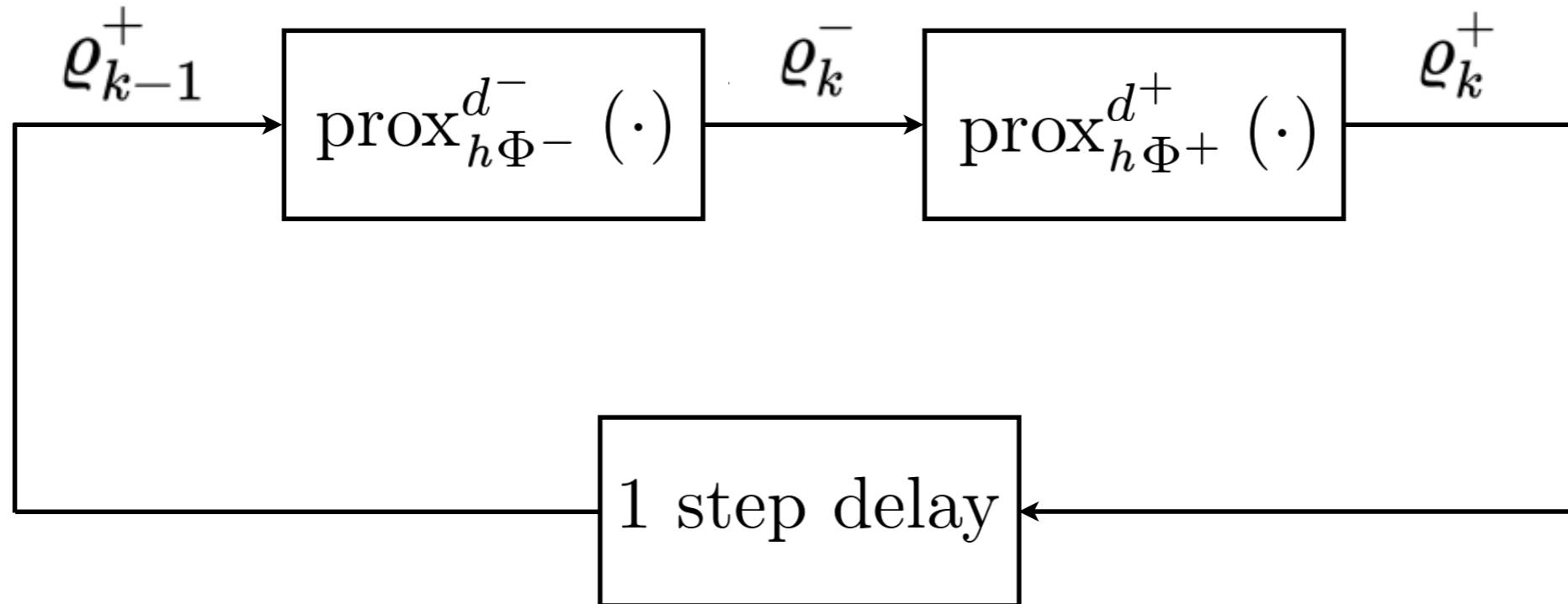
# Solving filtering as generalized gradient flow

# What's New?

Main idea: Solve the Kushner-Stratonovich SPDE

$$d\rho^+ = [\mathcal{L}_{\text{FP}} dt + \mathcal{L}(dz, dt, \rho^+)]\rho^+, \quad \rho(x, t=0) = \rho_0 \text{ as gradient flow in } \mathcal{P}_2(\mathcal{X})$$

Recursion of {deterministic  $\circ$  stochastic} proximal operators:



Convergence:  $\varrho_k^+(h) \rightarrow \rho^+(x, t = kh)$  as  $h \downarrow 0$

For prior, as before:  $d^- \equiv W^2, \quad \Phi^- \equiv \mathbb{E}_\varrho[\psi + \beta^{-1} \log \varrho]$

For posterior:  $d^+ \equiv d_{\text{FR}}^2 \text{ or } D_{\text{KL}}, \quad \Phi^+ \equiv \frac{1}{2} \mathbb{E}_{\varrho^+} [(y_k - h(x))^\top R^{-1} (y_k - h(x))]$

# Explicit Recovery of the Kalman-Bucy Filter

**Model:**

$$d\mathbf{x}(t) = \mathbf{A}\mathbf{x}(t)dt + \mathbf{B}d\mathbf{w}(t), \quad d\mathbf{w}(t) \sim \mathcal{N}(0, \mathbf{Q}dt)$$

$$d\mathbf{z}(t) = \mathbf{C}\mathbf{x}(t)dt + d\mathbf{v}(t), \quad d\mathbf{v}(t) \sim \mathcal{N}(0, \mathbf{R}dt)$$

**Given  $\mathbf{x}(0) \sim \mathcal{N}(\mu_0, \mathbf{P}_0)$ , want to recover:**

$$d\mu^+(t) = \mathbf{A}\mu^+(t)dt + \boxed{\mathbf{K}(t)}^\top (d\mathbf{z}(t) - \mathbf{C}\mu^+(t)dt),$$

$$\dot{\mathbf{P}}^+(t) = \mathbf{A}\mathbf{P}^+(t) + \mathbf{P}^+(t)\mathbf{A}^\top + \mathbf{B}\mathbf{Q}\mathbf{B}^\top - \mathbf{K}(t)\mathbf{R}\mathbf{K}(t)^\top.$$

— A.H. and T.T. Georgiou, Gradient Flows in Uncertainty Propagation and Filtering of Linear Gaussian Systems, *CDC 2017*.

— A.H. and T.T. Georgiou, Gradient Flows in Filtering and Fisher-Rao Geometry, *ACC 2018*.

# Explicit Recovery of the Wonham Filter

**Model:**

$$x(t) \sim \text{Markov}(Q), \\ dz(t) = h(x(t)) dt + \sigma_v(t) dv(t)$$

**State space:**  $\Omega := \{a_1, \dots, a_m\}$

**Posterior**  $\pi^+(t) := \{\pi_1^+(t), \dots, \pi_m^+(t)\}$  **solves the nonlinear SDE:**

$$d\pi^+(t) = \pi^+(t)Q dt + \frac{1}{(\sigma_v(t))^2} \pi^+(t) \left( H - \hat{h}(t)I \right) \left( dz(t) - \hat{h}(t)dt \right),$$

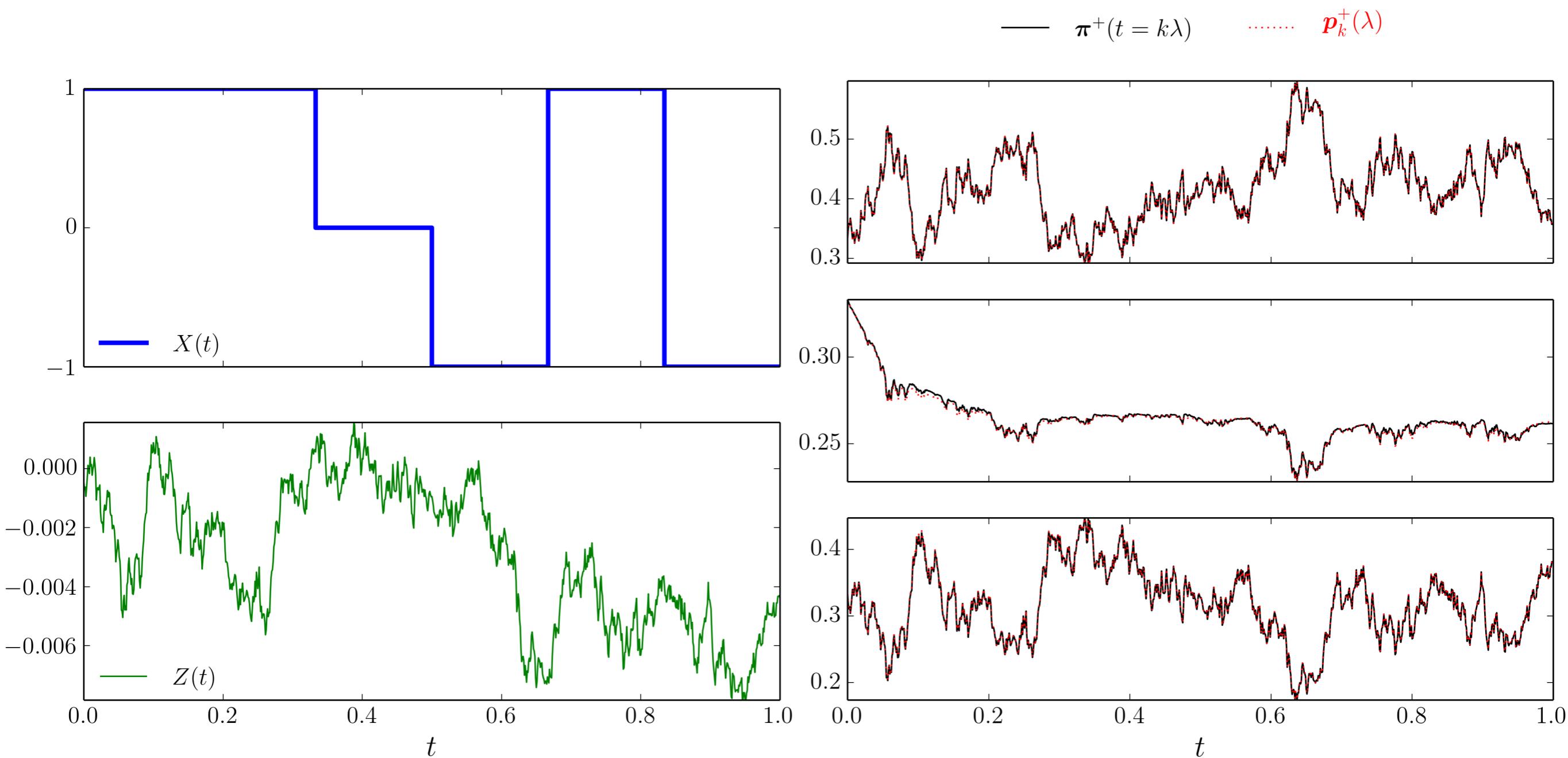
**where**  $H := \text{diag}(h(a_1), \dots, h(a_m)), \quad \hat{h}(t) := \sum_{i=1}^m h(a_i) \pi_i^+(t),$

**Initial condition:**  $\pi^+(t=0) = \pi_0,$

**By defn.**  $\pi^+(t) = \mathbb{P}(x(t) = a_i \mid z(s), 0 \leq s \leq t)$

— A.H. and T.T. Georgiou, Proximal Recursion for the Wonham Filter, *CDC 2019*.

