

---

# Theoretical Foundations of Wasserstein Policy Optimization

---

Anonymous AI Agent (first author)   Anonymous Human Co-author(s)

Affiliation

Address

email

## Abstract

1        We revisit Wasserstein Policy Optimization (WPO) as policy transport on action  
2        densities followed by a projection onto a parametric manifold. Evolving poli-  
3        cies by a 2-Wasserstein gradient flow and projecting in the Fisher/KL inner prod-  
4        uct yields a covariant natural step with a mixed-derivative cross-term. We make  
5        this projection-based view explicit, prove baseline invariance (via a constrained  
6        Gâteaux variation) and parameterization covariance (via Fisher pullbacks), and  
7        delineate when the step coincides with natural policy gradient (affine-in-action  
8        exponential families) versus when it departs (mixtures, squashings). For Gaus-  
9        sian policies we give mean and covariance updates, including a full-covariance  
10       Cholesky implementation that preserves SPD. We extend to c-Wasserstein dynam-  
11       ics to obtain principled stability via convex conjugates and state precise energy  
12       inequalities in the frozen-critic regime. Assumptions and weak-form conditions  
13       are spelled out, and connections to classic PG/DPG/NPG are established.

## 14    1   Introduction

15    Wasserstein Policy Optimization (WPO) offers a principled link between optimal transport and pol-  
16    icy updates through Wasserstein gradient flows. Despite elegant foundations, theoretical questions  
17    remain: the precise projection from infinite-dimensional flows to parametric updates; invariance to  
18    baselines and parameterizations; and stability when action-value gradients are large. We present a  
19    self-contained treatment addressing these points with label-anchored derivations, and we identify  
20    the conditions under which the WPO update coincides with, or departs from, classic policy gradi-  
21    ent/natural gradient directions.

## 22    2   Background

23    On the 2-Wasserstein manifold, steepest descent of a functional  $\mathcal{J}[\pi]$  follows the continuity equation  
24    with velocity  $v = -\nabla_{\mathbf{a}}(\delta\mathcal{J}/\delta\pi)$  [Ambrosio et al., 2008, Benamou and Brenier, 2000]. Taking  
25     $\mathcal{J}[\pi] = -\mathbb{E}[Q^\pi]$  yields  $\partial_t\pi = -\nabla_{\mathbf{a}} \cdot (\pi \nabla_{\mathbf{a}} Q^\pi)$ . Projecting this infinite-dimensional flow to a  
26    parametric family  $\{\pi_\theta\}$  leads to the natural-gradient form  $\Delta\theta = F_{\theta\theta}^{-1} \mathcal{F}_{t\theta}$ , where the cross term is  
27     $\mathbb{E}[\nabla_\theta \nabla_{\mathbf{a}} \log \pi \nabla_{\mathbf{a}} Q^\pi]$  under mild regularity and boundary conditions. We relate to natural policy  
28    gradient [Kakade, 2001] and to neural ODE views of transport [Chen et al., 2018].

29    **Weak form and boundary conditions.** We interpret the continuity equation in weak form, with  
30    no-flux boundary condition  $(\pi v) \cdot n = 0$  on bounded action domains or vanishing flux at infinity on  
31     $\mathbb{R}^d$ . Under  $C^1/C^2$  regularity and dominated convergence, integration by parts is valid and boundary  
32    terms vanish; all energy identities and projections below are stated in this weak-form sense.

### 3 Preliminaries and Notation

Let  $\mathcal{A} \subseteq \mathbb{R}^d$  denote the action domain and  $\mathcal{S}$  the state space. Policies are absolutely continuous densities  $\pi_\theta(\mathbf{a} \mid s)$  on  $\mathcal{A}$  for each  $s \in \mathcal{S}$ . We write expectations as  $\mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi_\theta(\cdot \mid s)}[\cdot]$  for a fixed reference weighting  $d^\pi(s)$  (e.g., discounted occupancy). In line with the policy gradient theorem, our first-order variations operate in a semi-gradient setting where  $d^\pi$  and  $Q^\pi$  are treated as fixed.

**Notation.**  $\mathcal{S}$  states,  $\mathcal{A} \subseteq \mathbb{R}^d$  actions;  $d^\pi$  discounted occupancy;  $\pi_\theta(\mathbf{a} \mid s)$  policy density;  $Q^\pi(s, \mathbf{a})$  critic;  $F_{\theta\theta}$  Fisher;  $\text{sym}(X) = \frac{1}{2}(X + X^\top)$ ;  $\mathbb{S}_{++}^d$  SPD cone. We use bold  $\mathbf{a}$  for vectors and  $a$  for 1D examples.

Throughout, expectations are with respect to the joint measure  $s \sim d^\pi$ ,  $\mathbf{a} \sim \pi_\theta(\cdot \mid s)$ . When we abbreviate  $\mathbb{E}_{s, \mathbf{a}}[\cdot]$ , this is the understood measure unless stated otherwise.

The Fisher matrix is

$$F_{\theta\theta} = \mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi_\theta(\cdot \mid s)} [\nabla_\theta \log \pi_\theta(\mathbf{a} \mid s) \nabla_\theta \log \pi_\theta(\mathbf{a} \mid s)^\top]. \quad (1)$$

We equip density variations with the inner product

$$\langle f, g \rangle = \mathbb{E}_{s \sim d^\pi} \left[ \int_{\mathcal{A}} \frac{f(\mathbf{a}, s) g(\mathbf{a}, s)}{\pi_\theta(\mathbf{a} \mid s)} d\mathbf{a} \right], \quad (2)$$

for which the parametric tangent directions are  $\pi \nabla_\theta \log \pi$ . Throughout, we assume sufficient smoothness ( $C^1/C^2$ ), dominated convergence conditions to exchange expectation and differentiation, and vanishing boundary terms so that integration by parts is valid.

## 4 Derivation of WPO

### 4.1 Problem Setup

We consider continuous control with policy  $\pi_\theta(\mathbf{a} \mid s)$  and critic  $Q^\pi(s, \mathbf{a})$ . We adopt a reference state distribution  $d^\pi(s)$  (e.g., discounted occupancy) and study policy updates that arise as projections of Wasserstein gradient flows on action densities.

