

Wasserstein Gradient Flow 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 Kenneth F. Caluya (UC Santa Cruz),
Tryphon T. Georgiou (UC Irvine), Walid Krichene (Google)

Controls, Autonomy and Robotics Seminar, UT Austin, TX, Nov. 18, 2020

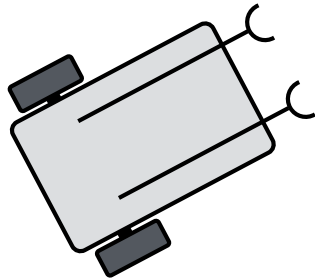


Overarching Theme

Systems-control theory for densities

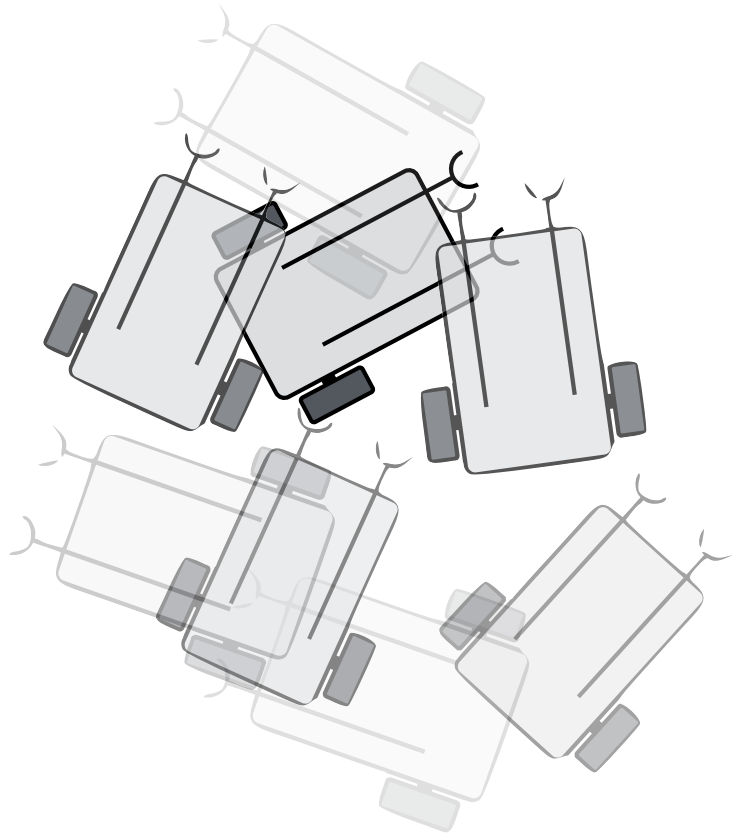
What is density?

Probability Density Fn.



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

Probability Density Fn.

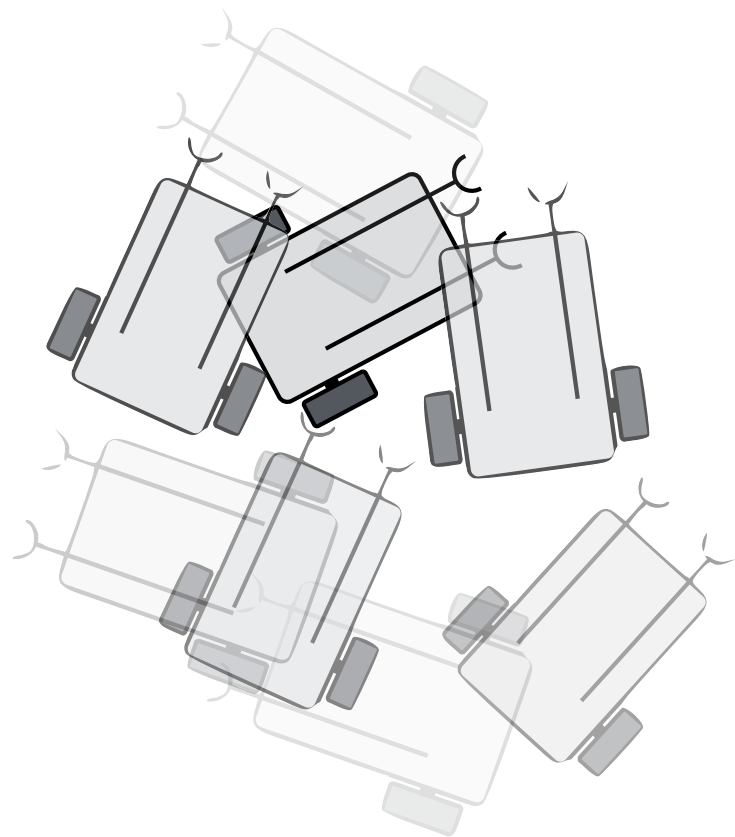


$$\mathbf{x}(t) \in \begin{pmatrix} x \\ y \\ \theta \end{pmatrix} \in \mathcal{X} \equiv \mathbb{R}^2 \times \mathbb{S}^1$$

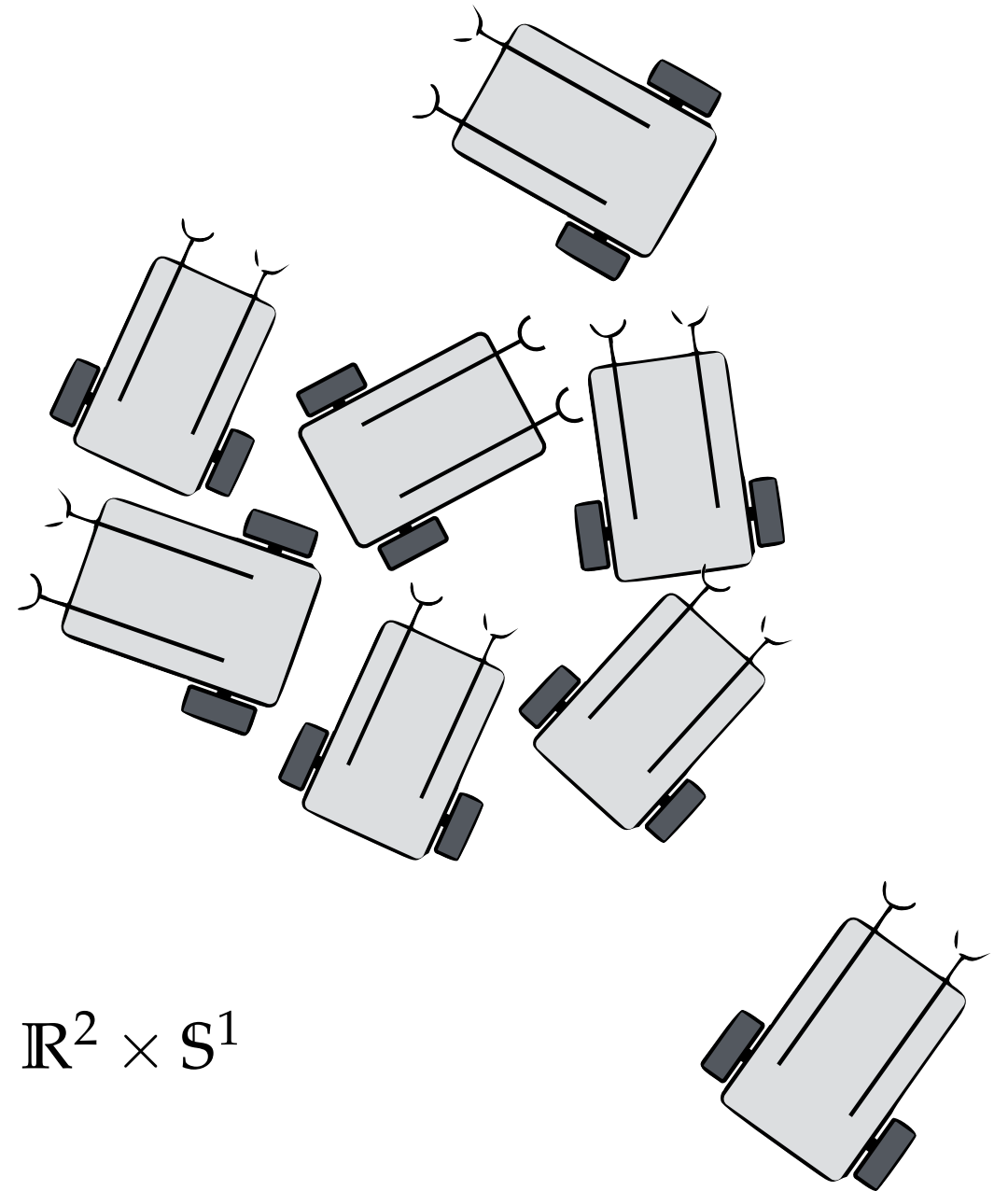
$$\rho(\mathbf{x}, t) : \mathcal{X} \times [0, \infty) \mapsto \mathbb{R}_{\geq 0}$$

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

Probability Density Fn.



Population Density Fn.



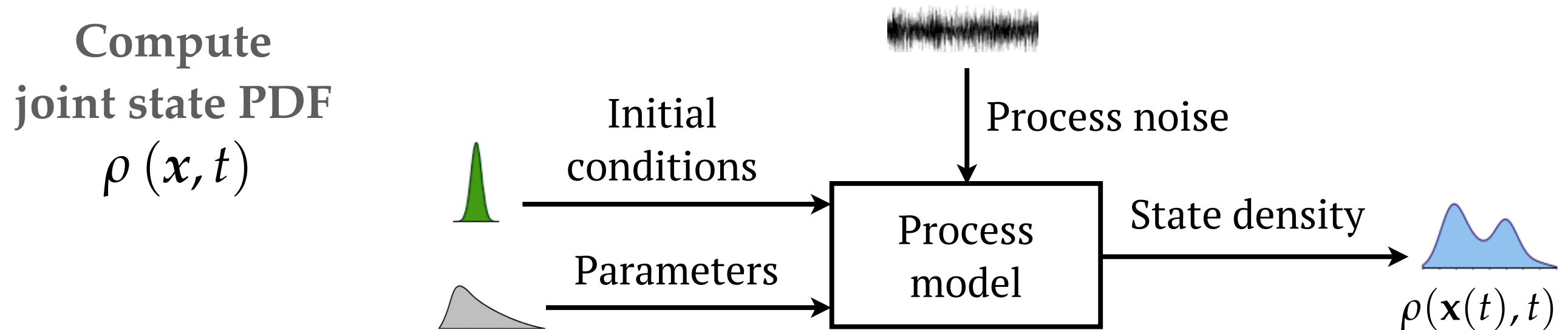
$$\mathbf{x}(t) \in \begin{pmatrix} x \\ y \\ \theta \end{pmatrix} \in \mathcal{X} \equiv \mathbb{R}^2 \times \mathbb{S}^1$$

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

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

Why care about densities?

Prediction Problem



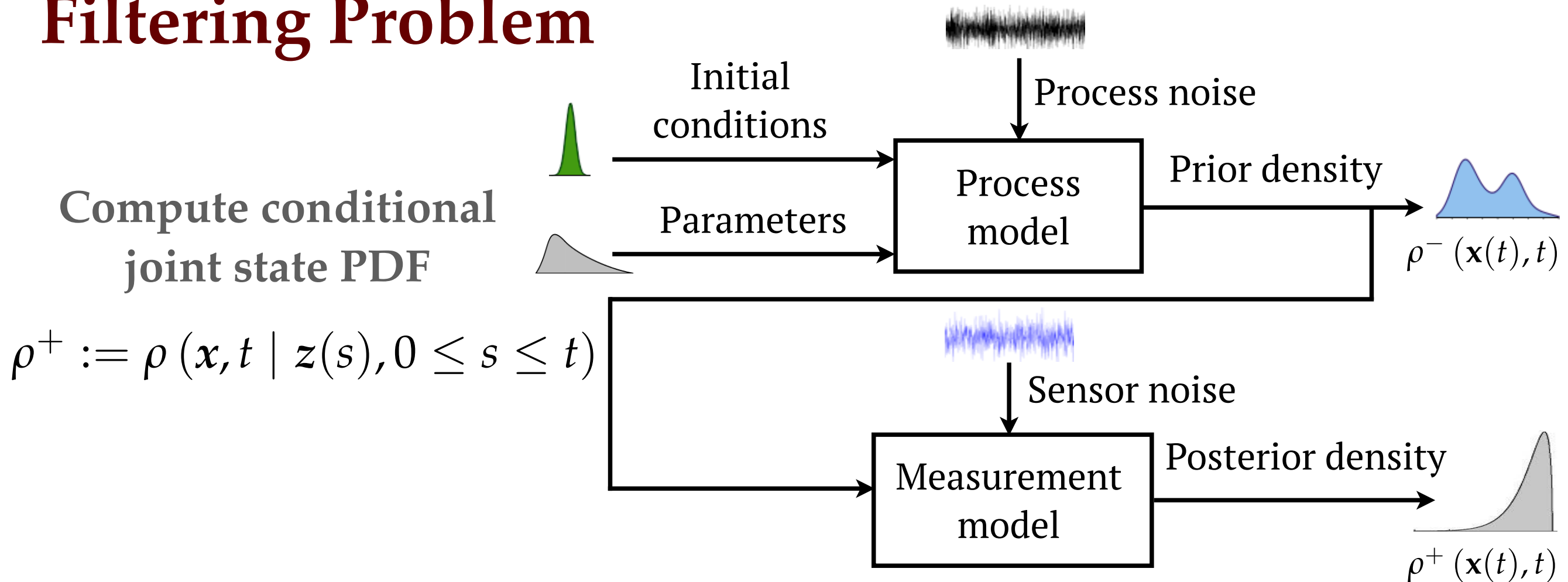
Trajectory flow:

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

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



Trajectory flow:

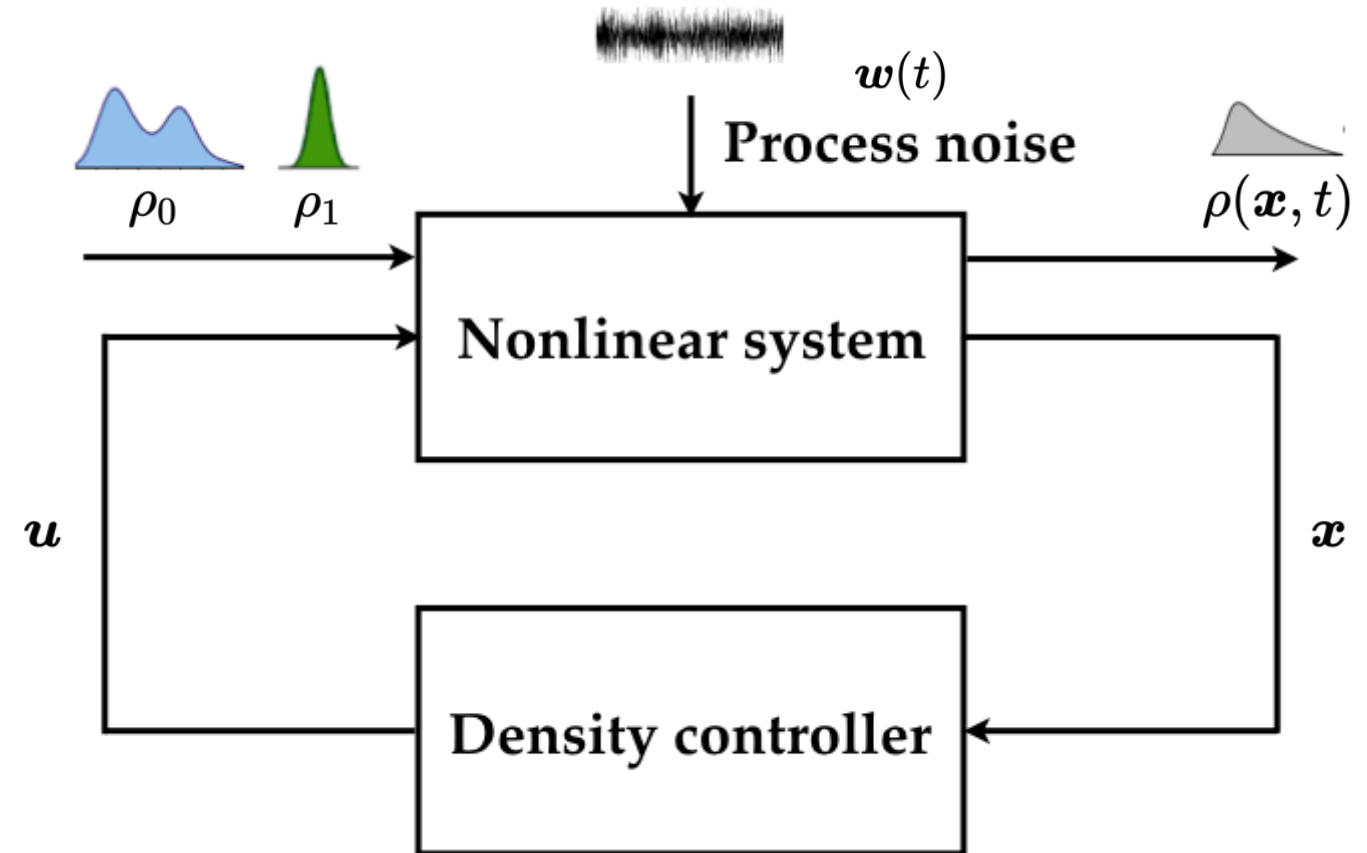
$$\begin{aligned} d\mathbf{X}(t) &= \mathbf{f}(\mathbf{X}, t) dt + \mathbf{g}(\mathbf{X}, t) d\mathbf{w}(t), & d\mathbf{w}(t) &\sim \mathcal{N}(0, \mathbf{Q}dt) \\ d\mathbf{Z}(t) &= \mathbf{h}(\mathbf{X}, t) dt + d\mathbf{v}(t), & d\mathbf{v}(t) &\sim \mathcal{N}(0, \mathbf{R}dt) \end{aligned}$$

