

Understanding Adaptive Optimizers with Central Flows

ELLIT Symposium 2026

Alex Damian

Kempner Institute → **MIT**

Central Flows: A Brief Recap

Central Flows: A Brief Recap

Challenge: optimizers operate at the **edge of stability** and their trajectories are oscillatory

Central Flows: A Brief Recap

Challenge: optimizers operate at the **edge of stability** and their trajectories are oscillatory

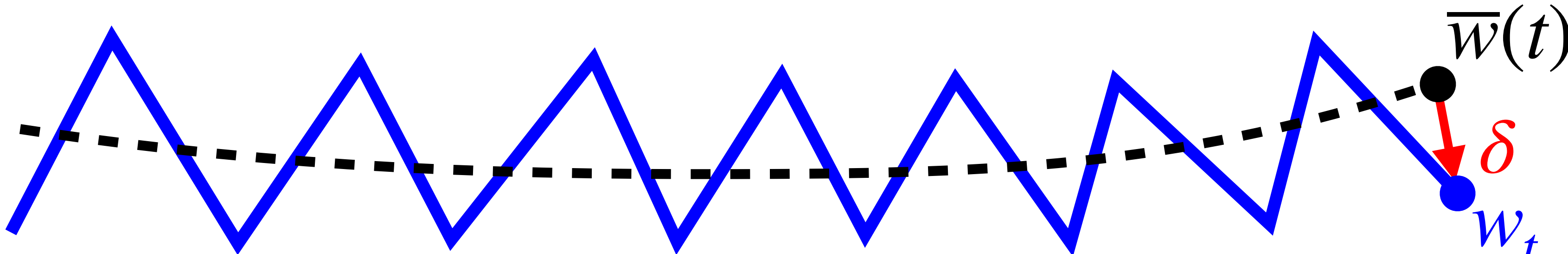
Approach: model the **time-averaged** optimization trajectory

Central Flows: A Brief Recap

Challenge: optimizers operate at the **edge of stability** and their trajectories are oscillatory

Approach: model the **time-averaged** optimization trajectory

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$

$$w_t = \bar{w}(t) + \delta_t \quad \mathbb{E}[\delta_t] = 0 \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$


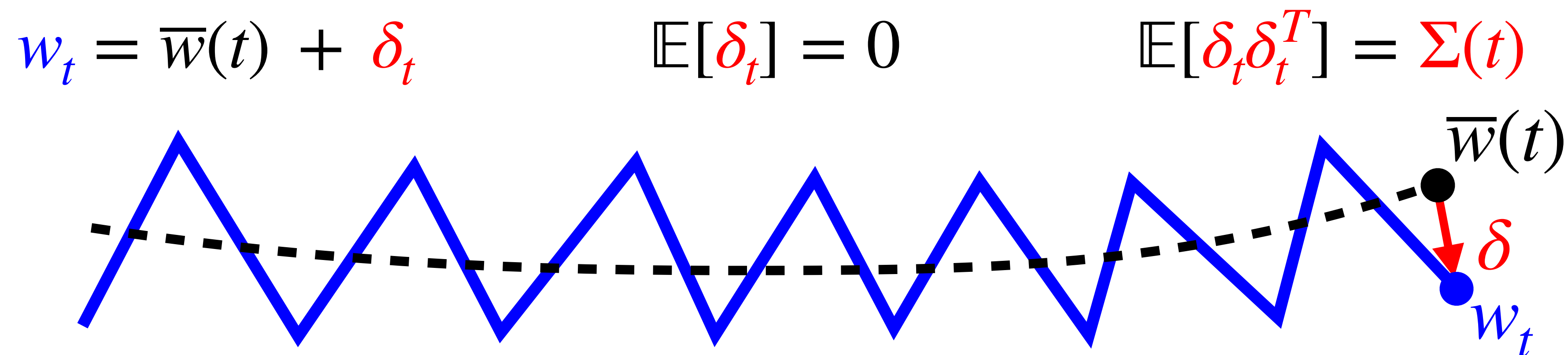
The diagram illustrates the decomposition of the weight vector w_t into a central flow $\bar{w}(t)$ and a deviation δ_t . A blue zigzag line represents the oscillatory trajectory w_t . A dashed black line represents the central flow $\bar{w}(t)$. A red arrow labeled δ points from the current point w_t (blue dot) to the central flow $\bar{w}(t)$ (black dot).

Central Flows: A Brief Recap

Challenge: optimizers operate at the **edge of stability** and their trajectories are oscillatory

Approach: model the **time-averaged** optimization trajectory

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$



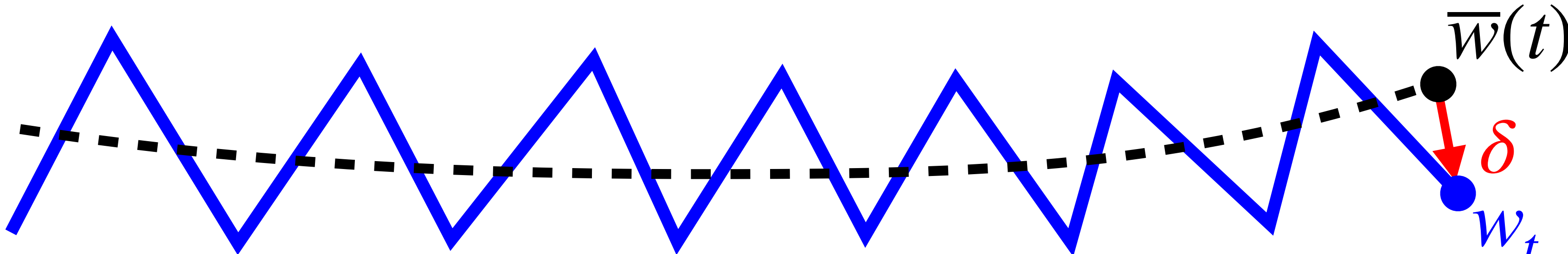
- ▶ Derive necessary conditions for $(\bar{w}(t), \Sigma(t))$ which allow us to solve for $\frac{d}{dt}\bar{w}(t)$

Central Flows: A Brief Recap

Challenge: optimizers operate at the **edge of stability** and their trajectories are oscillatory

Approach: model the **time-averaged** optimization trajectory

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$

$$w_t = \bar{w}(t) + \delta_t \quad \mathbb{E}[\delta_t] = 0 \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$


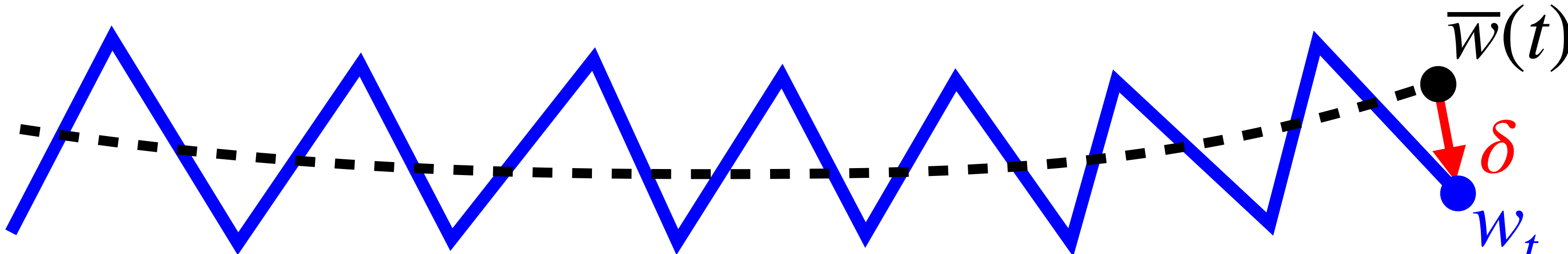
- ▶ Derive necessary conditions for $(\bar{w}(t), \Sigma(t))$ which allow us to solve for $\frac{d}{dt}\bar{w}(t)$
- ▶ Empirically verify these central flows in realistic deep learning settings

Central Flows: A Brief Recap

Challenge: optimizers operate at the **edge of stability** and their trajectories are oscillatory

Approach: model the **time-averaged** optimization trajectory

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$

$$w_t = \bar{w}(t) + \delta_t \quad \mathbb{E}[\delta_t] = 0 \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$


- ▶ Derive necessary conditions for $(\bar{w}(t), \Sigma(t))$ which allow us to solve for $\frac{d}{dt}\bar{w}(t)$
- ▶ Empirically verify these central flows in realistic deep learning settings
- ▶ Interpret the central flow to understand the implicit behavior of the optimizer

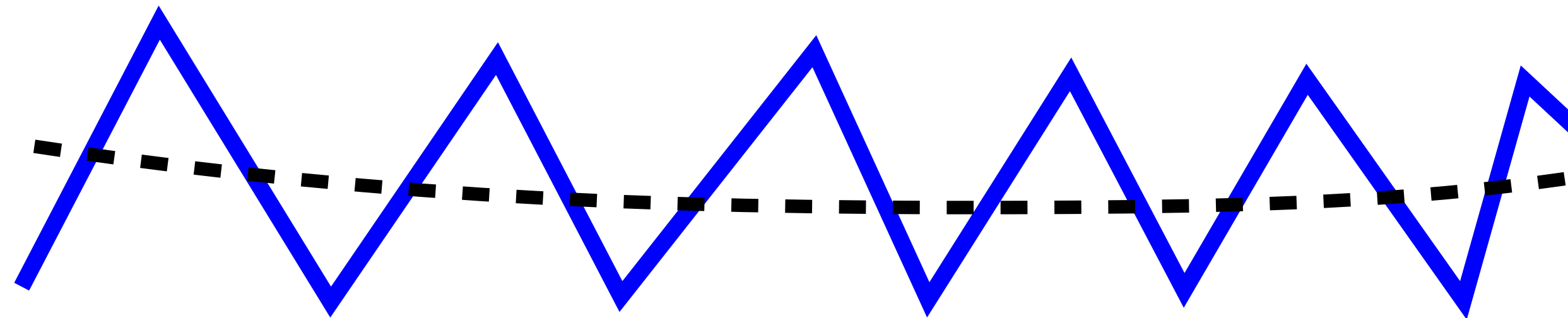
Central Flows: A Brief Recap

Challenge: optimizers operate at the **edge of stability** and their trajectories are oscillatory

Approach: model the **time-averaged** optimization trajectory

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance

$$w_t = \bar{w}(t) + \delta_t \quad \mathbb{E}[\delta_t] = 0 \quad \mathbb{E}[\delta_t \delta_s^T] = \Sigma(t) \delta_{ts}$$



- ▶ Derive necessary conditions for $(\bar{w}(t), \Sigma(t))$ which allow us to
- ▶ Empirically verify these central flows in realistic deep learning
- ▶ Interpret the central flow to understand the implicit behavior of the



but it is **very** predictive!

RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

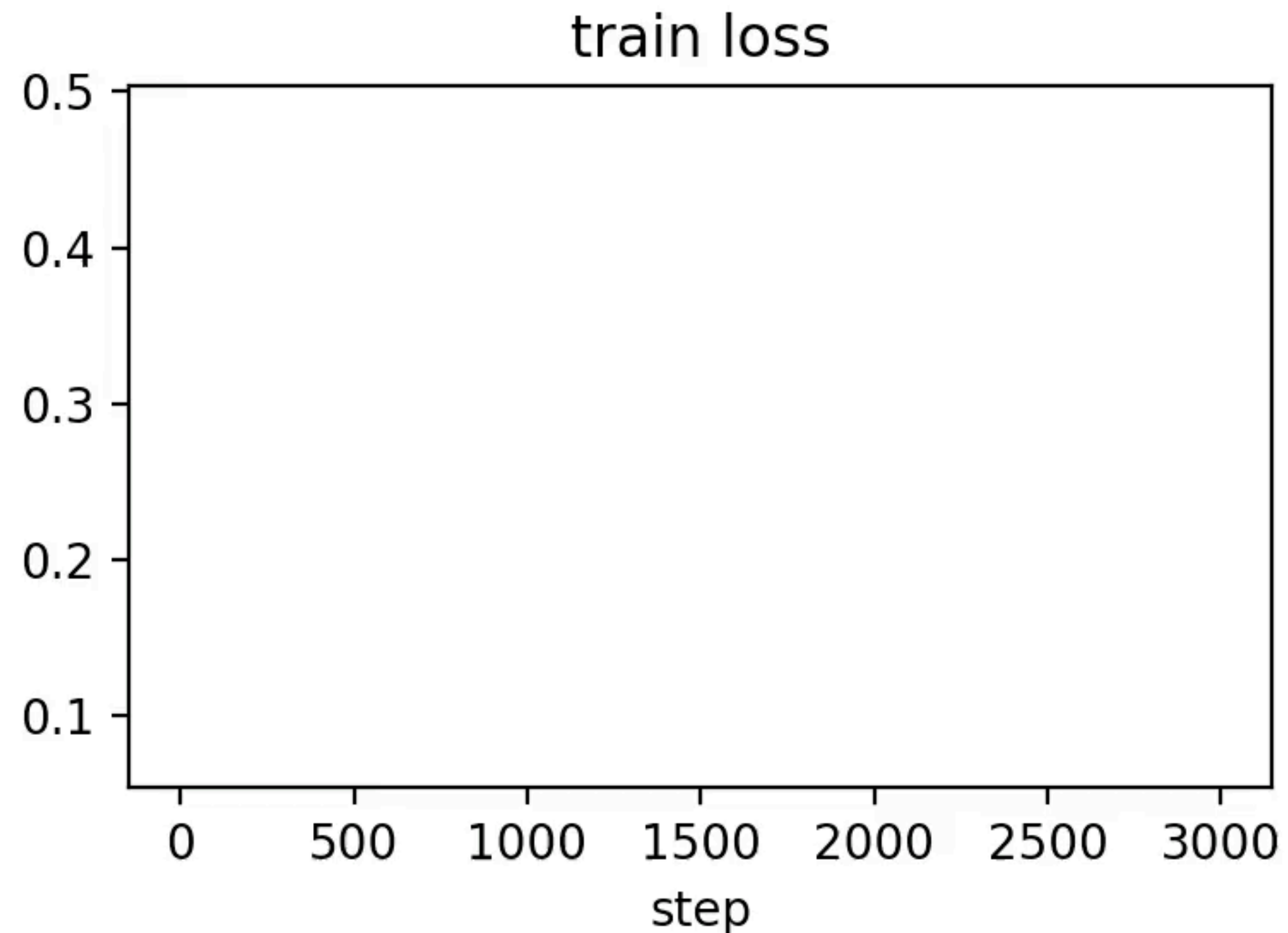
RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes



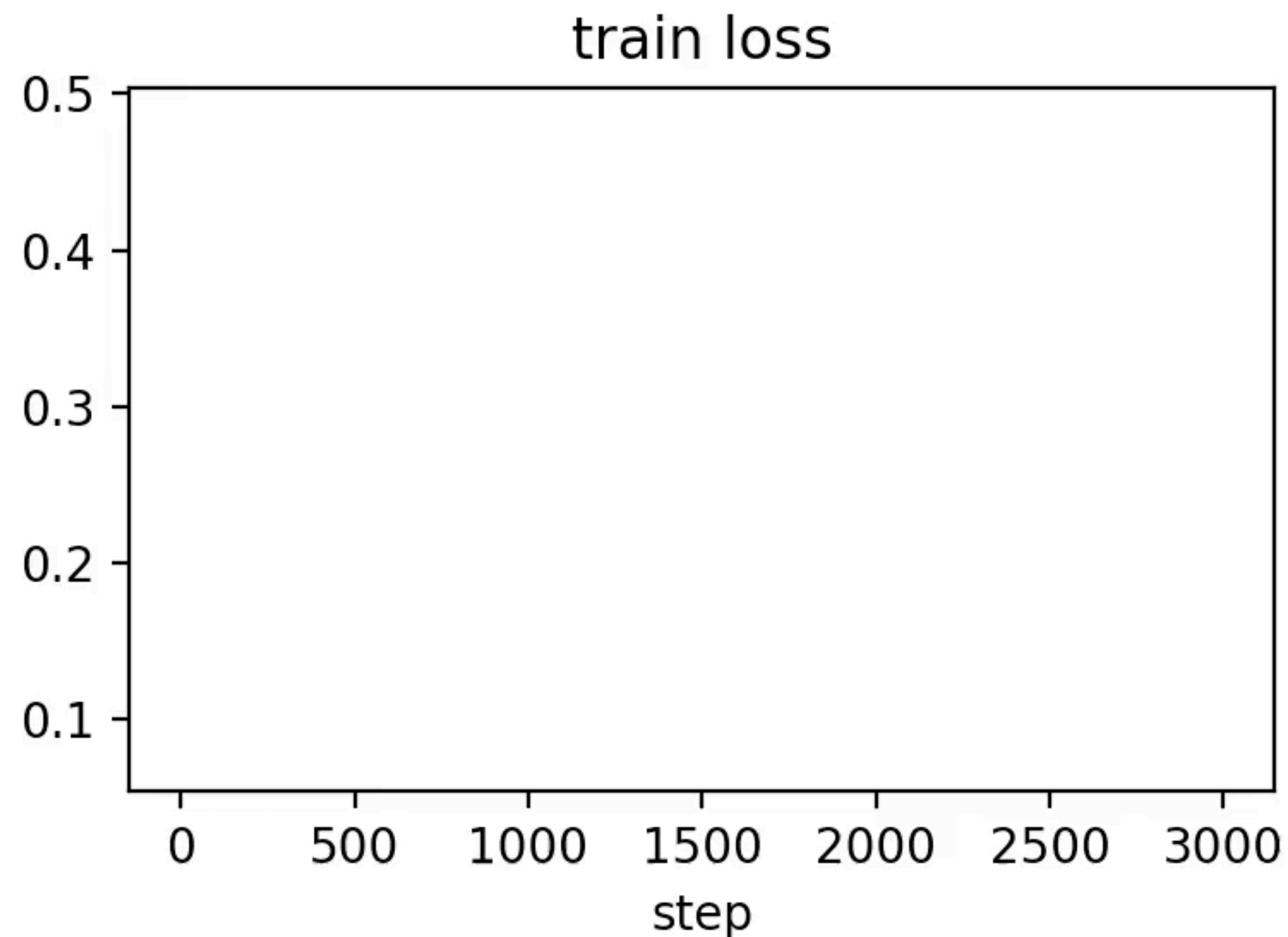
RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes



If Adam is adaptive,

- ▶ what does it adapt **to**?
- ▶ what does the learning rate do?

RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

- ▶ what does it adapt **to**?

RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

► what does it adapt **to**?

A: it adapts to the gradient scale... duh?

RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

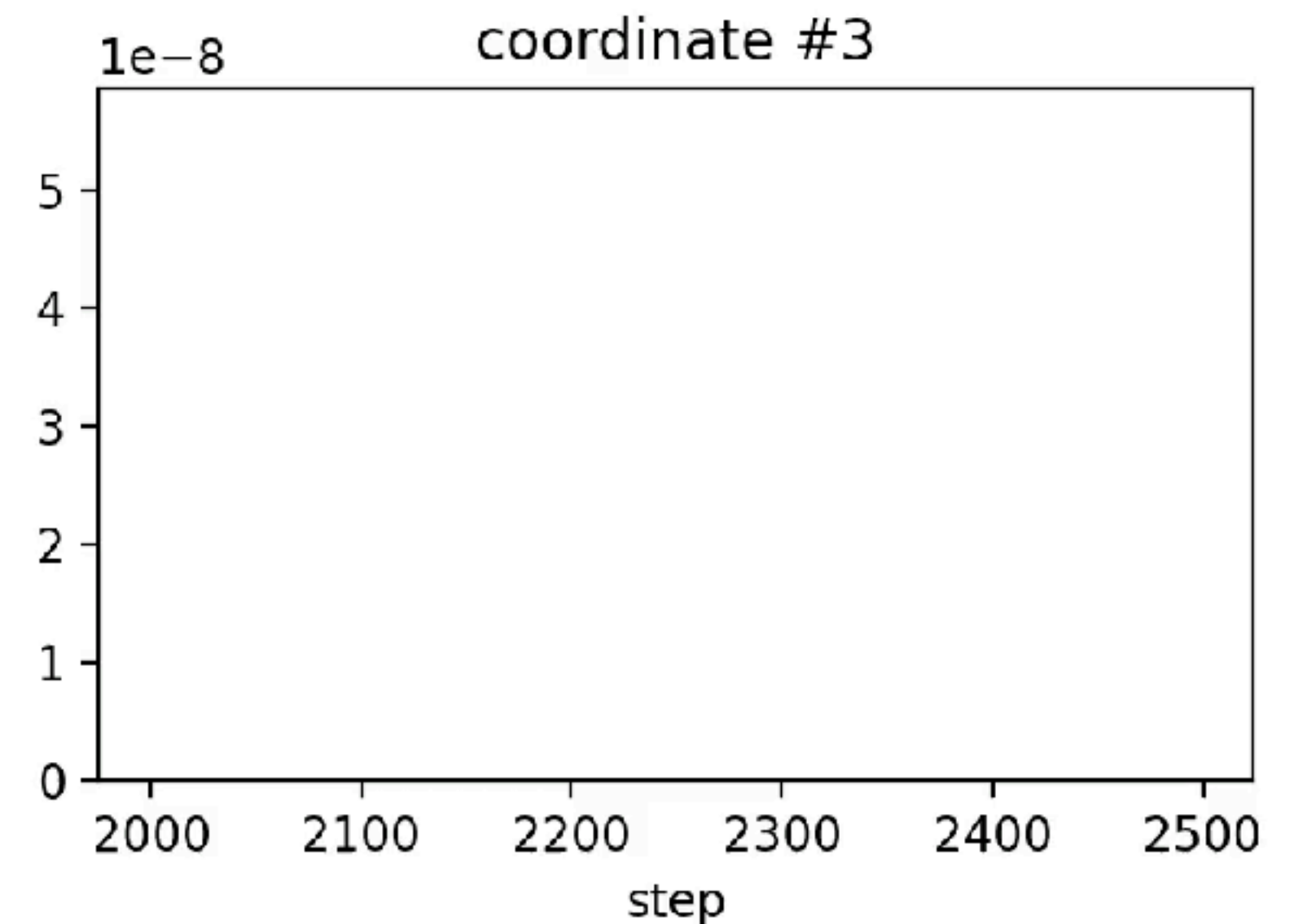
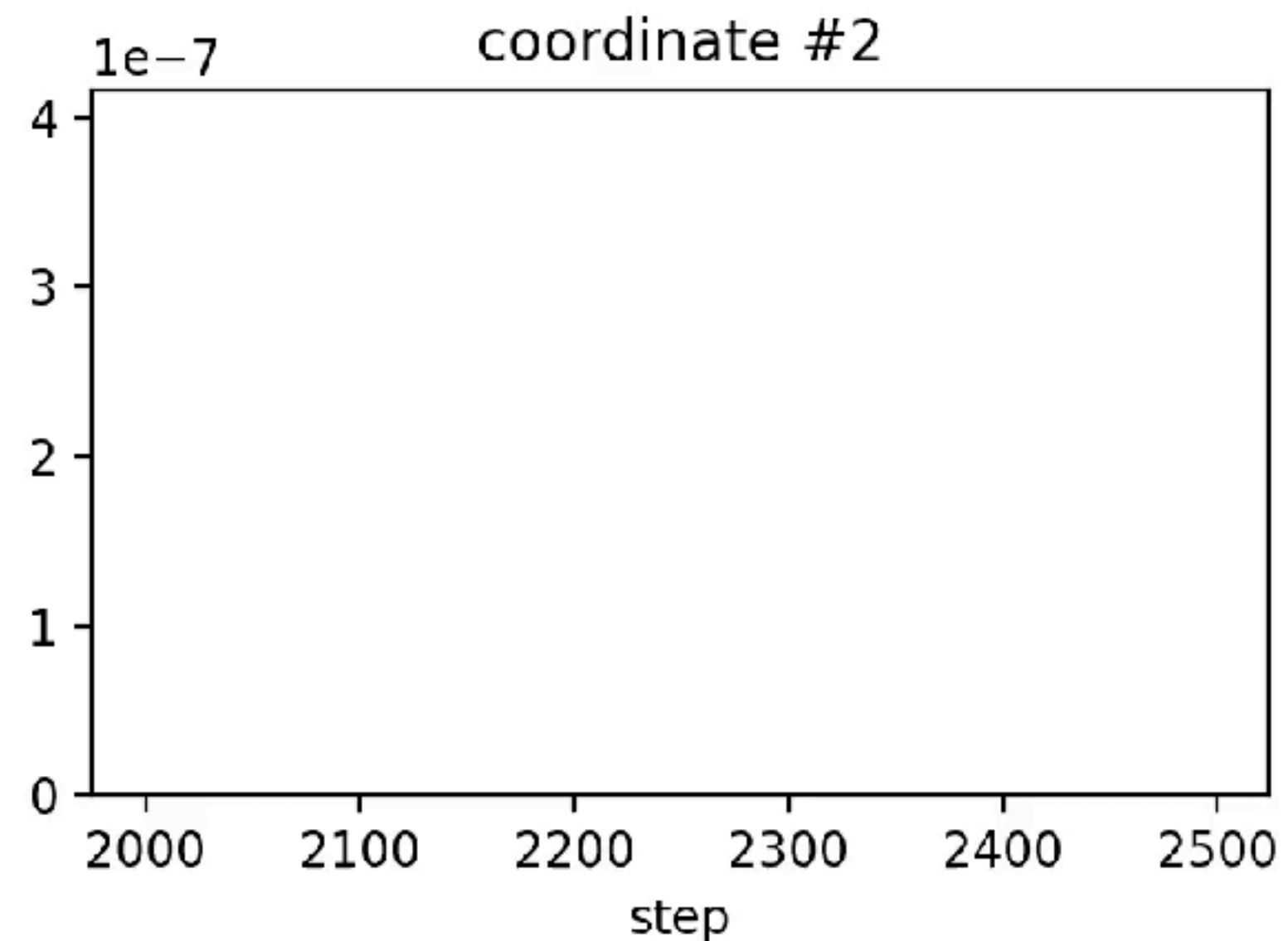
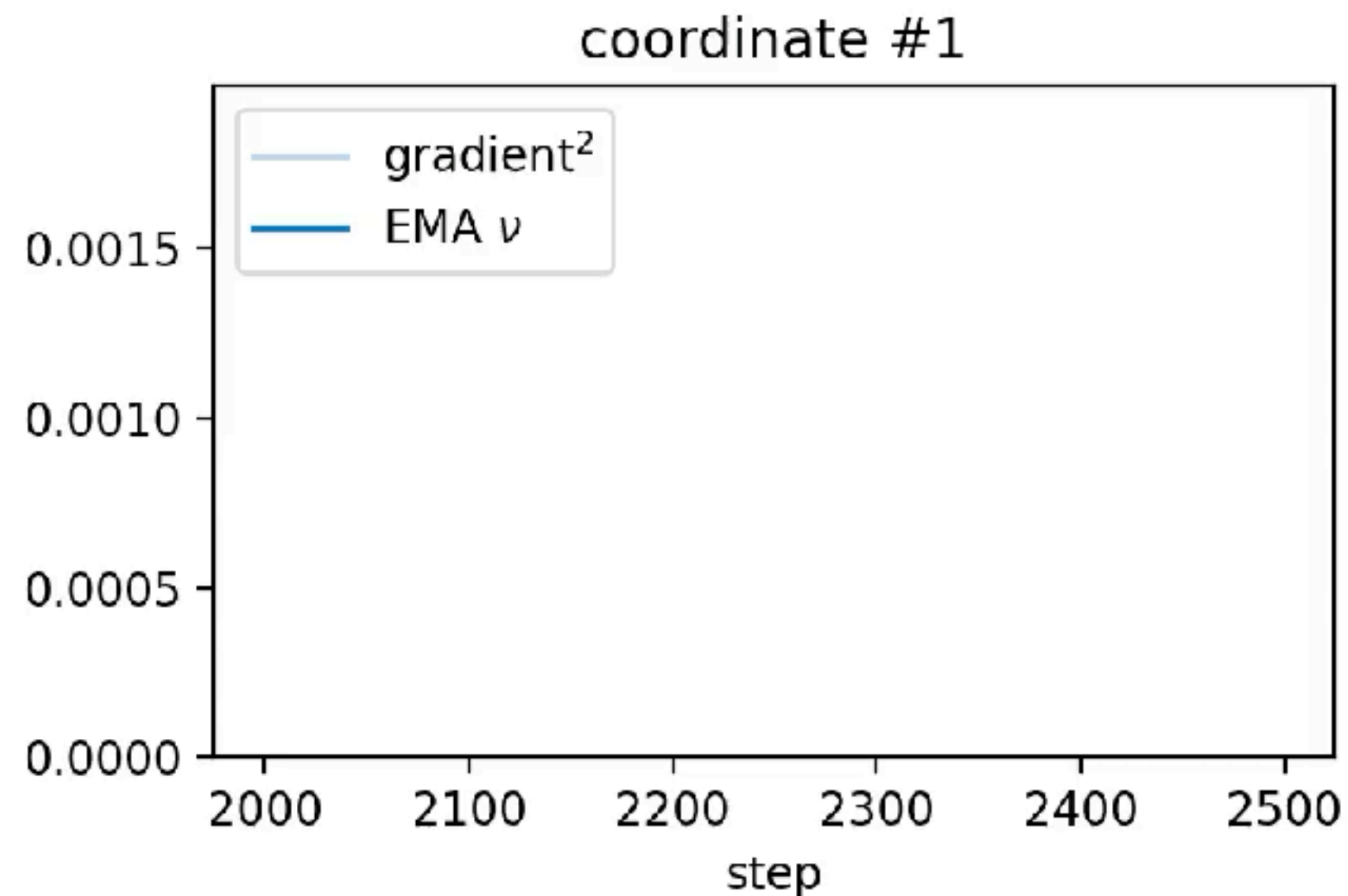
maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

► what does it adapt **to**?

A: it adapts to the gradient scale... duh?



RMSProp (Adam without momentum)

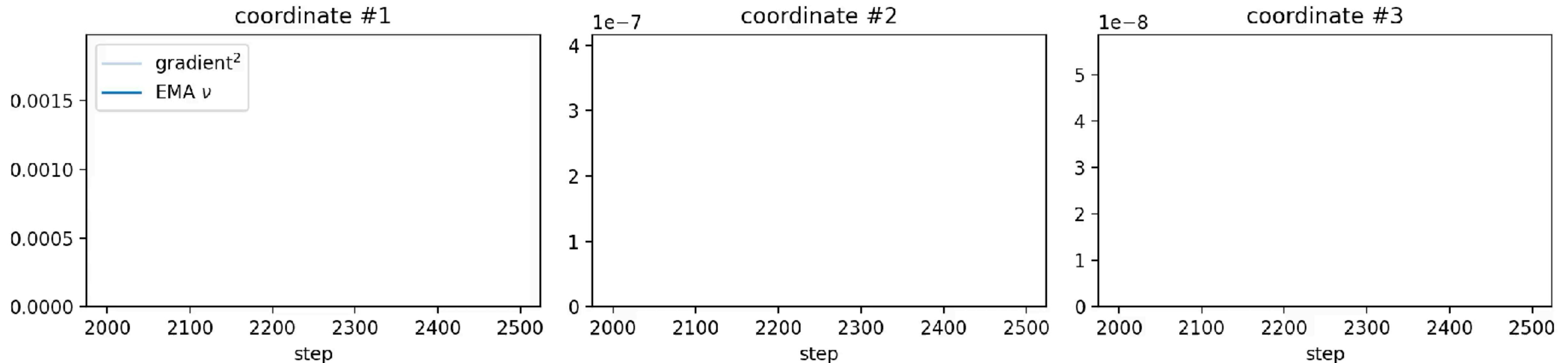
$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

Like Gradient Descent, RMSProp operates at the **edge of stability**



RMSProp (Adam without momentum)

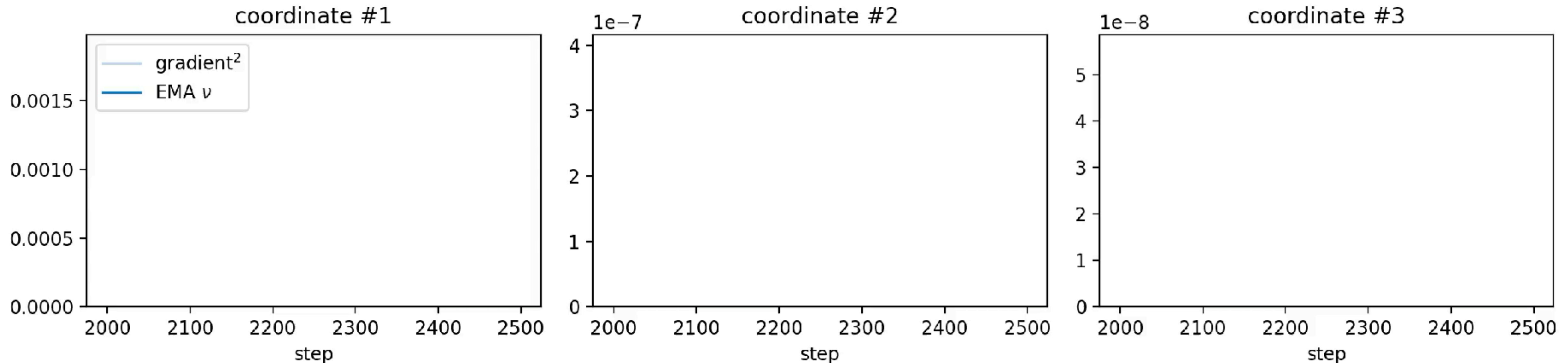
$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

Like Gradient Descent, RMSProp operates at the **edge of stability**



RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

► If $P_t^{-1} := \text{diag} \left(\frac{\eta}{\sqrt{\nu_t}} \right)$, RMSProp is just GD with an evolving preconditioner P_t :

$$w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$

RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

- ▶ If $P_t^{-1} := \text{diag} \left(\frac{\eta}{\sqrt{\nu_t}} \right)$, RMSProp is just GD with an evolving preconditioner P_t :

$$w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$

- ▶ On a quadratic $L(w) = \frac{1}{2} w^T H w$, the update is $w \leftarrow (I - P_t^{-1} H) w$

RMSProp (Adam without momentum)

$$\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \nabla L(w_t)^{\odot 2}$$

maintain an EMA of the squared gradient

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\nu_t}} \nabla L(w_t)$$

effective
step sizes

- ▶ If $P_t^{-1} := \text{diag} \left(\frac{\eta}{\sqrt{\nu_t}} \right)$, RMSProp is just GD with an evolving preconditioner P_t :

$$w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$

- ▶ On a quadratic $L(w) = \frac{1}{2} w^T H w$, the update is $w \leftarrow (I - P_t^{-1} H) w$

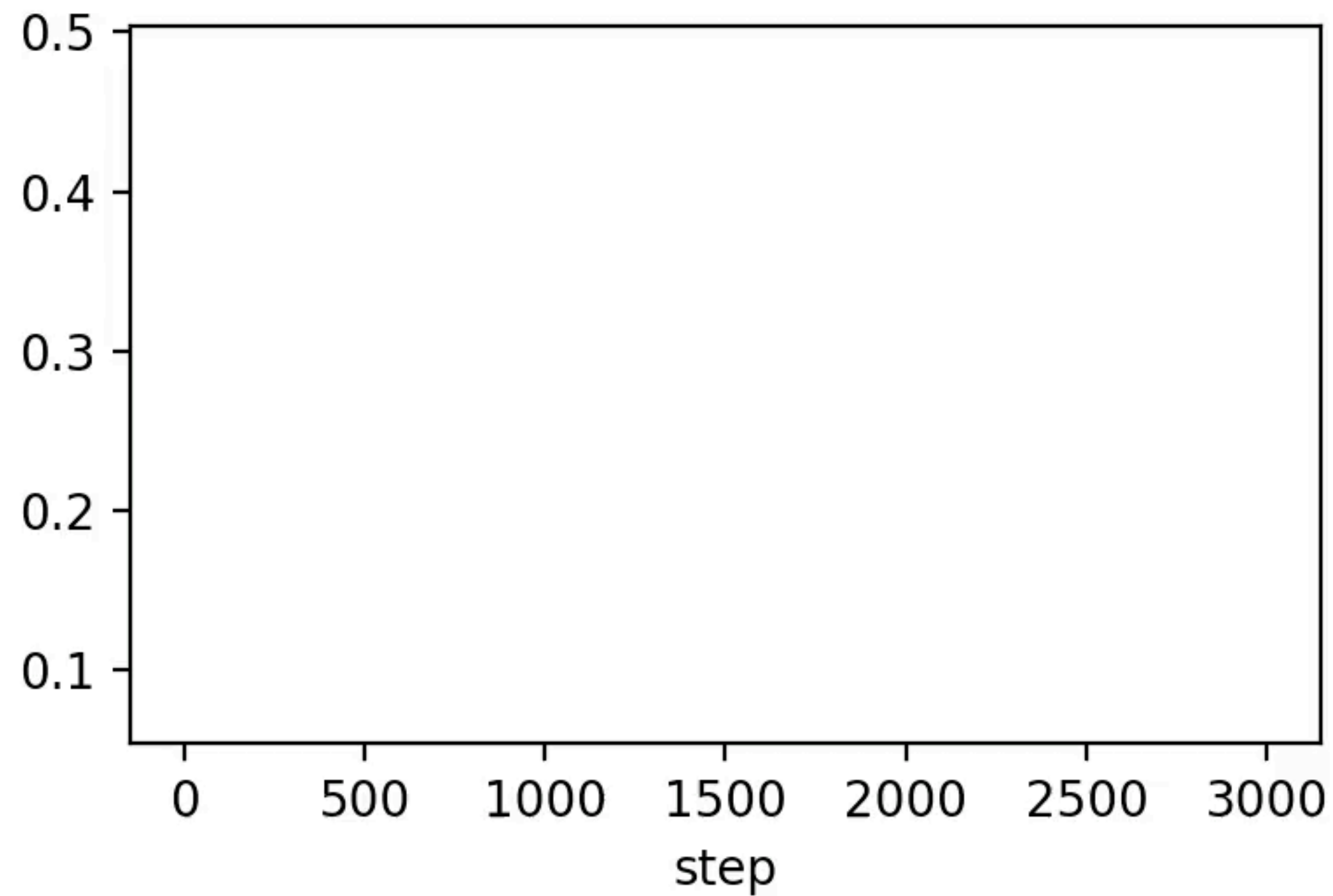
- ▶ In general, local stability is given by the eigenvalues of the **effective Hessian**:

$$\lambda_1(H_{\text{eff}}) \leq 2 \quad \text{where} \quad H_{\text{eff}} = P_t^{-1} H(w_t) \quad \text{and} \quad H = \nabla^2 L$$

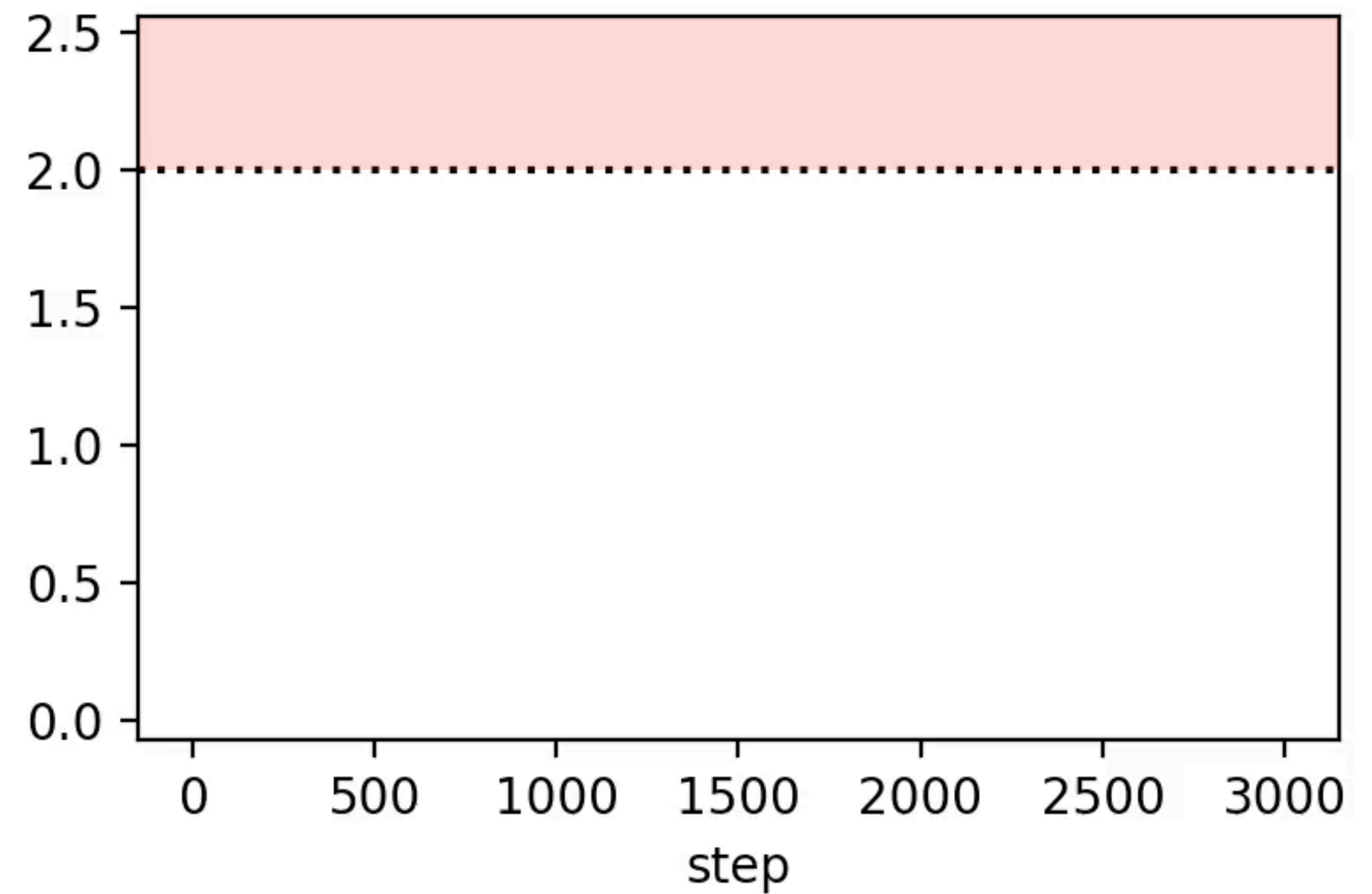
RMSProp in Deep Learning

RMSProp on a ResNet, trained on CIFAR10, $\eta = 2 \times 10^{-5}$, $\beta = 0.99$

Loss $L(w)$



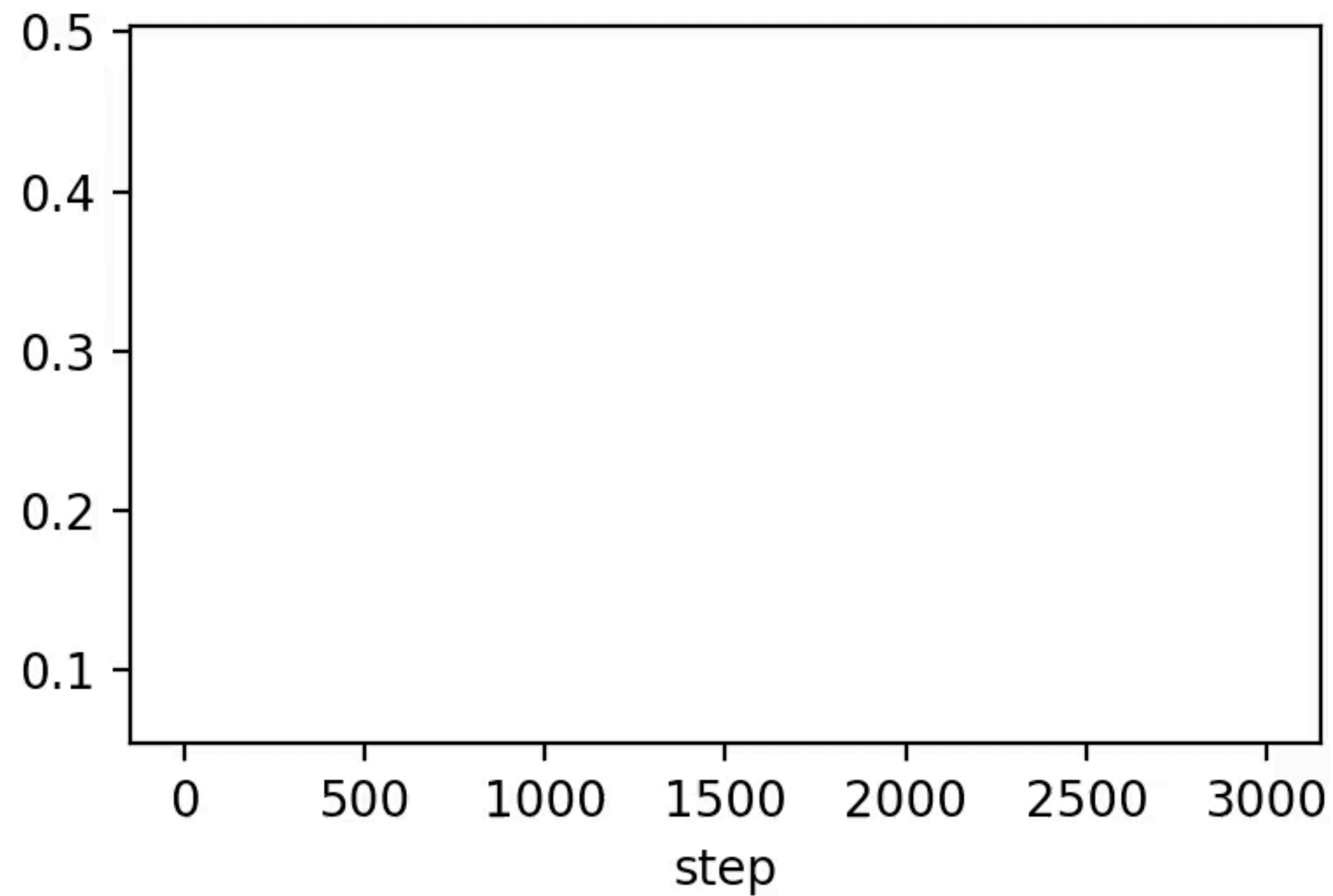
Top 4 evals of $H_{\text{eff}} := P_t^{-1}H(w_t)$



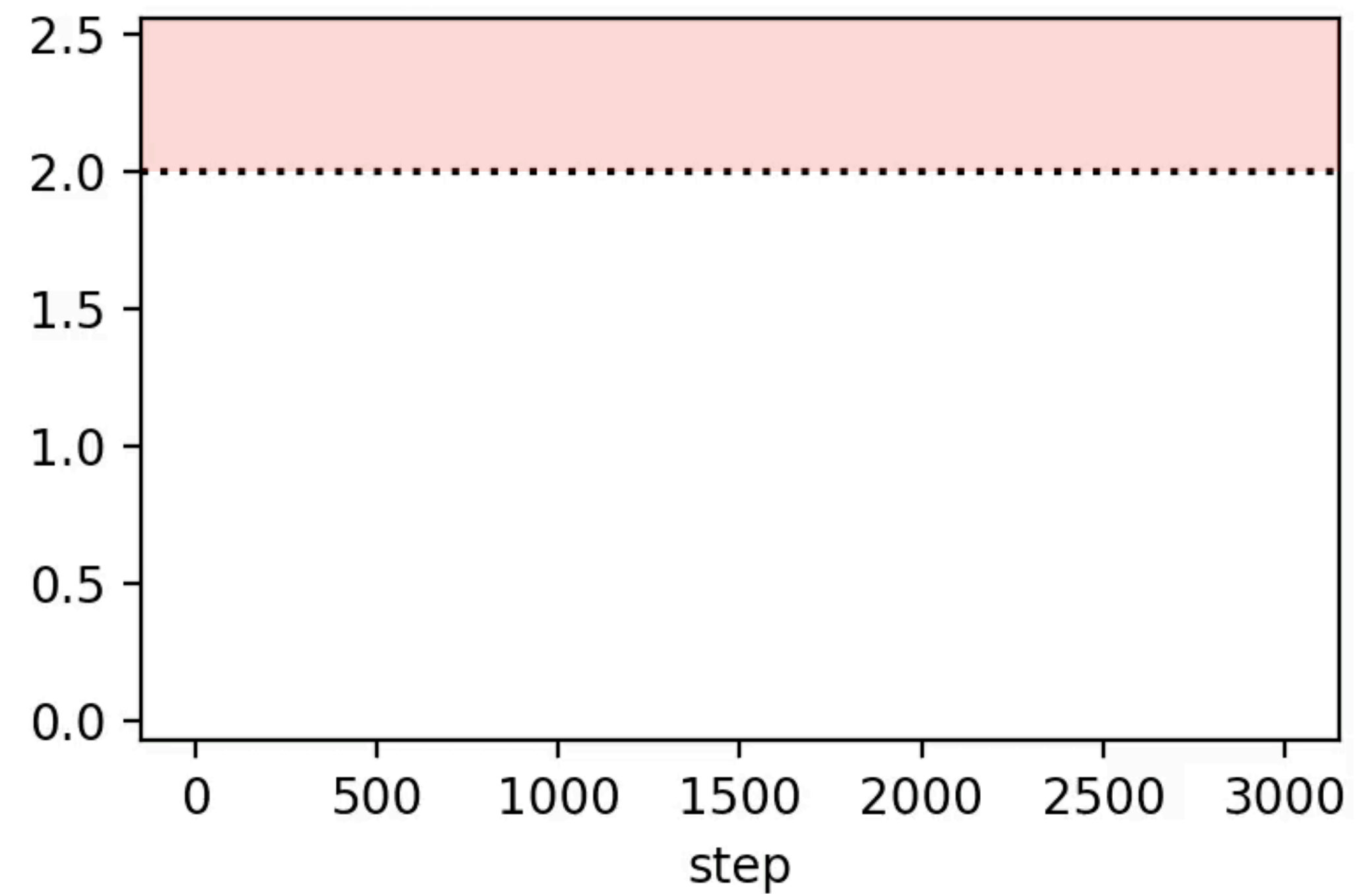
RMSProp in Deep Learning

RMSProp on a ResNet, trained on CIFAR10, $\eta = 2 \times 10^{-5}$, $\beta = 0.99$

Loss $L(w)$

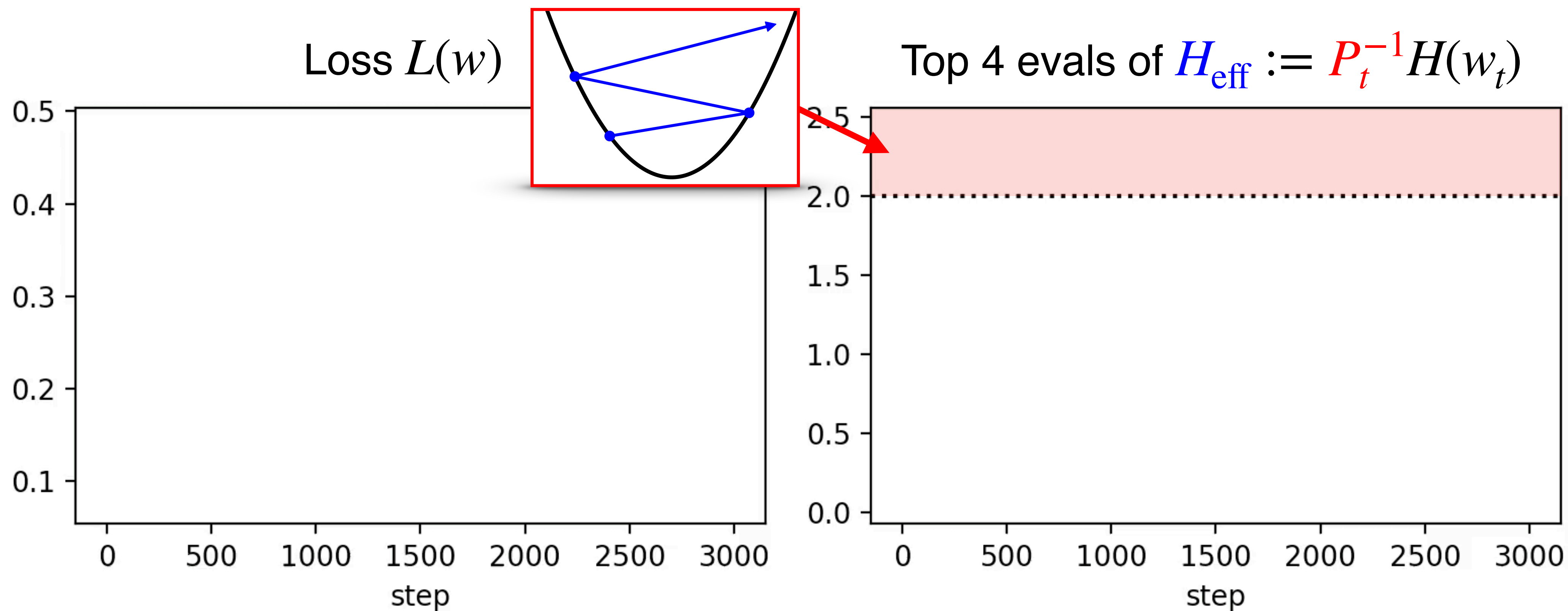


Top 4 evals of $H_{\text{eff}} := P_t^{-1}H(w_t)$



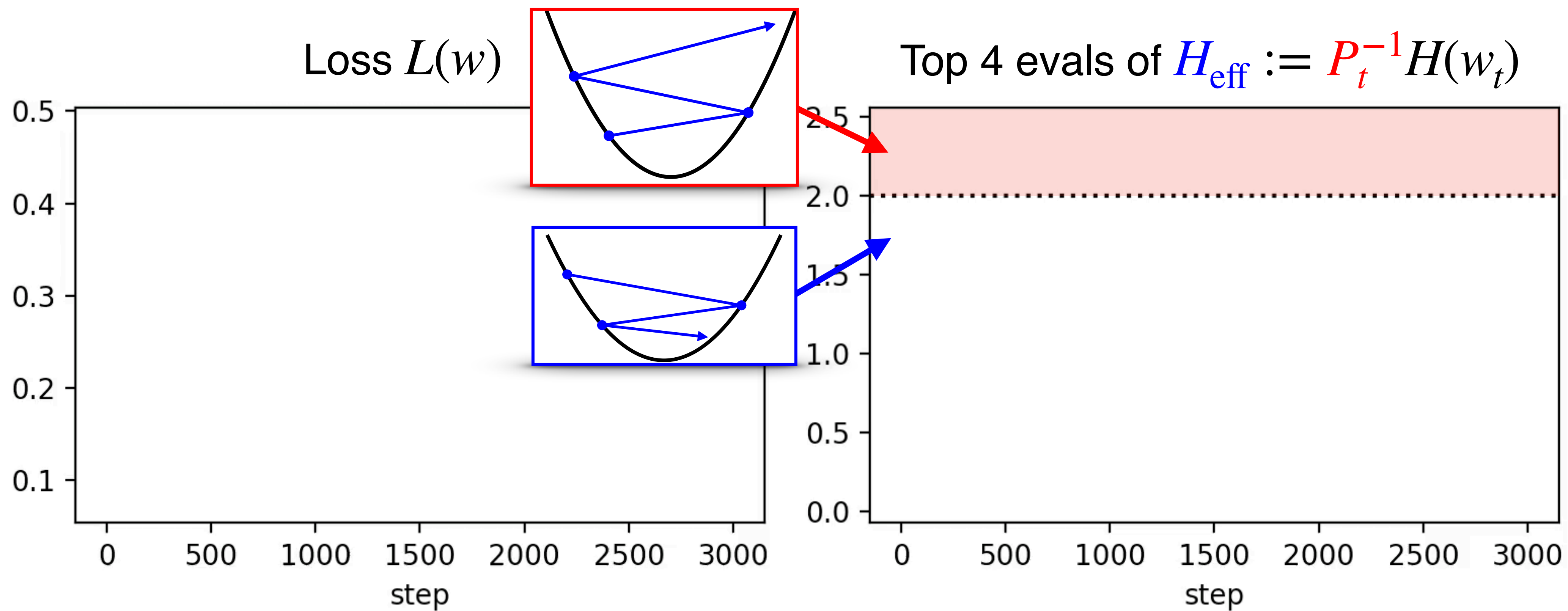
RMSProp in Deep Learning

RMSProp on a ResNet, trained on CIFAR10, $\eta = 2 \times 10^{-5}$, $\beta = 0.99$