### 4.2 From Wasserstein Flow to Parametric Update

Let  $\mathcal{J}[\pi] = -\mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi(\cdot \mid s)}[Q^\pi(s, \mathbf{a})]$ . We adopt a per-state time rescaling and work with the rescaled flow whose velocity uses  $v = \nabla_{\mathbf{a}} Q^\pi$  (absorbing the  $d^\pi$  factor). The 2-Wasserstein flow is

$$\partial_t \pi = -\nabla_{\mathbf{a}} \cdot (\pi \nabla_{\mathbf{a}} Q^\pi). \quad (3)$$

We note that embedding the discounted state weighting  $d^\pi(s)$  into  $\delta \mathcal{J} / \delta \pi$  scales the per-state velocity by  $d^\pi(s)$ . This is a per-state time reparameterization and preserves weak-form energy identities, but a Fisher-projected finite step may change unless additional conditions (e.g., state-wise collinearity of velocities or constant scaling) hold. To remove this ambiguity, we adopt the convention  $\delta \mathcal{J} / \delta \pi = -Q^\pi$  and carry  $d^\pi(s)$  only as an outer expectation in all cross terms and Fisher quantities. We project to  $\{\pi_\theta\}$  by minimizing, in the Fisher metric, the discrepancy between  $\partial_t \pi$  and  $\delta \pi_\theta = \pi (\nabla_\theta \log \pi) \cdot \Delta \theta$ , leading to  $F_{\theta\theta} \Delta \theta = \mathcal{F}_{t\theta}$  with

$$\mathcal{F}_{t\theta} = \mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi} [\nabla_\theta \nabla_{\mathbf{a}} \log \pi_\theta(\mathbf{a} \mid s) \nabla_{\mathbf{a}} Q^\pi(s, \mathbf{a})]. \quad (4)$$

### 4.3 Energy Dissipation

Under suitable regularity, and in the frozen-critic regime, the rescaled gradient flow (3) dissipates  $\mathcal{J}$ :

**Lemma 1** (Energy dissipation). *Along solutions  $\pi_t$  of (3), one has  $\frac{d}{dt} \mathcal{J}[\pi_t] \leq 0$  with  $\frac{d}{dt} \mathcal{J}[\pi_t] = -\mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi_t} [\|\nabla_{\mathbf{a}} Q^\pi(s, \mathbf{a})\|^2]$ . For the  $c$ -Wasserstein flow with velocity  $\nabla c^*(\nabla_{\mathbf{a}} \delta \mathcal{J} / \delta \pi)$  one has  $\frac{d}{dt} \mathcal{J}[\pi_t] = -\mathbb{E}_{s, \mathbf{a}} [\langle \nabla c^*(\nabla_{\mathbf{a}} Q^\pi), \nabla_{\mathbf{a}} Q^\pi \rangle] \leq 0$  by monotonicity of  $\nabla c^*$ . See Appendix (Energy Dissipation) for weak-form derivations and function-space assumptions.*

#### 70 4.4 Functional Derivative of $\mathcal{J}$

71 For discounted RL, under the semi-gradient convention (holding  $d^\pi$  and  $Q^\pi$  fixed), we adopt the  
72 rescaled functional derivative compatible with (3):

$$\frac{\delta \mathcal{J}}{\delta \pi}(s, \mathbf{a}) = -Q^\pi(s, \mathbf{a}), \quad d^\pi(s) = (1 - \gamma) \sum_t \gamma^t \Pr(s_t = s). \quad (5)$$

73 This corresponds to absorbing  $d^\pi(s)$  via a per-state time reparameterization. Under the alternative  
74 convention that embeds  $d^\pi(s)$  inside  $\delta \mathcal{J} / \delta \pi$ , the Fisher-projected finite step can differ unless state-  
75 wise collinearity or constant scaling holds; our adopted convention avoids this ambiguity while  
76 preserving the weak-form energy identities.

#### 77 4.5 KL/Fisher Projection and Normal Equations

78 Equivalently, minimizing  $\|\partial_t \pi - \delta \pi_\theta\|^2$  in the Fisher norm (2) yields the normal equations  $F_{\theta\theta} \Delta \theta =$   
79  $\mathcal{F}_{t\theta}$ . The cross term satisfies

$$\mathcal{F}_{t\theta} = \int \nabla_\theta \log \pi_\theta(\mathbf{a} | s) \partial_t \pi_\theta(\mathbf{a} | s) d\mathbf{a} \quad (6)$$

$$= - \int \nabla_\theta \log \pi_\theta \nabla_{\mathbf{a}} \cdot (\pi_\theta \nabla_{\mathbf{a}} Q^\pi) d\mathbf{a} \quad (7)$$

$$= \int (\nabla_{\mathbf{a}} \nabla_\theta \log \pi_\theta)^\top (\pi_\theta \nabla_{\mathbf{a}} Q^\pi) d\mathbf{a} \quad (8)$$

$$= \mathbb{E}_{\mathbf{a} \sim \pi_\theta} [\nabla_\theta \nabla_{\mathbf{a}} \log \pi_\theta \nabla_{\mathbf{a}} Q^\pi], \quad (9)$$

80 where boundary terms vanish and derivatives commute under the stated assumptions. Under  $\mathcal{C}^2$   
81 regularity and vanishing boundary terms, this follows from integration by parts. The natural-gradient  
82 step is  $\Delta \theta = F_{\theta\theta}^{-1} \mathcal{F}_{t\theta}$ .

#### 83 4.6 Projection Metric: Fisher vs. $W_2$

84 We project the  $W_2$  flow in the Fisher/KL inner product on densities, which induces a covariant  
85 natural step on the parametric manifold and yields simple estimators. A  $W_2$ -parametric projection  
86 would instead align the parametric velocity with the  $W_2$  tangent metric via a velocity potential on  
87 parameters.

88 *Operator-level relation (per state).* Let  $\mathcal{T}_\theta$  denote the parametric tangent span  $\{\pi \nabla_\theta \log \pi u : u \in$   
89  $\mathbb{R}^{\dim \theta}\}$ . The Fisher projector and the  $W_2$ -tangent projector are two positive self-adjoint operators  
90 on  $\mathcal{T}_\theta$  induced by different inner products. They coincide when  $\mathcal{T}_\theta$  is invariant and the per-state  
91 velocity is collinear with the tangent directions (e.g., exponential families with sufficient statistics  
92 affine in  $\mathbf{a}$  under a compatible parameterization, such as Gaussian means with fixed covariance). In  
93 general they differ by a positive operator on  $\mathcal{T}_\theta$ . We favor Fisher projection for statistical stability,  
94 parameterization covariance, and practicality.

#### 95 4.7 Stability via c-Wasserstein

96 For convex  $c$ , the c-Wasserstein flow uses the velocity  $v = \nabla c^*(\nabla_{\mathbf{a}}(\delta \mathcal{J} / \delta \pi)) = \nabla c^*(-\nabla_{\mathbf{a}} Q^\pi)$ .  
97 The corresponding flow is  $\partial_t \pi = -\nabla_{\mathbf{a}} \cdot (\pi \nabla c^*(\nabla_{\mathbf{a}}(\delta \mathcal{J} / \delta \pi)))$ . If  $c^*$  is even, then  $\nabla c^*(-x) =$   
98  $-\nabla c^*(x)$  and one may equivalently write  $\partial_t \pi = -\nabla_{\mathbf{a}} \cdot (\pi \nabla c^*(\nabla_{\mathbf{a}} Q^\pi))$ . For a separable choice  
99  $c(u) = \frac{1}{4} \sum_i |u_i|^4$ , one has  $c^*(x) = \frac{3}{4} \sum_i |x_i|^{4/3}$  and the elementwise cube-root map  $\nabla c^*(x)_i =$   
100  $\text{sign}(x_i) |x_i|^{1/3}$ , controlling large action-gradients while preserving descent/ascent directions. A  
101 rotationally invariant alternative  $c(u) = \frac{1}{4} \|u\|_2^4$  yields  $\nabla c^*(x) = \|x\|_2^{-2/3} x$  (radial shrinkage).  
102 These shrinkage maps are direction-preserving, monotone, and Hölder-continuous, ensuring the  
103 energy identity holds and often allowing larger stable steps than heuristic gradient clipping under  
104 the same frozen-critic assumptions.