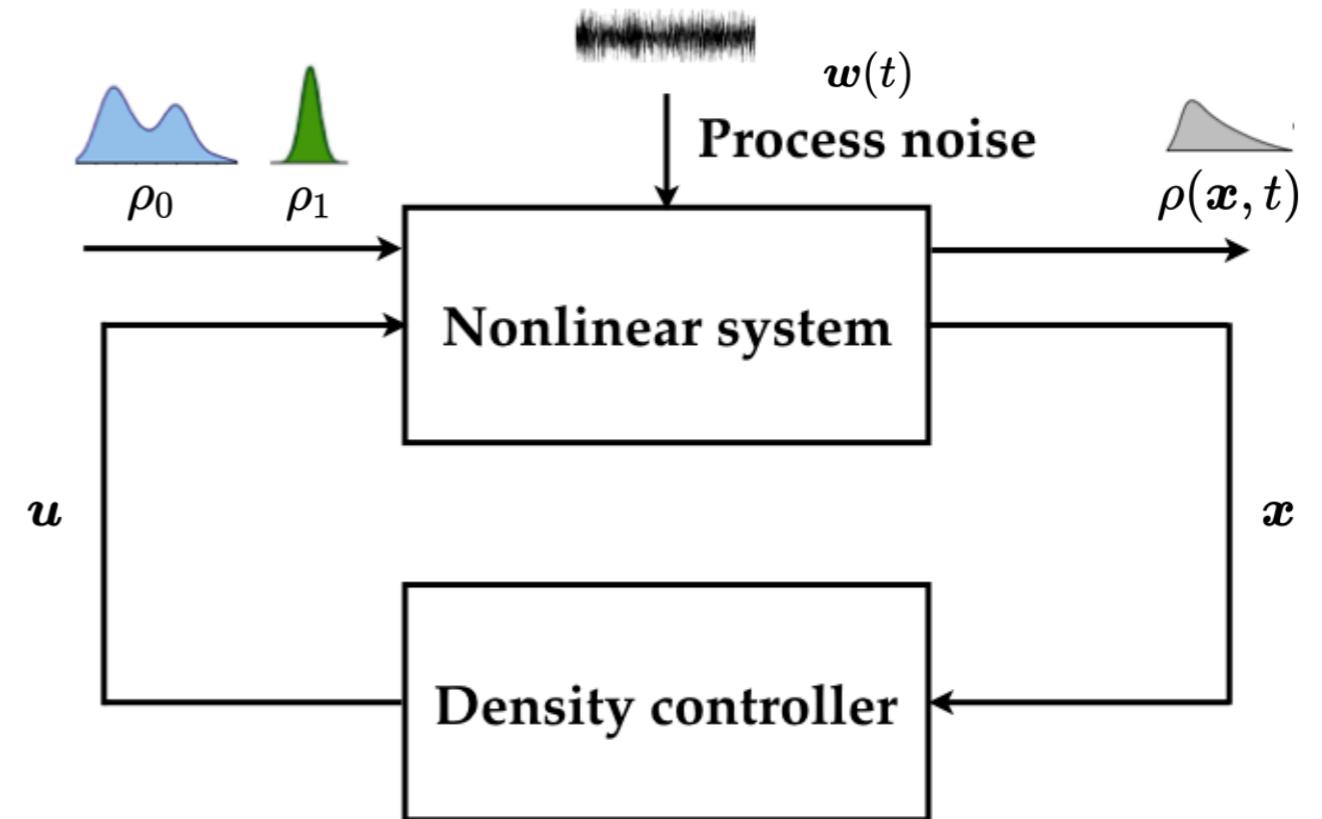
# Numerical Results for the Wonham Filter



# Solving density control as generalized gradient flow

# State Feedback Density Steering

Steer joint state PDF via feedback control over finite time horizon



Common scenario:  $G \equiv B$

$$\underset{u \in \mathcal{U}}{\text{minimize}} \quad \mathbb{E} \left[ \int_0^1 \left( \frac{1}{2} \|u(t, x_t^u)\|_2^2 + q(t, x_t^u) \right) dt \right]$$

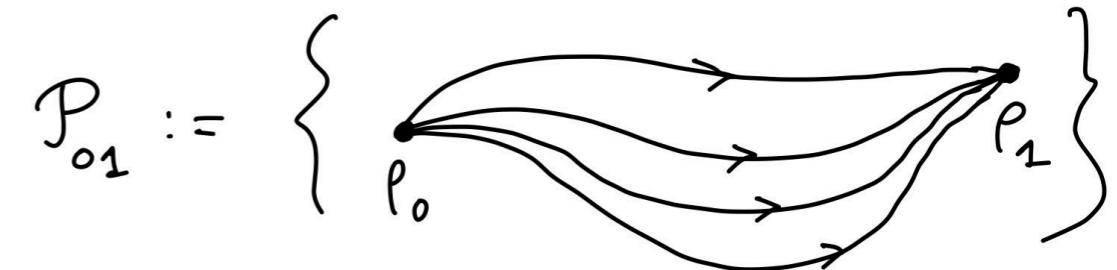
subject to

$$dx_t^u = \{f(t, x_t^u) + B(t, x_t^u)u\}dt + \sqrt{2}G(t, x_t^u)dw_t$$

$$x_0^u := x_t^u(t=0) \sim \rho_0, \quad x_1^u := x_t^u(t=1) \sim \rho_1$$

# Optimal Control Problem over PDFs

Diffusion tensor:  $D := GG^\top$



Hessian operator w.r.t. state: Hess

$$\inf_{(\rho, u) \in \mathcal{P}_{01} \times \mathcal{U}} \int_{\mathbb{R}^n} \int_0^1 \left( \frac{1}{2} \|u(t, x_t^u)\|_2^2 + q(t, x_t^u) \right) \rho(t, x_t^u) \, dt \, dx_t^u$$

subject to

$$\frac{\partial \rho}{\partial t} + \nabla \cdot ((f + Bu) \rho) = \langle \text{Hess}, D\rho \rangle$$

$$\rho(t=0, x_0^u) = \rho_0, \quad \rho(t=1, x_1^u) = \rho_1$$

# Optimal Control Problem over PDFs

Existence-uniqueness needs regularity assumptions

Are known to hold for many practical classes of nonlinearities

This talk: will focus on a few important classes

# Necessary Conditions of Optimality (Assuming $G \equiv B$ )

Coupled nonlinear PDEs + linear boundary conditions

Controlled Fokker-Planck or Kolmogorov's forward PDE

$$\frac{\partial \rho^{\text{opt}}}{\partial t} + \nabla \cdot ((f + D\nabla \psi) \rho^{\text{opt}}) = \langle \text{Hess}, D\rho \rangle$$

Hamilton-Jacobi-Bellman-like PDE

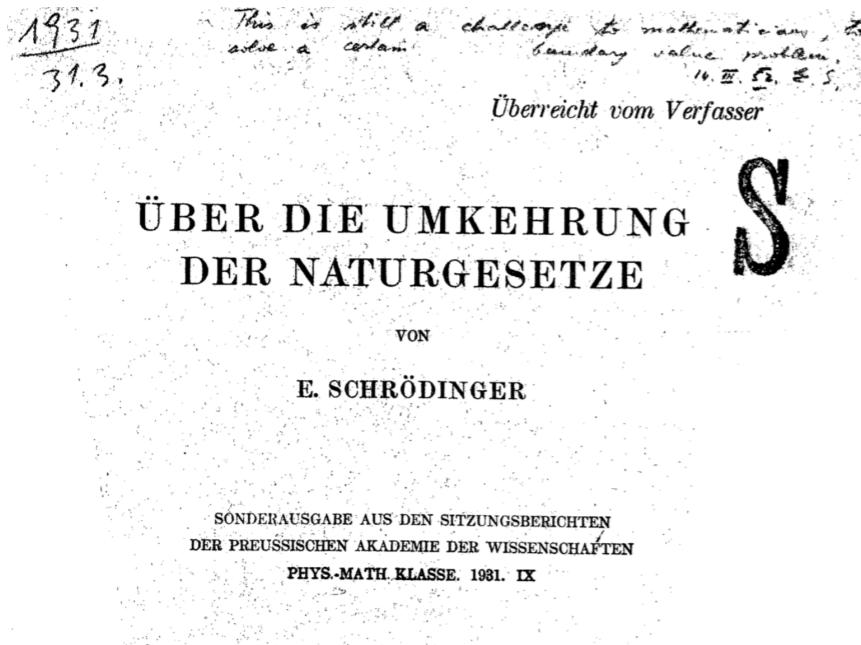
$$\frac{\partial \psi}{\partial t} + \langle \nabla \psi, f \rangle + \langle D, \text{Hess}(\psi) \rangle + \frac{1}{2} \langle \nabla \psi, D \nabla \psi \rangle = q$$