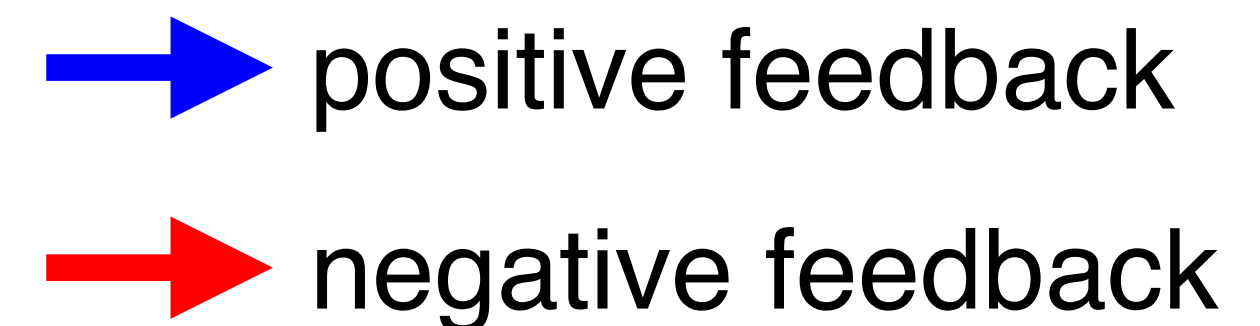
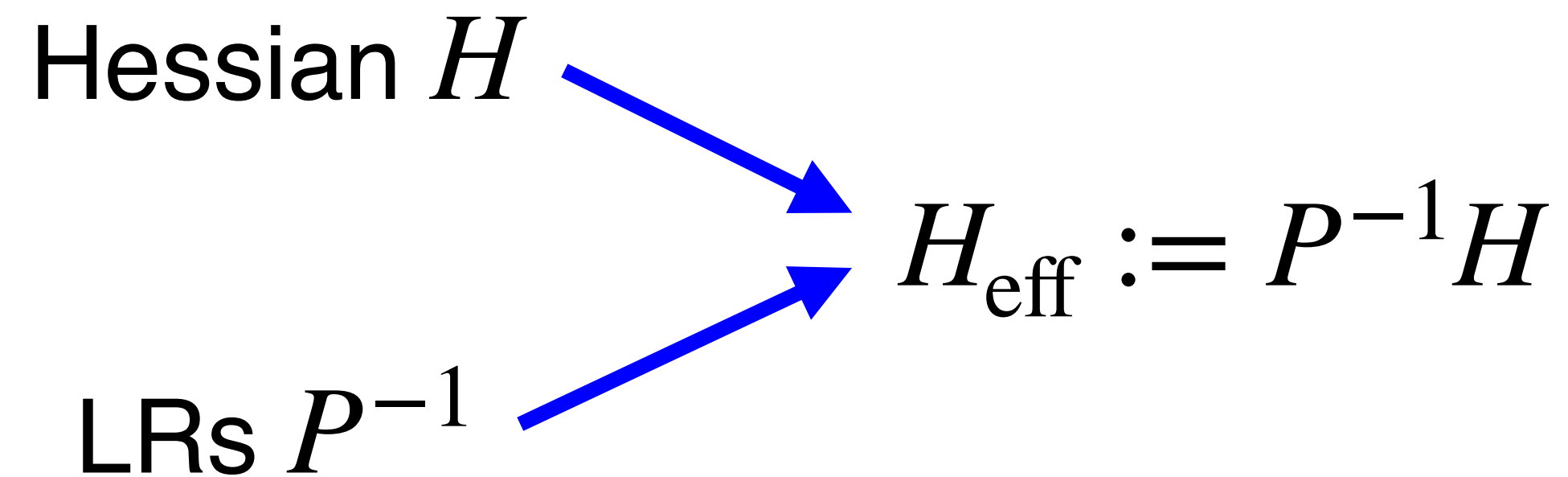
RMSProp in Deep Learning

RMSProp on a ResNet, trained on CIFAR10, $\eta = 2 \times 10^{-5}$, $\beta = 0.99$



RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\|\nabla L(w)\|^2) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$



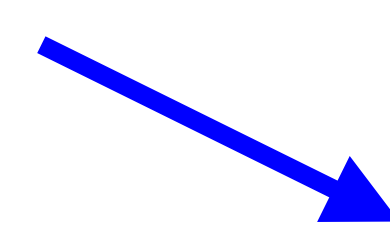
RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\nabla L(w)^{\odot 2}) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$

progressive
sharpening

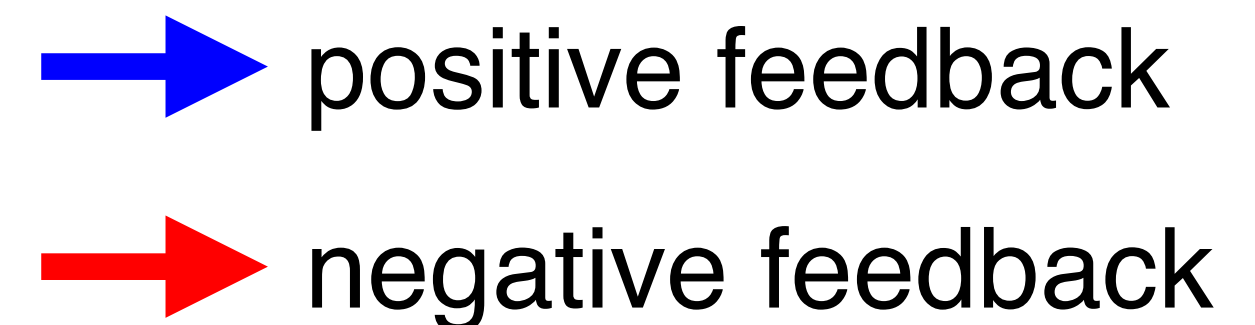
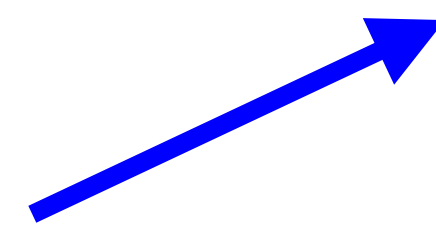


Hessian H



$$H_{\text{eff}} := P^{-1}H$$

LRs P^{-1}



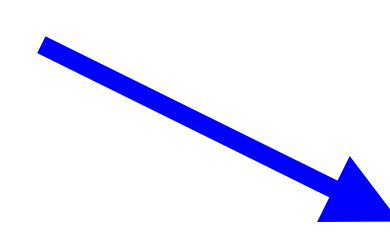
RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\nabla L(w)^{\odot 2}) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$

progressive
sharpening



Hessian H

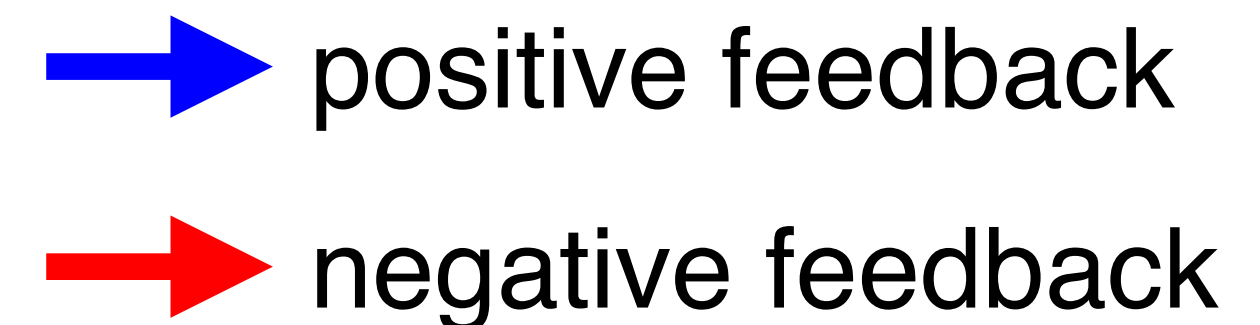
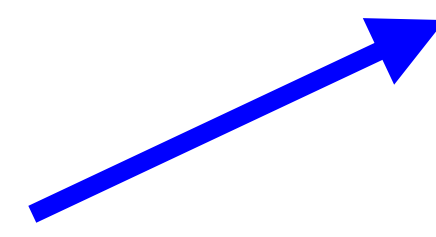


$H_{\text{eff}} := P^{-1}H$

$\nabla L(w)^{\odot 2}$

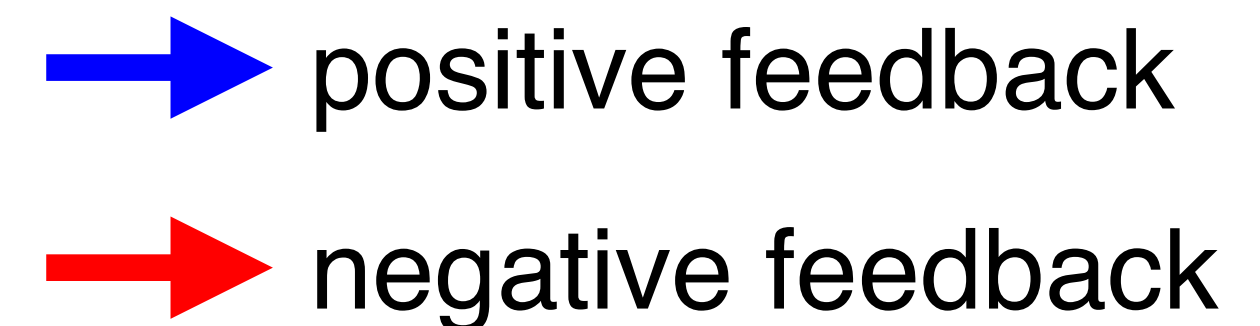
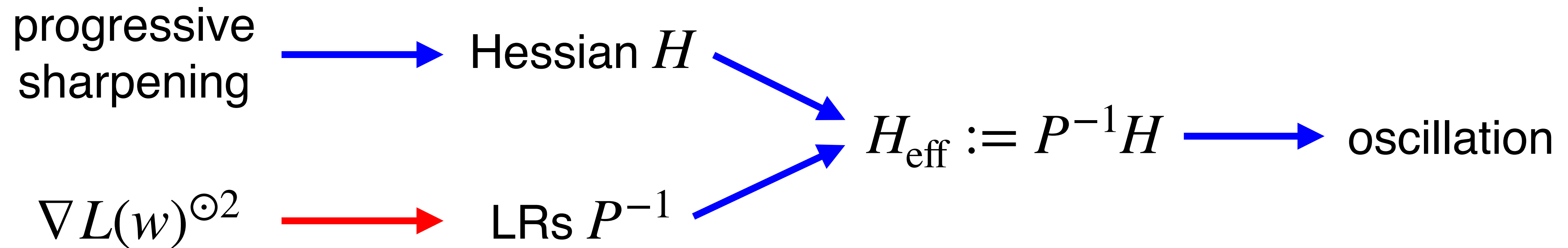


LRs P^{-1}



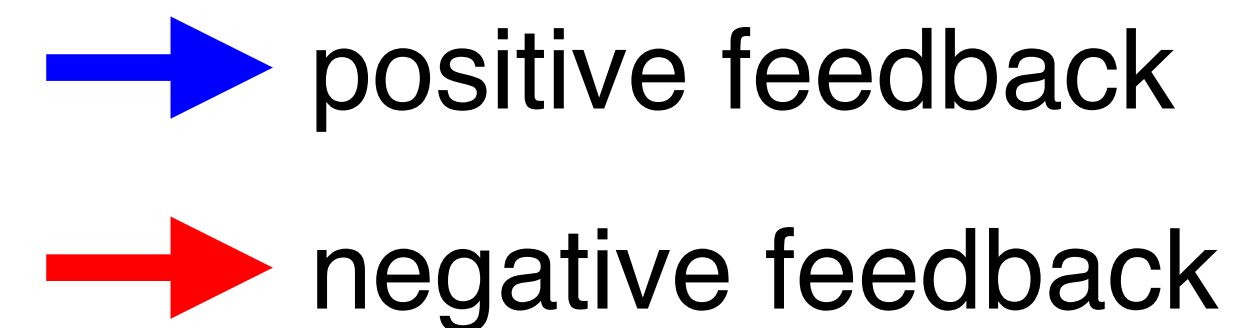
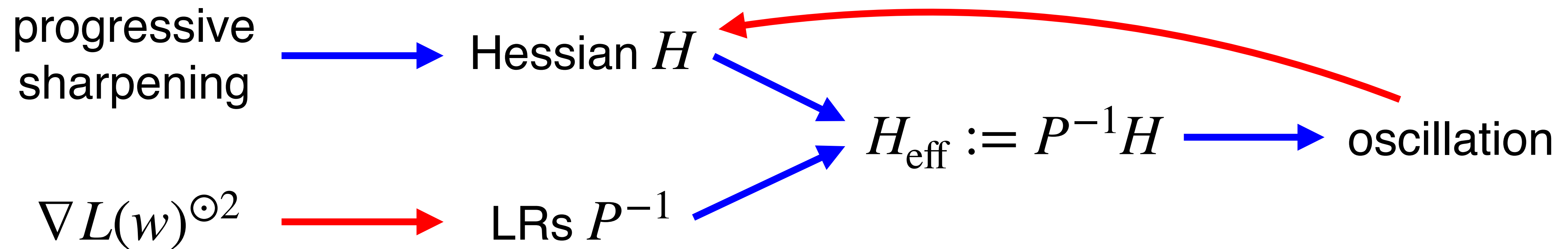
RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\nabla L(w)^{\odot 2}) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$



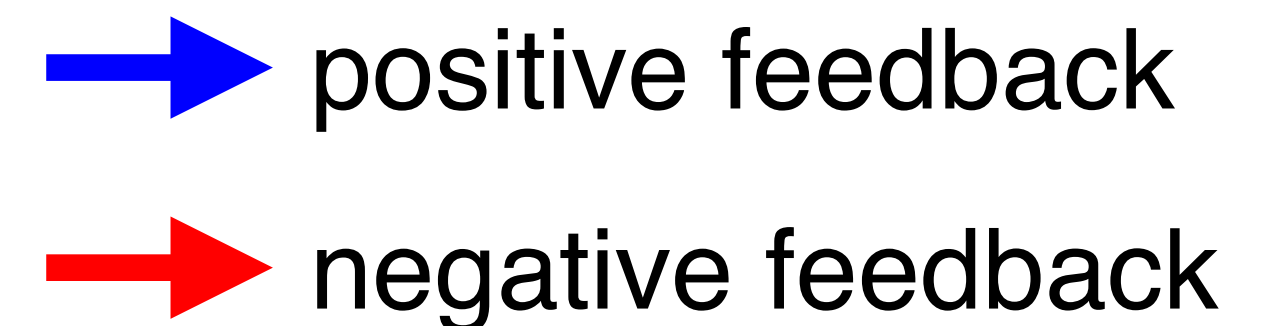
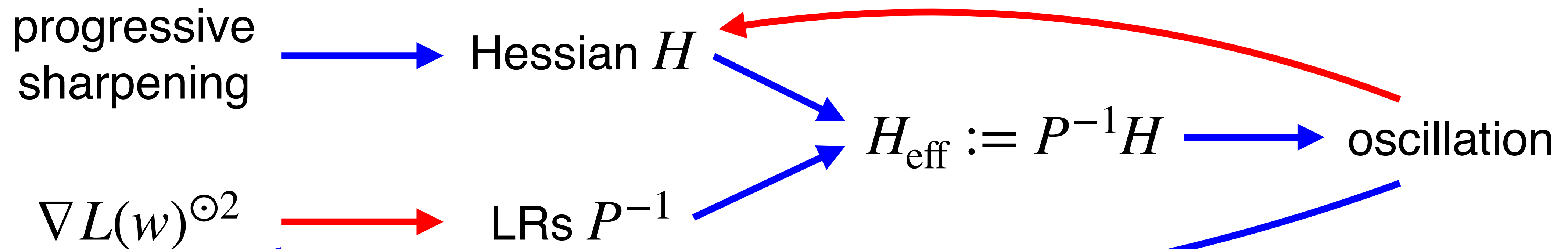
RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\nabla L(w)^{\odot 2}) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$



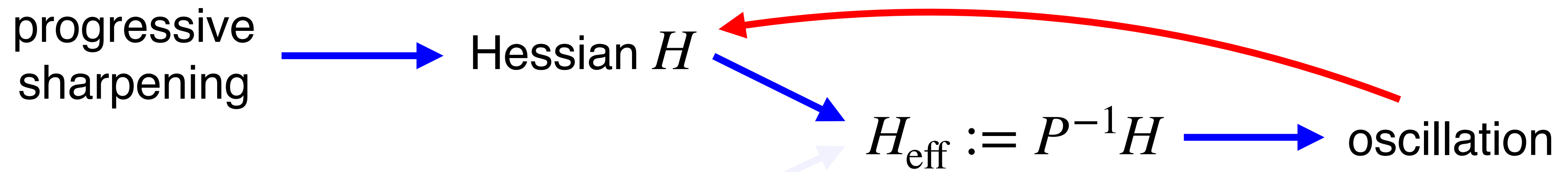
RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\nabla L(w)^{\odot 2}) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$



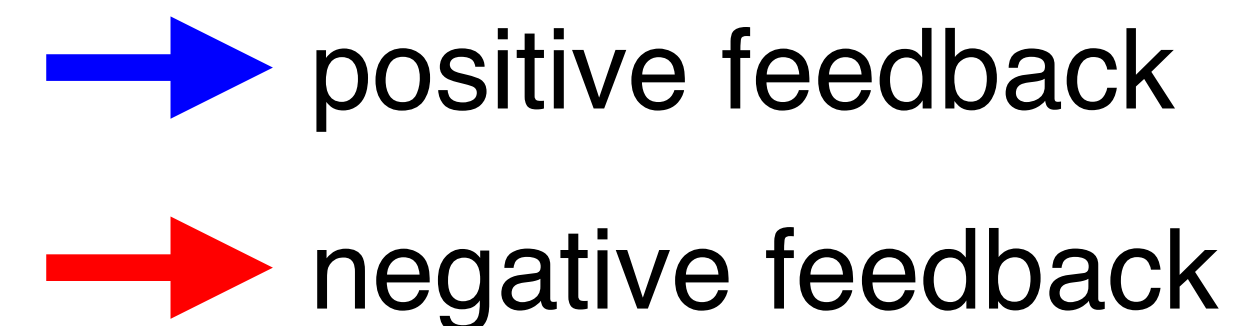
RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\nabla L(w)^{\odot 2}) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$



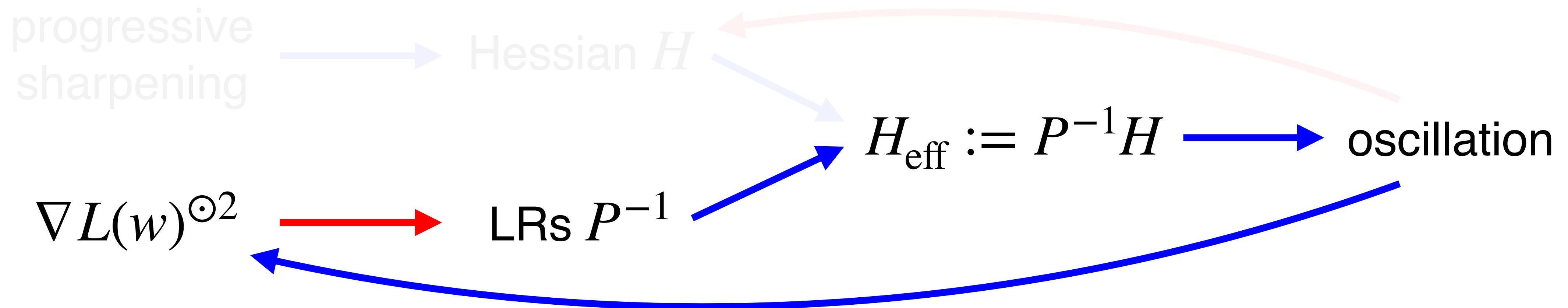
$\nabla L(w)^{\odot 2}$

LRs P^{-1}



RMSProp at the Edge of Stability

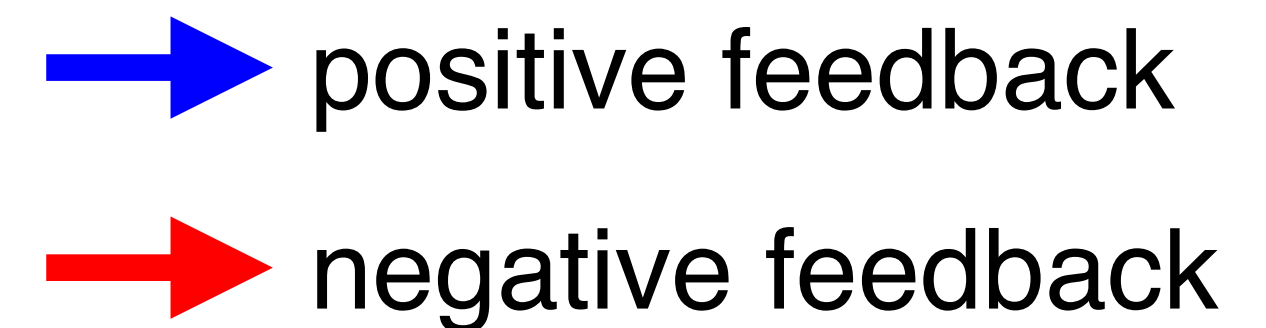
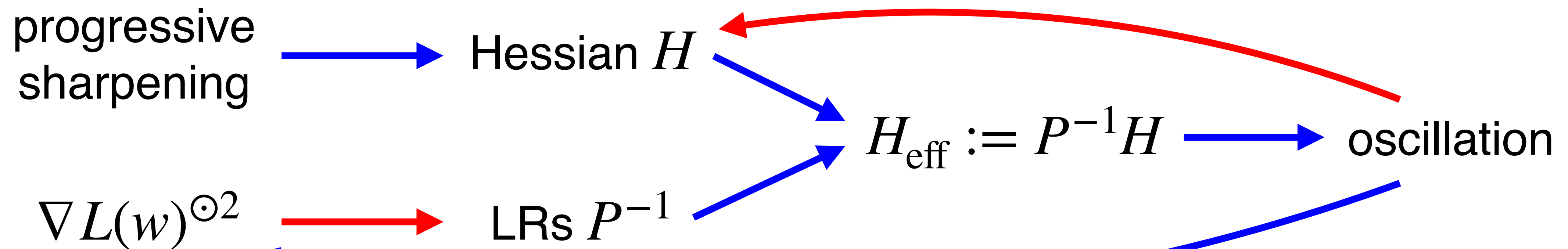
$$\nu = \text{EMA}(\nabla L(w)^{\odot 2}) \quad P_t^{-1} = \text{diag}\left(\frac{\eta}{\sqrt{\nu}}\right) \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$

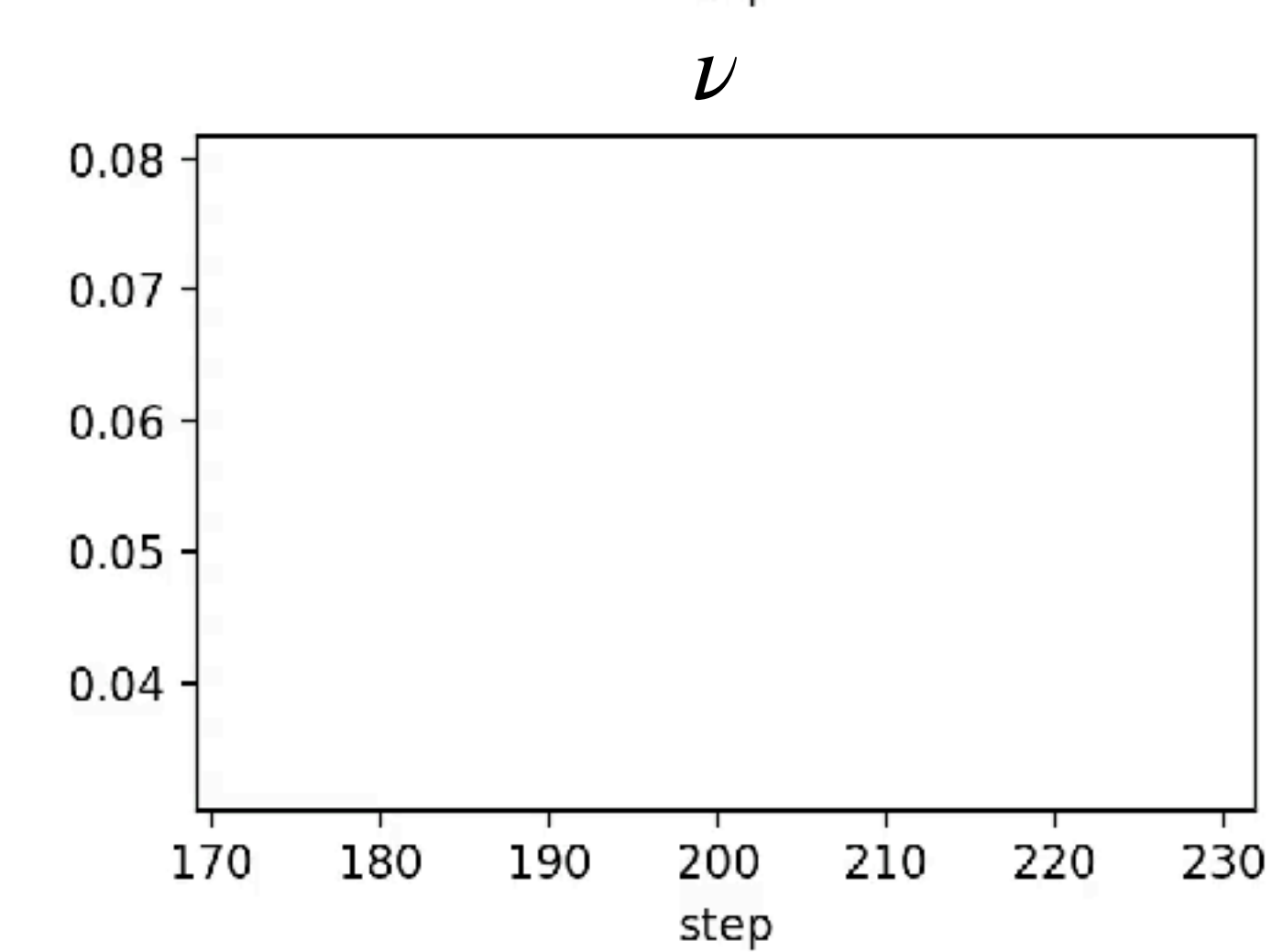
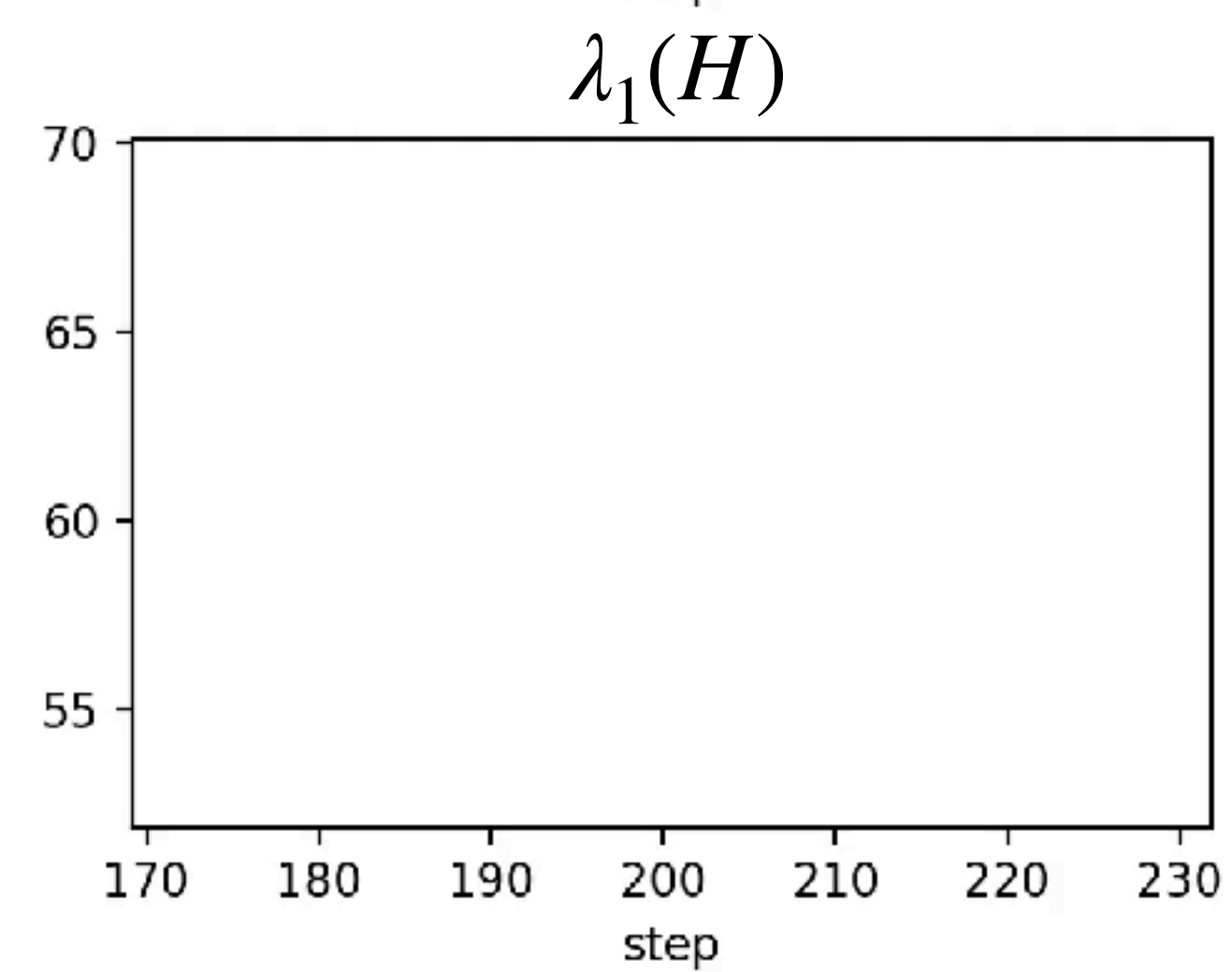
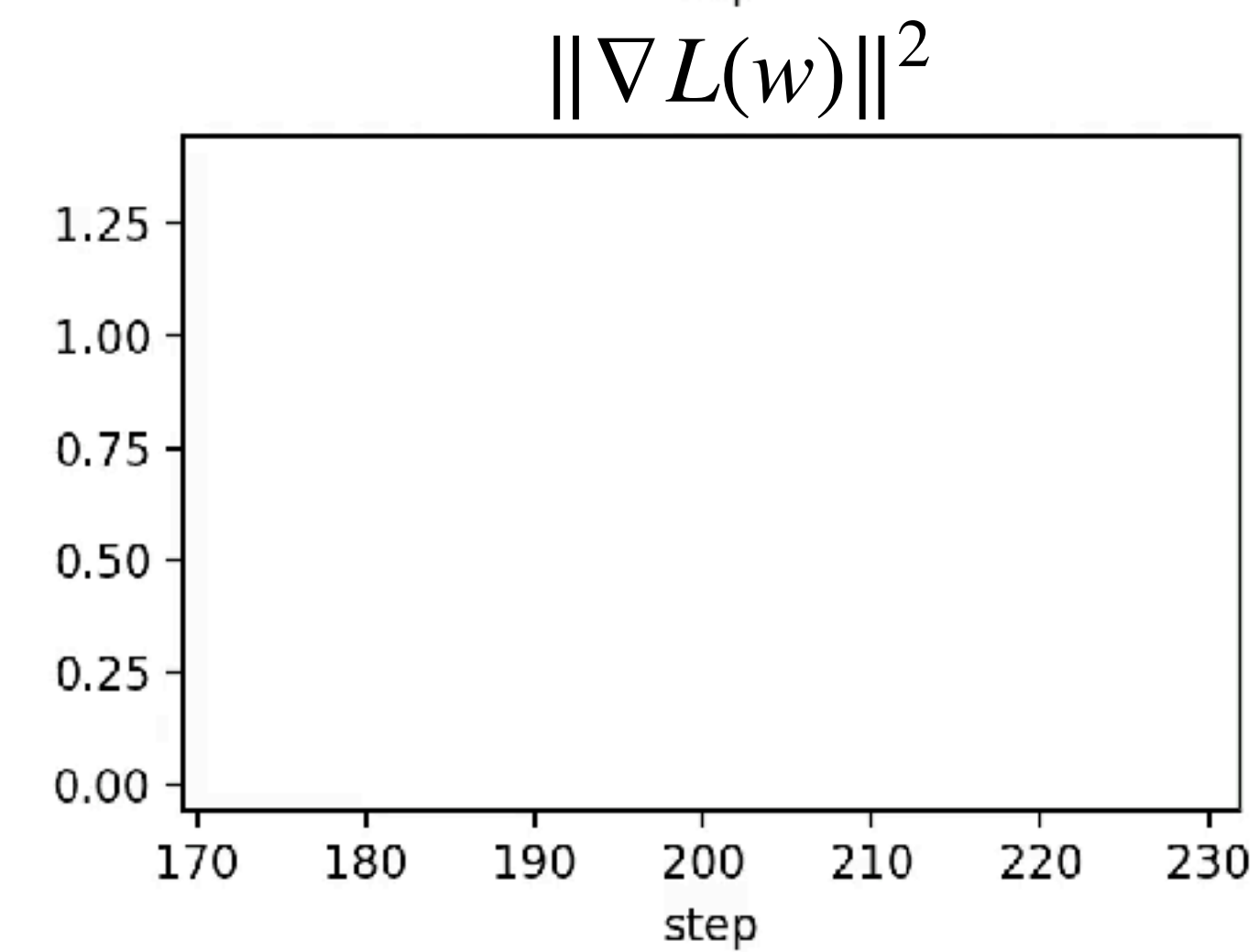
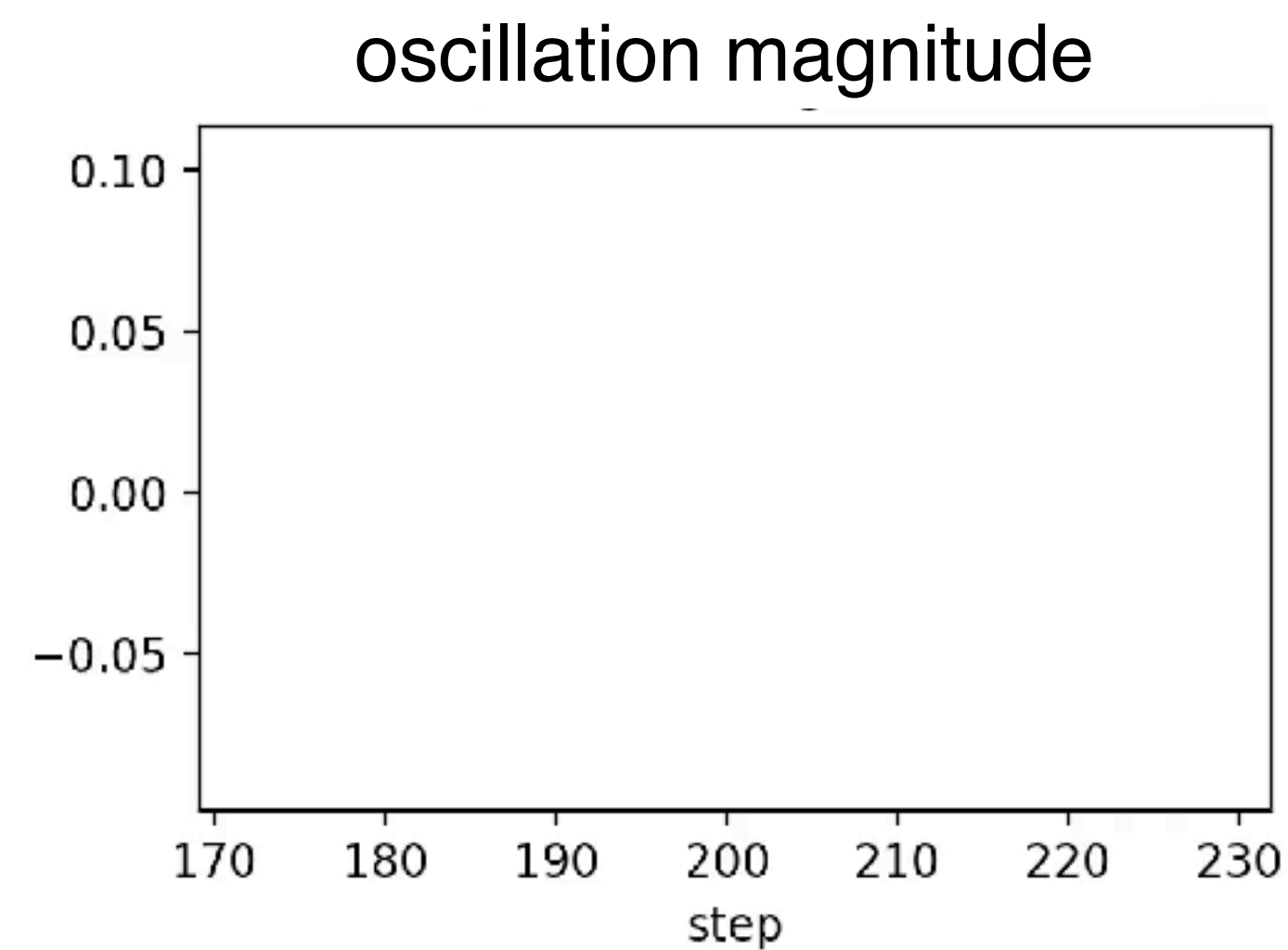
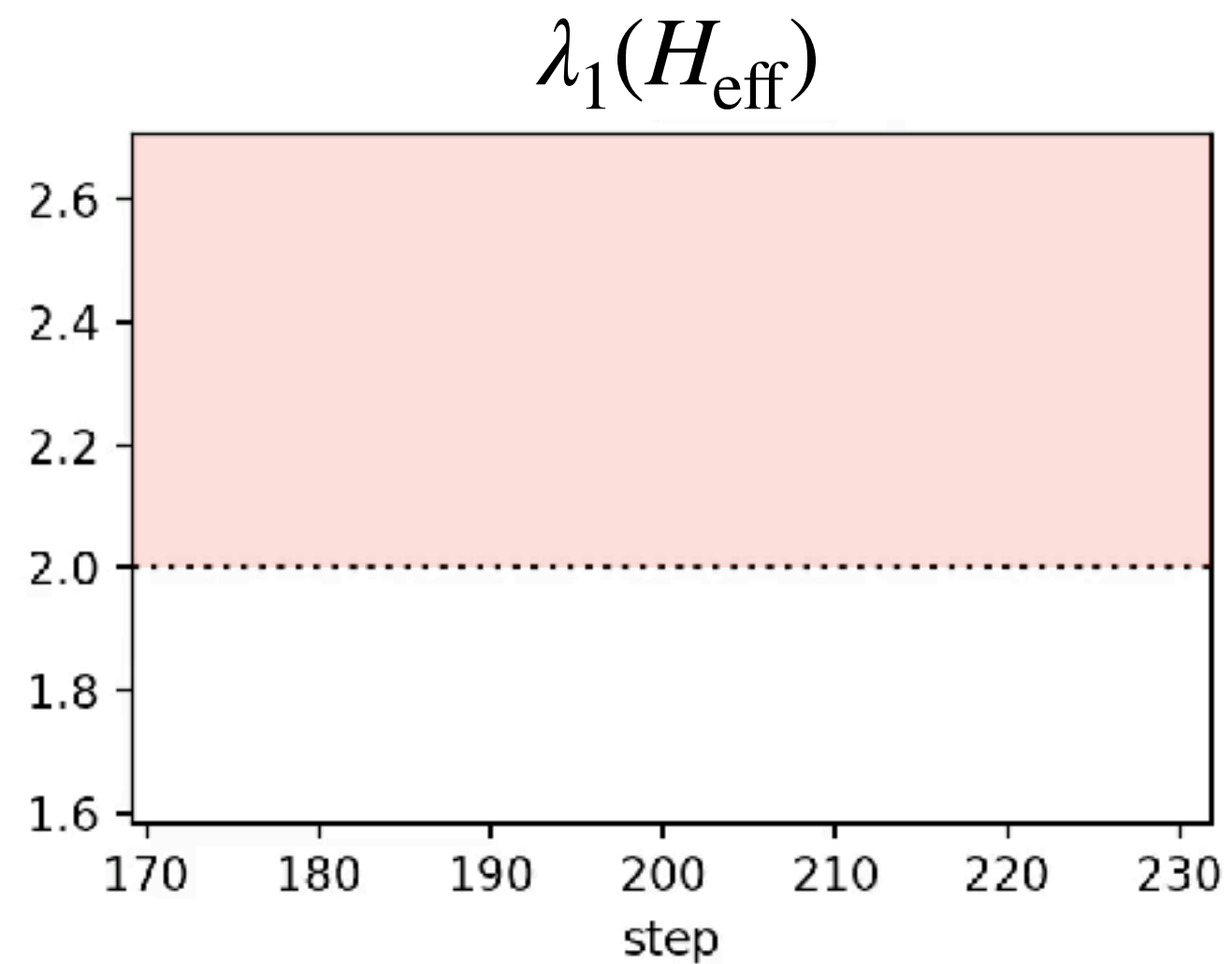
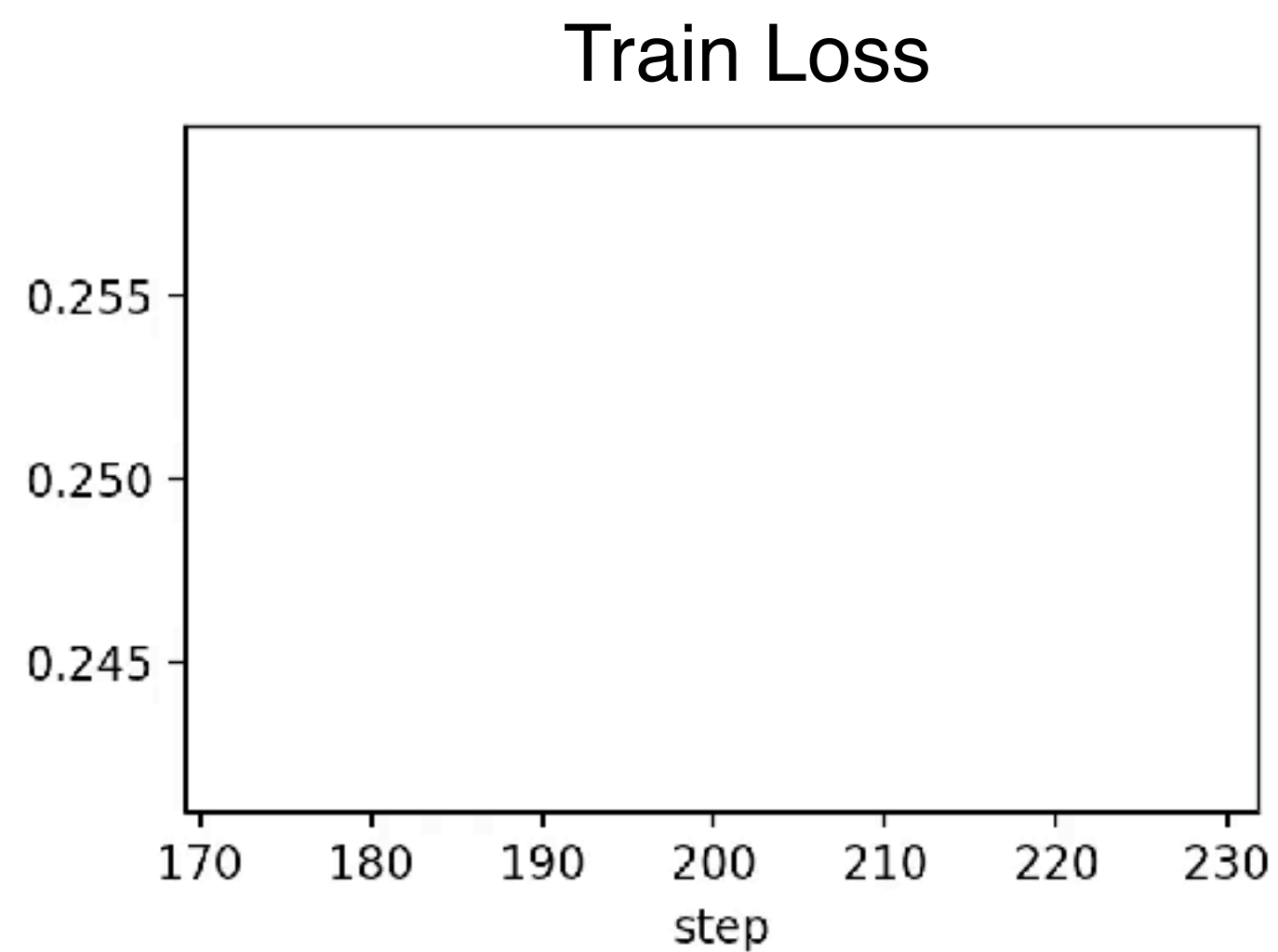
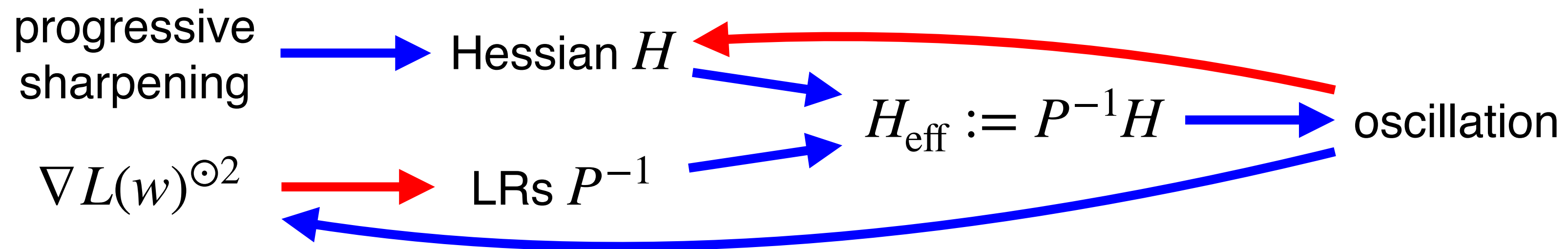


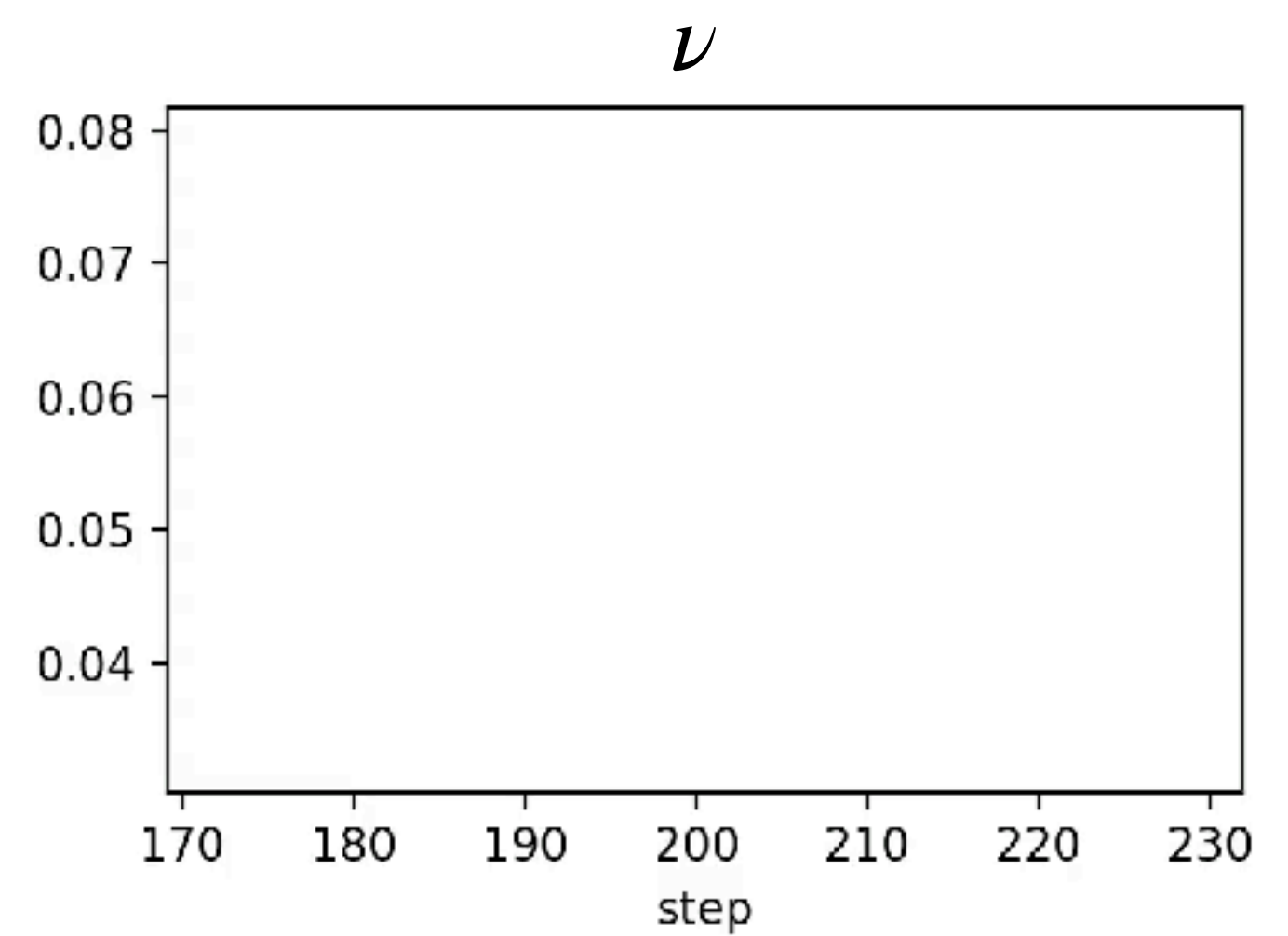
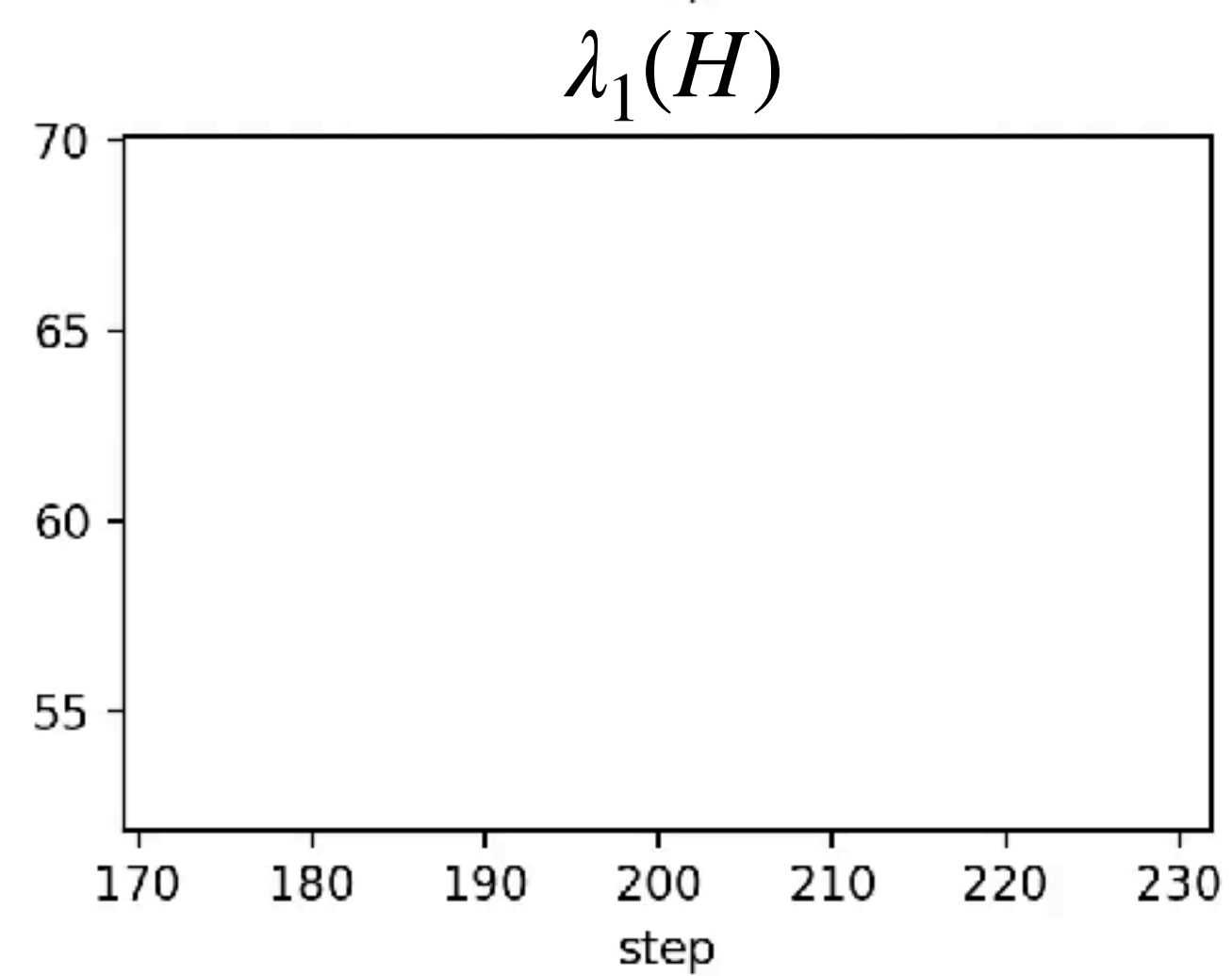
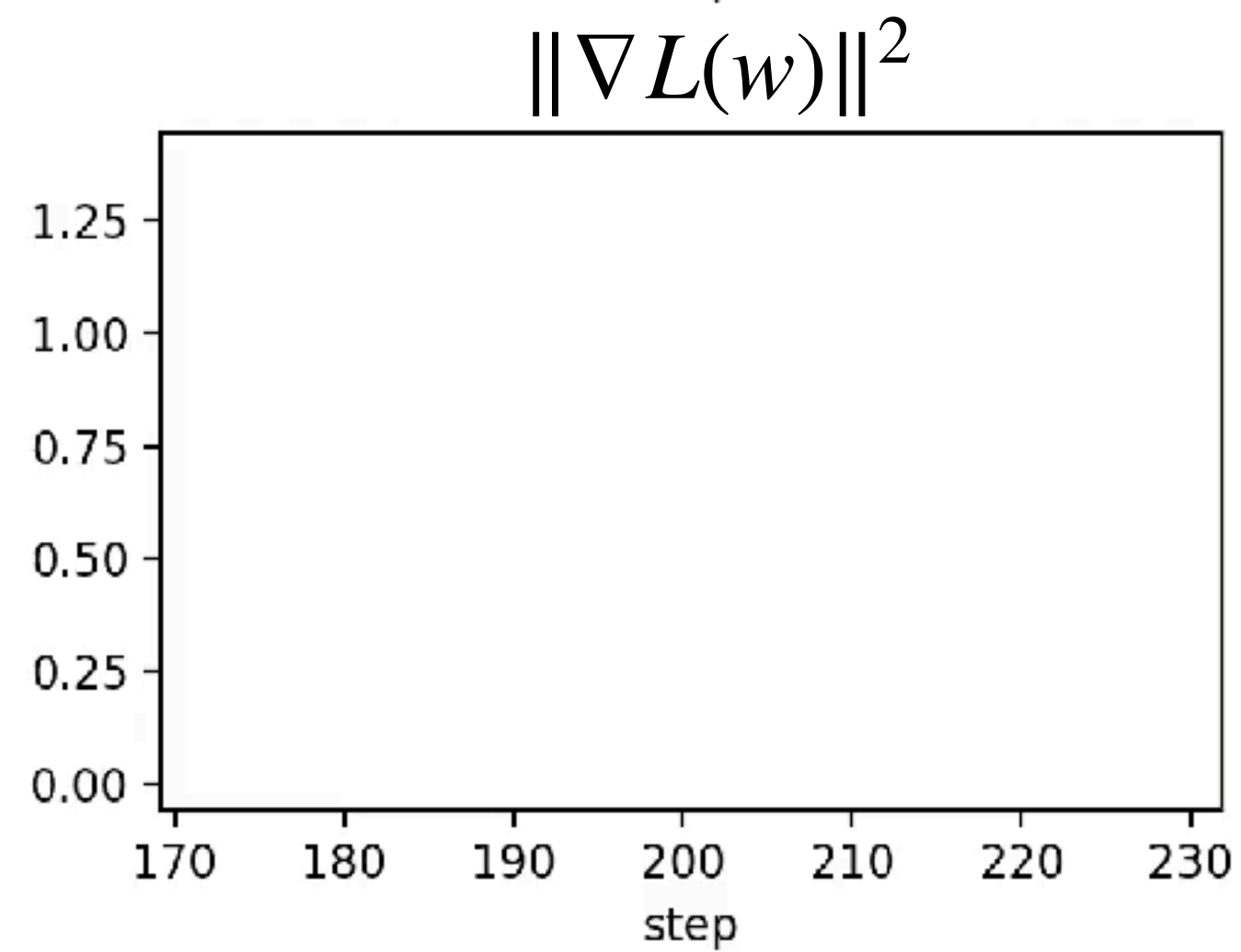
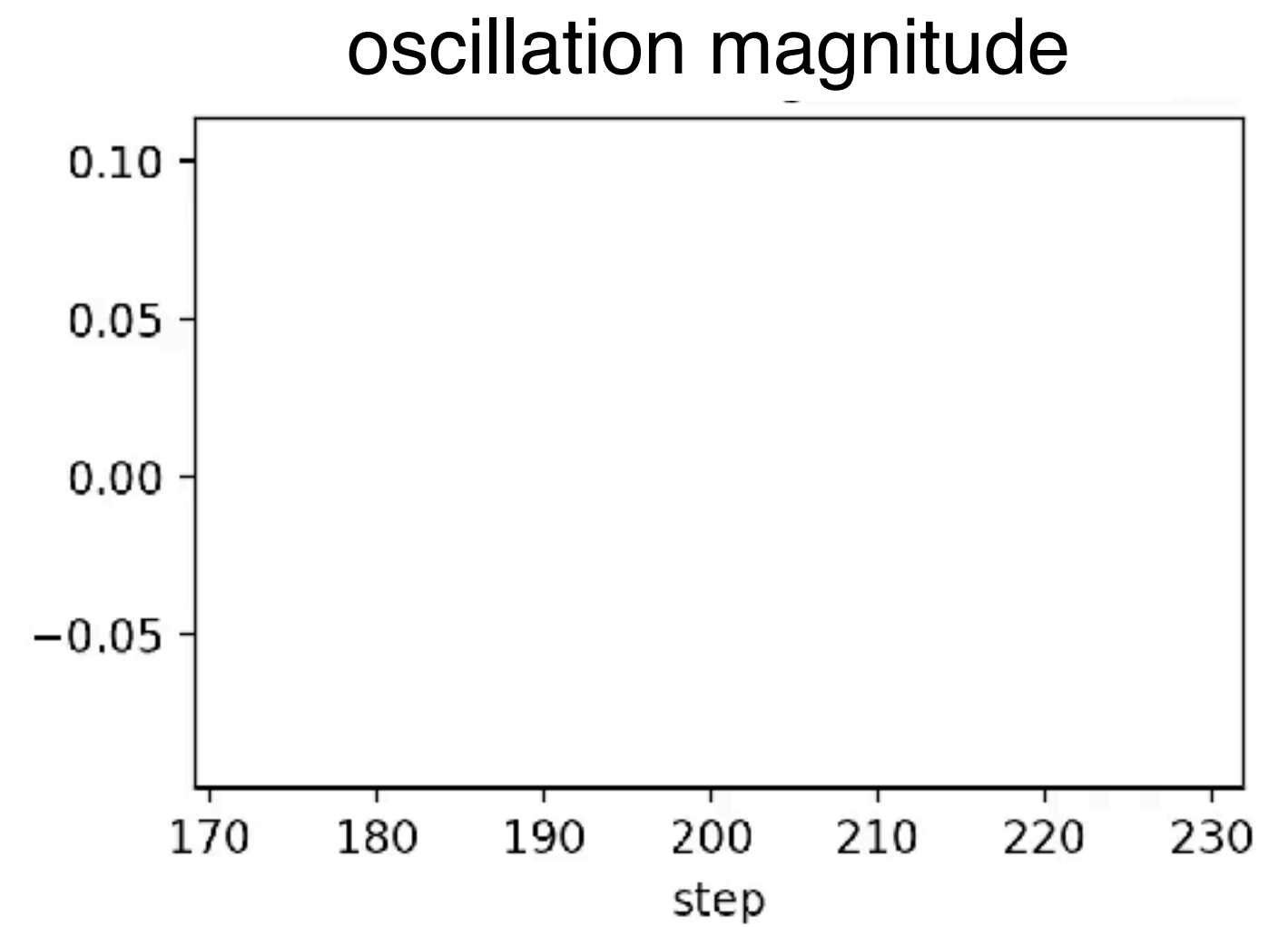
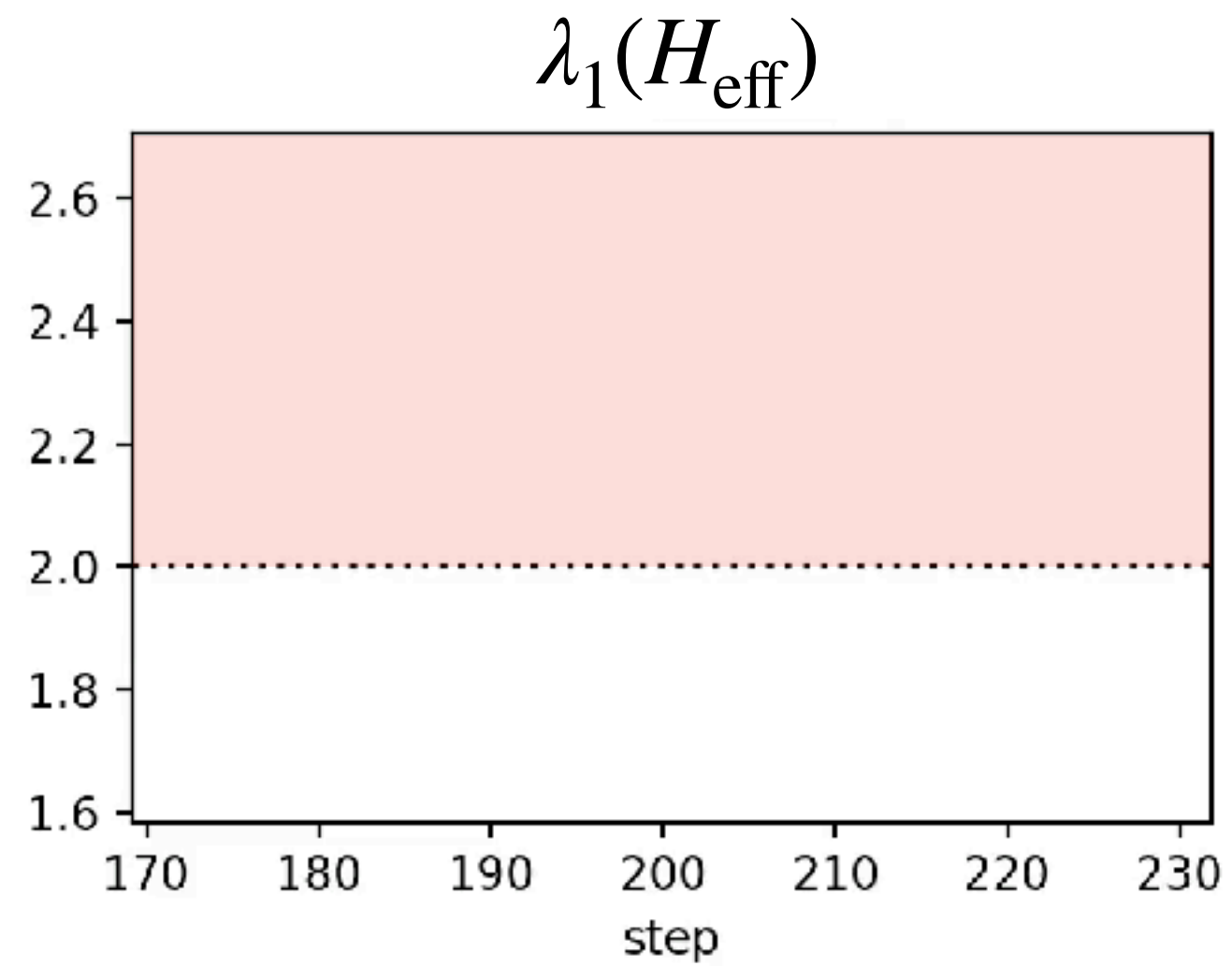
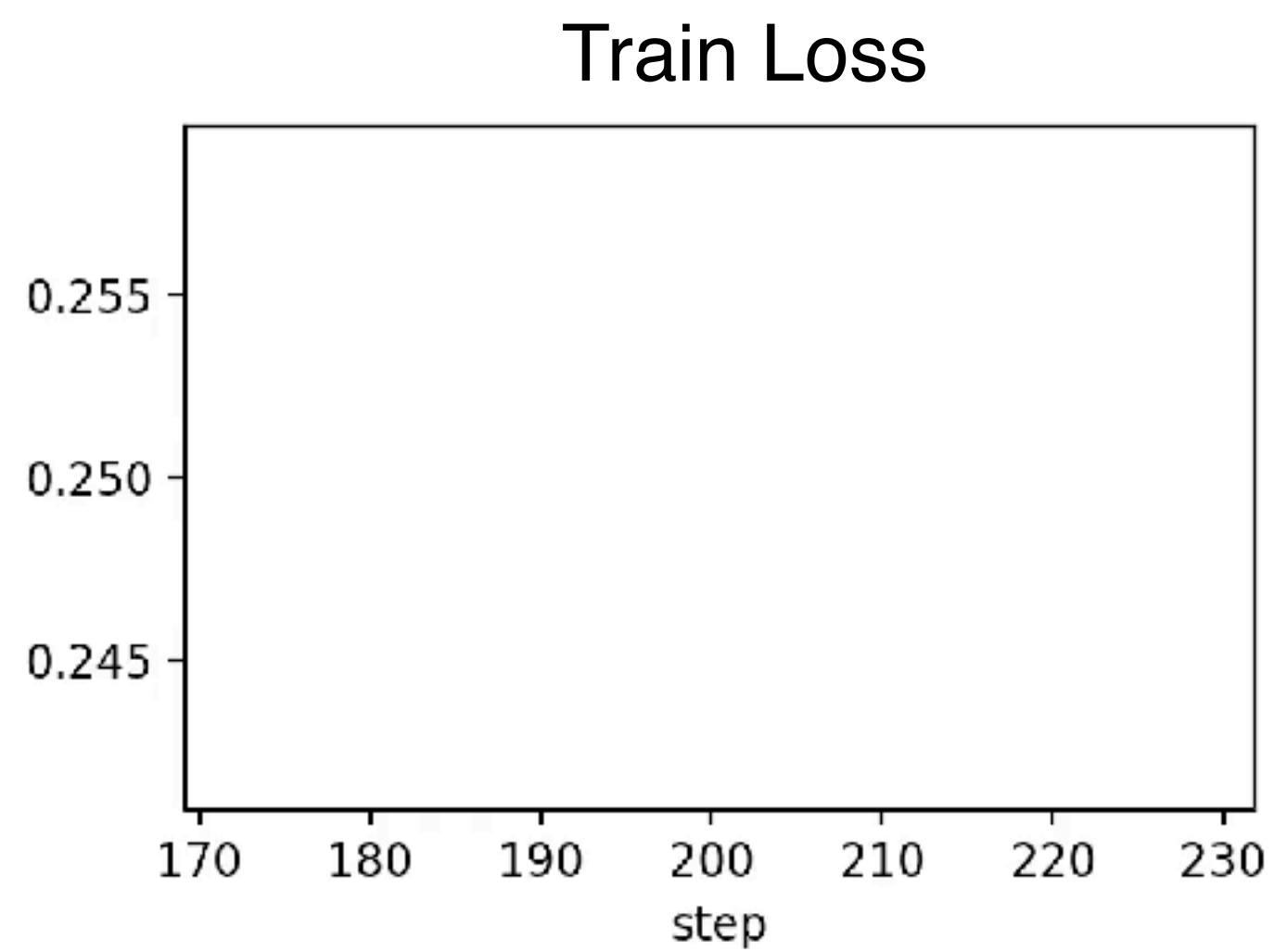
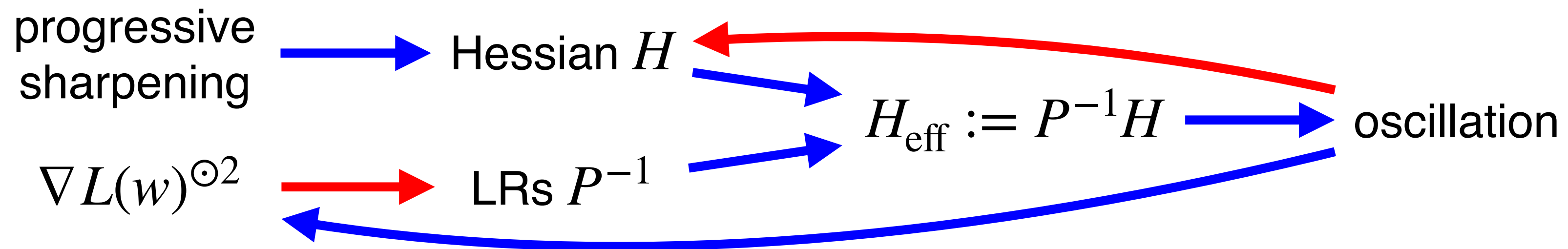
- ➡ positive feedback
- ➡ negative feedback