105 **Theorem 1** (Energy decay for c-Wasserstein flows). *Under Assumption 1, let  $c$  be proper, convex,*  
106 *l.s.c., with differentiable  $c^*$  and monotone  $\nabla c^*$ . The flow*

$$\partial_t \pi = -\nabla_{\mathbf{a}} \cdot \left( \pi \nabla c^*(\nabla_{\mathbf{a}}(\delta \mathcal{J} / \delta \pi)) \right)$$

107 satisfies the intrinsic energy identity

$$\frac{d}{dt} \mathcal{J}[\pi_t] = -\mathbb{E}_{s, \mathbf{a}} \left[ \langle \nabla c^*(\nabla_{\mathbf{a}}(\delta \mathcal{J}/\delta \pi)), \nabla_{\mathbf{a}}(\delta \mathcal{J}/\delta \pi) \rangle \right] \leq 0,$$

108 with equality iff  $\nabla_{\mathbf{a}}(\delta \mathcal{J}/\delta \pi) = 0$  almost everywhere. Specializing to  $\delta \mathcal{J}/\delta \pi = -Q^\pi$  gives  
 109  $-\mathbb{E}[\langle \nabla c^*(-\nabla_{\mathbf{a}} Q^\pi), -\nabla_{\mathbf{a}} Q^\pi \rangle] \leq 0$ ; if in addition  $c^*$  is even (hence  $\nabla c^*$  odd), this equals  
 110  $-\mathbb{E}[\langle \nabla c^*(\nabla_{\mathbf{a}} Q^\pi), \nabla_{\mathbf{a}} Q^\pi \rangle] \leq 0$ . For  $c(u) = \frac{1}{4} \sum_i |u_i|^4$  one has elementwise shrinkage  
 111  $\nabla c^*(x)_i = \text{sign}(x_i)|x_i|^{1/3}$ ; for  $c(u) = \frac{1}{4} \|u\|_2^4$  one has radial shrinkage  $\nabla c^*(x) = \|x\|_2^{-2/3} x$ .

112 Practical note: the induced map is direction-preserving and monotone, which yields energy decay  
 113 under frozen critics. Unlike heuristic gradient clipping, it retains a variational interpretation and  
 114 often permits larger stable steps in ill-conditioned regimes. See Theorem 1 and Appendix A.2.

#### 115 4.8 Estimators and variance reduction

116 With a reparameterization  $\mathbf{a} = f_\theta(\epsilon, s)$  and a frozen critic, an unbiased single-sample estimator of  
 117 the cross term is

$$\widehat{g}_\theta(\epsilon, s) = \nabla_\theta \nabla_{\mathbf{a}} \log \pi_\theta(f_\theta(\epsilon, s) \mid s) \widehat{\nabla_{\mathbf{a}} Q}(s, f_\theta(\epsilon, s)).$$

118 Mixed derivatives are computed efficiently via JVP/VJP (one forward- and one reverse-mode pass  
 119 suffice in most autodiff systems). The critic action-gradient  $\nabla_{\mathbf{a}} Q$  comes from differentiating the  
 120 critic wrt its action input. A score-aligned control variate subtracts  $c(s) \nabla_\theta \log \pi_\theta(\mathbf{a} \mid s)$ , with  
 121  $c(s)$  chosen by least-squares to minimize variance, leaving  $\mathbb{E}[\widehat{g}_\theta]$  unchanged. The optimal per-state  
 122 coefficient is

$$c^*(s) = \frac{\text{Cov}(\nabla_\theta \log \pi_\theta, \widehat{g}_\theta \mid s)}{\text{Var}(\nabla_\theta \log \pi_\theta \mid s)}, \quad (10)$$

123 estimated online via running moments.

#### 124 4.9 Gaussian Parameter Updates

125 For Gaussian policies, 1D mean/variance updates are

$$\Delta \mu = \mathbb{E}[\nabla_a Q], \quad (11)$$

$$\Delta(\sigma^2) = 2 \mathbb{E}[(a - \mu) \nabla_a Q]. \quad (12)$$

126 With diagonal covariance, updates act coordinate-wise. For full covariance  $\Sigma \in \mathbb{S}_{++}^d$ , the Fisher-  
 127 natural step simplifies to  $\Delta \Sigma = M + M^\top = 2 \text{sym}(M)$  with  $M = \mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi_\theta}[(\nabla_{\mathbf{a}} Q)(\mathbf{a} - \mu)^\top]$ .  
 128 Sampling formulas follow directly from the expectation definitions. The full-covariance update  
 129 entails  $\mathcal{O}(d^3)$  algebra (linear solves), whereas diagonal updates are  $\mathcal{O}(d)$ . In practice, full covariance  
 130 helps under strong anisotropy/ill-conditioning; diagonal is preferable when samples or compute are  
 131 limited.

132 **Full-covariance via Cholesky.** Write  $\Sigma = LL^\top$  with  $L$  lower-triangular and  $\text{diag}(L) > 0$ . To  
 133 implement the natural step stably, solve for  $\Delta L$  in the triangular Sylvester equation

$$L \Delta L^\top + \Delta L L^\top = S, \quad S = M + M^\top,$$

134 by forward substitution by columns  $j = 1, \dots, d$ :

$$(2L_{jj})(\Delta L)_{jj} = S_{jj} - 2 \sum_{k < j} L_{jk}(\Delta L)_{jk},$$

$$(\Delta L)_{ij} = \frac{S_{ij} - \sum_{k < j} (L_{ik}(\Delta L)_{jk} + (\Delta L)_{ik} L_{jk}) - L_{ij}(\Delta L)_{jj}}{L_{jj}}, \quad i = j + 1, \dots, d.$$

135 Then update  $L \leftarrow L + \eta \Delta L$  with a line search to keep  $\text{diag}(L) > 0$ . This realizes  $\Delta \Sigma = S$  while  
 136 preserving  $\Sigma \in \mathbb{S}_{++}^d$ . In practice, shrink  $\eta$  (backtracking) until  $\min_i (L + \eta \Delta L)_{ii} > 0$ ; this prevents  
 137 loss of SPD under large steps.

138 **Deterministic limit.** As  $\Sigma \rightarrow 0$ , the mean update reduces to the deterministic policy-gradient  
 139 direction evaluated at the mean, while the covariance update vanishes. The Fisher metric provides a  
 140 well-defined limit.

141 **Quadratic toy.** If locally  $\nabla_{\mathbf{a}} Q^\pi(s, \mathbf{a}) = H(s)(\mathbf{a} - \mu)$  with  $H(s) \succeq 0$ , then  $M = H(s)\Sigma$  and  
 142  $\Delta\Sigma = H(s)\Sigma + \Sigma H(s)$ , while  $\Delta\mu = 0$  at  $\mathbf{a} = \mu$ , consistent with the deterministic limit.