Boundary conditions:

$$\rho^{\text{opt}}(\cdot, t=0) = \rho_0, \quad \rho^{\text{opt}}(\cdot, t=1) = \rho_1$$

Optimal control:  $u^{\text{opt}} = B^\top \nabla \psi$

# Feedback Synthesis via the Schrödinger System



Sur la théorie relativiste de l'électron  
et l'interprétation de la mécanique quantique

PAR  
E. SCHRÖDINGER

## I. — Introduction

J'ai l'intention d'exposer dans ces conférences diverses idées concernant la mécanique quantique et l'interprétation qu'on en donne généralement à l'heure actuelle ; je parlerai principalement de la théorie quantique relativiste du mouvement de l'électron. Autant que nous pouvons nous en rendre compte aujourd'hui, il semble à peu près sûr que la mécanique quantique de l'électron, sous sa forme idéale, *que nous ne possédons pas encore*, doit former un jour la base de toute la physique. A cet intérêt tout à fait général, s'ajoute, ici à Paris, un intérêt particulier : vous savez tous que les bases de la théorie moderne de l'électron ont été posées à Paris par votre célèbre compatriote Louis de BROGLIE.



Hopf-Cole a.k.a. Fleming's logarithmic transform:

$$(\rho^{\text{opt}}, \psi) \mapsto (\hat{\varphi}, \varphi) \quad \text{— Schrödinger factors}$$

$$\hat{\varphi}(x, t) = \rho^{\text{opt}}(x, t) \exp(-\psi(x, t))$$

$$\varphi(x, t) = \exp(\psi(x, t)) \quad \text{for all } (x, t) \in \mathbb{R}^n \times [0, 1]$$

# Feedback Synthesis via the Schrödinger System

2 coupled nonlinear PDEs → boundary-coupled linear PDEs!!

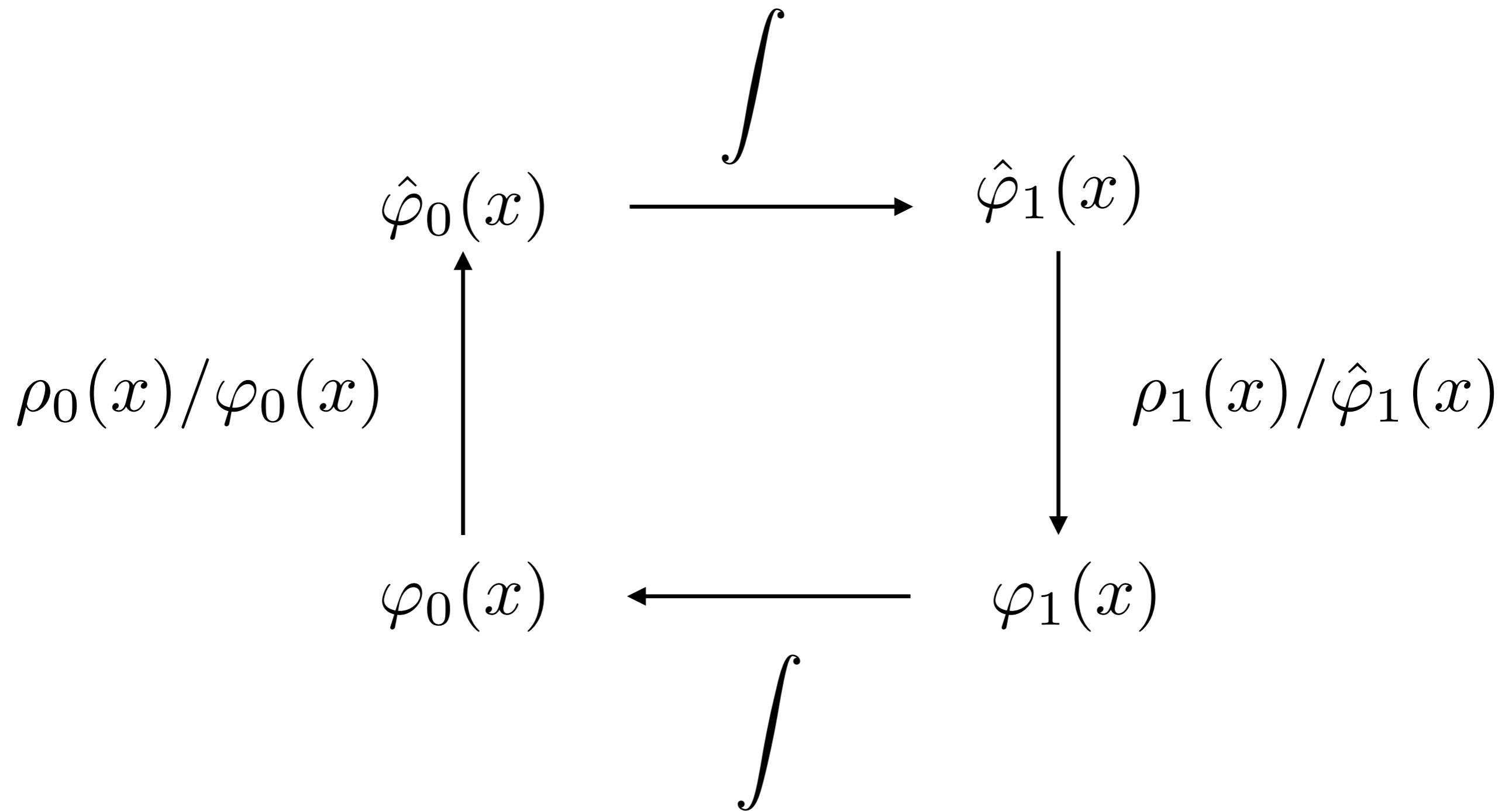
**Uncontrolled forward-backward Kolmogorov PDEs:**

$$\begin{aligned}\frac{\partial \hat{\varphi}}{\partial t} &= -\nabla \cdot (\hat{\varphi} f) + \langle \text{Hess}, D\hat{\varphi} \rangle - q\hat{\varphi}, & \hat{\varphi}_0 \varphi_0 &= \rho_0, \\ \frac{\partial \varphi}{\partial t} &= -\langle \nabla \varphi, f \rangle - \langle \text{Hess}(\varphi), D \rangle + q\varphi, & \hat{\varphi}_1 \varphi_1 &= \rho_1,\end{aligned}$$

Optimal controlled joint state PDF:  $\rho^{\text{opt}}(x, t) = \hat{\varphi}(x, t)\varphi(x, t)$

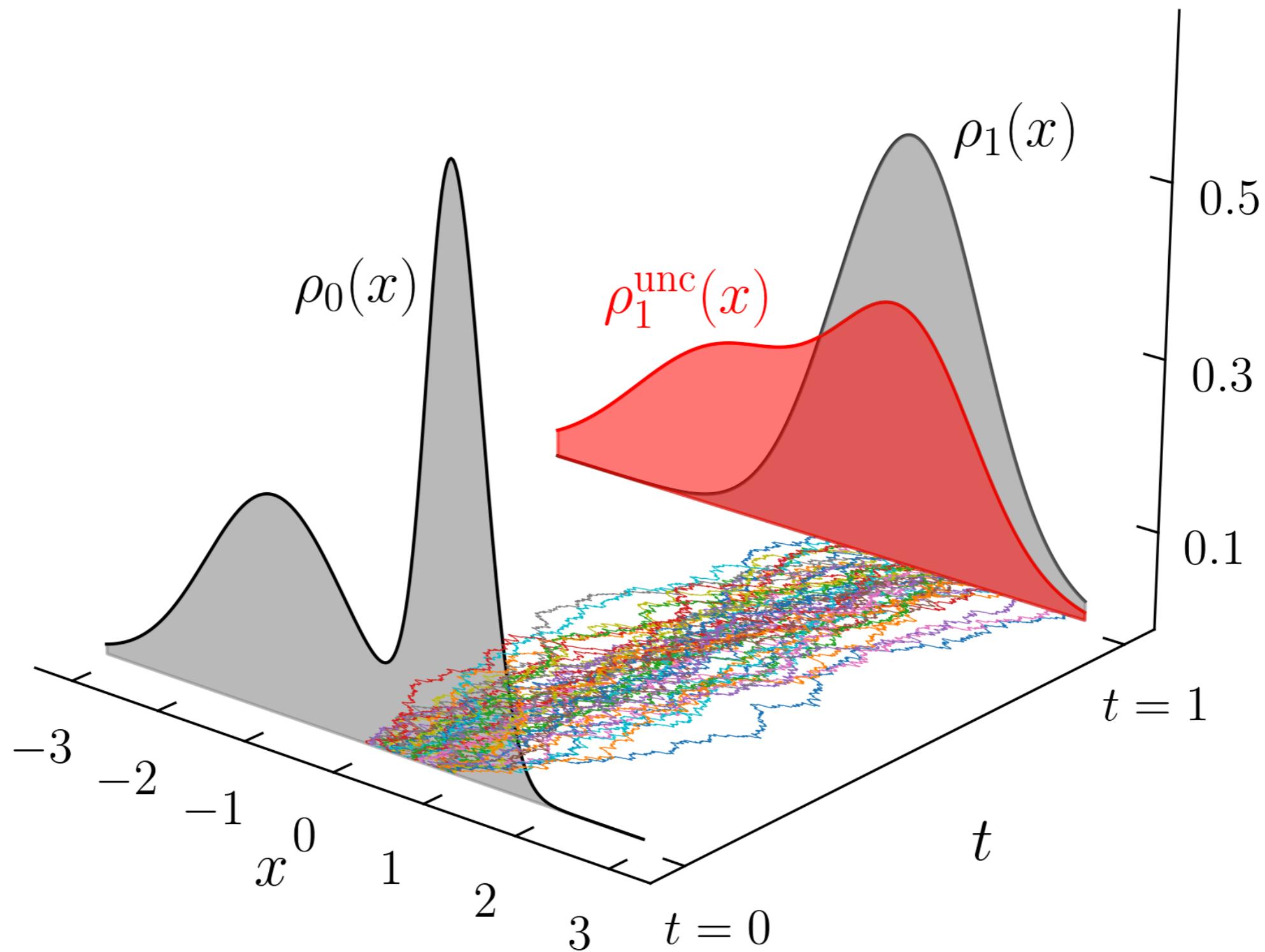
Optimal control:  $u^{\text{opt}}(x, t) = 2B^\top \nabla_x \log \varphi(x, t)$