Density flow:

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

Control Problem

Steer joint state PDF via feedback control over finite time horizon



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

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

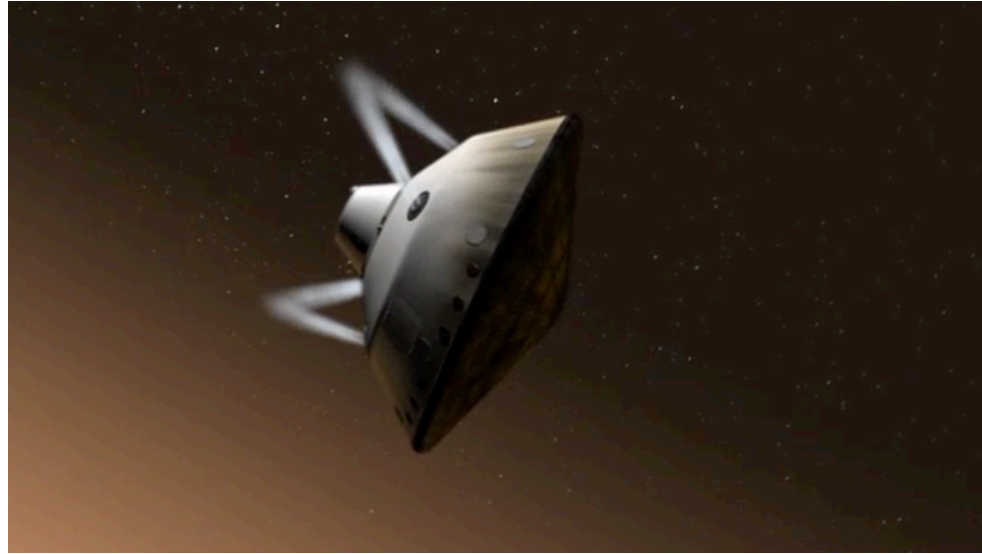
PDF dynamics:

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

PDFs in Mars Entry-Descent-Landing

Prediction problem

Filtering problem



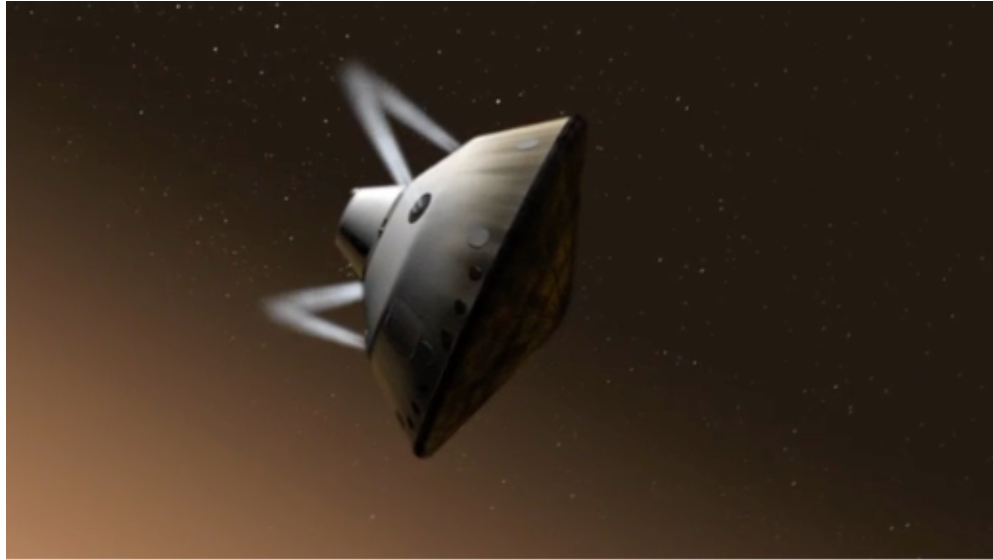
Predict heating rate uncertainty

Control problem

Learning problem

PDFs in Mars Entry-Descent-Landing

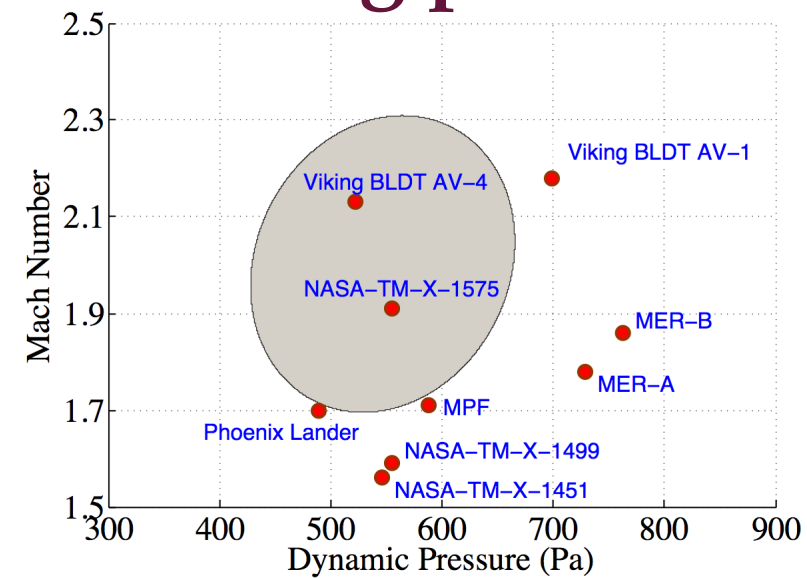
Prediction problem



Predict heating rate uncertainty

Control problem

Filtering problem

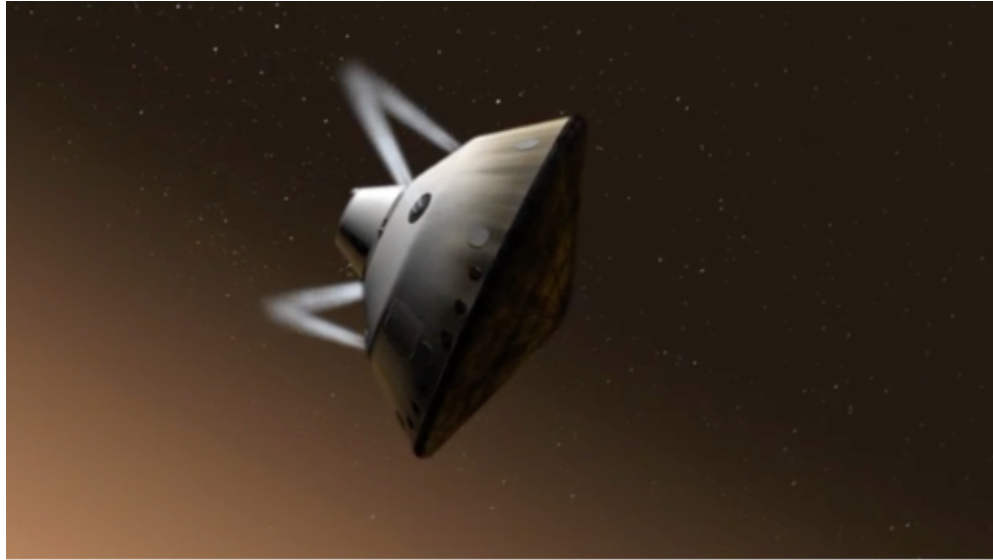


Estimate state to deploy parachute

Learning problem

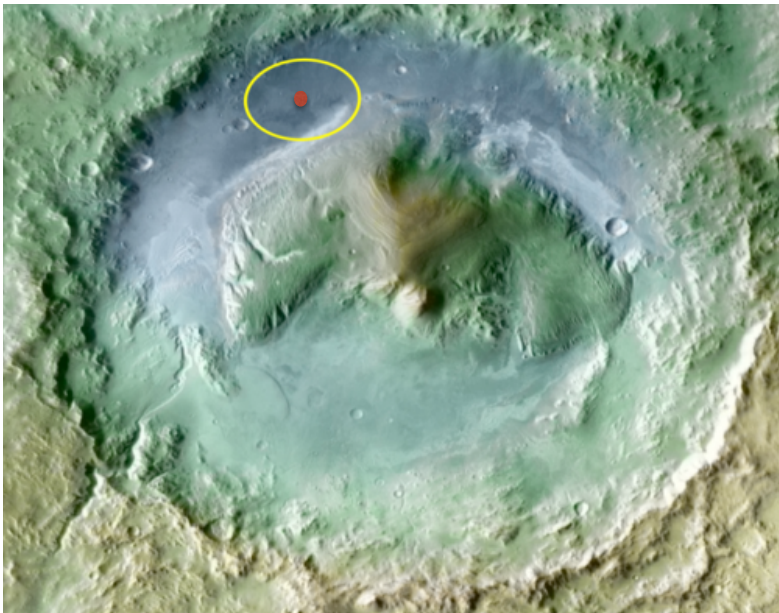
PDFs in Mars Entry-Descent-Landing

Prediction problem



Predict heating rate uncertainty

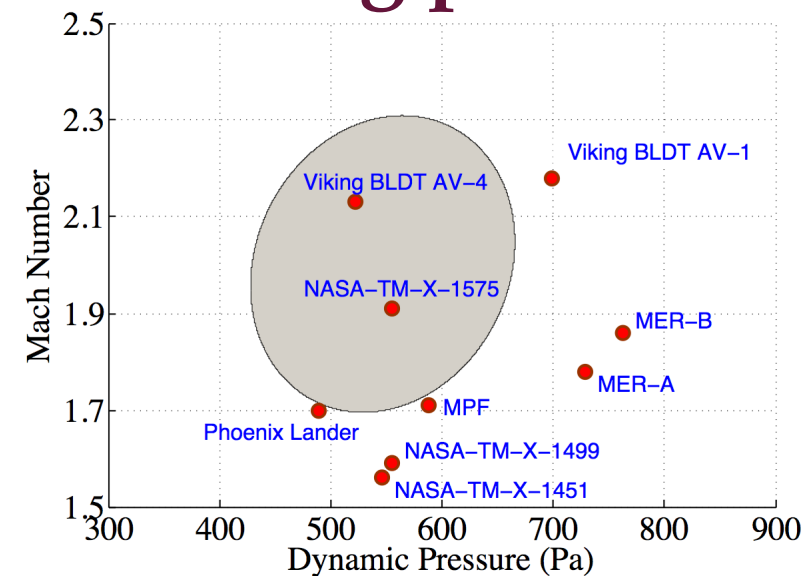
Control problem



Gale Crater (4.49S, 137.42E)

Steer state PDF to achieve
desired landing footprint accuracy

Filtering problem

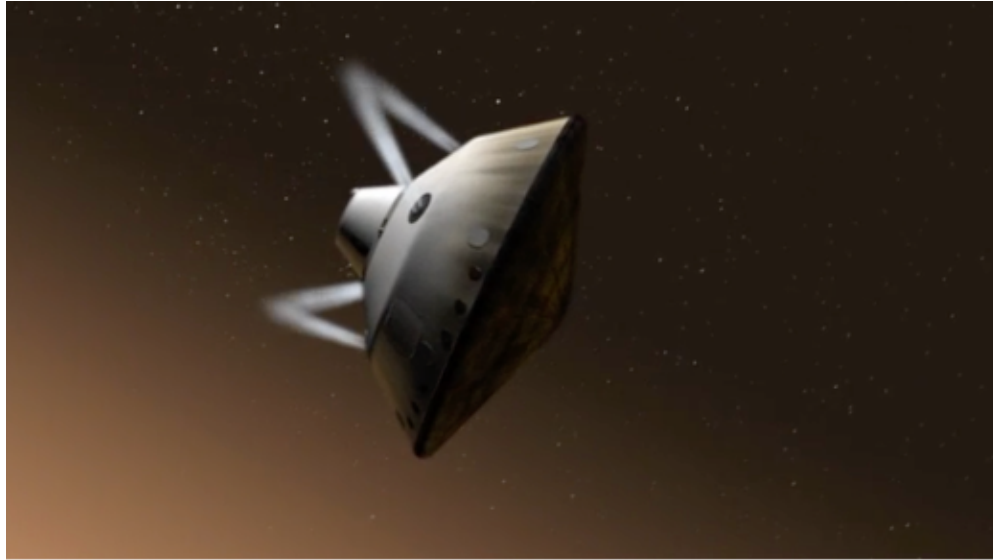


Estimate state to deploy parachute

Learning problem

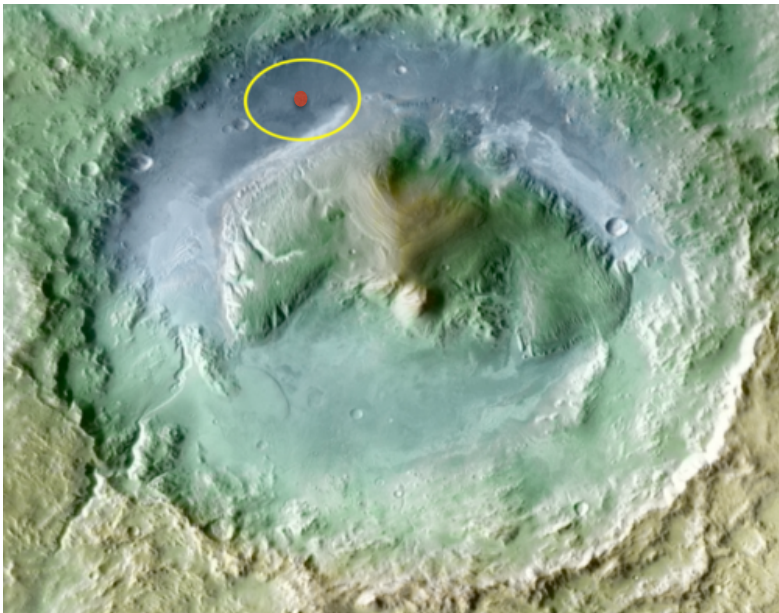
PDFs in Mars Entry-Descent-Landing

Prediction problem



Predict heating rate uncertainty

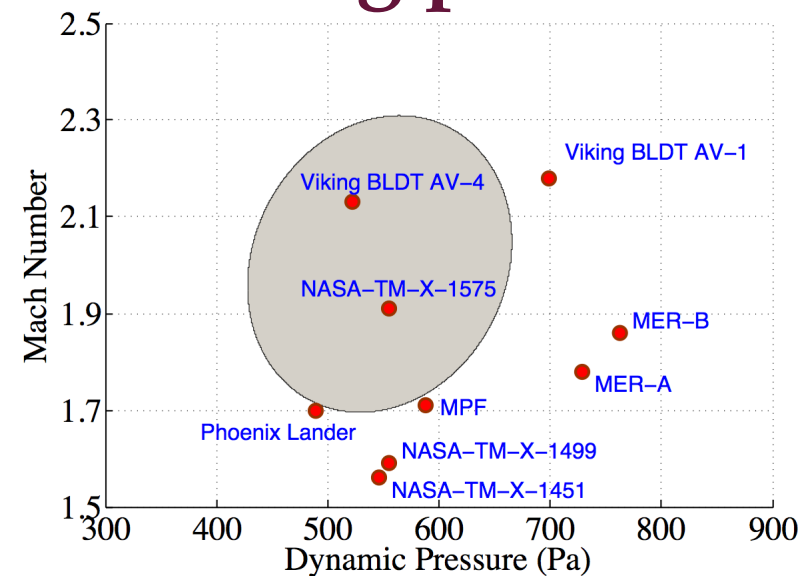
Control problem



Gale Crater (4.49S, 137.42E)

