



AdamO: A Collapse-Suppressed Optimizer for Offline RL

Nan Qiao, Sheng Yue, Shuning Wang, Ju Ren



ABSTRACT

Offline RL critics fail spectacularly when bootstrapped temporal-difference (TD) updates amplify their own errors, driving Q-values to extreme, unusable magnitudes. Prior work blames the backup rule, the network architecture, or the squared-TD loss. **We show a more basic cause: the optimizer dynamics themselves can trigger or suppress collapse.** Treating offline TD learning as a feedback control system, we analyze Adam-based critic updates and derive a necessary and sufficient stability condition: the linearized critic is stable **iff** the spectral radius of an explicit augmented update operator stays strictly below one. A tractable sufficient form of this condition splits into a scale term (handled by standard normalization) and a geometry term governed by how far the weights deviate from isometry — motivating parameter orthogonality. Because a loss-based orthogonality penalty contaminates Adam's adaptive moments, orthogonality must instead be enforced as a **decoupled, optimizer-level correction**. We propose **AdamO**: Adam plus a budgeted orthogonality drift on selected weight blocks. AdamO provably never harms worst-case task descent and preserves Adam's dissipative continuous-time dynamics, yet is a drop-in replacement — just swap the critic's optimizer. Across D4RL it consistently improves stability and returns, with the largest gains where Adam collapses outright.

MOTIVATION

Minimizing the second moment of the temporal difference (TD) residual δ under the dataset distribution, i.e. $\min \mathbb{E}[\delta^2]$.

$$\delta \triangleq r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\cdot | \mathbf{s}')} Q_{\phi'}(\mathbf{s}', \mathbf{a}') - Q_{\phi}(\mathbf{s}, \mathbf{a})$$

Why Does Adam Trigger Value Collapse?

Adam is not "fast SGD." Its momentum + variance adaptation turn TD learning into a second-order difference equation, so the right stability object is a spectral radius, not SGD intuition.

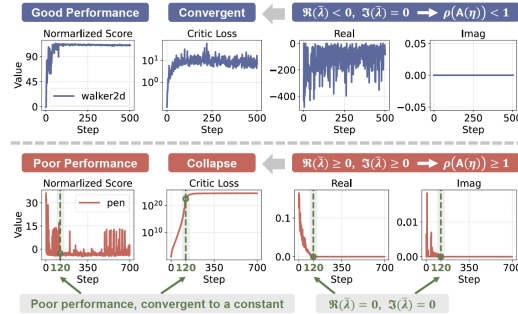
Linearized TD-error dynamics (frozen/terminal regime). With the Adam-preconditioned Gram operator $\mathbf{K}(X_1, X_2) \triangleq Z(X_1)^T D Z(X_2)$, the bootstrapped update propagates through the greedy targets, giving the TD dynamics operator

$$\mathbf{S} \triangleq \gamma \mathbf{K}(X', X) - \mathbf{K}(X, X)$$

In supervised learning, \mathbf{S} is negative semidefinite, the symmetric, marginally stable case of the Hurwitz condition, so no positive feedback, **no collapse**.

Coupling this with Adam's momentum EMA yields a closed-form second-order recurrence, summarized by the augmented operator

$$\mathbf{A}(\eta) \triangleq \begin{bmatrix} (1 + \beta_1)I + \eta(1 - \beta_1)\mathbf{S} & -\beta_1 I \\ I & 0 \end{bmatrix}$$



Theorem 4.1 (necessary & sufficient). The frozen linearized TD error decays to 0 iff $\rho(\mathbf{A}(\eta)) < 1$.

Theorem 4.2 + Remark 4.3 (mechanism). If \mathbf{S} is Hurwitz, a small enough stepsize guarantees $\rho < 1$. Collapse occurs exactly when bootstrapping turns into **positive feedback** that amplifies TD errors faster than the optimizer can damp them; at the boundary $(\Re(\lambda) \rightarrow 0)$ a unit root appears and training **stagnates** instead of recovering.

So the entire fix reduces to one question answered next: how do we keep \mathbf{S} Hurwitz? That is

How to Suppress It Without Breaking Adam?

A tractable sufficient condition for Hurwitzness (Prop. 5.1) cleanly separates two effects. With $\Phi = D^{1/2}Z(X)$, $\Phi_* = D^{1/2}Z(X')$:

$$\gamma \|\Phi\|_2 \|\Phi_*\|_2 + \|\Phi^T \Phi - I\|_2 < 1 \implies \mathbf{S} \text{ is Hurwitz.}$$

(i) bootstrapped scale (ii) feature geometry

- (i) **Scale term** \rightarrow already controlled by ubiquitous practice: input normalization / clipping + spectral-norm constraints.
- (ii) **Geometry term** \rightarrow reduces to a *parameter-only* near-isometry defect $\varepsilon = \|\Psi^T \Psi - I\|_2$, i.e. how far the weight blocks deviate from orthonormal. **This is the missing knob.**

Why not just add an orthogonality penalty to the loss? Running Adam on $\nabla \tilde{L}_t = g_t + \lambda r_t$ feeds the orthogonality gradient into Adam's *second-moment* recursion

$$v_{t+1} = \beta_2 v_t + (1 - \beta_2)(g_t + \lambda r_t)^{\otimes 2},$$

so λr_t leaks into (m_t, v_t) and **distorts future task steps even after it vanishes** (the same pathology AdamW fixes for weight decay).

