# Exponential weight averaging as damped harmonic motion

Jonathan Patsenker*[1], Henry Li*[1], Yuval Kluger[1]

[1]Applied Mathematics Program, Yale University
*Equal contribution; order decided by coin toss

## Motivations

◦ The exponential moving average (EMA) of neural network weights is a commonly used in deep learning optimization, especially in generative models

◦ EMA improves the stability of the inference model during and after training.

◦ Benefits *after* training have been studied

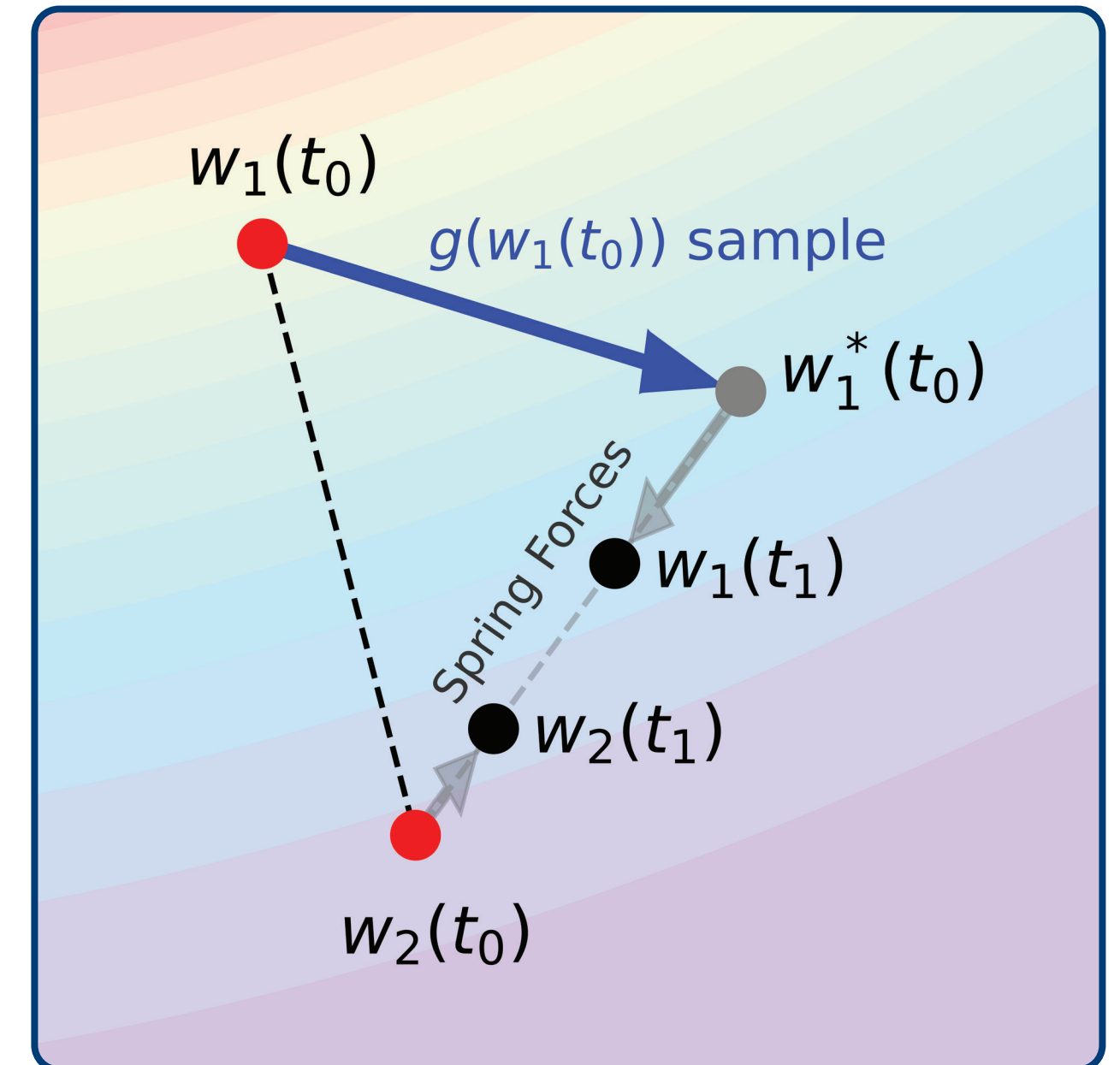◦ Benefits *during* training not well understood.