# Fixed Point Recursion over $(\hat{\varphi}_0, \varphi_1)$



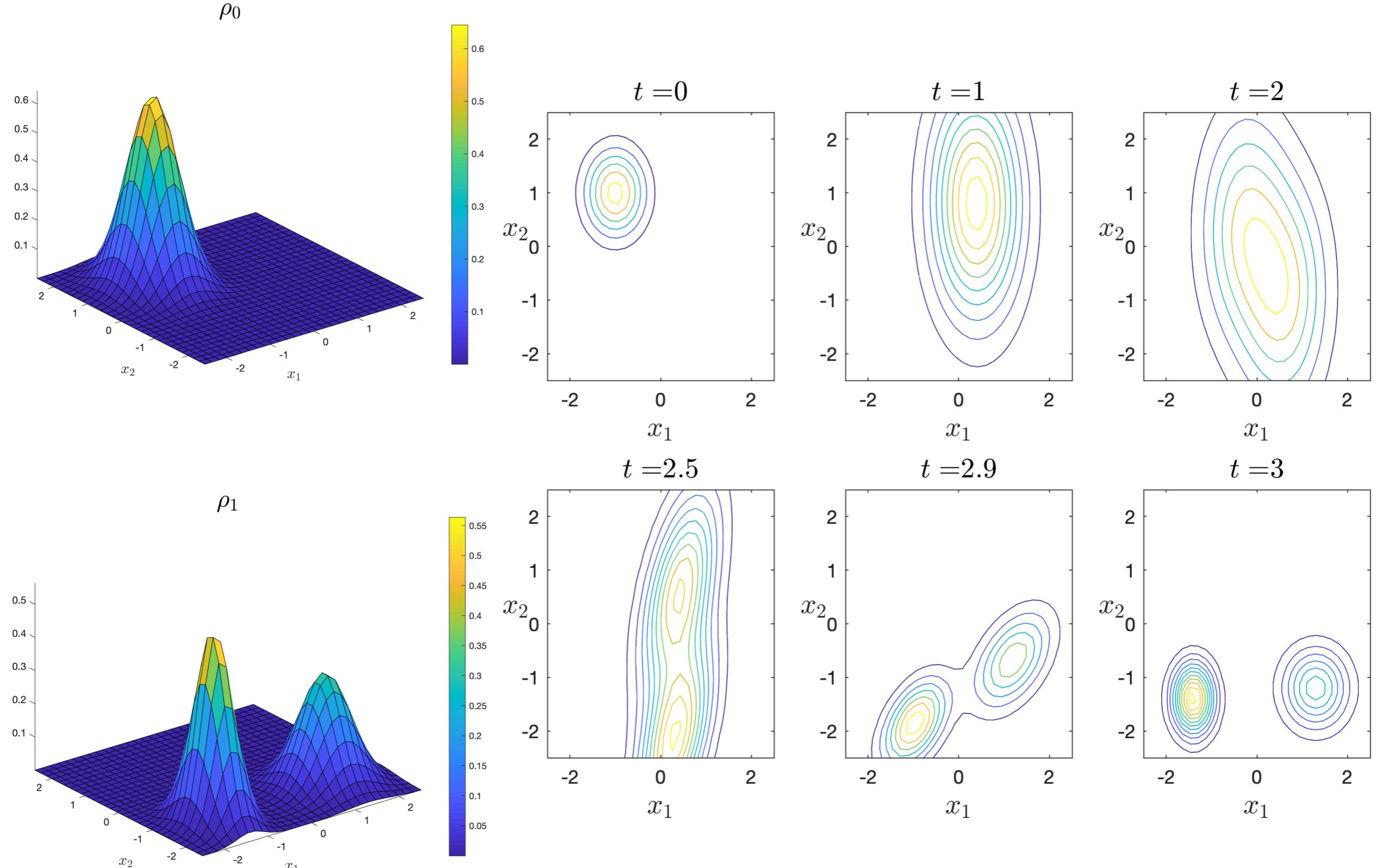
**This recursion is contractive in the Hilbert metric!!**

# Feedback Density Control: $f \equiv 0, B = G \equiv I, q \equiv 0$



Zero prior dynamics

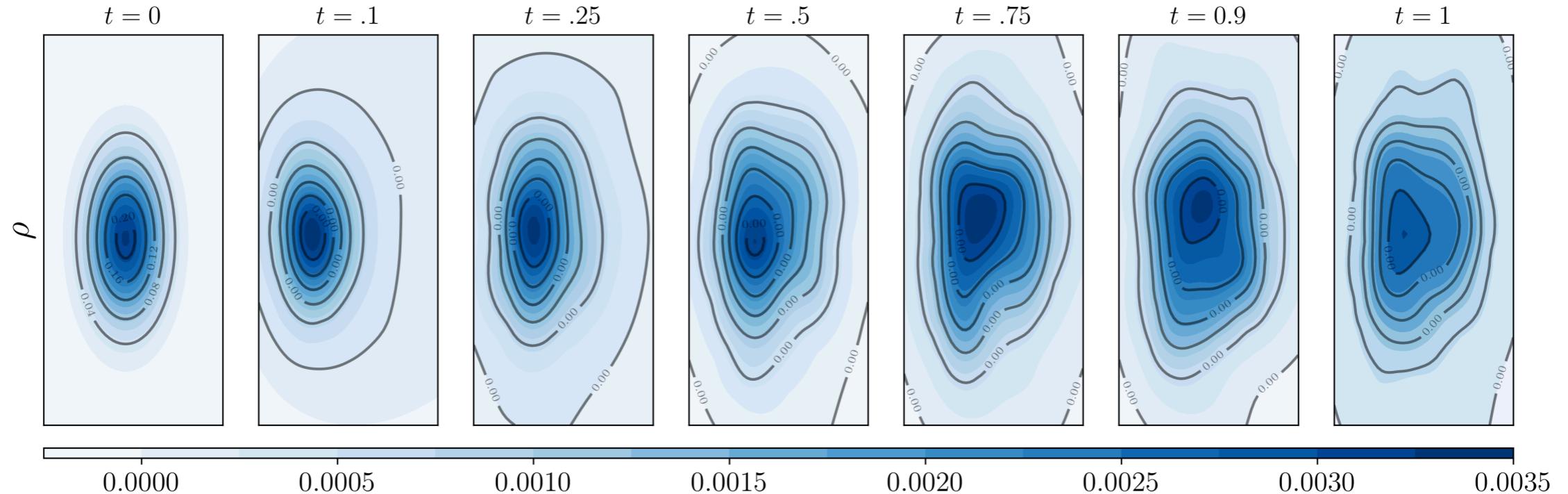
# Feedback Density Control: $f \equiv Ax, B = G, q \equiv 0$