(Scalar) RMSProp at the Edge of Stability

$$\nu = \text{EMA}(\|\nabla L(w)\|^2) \quad P_t^{-1} = \frac{\eta}{\sqrt{\nu}} \quad w_{t+1} = w_t - P_t^{-1} \nabla L(w_t)$$







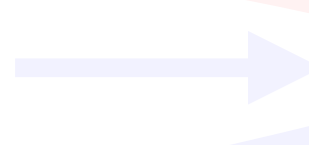
progressive
sharpening



Hessian H

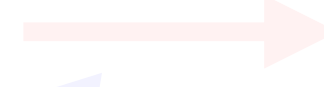


$H_{\text{eff}} := P^{-1}H$



oscillation

$\nabla L(w)^{\odot 2}$



LRs P^{-1}



Train Loss



$\lambda_1(H_{\text{eff}})$



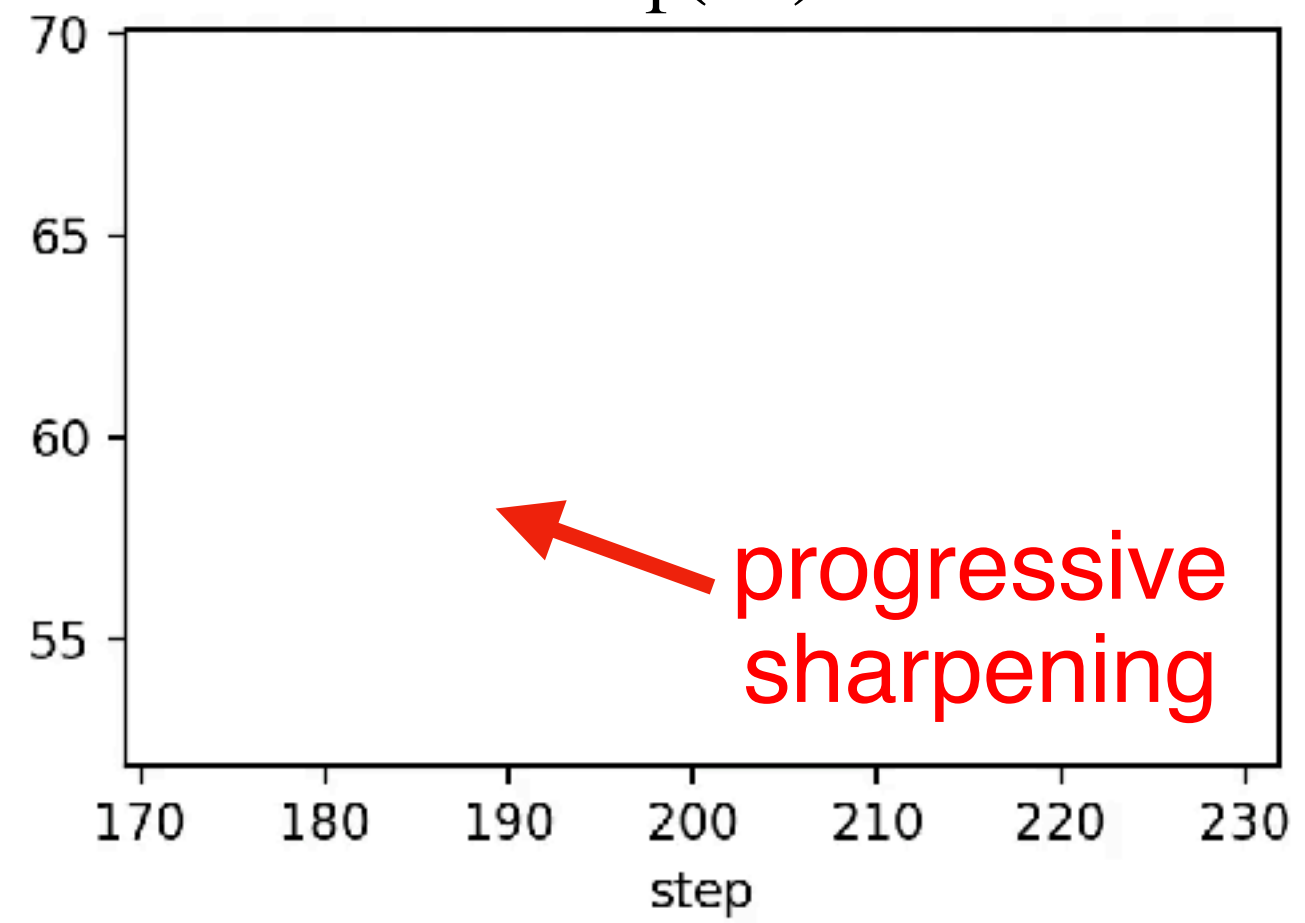
oscillation magnitude



$\|\nabla L(w)\|^2$



$\lambda_1(H)$



ν



progressive
sharpening



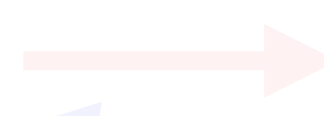
Hessian H



$H_{\text{eff}} := P^{-1}H$

oscillation

$\nabla L(w)^{\odot 2}$



LRs P^{-1}



Train Loss



$\lambda_1(H_{\text{eff}})$



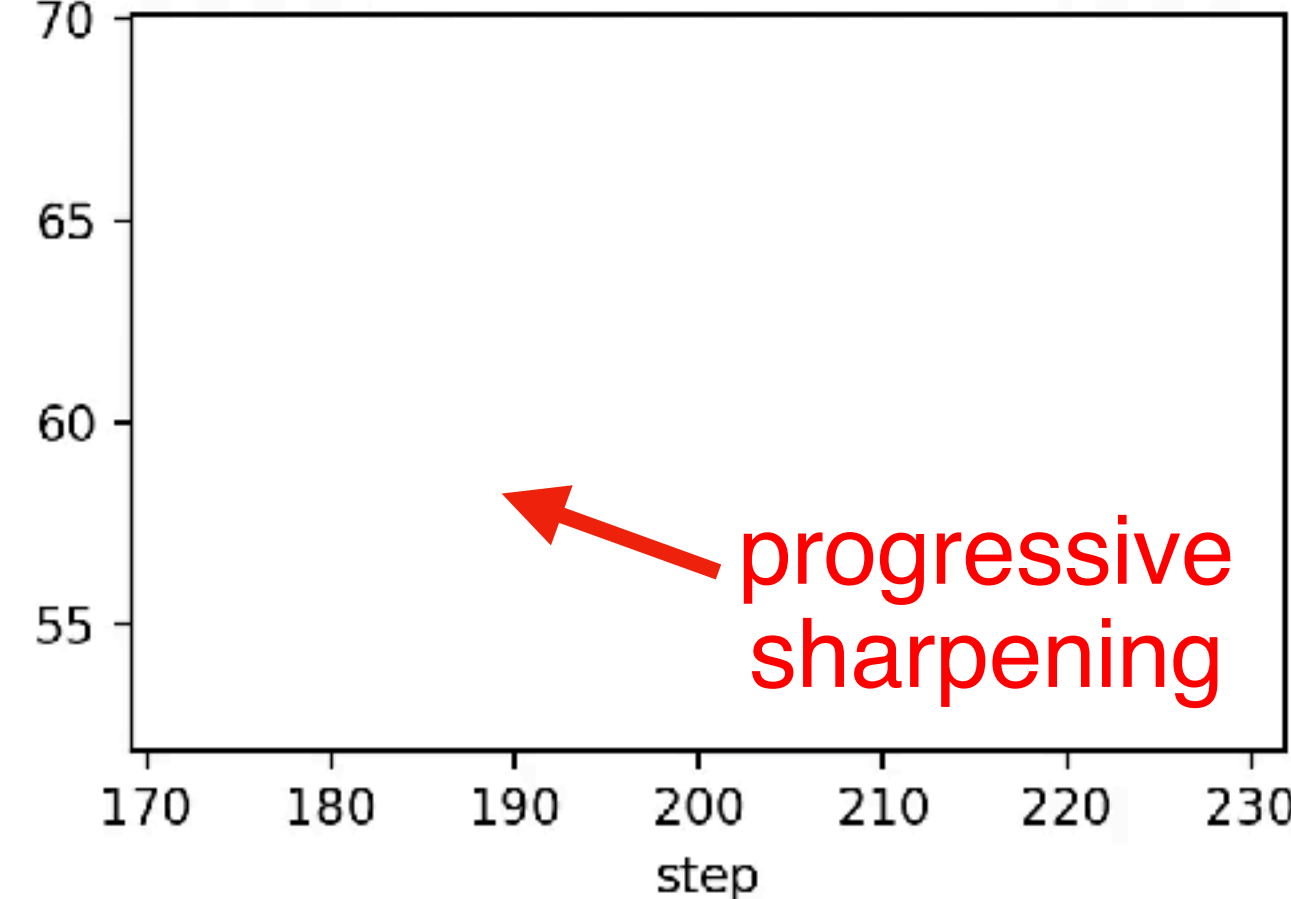
oscillation magnitude



$\|\nabla L(w)\|^2$

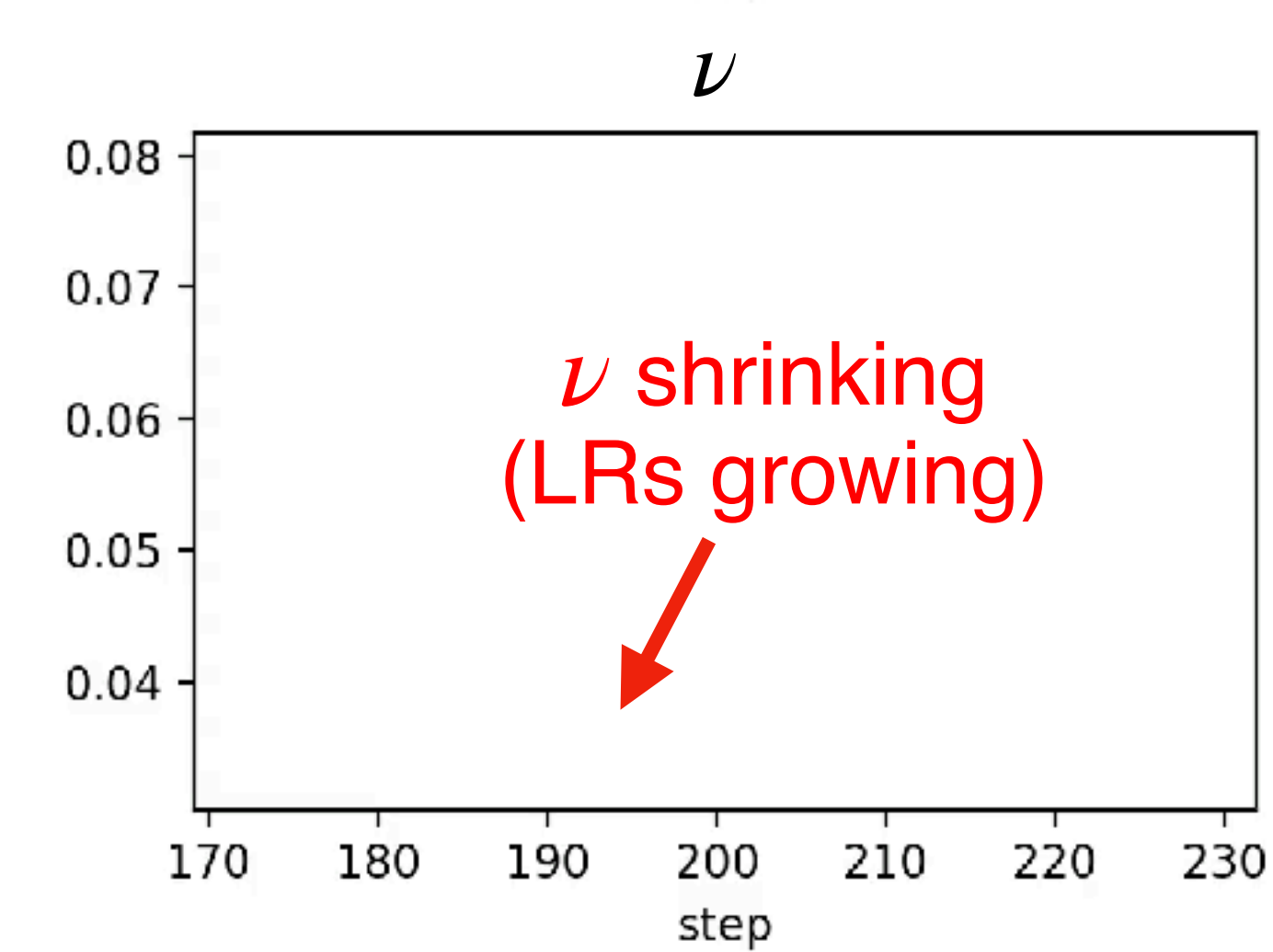
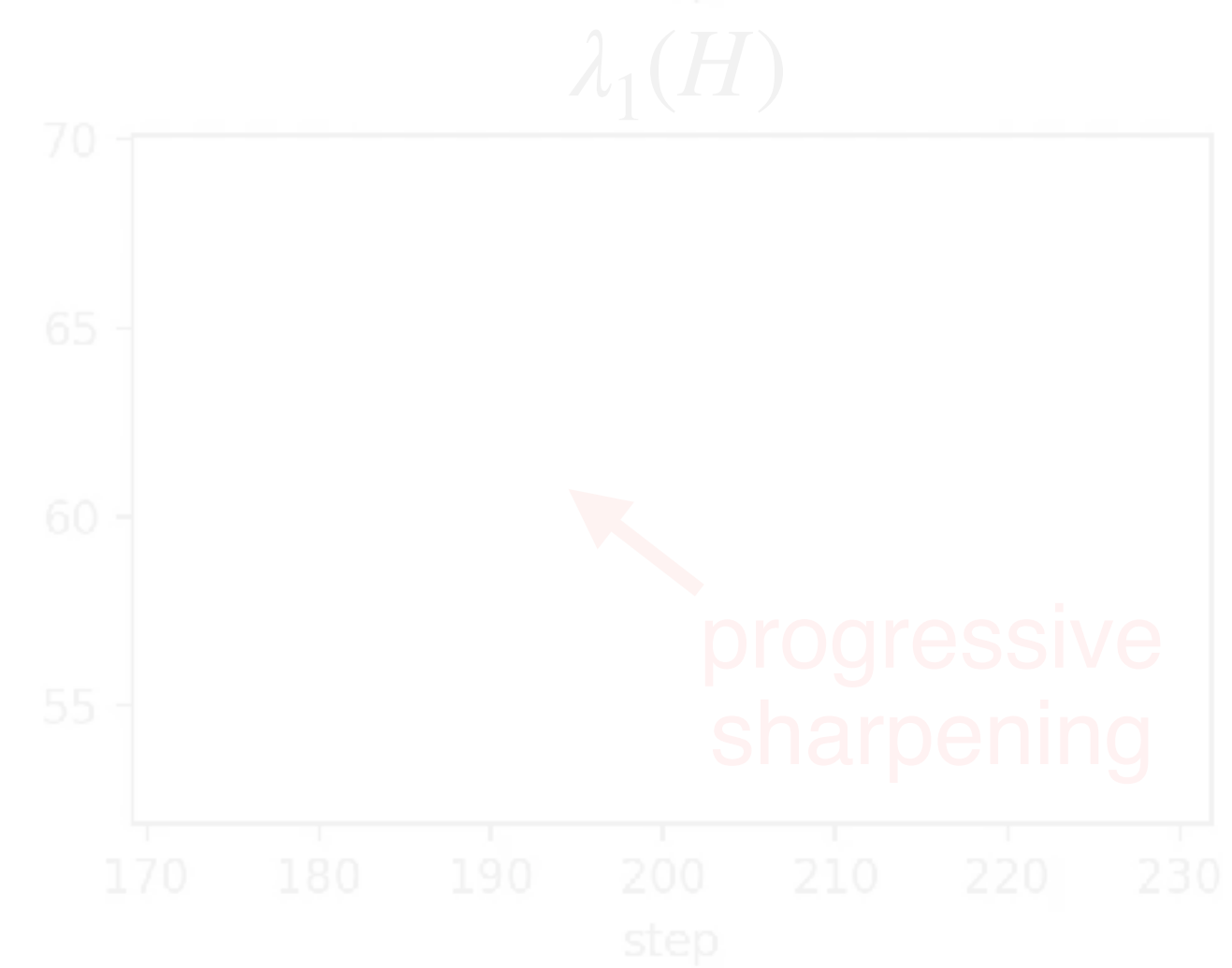
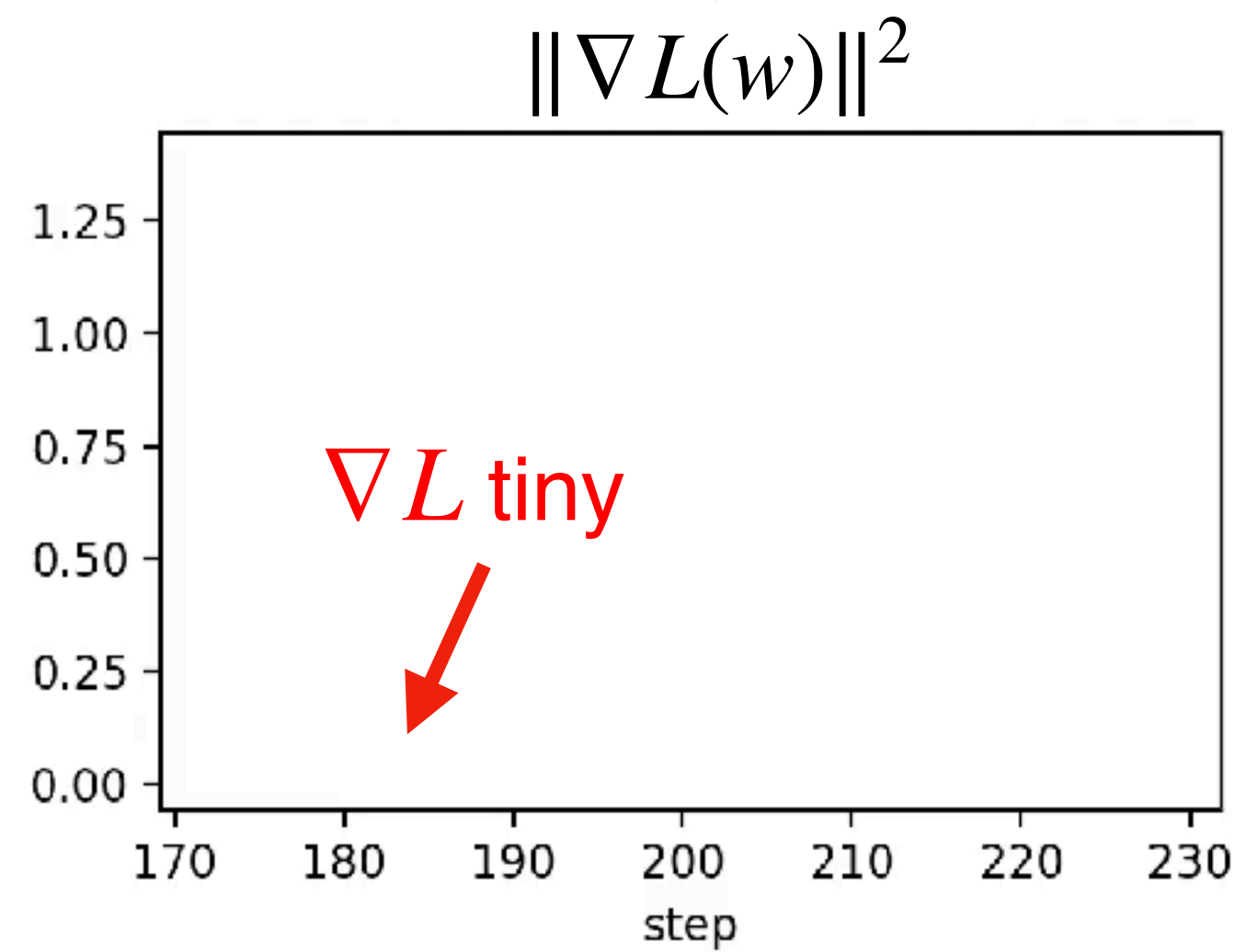
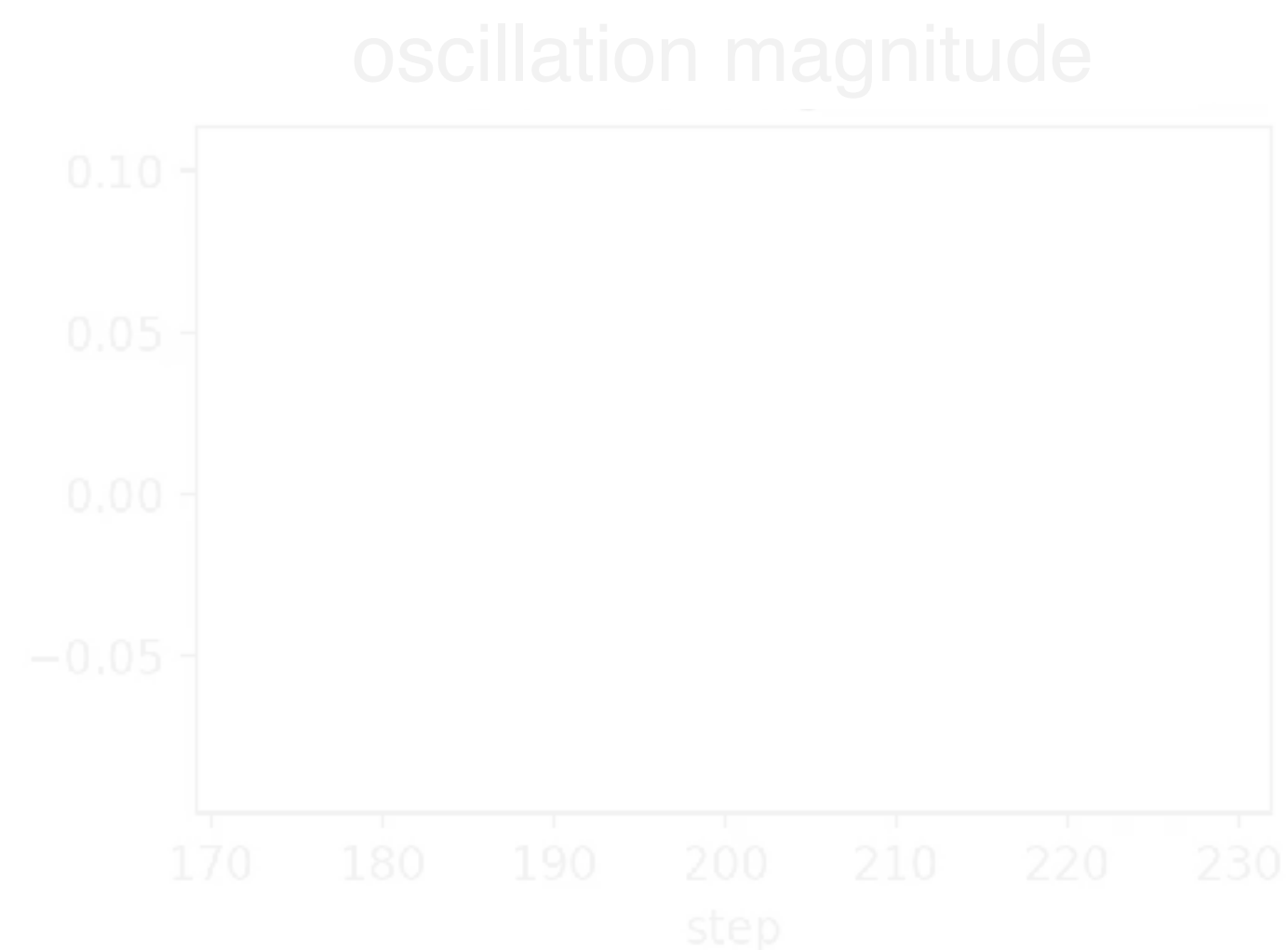
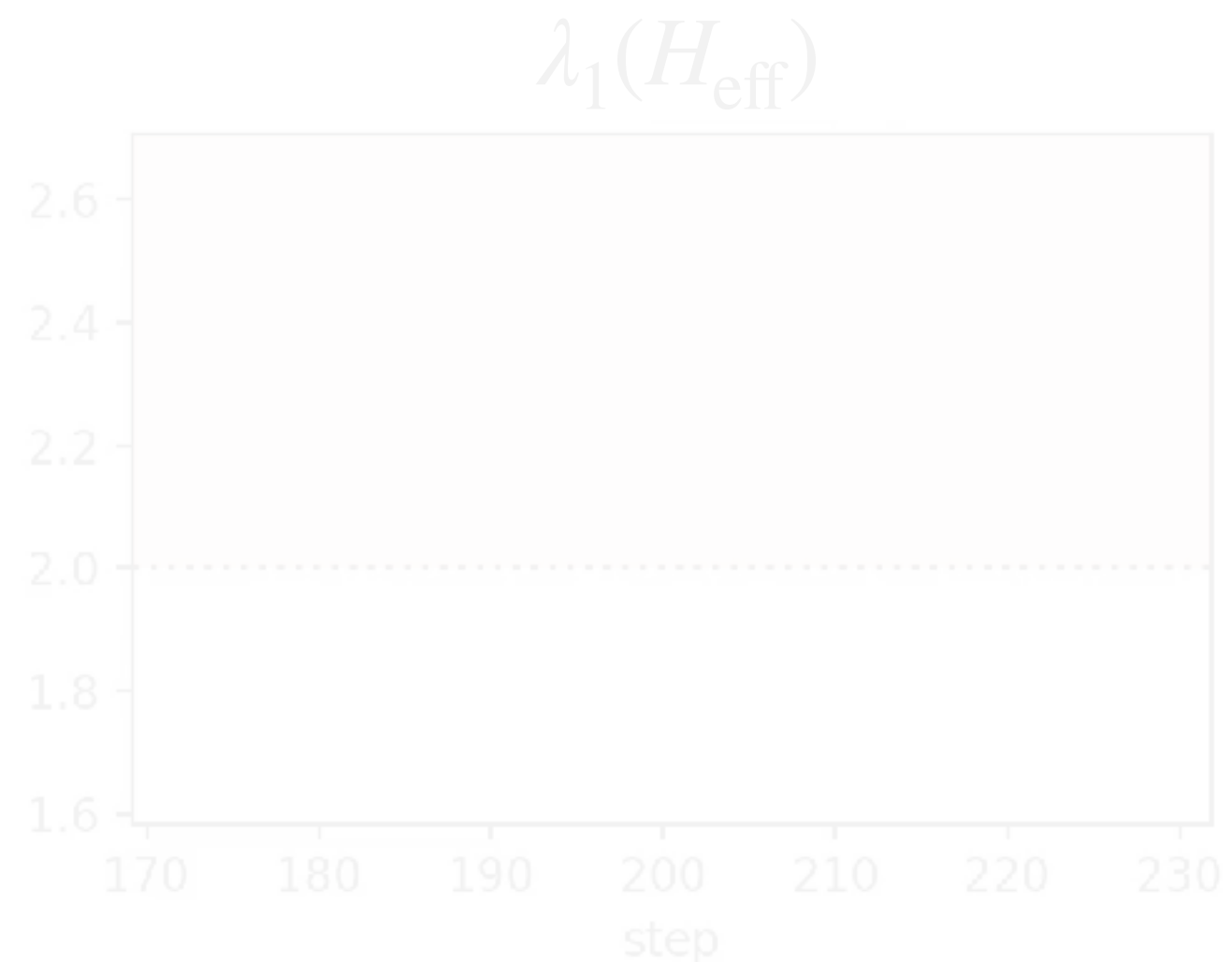
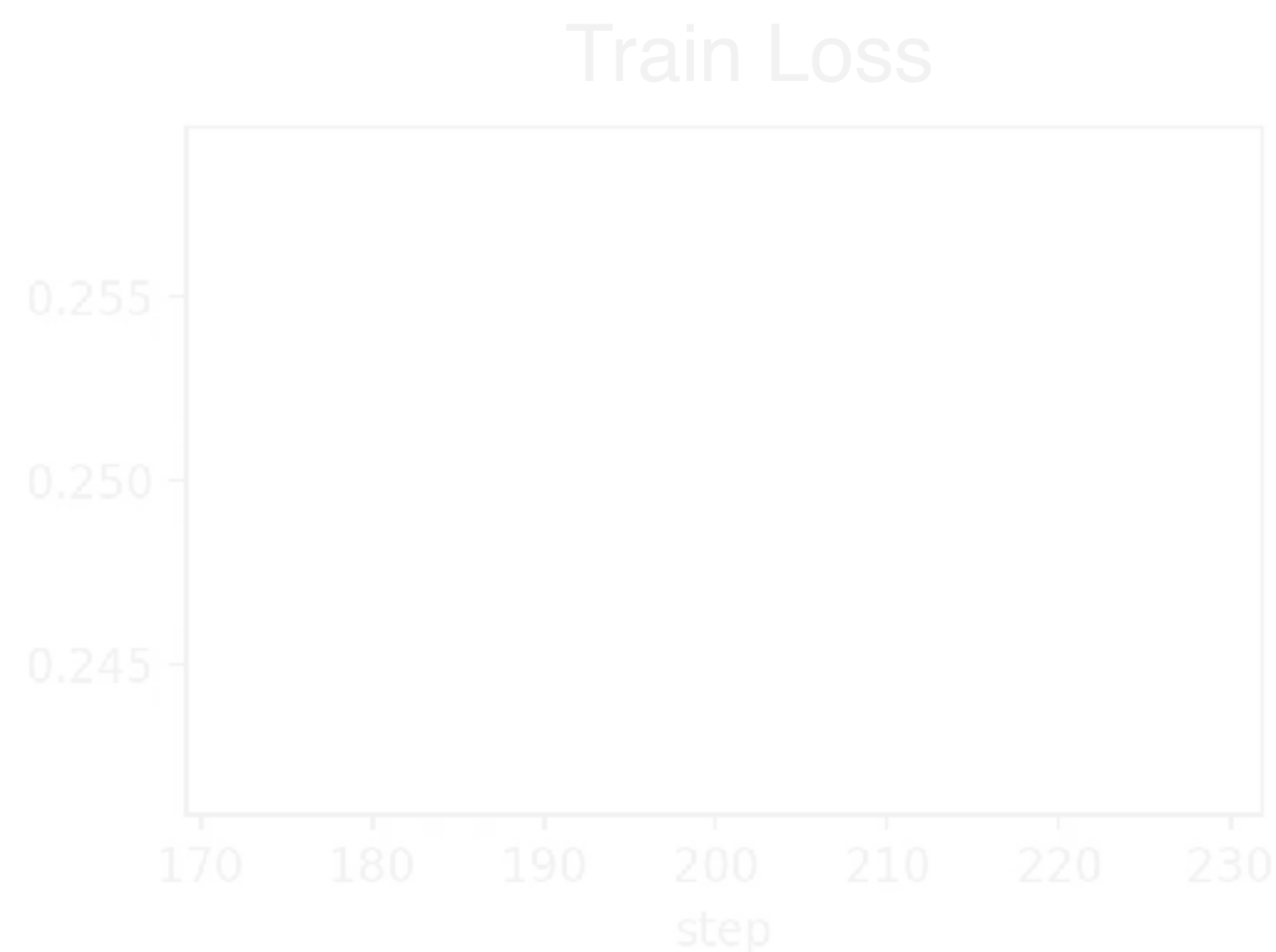
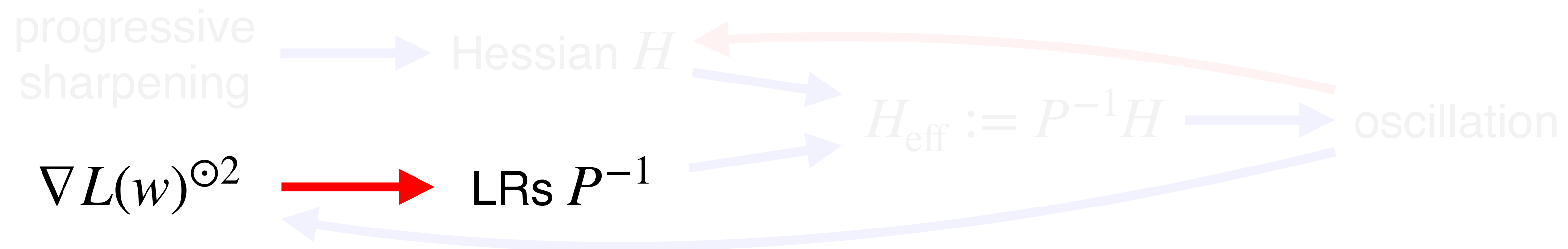


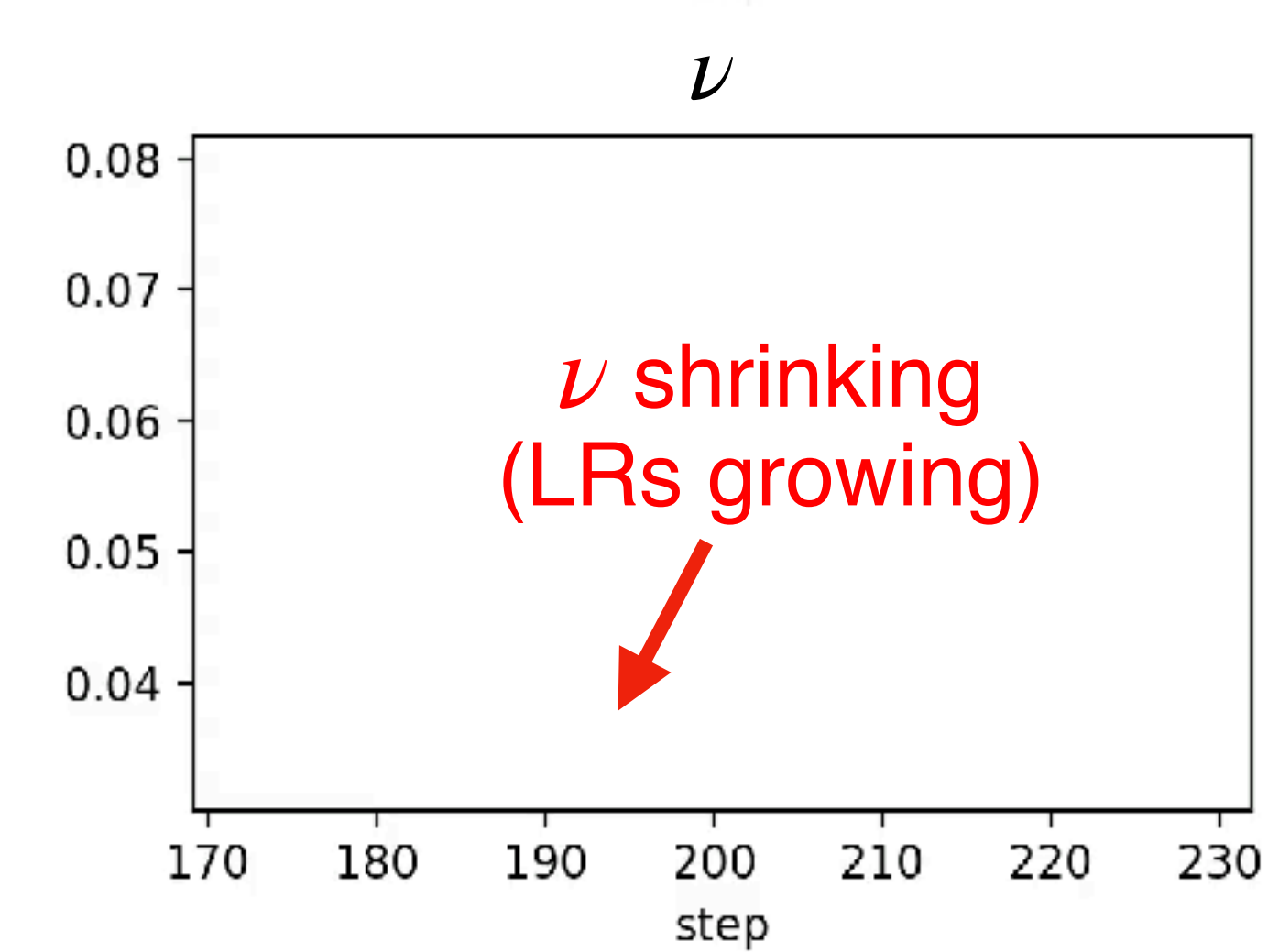
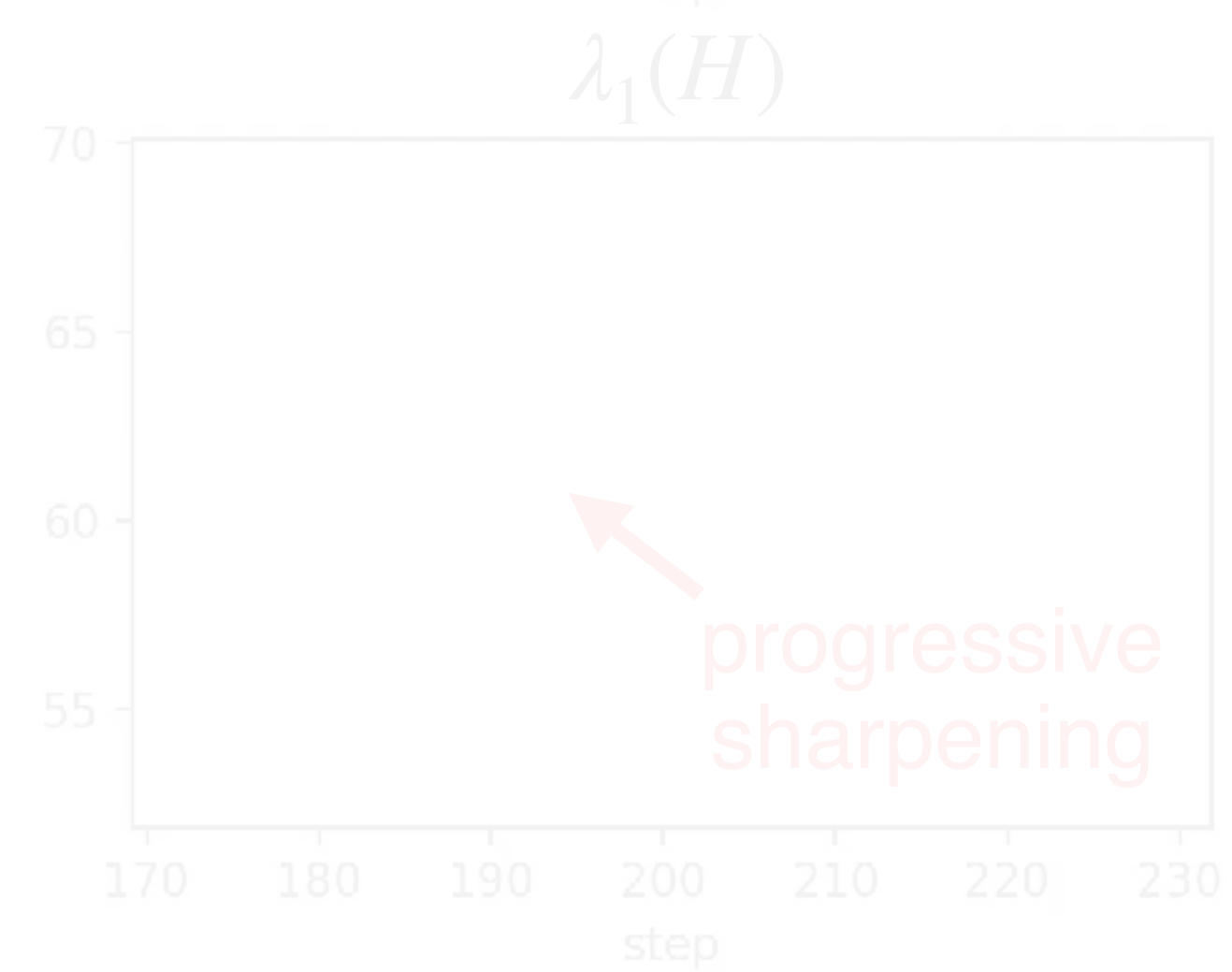
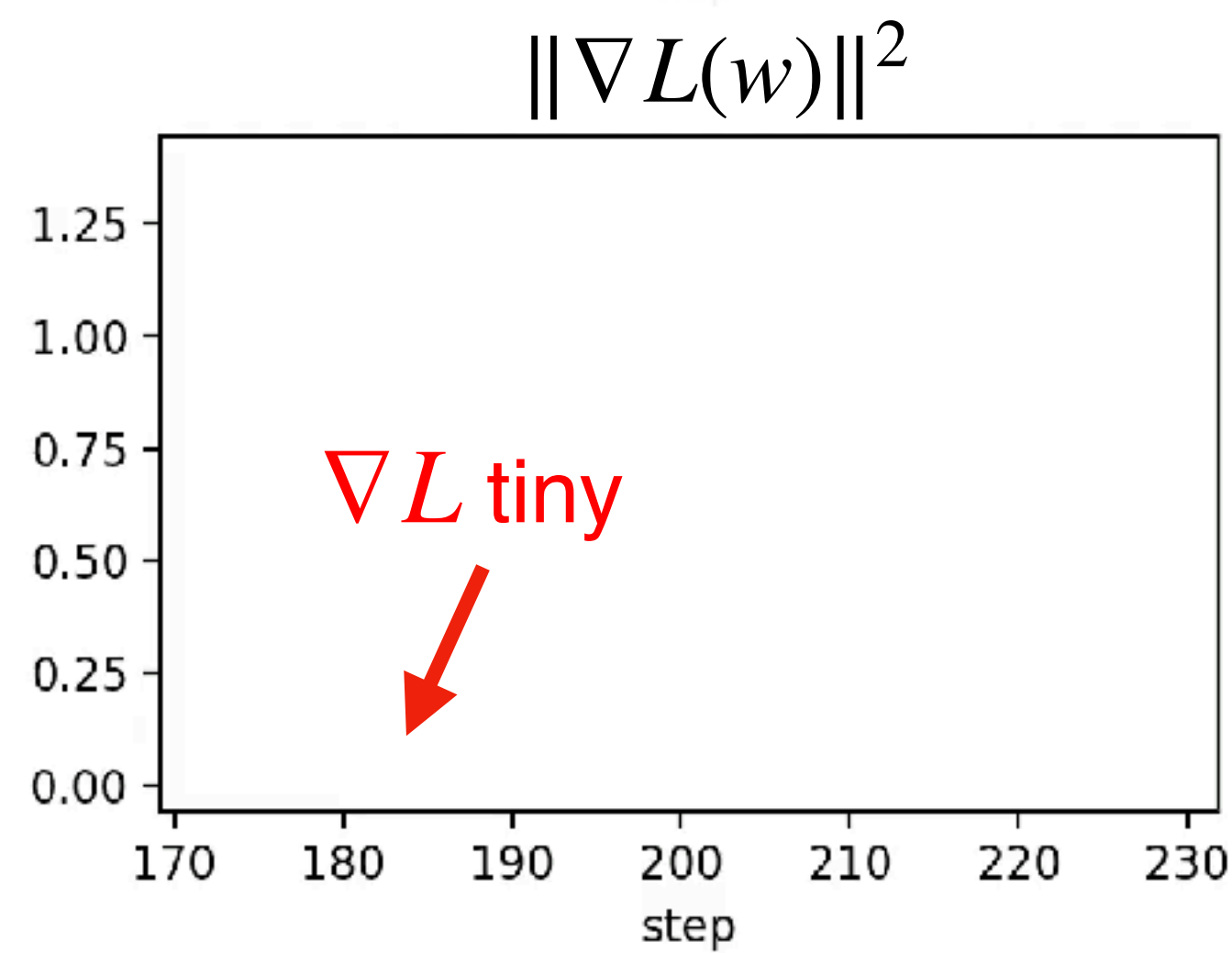
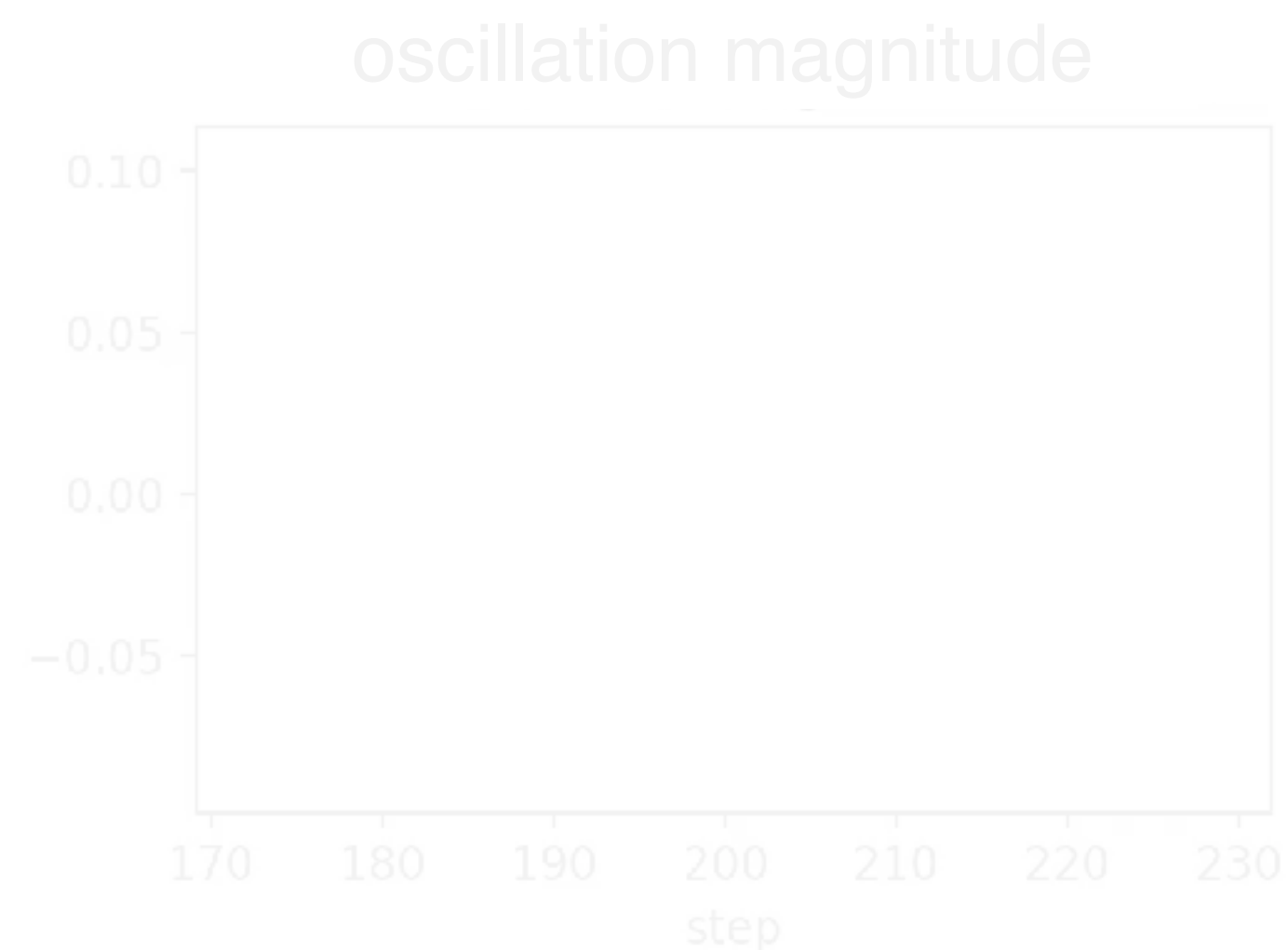
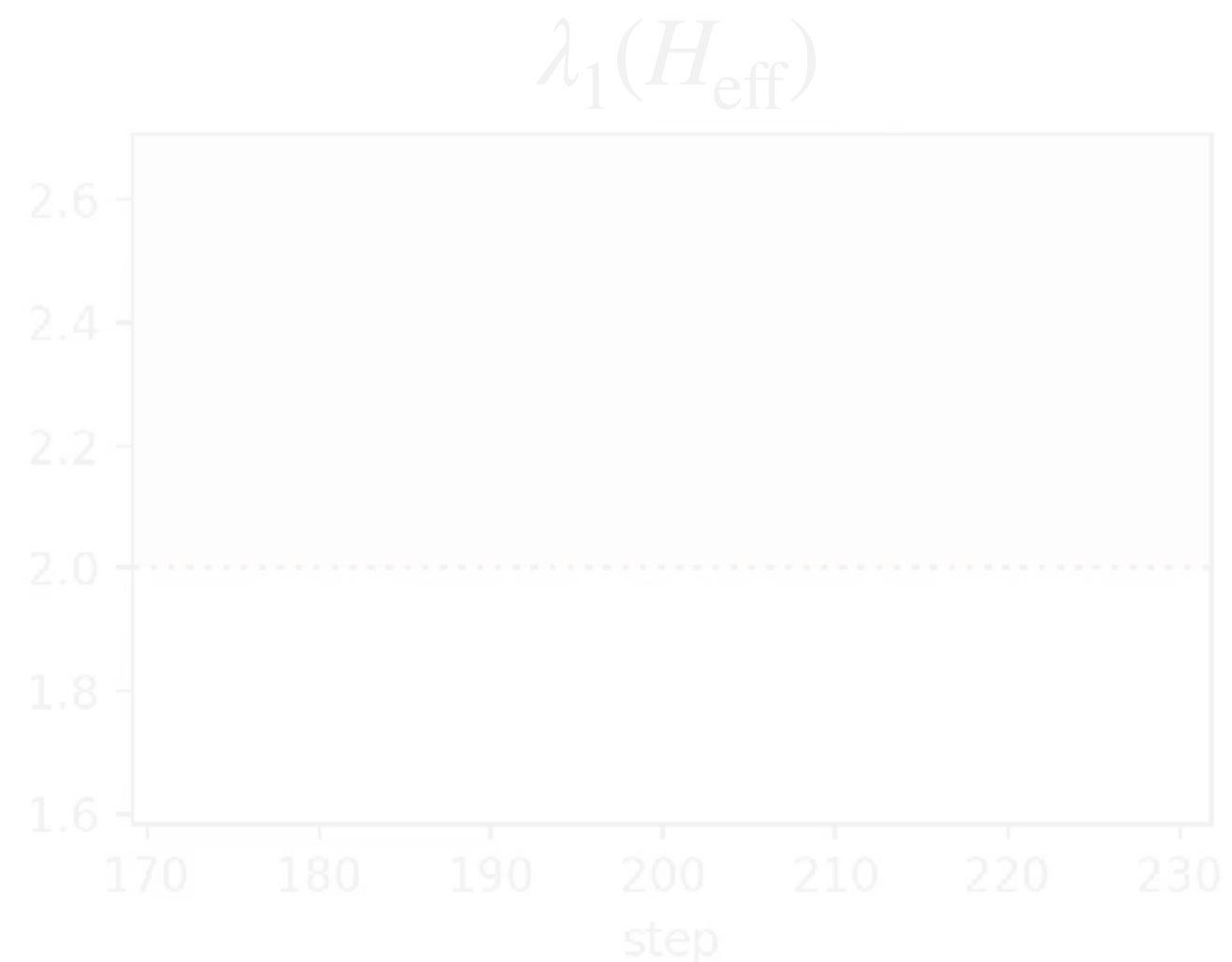
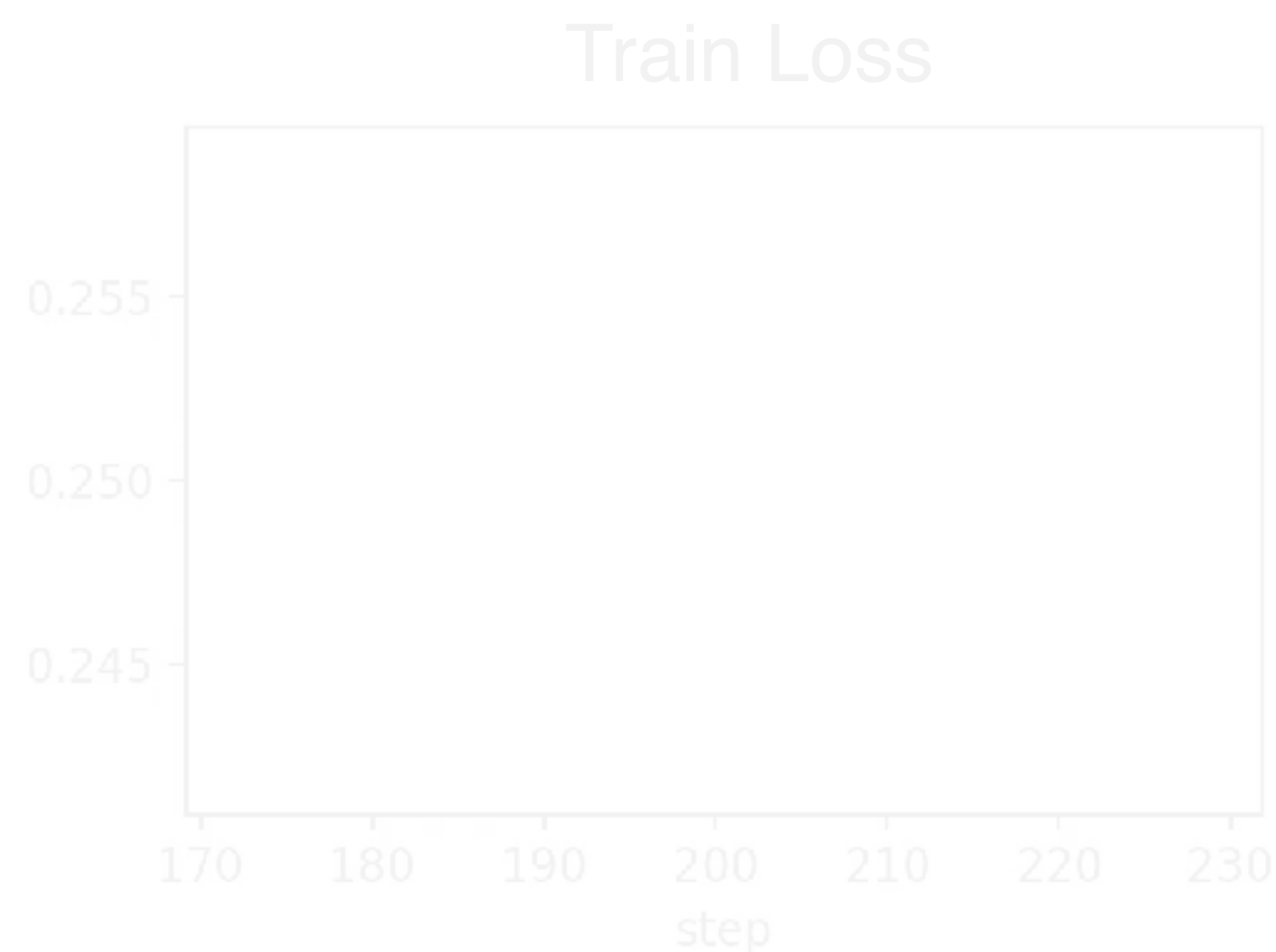
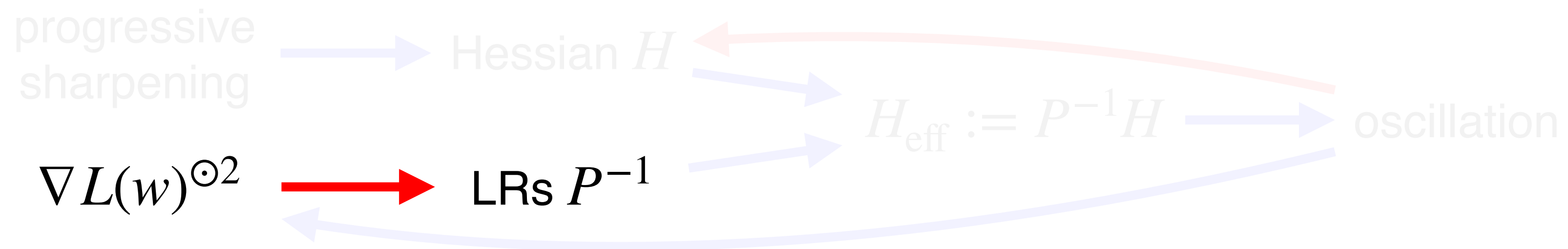
$\lambda_1(H)$