#### 143 4.10 Main Results

144 We state the primary results under the hypotheses above; proofs are provided in the appendix. Scope:  
 145 results are stated with fixed critic  $Q^\pi$  and state weighting  $d^\pi$  (semi-gradient convention). When both  
 146 evolve, additional residual terms appear.

147 **Assumption 1** (Standing hypotheses). (i)  $\log \pi_\theta(\mathbf{a} \mid s) \in C^2$  in  $(\mathbf{a}, \theta)$  with mixed partials com-  
 148 muting almost everywhere, and  $\nabla_\theta \nabla_{\mathbf{a}} \log \pi_\theta \in L^1(d^\pi \otimes \pi_\theta)$ ; (ii)  $Q^\pi(\cdot, s) \in C^1$  with  $\nabla_{\mathbf{a}} Q^\pi \in$   
 149  $L^2(\pi_\theta(\cdot \mid s))$  uniformly in  $s$ ; (iii) either  $\mathcal{A} = \mathbb{R}^d$  with vanishing flux at infinity or a no-flux bound-  
 150 ary on  $\partial\mathcal{A}$ ; (iv) dominated convergence justifies swapping expectations and derivatives. All energy  
 151 statements hold in the frozen-critic regime.

152 **Theorem 2** (Projection to Natural Gradient). Let  $\partial_t \pi$  satisfy (3) for  $\mathcal{J}[\pi] = -\mathbb{E}[Q^\pi]$ . With the inner  
 153 product (2), the Galerkin orthogonality conditions  $\langle \partial_t \pi - \delta \pi_\theta, \pi \nabla_\theta \log \pi \rangle = 0$  yield the normal  
 154 equations  $F_{\theta\theta} \Delta\theta = \mathcal{F}_{t\theta}$  with  $\mathcal{F}_{t\theta}$  in (4). Hence  $\Delta\theta = F_{\theta\theta}^{-1} \mathcal{F}_{t\theta}$ .

155 Finite steps  $\Delta\theta = \eta F_{\theta\theta}^{-1} \mathcal{F}_{t\theta}$  need not monotonically decrease  $\mathcal{J}$  unless  $\eta$  is sufficiently small; a  
 156 global line search (Armijo/backtracking) can ensure descent. Locally, if  $g(\theta) := \mathcal{F}_{t\theta}$  is  $L$ -Lipschitz  
 157 and  $F_{\theta\theta} \succeq mI$  in a neighborhood, then choosing  $\eta \leq m/L$  yields a guaranteed decrease for suffi-  
 158 ciently small neighborhoods.

159 **Corollary 1** (Normal equations). Under Assumption 1, the Fisher–Galerkin projection yields  
 160  $F_{\theta\theta} \Delta\theta = \mathcal{F}_{t\theta}$  with  $\mathcal{F}_{t\theta} = \mathbb{E}[\nabla_\theta \nabla_{\mathbf{a}} \log \pi \nabla_{\mathbf{a}} Q^\pi]$ .

161 **Proposition 1** (Baseline Invariance). For any baseline  $b(s)$ , replacing  $Q^\pi$  by  $A^\pi = Q^\pi - b(s)$   
 162 leaves both the PDE and the parametric cross term  $\mathcal{F}_{t\theta}$  unchanged. Equivalently, the constrained  
 163 first variation under per-state normalization satisfies  $\delta\mathcal{J}/\delta\pi = -(Q^\pi - b(s))$ .

164 Caution: action-dependent adjustments  $b(s, \mathbf{a})$  are not baselines in this sense; they alter  $\nabla_{\mathbf{a}} Q$  and  
 165 hence change both the PDE velocity and the projected update.

166 **Theorem 3** (Parameterization Covariance). Let  $\phi = \phi(\theta)$  be a local diffeomorphism with Jacobian  
 167  $J = \partial\phi/\partial\theta$  and pullback Fisher  $F_\phi = J^{-\top} F_\theta J^{-1}$ . Define  $g_\theta = \mathbb{E}[\nabla_\theta \nabla_{\mathbf{a}} \log \pi \nabla_{\mathbf{a}} Q]$  and  $g_\phi =$   
 168  $J^{-\top} g_\theta$ . Then the natural-gradient step is covariant:  $\Delta\phi = F_\phi^{-1} g_\phi = J F_\theta^{-1} g_\theta = J \Delta\theta$ .

169 **Lemma 2** (Gaussian Family). For a 1D Gaussian policy with fixed variance,  $\Delta\mu = \mathbb{E}[\nabla_{\mathbf{a}} Q]$ . With  
 170 log-variance  $\lambda = \log \sigma$ ,  $\Delta\lambda = \sigma^{-2} \mathbb{E}[(a - \mu) \nabla_{\mathbf{a}} Q]$ , i.e.,  $\Delta(\sigma^2) = 2 \mathbb{E}[(a - \mu) \nabla_{\mathbf{a}} Q]$ . For full  
 171 covariance  $\Sigma$ , with  $M = \mathbb{E}[(\nabla_{\mathbf{a}} Q)(\mathbf{a} - \mu)^\top]$  and  $G = \text{sym}(\Sigma^{-1} M \Sigma^{-1})$ , the affine-invariant  
 172 Fisher yields  $\Delta\Sigma = 2 \Sigma G \Sigma$ , which reduces to  $\Delta\Sigma = M + M^\top$ .

173 **Proposition 2** (c-Wasserstein Gradient Flow). For convex  $c$  with conjugate  $c^*$ , replacing  $\nabla_{\mathbf{a}} Q$  by  
 174  $\nabla c^*(\nabla_{\mathbf{a}} Q)$  defines the  $c$ -Wasserstein gradient flow of  $\mathcal{J}$  (in the sense of Ambrosio et al., 2008).  
 175 Under standard regularity,  $\mathcal{J}$  decreases along solutions. The elementwise cube-root arises from  
 176  $c^*(x) = \frac{3}{4}|x|^{4/3}$ .

177 **Proposition 3** (Alignment with NPG for Gaussian means). Let  $\pi_\theta(\cdot \mid s) = \mathcal{N}(\mu_\theta(s), \Sigma)$  with fixed  
 178  $\Sigma$ . Suppose  $\nabla_{\mathbf{a}} Q^\pi(s, \mathbf{a}) \approx H(s)(\mathbf{a} - \mu_\theta(s))$  in a neighborhood of  $\mu_\theta(s)$  with  $H(s) \succeq 0$ . Then the  
 179 WPO mean step  $\Delta\theta_\mu = F_{\theta_\mu}^{-1} \mathbb{E}[\nabla_{\theta_\mu} \nabla_{\mathbf{a}} \log \pi \nabla_{\mathbf{a}} Q^\pi]$  and the NPG mean step  $F_{\theta_\mu}^{-1} \mathbb{E}[\nabla_{\theta_\mu} \log \pi A^\pi]$   
 180 are collinear for each  $s$ , and hence globally after averaging over  $d^\pi$ .