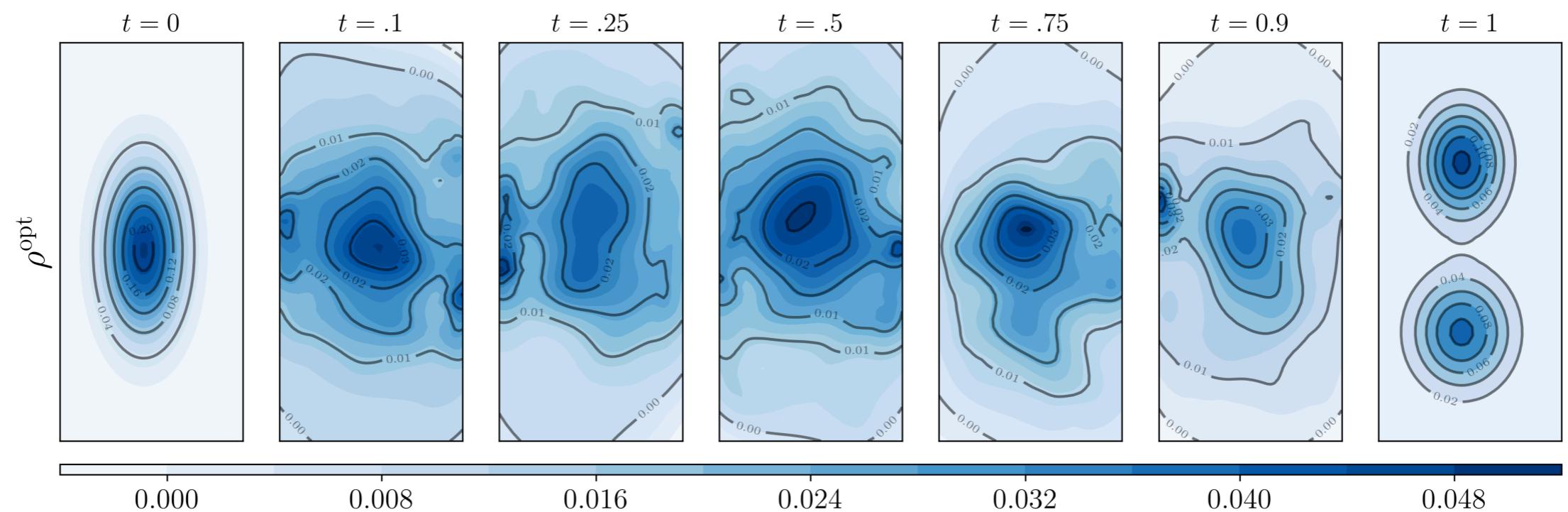
Linear prior dynamics

# Feedback Density Control: Nonlinear Grad. Drift

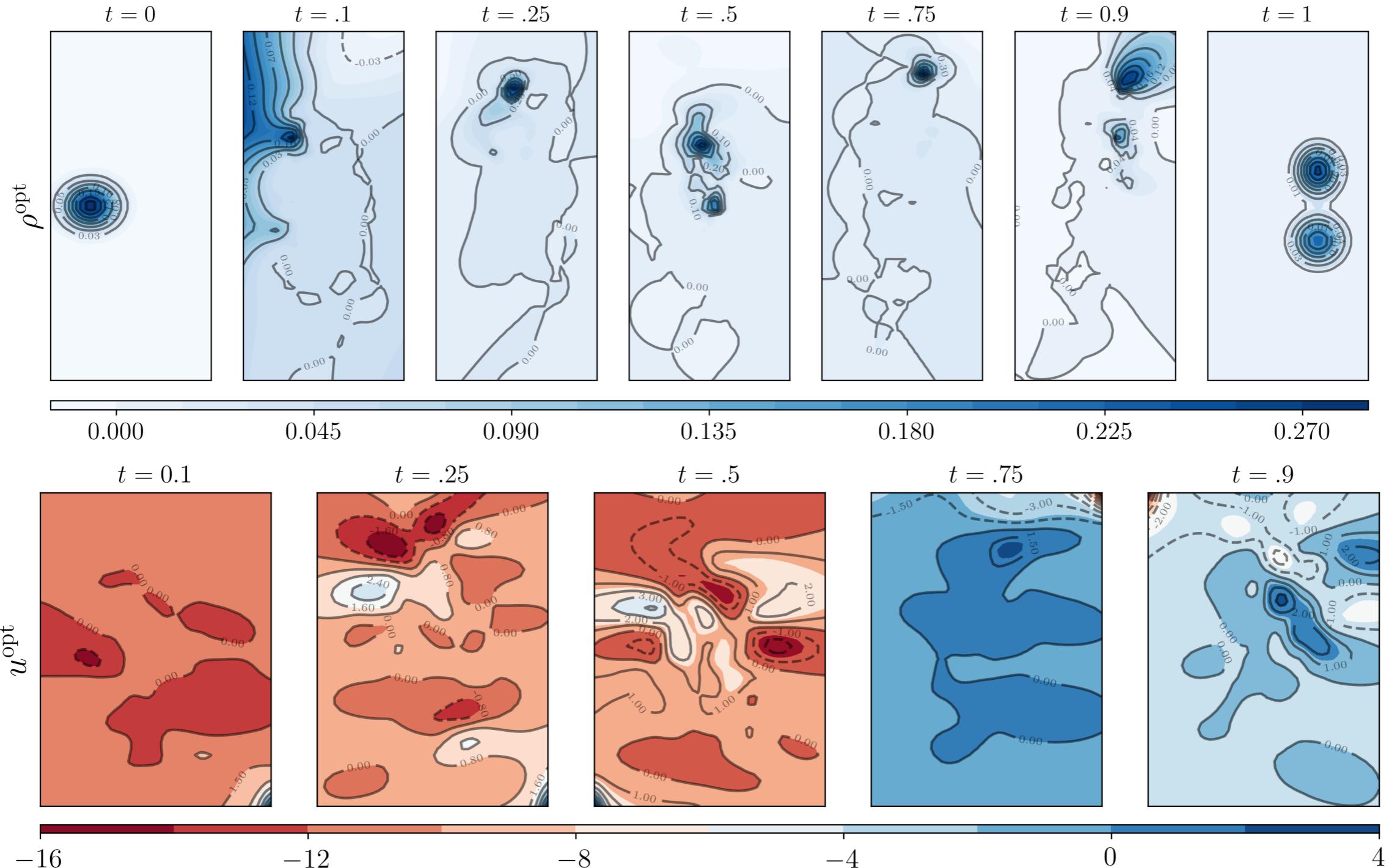
Uncontrolled joint PDF evolution:



Optimal controlled joint PDF evolution:

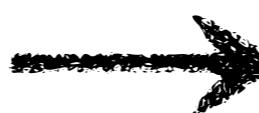


# Feedback Density Control: Mixed Conservative-Dissipative Drift



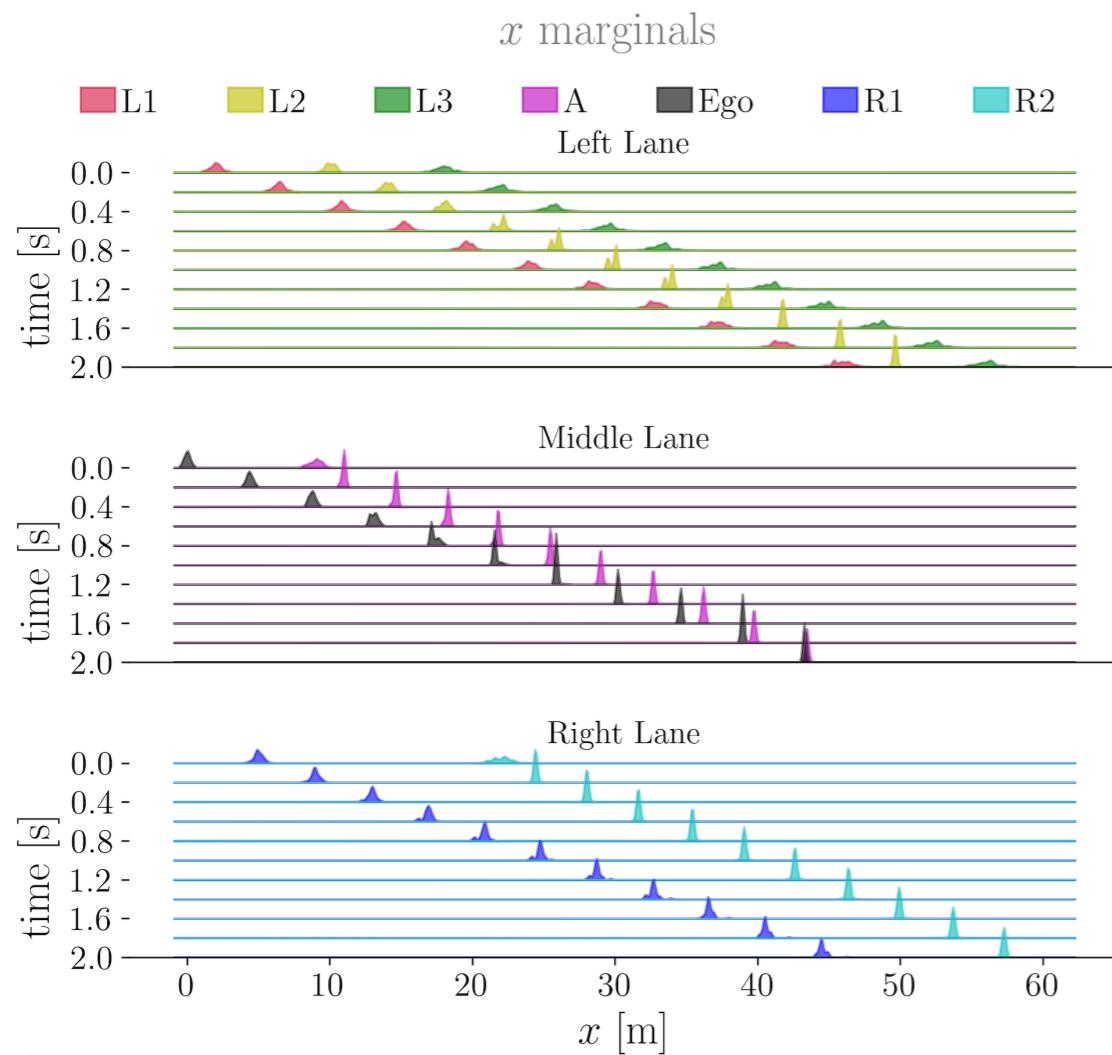
— K.F. Caluya and A.H., Wasserstein proximal algorithms for the Schrodinger bridge problem: density control with nonlinear drift, *IEEE TAC* 2021.