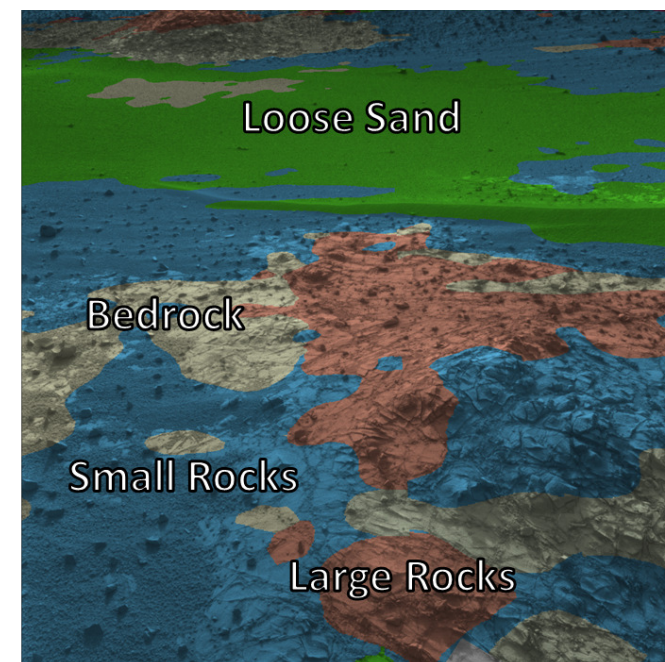
Steer state PDF to achieve desired landing footprint accuracy

Filtering problem



Estimate state to deploy parachute

Learning problem



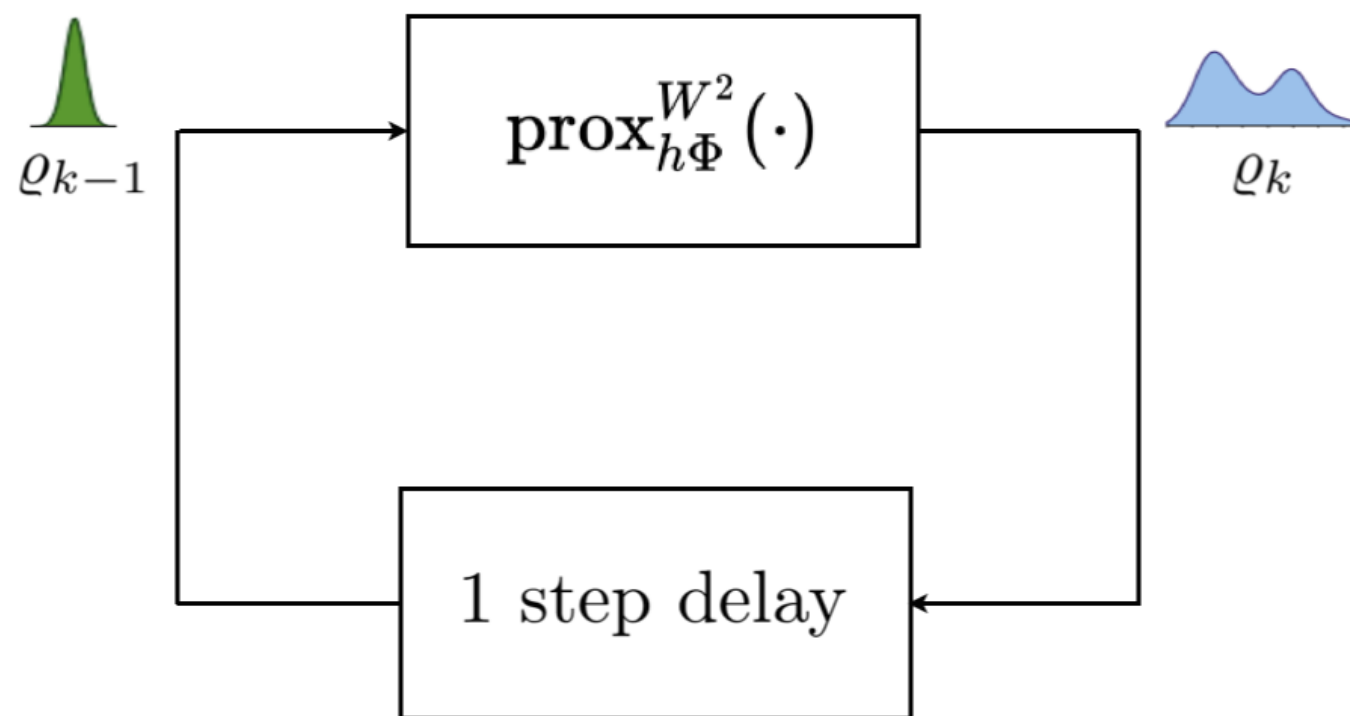
Learn surface feature from data

Solving prediction problem as Wasserstein 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: $\varrho_k = \text{prox}_{h\Phi}^{W^2}(\varrho_{k-1}) := \arg \inf_{\varrho \in \mathcal{P}_2(\mathcal{X})} \left\{ \frac{1}{2} W^2(\varrho, \varrho_{k-1}) + h\Phi(\varrho) \right\}$

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

Free energy functional: $\Phi(\varrho) := \int_{\mathcal{X}} \psi \varrho \, \mathrm{d}x + \beta^{-1} \int_{\mathcal{X}} \varrho \log \varrho \, \mathrm{d}x$

Geometric Meaning of Gradient Flow

Gradient Flow in \mathcal{X}

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

Recursion:

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

Convergence:

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

φ as Lyapunov function:

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

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:

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

Convergence:

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

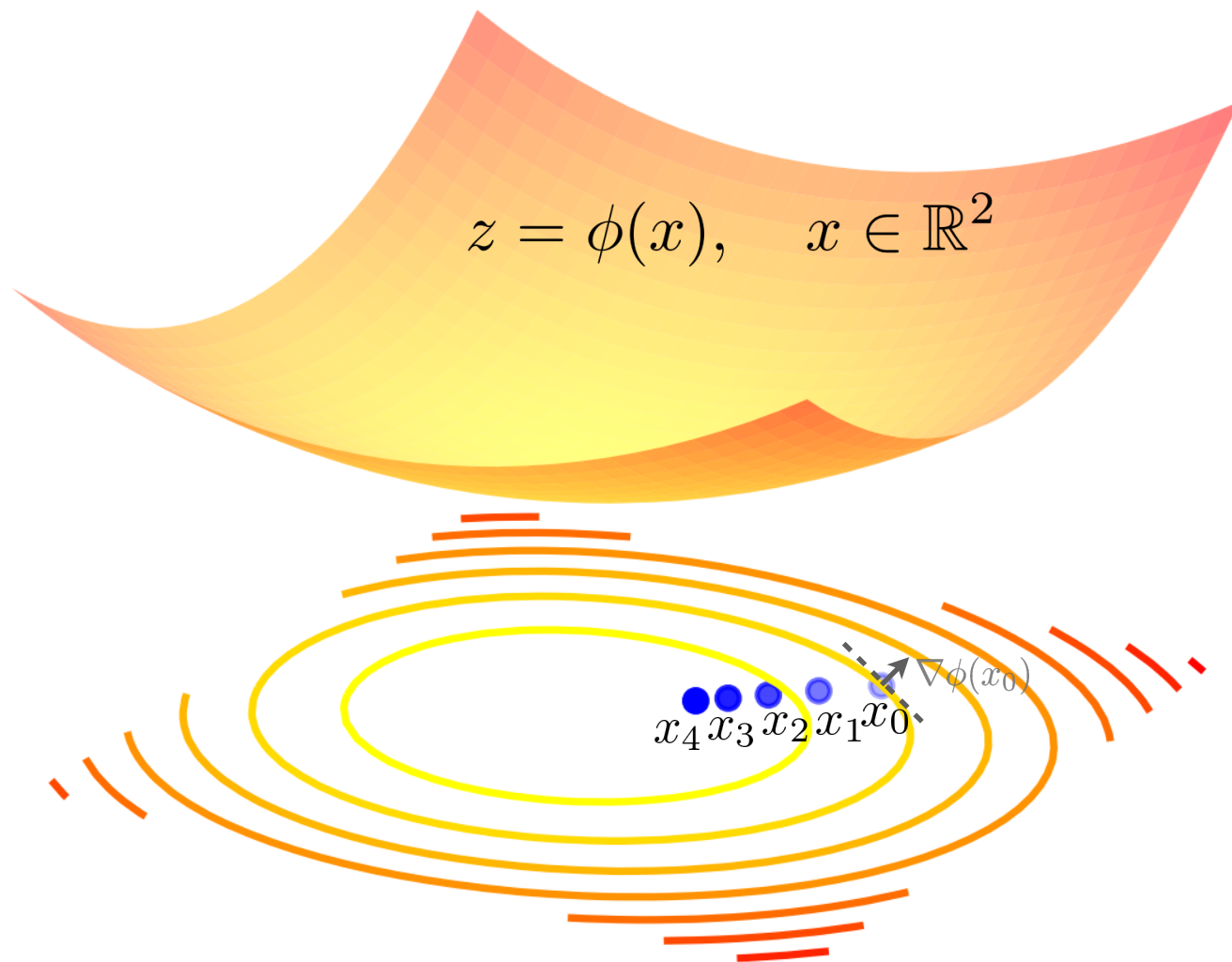
Φ as Lyapunov functional:

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

Geometric Meaning of Gradient Flow

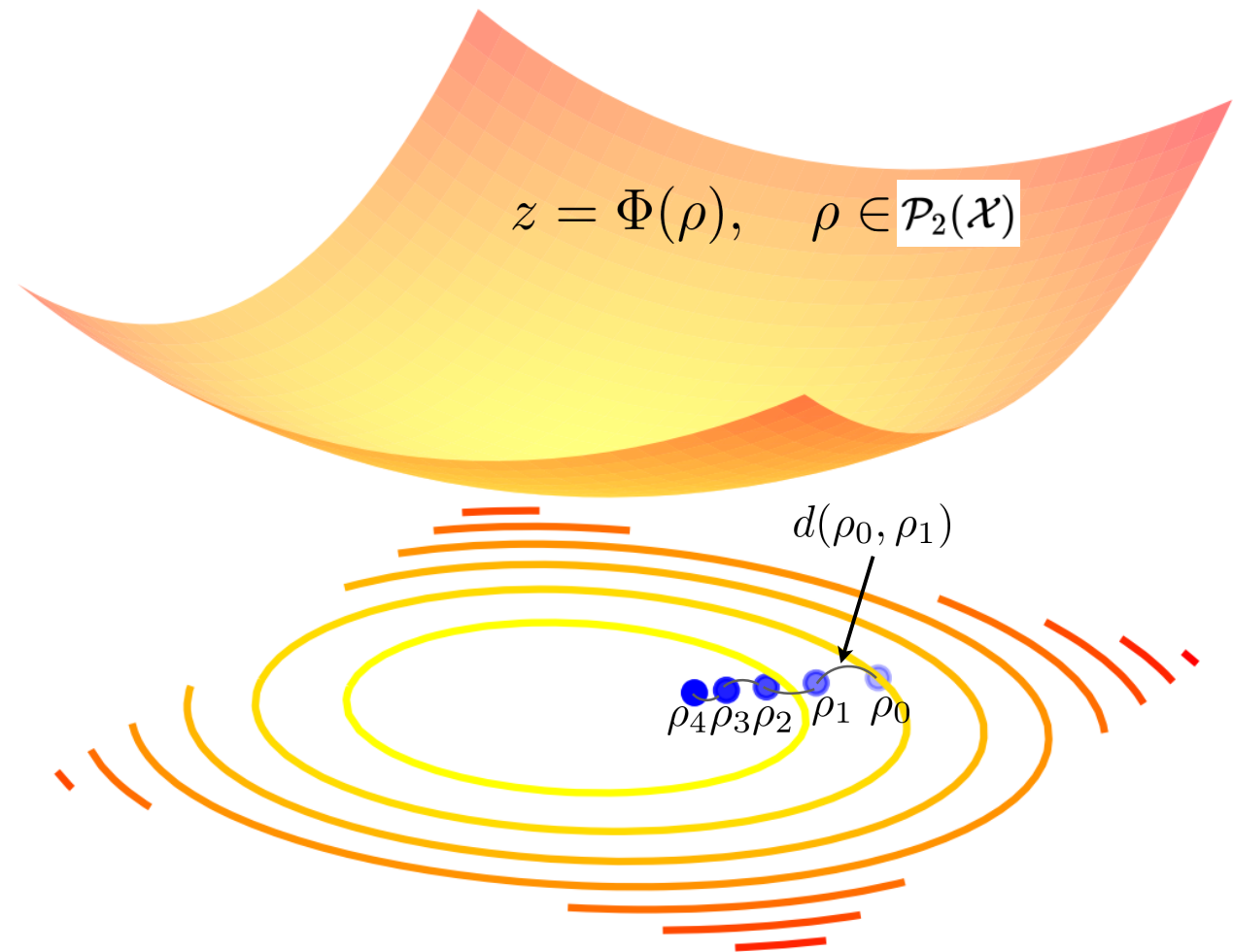
Gradient Flow in \mathcal{X}

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



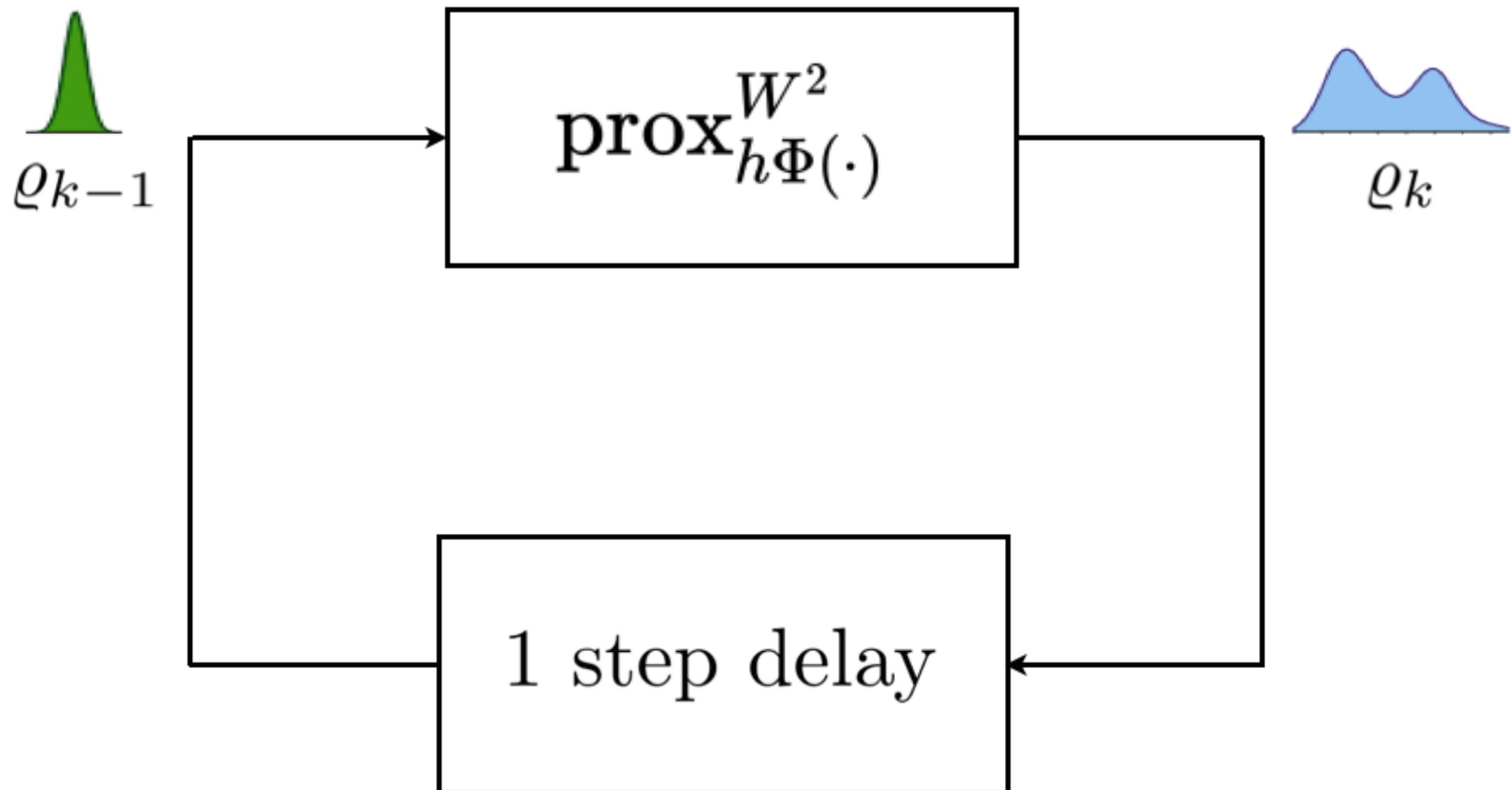
Gradient Flow in $\mathcal{P}_2(\mathcal{X})$