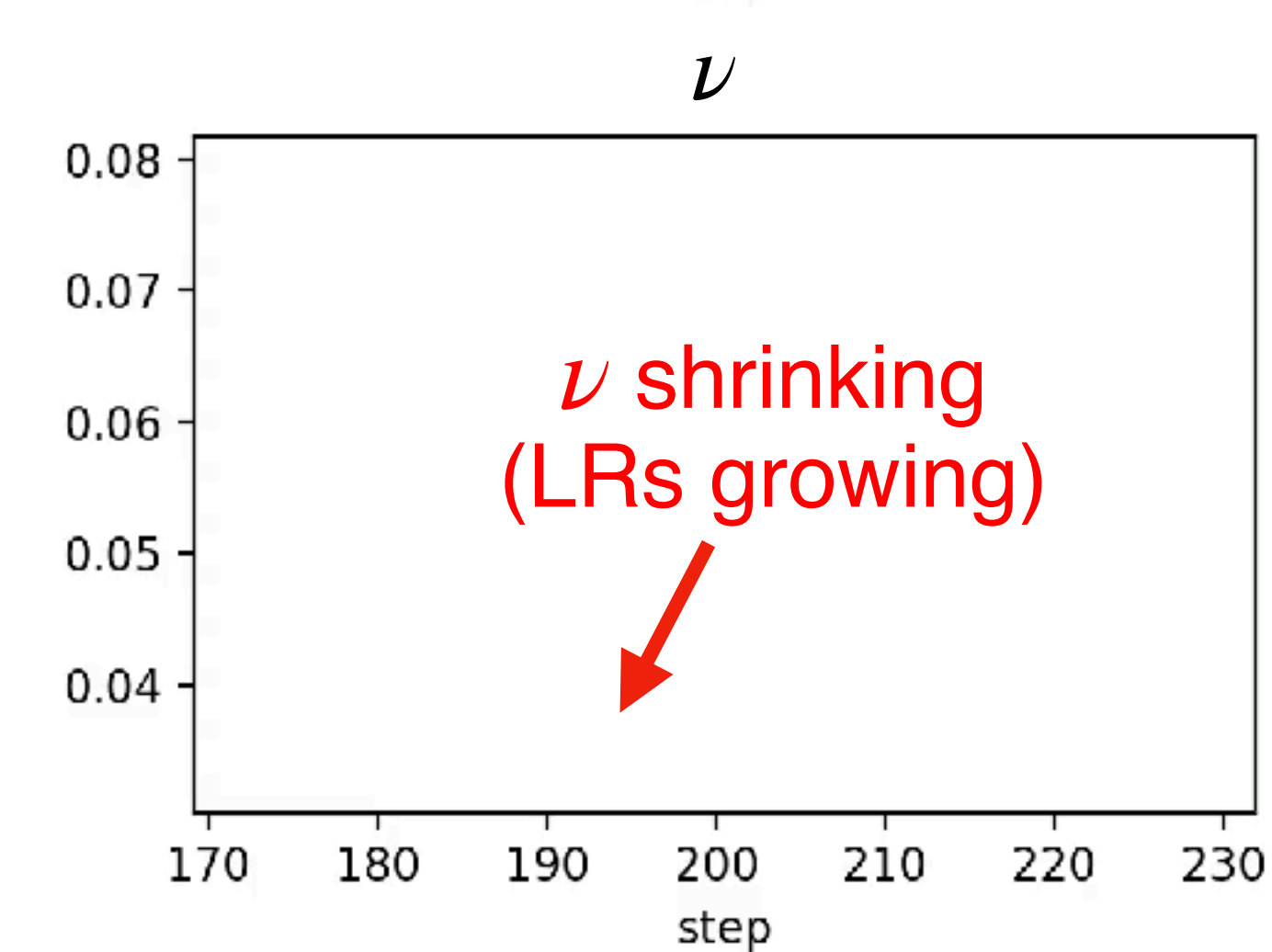
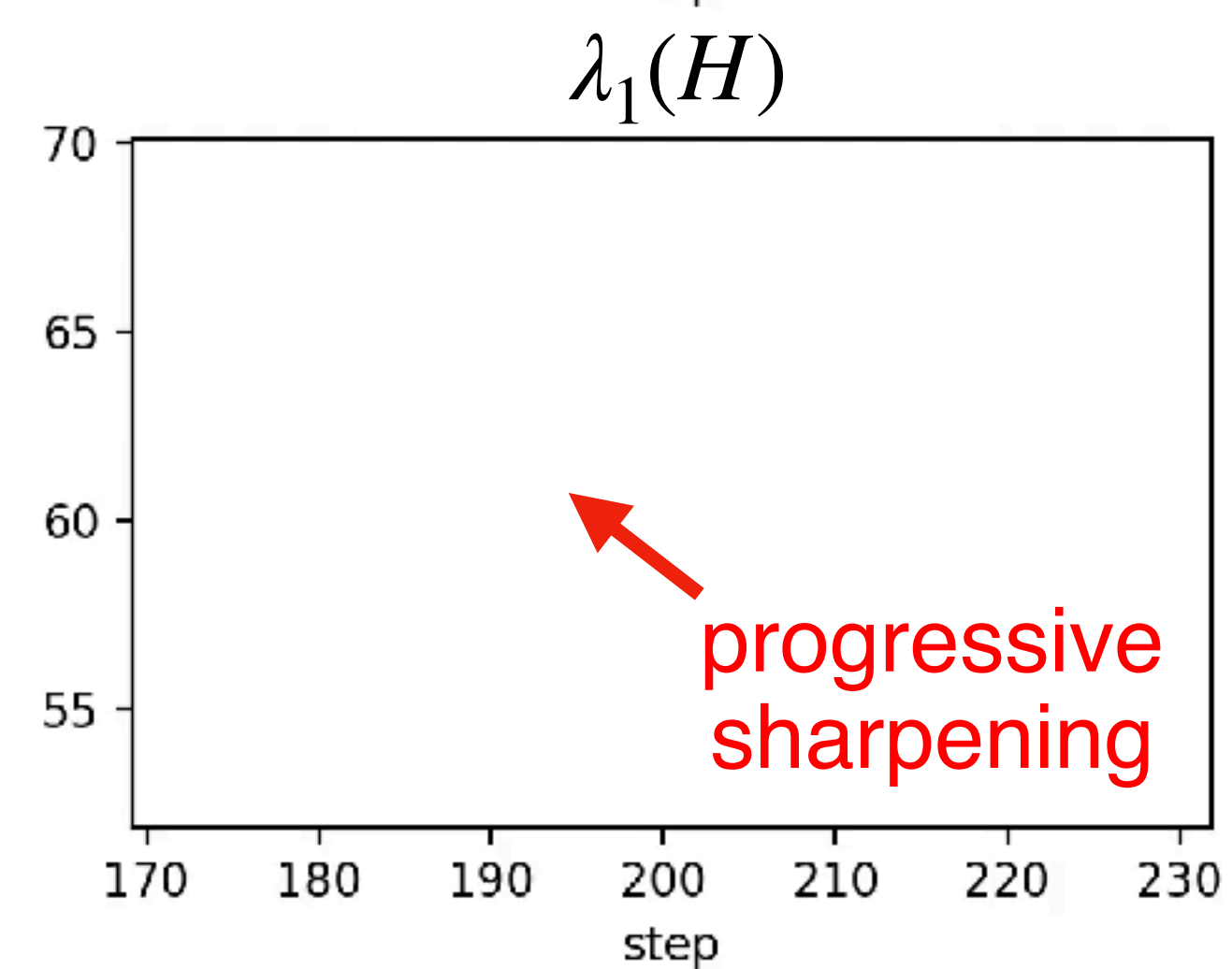
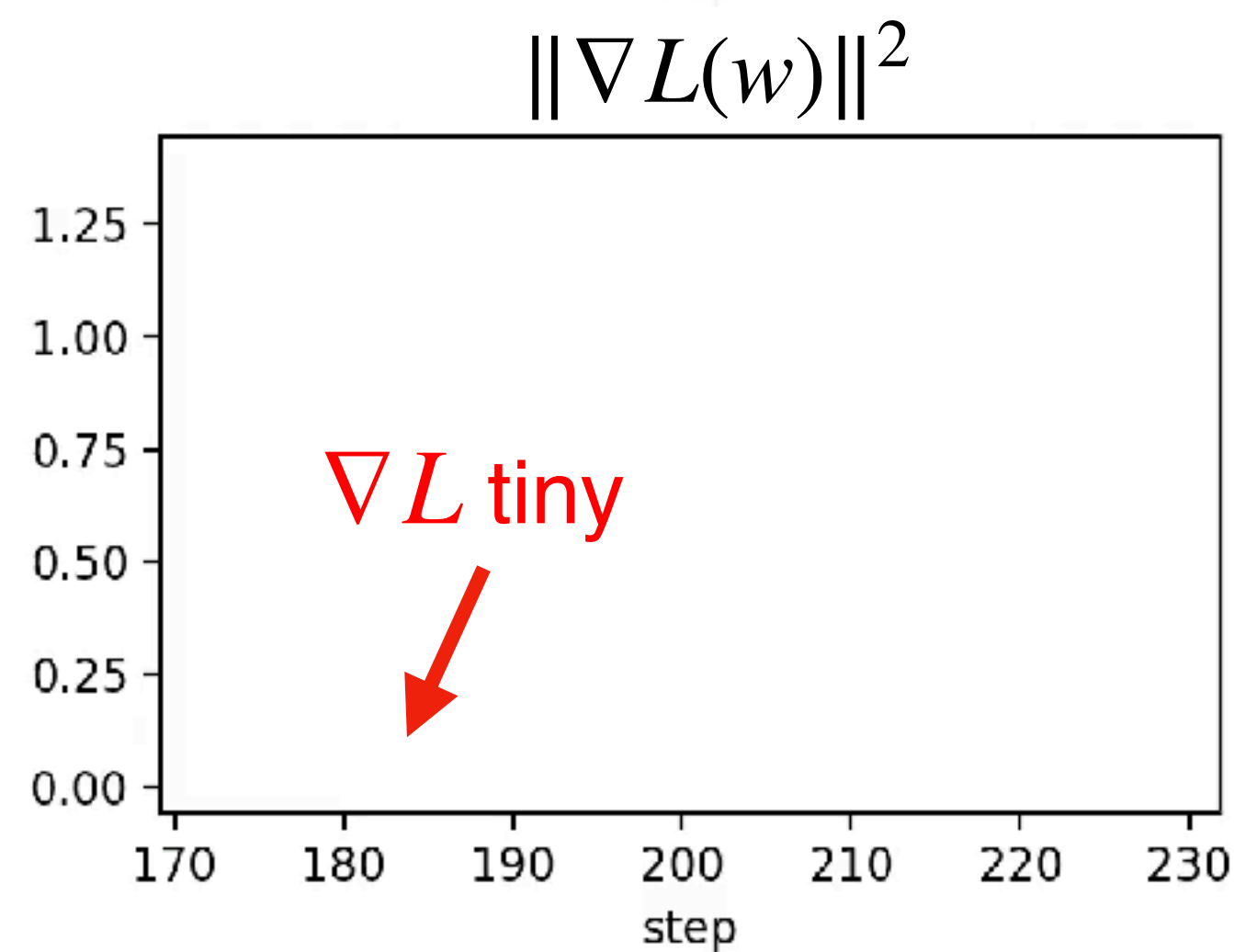
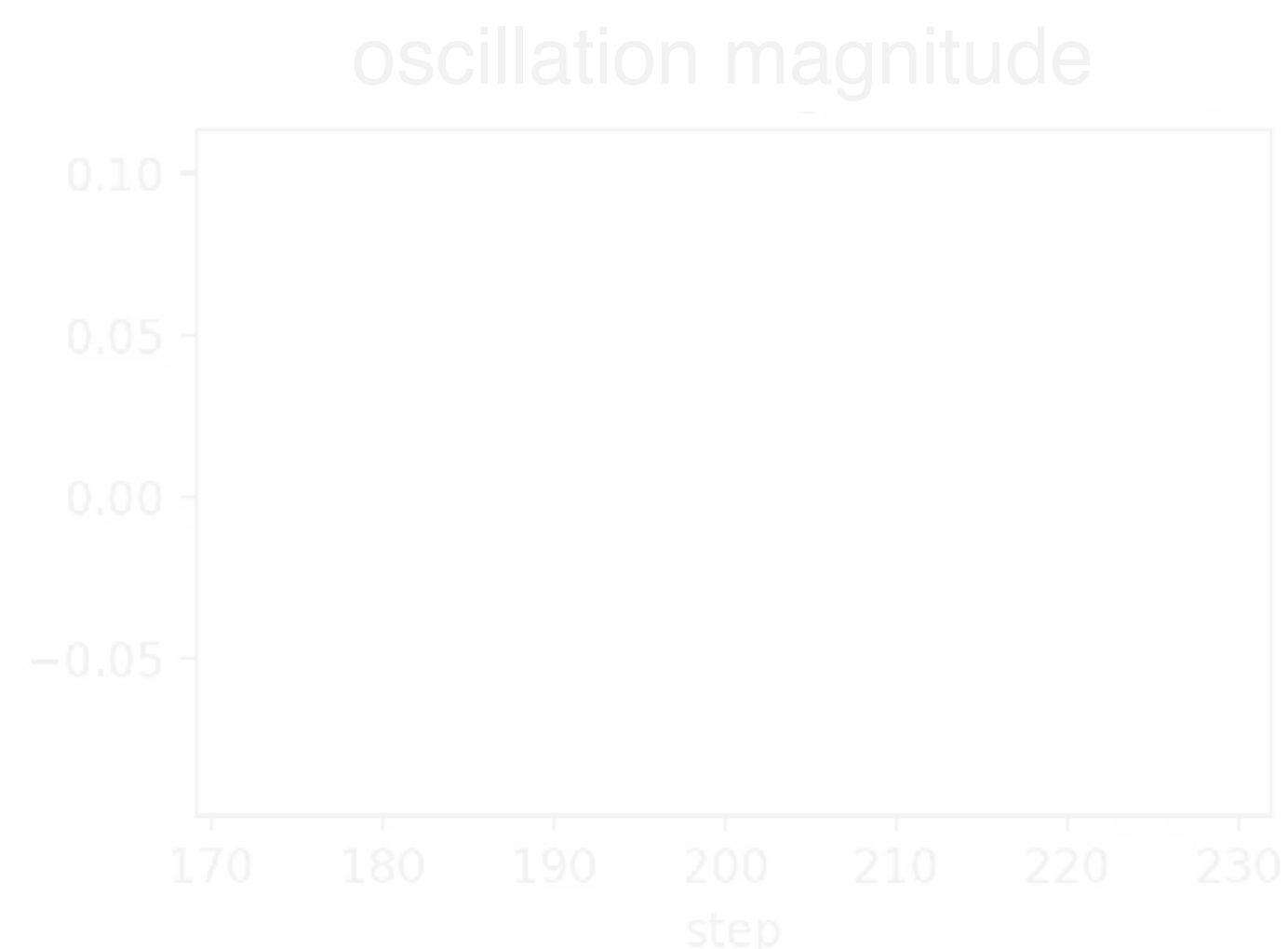
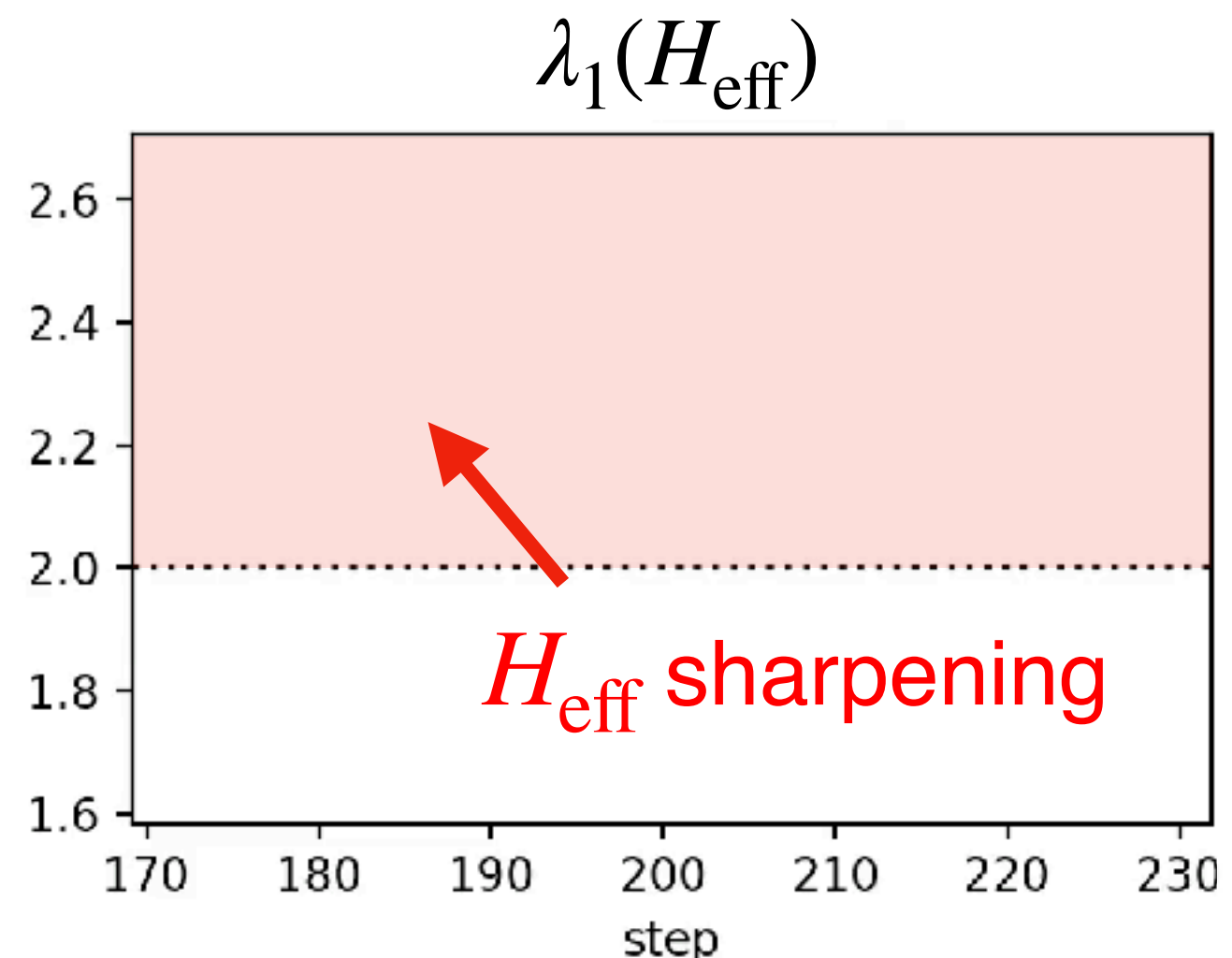
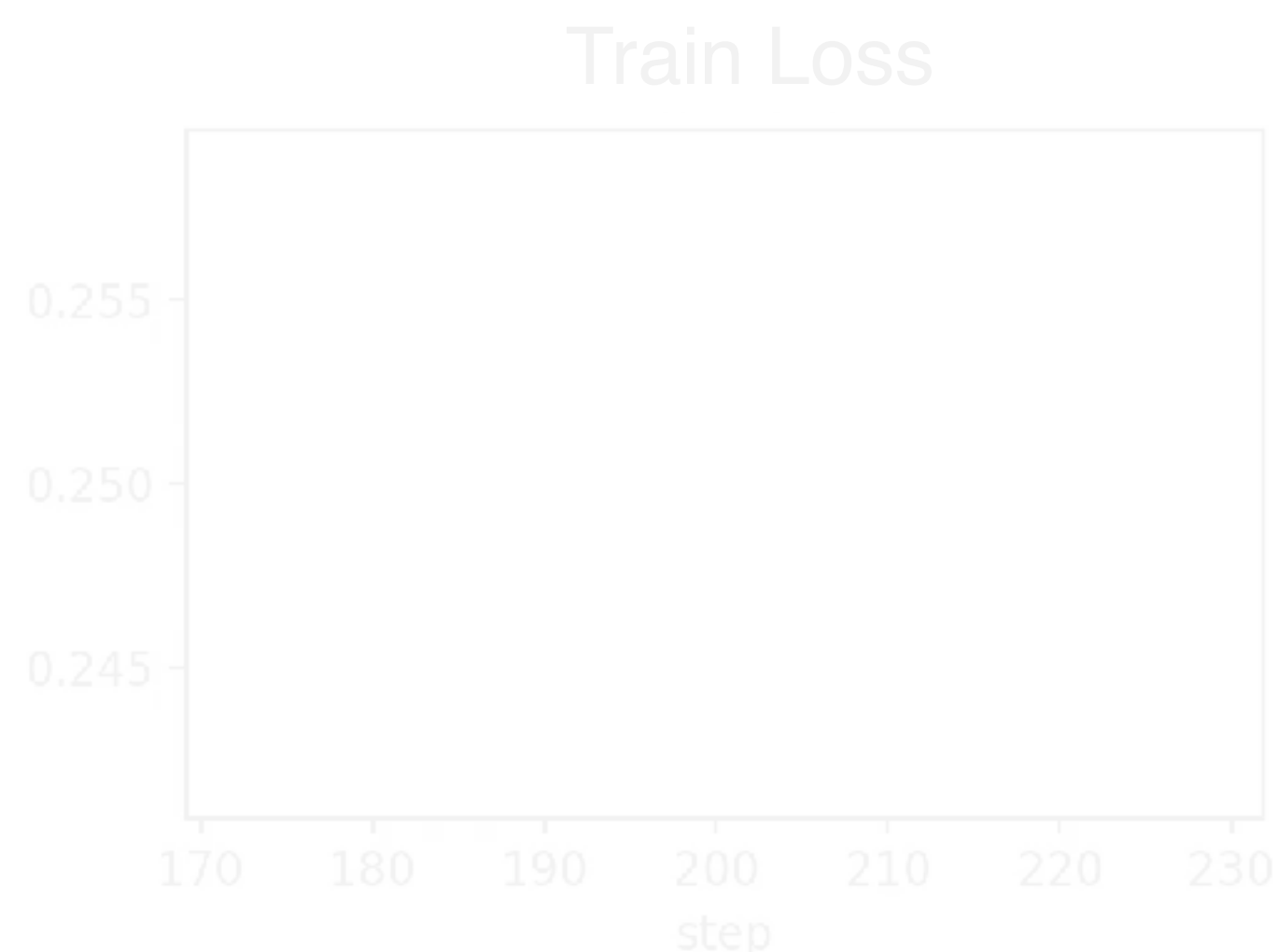
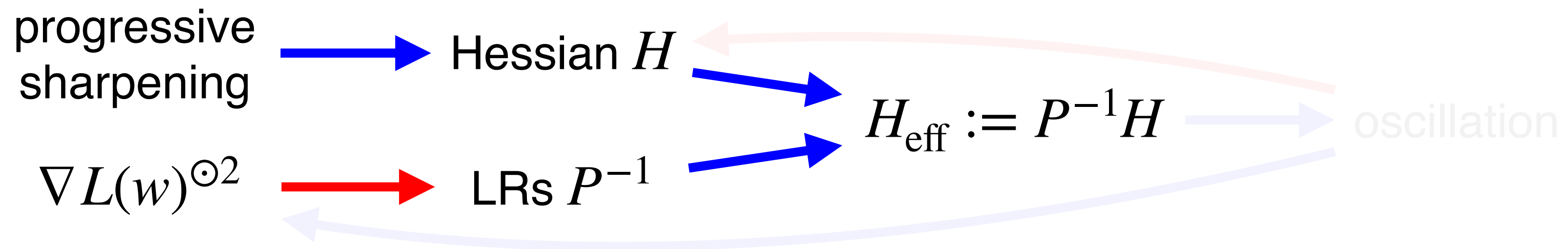


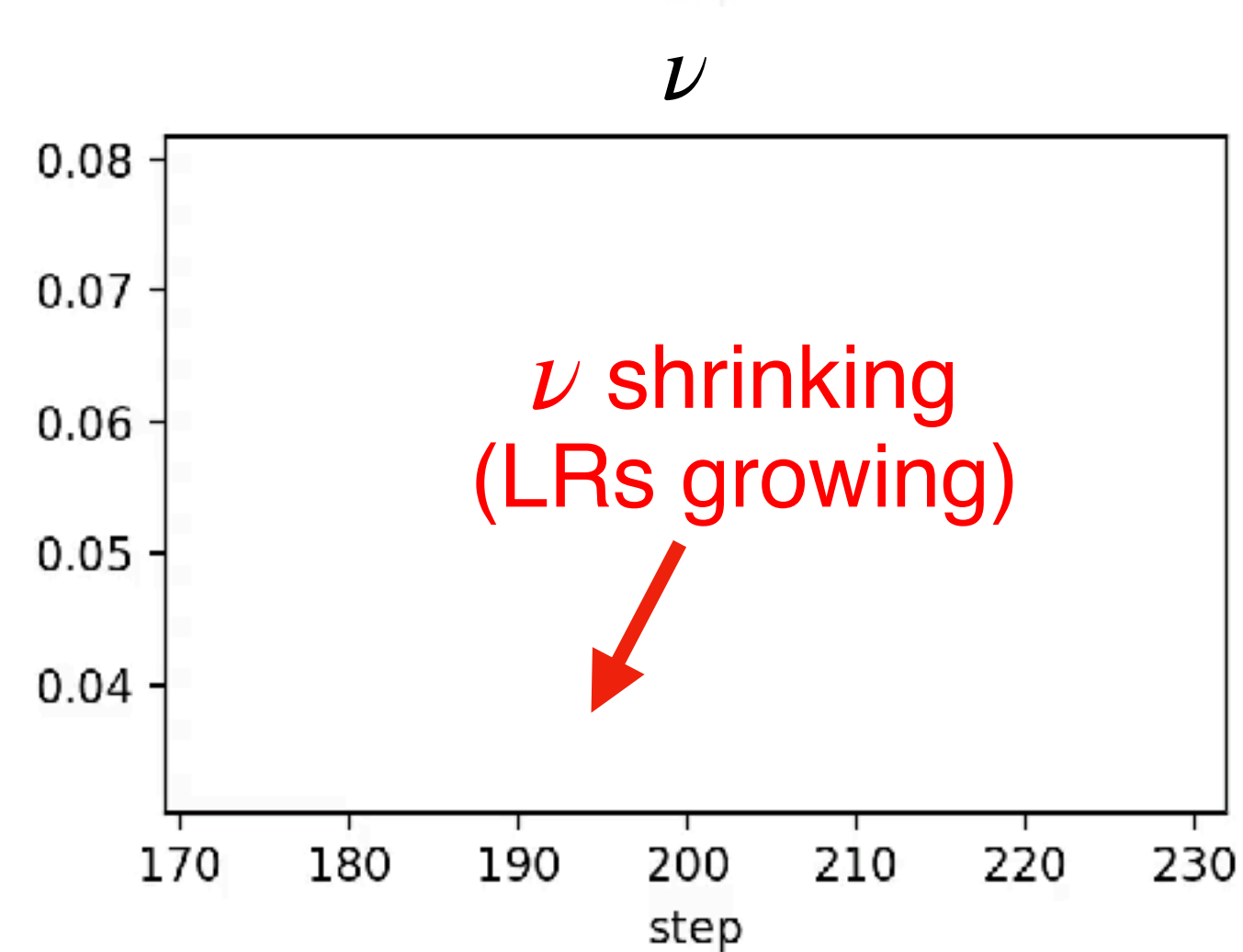
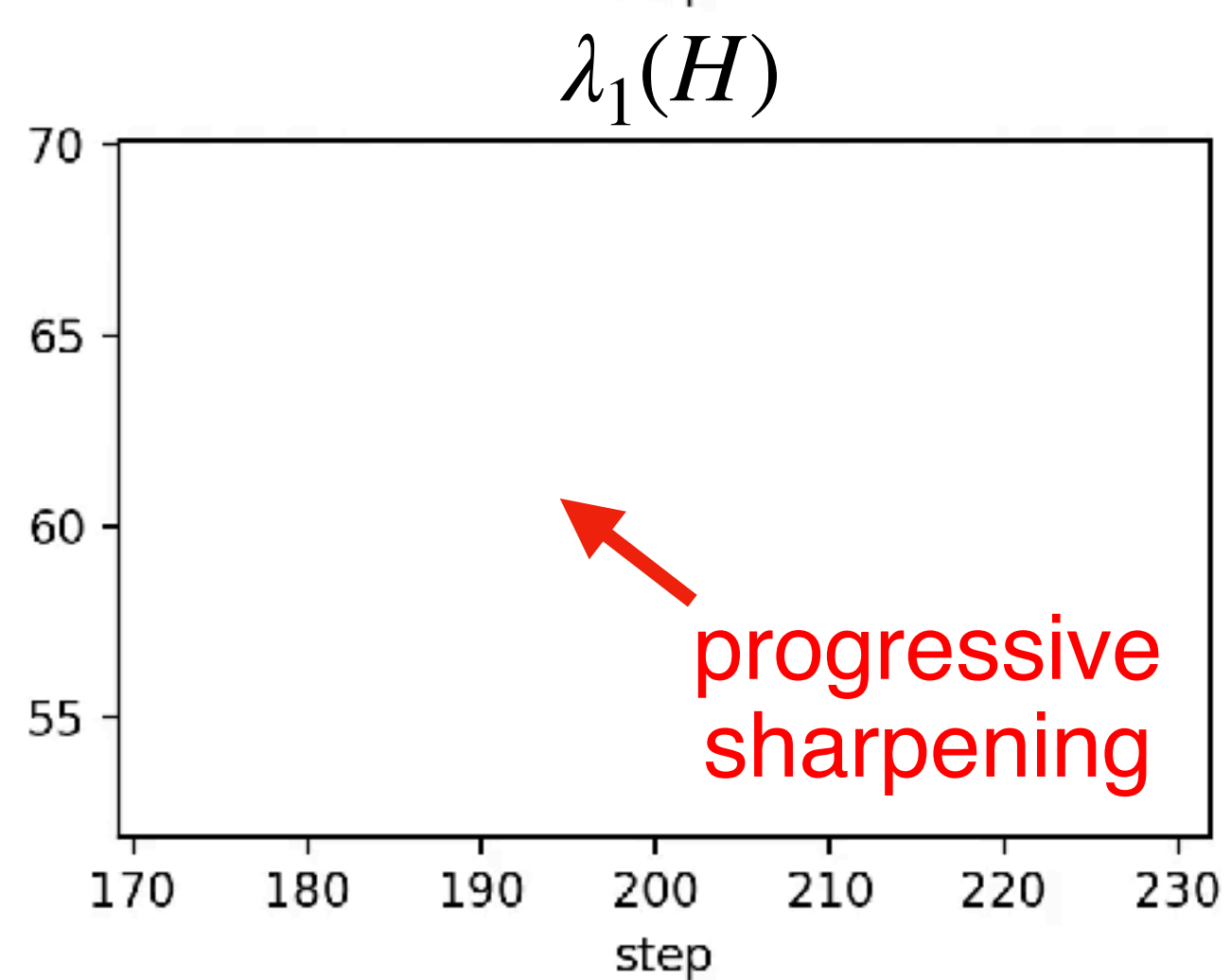
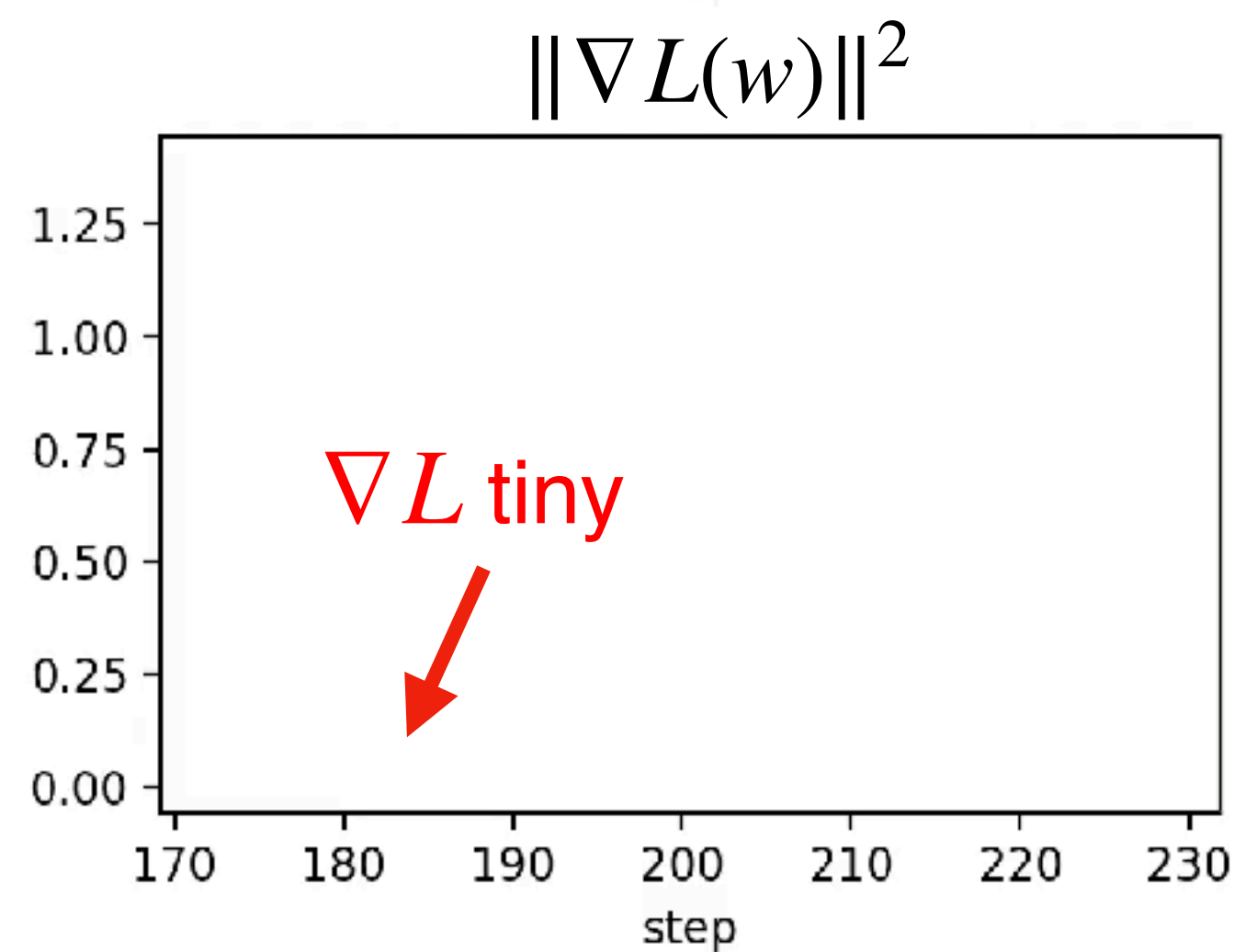
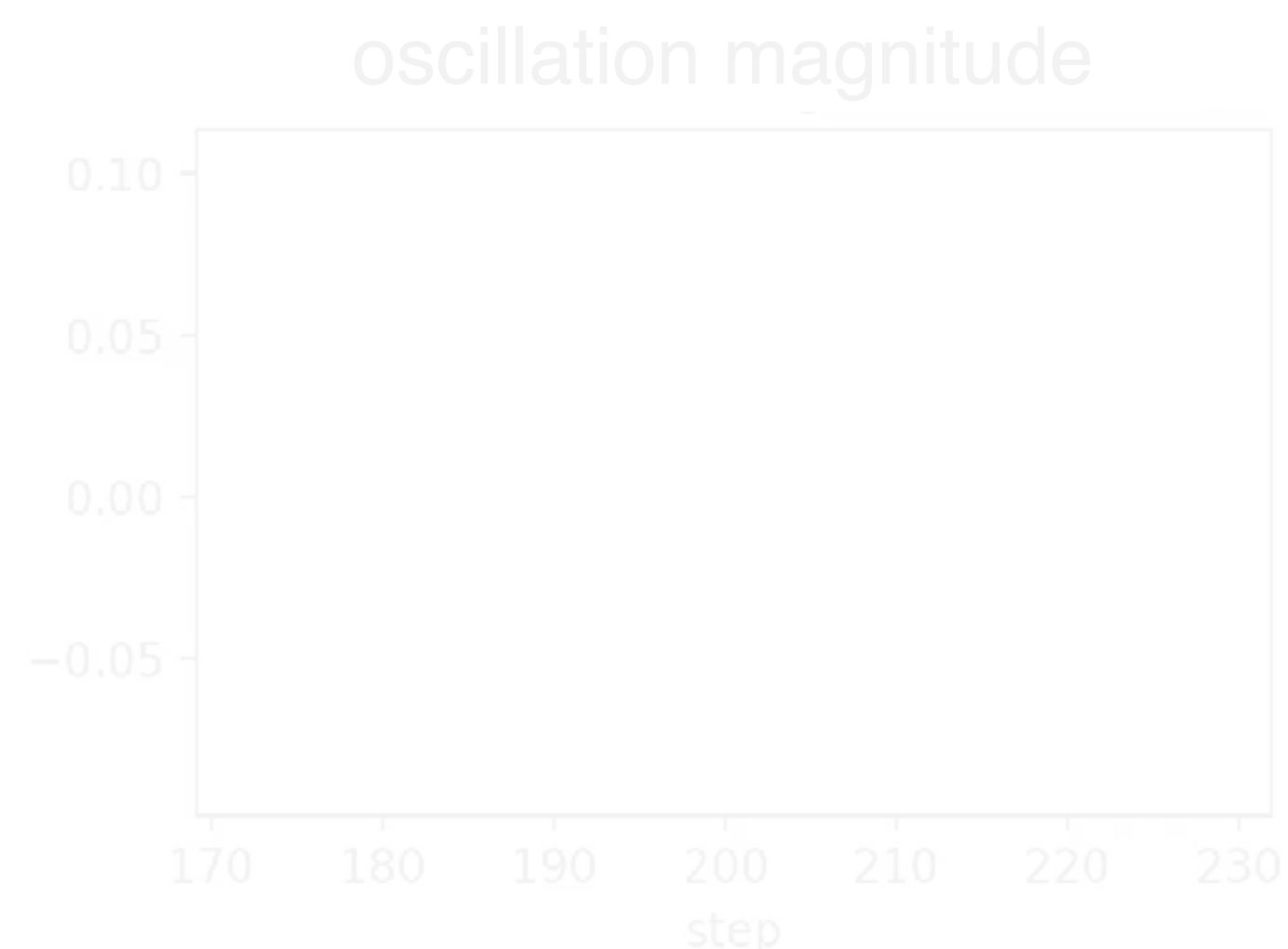
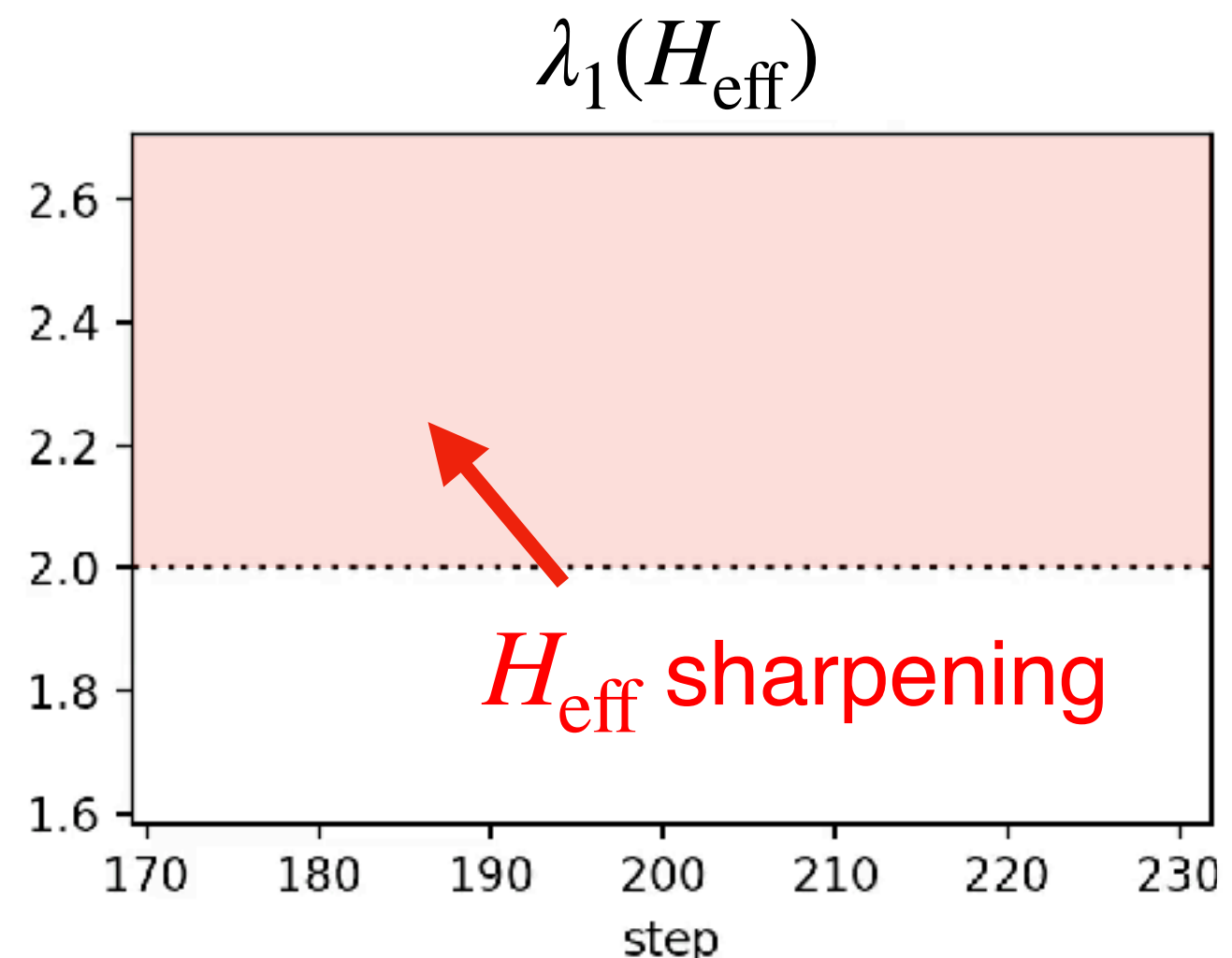
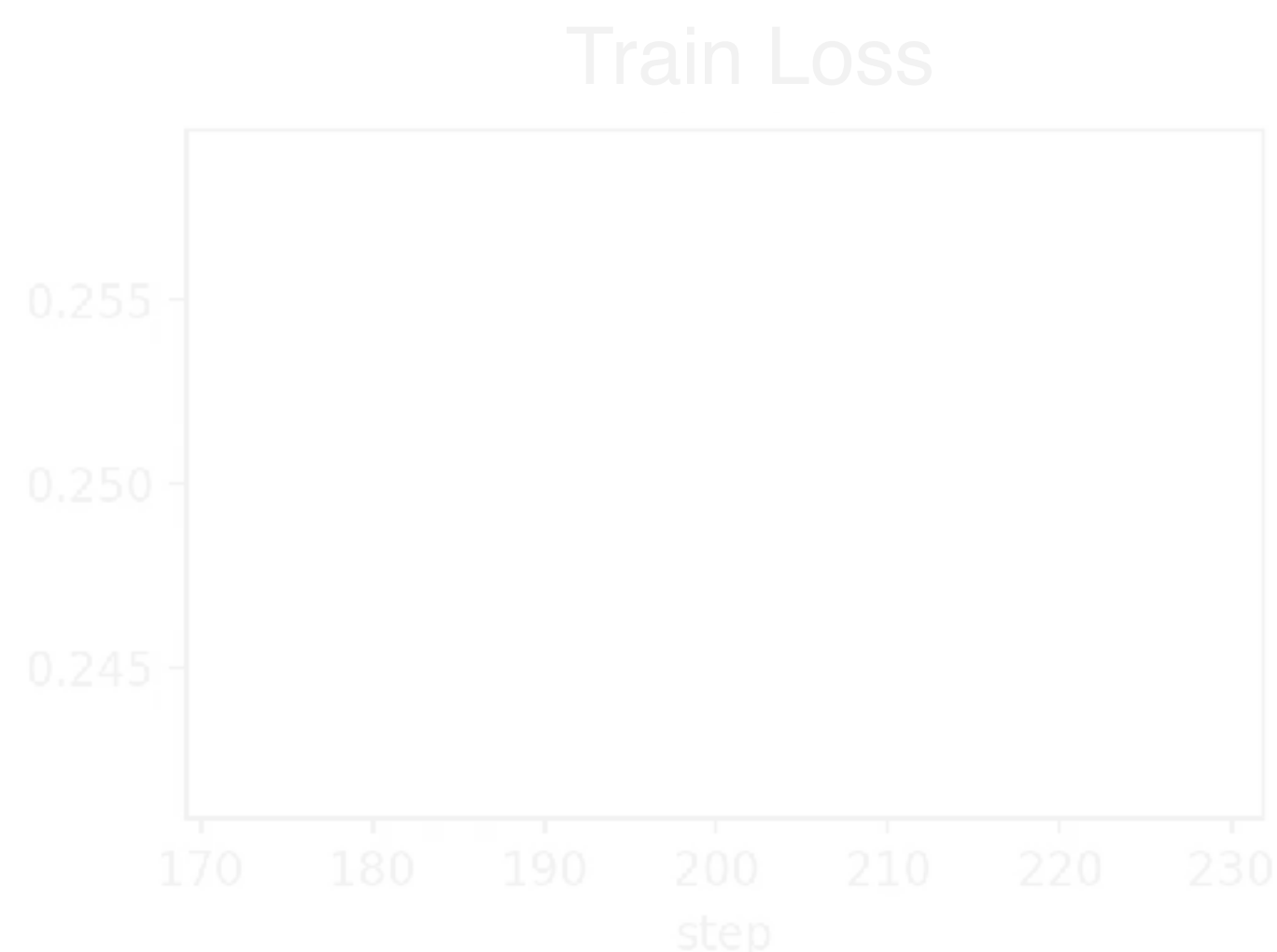
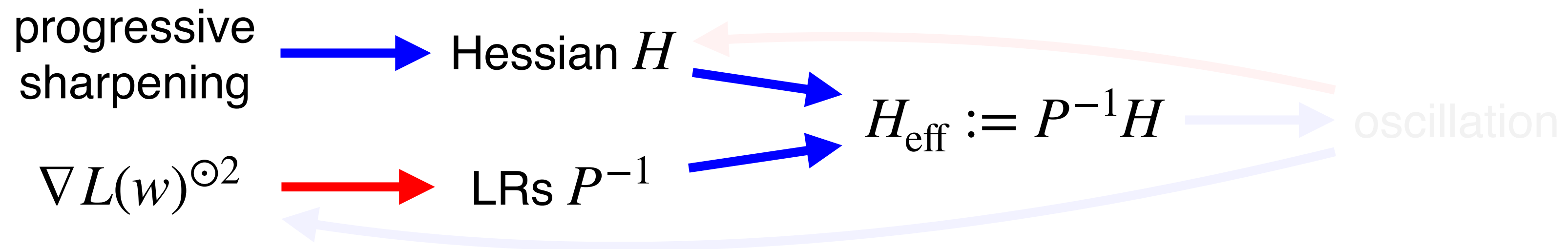
ν





