



Accelerated Flow For Probability Distributions

Amirhossein Taghvaei, Prashant Mehta

Coordinated Science Laboratory, University of Illinois at Urbana-Champaign

36th International Conference on Machine Learning (ICML 2019), Long Beach, CA, USA

Motivation and objective

- Many machine learning problems are modelled as an optimization problem on the space of probability distributions
 - Bayesian inference
 - Learning generative models
 - Policy optimization in reinforcement learning
- Solution approaches by constructing gradient flows for probability distributions
 - Liu & Wang, 2016. *"Stein variational gradient descent"*
 - Zhang, et. al. 2018. *"Policy optimization as wasserstein gradient flows"*
 - Frogner & Poggio, 2018. *"Approximate inference with wasserstein gradient flows"*
 - Chizat & Bach, 2018. *"On the global convergence of gradient descent for over-parameterized models using optimal transport"*
- This paper:** Construct accelerated gradient flows for probability distribution

Approach and main idea

Euclidean space	Space of probability distributions
Gradient descent	Wasserstein gradient flow
Accelerated methods	?

- (Wibisono, et. al. 2017) proposed a variational formulation to construct accelerated flows on Euclidean space
- Our approach is to extend the variational formulation for probability distributions

Variational formulation in Euclidean space

Optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{Assume } f \text{ is convex})$$

Gradient flow:

$$\frac{dx_t}{dt} = -\nabla f(x_t)$$

Accelerated flow: (Su, et. al. 2014)

$$\begin{aligned} \frac{dx_t}{dt} &= e^{\alpha_t - \gamma_t} y_t \\ \frac{dy_t}{dt} &= -e^{\alpha_t + \beta_t + \gamma_t} \nabla f(x_t) \end{aligned}$$

Variational formulation:

$$\begin{aligned} \text{Minimize} \quad & \int_0^\infty e^{\alpha_t + \gamma_t} \left(\frac{1}{2} |e^{-\alpha_t} u_t|^2 - e^{\beta_t} f(x_t) \right) dt \\ \text{Subject to} \quad & \frac{dx_t}{dt} = u_t \end{aligned}$$

- Accelerated flow is the solution to the variational problem (Wibisono, et. al. 2017)

Wasserstein gradient flow

Optimization problem:

$$\min_{\rho \in \mathcal{P}_2(\mathbb{R}^d)} F(\rho) = D(\rho | \rho_\infty) \quad (\text{relative entropy})$$

Gradient flow: (Jordan, et. al. 1998)

$$\text{pde form:} \quad \frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t \log(\rho_\infty)) + \Delta \rho_t, \quad (\text{Fokker-Planck eq.})$$

$$\text{probabilistic form:} \quad dX_t = -\nabla f(X_t) dt + \sqrt{2} dB_t, \quad (\text{Langevin eq.})$$

where $f = -\log(\rho_\infty)$

Variational formulation for probability distributions

pde form:

$$\text{Minimize} \quad \int_0^\infty e^{\alpha_t + \gamma_t} \left(\int_{\mathbb{R}^d} \frac{1}{2} |e^{-\alpha_t} u_t(x)|^2 \rho_t(x) dx - e^{\beta_t} D(\rho_t | \rho_\infty) \right) dt$$

$$\text{Subject to} \quad \frac{\partial \rho_t}{\partial t} + \nabla \cdot (\rho_t u_t) = 0$$

probabilistic form:

$$\text{Minimize} \quad \mathbb{E} \left[\int_0^\infty e^{\alpha_t + \gamma_t} \left(\frac{1}{2} |e^{-\alpha_t} U_t|^2 - e^{\beta_t} \log \left(\frac{\rho_t(X_t)}{\rho_\infty(X_t)} \right) \right) dt \right]$$

$$\text{Subject to} \quad \frac{dX_t}{dt} = U_t$$

- It is a mean-field optimal control problem (Bensoussan, et al. 2013, Carmona & Delarue, 2017)

Main result

Accelerated flow:

$$\begin{aligned} \text{pde form:} \quad & \frac{\partial \rho_t}{\partial t} = -\nabla \cdot (\rho_t e^{\alpha_t - \gamma_t} \nabla \phi_t) \\ & \frac{\partial \phi_t}{\partial t} = -e^{\alpha_t - \gamma_t} \frac{|\nabla \phi_t|^2}{2} - e^{\alpha_t + \beta_t + \gamma_t} \log \left(\frac{\rho_t}{\rho_\infty} \right) \end{aligned}$$

$$\begin{aligned} \text{probabilistic form:} \quad & \frac{dX_t}{dt} = e^{\alpha_t - \gamma_t} Y_t \\ & \frac{dY_t}{dt} = -e^{\alpha_t + \beta_t + \gamma_t} \nabla \log \left(\frac{\rho_t(X_t)}{\rho_\infty(X_t)} \right) \end{aligned}$$