PRACTICAL ALGORITHM

AdamO keeps Adam's moments driven only by the task gradient, and adds a per-layer orthogonality correction that is (a) scale-matched to the Adam step, (b) capped by a task-alignment budget, and (c) exactly Adam when $\kappa=0$.

For each constrained weight block W (regularizer $R(W) = \frac{1}{4} \|W^T W - I\|_F^2$, closed-form gradient r_t):

$$\text{scale-match: } \delta_{t,0} = \kappa \frac{\|u_t\|_F}{\|r_t\|_F + \varepsilon_r} r_t, \quad \|\delta_t\| \leq \kappa \|u_t\|.$$

$$\text{budget: } (g_t, \delta_t)_F \geq -(g_t, u_t)_F, \quad \tau \in [0, 1],$$

with the largest feasible step taken in closed form (Eq. 19), and finally

$$\omega_{t+1} = \omega_t - \eta(u_t + \delta_t)$$

Algorithm 1 AdamO: Adam with Orthogonality Correction

- 1: Init $\omega_0, m_0 = v_0 = 0, \eta, \kappa, \tau$ and Adam params
- 2: for $t = 0, 1, \dots$ do
- 3: $g_t \leftarrow \nabla L_t(\omega_t)$
- 4: Compute standard Adam update u_t
- 5: Compute δ_t by applying (17)–(19) layer-wise **Orthogonality Correction**
- 6: $\omega_{t+1} \leftarrow \omega_t - \eta(u_t + \delta_t)$
- 7: end for

EXPERIMENTAL RESULTS

Scarce-data regime: D4RL with 10k state–action pairs. AdamO is dropped into six offline RL algorithms, TD3+BC, IQL, ReBRAC, ACTIVE, PARS, SQOG, replacing only the critic optimizer.