Fig 1. A visualization of the BELAY update step. The background color corresponds to the true full-batch loss function, and $g$ is sampled using an optimizer on a minibatch.

## BELAY: Physical EMA

Let $w_1$, $w_2$ represent point-particles with masses $m_1$, $m_2$, attached by a 0-length spring with spring constant k. The particles are subject to damping with constants $c_1$, $c_2$ respectively. External forces notated by $f(w_1,t)$ are exerted upon $w_1$ but not $w_2$. We break down the total forces ($F_1$,$F_2$) exerted on $w_1$, $w_2$.

$$F_1 = \underbrace{k(w_2 - w_1)}_{\text{Hookean}} - \underbrace{c_1 \dot{w}_1}_{\text{Damping}} + \underbrace{f(w_1, t)}_{\text{External}} = \underbrace{m_1 \ddot{w}_1}_{\text{Newton's 2nd Law}}$$

$$F_2 = k(w_1 - w_2) - c_2 \dot{w}_2$$

$$\ddot{w}_1 = \frac{k}{m_1}(w_2 - w_1) - \frac{c_1}{m_1}\dot{w}_1 + \frac{1}{m_1}f(w_1,t)$$

$$\ddot{w}_2 = \frac{k}{m_2}(w_1 - w_2) - \frac{c_2}{m_2}\dot{w}_2$$

*Harmonic Oscillator: motion of spring system*

*Discretization with Kinematics*

$$w_1(t+1) = w_1(t) + \dot{w}_1(t) + \frac{k}{2m_1}(w_2(t) - w_1(t))$$
$$- \frac{c_1}{2m_1}\dot{w}_1 + \frac{1}{2m_1}f(w_1,t)$$

$$= (1-\beta)\underbrace{w_1^*(t)}_{w_1(t)+\eta f(w_1,t)} + \beta w_2(t) \overset{\beta \to 0}{=} w_1^*(t)$$

*Optimizer Update*

$$w_2(t+1) = w_2(t) + \dot{w}_2(t) + \frac{k}{2m_2}(w_1(t) - w_2(t)) - \frac{c_2}{2m_2}\dot{w}_2$$

$$= (1-\alpha)w_2(t) + \alpha w_1(t) \to \mathbf{EMA}(w_1)$$

When $c_1 = 2m_1$, $c_2 = 2m_2$, for constants $\alpha, \beta, \eta$.