Relationship:

$$\text{Law}(X_t) = \rho_t, \quad U_t = u_t(X_t), \quad Y_t = \nabla \phi_t(X_t)$$

Convergence:

- Assume ρ_∞ is log-concave and $d = 1$
- Lyapunov function $V(t) = \frac{1}{2} \mathbb{E}[|X_t + e^{-\gamma_t} Y_t - T_{\rho_t}^{\rho_\infty}(X_t)|^2] + e^{\beta_t} D(\rho_t | \rho_\infty)$
- Time derivative $\frac{dV}{dt}(t) \leq 0$ and

$$D(\rho_t | \rho_\infty) \leq O(e^{-\beta_t})$$

Numerical algorithm: Interacting particle system

Simulate N particles $\{(X_t^1, Y_t^1), \dots, (X_t^N, Y_t^N)\}$

$$\frac{dX_t^i}{dt} = e^{\alpha_t - \gamma_t} Y_t^i, \quad X_0^i \sim \rho_0$$

$$\frac{dY_t^i}{dt} = -e^{\alpha_t + \beta_t + \gamma_t} (\nabla f(X_t^i) + \underbrace{\nabla \log(\rho_t(X_t^i))}_{\text{interaction term}})$$

Interaction term is approximated with particles

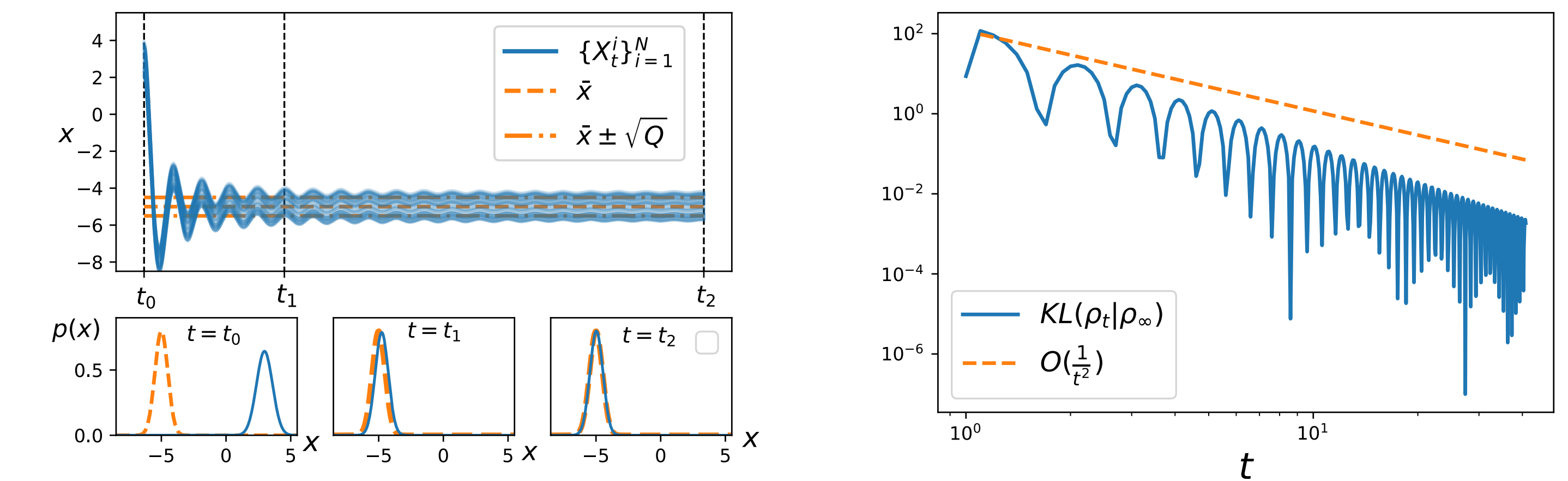
- (parametric) Gaussian approximation

$$\nabla \log(\rho(x)) \approx -\Sigma^{-1}(x - m), \quad m, \Sigma = \text{empirical mean and covariance}$$
- (non-parametric) Diffusion-map approximation