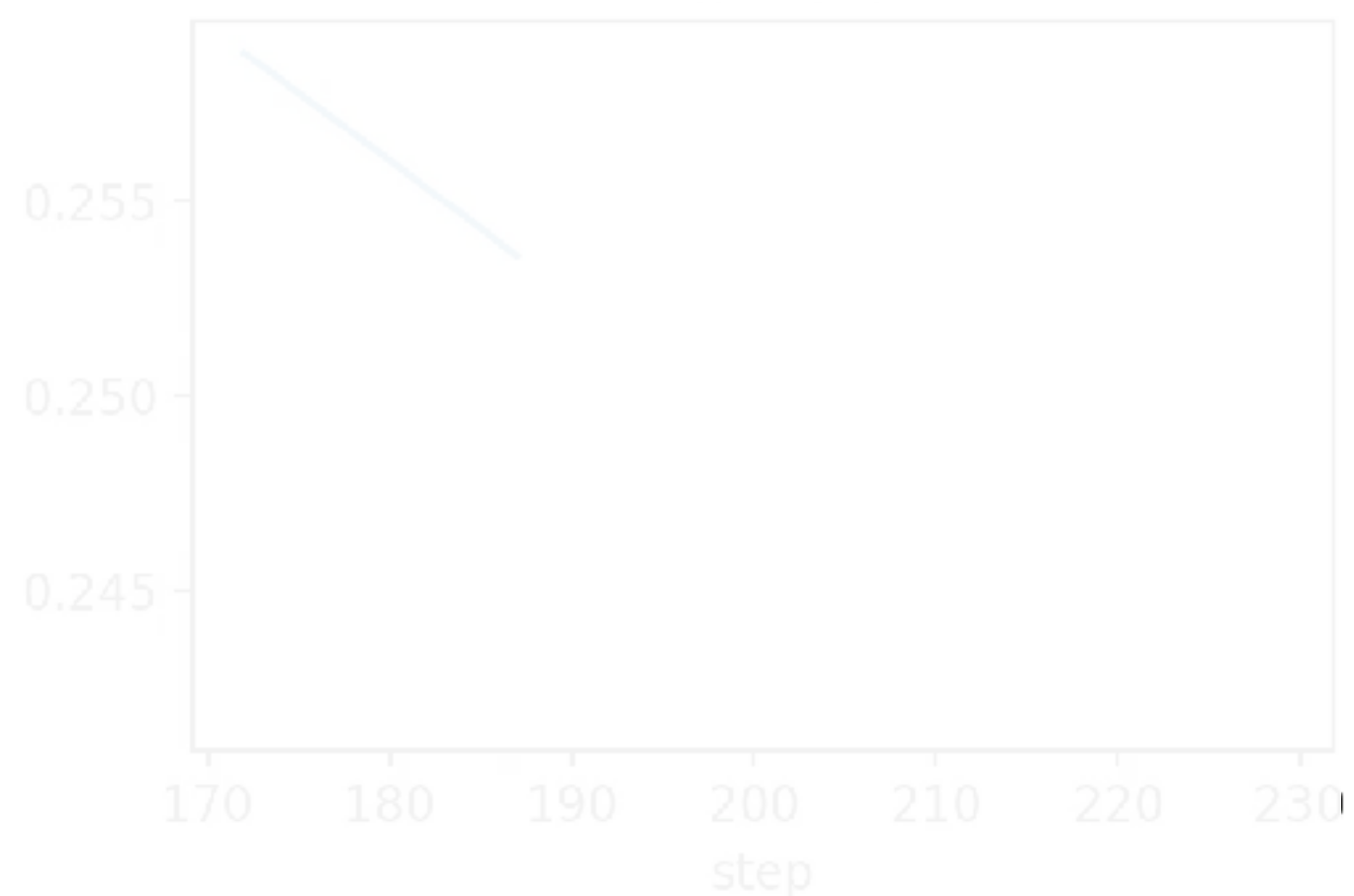




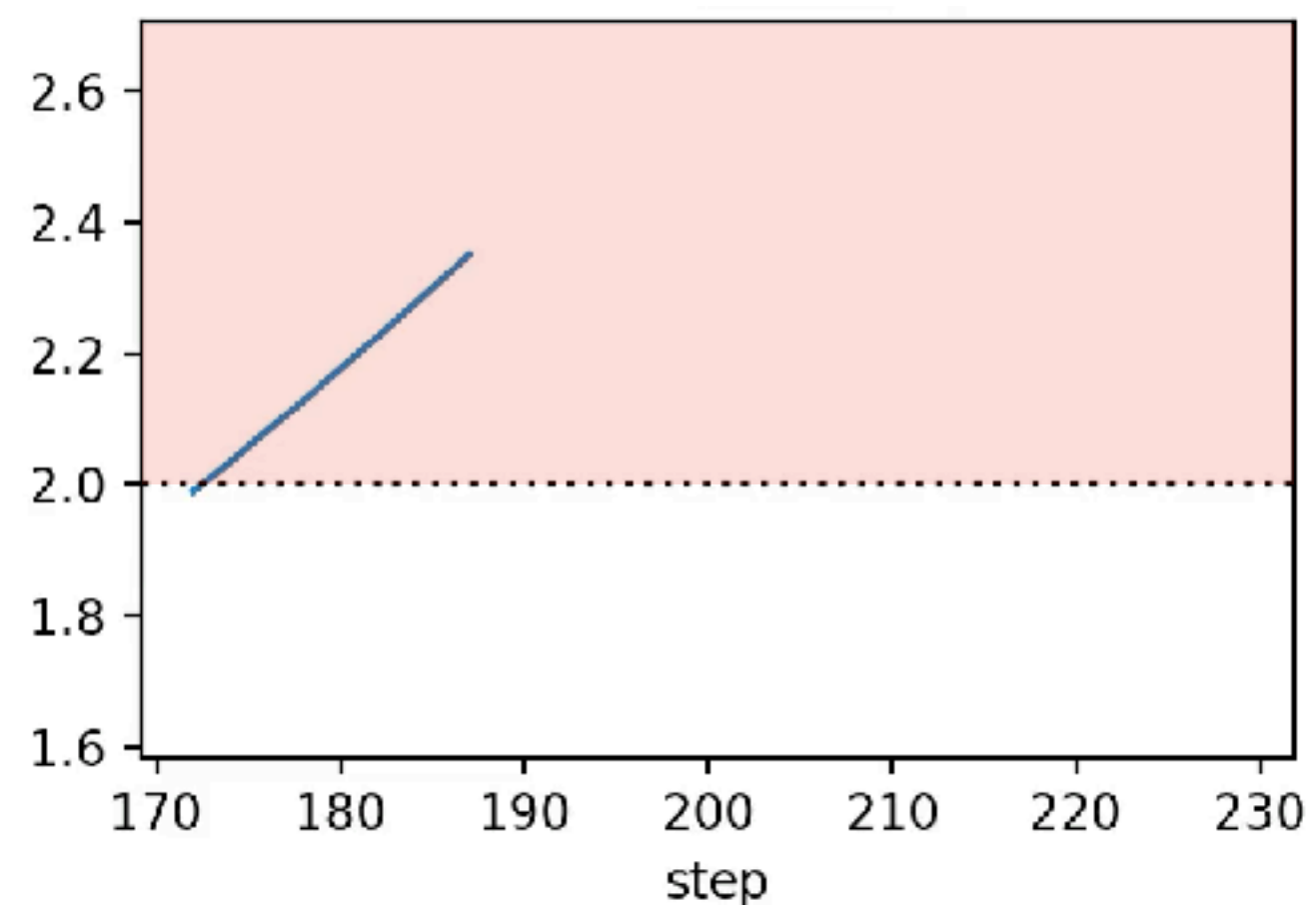




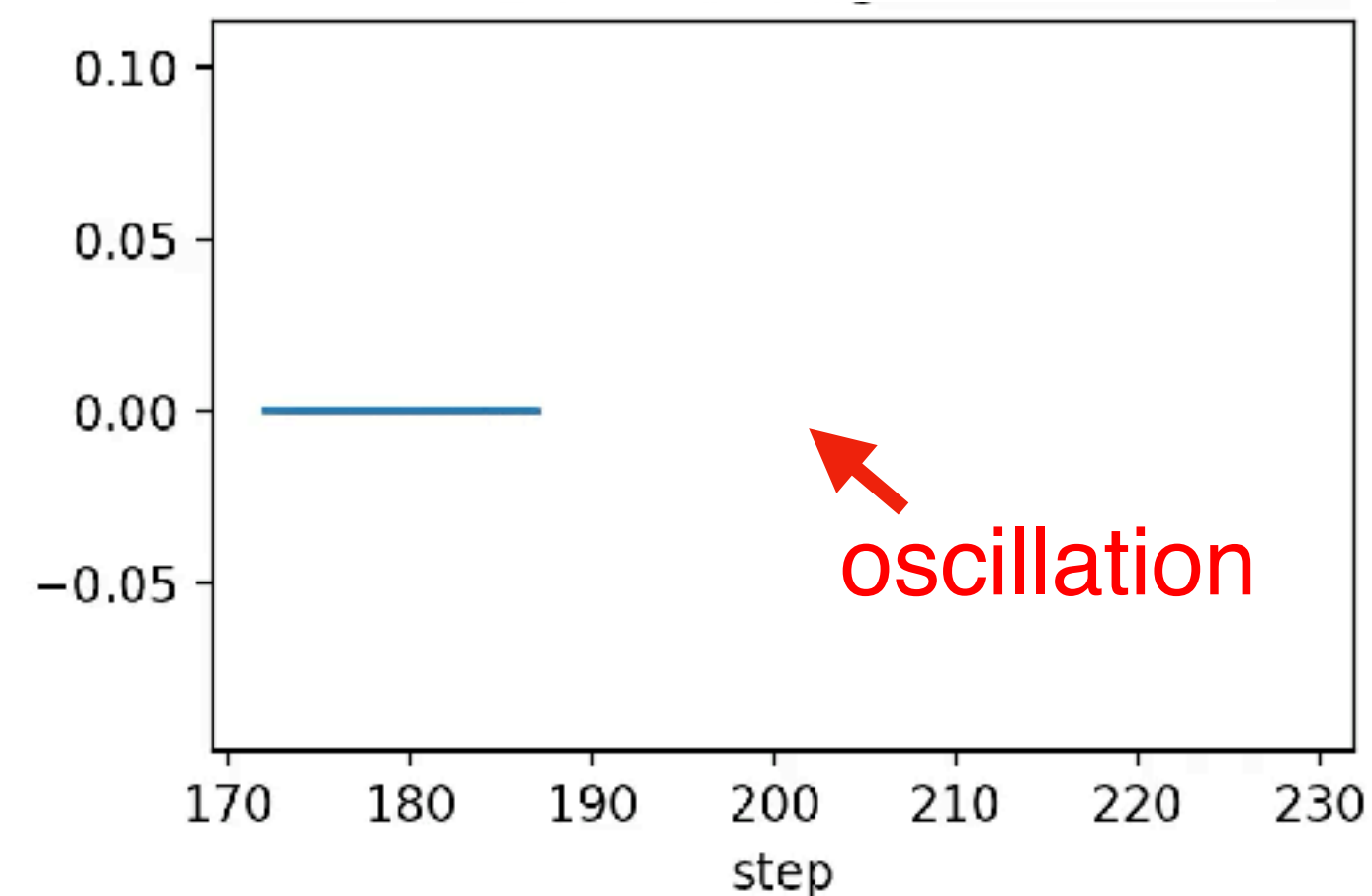
Train Loss



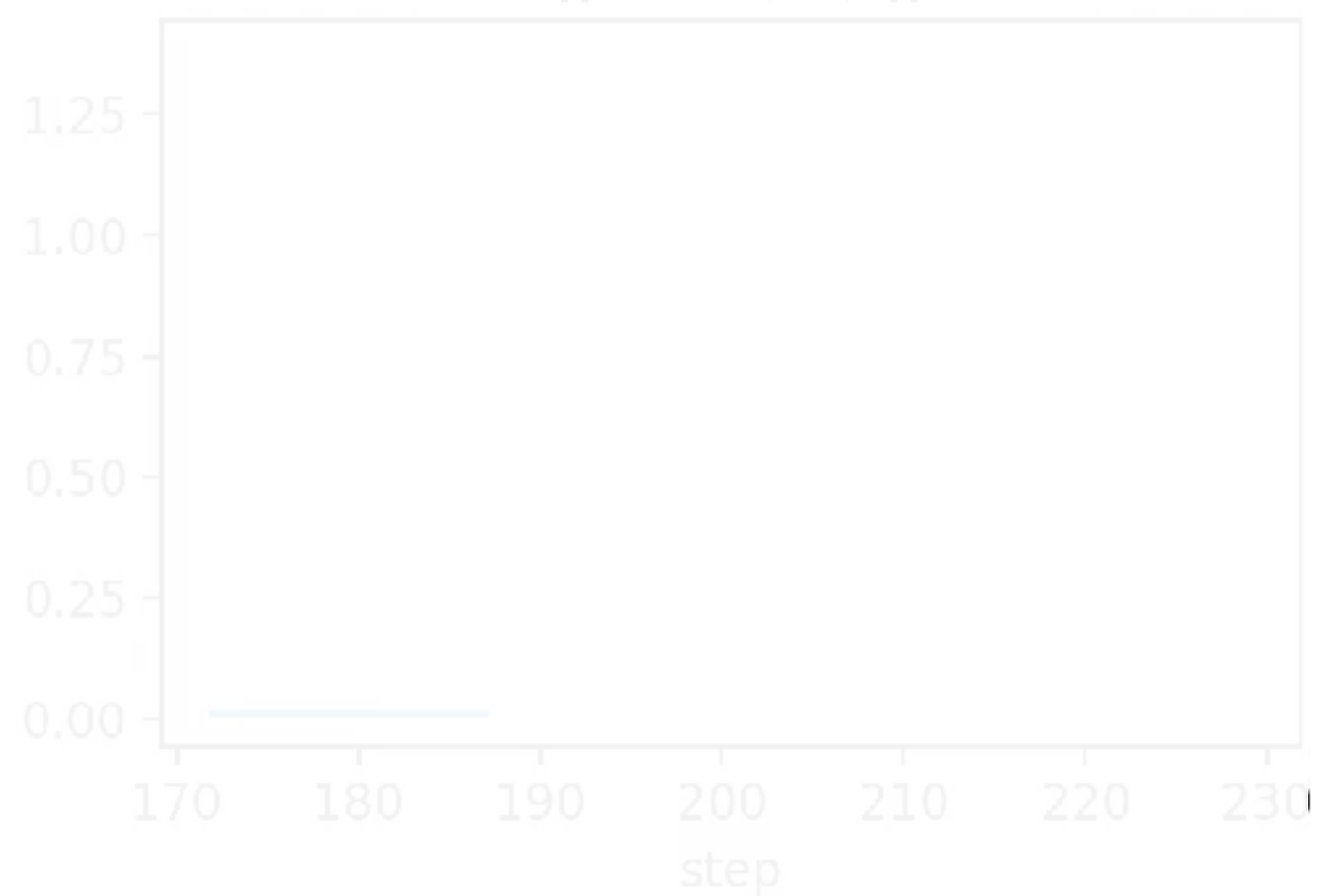
$\lambda_1(H_{\text{eff}})$



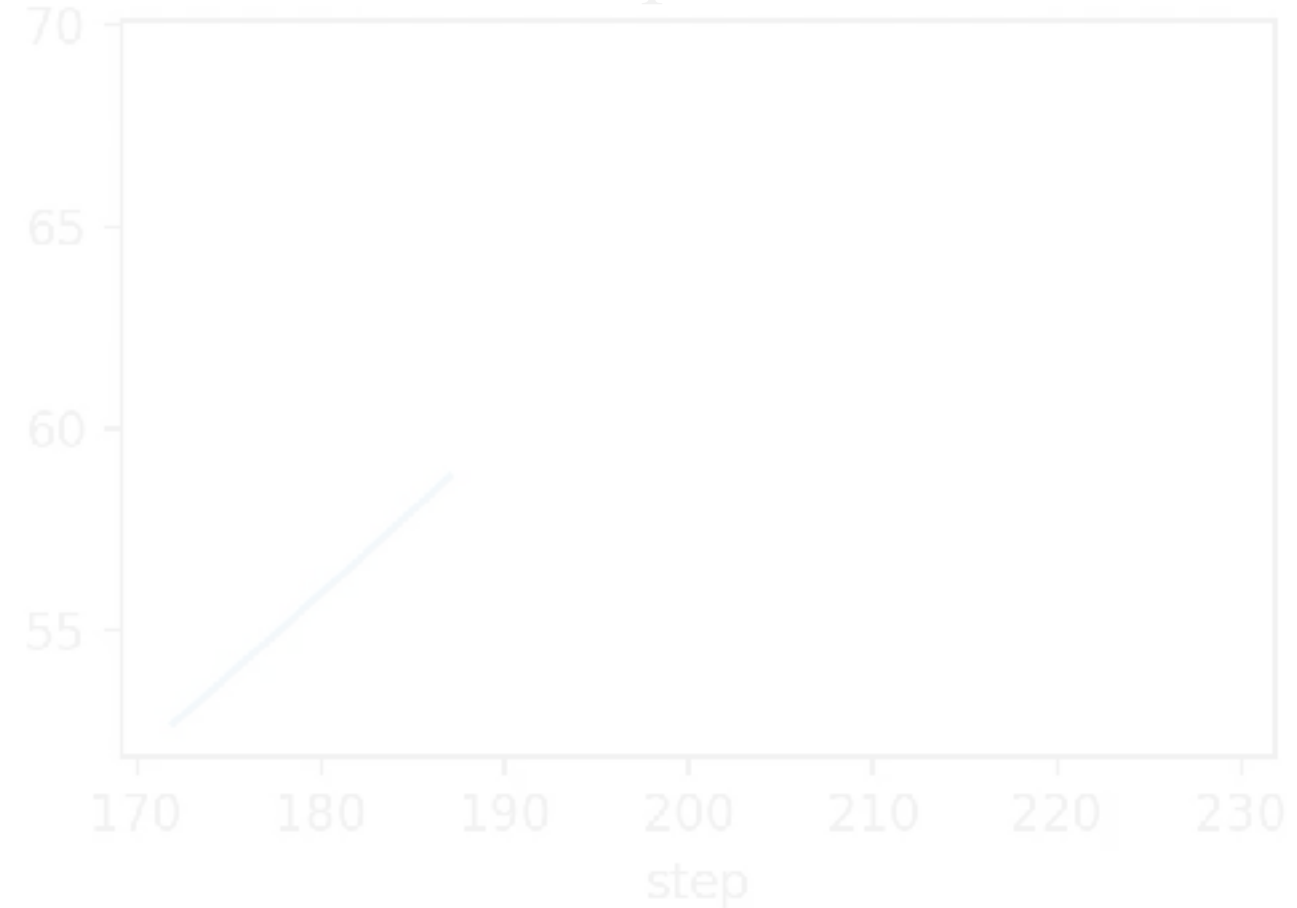
oscillation magnitude



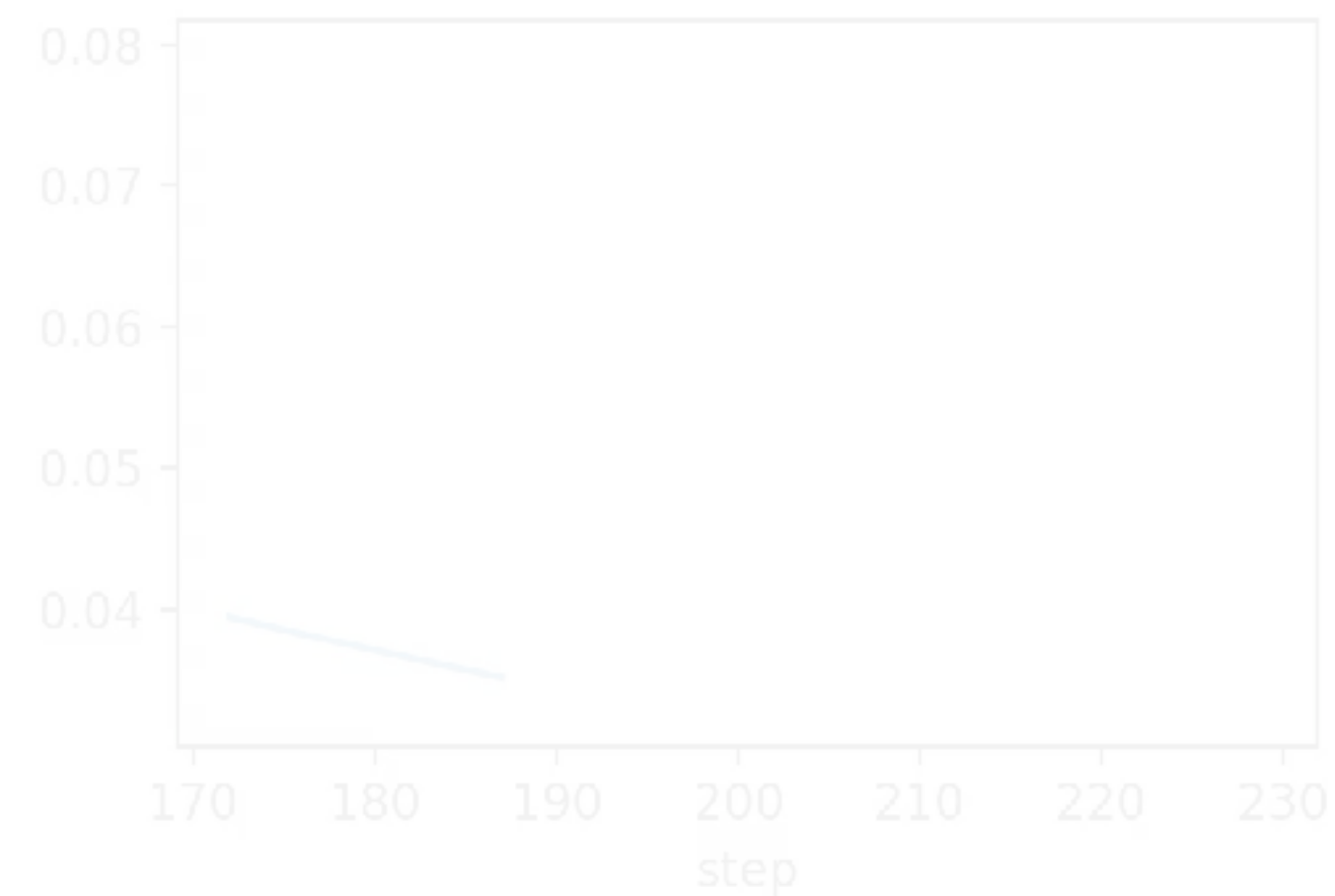
$\|\nabla L(w)\|^2$



$\lambda_1(H)$

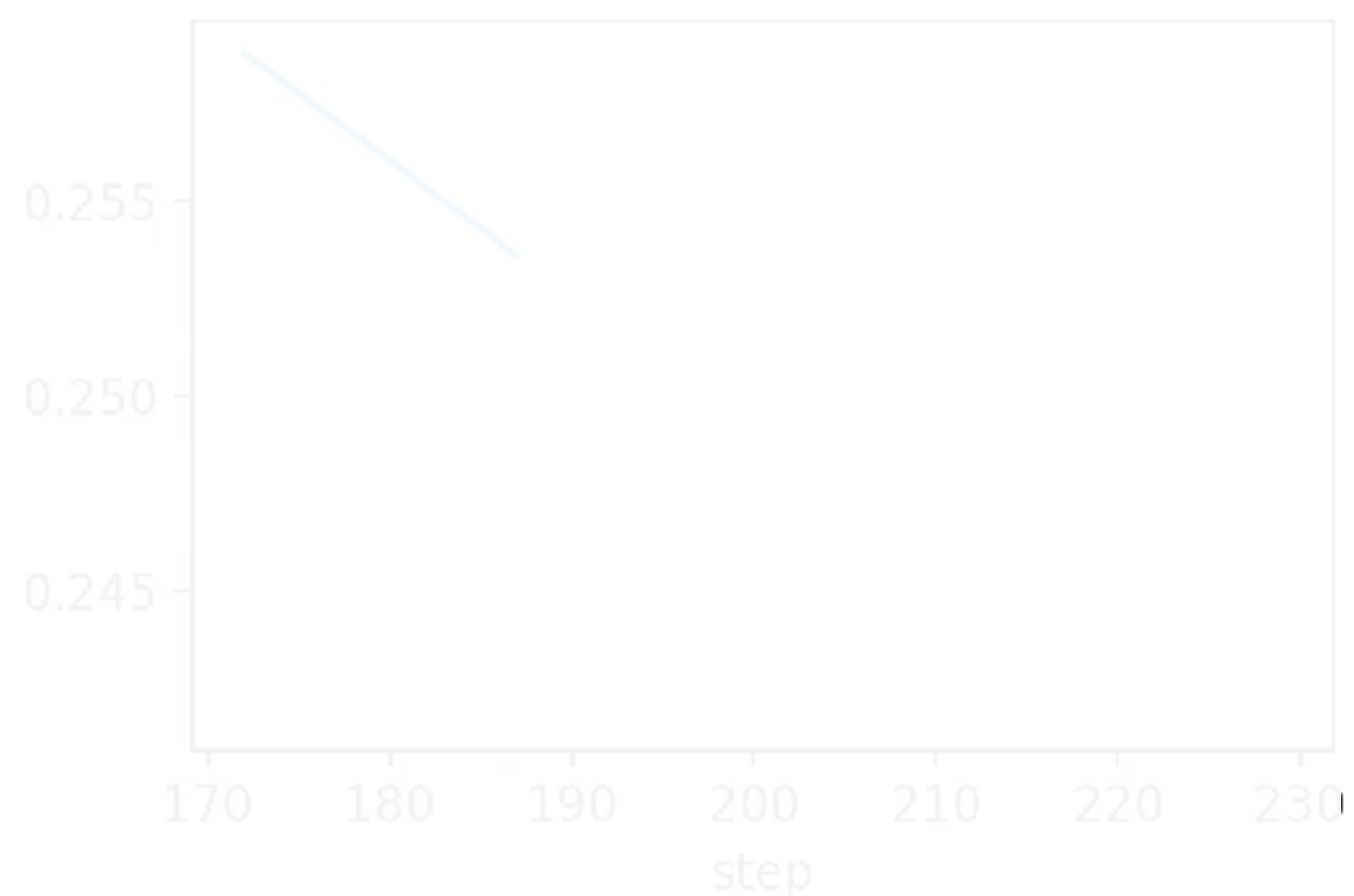


ν

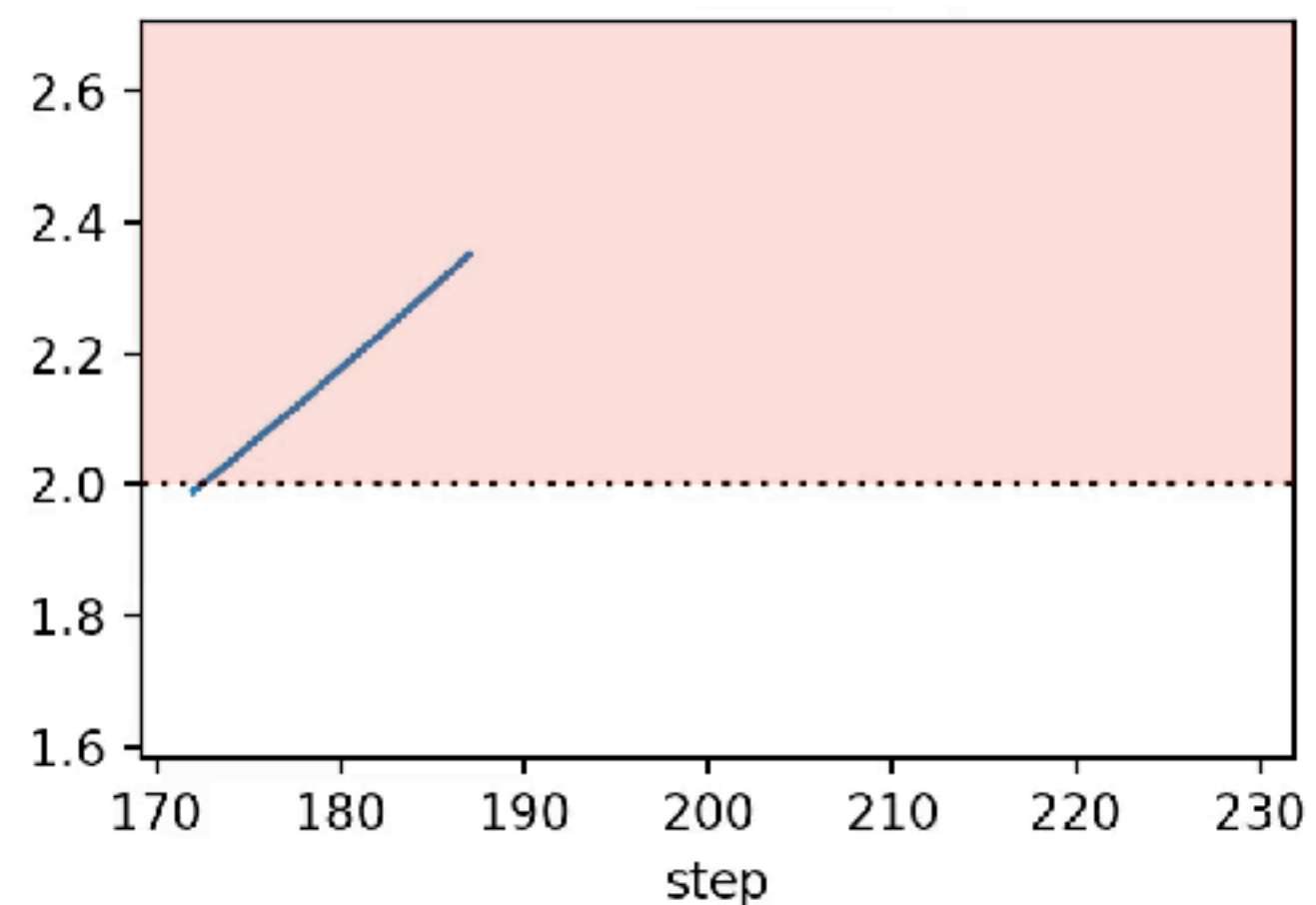




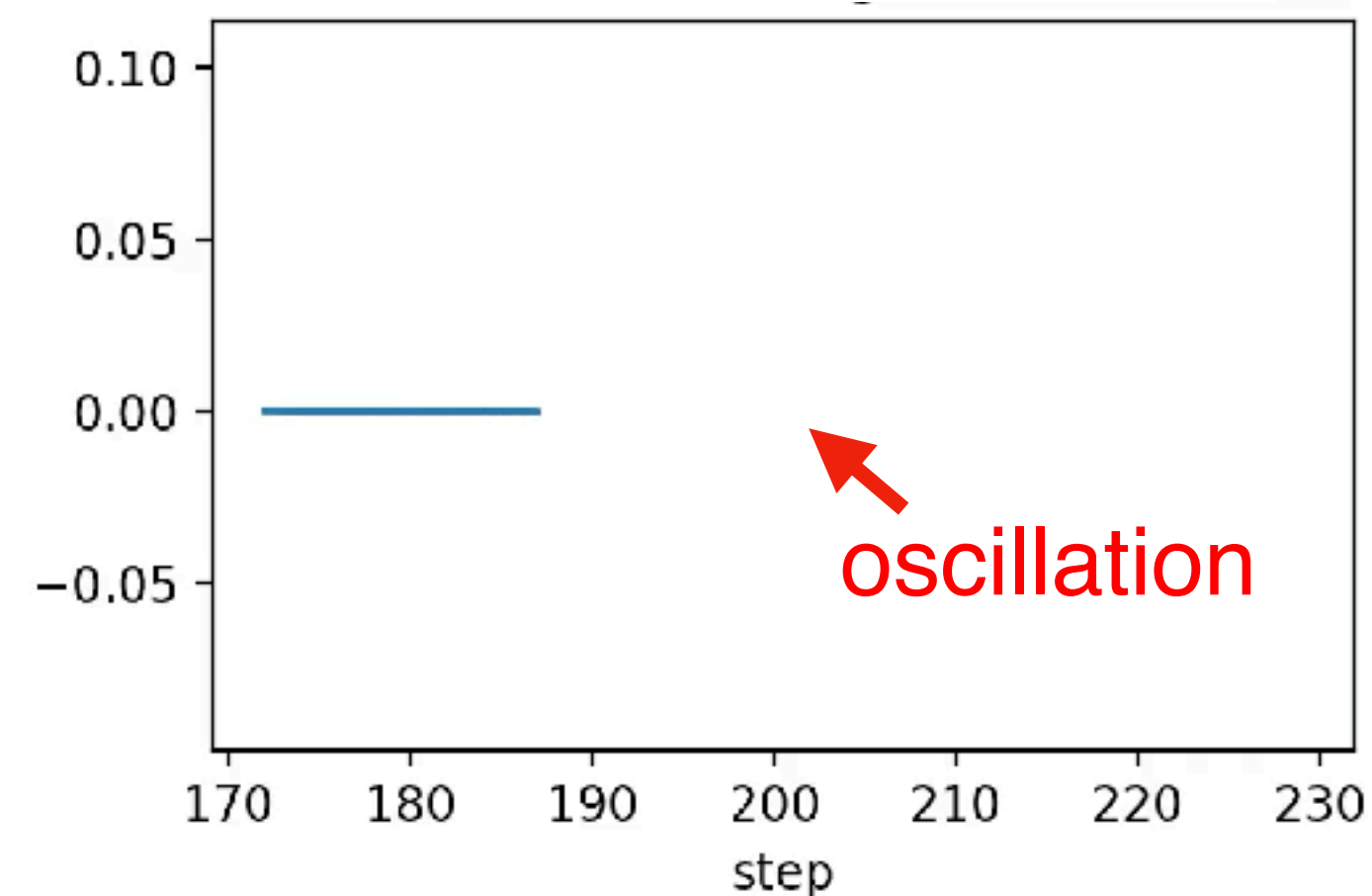
Train Loss



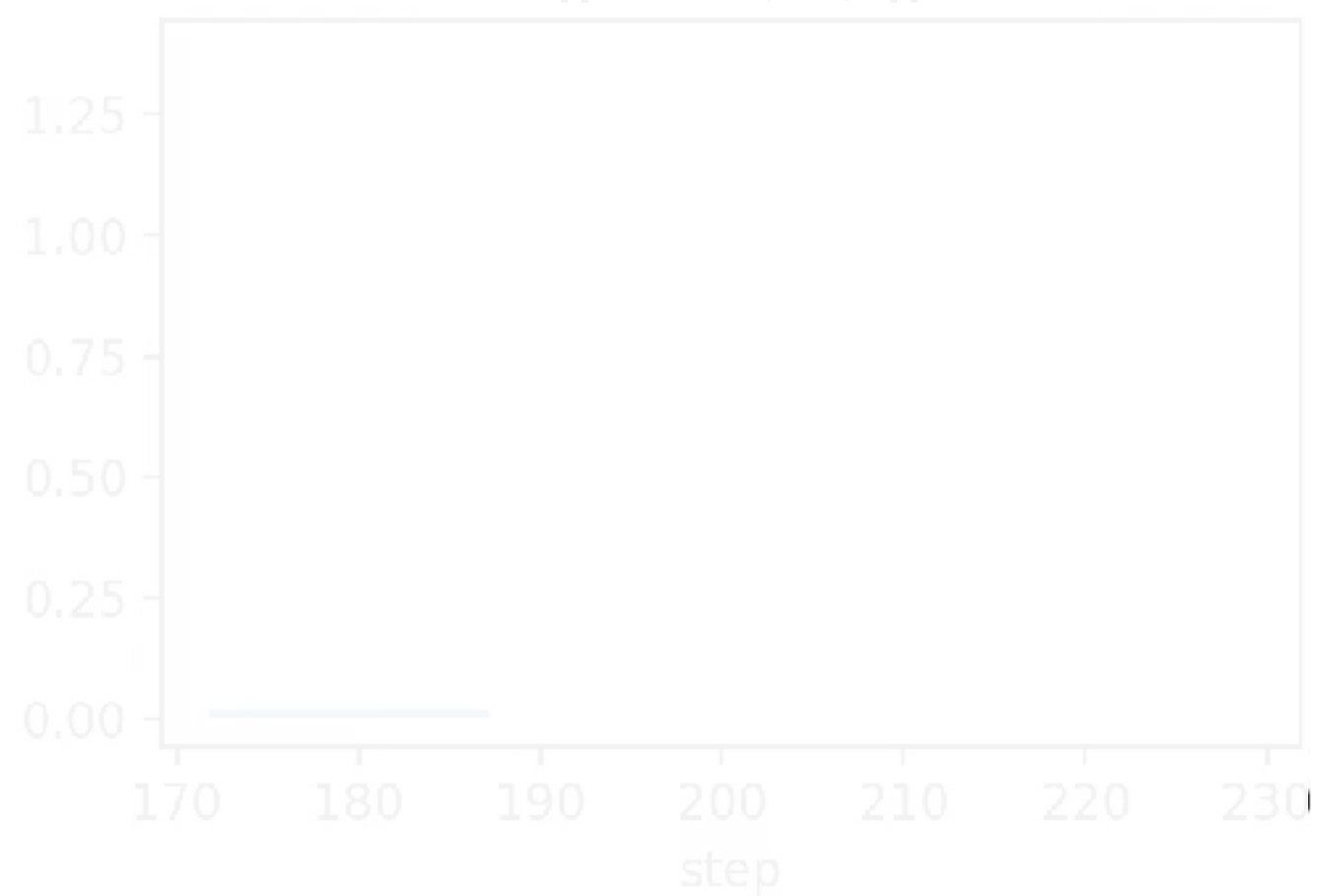
$\lambda_1(H_{\text{eff}})$



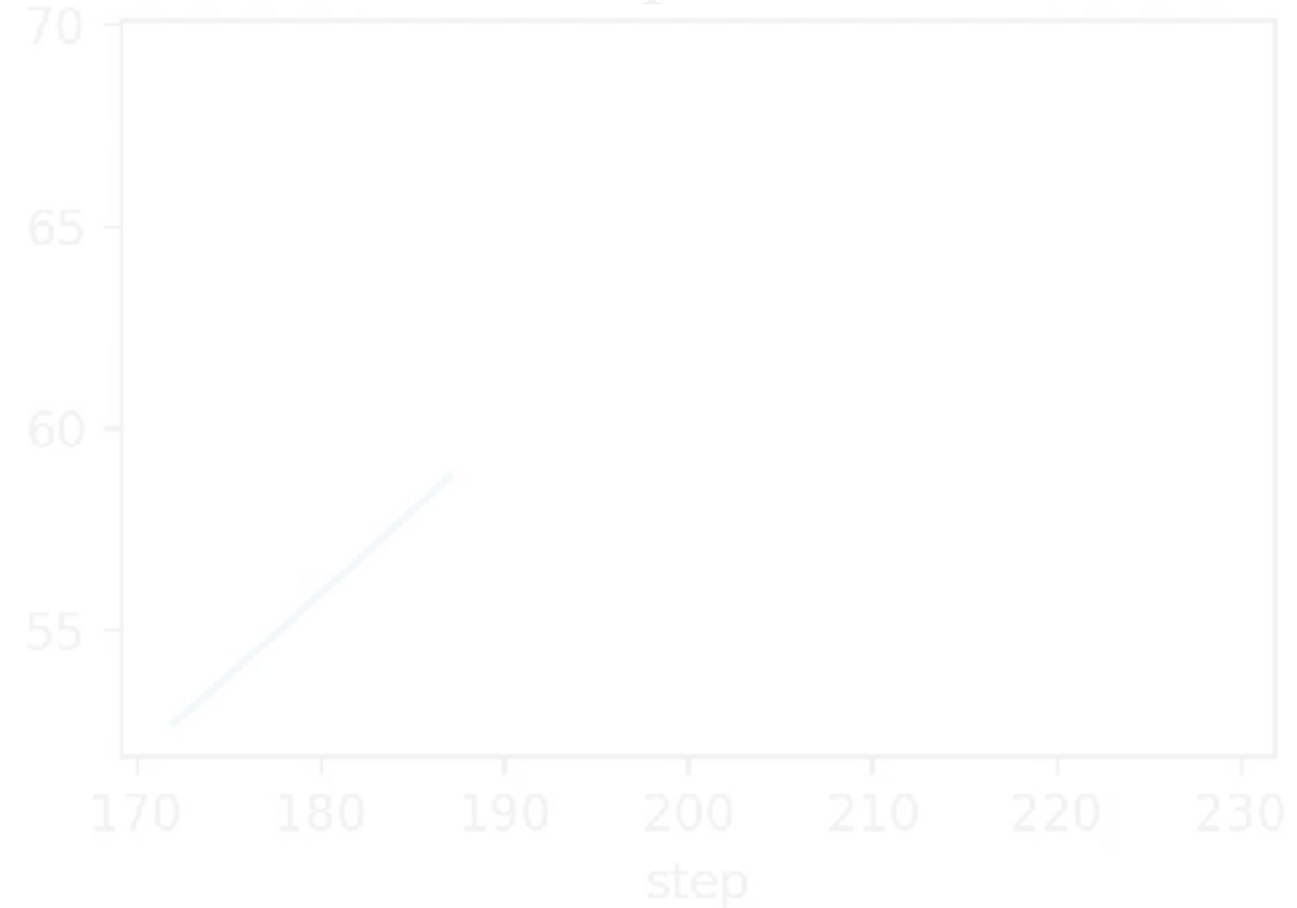
oscillation magnitude



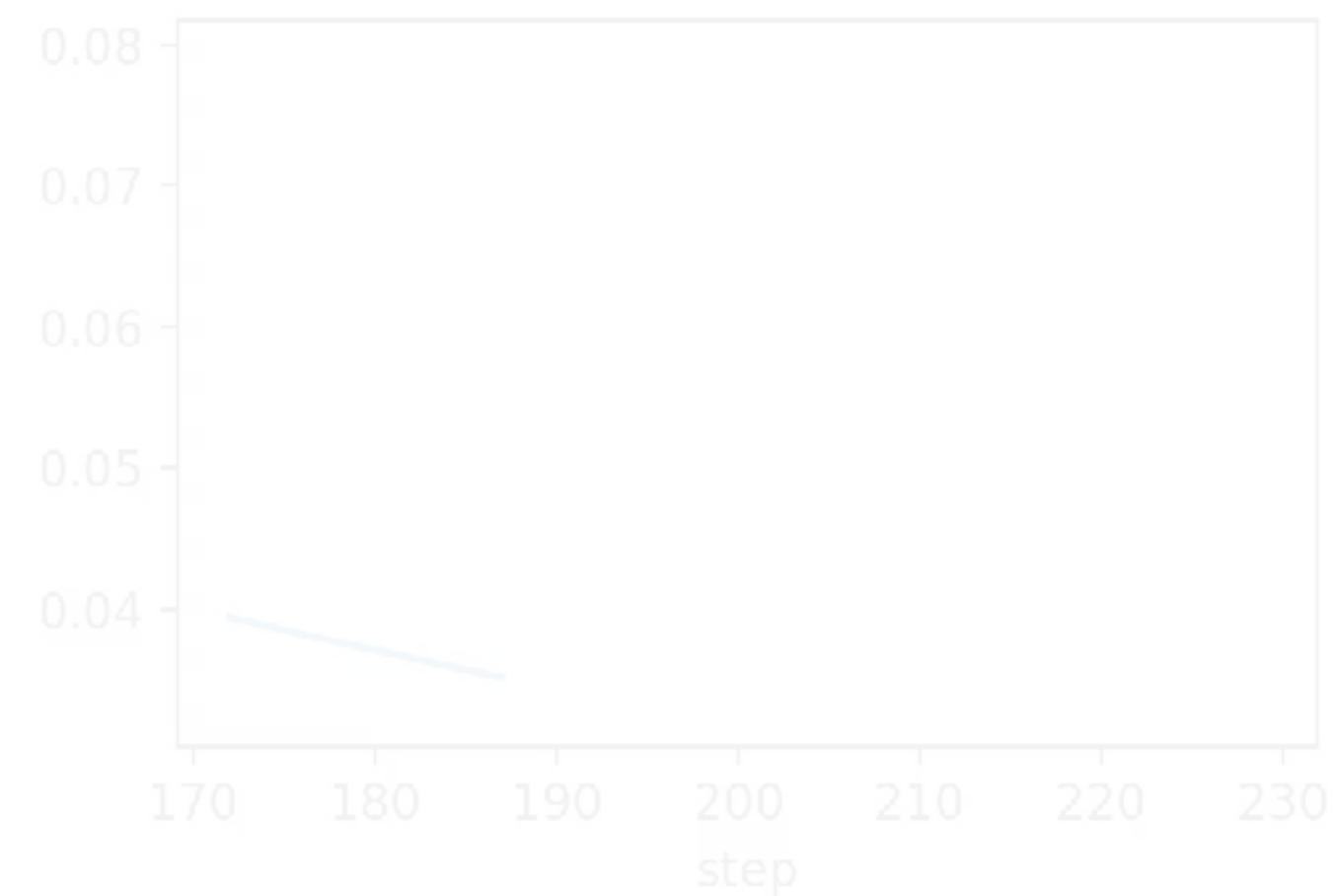
$\|\nabla L(w)\|^2$

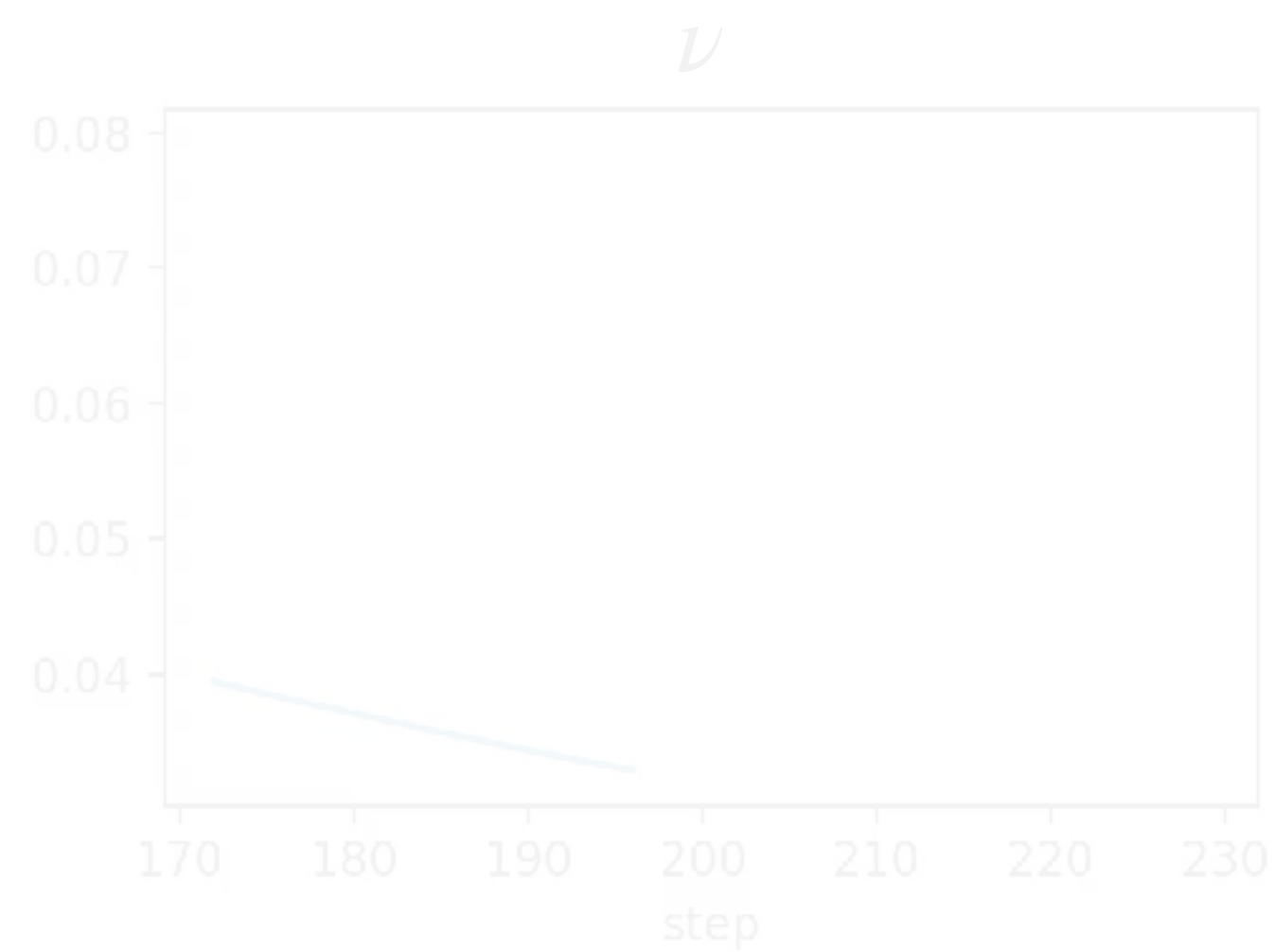
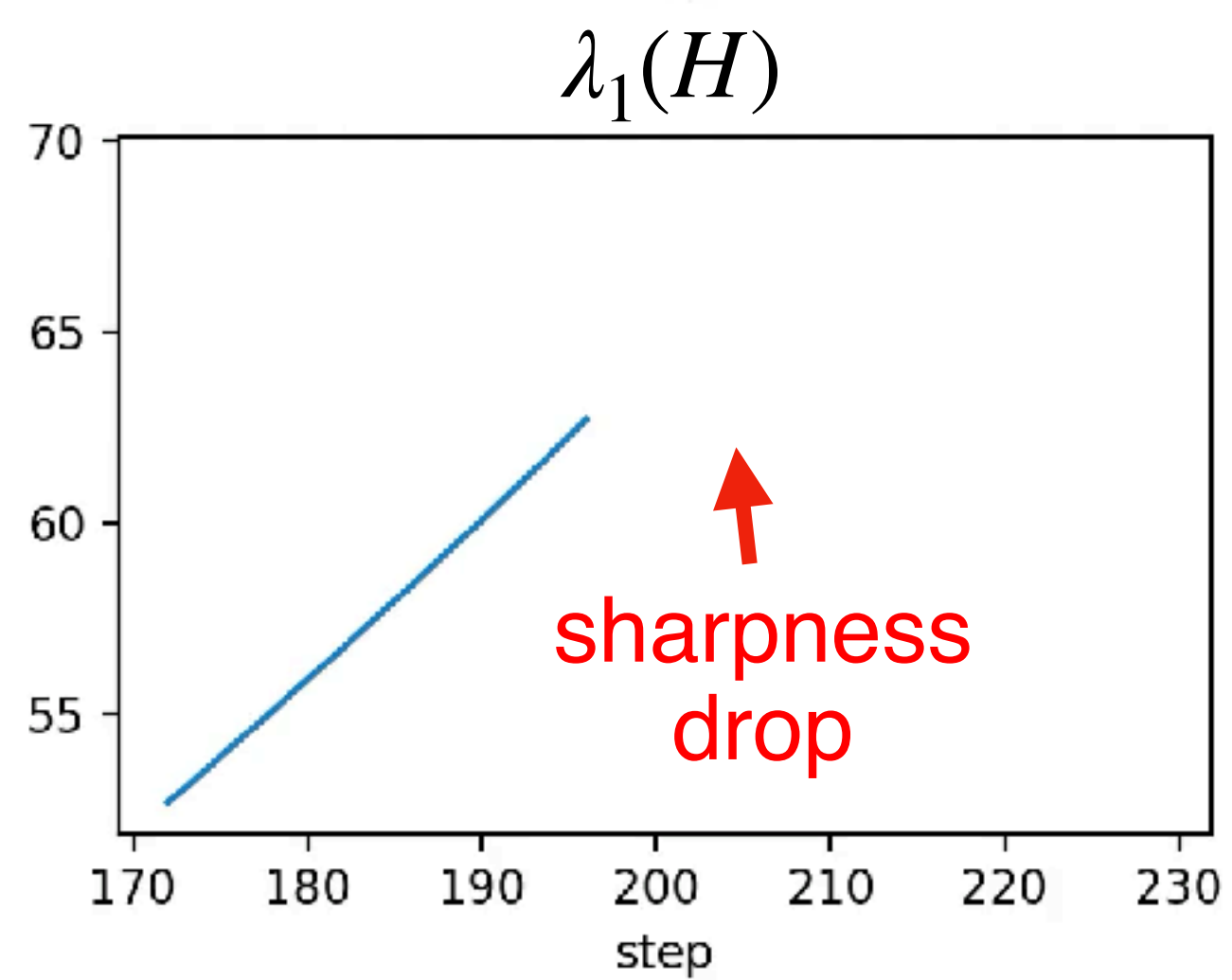
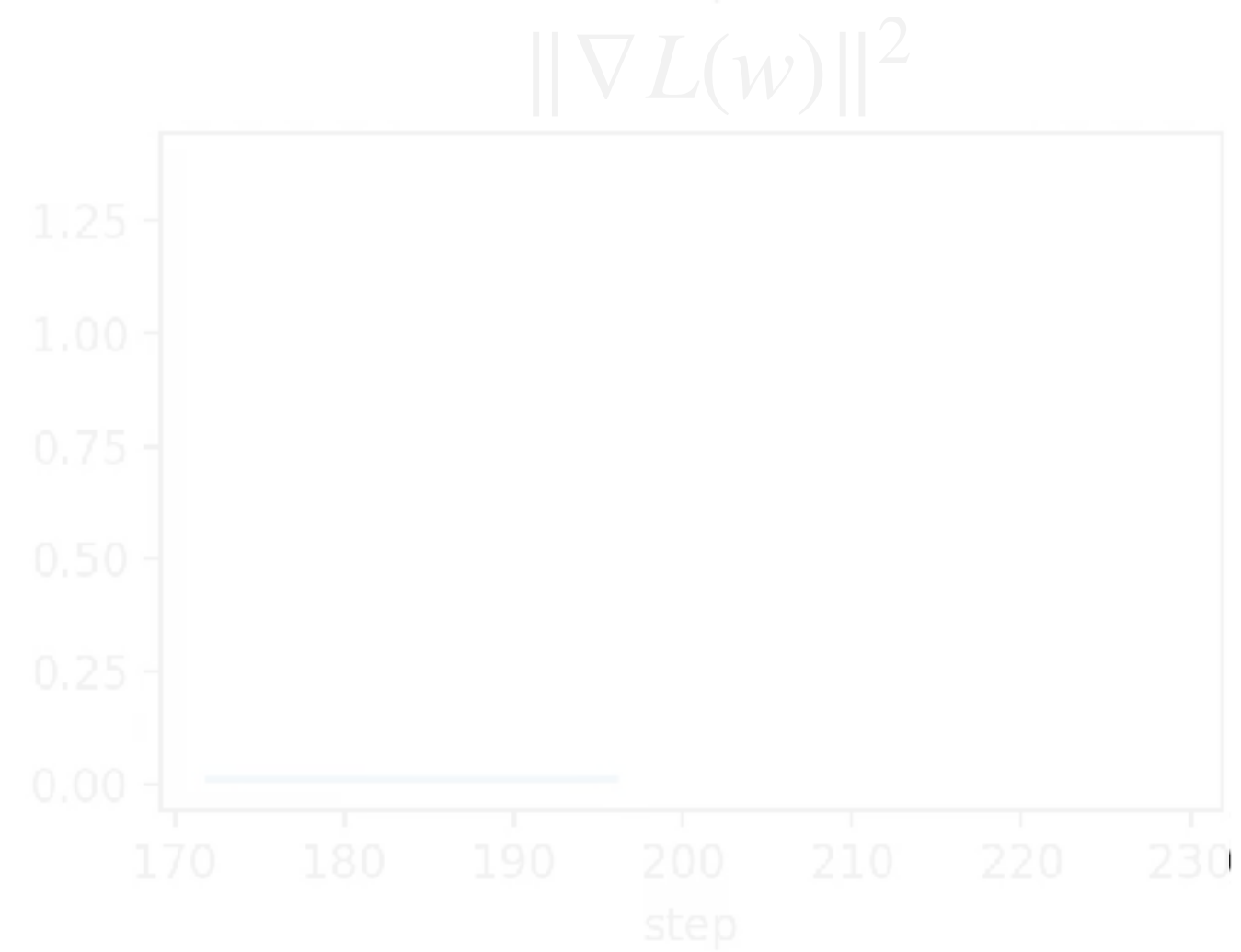
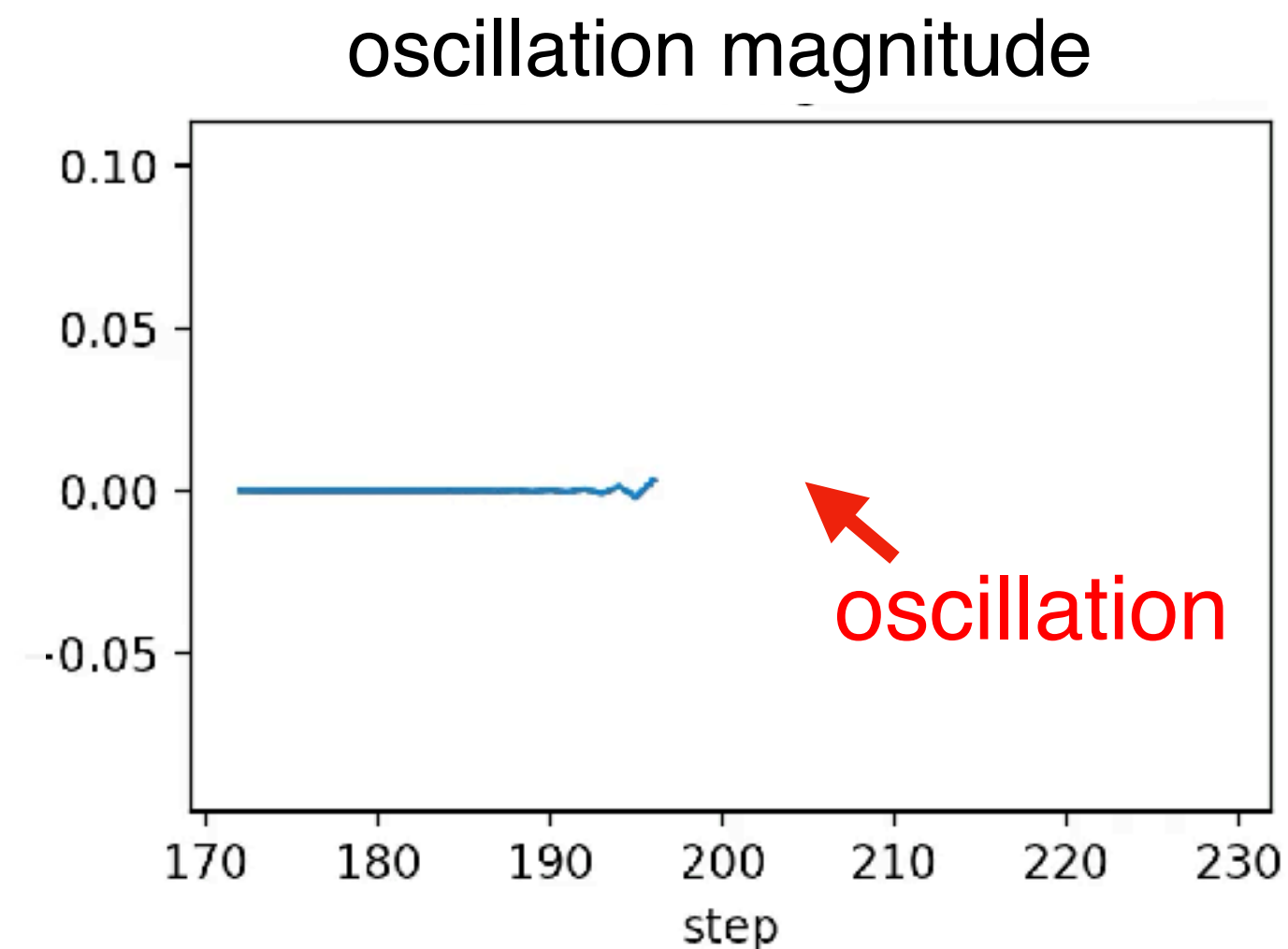
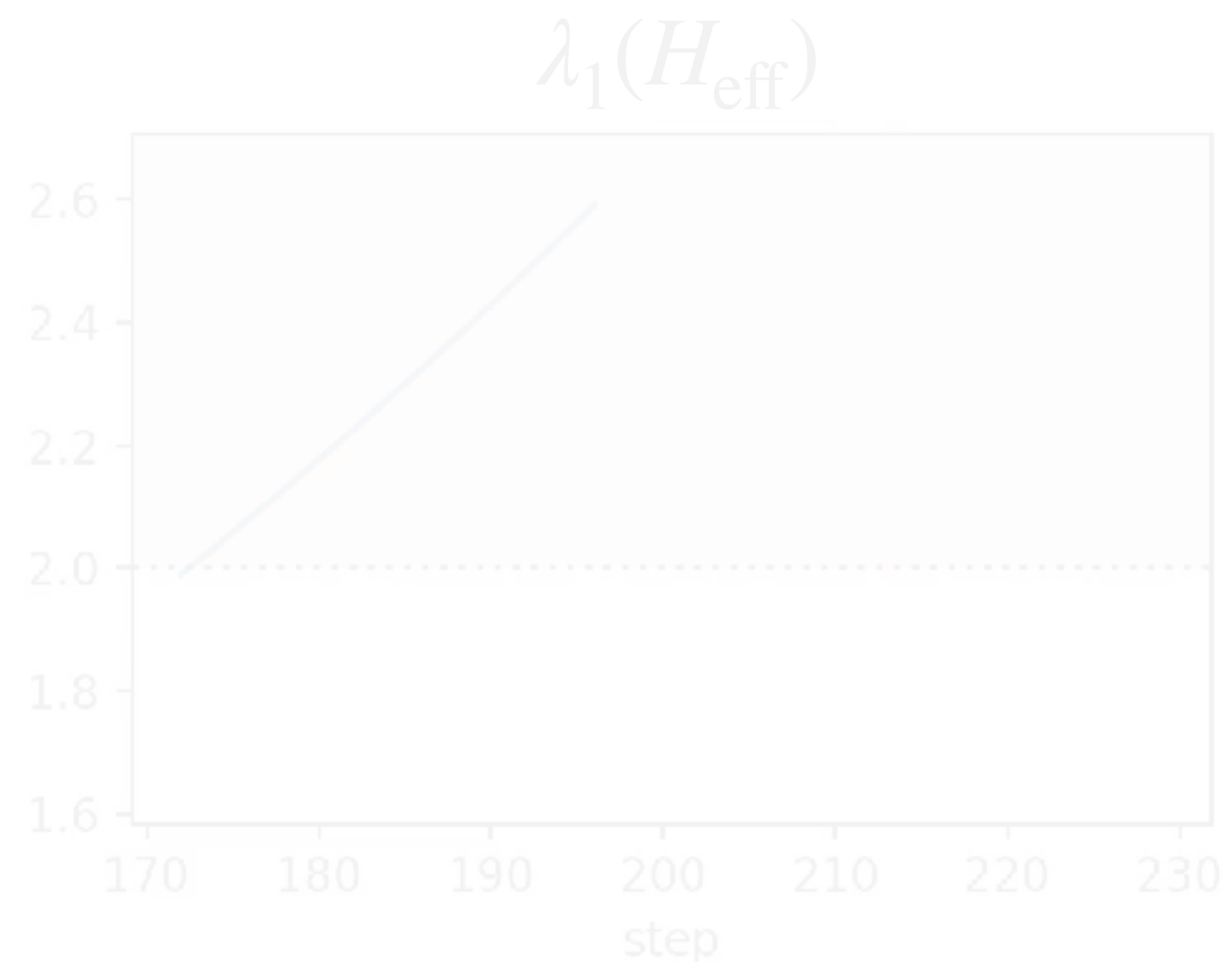
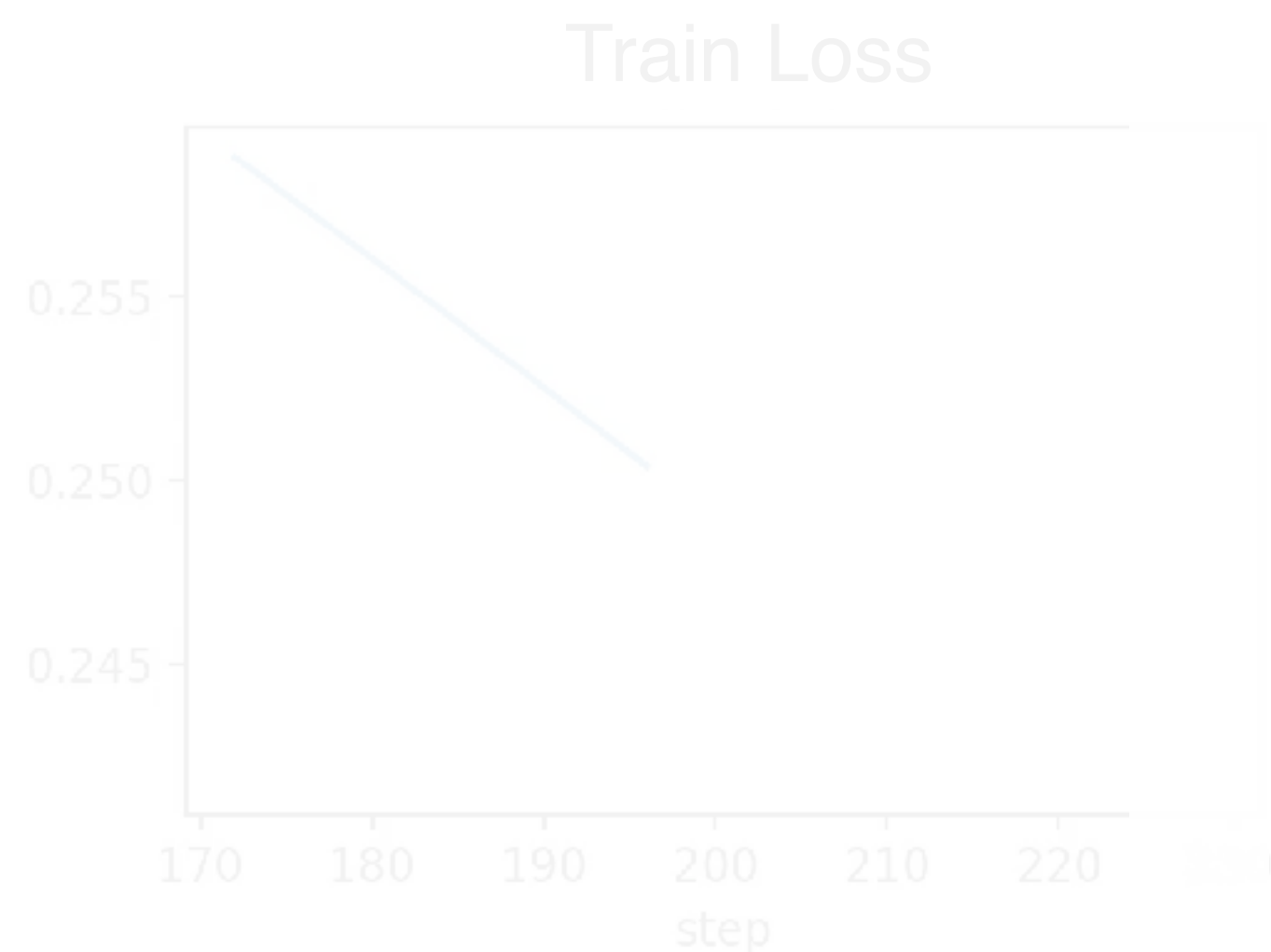
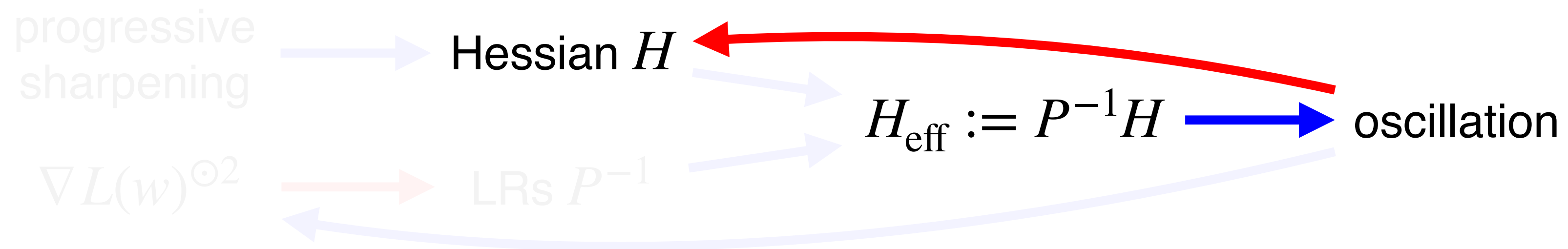