## 181 5 Theoretical Discussion

182 **Equivalence vs. departure.** Proposition 3 formalizes alignment with natural policy gradient for  
 183 Gaussian means with fixed covariance under a local quadratic assumption on  $Q^\pi$ . More generally,  
 184 in exponential families with sufficient statistics affine in action and compatible parameterizations,  
 185 the WPO and NPG directions can be collinear per state. Departures arise for non-Gaussian families  
 186 (e.g., mixtures, tanh-squashed Gaussians) through the mixed derivative  $\nabla_\theta \nabla_{\mathbf{a}} \log \pi$ , which reflects  
 187 policy-manifold geometry and can rotate the update relative to NPG.

188 **Non-Gaussian example.** For a two-component Gaussian mixture with shared covariance, com-  
 189 ponent responsibilities  $\phi_i(\mathbf{a}, s) = \frac{\rho_i \mathcal{N}(\mathbf{a}|\mu_i, \Sigma)}{\sum_k \rho_k \mathcal{N}(\mathbf{a}|\mu_k, \Sigma)}$  enter  $\nabla_\theta \nabla_{\mathbf{a}} \log \pi$  and weight the cross-moment  
 190  $\mathbb{E}[\phi_i(\mathbf{a}, s) \nabla_{\mathbf{a}} Q(\mathbf{a} - \mu_i)^\top]$ . These curvature terms change the step direction even after Fisher pre-  
 191 conditioning.

## 192 6 Related Work

193 Policy gradient methods include REINFORCE [Williams, 1992], natural policy gradient [Kakade,  
 194 2001, Pascanu and Bengio, 2013], deterministic policy gradients [Silver et al., 2014], and successors  
 195 such as DDPG [Lillicrap et al., 2015]. Wasserstein geometry has informed optimization and learning  
 196 [Ambrosio et al., 2008, Benamou and Brenier, 2000], and several works explored Wasserstein in RL:  
 197 robust formulations [Abdullah et al., 2019], natural gradients [Moskovitz et al., 2020], and policy  
 198 optimization views [Zhang et al., 2018]. MPO [Abdolmaleki et al., 2018] relates as a policy iteration  
 199 method with KL-regularized exponentiated targets. Score-based transport (SVGD) [Liu and Wang,  
 200 2016] and its policy variants differ in using score-driven particle flows rather than density flows plus  
 201 projection. For covariance updates, affine-invariant geometry on SPD matrices [Absil et al., 2008]  
 202 and second-order approximations like K-FAC [Martens and Grosse, 2015] provide complementary  
 203 perspectives.

## 204 7 Limitations

205 Our energy and dissipation statements rely on the frozen-critic regime and matching state weight-  
 206 ing in the projection metric; fully coupled policy–critic dynamics introduce residual terms we do  
 207 not analyze. Non-smooth architectures may require weak-form interpretations. We do not ana-  
 208 lyze the variance/bias of mixed-derivative estimators beyond qualitative remarks. The Fisher vs.  
 209  $W_2$  parametric projection gap is only closed under affine-in-action exponential families. Extend-  
 210 ing equivalence results beyond Gaussians and characterizing regularity under which c-Wasserstein  
 211 squashing preserves optimality are important directions.

## 212 8 Conclusion

213 We clarified the WPO foundation with label-anchored derivations: Fisher–Galerkin projection to  
 214 a parametric update, baseline invariance, parameterization covariance, c-Wasserstein stability, and  
 215 Gaussian family updates. Future work includes extending equivalence conditions beyond Gaussians  
 216 and analyzing stronger regularity requirements for non-smooth architectures.

## 217 A Derivations and Assumptions

### 218 A.1 Convention (Rescaled Flow)

219 We work throughout with a per-state time-rescaled flow: the Eulerian velocity uses  $v = \nabla_{\mathbf{a}} Q^\pi$   
 220 and the functional derivative is  $\delta \mathcal{J} / \delta \pi = -Q^\pi$ . The state weighting  $d^\pi(s)$  enters only as an outer  
 221 expectation over  $s$ . All projection inner products and Fisher expectations use the same  $d^\pi(s)$ , so the  
 222 projected direction is invariant under this rescaling. Using the unrescaled convention instead would  
 223 multiply the fiberwise energy identities by  $d^\pi(s)^2$  inside the state integral.

### 224 A.2 Energy Dissipation (Proofs)

225 **Weak form.** Fix a state  $s$ . A curve  $t \mapsto \pi_t(\cdot | s)$  is a weak solution of the continuity equation  
 226  $\partial_t \pi = -\nabla_{\mathbf{a}} \cdot (\pi v)$  if, for every  $\varphi \in C_c^\infty(\mathcal{A})$ , the map  $t \mapsto \int \varphi \pi_t d\mathbf{a}$  is absolutely continuous and

$$\frac{d}{dt} \int \varphi \pi_t d\mathbf{a} = \int \nabla_{\mathbf{a}} \varphi(\mathbf{a}) \cdot v_t(\mathbf{a}) \pi_t(\mathbf{a}) d\mathbf{a}.$$

227 Assume either vanishing flux at infinity or a no-flux boundary condition on  $\partial \mathcal{A}$ .

228 **Function spaces.** Assume  $Q^\pi(\cdot, s) \in C^1$ ,  $\nabla_{\mathbf{a}} Q^\pi \in L^2(\pi_t(\cdot | s))$  uniformly in  $t$ , and  $\pi_t v_t \in L^1$ .  
 229 Fubini/Tonelli applies to interchange the  $s$  and  $\mathbf{a}$  integrals under  $d^\pi(s) \pi_t(\mathbf{a} | s)$ .

230 **W<sub>2</sub> case.** With the rescaled convention  $v = \nabla_{\mathbf{a}} Q^\pi$  and  $\delta \mathcal{J} / \delta \pi = -Q^\pi$ , take  $\varphi = \delta \mathcal{J} / \delta \pi$  as a  
 231 test function and use the chain rule (justified by dominated convergence under the  $L^2$  bounds):

$$\frac{d}{dt} \mathcal{J}[\pi_t] = \int \frac{\delta \mathcal{J}}{\delta \pi} \partial_t \pi_t d\mathbf{a} = - \int \frac{\delta \mathcal{J}}{\delta \pi} \nabla_{\mathbf{a}} \cdot (\pi_t v_t) d\mathbf{a} = \int \nabla_{\mathbf{a}} \left( \frac{\delta \mathcal{J}}{\delta \pi} \right) \cdot v_t \pi_t d\mathbf{a},$$

232 where the boundary term vanishes. Substituting  $\nabla_{\mathbf{a}}(\delta \mathcal{J} / \delta \pi) = -\nabla_{\mathbf{a}} Q^\pi$  and  $v_t = \nabla_{\mathbf{a}} Q^\pi$  gives  
 233  $\frac{d}{dt} \mathcal{J}[\pi_t] = - \int \pi_t \|\nabla_{\mathbf{a}} Q^\pi\|^2 d\mathbf{a}$ . Taking the expectation over  $s \sim d^\pi$  yields Lemma 1.