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



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

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

\Updownarrow

Proximal Recursion

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

\Downarrow

Discrete Primal Formulation

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

\Downarrow

Entropic Regularization

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

\Updownarrow

Dualization

$$\boldsymbol{\lambda}_0^{\text{opt}}, \boldsymbol{\lambda}_1^{\text{opt}} = \arg \max_{\boldsymbol{\lambda}_0, \boldsymbol{\lambda}_1 \geq 0} \left\{ \langle \boldsymbol{\lambda}_0, \boldsymbol{\varrho}_{k-1} \rangle - F^*(-\boldsymbol{\lambda}_1) - \frac{\epsilon}{h} \left(\exp(\boldsymbol{\lambda}_0^\top h / \epsilon) \exp(-\mathbf{C}_k / 2\epsilon) \exp(\boldsymbol{\lambda}_1 h / \epsilon) \right) \right\}$$

Recursion on the Cone

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

Coupled Transcendental Equations in \mathbf{y} and \mathbf{z}

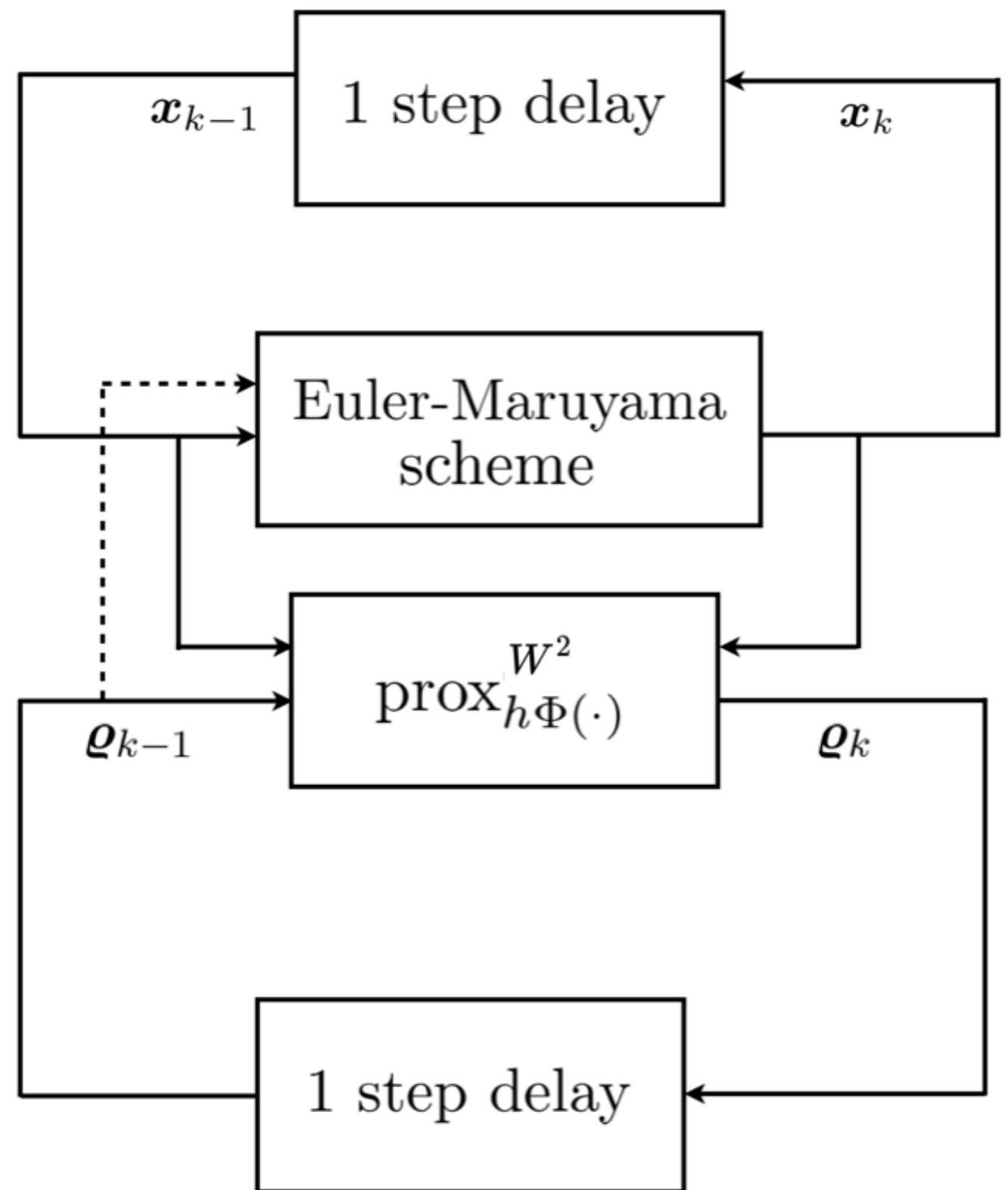
$$\begin{array}{l} \Gamma_k = e^{\frac{-\mathbf{c}_k}{2\epsilon}} \\ \varrho_{k-1} \\ \xi_{k-1} = \frac{e^{-\beta\psi_{k-1}}}{e} \end{array} \begin{array}{l} \longrightarrow \\ \longrightarrow \\ \longrightarrow \end{array} \boxed{\begin{array}{l} \mathbf{y} \odot \Gamma_k \mathbf{z} = \varrho_{k-1} \\ \mathbf{z} \odot \Gamma_k^\top \mathbf{y} = \xi_{k-1} \odot \mathbf{z}^{-\beta\epsilon/2h} \end{array}} \longrightarrow \varrho_k = \mathbf{z} \odot \Gamma_k^\top \mathbf{y}$$

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

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

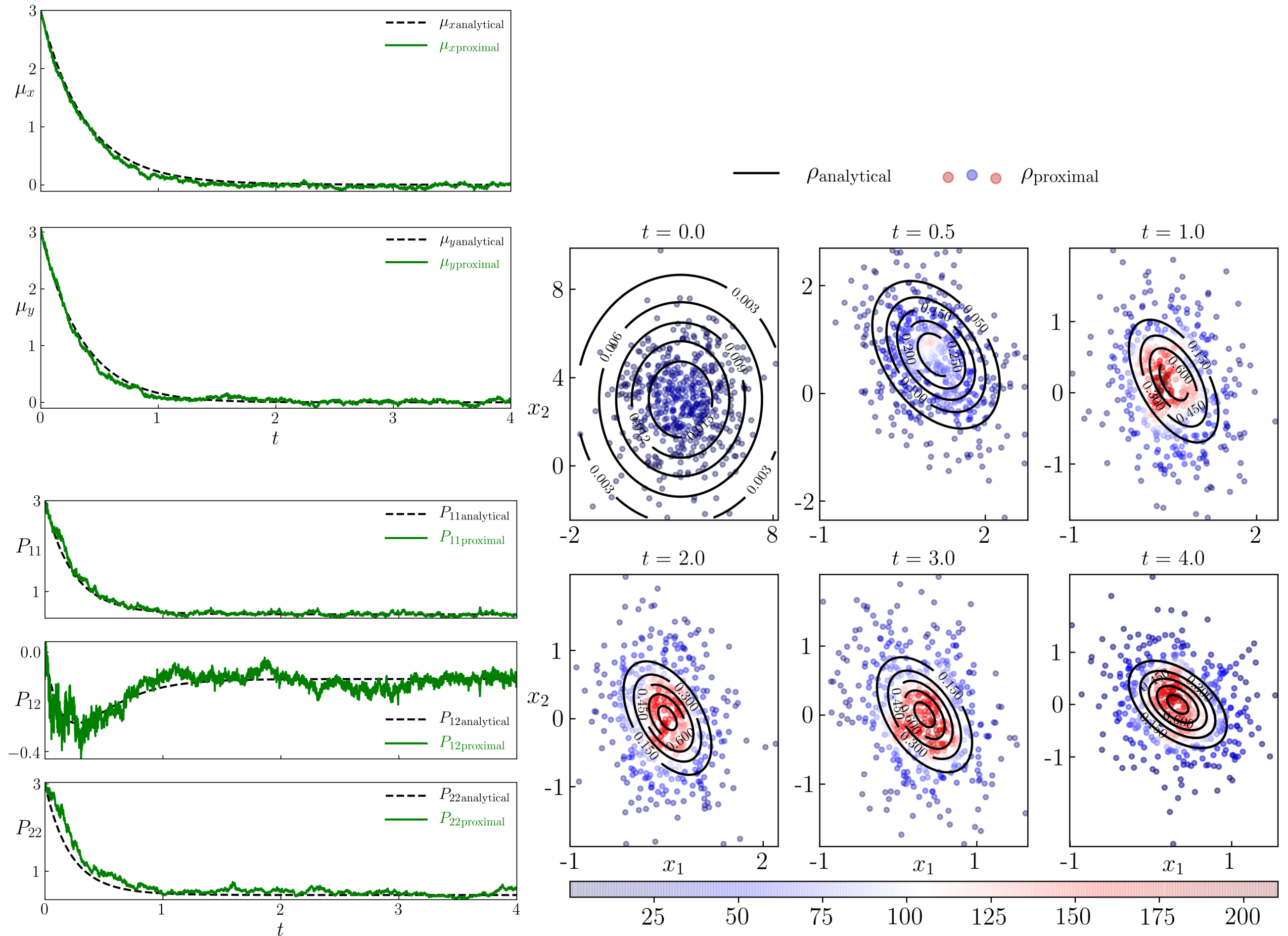
Then the solution $(\mathbf{y}^*, \mathbf{z}^*)$ gives the proximal update $\varrho_k = \mathbf{z}^* \odot (\Gamma_k^\top \mathbf{y}^*)$

Algorithmic Setup

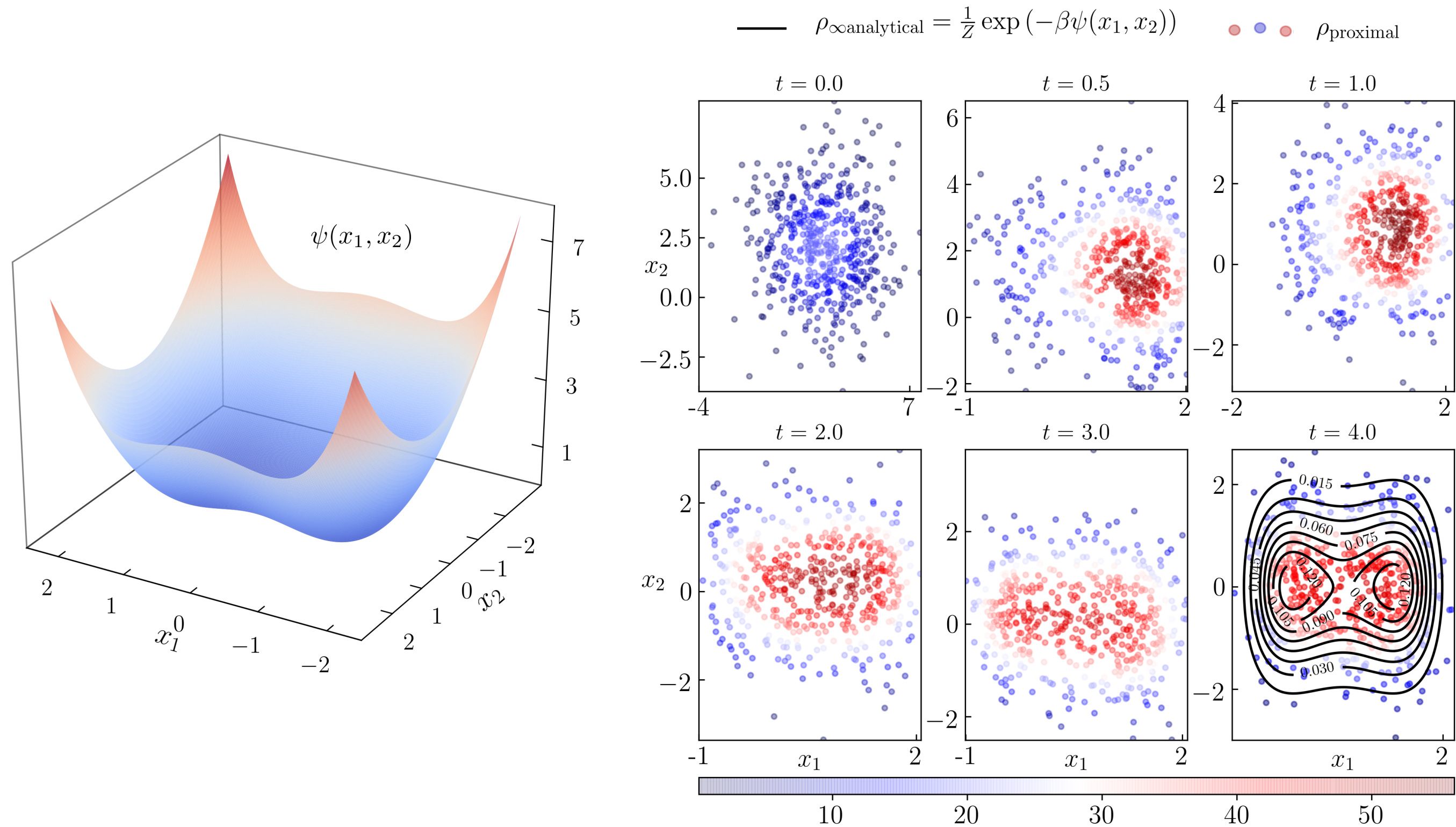


Theorem: Block co-ordinate iteration of (\mathbf{y}, \mathbf{z}) recursion is contractive on $\mathbb{R}_{>0}^n \times \mathbb{R}_{>0}^n$.