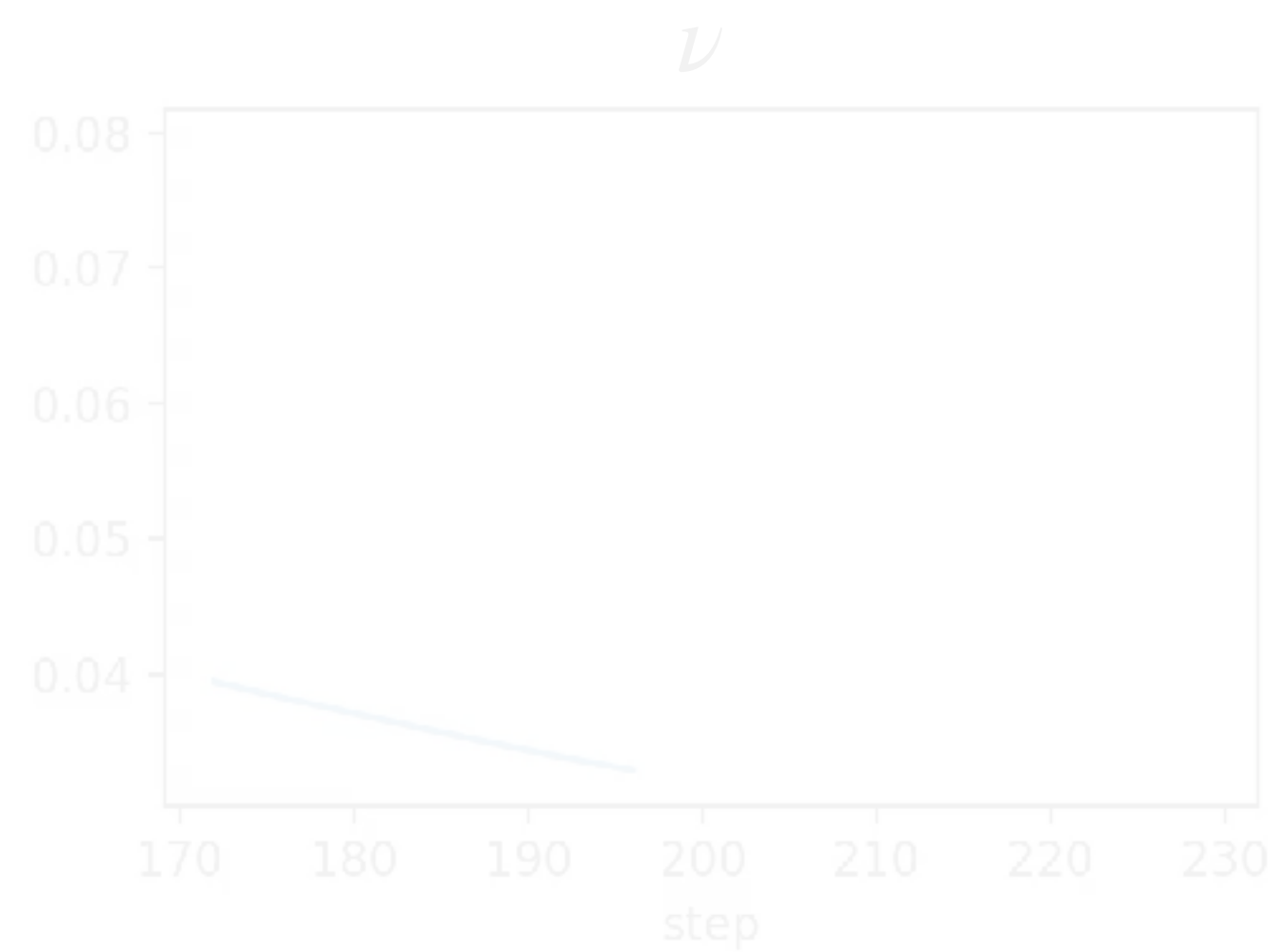
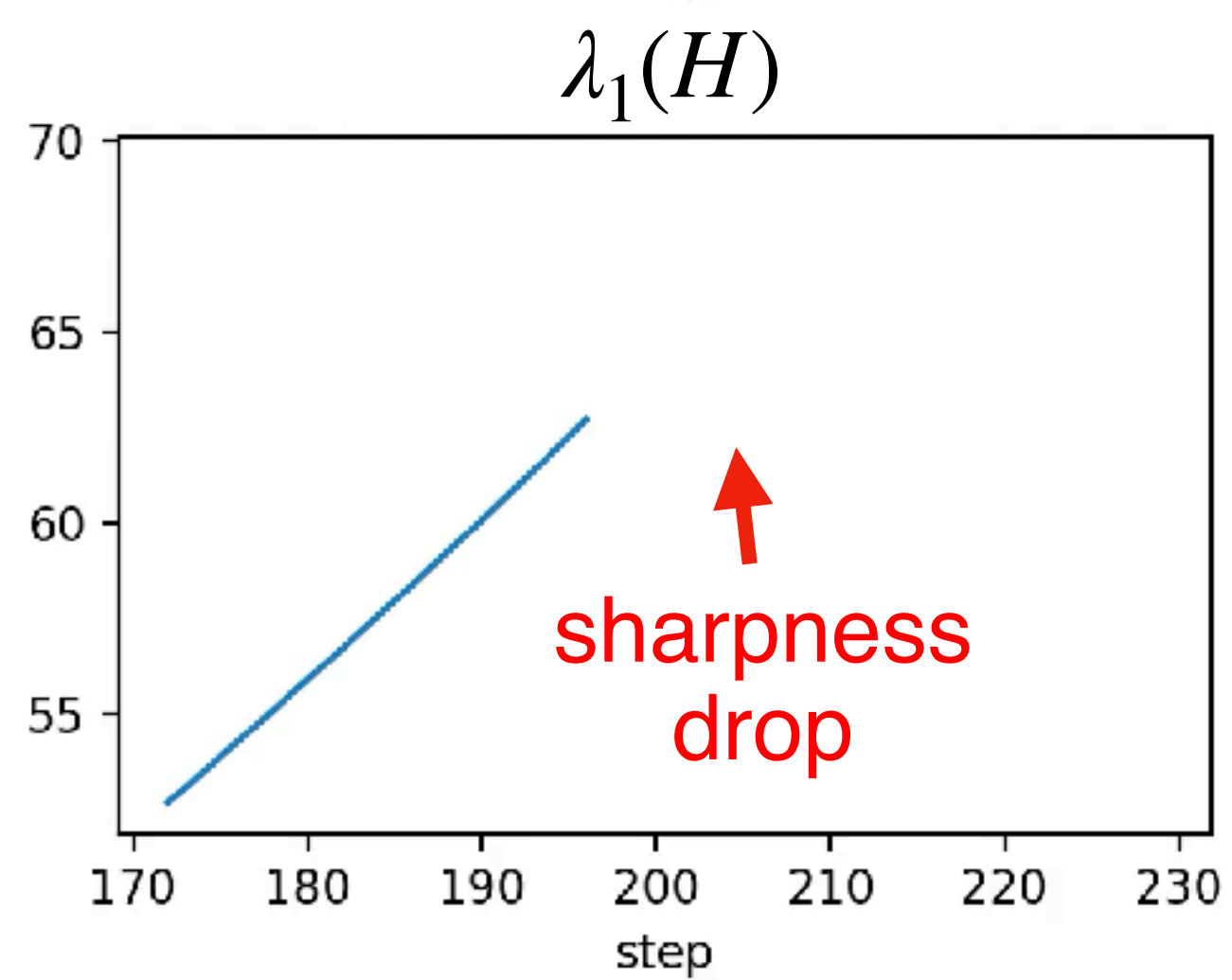
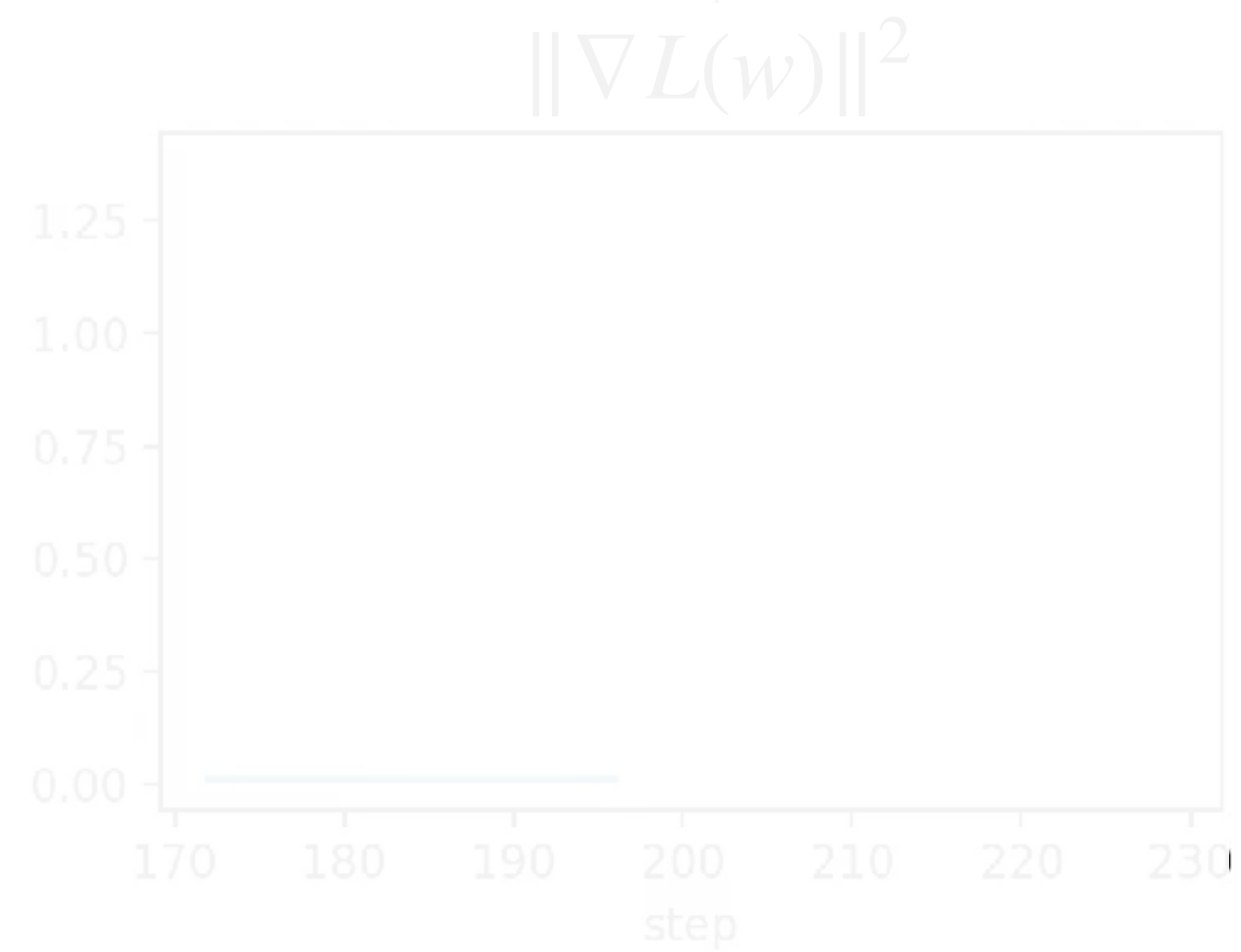
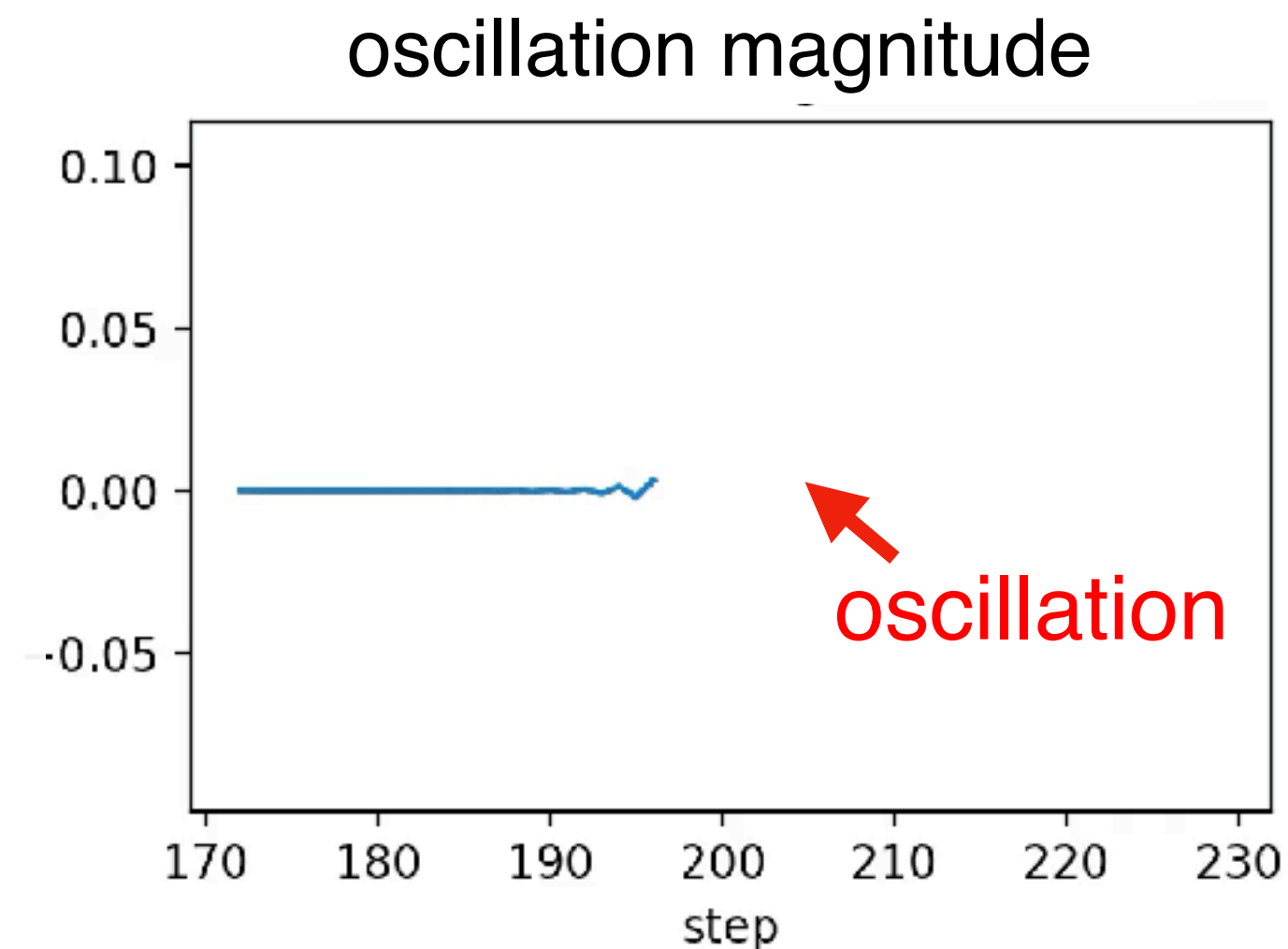
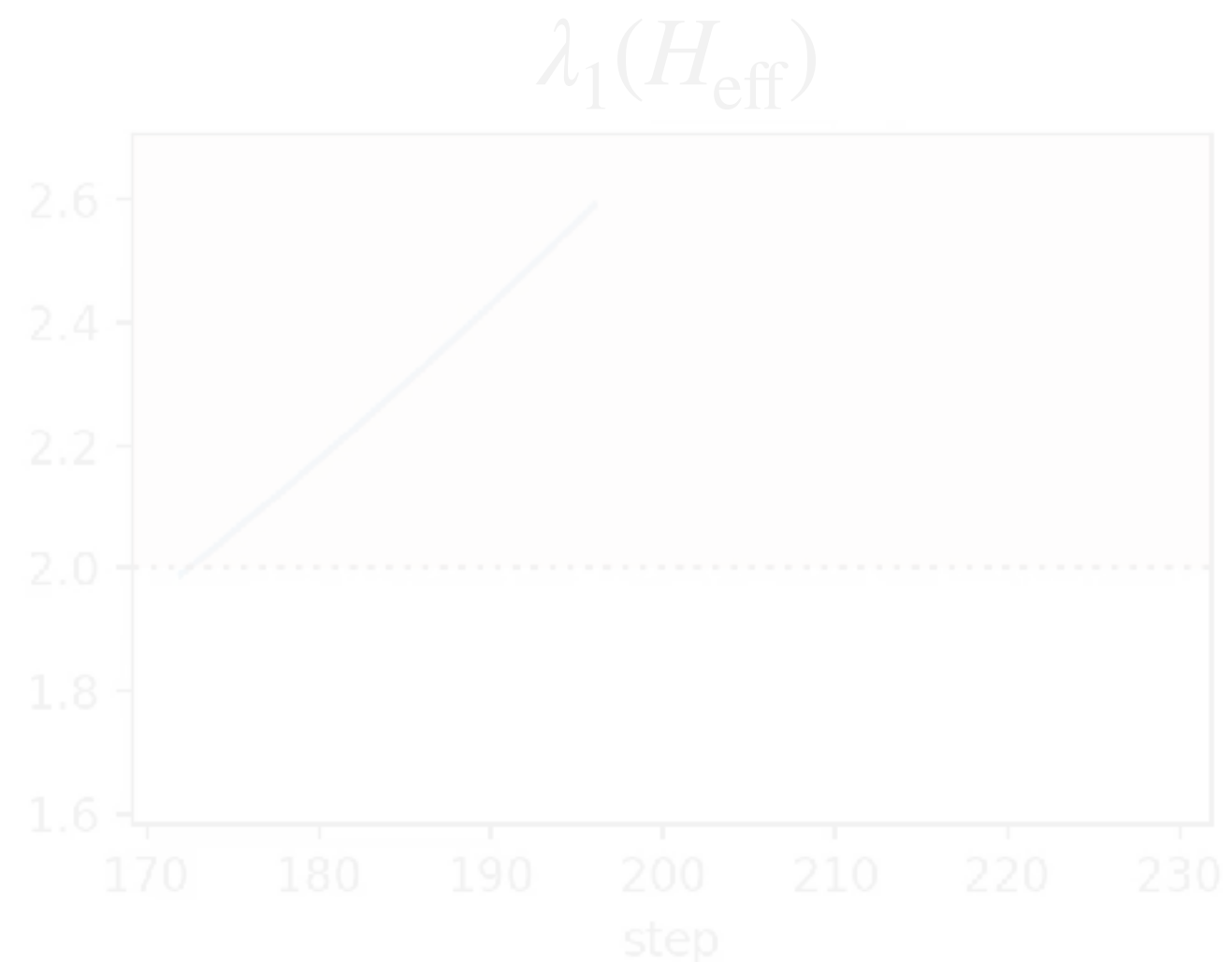
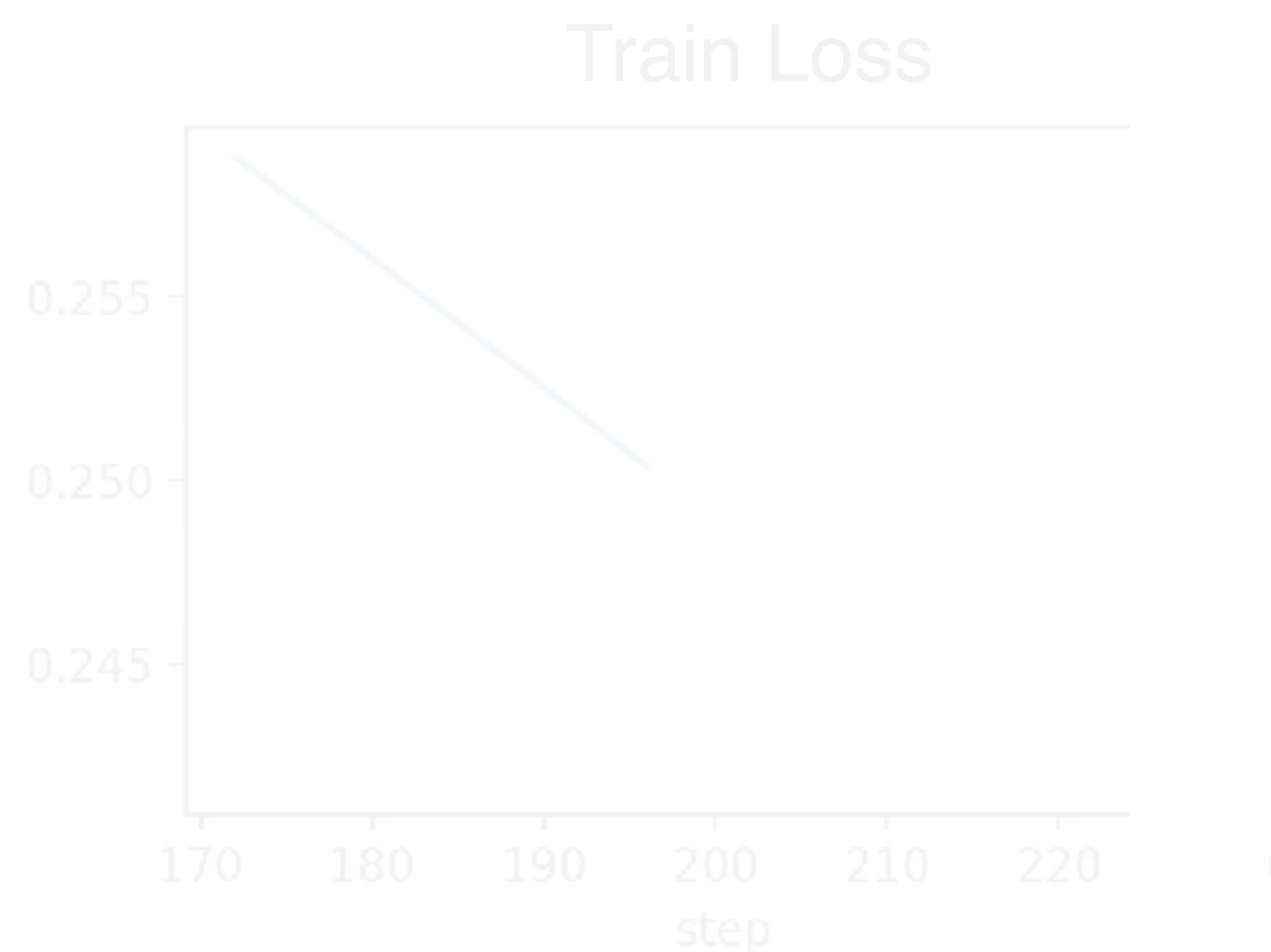
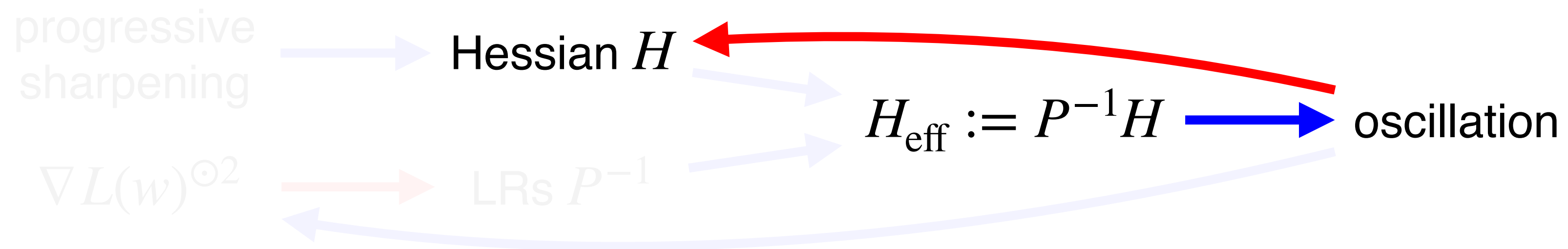
$\lambda_1(H)$

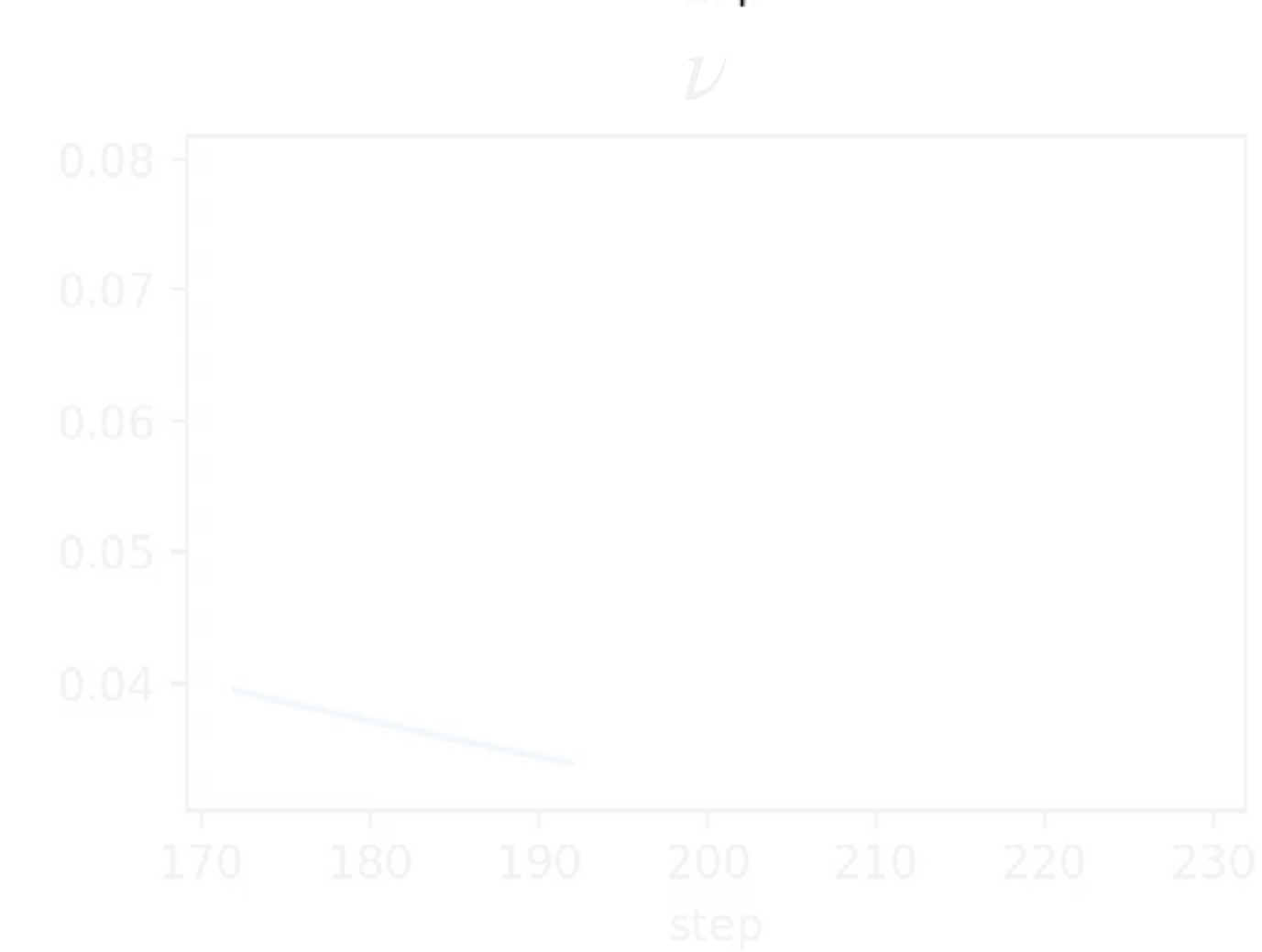
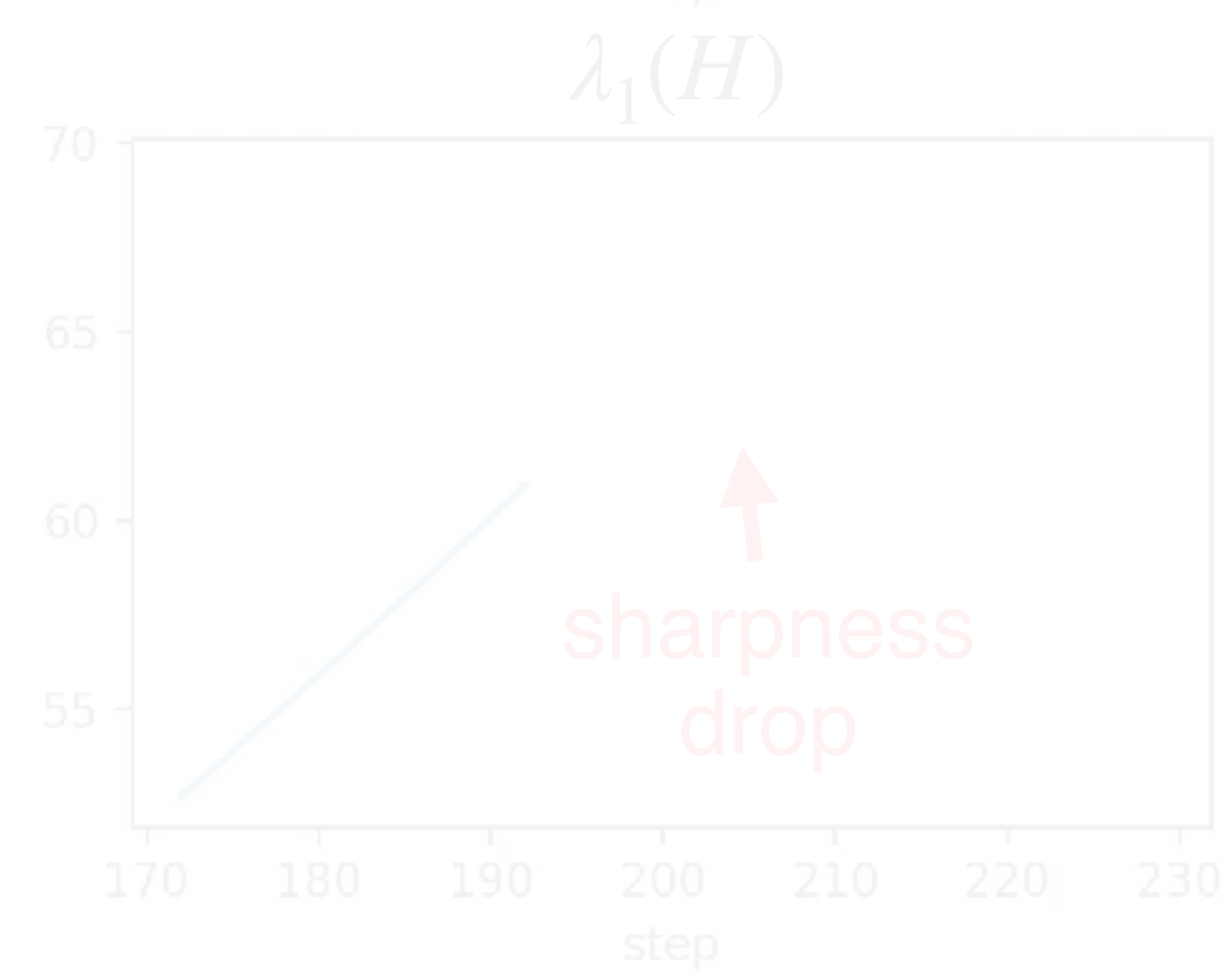
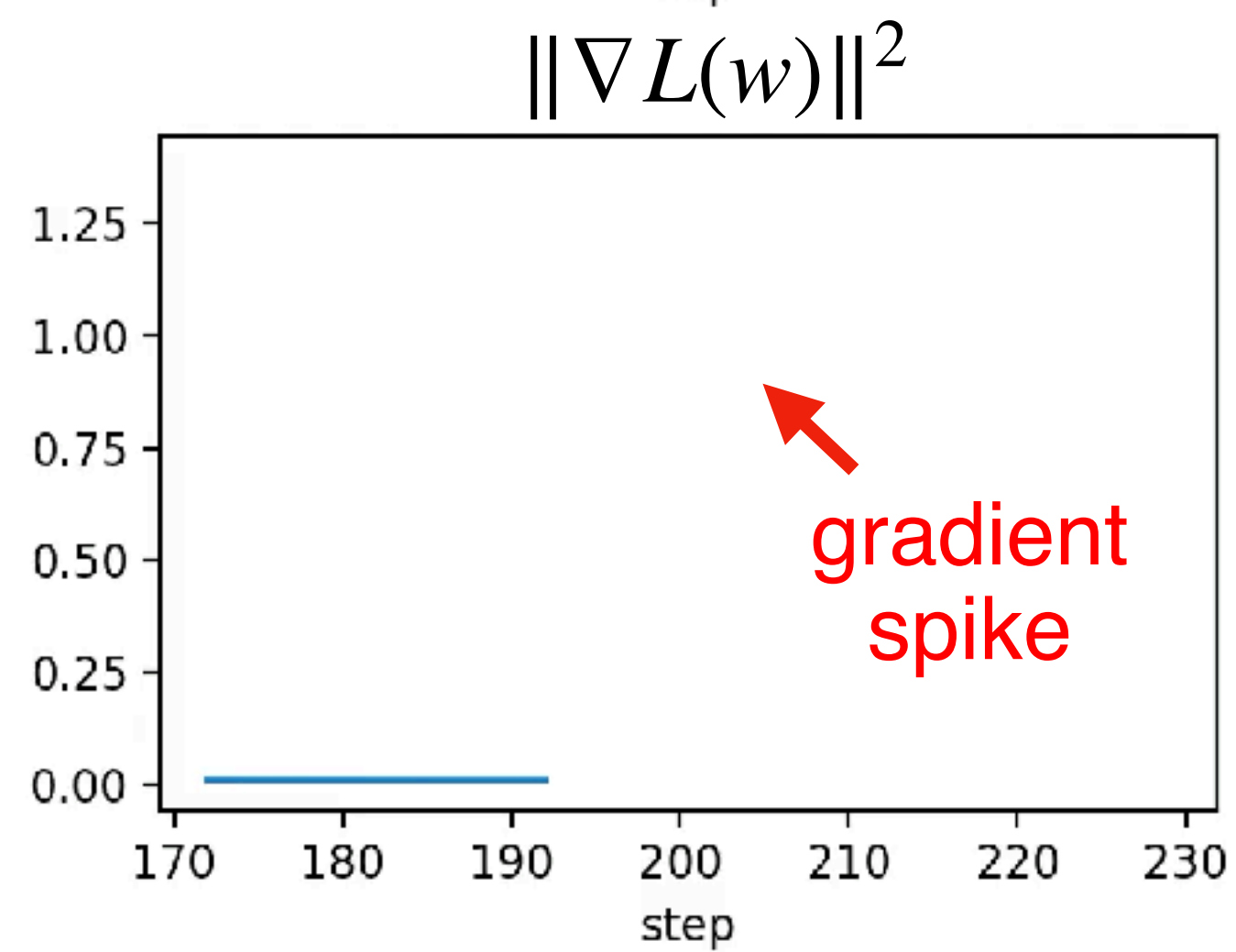
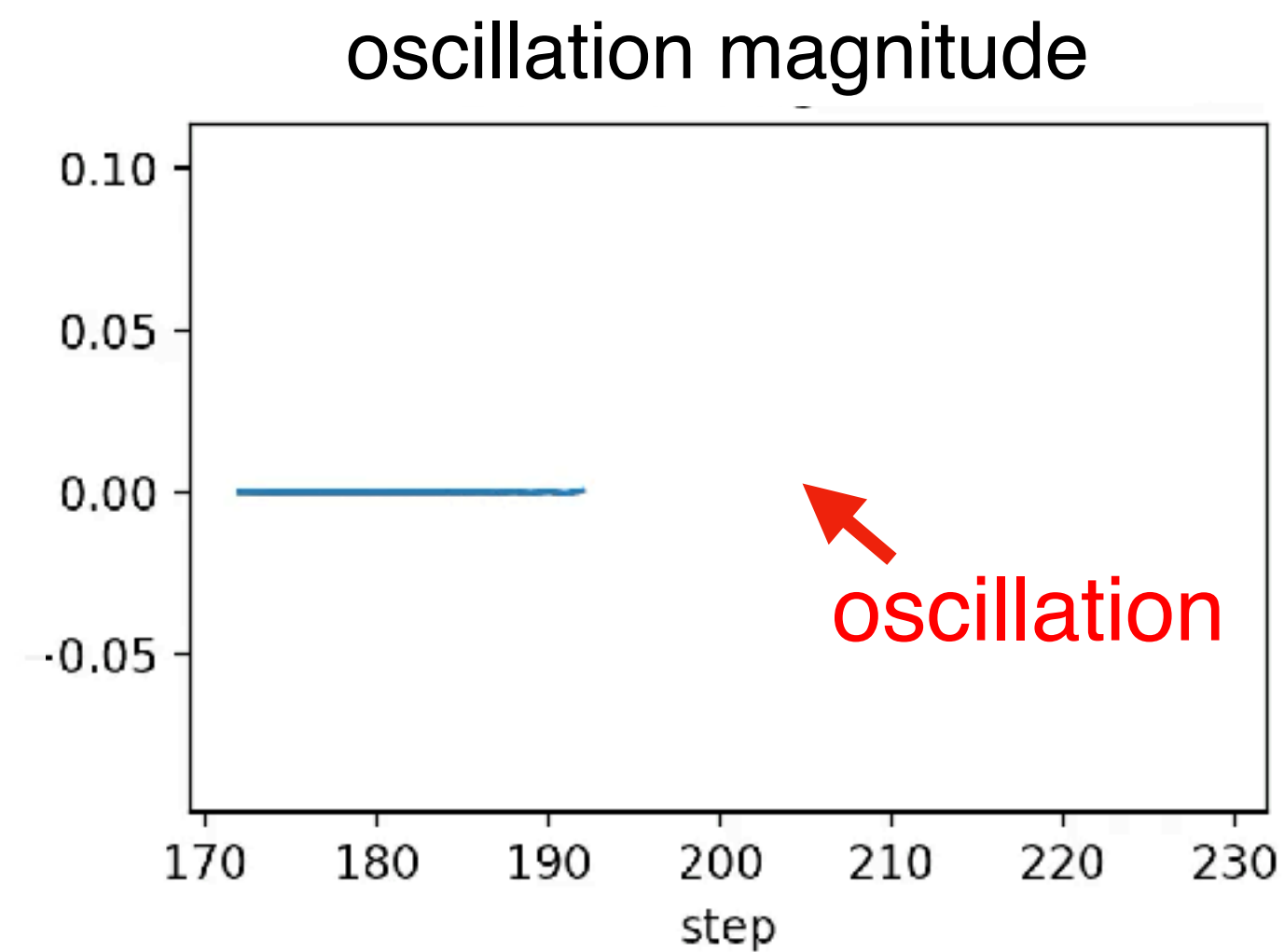
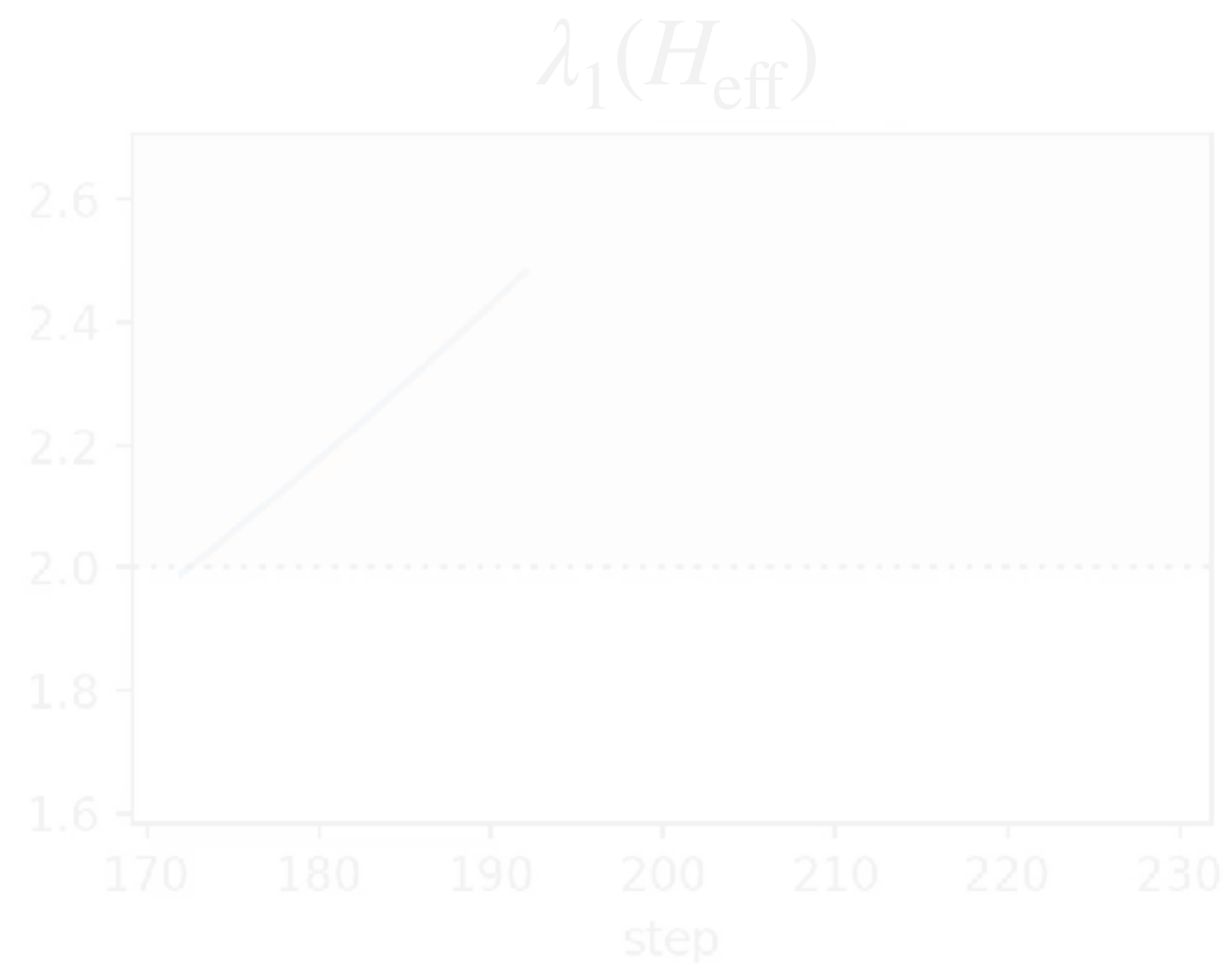
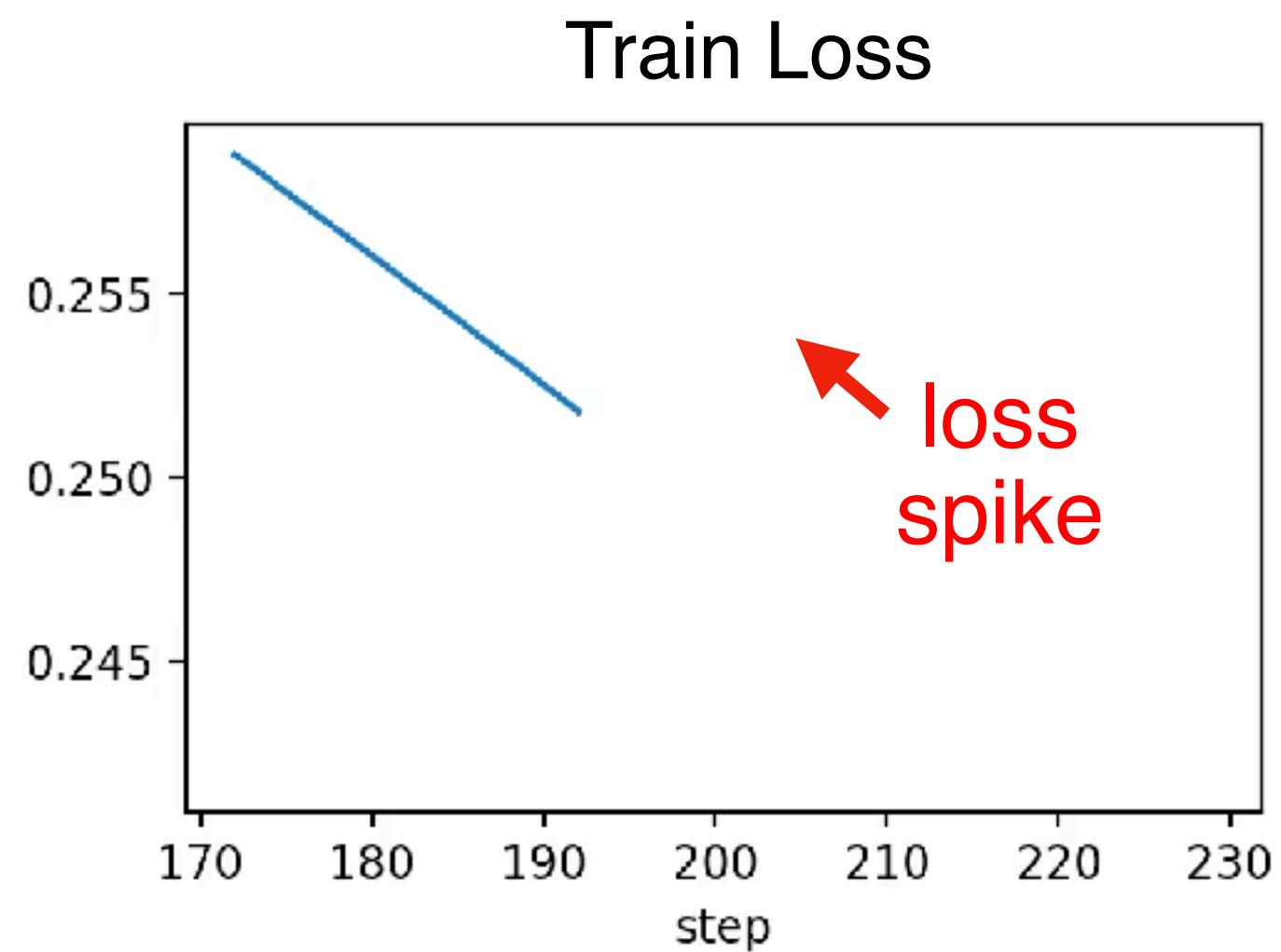