# Density Prediction for Safe Automated Driving

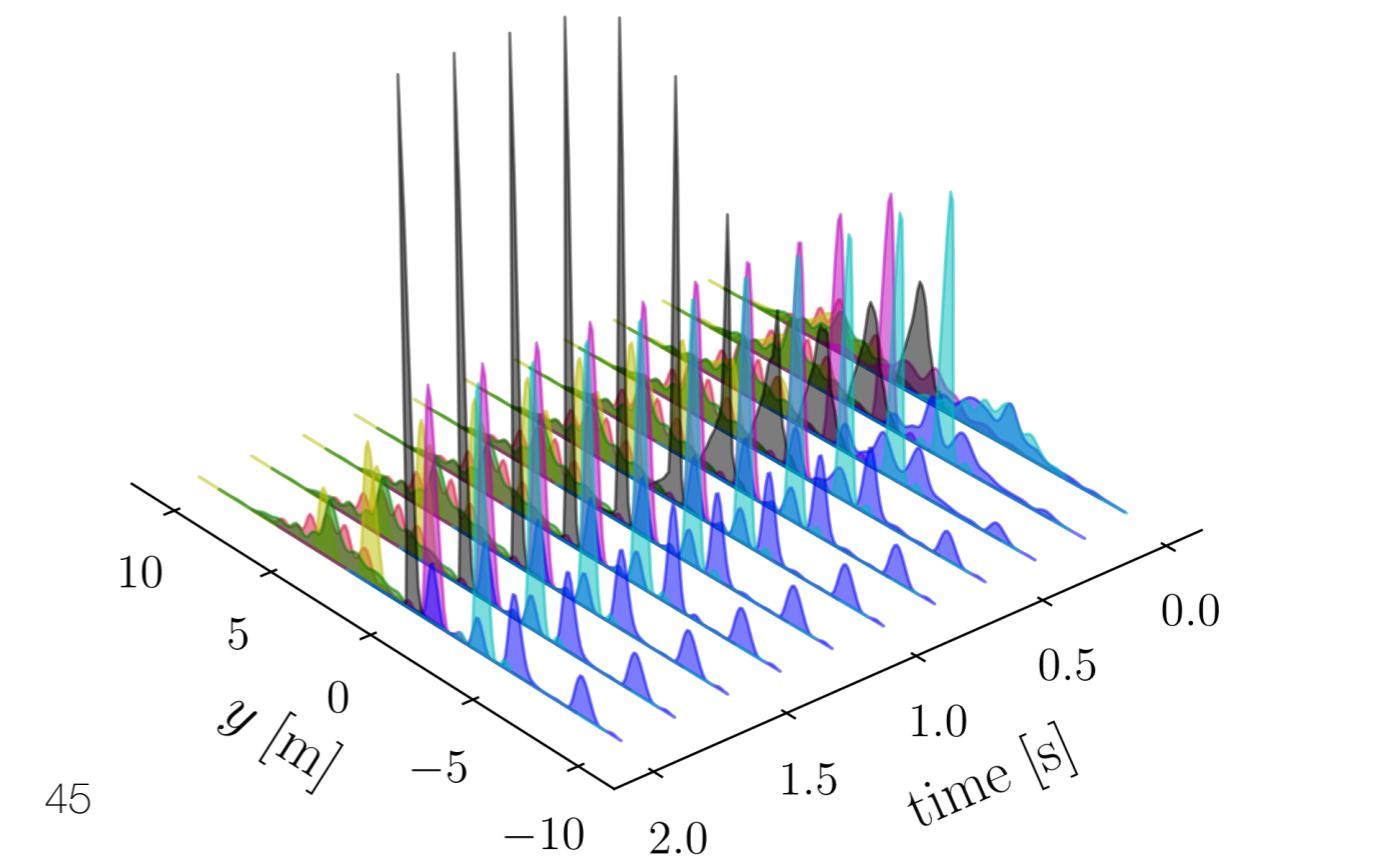


$t_1$

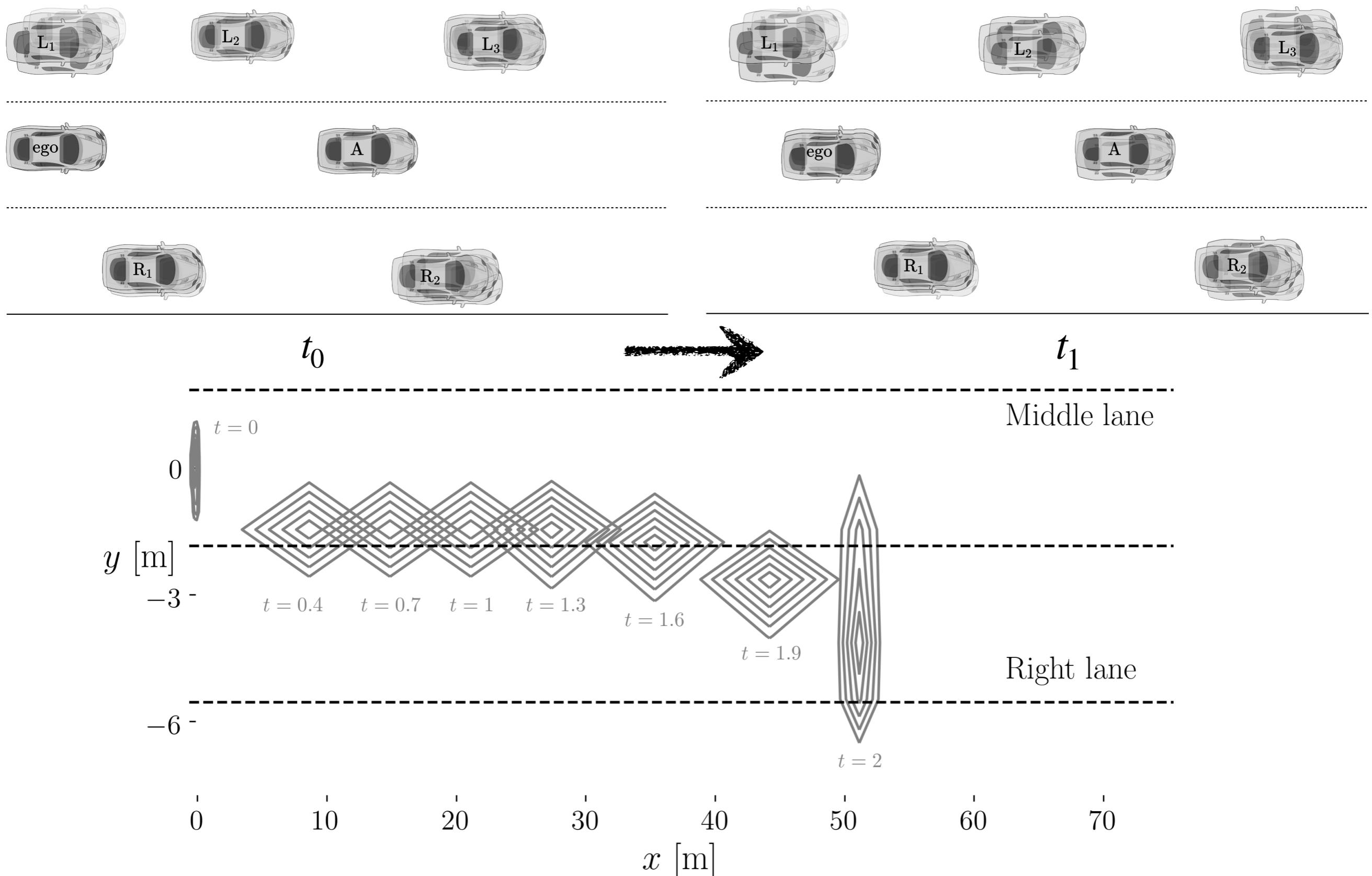
$y$  marginals



L1 L2 L3 A Ego R1 R2



# Density Control for Safe Automated Driving



# Application to Safe Automated Driving

S. Haddad, A.H., and B. Singh, Density-based stochastic reachability computation for occupancy prediction in automated driving, *IEEE Transactions on Control Systems Technology*, 2022.

S. Haddad, K.F. Caluya, A.H., and B. Singh, Prediction and optimal feedback steering of probability density functions for safe automated driving, *IEEE Control Systems Letters*, 2021.