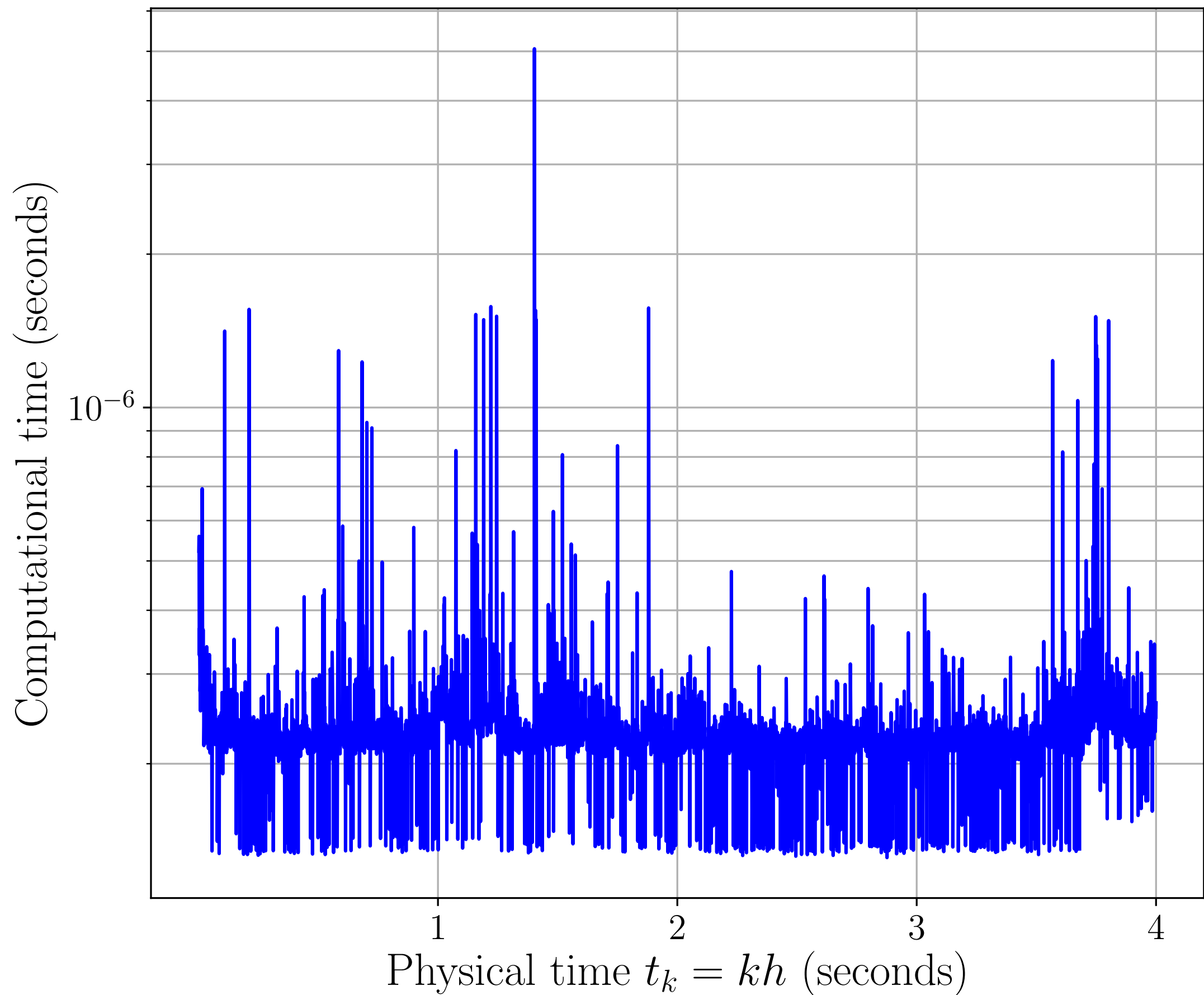
Proximal Prediction: 2D Linear Gaussian



Proximal Prediction: Nonlinear Non-Gaussian



Computational Time: Nonlinear Non-Gaussian



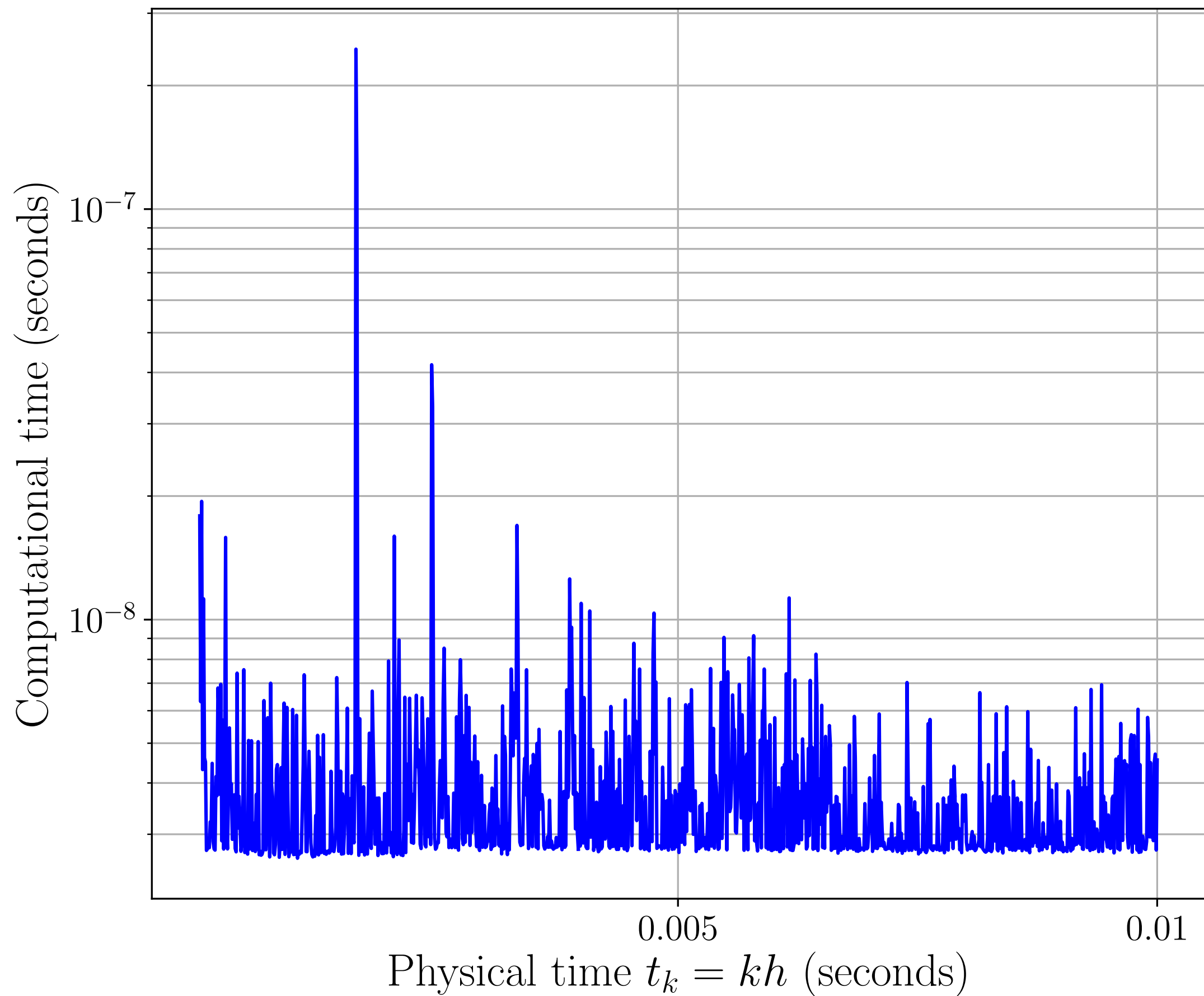
Proximal Prediction: Satellite in Geocentric Orbit

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

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

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

Computational Time: Satellite in Geocentric Orbit



Extensions: Nonlocal Interactions

PDF dependent sample path dynamics:

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

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

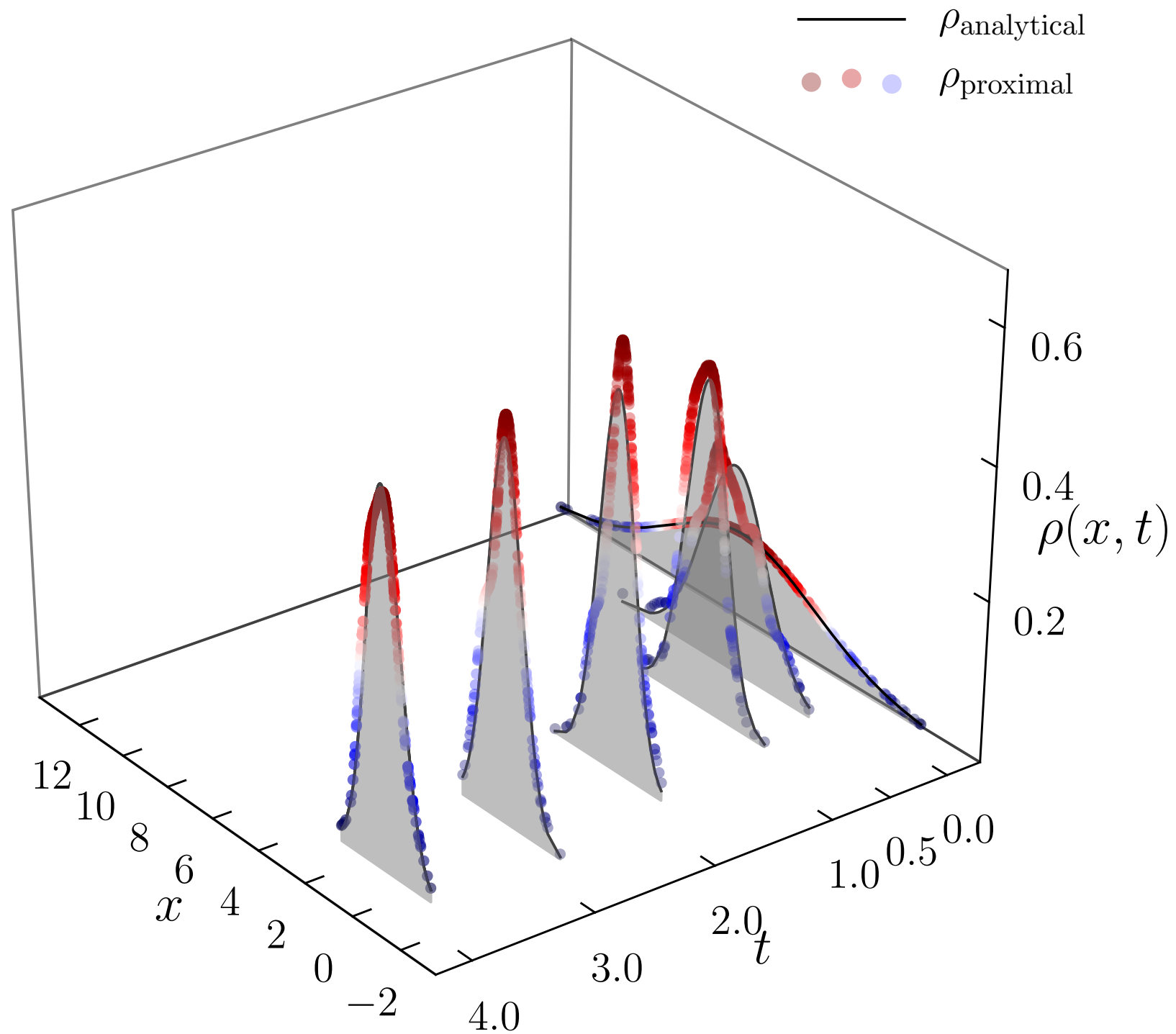
$$\frac{\partial \rho}{\partial t} = \nabla \cdot (\rho \nabla (U + \rho * V)) + \beta^{-1} \Delta \rho$$

Free energy:

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

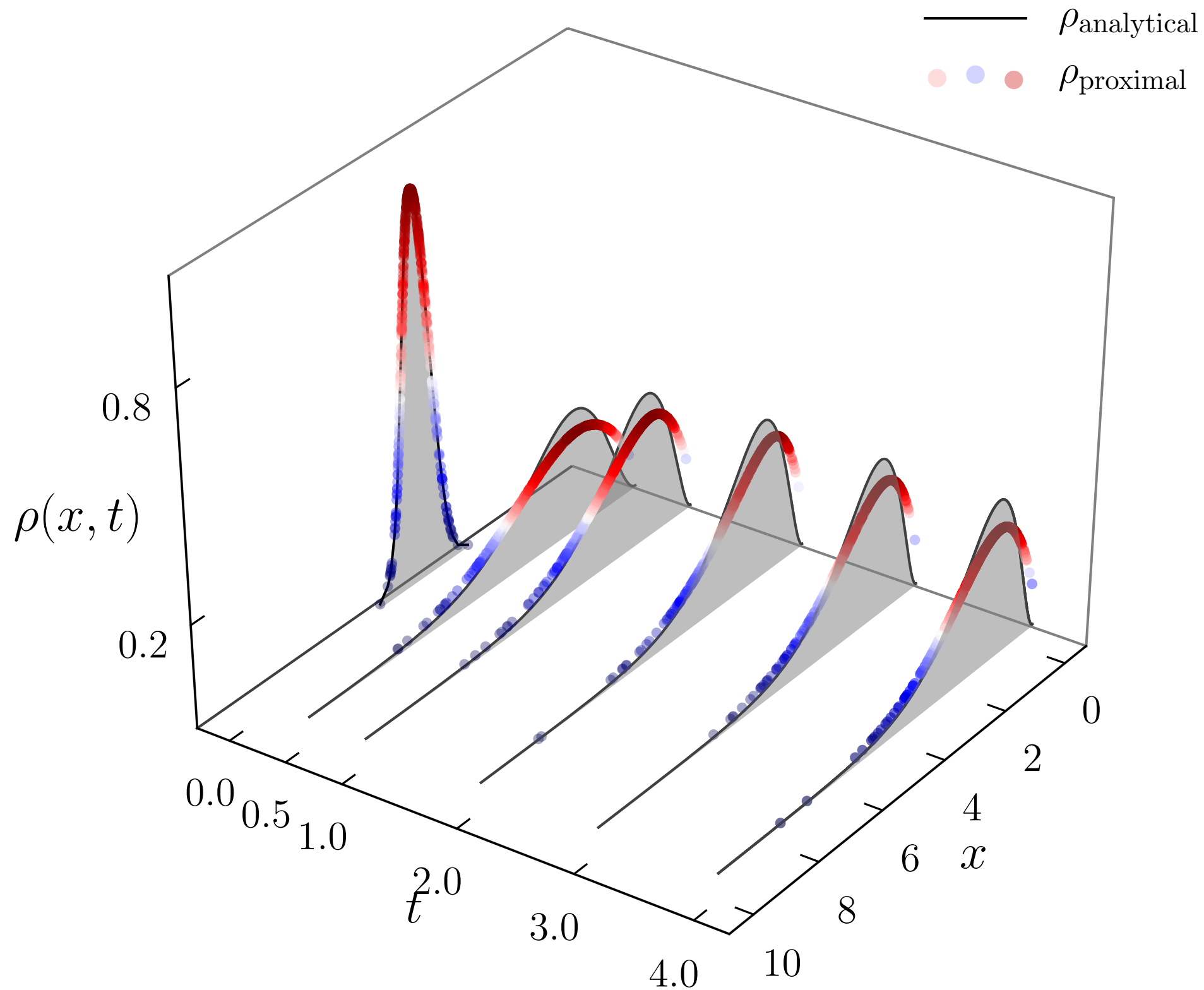
Extensions: Nonlocal Interactions

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



Extensions: Multiplicative Noise

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



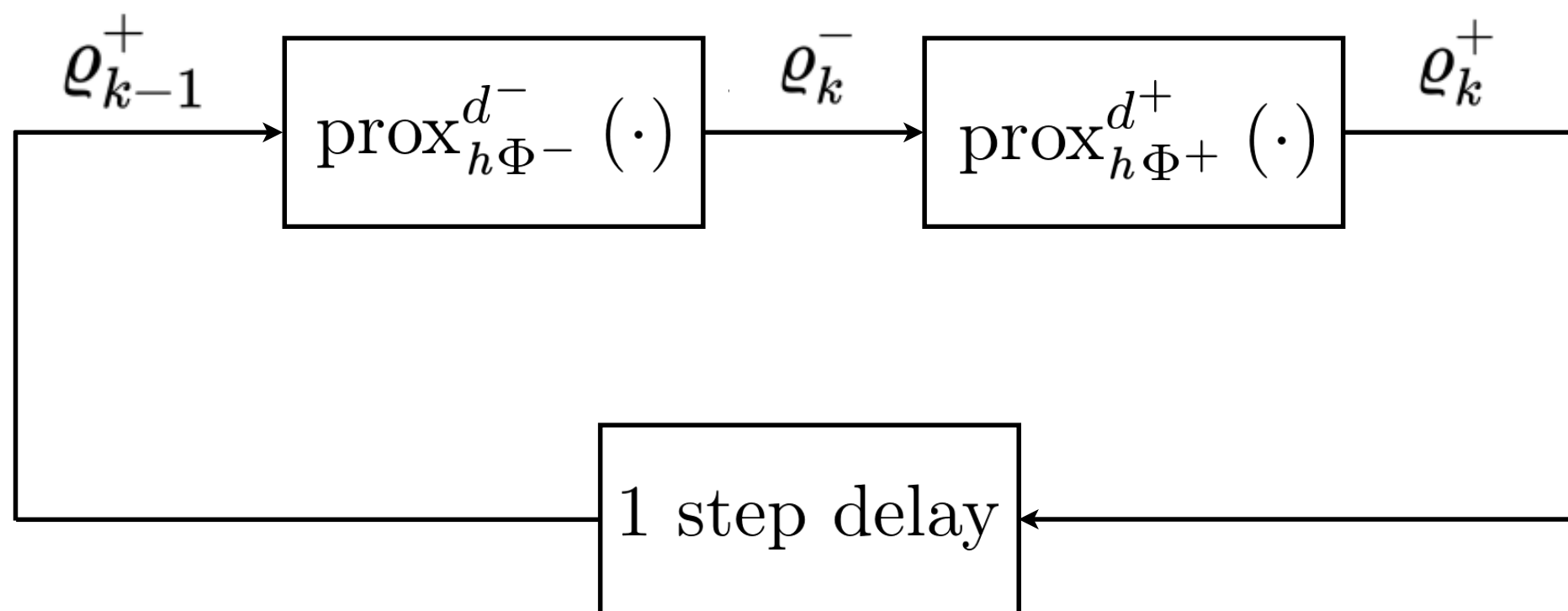
Solving filtering as Wasserstein gradient flow

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 ◦ stochastic} proximal operators:



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

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

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

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:

$$d\mu^+(t) = \mathbf{A}\mu^+(t)dt + \overset{\mathbf{P}^+\mathbf{C}\mathbf{R}^{-1}}{\underset{\text{I}}{\mathbf{K}(t)}} (d\mathbf{z}(t) - \mathbf{C}\mu^+(t)dt),$$