234 **c-Wasserstein case.** Let  $c$  be proper, convex, l.s.c., with conjugate  $c^*$  differentiable and  $\nabla c^*$   
 235 monotone. With  $v_t = \nabla c^*(-\nabla_{\mathbf{a}} \delta \mathcal{J} / \delta \pi)$  the same argument yields

$$\frac{d}{dt} \mathcal{J}[\pi_t] = \int \nabla_{\mathbf{a}} \left( \frac{\delta \mathcal{J}}{\delta \pi} \right) \cdot \nabla c^* \left( -\nabla_{\mathbf{a}} \left( \frac{\delta \mathcal{J}}{\delta \pi} \right) \right) \pi_t d\mathbf{a} = - \int \langle \nabla c^*(\nabla_{\mathbf{a}} Q^\pi), \nabla_{\mathbf{a}} Q^\pi \rangle \pi_t d\mathbf{a} \leq 0,$$

236 with equality only when  $\nabla_{\mathbf{a}} Q^\pi = 0$  almost everywhere (or when  $\nabla c^*(\cdot) = 0$  at that argument).

### 237 A.3 Assumptions and Boundary Conditions

238 We assume  $C^1/C^2$  smoothness as required and vanishing boundary terms in integration by parts;  
 239 interchange of expectation and differentiation follows from dominated convergence under integrable  
 240 bounds. We assume either (i)  $\mathcal{A} = \mathbb{R}^d$  with tails making the flux vanish at infinity, or (ii) bounded  
 241  $\mathcal{A}$  with no-flux boundary condition  $(\pi v) \cdot n = 0$  on  $\partial \mathcal{A}$ . We assume  $Q^\pi(\cdot, s) \in C^1$  in  $\mathbf{a}$  with  $\nabla_{\mathbf{a}} Q^\pi$   
 242 locally Lipschitz (uniformly in  $s$ ), ensuring the weak formulation and energy identity.

243 Sufficient boundary/decay conditions include: (i) bounded  $\mathcal{A}$  with  $(\pi v) \cdot n = 0$ ; or (ii)  $\mathcal{A} = \mathbb{R}^d$  and  
 244  $\pi(\mathbf{a} | s) \|\nabla_{\mathbf{a}} Q^\pi(s, \mathbf{a})\| \rightarrow 0$  as  $\|\mathbf{a}\| \rightarrow \infty$  (uniformly in  $s$ ).

245 Scope of descent (frozen critic). Energy decay statements apply to the proxy functional with  $Q^\pi$   
 246 and  $d^\pi$  held fixed during the inner flow/projection step (semi-gradient setting). They do not by  
 247 themselves imply monotone improvement of the true return when the critic or occupancy evolves  
 248 between steps.

### 249 A.4 Fisher–Galerkin Projection

250 Minimizing the Fisher-weighted squared error between  $\partial_t \pi$  and  $\delta \pi_\theta$  yields the orthogonality con-  
 251 ditions  $\langle \partial_t \pi - \delta \pi_\theta, \pi \nabla_{\theta_k} \log \pi \rangle = 0$  for each coordinate  $k$ , i.e.,  $F_{\theta\theta} \Delta \theta = \mathcal{F}_{t\theta}$  with  $F_{\theta\theta} =$   
 252  $\mathbb{E}[\nabla_\theta \log \pi \nabla_\theta \log \pi^\top]$ . Using the continuity equation and integrating by parts in  $\mathbf{a}$  (boundary terms  
 253 vanish),

$$\mathcal{F}_{t\theta_k} = \mathbb{E} \left[ \int \nabla_{\theta_k} \log \pi \partial_t \pi d\mathbf{a} \right] = -\mathbb{E} \left[ \int \nabla_{\theta_k} \log \pi \nabla_{\mathbf{a}} \cdot (\pi \nabla_{\mathbf{a}} Q) d\mathbf{a} \right] = \mathbb{E} \left[ \int (\nabla_{\mathbf{a}} \nabla_{\theta_k} \log \pi)^\top (\pi \nabla_{\mathbf{a}} Q) d\mathbf{a} \right],$$

254 which gives  $\mathcal{F}_{t\theta} = \mathbb{E}[\nabla_\theta \nabla_{\mathbf{a}} \log \pi \nabla_{\mathbf{a}} Q]$ . Commutation of  $\nabla_\theta$  with the integral follows by domi-  
 255 nated convergence if  $\nabla_\theta \nabla_{\mathbf{a}} \log \pi \in L^1$  uniformly.

### 256 A.5 Gâteaux Variation and Baselines

257 Consider the Lagrangian  $\mathcal{L}(\pi, \lambda) = -\mathbb{E}_s \mathbb{E}_{\mathbf{a} \sim \pi(\cdot | s)}[Q^\pi(s, \mathbf{a})] + \mathbb{E}_s[\lambda(s)(\int \pi(\mathbf{a} | s) d\mathbf{a} - 1)]$ . The  
 258 first variation in a direction  $h$  with  $\int h d\mathbf{a} = 0$  per state is  $\delta \mathcal{L}(\pi; h) = -\mathbb{E}_s \int Q^\pi(s, \mathbf{a}) h(\mathbf{a}, s) d\mathbf{a}$ .  
 259 Stationarity yields  $\delta \mathcal{J} / \delta \pi = -Q^\pi + \lambda(s)$ . Replacing  $Q^\pi$  by  $Q^\pi - b(s)$  shifts  $\lambda(s)$  to  $\lambda(s) + b(s)$ ;  
 260 since  $\nabla_{\mathbf{a}} b(s) = 0$ , both the Eulerian velocity and the cross term  $\mathcal{F}_{t\theta}$  are unchanged.

### 261 A.6 Gaussian Updates

262 For 1D Gaussian,  $\partial_\mu \log \pi = (a - \mu) / \sigma^2$  and  $\partial_\lambda \log \pi = (a - \mu)^2 / \sigma^2 - 1$  with  $\lambda = \log \sigma$ .  
 263 Fisher blocks are  $F_{\mu\mu} = 1 / \sigma^2$ ,  $F_{\lambda\lambda} = 2$ ,  $F_{\mu\lambda} = 0$ . The natural step gives  $\Delta \mu = \mathbb{E}[\nabla_a Q]$  and  
 264  $\Delta \lambda = \frac{1}{\sigma^2} \mathbb{E}[(a - \mu) \nabla_a Q]$ .

265 For full covariance, work on  $\text{SPD}(d)$  with the affine-invariant inner product  $\langle U, V \rangle_\Sigma =$   
 266  $\text{tr}(\Sigma^{-1} U \Sigma^{-1} V)$ . Let  $M = \mathbb{E}[(\nabla_{\mathbf{a}} Q)(\mathbf{a} - \mu)^\top]$ . A standard calculation gives the Riemannian

267 gradient  $\text{grad } \mathcal{J}(\Sigma) = \text{sym}(\Sigma^{-1} M \Sigma^{-1})$ . The natural step is  $\Delta \Sigma = -\alpha \text{grad } \mathcal{J}$  pulled back to  
 268 Euclidean coordinates, i.e.,

$$\Delta \Sigma = 2 \Sigma \text{sym}(\Sigma^{-1} M \Sigma^{-1}) \Sigma = M + M^\top,$$

269 which matches the Fisher-preconditioned update.