# Learning a neural network as generalized gradient flow

# Learning Neural Network from Data

(feature vector, label) =  $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ ,  $i = 1, \dots, n$

Consider shallow NN: 1 hidden layer with  $n_H$  neurons

NN parameter vector  $\boldsymbol{\theta} := (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_{n_H})^\top \in \mathbb{R}^{pn_H}$

Approximating function:

$$\hat{f}(\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{n_H} \sum_{i=1}^{n_H} \Phi(\mathbf{x}, \boldsymbol{\theta}_i), \text{ example: } \Phi(\mathbf{x}, \boldsymbol{\theta}_i) = a_i \sigma(\mathbf{w}_i^\top \mathbf{x} + b_i)$$

Population risk functional:

$$R(\hat{f}) = \mathbb{E}_{(\mathbf{x}, y)} \left[ (y - \hat{f}(\mathbf{x}, \boldsymbol{\theta}))^2 \right] \approx \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i, \boldsymbol{\theta}))^2$$

Learning problem: minimize  $R(\hat{f})$   
 $\boldsymbol{\theta} \in \mathbb{R}^{pn_H}$

# Mean Field Density Dynamics of SGD

Free energy functional:  $F(\rho) := R(\hat{f}(\mathbf{x}, \rho))$

For quadratic loss:

$$F(\rho) = \underbrace{F_0}_{\text{independent of } \rho} + \underbrace{\int_{\mathbb{R}^p} V(\boldsymbol{\theta}) \rho(\boldsymbol{\theta}) d\boldsymbol{\theta}}_{\text{advection potential energy, linear in } \rho} + \underbrace{\int_{\mathbb{R}^p} \int_{\mathbb{R}^p} U(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \rho(\boldsymbol{\theta}) \rho(\tilde{\boldsymbol{\theta}}) d\boldsymbol{\theta} d\tilde{\boldsymbol{\theta}}}_{\text{interaction potential energy, nonlinear in } \rho} ,$$

where

$$F_0 := \mathbb{E}_{(\mathbf{x}, y)} [y^2], \quad V(\boldsymbol{\theta}) := \mathbb{E}_{(\mathbf{x}, y)} [-2y\Phi(\mathbf{x}, \boldsymbol{\theta})],$$

PDF dynamics for SGD:

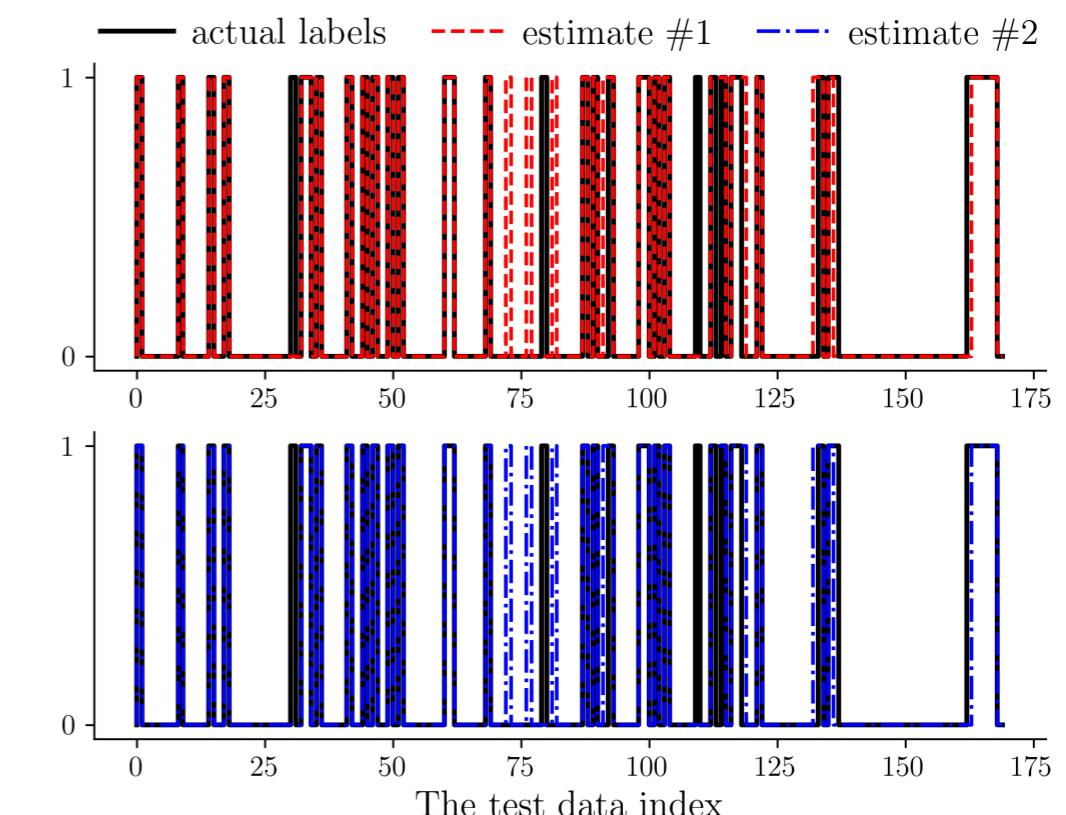
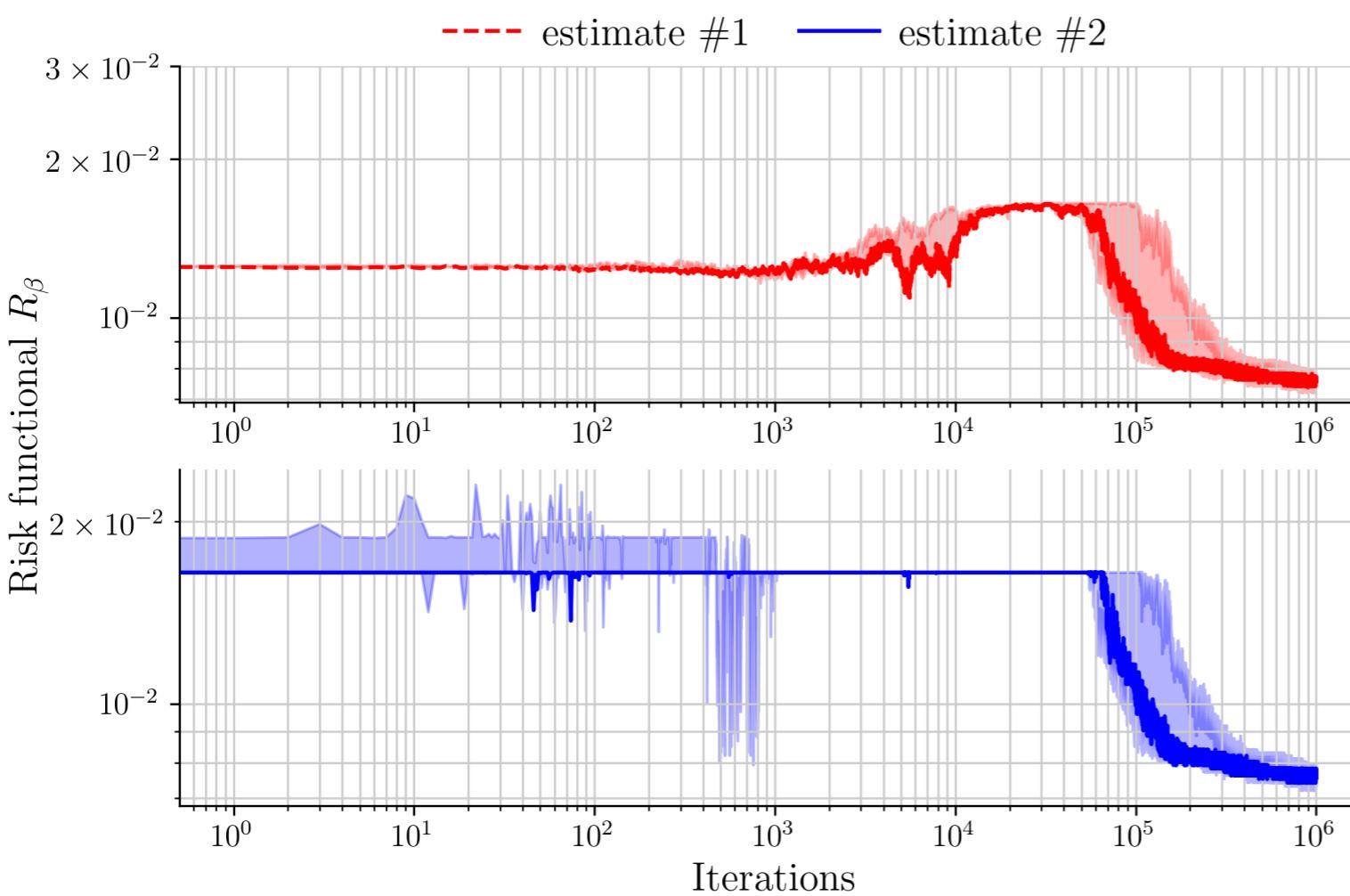
$$\frac{\partial \rho}{\partial t} = \nabla \cdot (\rho \nabla (\underbrace{V + U \circledast \rho}_{\frac{\delta F}{\delta \rho}})), \text{ where } (U \circledast \rho)(\boldsymbol{\theta}) := \int_{\mathbb{R}^p} U(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \rho(\tilde{\boldsymbol{\theta}}) d\tilde{\boldsymbol{\theta}}$$

This PDE is the gradient flow of functional  $F$  w.r.t. the Wasserstein metric  $W$