**BELAY**: *modified EMA as Harmonic Oscillator*



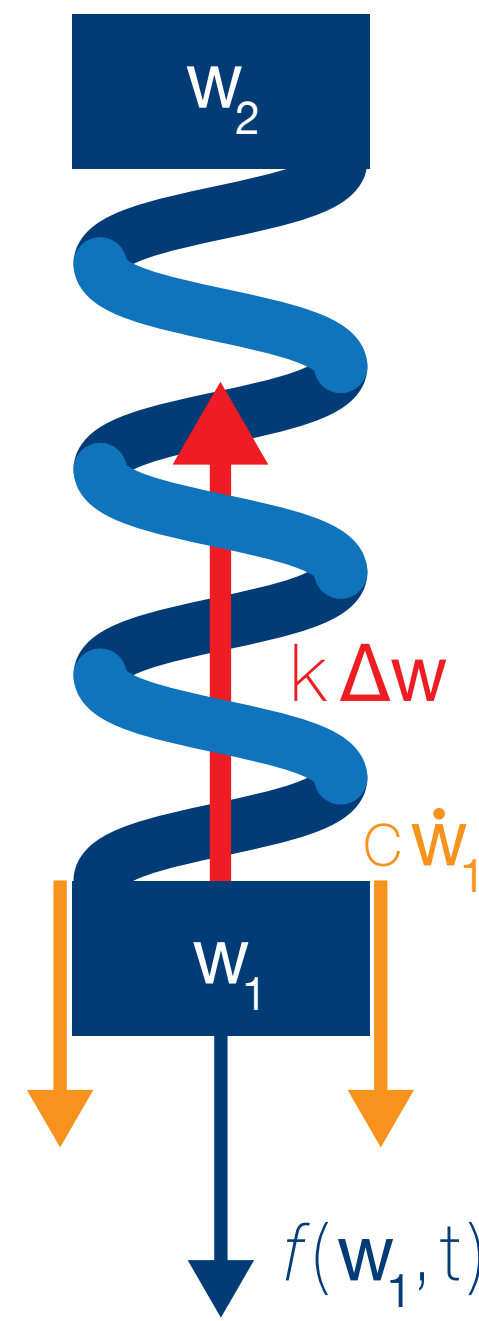Fig P2. An illustration of a classical Hookean spring system in the scenario described. Forces illustrated on $w_1$. Forces on $w_2$ are not pictured, but would be mirror images, with no external force applied.
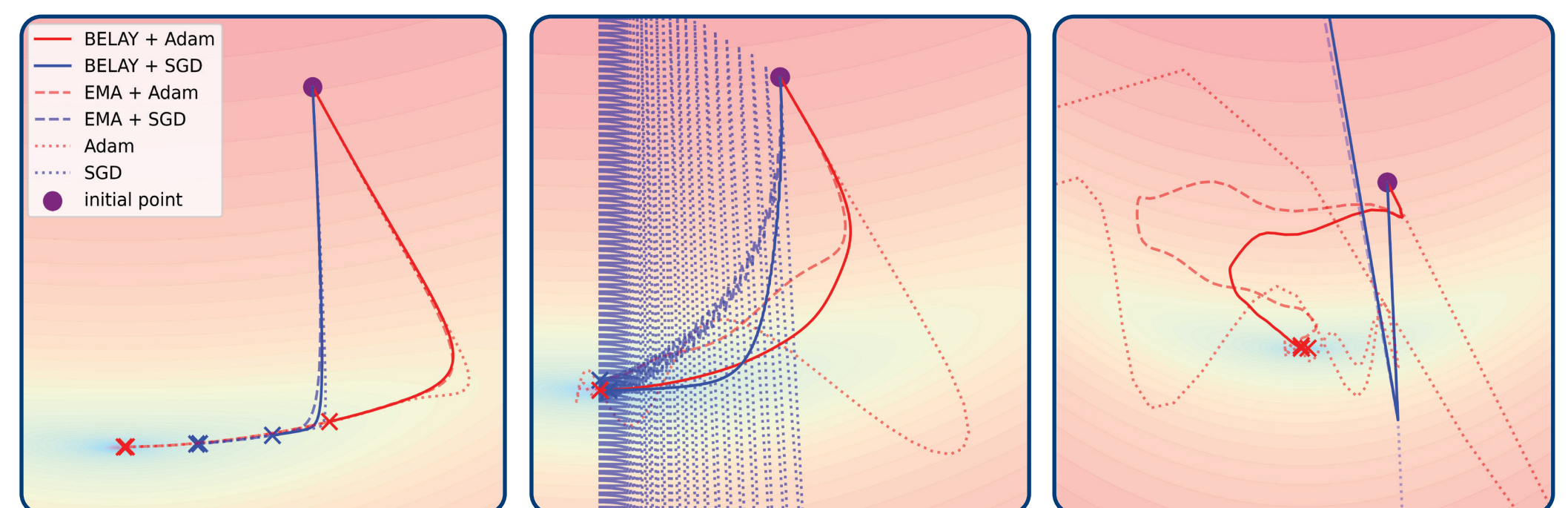