Task Name	TD3+BC		IQL		ReBRAC		ACTIVE		PARS		SQOG	
	Adam	AdamO	Adam	AdamO	Adam	AdamO	Adam	AdamO	Adam	AdamO	Adam	AdamO
AntMaze-umaze	72.2	92.5 _{±0.0%}	80.5	83.8 _{±0.6%}	94.3	96.5 _{±0.6%}	92.4	95.1 _{±0.6%}	97.3	93.5 _{±0.1%}	89.6	93.1 _{±0.1%}
AntMaze-umaze-div	47.0	82.2 _{±0.0%}	55.8	68.2 _{±0.0%}	87.2	91.5 _{±0.0%}	75.9	78.2 _{±0.0%}	93.2	92.4 _{±0.0%}	72.8	87.1 _{±0.0%}
AntMaze-med-play	0.3	28.5 _{±0.0%}	70.4	70.1 _{±0.4%}	85.2	89.4 _{±0.0%}	73.7	86.5 _{±0.0%}	91.5	92.8 _{±0.0%}	60.9	65.4 _{±0.0%}
AntMaze-med-div	0.2	16.4 _{±0.0%}	66.9	75.5 _{±0.0%}	79.8	85.1 _{±0.0%}	77.8	82.3 _{±0.0%}	87.1	90.6 _{±0.0%}	65.8	81.2 _{±0.0%}
AntMaze-large-play	0.0	13.7 _{±0.0%}	38.5	41.8 _{±0.0%}	55.4	61.8 _{±0.0%}	59.9	56.5 _{±0.0%}	46.6	52.9 _{±0.0%}	32.6	37.4 _{±0.0%}
AntMaze-large-div	0.0	16.5 _{±0.0%}	32.4	36.1 _{±0.0%}	59.5	56.2 _{±0.0%}	46.6	61.8 _{±0.0%}	50.8	56.5 _{±0.0%}	47.5	52.8 _{±0.0%}
AntMaze-ultra-div	0.0	2.4 _{±0.0%}	19.0	31.1 _{±0.0%}	5.7	9.4 _{±0.0%}	10.0	19.8 _{±0.0%}	42.1	48.6 _{±0.0%}	0.0	4.2 _{±0.0%}
AntMaze-ultra-play	0.0	1.8 _{±0.0%}	21.0	29.9 _{±0.0%}	20.6	26.8 _{±0.0%}	13.2	11.5 _{±0.0%}	55.9	62.4 _{±0.0%}	0.0	2.1 _{±0.0%}
AntMaze Avg.	15.0	31.8 _{±0.0%}	48.1	54.6 _{±0.0%}	59.8	64.6 _{±0.0%}	55.1	61.4 _{±0.0%}	70.6	73.7 _{±0.0%}	48.7	55.4 _{±0.0%}
HalfCheetah-m	35.9	53.5 _{±0.0%}	29.7	35.2 _{±0.0%}	42.8	49.4 _{±0.0%}	42.9	48.5 _{±0.0%}	44.9	51.2 _{±0.0%}	36.5	41.8 _{±0.0%}
HalfCheetah-mr	39.1	46.2 _{±0.0%}	32.7	38.4 _{±0.0%}	40.7	47.1 _{±0.0%}	34.9	40.3 _{±0.0%}	45.4	50.8 _{±0.0%}	30.2	35.5 _{±0.0%}
HalfCheetah-me	33.5	60.1 _{±0.0%}	48.1	54.6 _{±0.0%}	63.9	71.5 _{±0.0%}	55.7	62.3 _{±0.0%}	75.7	82.9 _{±0.0%}	58.9	65.4 _{±0.0%}
Hopper-m	40.7	75.5 _{±0.0%}	38.9	45.1 _{±0.0%}	78.6	86.2 _{±0.0%}	26.2	31.8 _{±0.0%}	73.7	80.4 _{±0.0%}	45.8	51.3 _{±0.0%}
Hopper-mr	21.3	55.8 _{±0.0%}	46.6	52.9 _{±0.0%}	64.2	70.8 _{±0.0%}	62.5	68.7 _{±0.0%}	67.5	74.1 _{±0.0%}	61.4	67.2 _{±0.0%}
Hopper-me	32.6	84.2 _{±0.0%}	66.5	73.2 _{±0.0%}	78.5	85.9 _{±0.0%}	62.7	69.1 _{±0.0%}	75.5	81.6 _{±0.0%}	72.9	78.5 _{±0.0%}
Walker2d-m	21.2	68.4 _{±0.0%}	54.9	61.5 _{±0.0%}	62.2	69.4 _{±0.0%}	53.6	59.2 _{±0.0%}	68.0	74.5 _{±0.0%}	49.8	55.6 _{±0.0%}
Walker2d-mr	19.3	52.5 _{±0.0%}	51.4	57.8 _{±0.0%}	69.4	76.1 _{±0.0%}	45.8	51.4 _{±0.0%}	63.4	69.2 _{±0.0%}	58.5	64.3 _{±0.0%}
Walker2d-me	22.4	98.5 _{±0.0%}	57.3	63.7 _{±0.0%}	74.0	80.8 _{±0.0%}	58.6	64.9 _{±0.0%}	81.6	88.3 _{±0.0%}	68.7	74.9 _{±0.0%}
Locomotion Avg.	29.6	66.1 _{±0.0%}	47.3	53.6 _{±0.0%}	63.8	70.8 _{±0.0%}	49.2	55.2 _{±0.0%}	66.2	72.6 _{±0.0%}	53.6	59.4 _{±0.0%}
Pen-human	-4.1	83.1 _{±0.0%}	79.1	89.8 _{±0.0%}	105.4	102.2 _{±0.0%}	106.2	109.8 _{±0.0%}	88.1	94.5 _{±0.0%}	77.0	75.0 _{±0.0%}
Pen-cloned	5.6	82.4 _{±0.0%}	46.5	82.2 _{±0.0%}	98.5	92.4 _{±0.0%}	96.5	100.2 _{±0.0%}	107.1	105.8 _{±0.0%}	73.6	87.4 _{±0.0%}
Door-human	-0.3	0.2 _{±0.0%}	3.5	4.8 _{±0.0%}	-0.1	0.1 _{±0.0%}	0.0	0.2 _{±0.0%}	0.1	0.3 _{±0.0%}	-0.1	0.1 _{±0.0%}
Door-cloned	-0.3	0.1 _{±0.0%}	3.3	4.5 _{±0.0%}	0.1	0.5 _{±0.0%}	0.1	0.4 _{±0.0%}	0.1	0.3 _{±0.0%}	0.1	0.5 _{±0.0%}
Hammer-human	1.1	0.5 _{±0.0%}	1.9	3.1 _{±0.0%}	0.3	0.8 _{±0.0%}	0.3	0.9 _{±0.0%}	0.3	0.8 _{±0.0%}	0.3	0.8 _{±0.0%}
Hammer-cloned	0.3	0.4 _{±0.0%}	1.7	3.2 _{±0.0%}	5.4	7.5 _{±0.0%}	5.9	7.8 _{±0.0%}	4.6	7.0 _{±0.0%}	5.1	6.8 _{±0.0%}
Relocate-human	-0.3	0.2 _{±0.0%}	0.1	0.5 _{±0.0%}	0.2	0.6 _{±0.0%}	0.2	0.7 _{±0.0%}	0.1	0.7 _{±0.0%}	0.2	0.6 _{±0.0%}
Relocate-cloned	0.1	0.0 _{±0.0%}	0.0	0.0 _{±0.0%}	0.4 _{±0.0%}	0.2 _{±0.0%}	0.1 _{±0.0%}	0.1 _{±0.0%}	0.1 _{±0.0%}	0.0 _{±0.0%}	-0.2 _{±0.0%}	0.6 _{±0.0%}
Adopt Avg.	0.3	20.9 _{±0.0%}	17.0	23.5 _{±0.0%}	26.2	25.6 _{±0.0%}	26.2	27.5 _{±0.0%}	25.1	26.2 _{±0.0%}	19.5	21.6 _{±0.0%}