## 270 A.7 Cholesky Triangular Sylvester Solve

271 Linearizing  $\Sigma = LL^\top$  yields  $\Delta \Sigma = L \Delta L^\top + \Delta L L^\top$ . Given a symmetric target  $S$ , the  
 272 lower-triangular recursion stated in the main text solves  $L \Delta L^\top + \Delta L L^\top = S$  uniquely by forward  
 273 substitution (prove by induction on columns). Hence the realized increment equals  $S = M + M^\top$ .  
 274 For SPD preservation, a sufficient condition is to backtrack  $\eta$  until  $\min_i (L + \eta \Delta L)_{ii} > 0$ ; for  
 275 example, if  $\eta \|L^{-1} \Delta L\|_\infty < 1$ , the updated diagonal remains positive.  
 276

277 **Algorithm (Cholesky SPD update).** Given  $L$  and a symmetric target  $S = M + M^\top$ :

- 278 1. Solve  $L \Delta L^\top + \Delta L L^\top = S$  for lower-triangular  $\Delta L$  by forward substitution (column-  
 279 wise recursion above).
- 280 2. Line search  $\eta > 0$  (e.g., backtracking Armijo) until  $\min_i (L + \eta \Delta L)_{ii} > 0$ .
- 281 3. Update  $L \leftarrow L + \eta \Delta L$  and return  $\Sigma \leftarrow LL^\top$ .

282 Cost:  $\mathcal{O}(d^3)$  per update; diagonal  $\Sigma$  reduces to  $\mathcal{O}(d)$ .

## 283 A.8 Parameterization Covariance

284 Let  $\phi = \phi(\theta)$  be a local diffeomorphism with Jacobian  $J = \partial \phi / \partial \theta$ . The Fisher transforms as  $F_\phi =$   
 285  $\mathbb{E}[\nabla_\phi \log \pi \nabla_\phi \log \pi^\top] = J^{-\top} F_\theta J^{-1}$ . The cross term obeys  $g_\phi = \mathbb{E}[\nabla_\phi \nabla_a \log \pi \nabla_a Q] = J^{-\top} g_\theta$   
 286 by the chain rule  $\nabla_\phi = J^{-\top} \nabla_\theta$ . Therefore  $\Delta \phi = F_\phi^{-1} g_\phi = J F_\theta^{-1} g_\theta = J \Delta \theta$ .

## 287 A.9 c-Wasserstein

288 Assume  $c$  is proper, convex, l.s.c., with conjugate  $c^*$  differentiable and  $\nabla c^*$  monotone. The  
 289 c-Wasserstein flow uses velocity  $\nabla c^*(-\nabla_a \delta \mathcal{J} / \delta \pi)$ . The elementwise cube-root mapping used  
 290 in the main text arises from the *separable* choice  $c(u) = \frac{1}{4} \sum_i |u_i|^4$ , whose conjugate has  
 291  $c^*(x) = \frac{3}{4} \sum_i |x_i|^{4/3}$  and  $\nabla c^*(x)_i = \text{sign}(x_i) |x_i|^{1/3}$ . This map is direction-preserving and Hölder-  
 292 continuous, supporting the stated stability. A rotationally invariant alternative  $c(u) = \frac{1}{4} \|u\|_2^4$  yields  
 293  $\nabla c^*(x) = \|x\|_2^{-2/3} x$ .

## 294 B Proofs of Main Results

295 *Proof of Theorem 2.* Projecting  $\partial_t \pi$  onto  $\delta \pi_\theta$  in the Fisher inner product yields normal equations  
 296  $\langle \partial_t \pi - \delta \pi_\theta, \pi \nabla_\theta \log \pi \rangle = 0$ , which rearrange to  $F_{\theta\theta} \Delta \theta = \mathcal{F}_{t\theta}$ . Integration by parts under the  
 297 stated regularity gives the cross term in (4).  $\square$

298 *Proof of Proposition 1.* Constrained variation with per-state normalization introduces a Lagrange  
 299 multiplier  $\lambda(s)$ , yielding  $\delta \mathcal{J} / \delta \pi = -Q^\pi + \lambda(s)$ . Any baseline  $b(s)$  can be absorbed into  $\lambda$ , so both  
 300 the PDE and  $\mathcal{F}_{t\theta}$  are invariant to  $Q^\pi \mapsto Q^\pi - b(s)$ .  $\square$

301 *Proof of Theorem 3.* Under the pullback metric,  $F_\phi^{-1} = J F_\theta^{-1} J^\top$  and  $g_\phi = J^{-\top} g_\theta$ . Thus  
 302  $F_\phi^{-1} g_\phi = J F_\theta^{-1} g_\theta = J \Delta \theta$ .  $\square$

303 *Proof of Lemma 2.* Use  $\nabla_a \log \pi = -(a - \mu) / \sigma^2$  and  $\nabla_\theta \nabla_a \log \pi = (\nabla_\theta \mu) / \sigma^2$  to compute  $g_\theta$  and  
 304 Fisher blocks; apply  $\Delta \theta = F_\theta^{-1} g_\theta$  to obtain the stated updates. For full  $\Sigma$ , the affine-invariant Fisher  
 305 yields  $\Delta \Sigma = M + M^\top$ .  $\square$



306 *Proof of Proposition 2.* For convex  $c$ ,  $c^*$  is convex and  $\nabla c^*$  is monotone. Replacing  $\nabla_{\mathbf{a}} Q$  by  
307  $\nabla c^*(\nabla_{\mathbf{a}} Q)$  corresponds to the Eulerian form of the c-Wasserstein gradient flow, preserving de-  
308 scent/ascent while taming large gradients.  $\square$

309 *Proof of Proposition 3. Proof sketch.* For fixed covariance, the mean-block Fisher is  $F_{\theta_\mu, \theta_\mu} =$   
310  $\mathbb{E}[\nabla_{\theta_\mu} \log \pi \nabla_{\theta_\mu} \log \pi^\top]$ . Under the local quadratic model  $\nabla_{\mathbf{a}} Q^\pi \approx H(s)(\mathbf{a} - \mu)$  and Gaussian  $\pi$ ,  
311 one has  $\mathbb{E}[\nabla_{\mathbf{a}} Q^\pi (\mathbf{a} - \mu)^\top] = H(s) \Sigma$ . Using  $\nabla_{\theta_\mu} \nabla_{\mathbf{a}} \log \pi = (\partial \mu / \partial \theta_\mu) \Sigma^{-1}$ , the WPO cross term  
312 becomes  $g_{\theta_\mu} = \mathbb{E}[(\partial \mu / \partial \theta_\mu) \Sigma^{-1} H(s) \Sigma]$ . The NPG mean direction with advantage  $A^\pi$  linearized  
313 around  $\mu$  yields the same factor  $H(s)$  multiplying  $\nabla_{\theta_\mu} \log \pi$ , so after preconditioning by  $F_{\theta_\mu, \theta_\mu}^{-1}$  both  
314 steps are collinear. Averaging over  $d^\pi$  preserves collinearity.  $\square$

315 **Acknowledgments** Omitted for double-blind review.