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

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

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

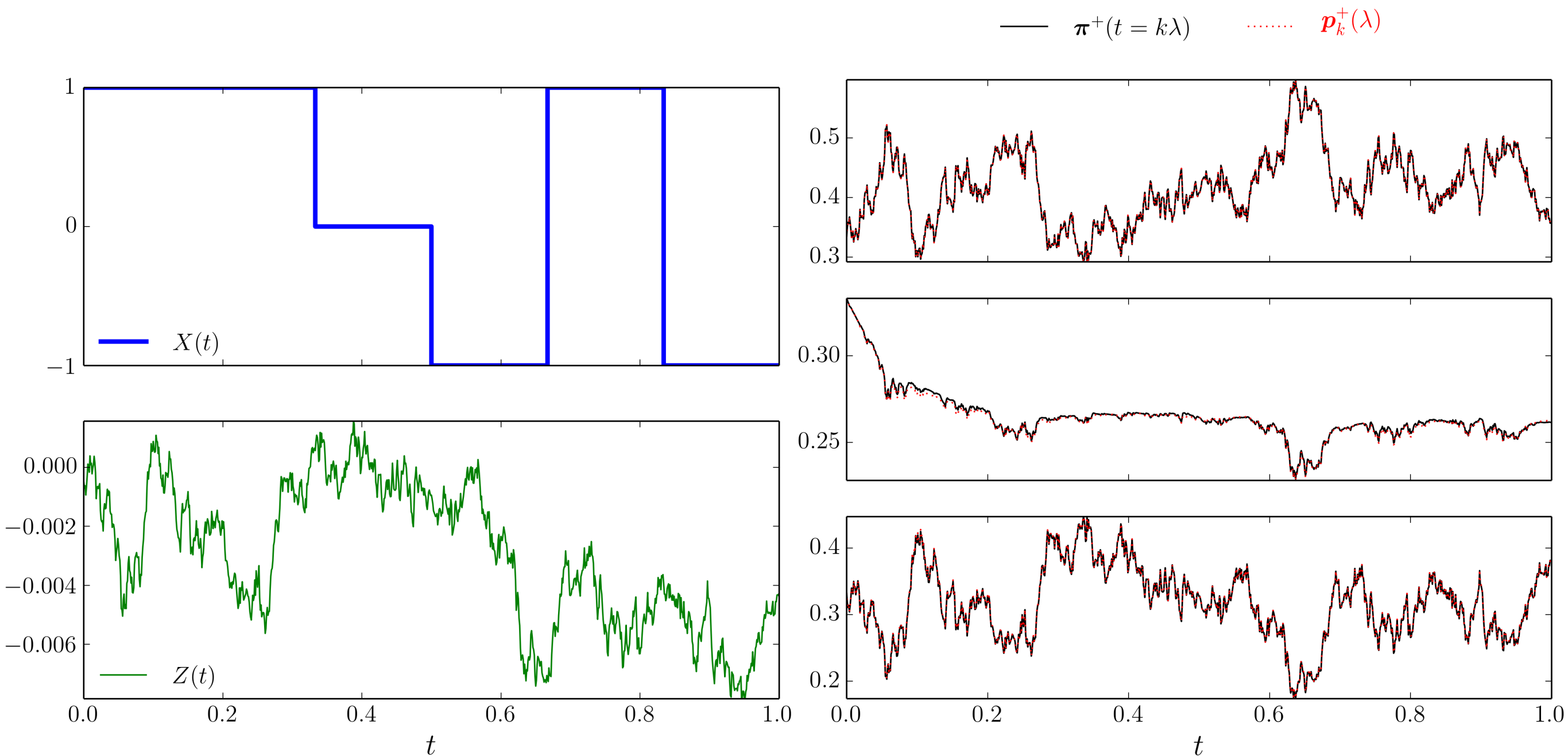
By defn. $\pi^+(t) = \mathbb{P}(x(t) = a_i \mid z(s), 0 \leq s \leq t)$

J.SIAM CONTROL
Ser. A, Vol. 2, No. 3
Printed in U.S.A., 1965

SOME APPLICATIONS OF STOCHASTIC DIFFERENTIAL
EQUATIONS TO OPTIMAL NONLINEAR FILTERING*

W. M. WONHAM†

Numerical Results for the Wonham Filter



Solving density control as Wasserstein gradient flow

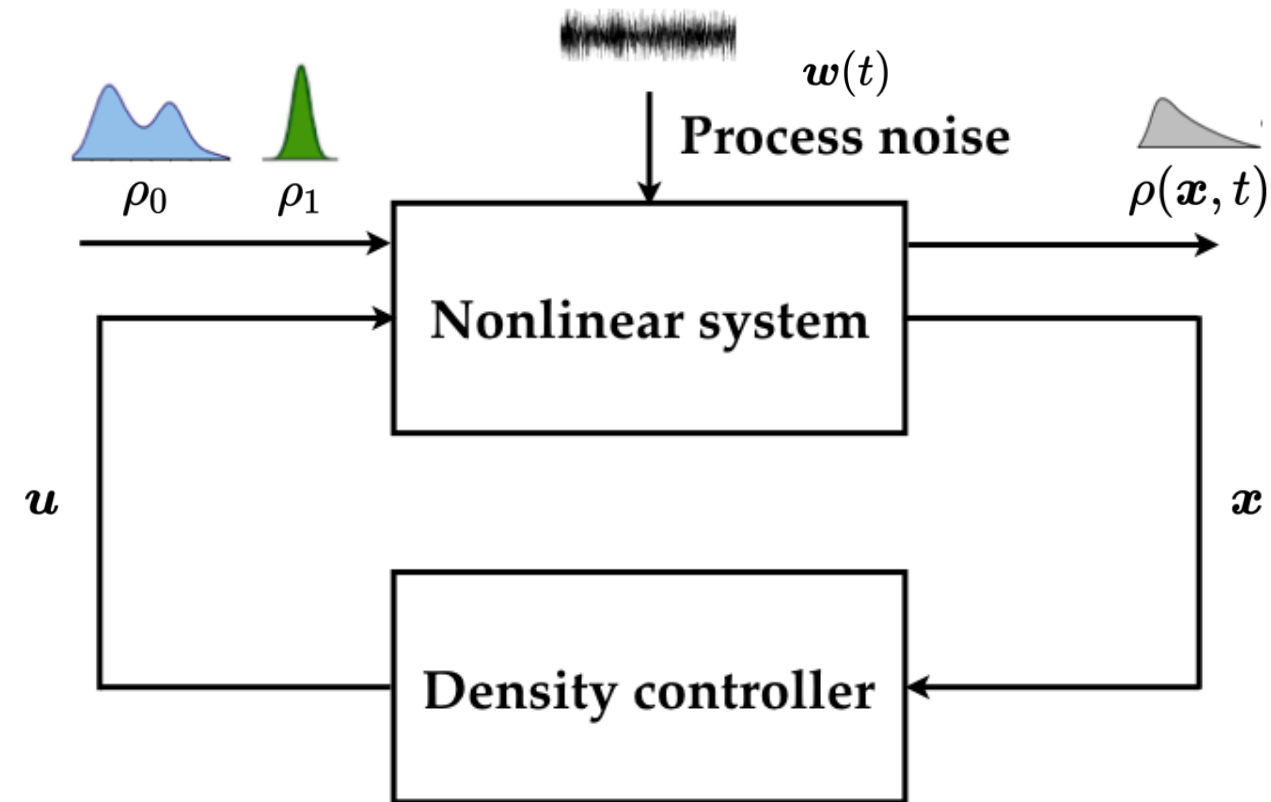
Finite Horizon Feedback Density Control

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

subject to

$$dx = \left\{ f(x, t) + B(t)u(x, t) \right\} dt + \sqrt{2\epsilon} B(t) dw,$$

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



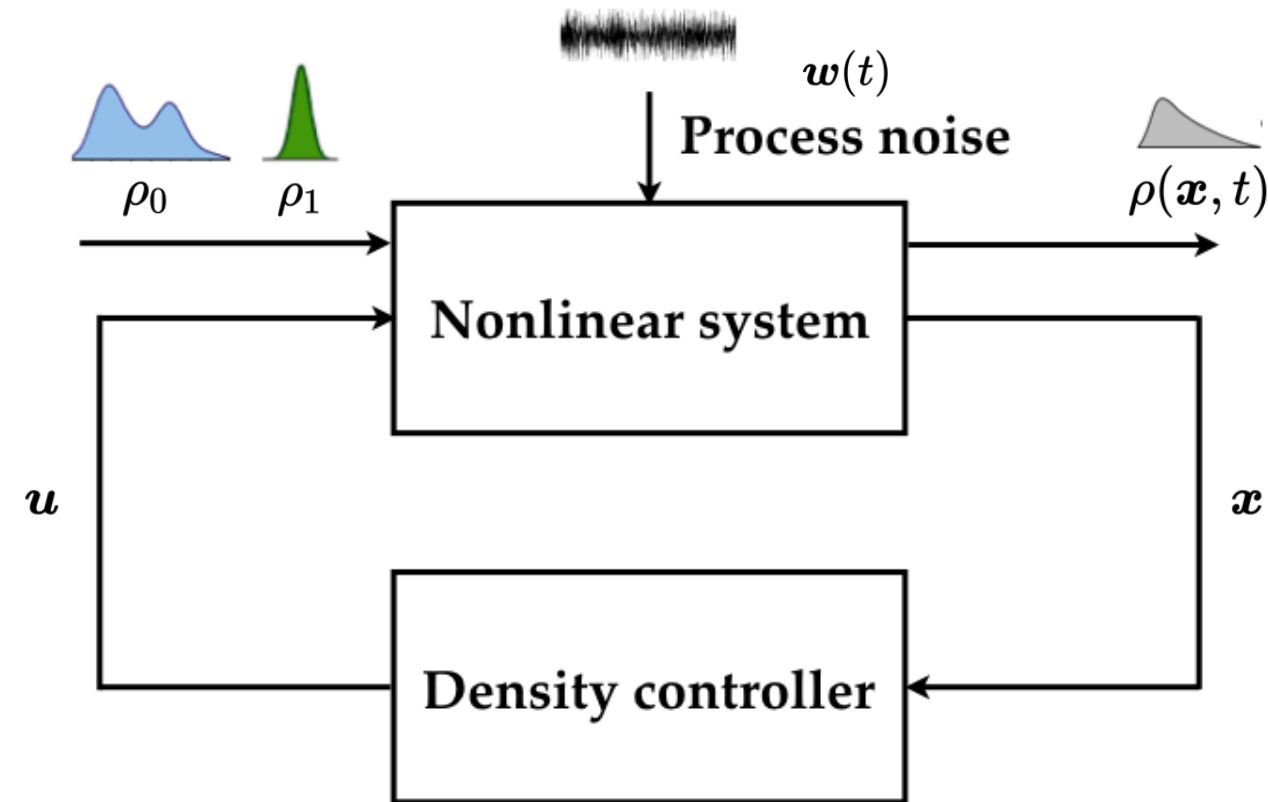
Finite Horizon Feedback Density Control

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

subject to

$$dx = \left\{ f(x, t) + B(t)u(x, t) \right\} dt + \sqrt{2\epsilon} B(t) dw,$$

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



Necessary conditions for optimality: coupled nonlinear PDEs (FPK + HJB)

$$\frac{\partial \rho^{\text{opt}}}{\partial t} + \nabla \cdot \left(\rho^{\text{opt}} \left(f + B(t)^\top \nabla \psi \right) \right) = \epsilon \mathbf{1}^\top \left(D(t) \odot \text{Hess}(\rho^{\text{opt}}) \right) \mathbf{1},$$

$$\frac{\partial \psi}{\partial t} + \frac{1}{2} \|B(t)^\top \nabla \psi\|_2^2 + \langle \nabla \psi, f \rangle = -\epsilon \langle D(t), \text{Hess}(\psi) \rangle$$

Boundary conditions:

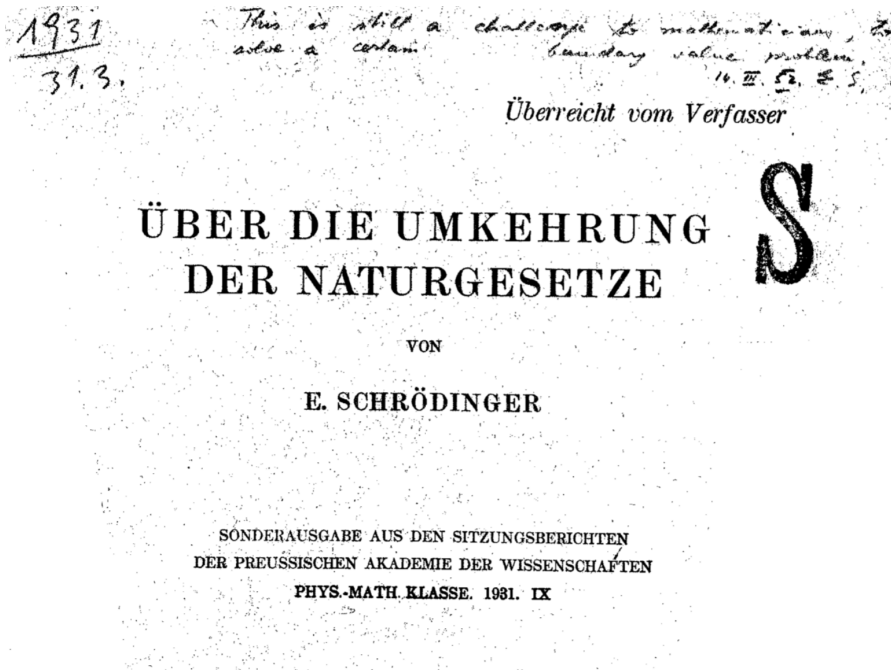
$$\rho^{\text{opt}}(x, 0) = \rho_0(x), \quad \rho^{\text{opt}}(x, 1) = \rho_1(x)$$

Optimal control:

$$u^{\text{opt}}(x, t) = B(t)^\top \nabla \psi$$

Feedback Synthesis via the Schrödinger System

Schrödinger's (until recently) forgotten papers:



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

PAF

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 transform: $(\rho^{\text{opt}}, \psi) \mapsto (\varphi, \hat{\varphi})$

$$\begin{aligned}\varphi(\boldsymbol{x}, t) &= \exp\left(\frac{\psi(\boldsymbol{x}, t)}{2\epsilon}\right), \\ \hat{\varphi}(\boldsymbol{x}, t) &= \rho^{\text{opt}}(\boldsymbol{x}, t) \exp\left(-\frac{\psi(\boldsymbol{x}, t)}{2\epsilon}\right),\end{aligned}$$

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

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

Feedback Synthesis via the Schrödinger System

2 coupled nonlinear PDEs \rightarrow boundary-coupled linear PDEs!!

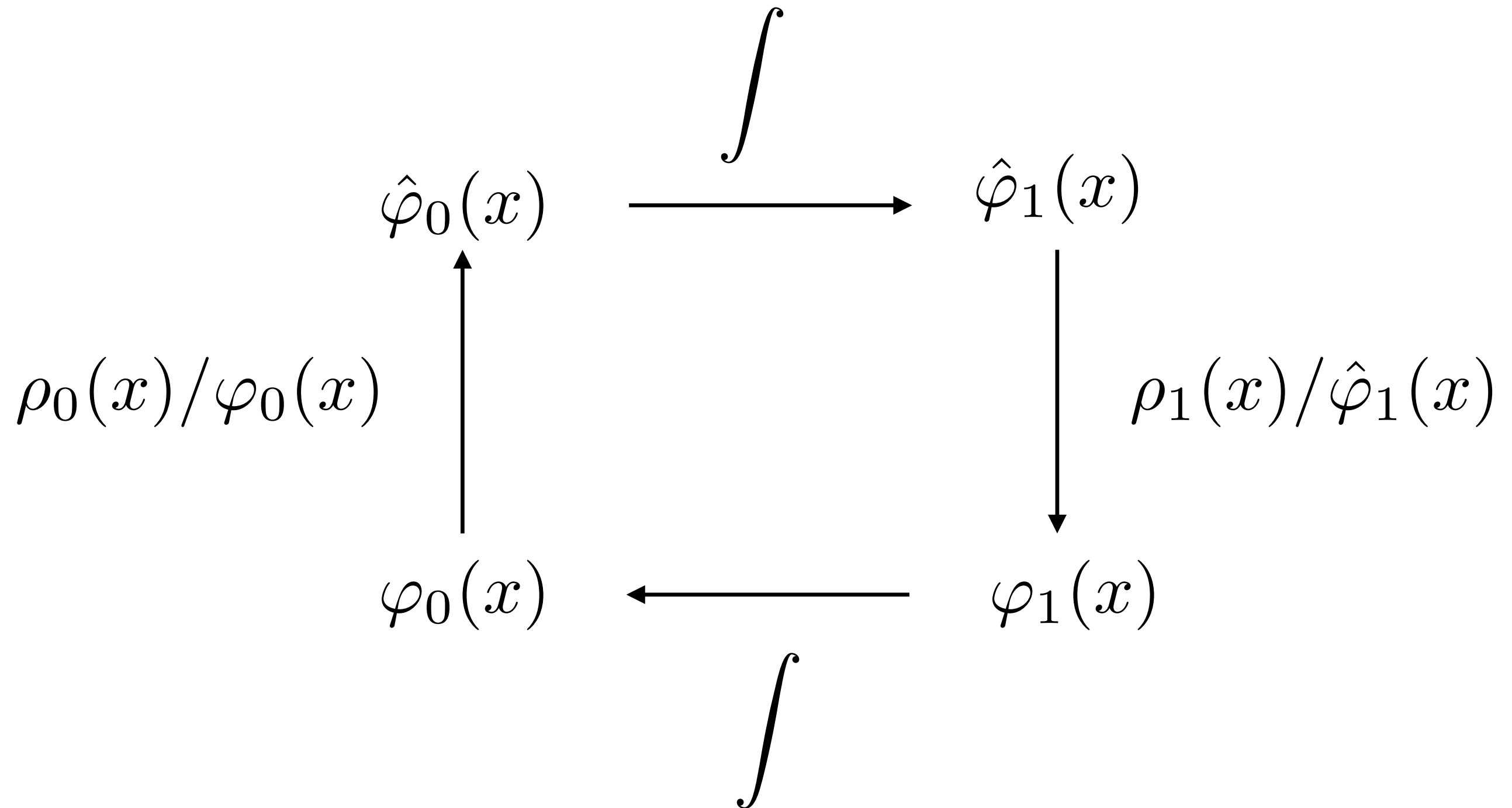
$$\underbrace{\frac{\partial \hat{\phi}}{\partial t} = -\nabla \cdot (\hat{\phi} \mathbf{f}) + \epsilon \mathbf{1}^\top (\mathbf{D}(t) \odot \text{Hess}(\hat{\phi})) \mathbf{1}}_{\text{forward Kolmogorov PDE}}, \quad \varphi_0 \hat{\phi}_0 = \rho_0,$$

$$\underbrace{\frac{\partial \varphi}{\partial t} = -\langle \nabla \varphi, \mathbf{f} \rangle - \epsilon \langle \mathbf{D}(t), \text{Hess}(\varphi) \rangle}_{\text{backward Kolmogorov PDE}}, \quad \varphi_1 \hat{\phi}_1 = \rho_1.$$