Fig 2. Comparison of BELAY against EMA and a control, using both Adam, and SGD on the Rosenbrock function across learning rates. Robustness to learning rate (η) is related to robustness across varying function smoothness.
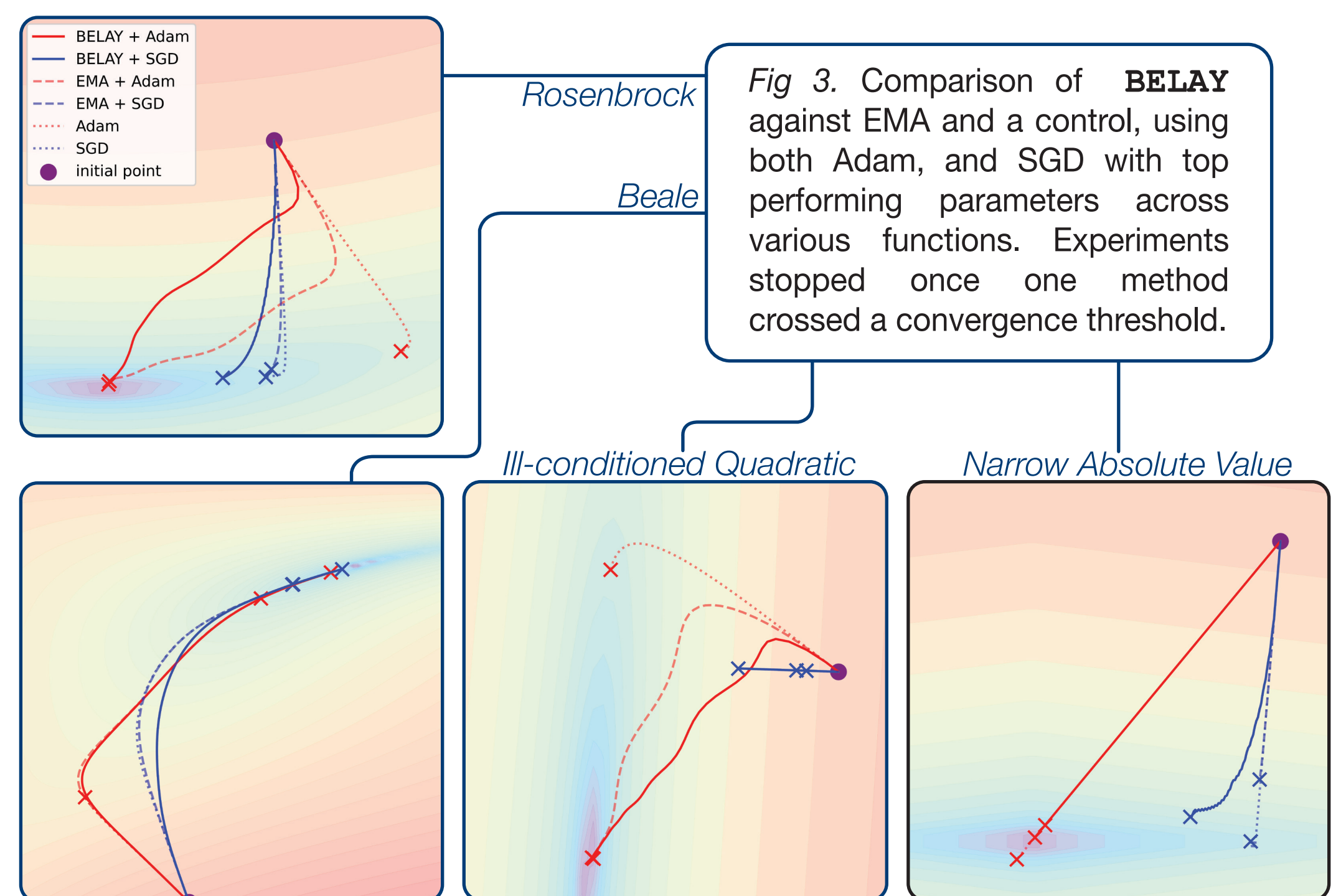


Fig 3. Comparison of BELAY against EMA and a control, using both Adam, and SGD with top performing parameters across various functions. Experiments stopped once one method crossed a convergence threshold.