## 316 Responsible AI Statement

317 This is a theoretical contribution. It contains no human or animal subjects, no personal or sensi-  
318 tive data, and no deployed systems. All results are formal statements with explicit assumptions  
319 and proofs; we discuss limitations, scope, and failure modes where applicable. The work adheres  
320 to the Agents4Science Code of Ethics (and the NeurIPS Code of Ethics in spirit): we avoid pro-  
321 hibited practices, dual-use concerns, and undisclosed human subject data; environmental impact is  
322 negligible as no large-scale compute is used.

## 323 Reproducibility Statement

324 All claims are formal theorems, propositions, or lemmas, each with clearly stated assumptions and  
325 proofs. Key equations are label-anchored for unambiguous cross-referencing (e.g., (3), (1), (2), (4)).  
326 Weak-form and boundary assumptions are made explicit; derivations are given both in the main text  
327 and in the appendix. No datasets or empirical experiments are involved. A minimal prototype used  
328 in separate notes mirrors the Gaussian mean/covariance updates but is not required to verify the  
329 theoretical results.

## 330 References

- 331 Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin Ried-  
332 miller. Maximum a posteriori policy optimisation. In *International Conference on Learning Representations*,  
333 2018.
- 334 Mohammed Amin Abdullah, Hang Ren, Haitham Bou Ammar, Vladimir Milenkovic, Rui Luo, Mingtian  
335 Zhang, and Jun Wang. Wasserstein robust reinforcement learning, 2019.
- 336 P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University  
337 Press, 2008.
- 338 Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of*  
339 *probability measures*. Springer, 2008.
- 340 Jean-David Benamou and Yann Brenier. A computational fluid mechanics solution to the monge-kantorovich  
341 mass transfer problem. *Numerische Mathematik*, 84(3):375–393, 2000.
- 342 Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential  
343 equations. *Advances in Neural Information Processing Systems*, 31, 2018.
- 344 Sham M Kakade. A natural policy gradient. In *Advances in Neural Information Processing Systems*, 2001.
- 345 Timothy Lillicrap et al. Continuous control with deep reinforcement learning. *arXiv preprint*  
346 *arXiv:1509.02971*, 2015.
- 347 Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm.  
348 In *Advances in Neural Information Processing Systems*, 2016.

- 349 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature.  
350 In *International Conference on Machine Learning*. PMLR, 2015.
- 351 Ted Moskovitz, Michael Arbel, Ferenc Huszar, and Arthur Gretton. Efficient wasserstein natural gradients for  
352 reinforcement learning. *arXiv preprint arXiv:2010.05380*, 2020.
- 353 Razvan Pascanu and Yoshua Bengio. Revisiting natural gradient for deep networks. *arXiv preprint*  
354 *arXiv:1301.3584*, 2013.
- 355 David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic  
356 policy gradient algorithms. In *International Conference on Machine Learning*, pages 387–395. PMLR, 2014.
- 357 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning.  
358 *Machine learning*, 8:229–256, 1992.
- 359 Ruiyi Zhang, Changyou Chen, Chunyuan Li, and Lawrence Carin. Policy optimization as wasserstein gradient  
360 flows. In *International Conference on Machine Learning*, pages 5737–5746. PMLR, 2018.

## Agents4Science AI Involvement Checklist

1. **Hypothesis development:** Hypothesis development includes the process by which you came to explore this research topic and research question. This can involve the background research performed by either researchers or by AI. This can also involve whether the idea was proposed by researchers or by AI.

Answer: [A]

Explanation: The theoretical questions, hypotheses, and scope (projection of Wasserstein flows, invariance properties, and c-Wasserstein stability) were developed by the authors without AI assistance beyond standard search/reading tools.

2. **Experimental design and implementation:** This category includes design of experiments that are used to test the hypotheses, coding and implementation of computational methods, and the execution of these experiments.

Answer: [NA]

Explanation: The paper presents mathematical derivations and proofs only; there are no empirical experiments or released implementations.

3. **Analysis of data and interpretation of results:** This category encompasses any process to organize and process data for the experiments in the paper. It also includes interpretations of the results of the study.

Answer: [NA]

Explanation: No datasets were used; results are theoretical statements with formal proofs and stated assumptions.

4. **Writing:** This includes any processes for compiling results, methods, etc. into the final paper form. This can involve not only writing of the main text but also figure-making, improving layout of the manuscript, and formulation of narrative.

Answer: [A]

Explanation: Writing and editing were performed by the authors using LaTeX. No generative AI systems were used to draft or edit the scientific content.

5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or lead author?

Description: Not applicable (no AI systems were used in ideation, analysis, coding, or writing).

## Agents4Science Paper Checklist

### 1. Claims

Question: Are the claims made in the abstract and introduction supported by the results?

Answer: [Yes]

Justification: Claims (projection to natural gradient, baseline invariance, parameterization covariance, Gaussian/covariance updates, and c-Wasserstein energy decay) are supported by Theorems/Propositions with proofs and stated assumptions (see Theorem 2, Proposition 1, Theorem 3, Lemma 2, Theorem 1).

### 2. Tasks and baselines

Question: Did you describe the limitations of your work?

Answer: [Yes]

Justification: A dedicated Limitations section enumerates scope and assumptions (e.g., frozen-critic regime, smoothness/weak-form conditions, and when equivalences fail).

Question: Did you discuss any potential negative societal impacts of your work?

Answer: [NA]

Justification: The work is purely theoretical and does not involve deployment, datasets, or application domains with direct societal impact.

Question: Did you discuss the failure modes of your method?

Answer: [Yes]

Justification: The paper discusses conditions where alignment with NPG holds vs. fails (e.g., mixtures/squashing), boundary/regularity requirements, and stability caveats beyond the frozen-critic setting.

### 3. Reproducibility

Question: Is your code and data publicly available?

Answer: [NA]

Justification: No code or datasets were produced; the paper contains only theoretical results.

Question: Does the paper provide sufficient details to reproduce the main results (either in the text or as supplemental material)?

Answer: [Yes]

Justification: Full derivations with assumptions, weak-form statements, and proofs are provided in the main text and appendix.

### 4. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: Not applicable; there are no experiments.

### 5. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Not applicable; there are no experiments.

### 6. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: Not applicable; there are no experiments.

### 7. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the Agents4Science Code of Ethics (see conference website)?

Answer: [Yes]

441 Justification: Theoretical work with no human/animal subjects, sensitive data, or prohibited practices.  
442 8. **Broader impacts**  
443 Question: Does the paper discuss both potential positive societal impacts and negative societal im-  
444 pacts of the work performed?  
445 Answer: [NA]  
446 Justification: Not applicable to this purely theoretical contribution; no deployment or application  
447 domain is studied.