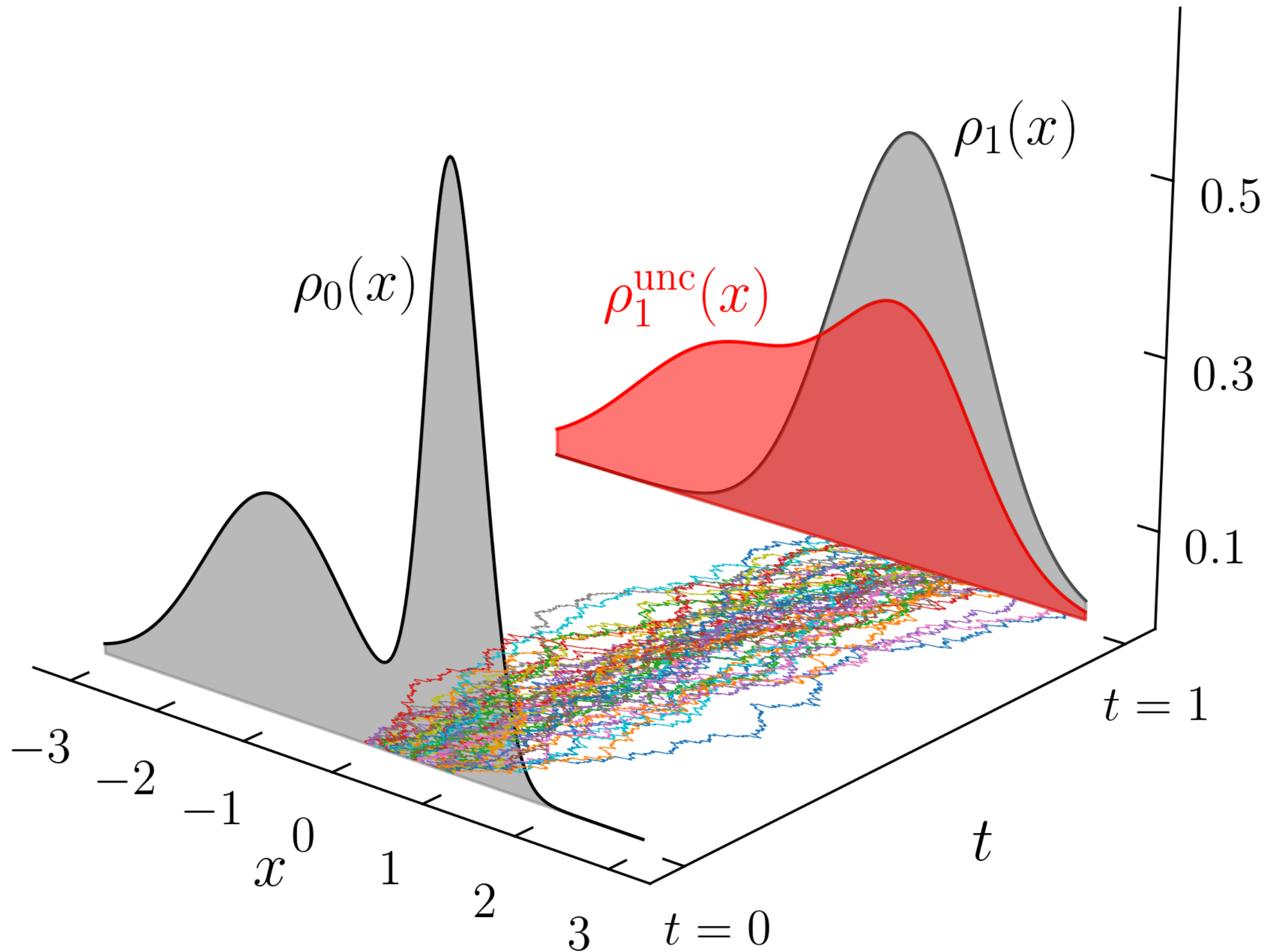
Wasserstein proximal algorithm \rightarrow fixed point recursion over $(\hat{\phi}_0, \varphi_1)$

(Contractive in Hilbert metric)

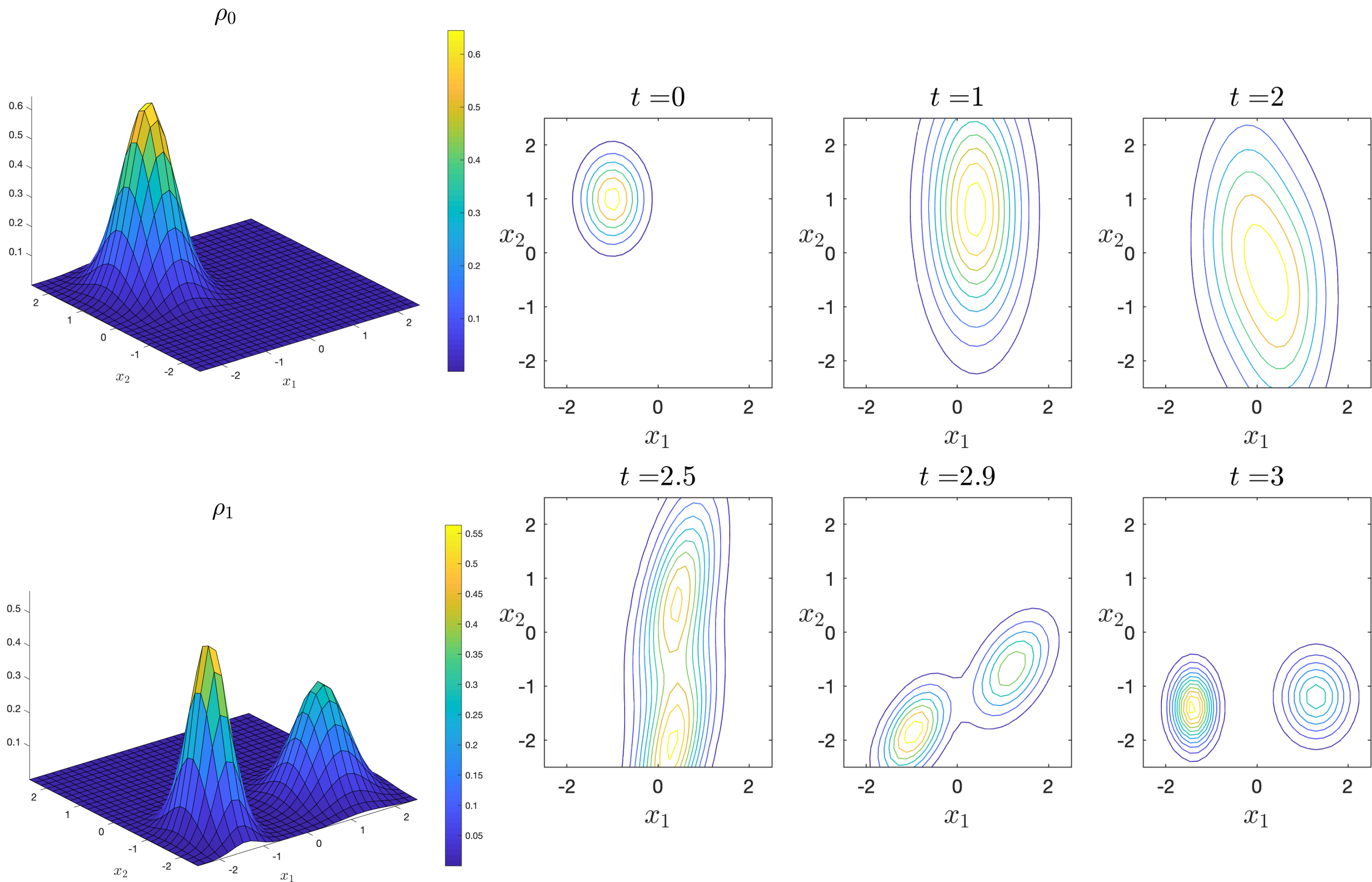
Fixed Point Recursion over $(\hat{\varphi}_0, \varphi_1)$



Feedback Density Control: Zero Prior Dynamics

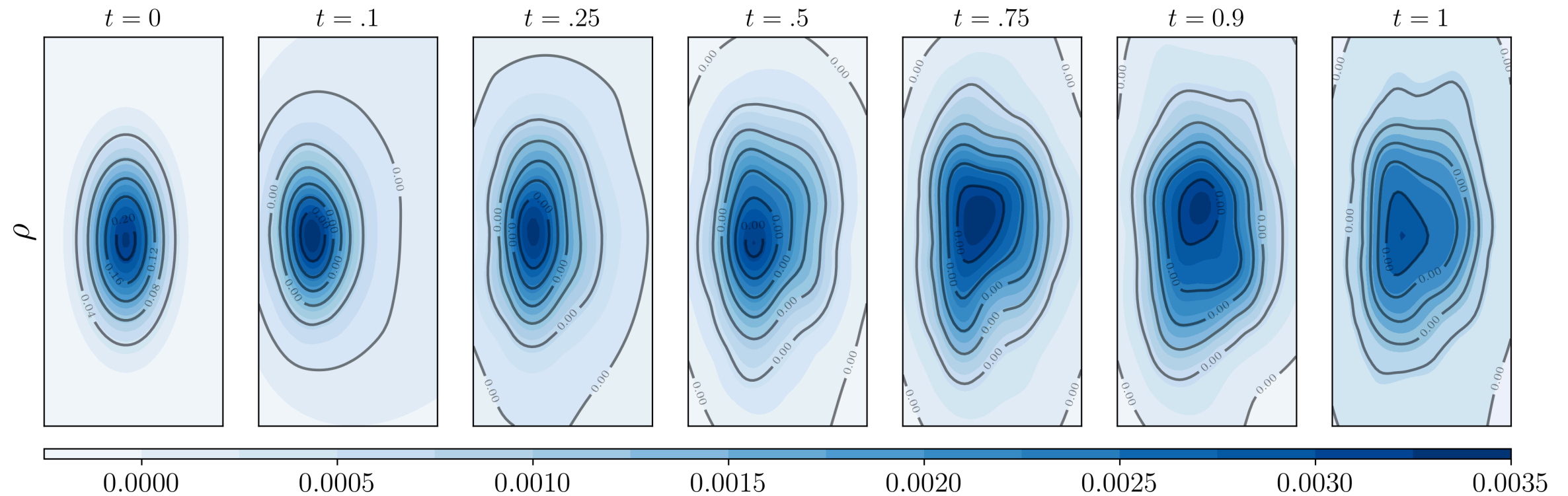


Feedback Density Control: LTI Prior Dynamics

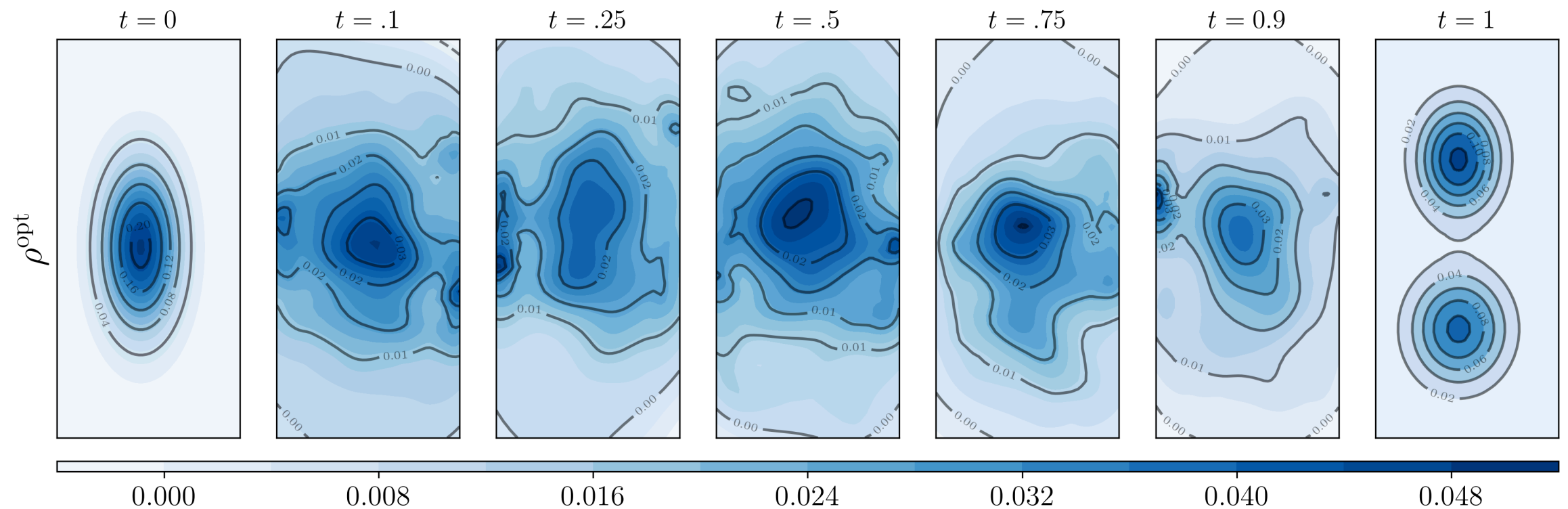


Feedback Density Control: Nonlinear Grad. Drift

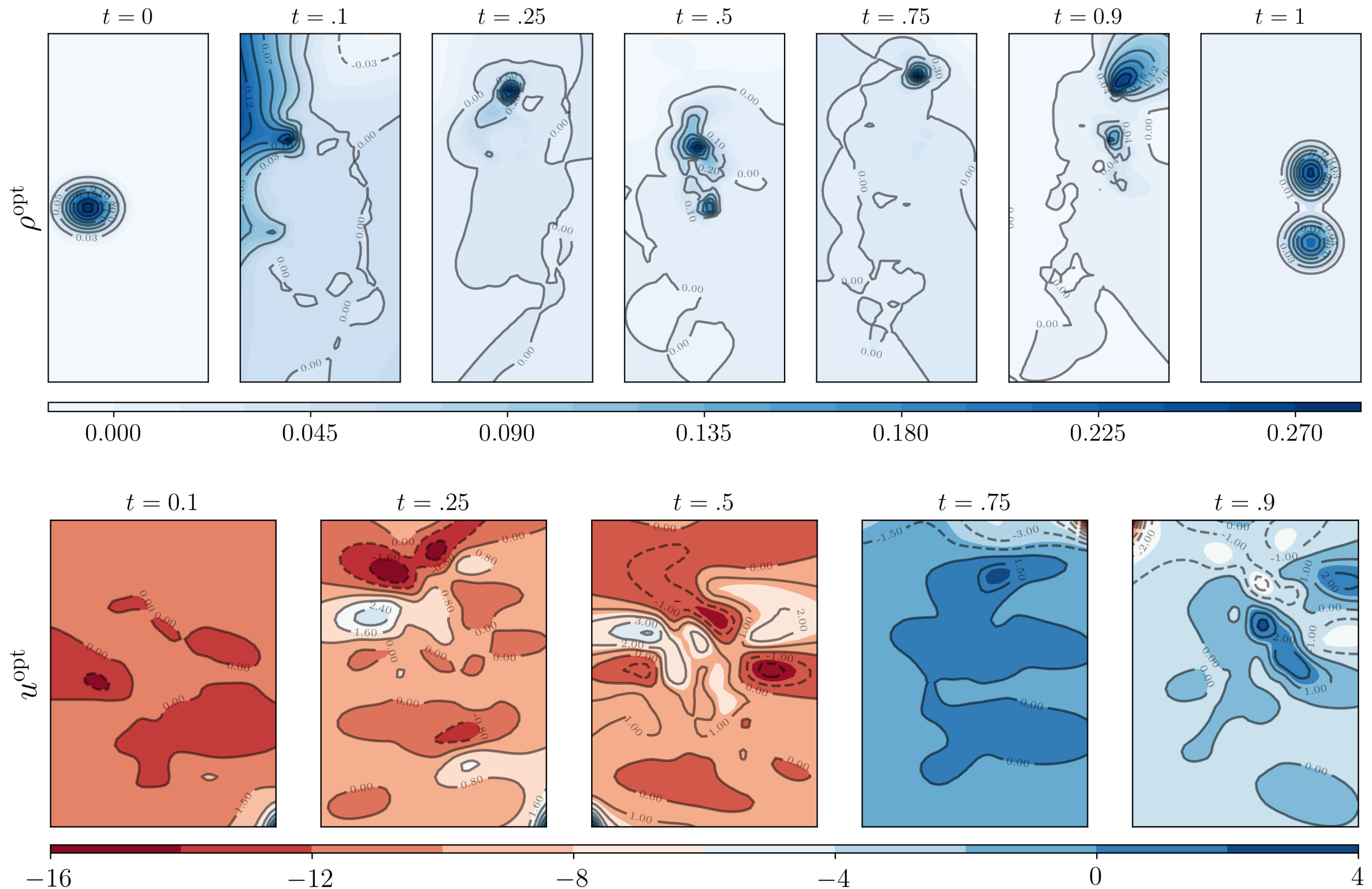
Uncontrolled joint PDF evolution:



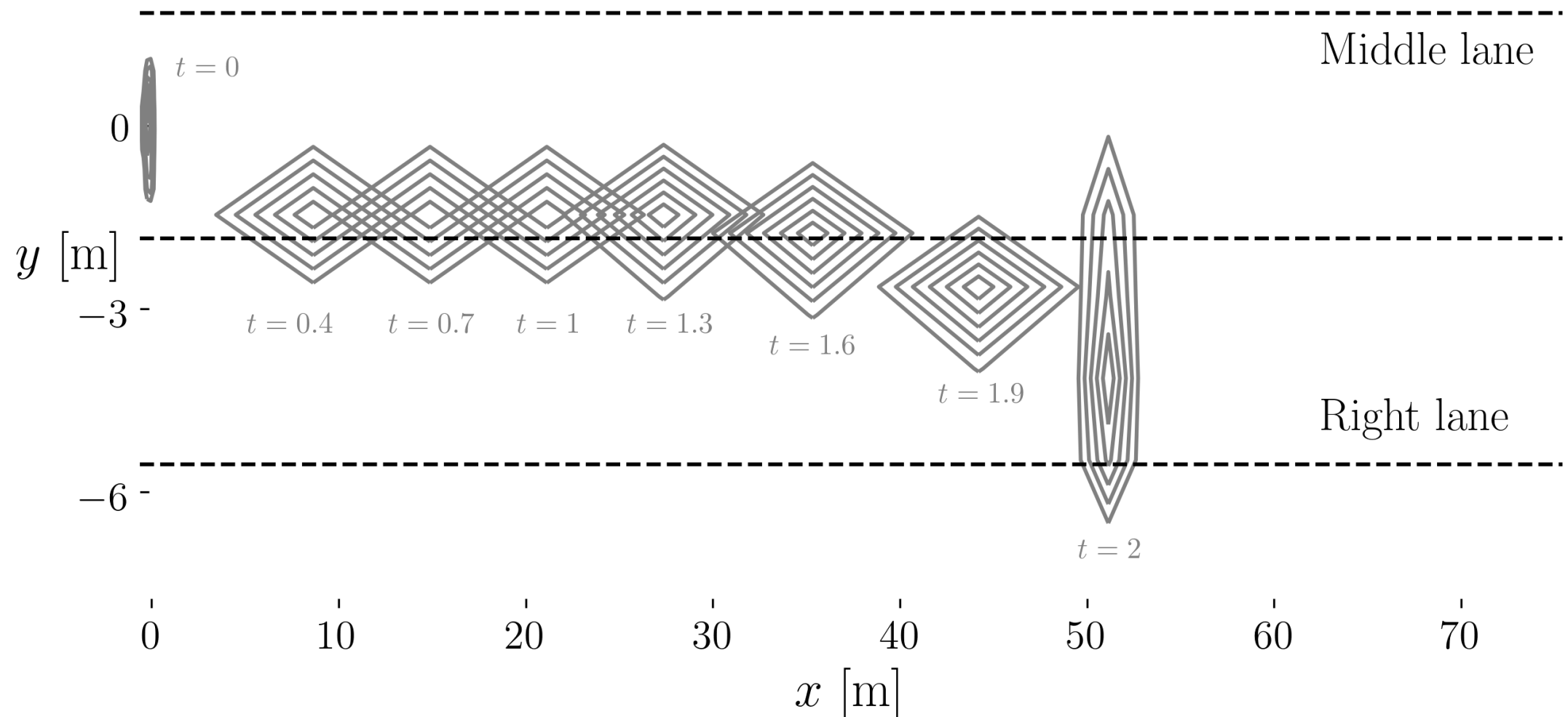
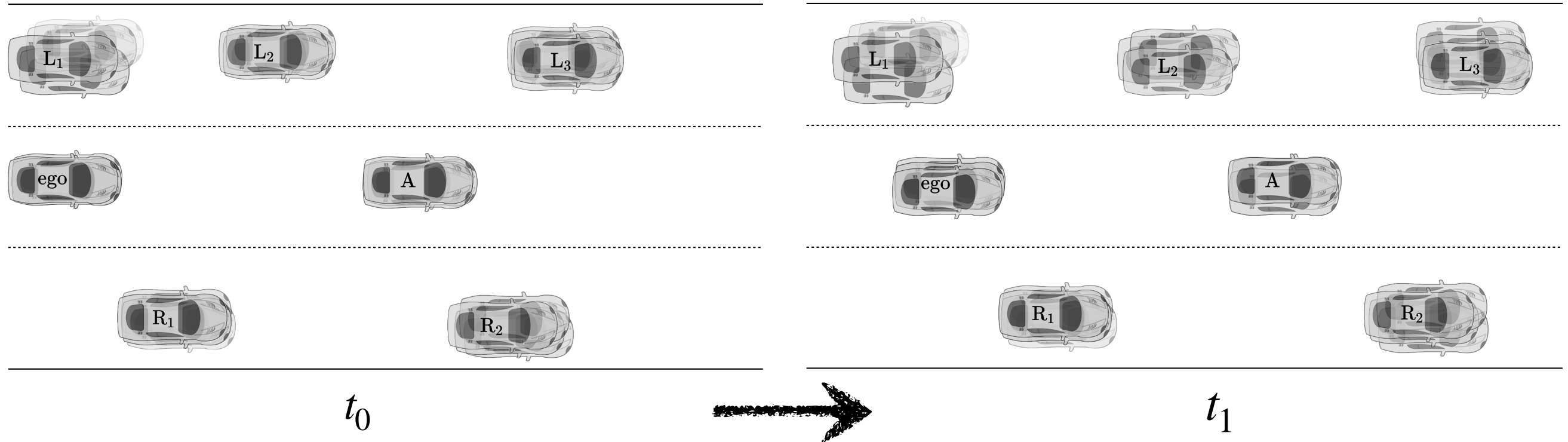
Optimal controlled joint PDF evolution:



Feedback Density Control: Mixed Conservative-Dissipative Drift



Density Control for Safe Automated Driving



Learning a neural network as Wasserstein 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:

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

Population risk functional:

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

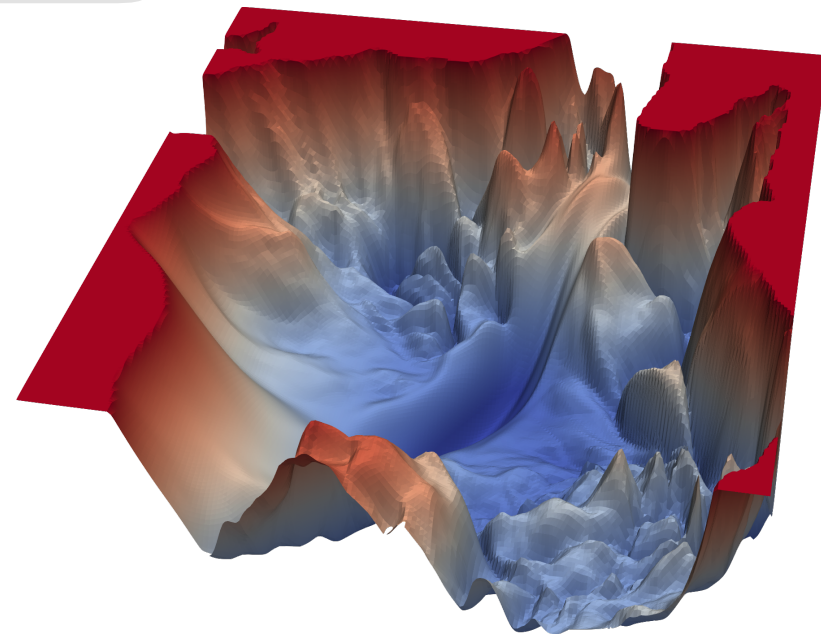
Learning problem: minimize $R(\hat{f})$
 $\boldsymbol{\theta} \in \mathbb{R}^{pn_H}$

Learning Neural Network from Data

Learning problem: minimize $R(\hat{f})$
 $\theta \in \mathbb{R}^{p_{\text{H}}}$

Challenge: highly non-convex (many local minima)

Surprise: SGD and its variants work in practice!!



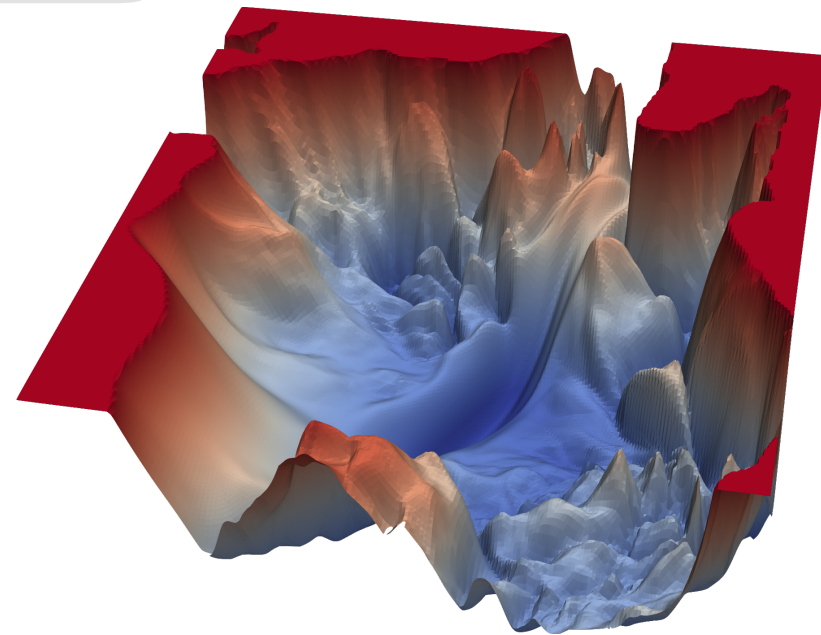
Learning Neural Network from Data

Learning problem: minimize $R(\hat{f})$
 $\theta \in \mathbb{R}^{p_{\text{H}}}$

Challenge: highly non-convex (many local minima)

Surprise: SGD and its variants work in practice!!

Good news: emerging theory (starting in 2018!!)



Chizat and Bach (NIPS 2018), Mei, Montanari and Nguyen (PNAS 2018), Rotskoff and Vanden-Eijnden (arXiv:1805.00915, 2018), Williams et al (arXiv:1906.07842, 2019)

Idea: Think of the mean field, i.e., infinite width ($n_{\text{H}} \rightarrow \infty$) limit

$$\hat{f} \equiv \hat{f}(\mathbf{x}, \rho) = \int_{\mathbb{R}^p} \Phi(\mathbf{x}, \boldsymbol{\theta}) \rho(\boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta}$$

Then, learning problem: minimize $R(\hat{f})$
 $\rho \in \mathcal{P}_2(\mathbb{R}^p)$

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

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

PDF dynamics for SGD:

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

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

Wasserstein Proximal Recursion for Training NN

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

Classifying two Gaussians:

$$d = 40, n = 100,$$

$$a = 1, b = 0, \sigma(\cdot) = \tanh(\cdot),$$

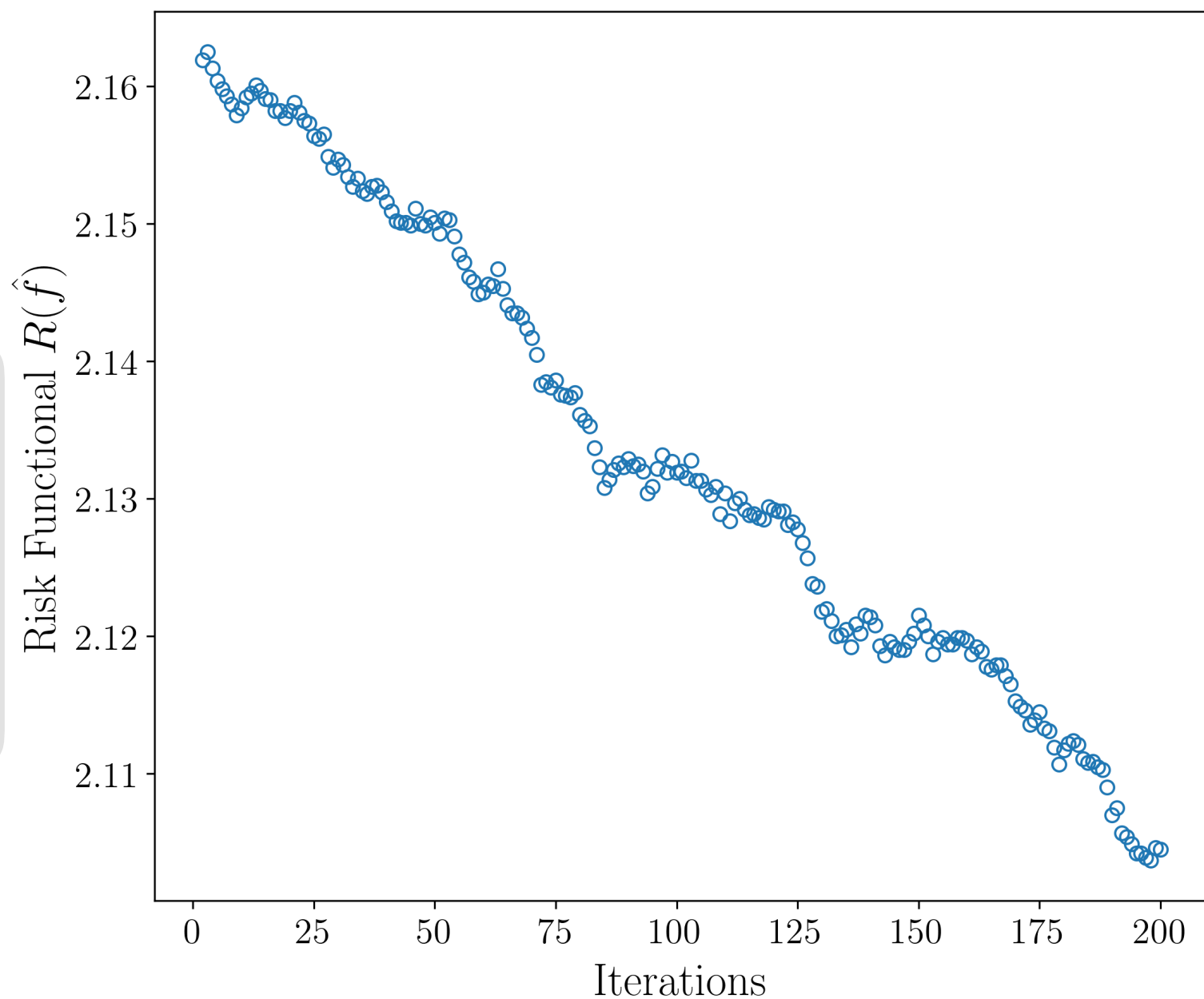
Joint law of $(\mathbf{x}, y) \in \mathbb{R}^d \times \mathbb{R}$:

$$\text{Prob}(y = +1, \mathbf{x} \sim \mathcal{N}(\mathbf{0}, (1 + \Delta)^2 \mathbf{I}_d)) = \frac{1}{2},$$

$$\text{Prob}(y = -1, \mathbf{x} \sim \mathcal{N}(\mathbf{0}, (1 - \Delta)^2 \mathbf{I}_d)) = \frac{1}{2},$$

$$\tau = 10^{-3}, n_{\text{sample}} = 100, \Delta = 0.2,$$

$$\text{Noisy SGD with } \beta = \frac{1}{3}$$



Take Home Message



Thank You

Support:



CITRIS
PEOPLE AND
ROBOTS