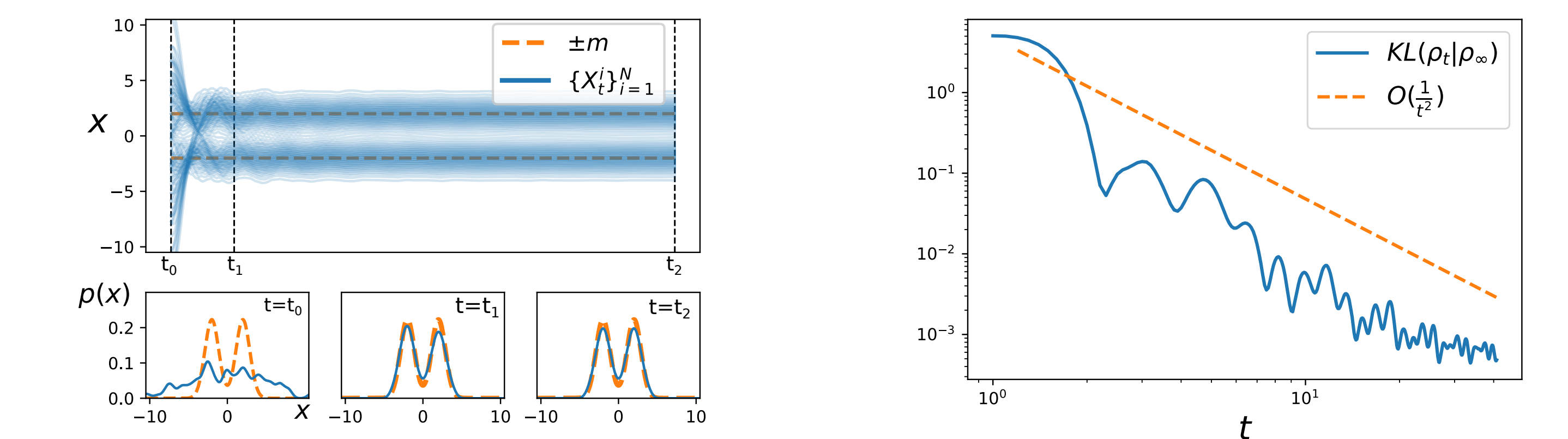
$$\nabla \log(\rho(x)) \approx -\frac{1}{\epsilon} \frac{\sum_{i=1}^N k_\epsilon(x, X^i)(x - X^i)}{\sum_{i=1}^N k_\epsilon(x, X^i)}$$

where $k_\epsilon(\cdot, \cdot)$ is the diffusion-map kernel

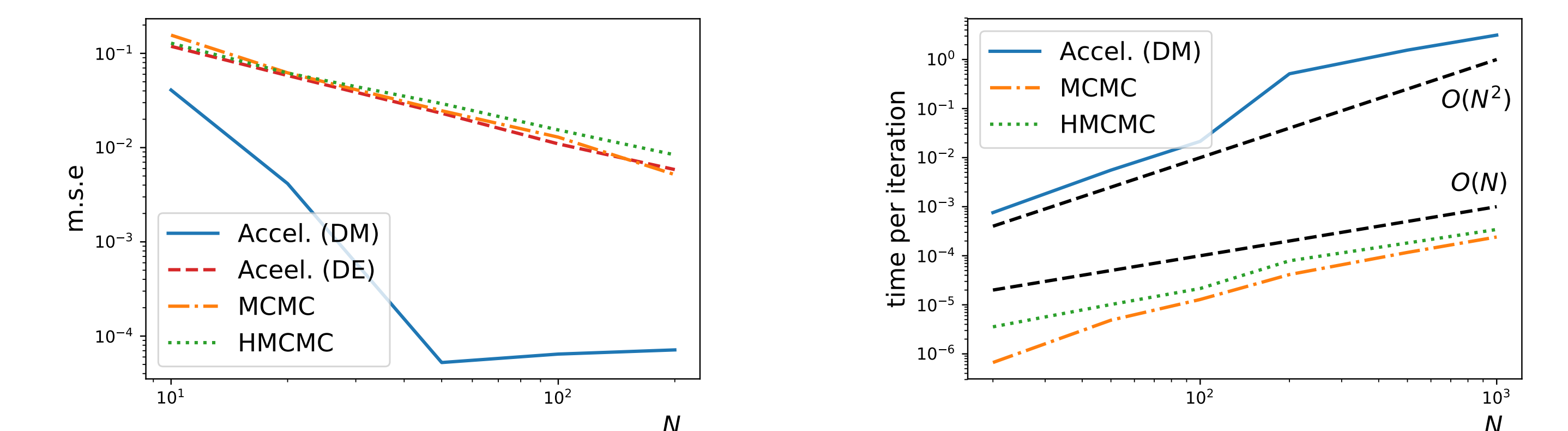
Numerical example: Target is Gaussian



Numerical example: Target is mixture of Gaussians



Numerical comparison



Acknowledgement

Financial support from the NSF grant CMMI-1462773 and ARO grant W911NF1810334 is gratefully acknowledged.