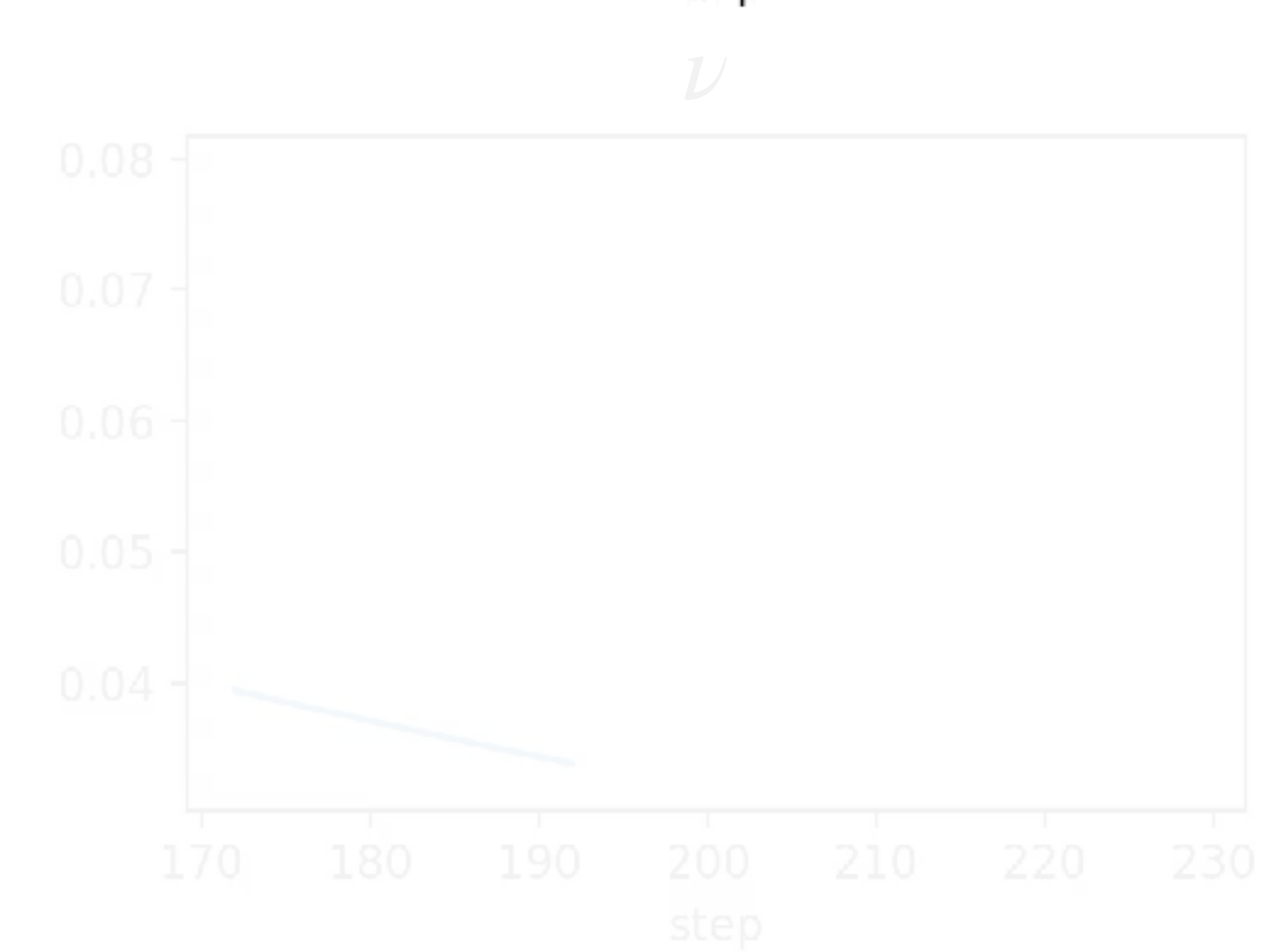
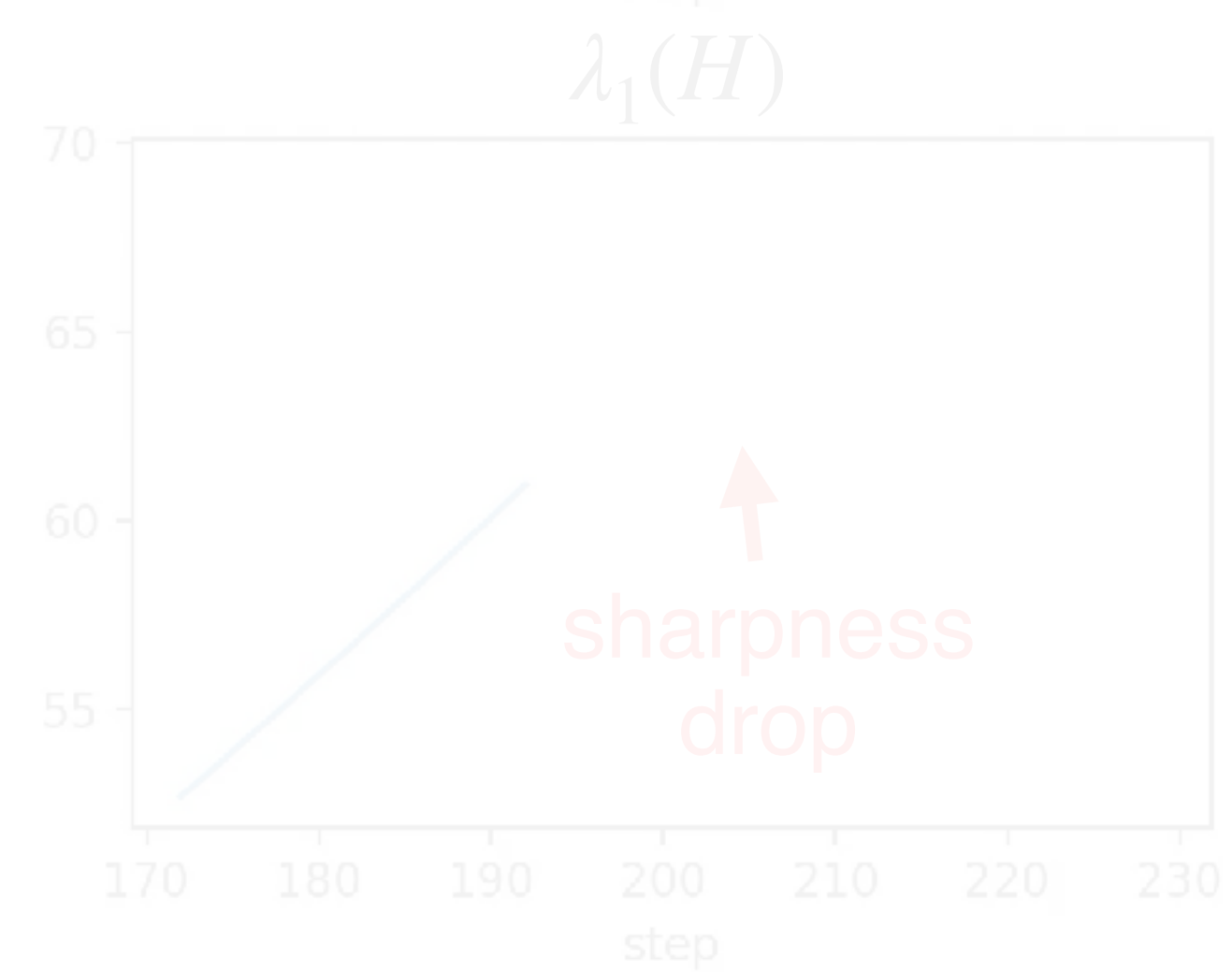
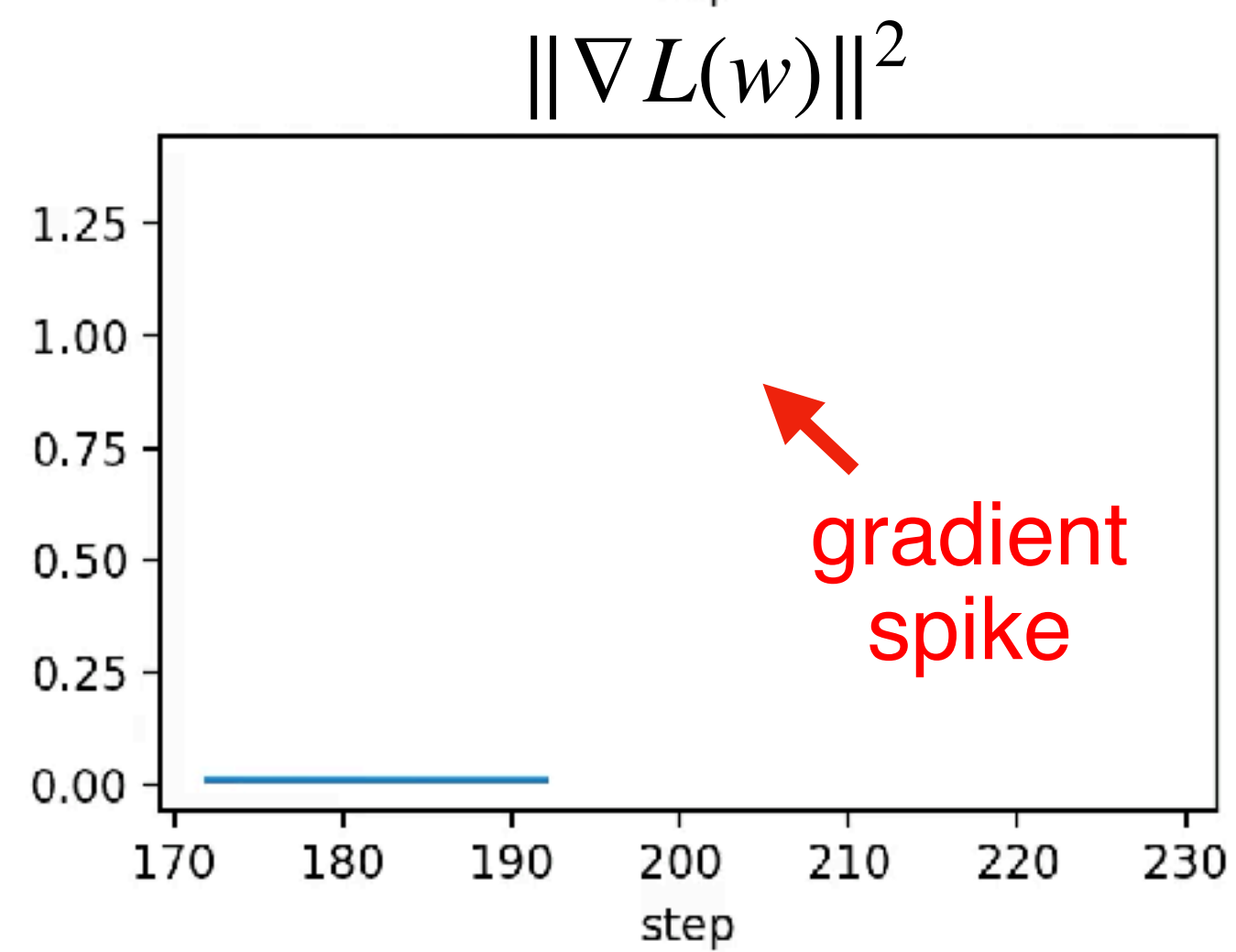
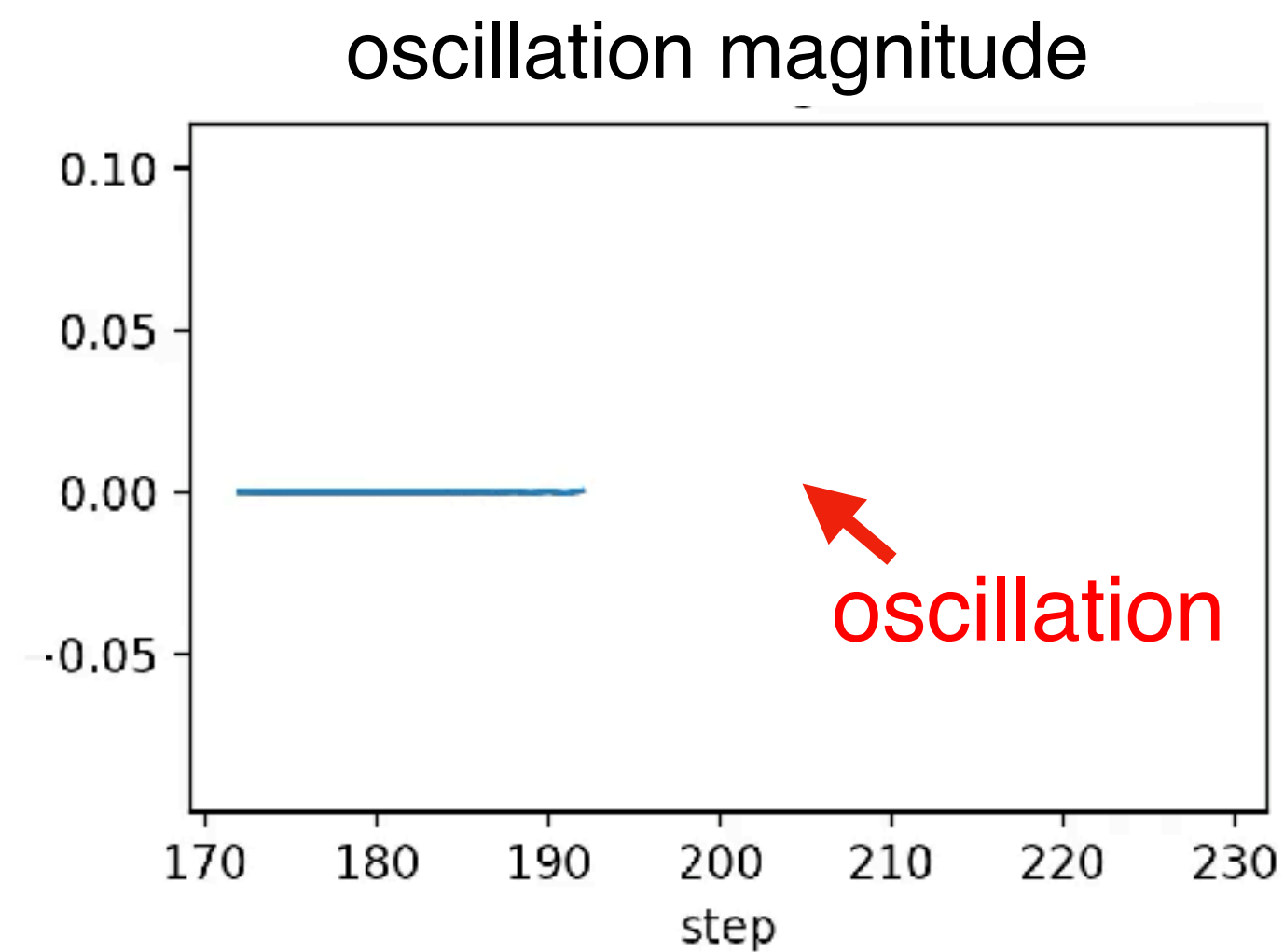
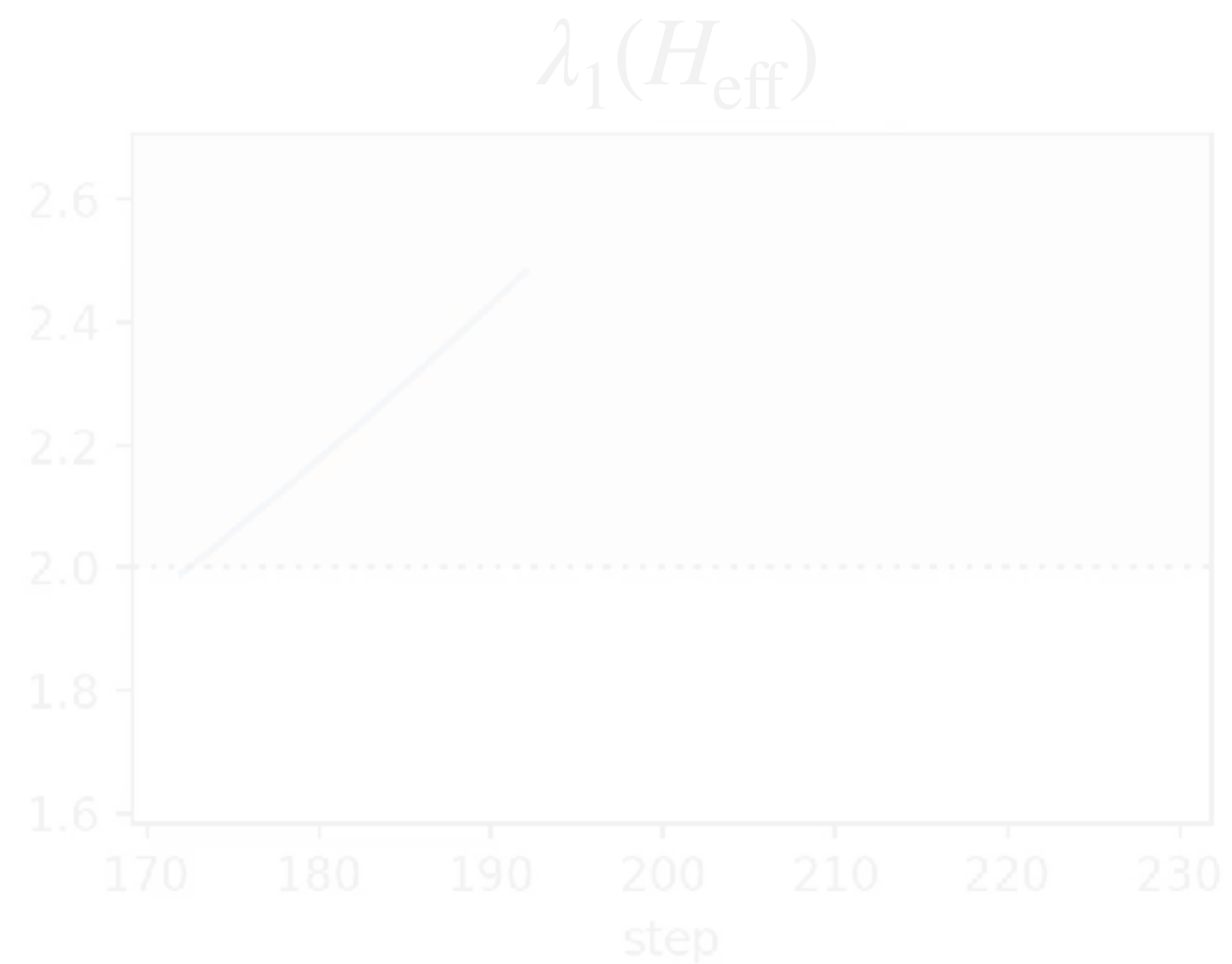
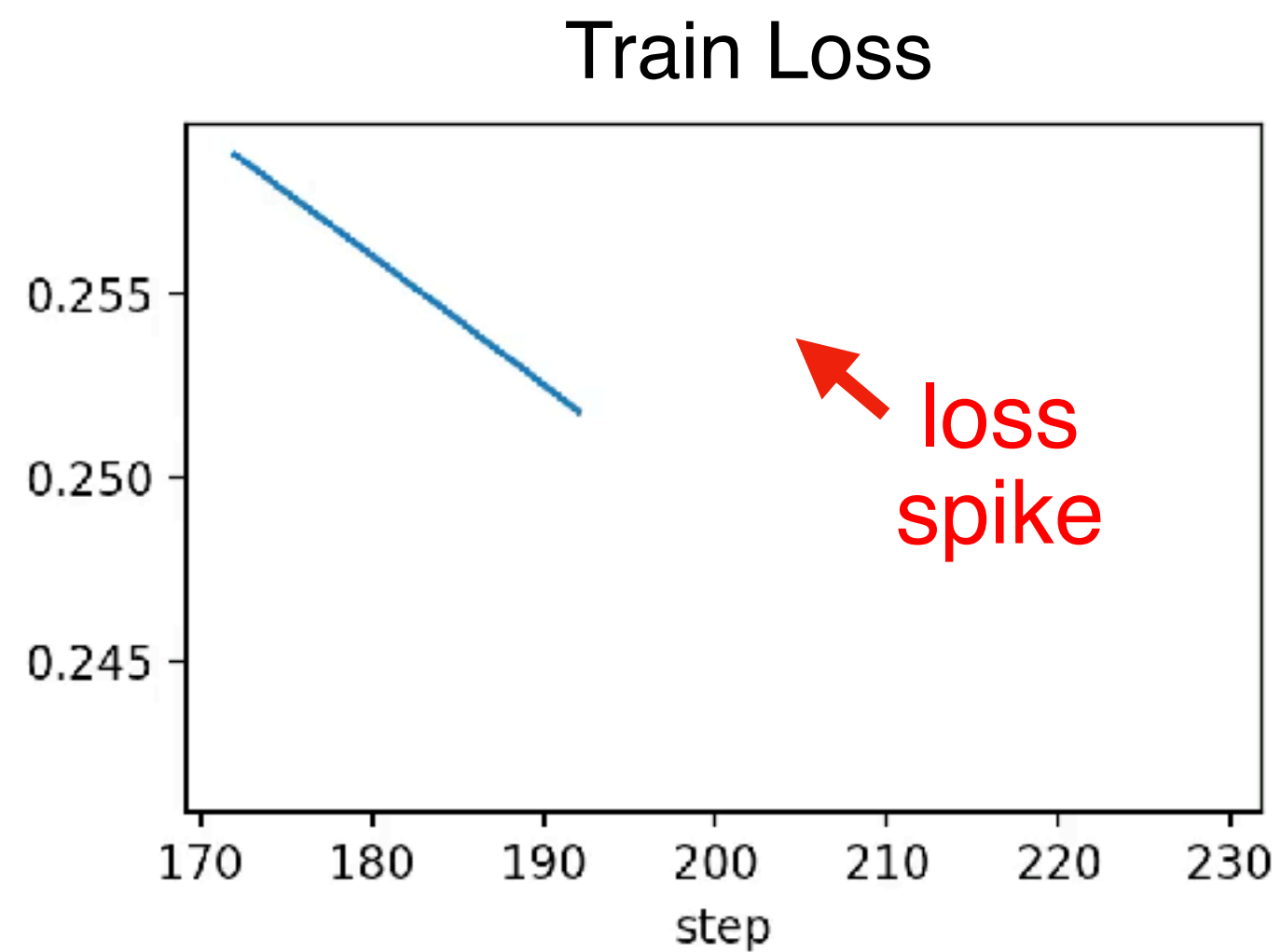
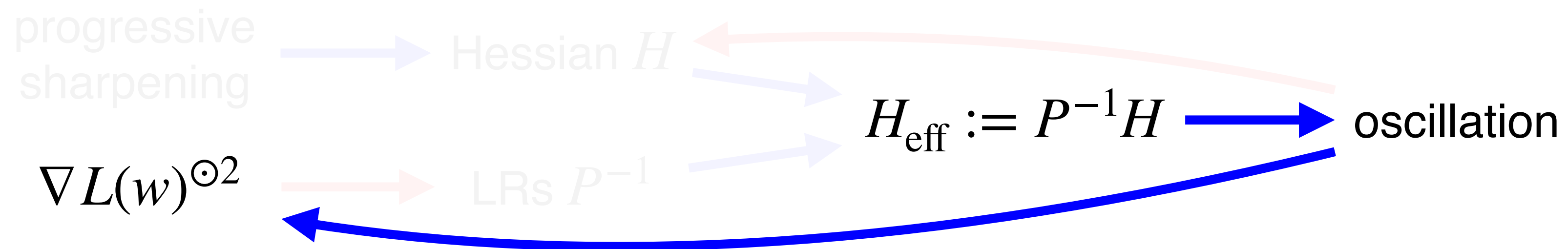
ν

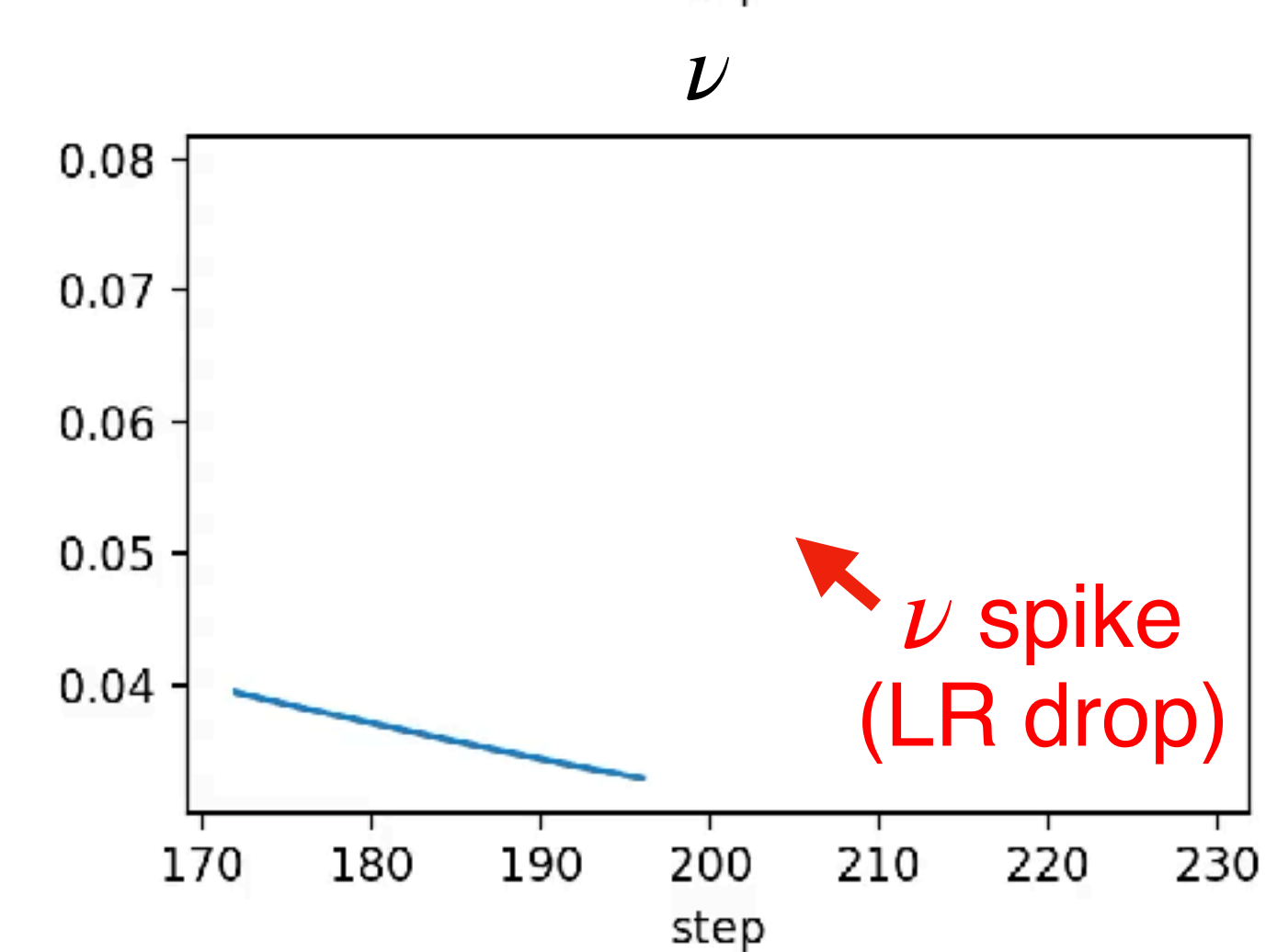
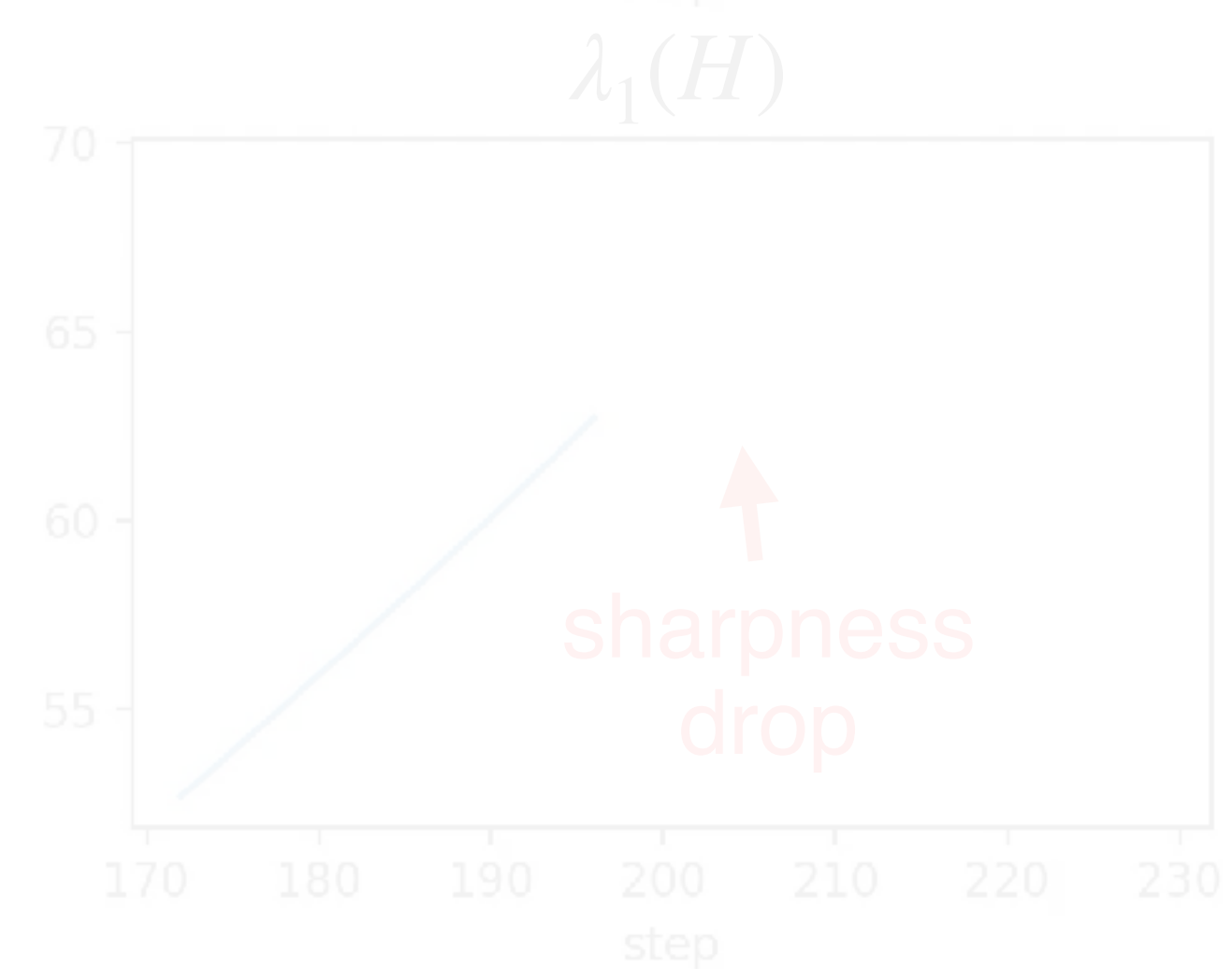
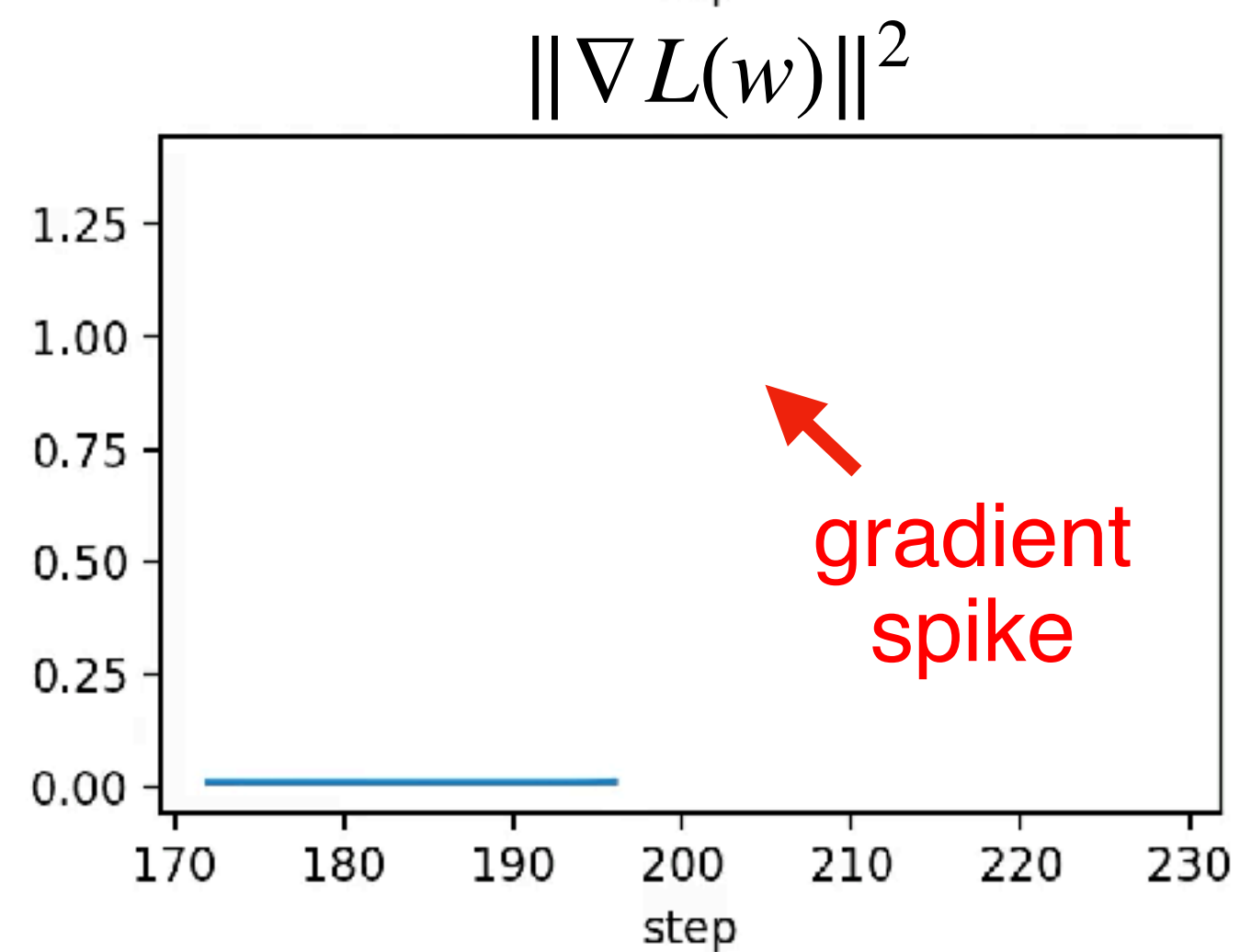
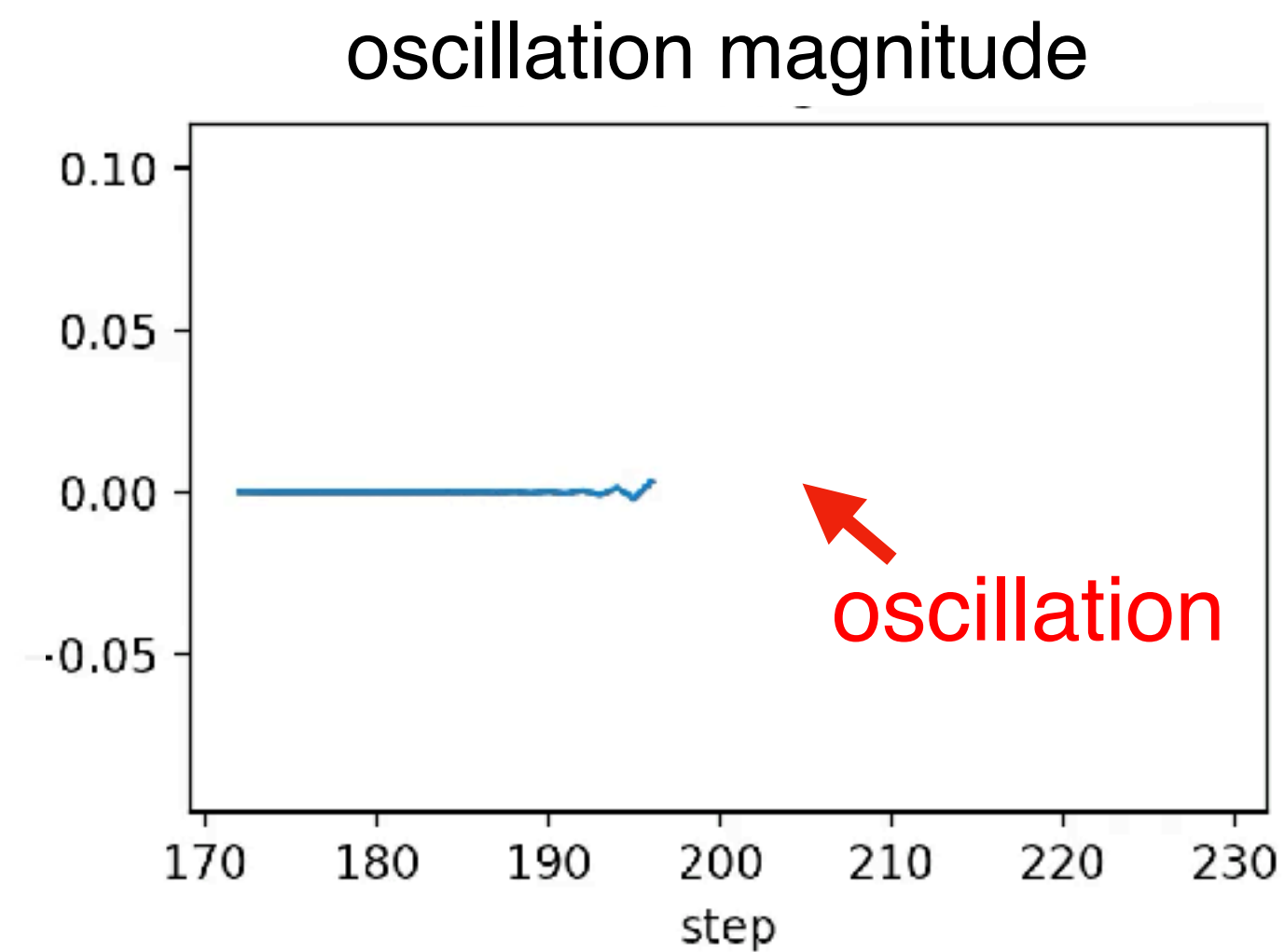
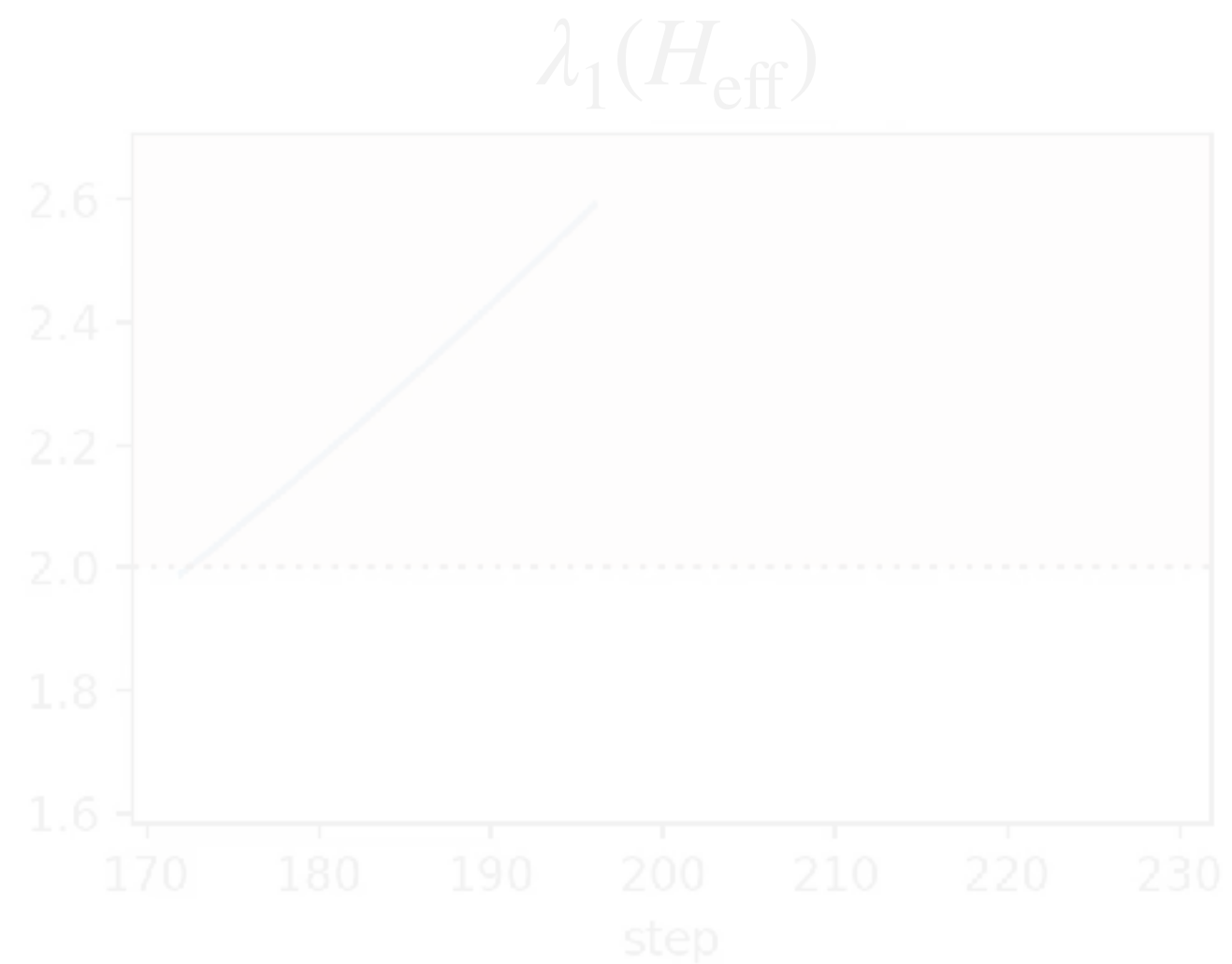
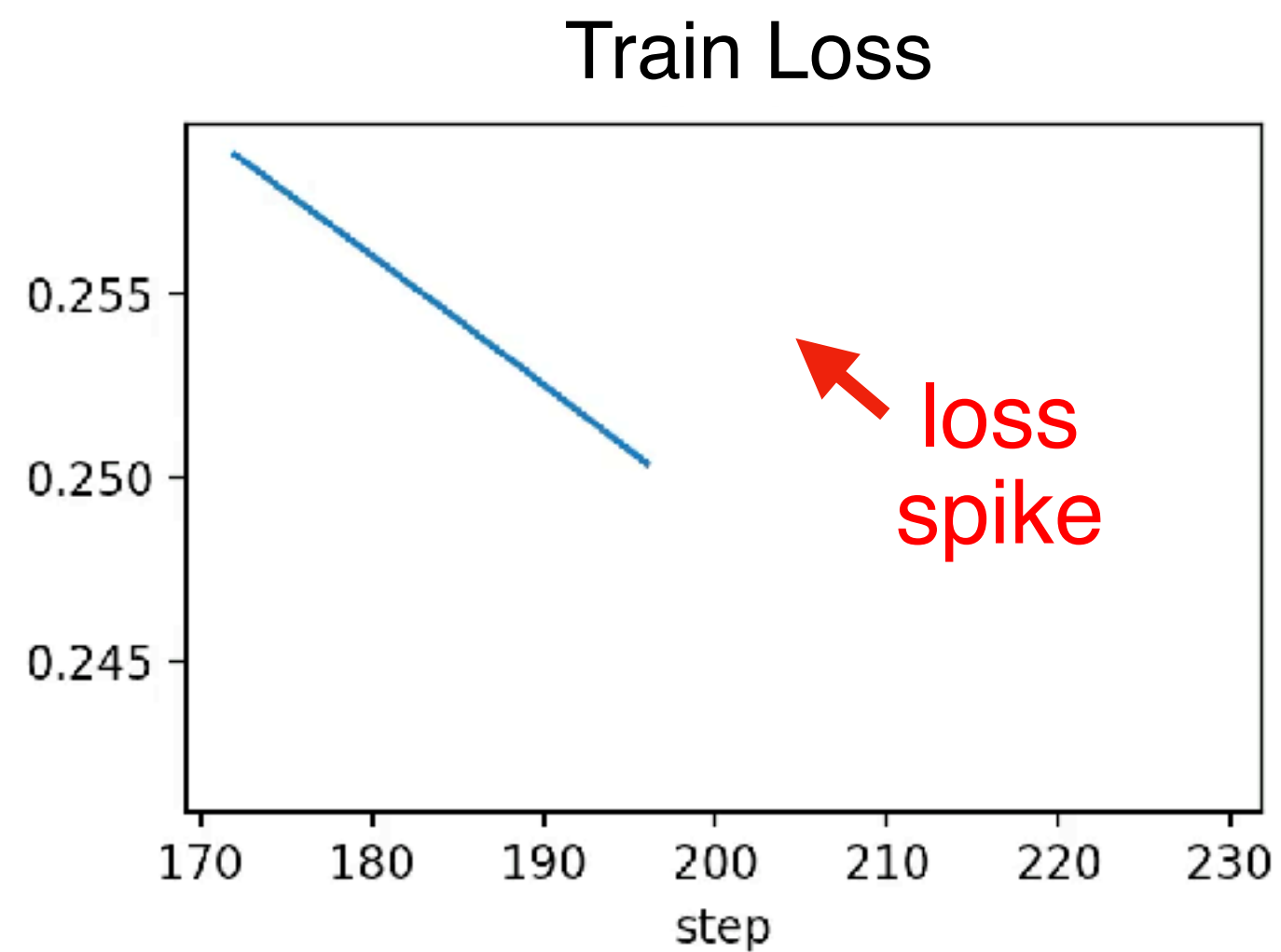
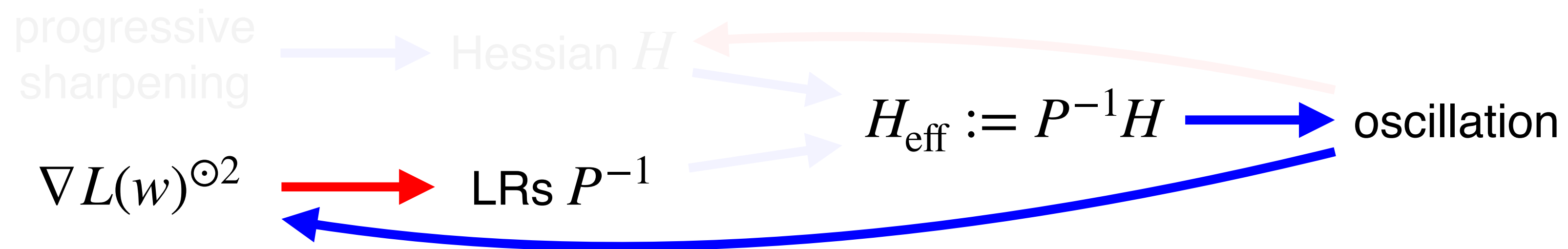


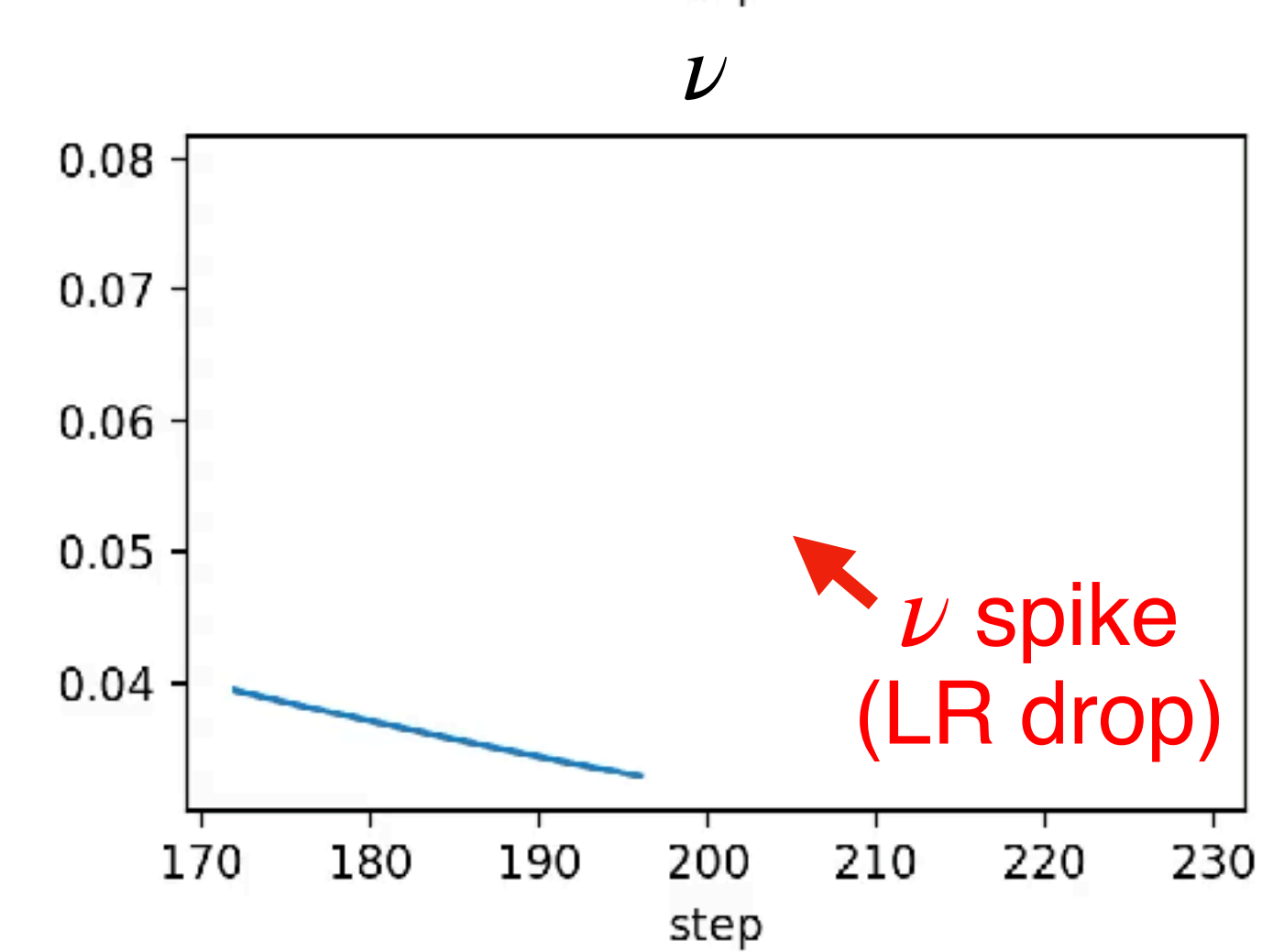
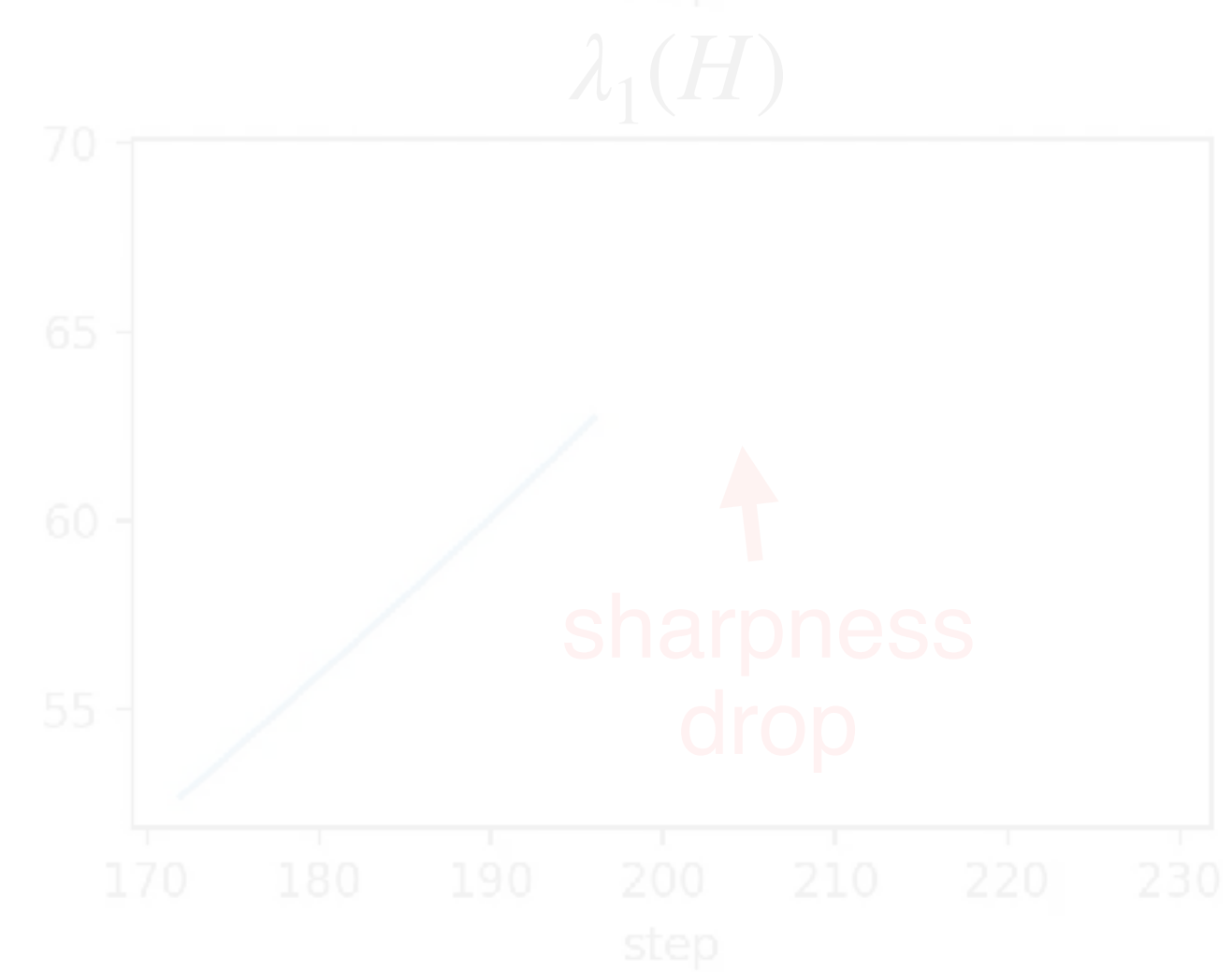
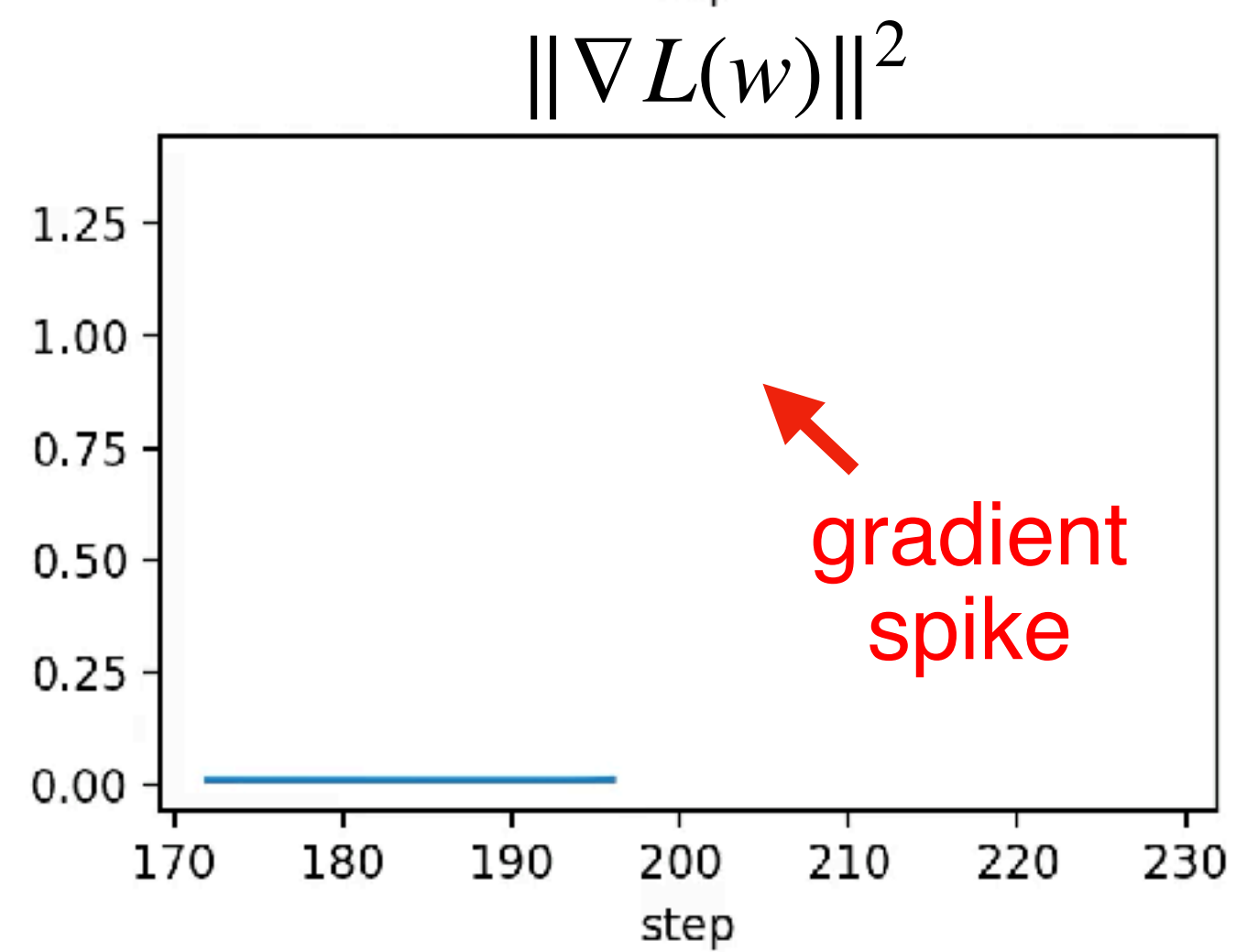
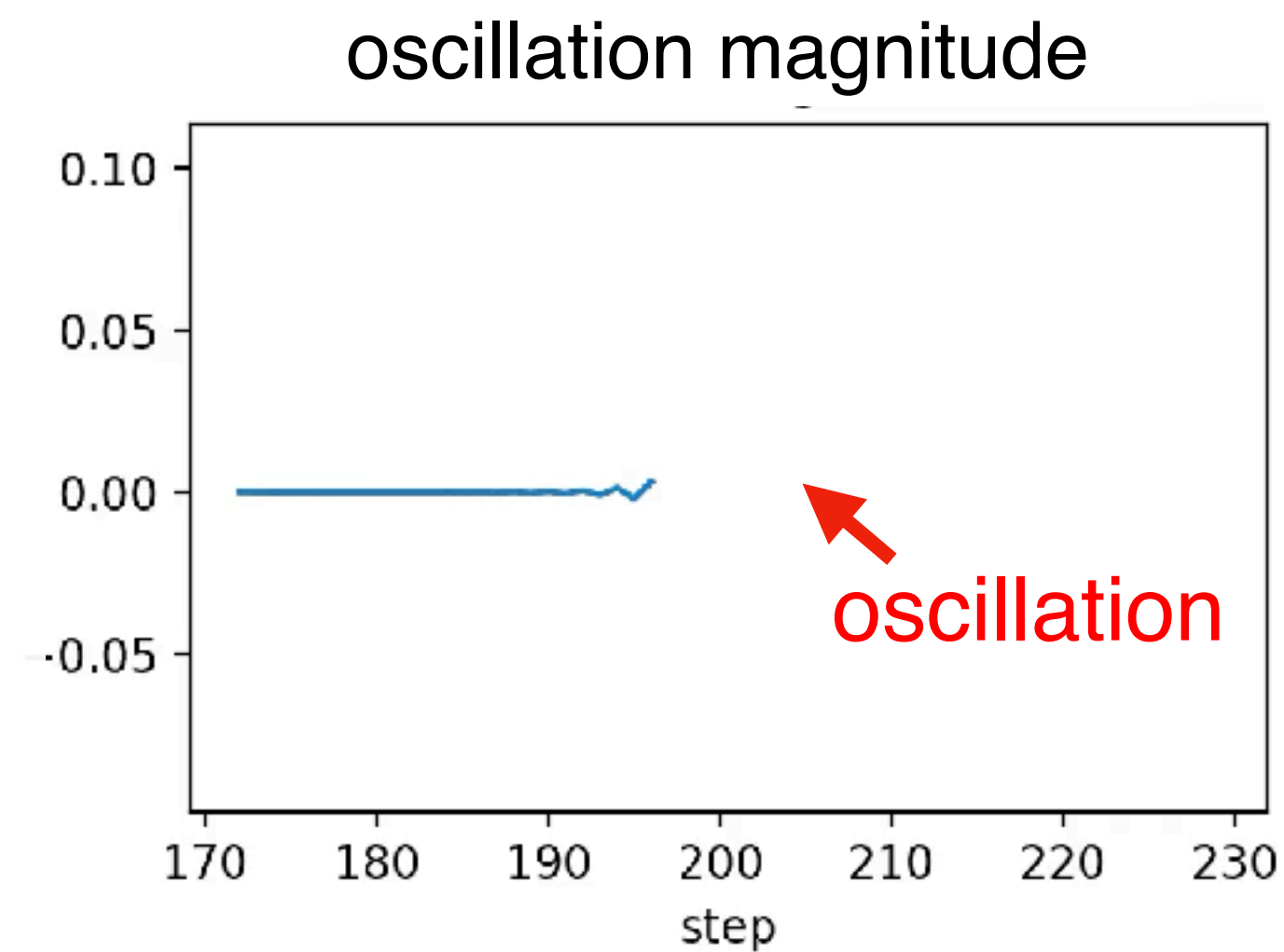
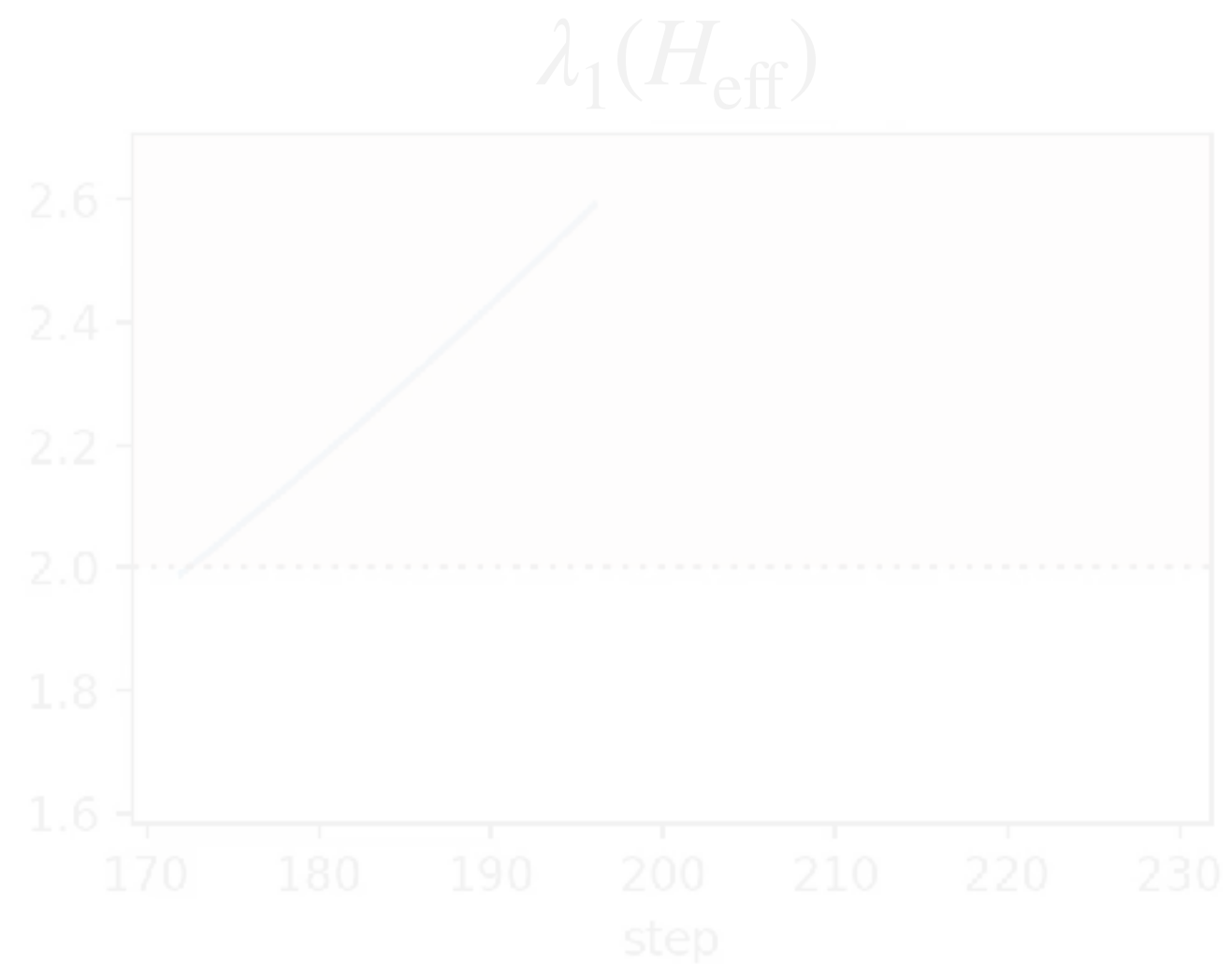
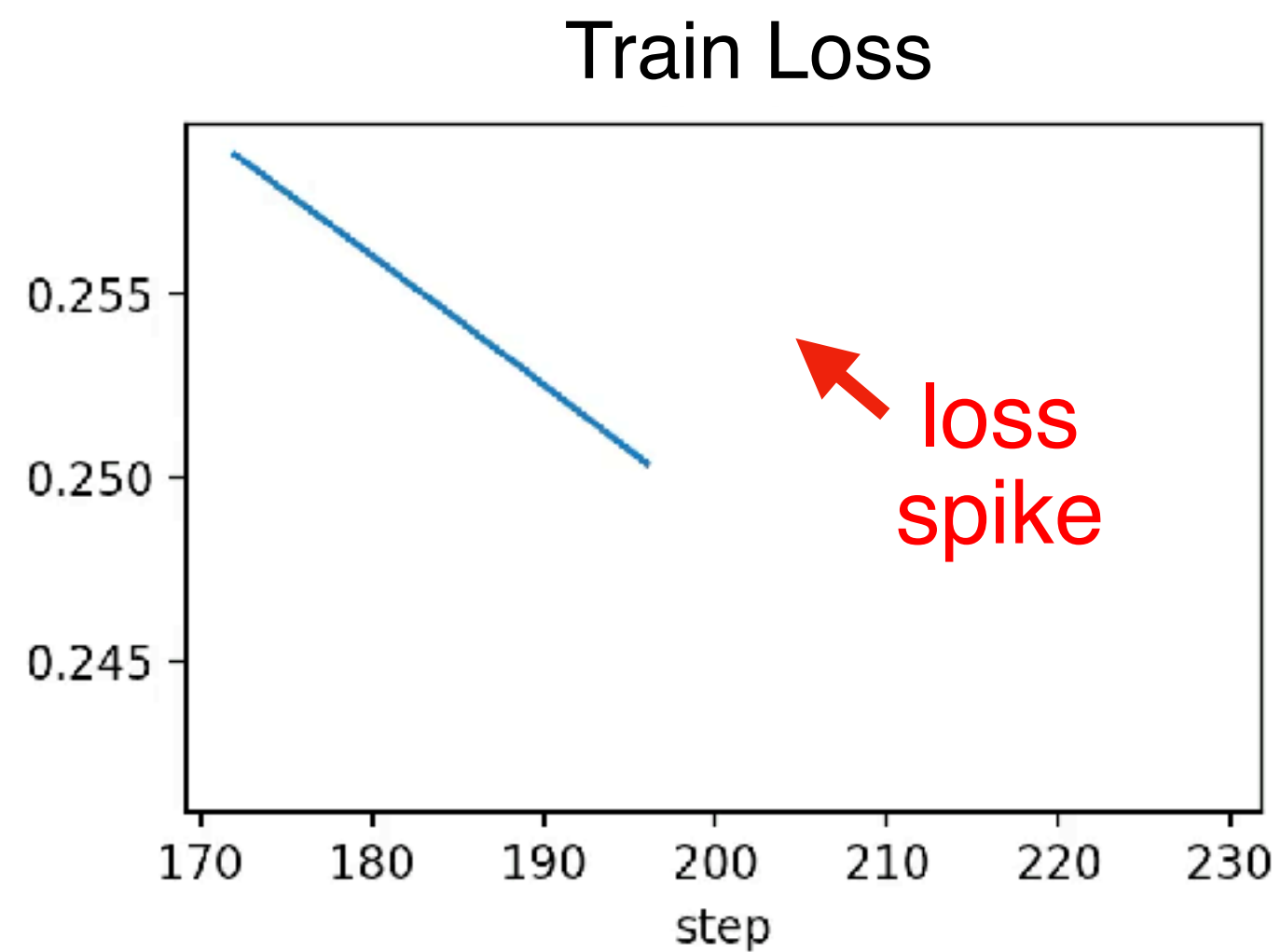
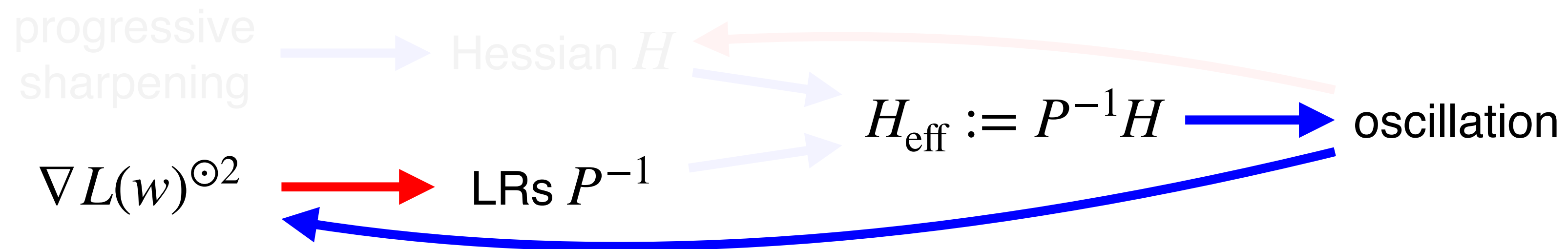