# Proximal Recursion for SGD Training of NN

$$\begin{aligned} \varrho_k(\tau, \theta) &= \arg \min_{\varrho \in \mathcal{P}(\mathbb{R}^p)} \frac{1}{2} (W(\varrho(\theta), \varrho_{k-1}(\tau, \theta)))^2 + \tau F(\varrho(\theta)) \\ &= \text{prox}_{\tau F}^W (\varrho_{k-1}) \end{aligned}$$

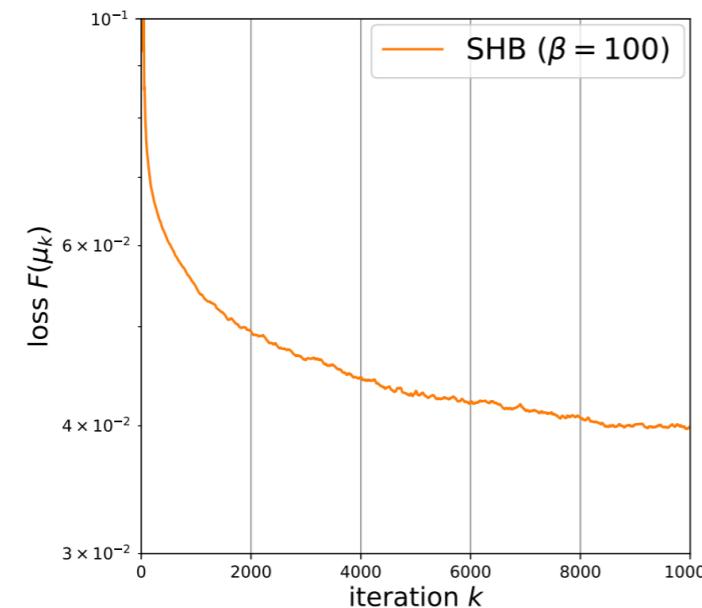
## Case study: Wisconsin Breast Cancer (Diagnostic) Data Set



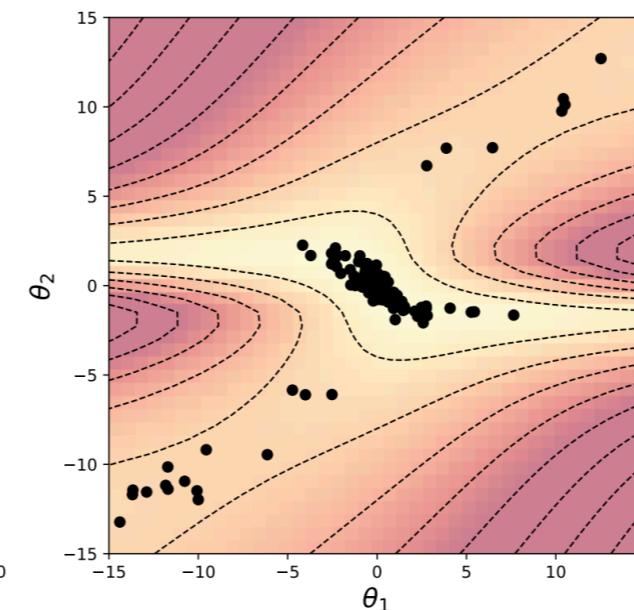
Classification accuracy for the WBDC dataset		
$\beta$	Estimate #1	Estimate #2
0.03	91.17%	92.35%
0.05	92.94%	92.94%
0.07	78.23%	92.94%

# Mean Field Density Dynamics of Stoc. Heavy Ball

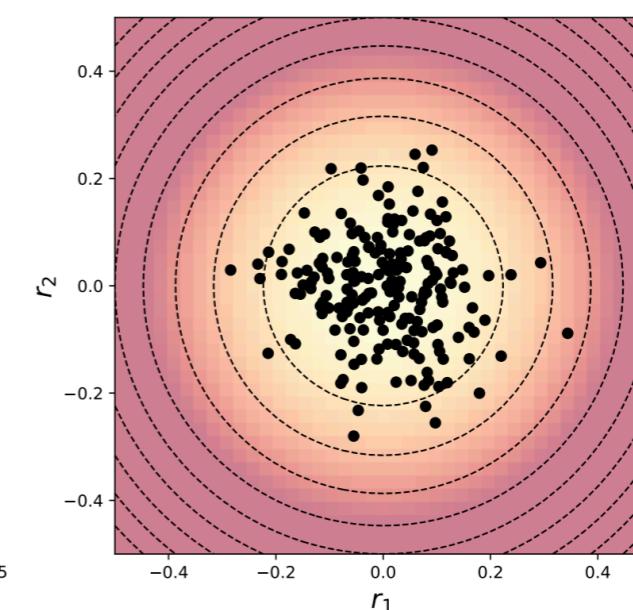
$$\partial_t \mu_t = -\nabla \cdot \left[ \mu_t \cdot \begin{pmatrix} r \\ -\nabla F'([\mu_t]^\theta) - \gamma r \end{pmatrix} \right] + \gamma \beta^{-1} \Delta_r \mu_t$$



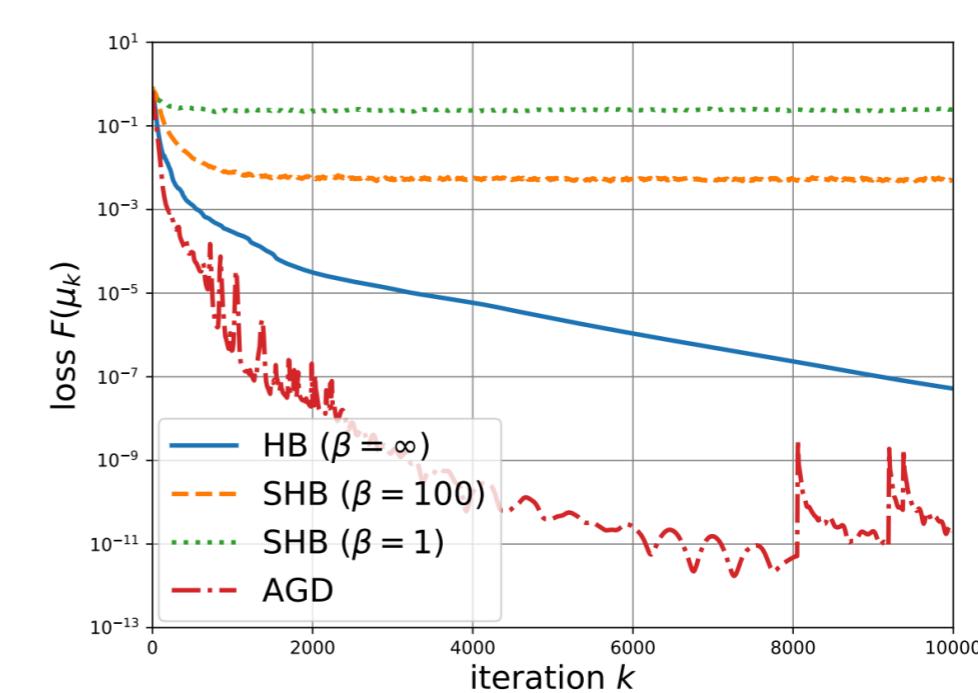
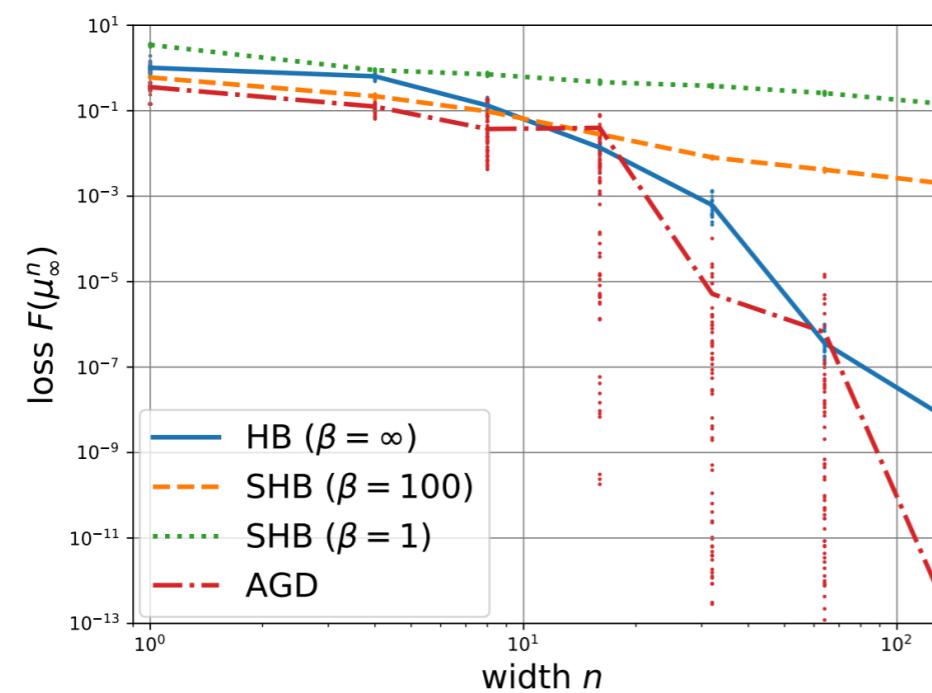
(a) Loss  $F(\mu_k^n)$



(b)  $\theta$  marginal



(c)  $r$  marginal



# Summary



# Thank You

Support:



CITRIS  
PEOPLE AND  
ROBOTS