## Further Insights
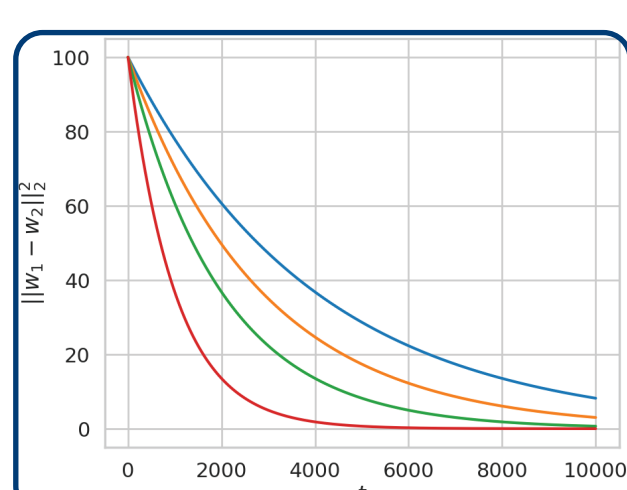
### Connections to Momentum-based Methods

$$\text{Momentum} \begin{cases} v(t) = \lambda g(w(t)) + (1-\lambda)v(t-1) = (1-\lambda)^s \lambda g(w(t-s)) \\ w(t+1) = w(t) + \alpha v(t) = w(t) + \alpha \sum_{s=0}^{t} a_s g(w(t-s)) \end{cases}$$

*↓ Linear g ↓*

$$= w(t) + \alpha g\left(\sum_{s=0}^{t} a_s w(t-1)\right) = w(t) + \alpha g(w^{EMA}(t)) \quad \text{BELAY}$$

### Physically-based Spring Parameterization

$$w(t) = C_1 e^{\left(-\delta + \sqrt{\delta^2 - \frac{k}{m}}\right)t} + C_2 e^{\left(-\delta - \sqrt{\delta^2 - \frac{k}{m}}\right)t}$$
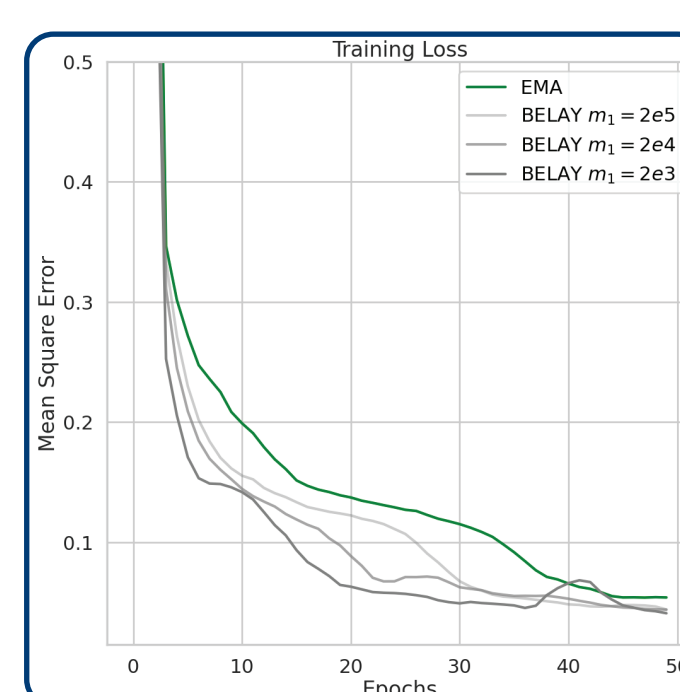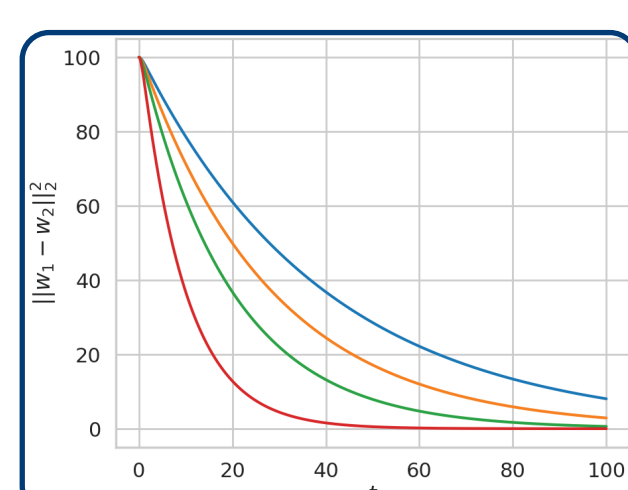
*Spring System Solution*

*Time-invariant Dynamics*



Fig 4. Comparison of BELAY against EMA on the MNIST dataset. The standard EMA algorithm is compared against BELAY with various settings of the model mass $m_2$.