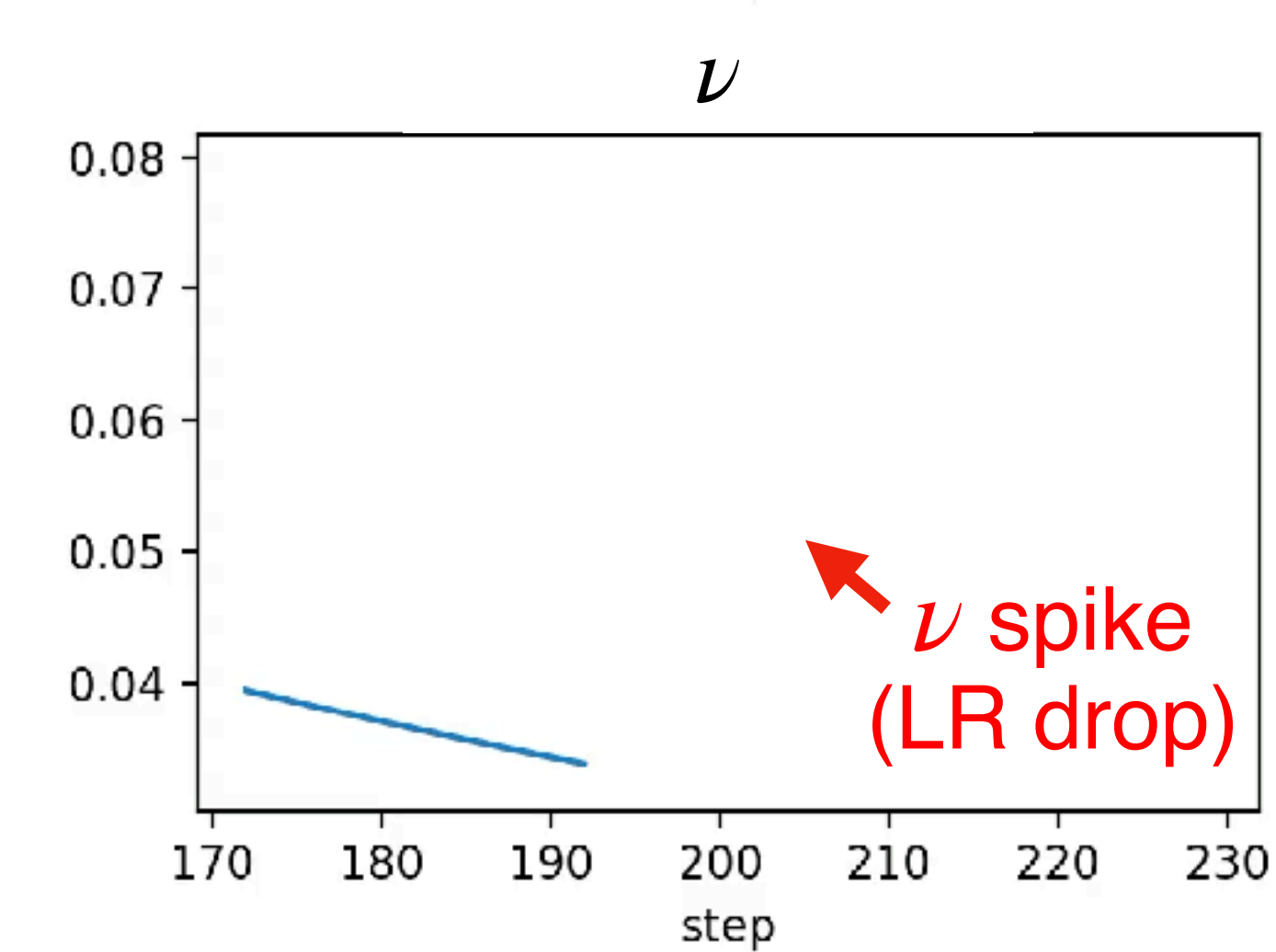
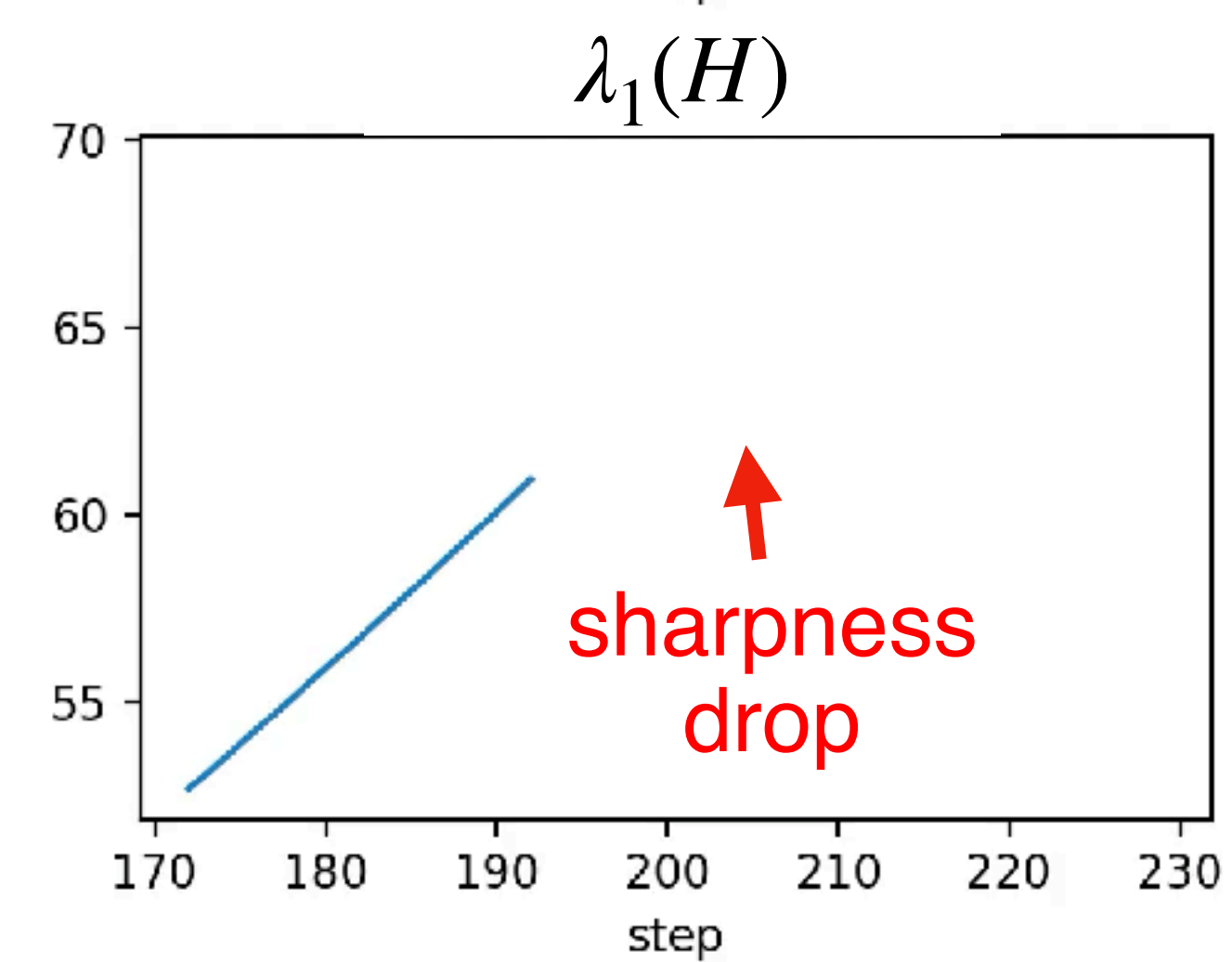
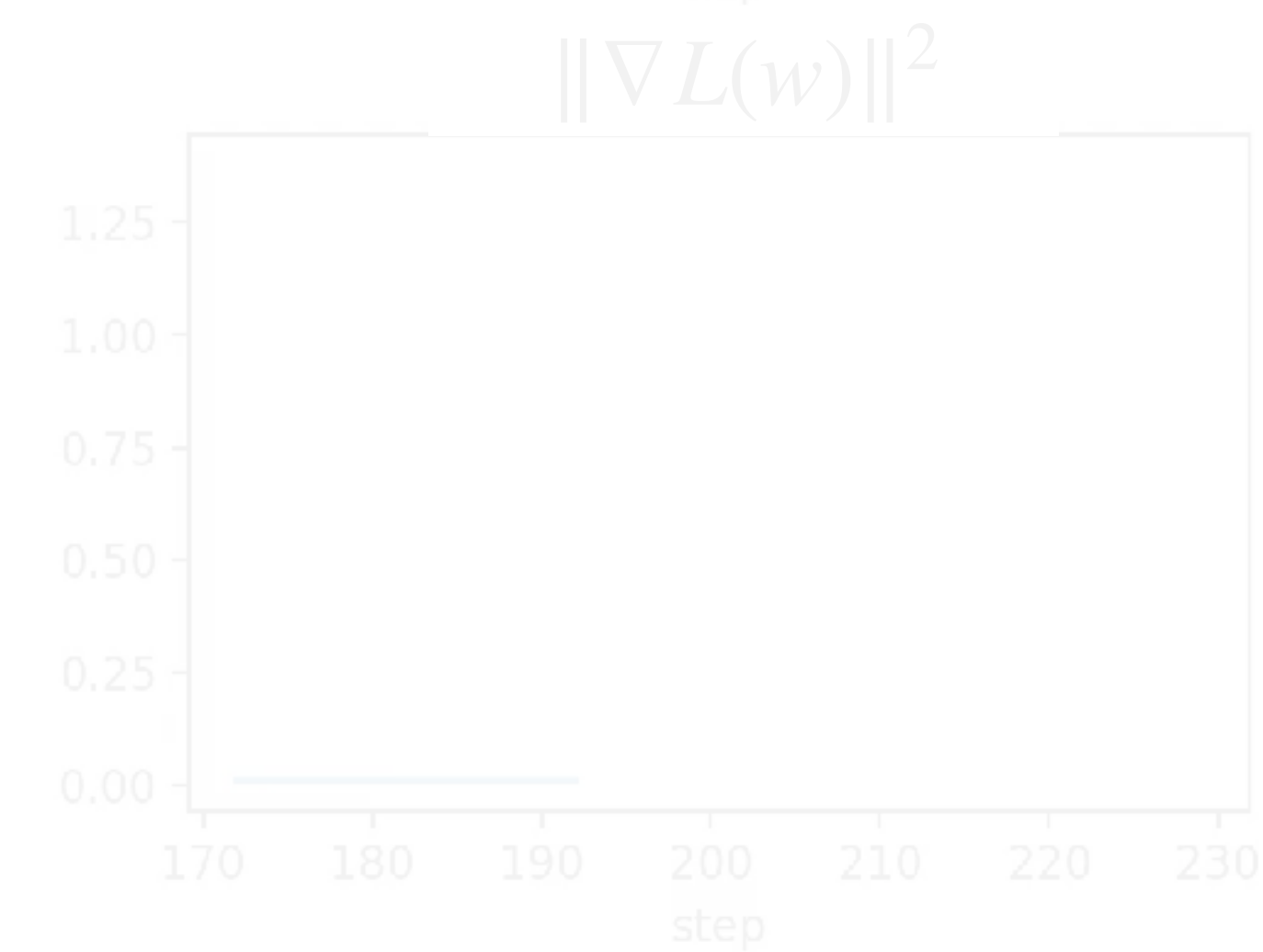
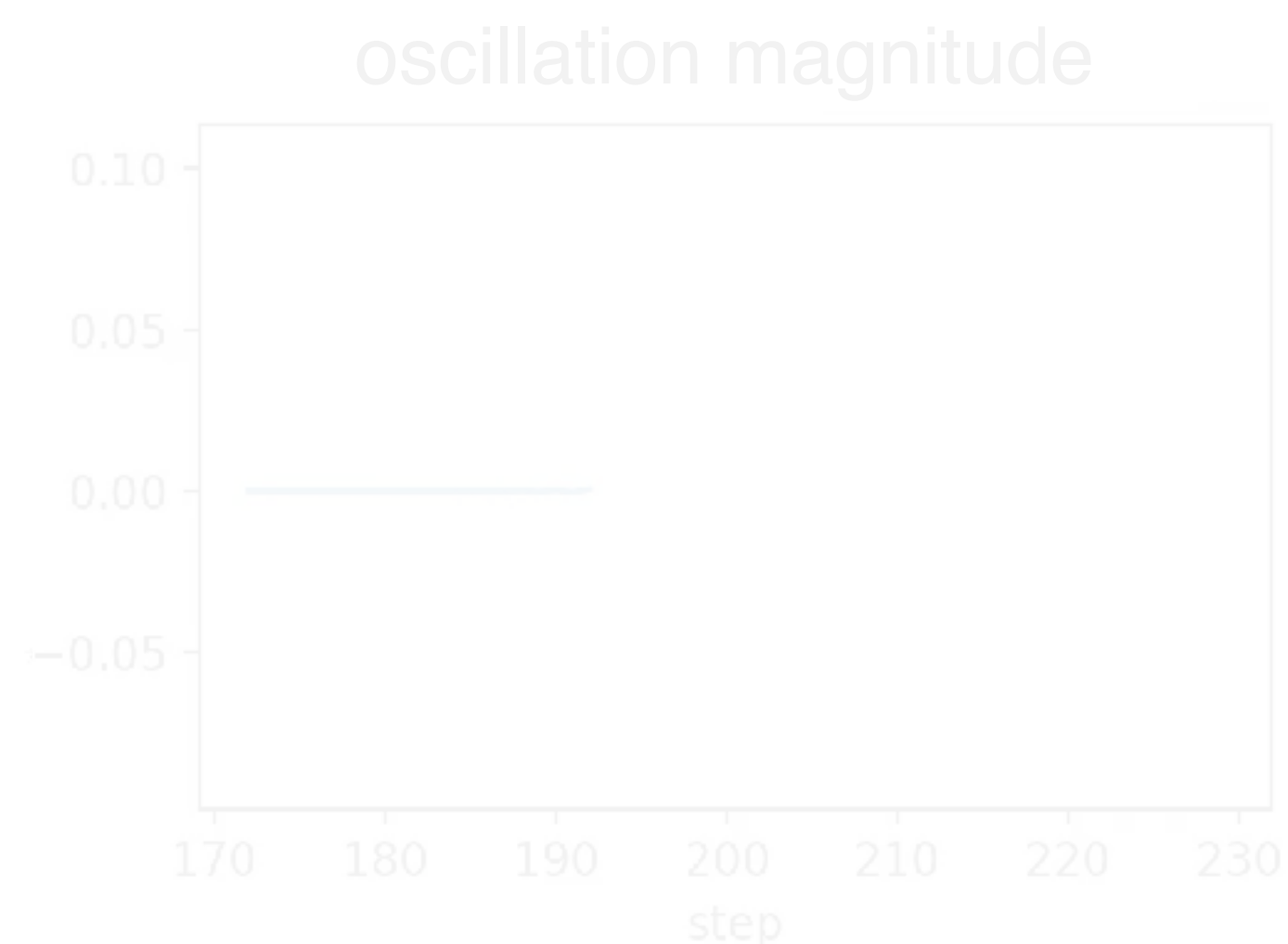
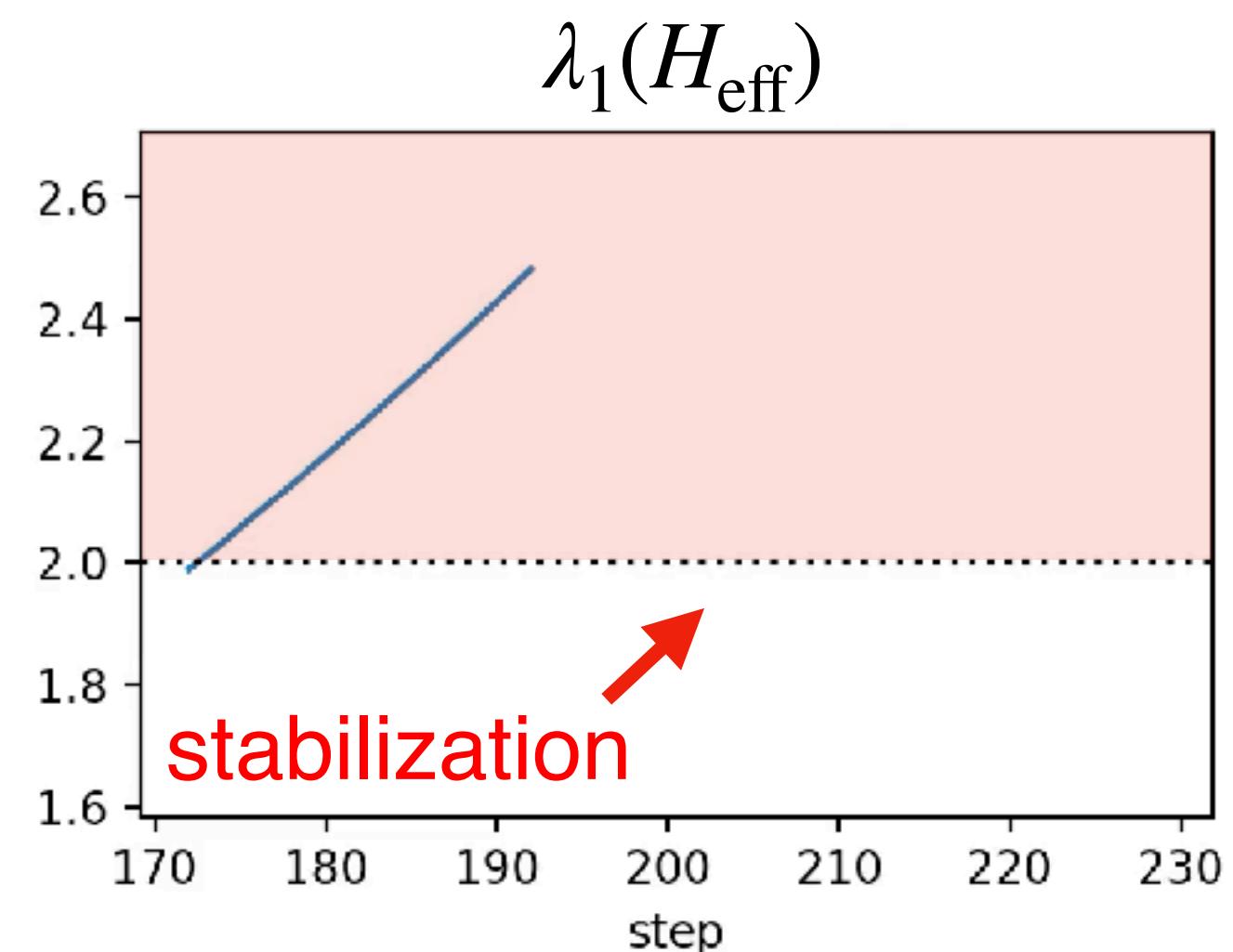
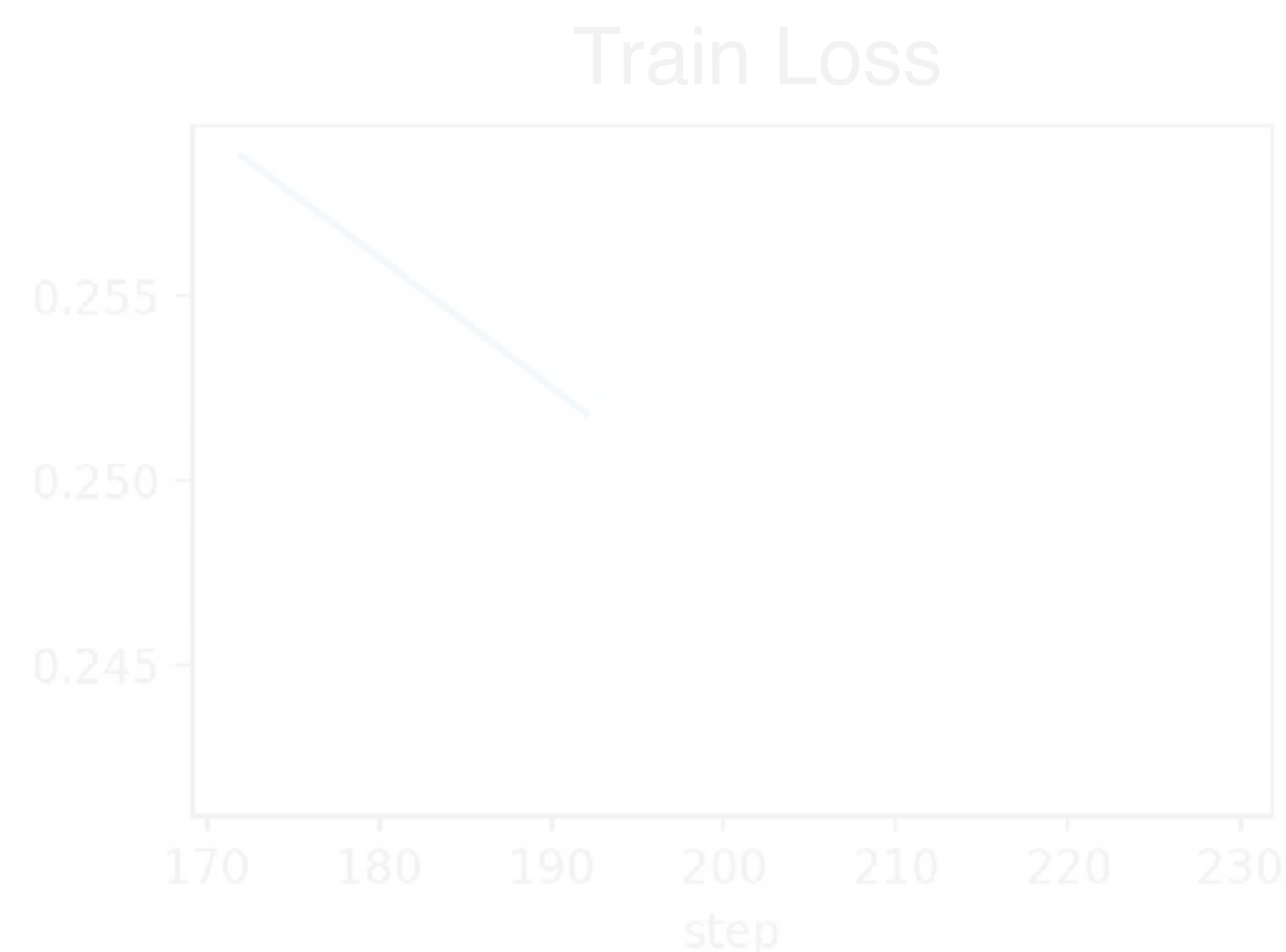
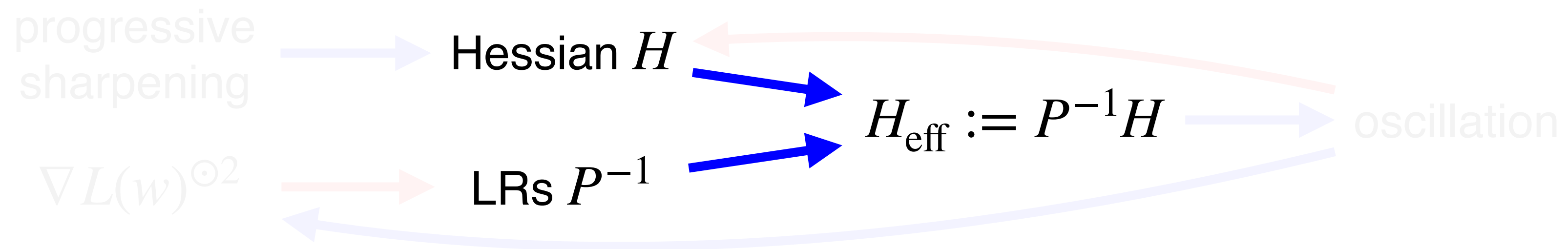


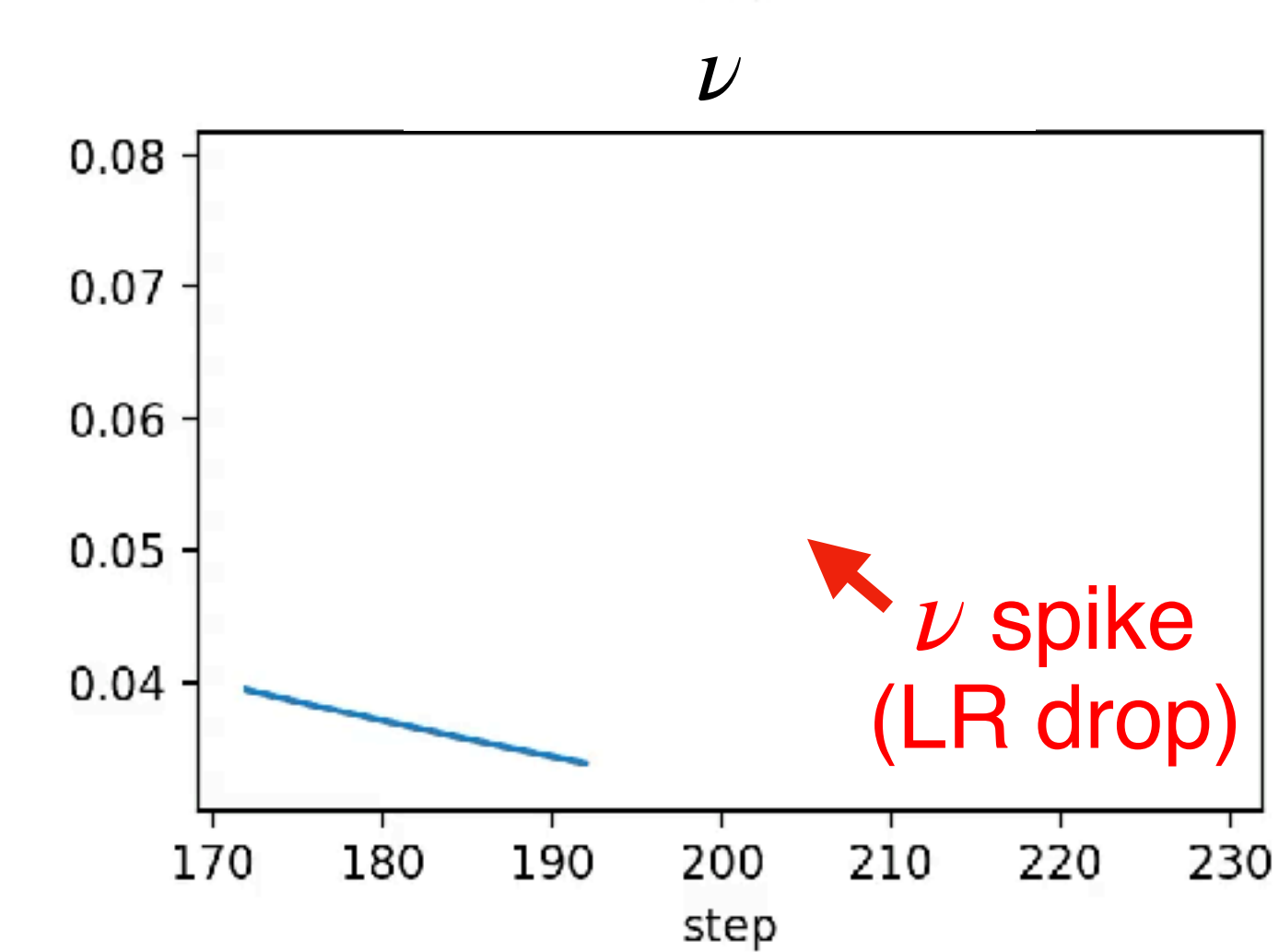
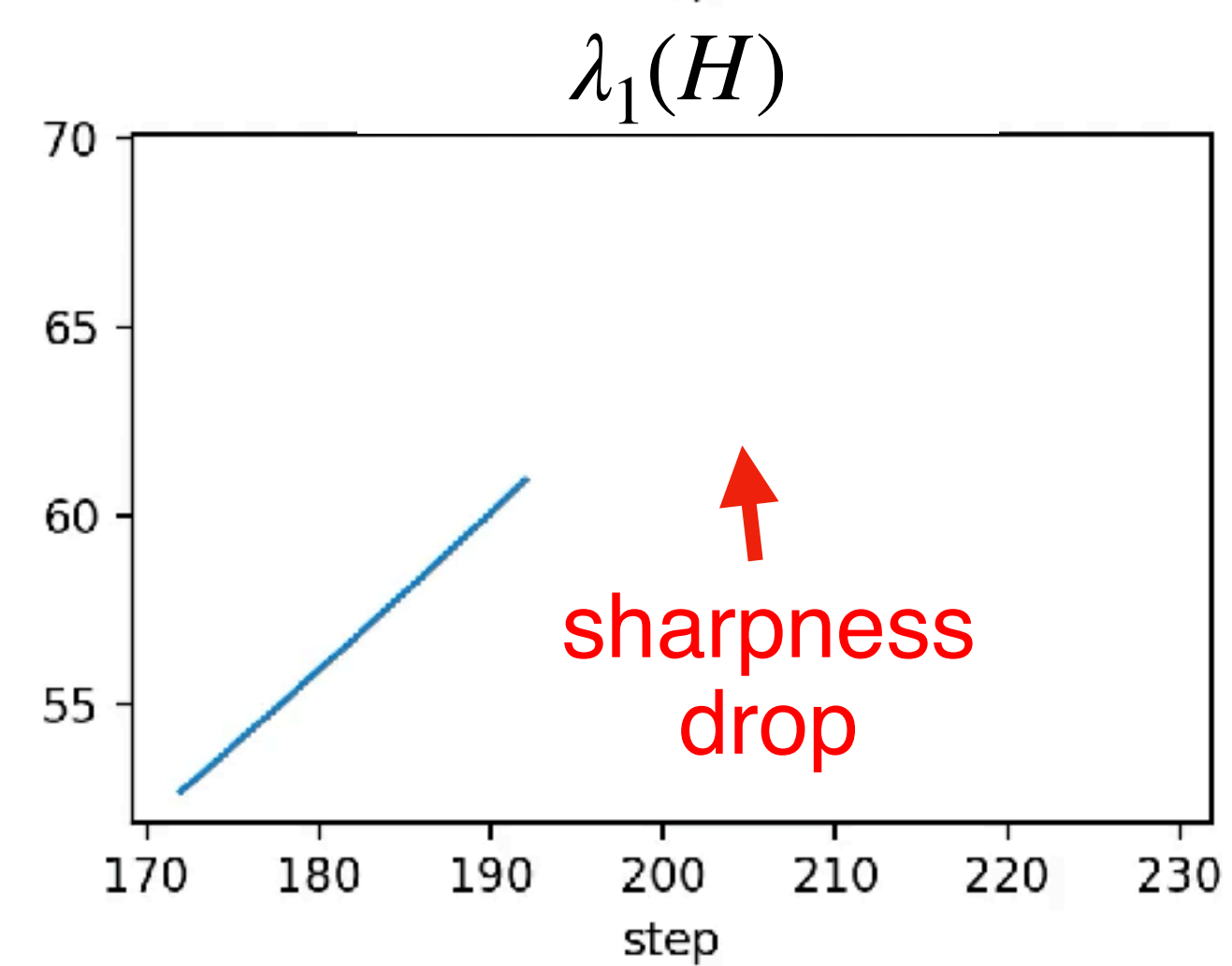
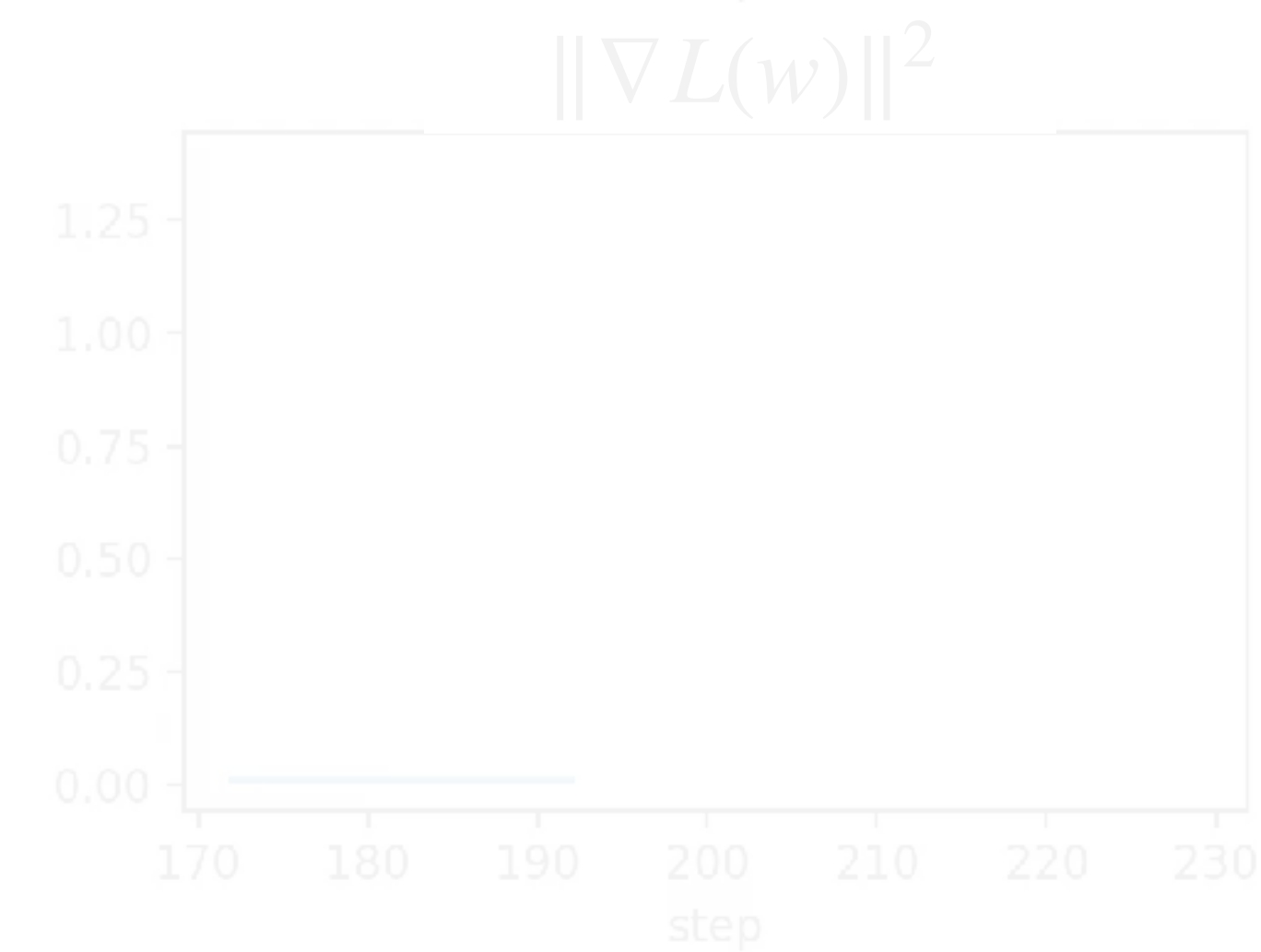
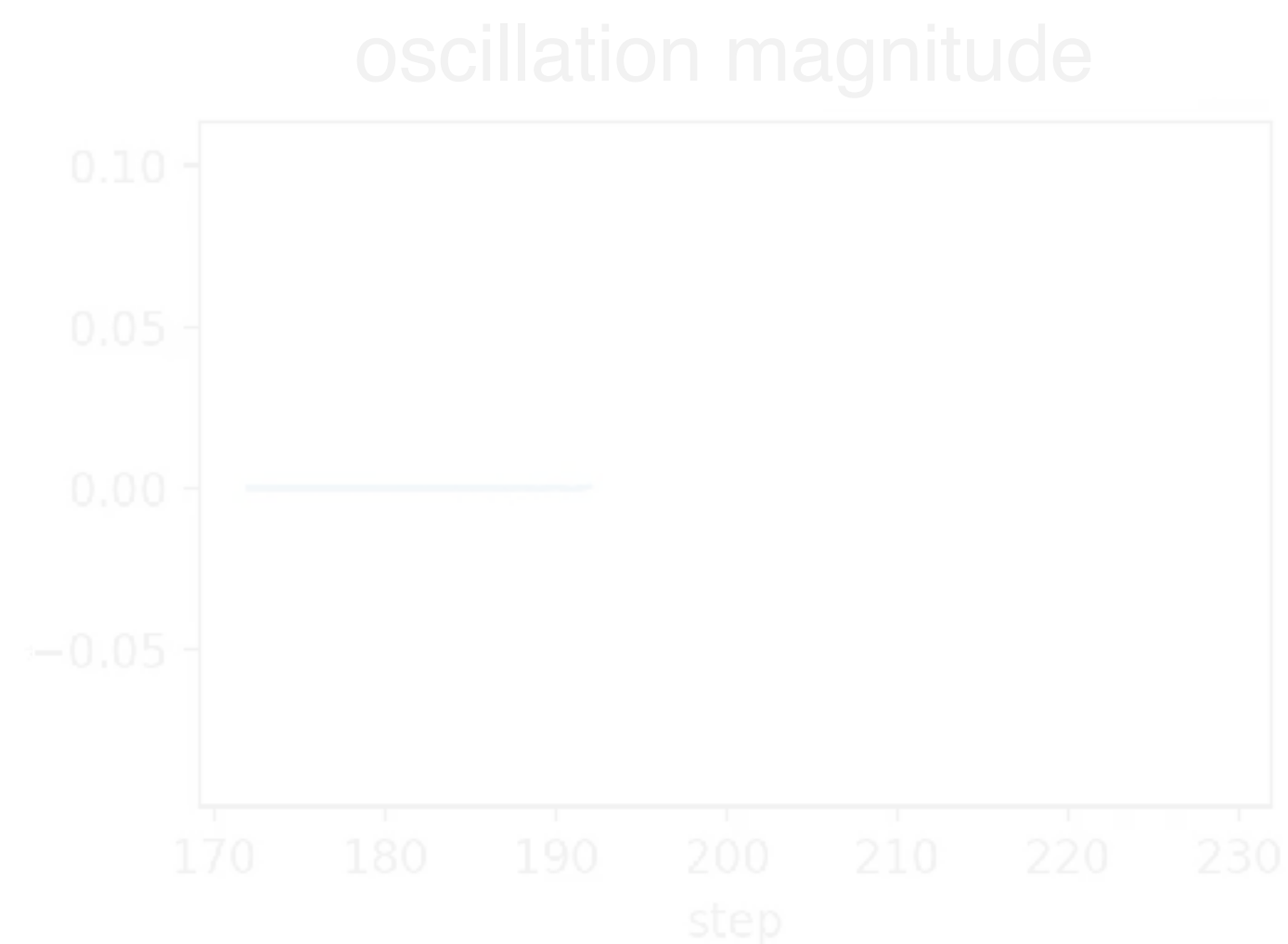
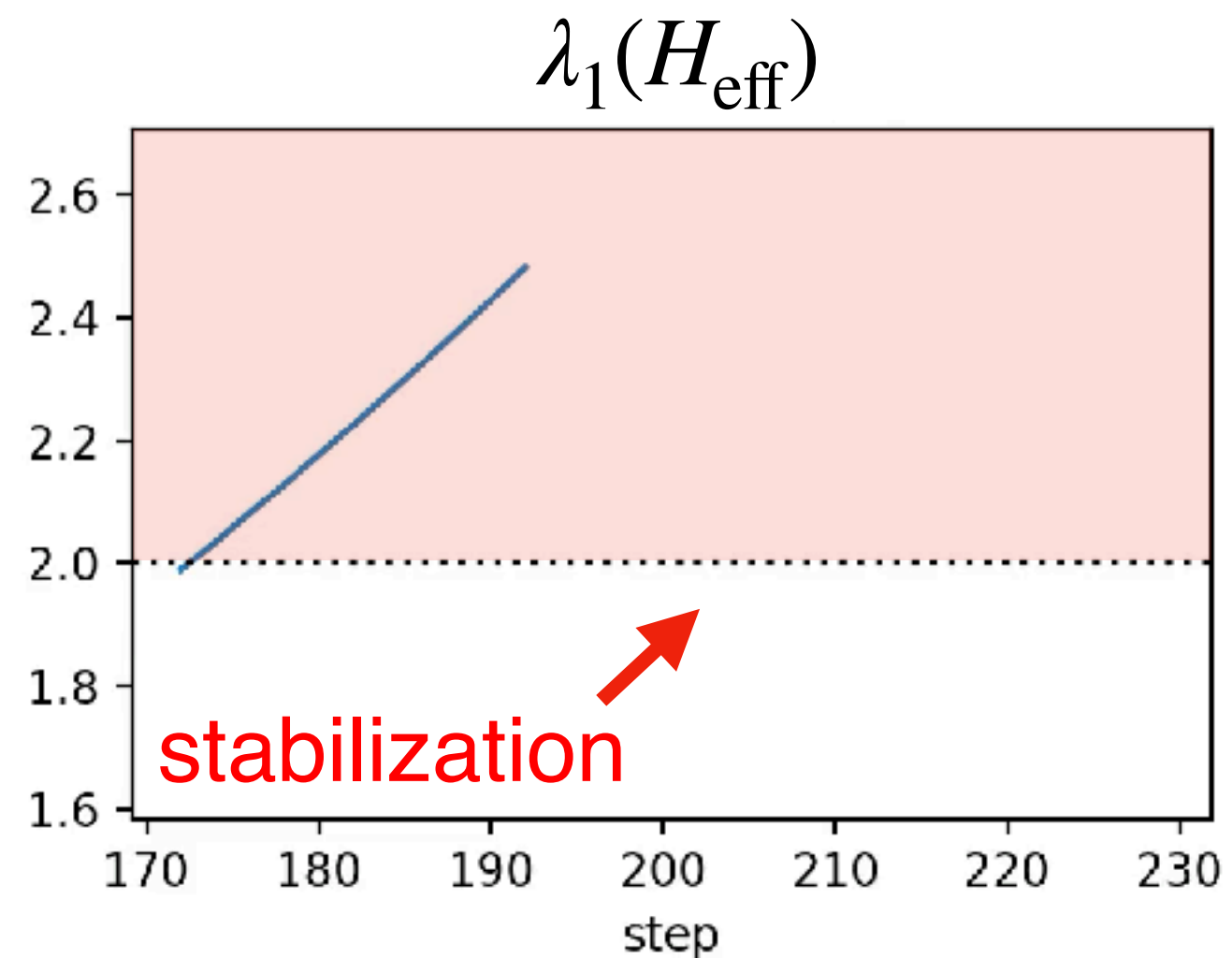
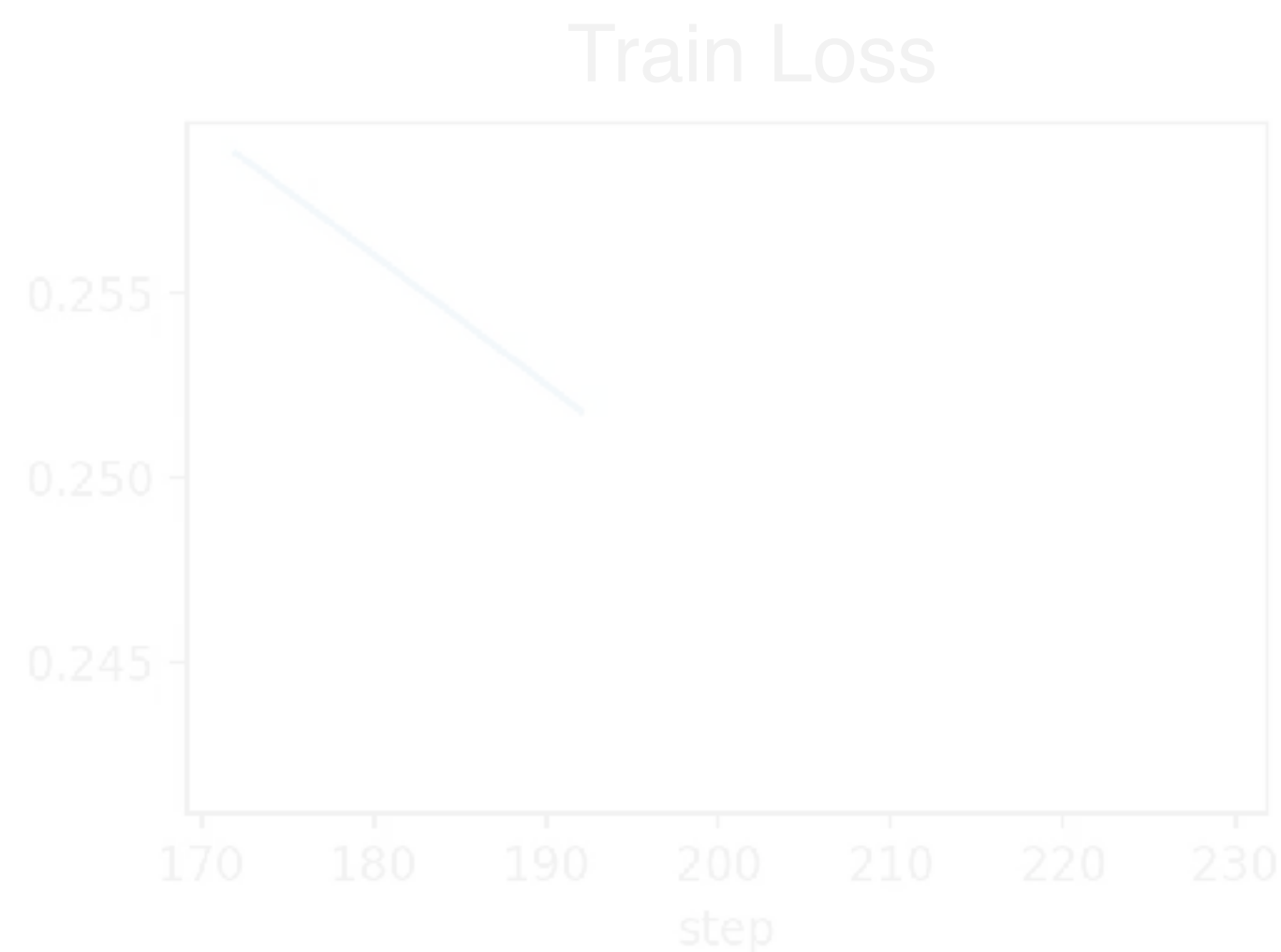
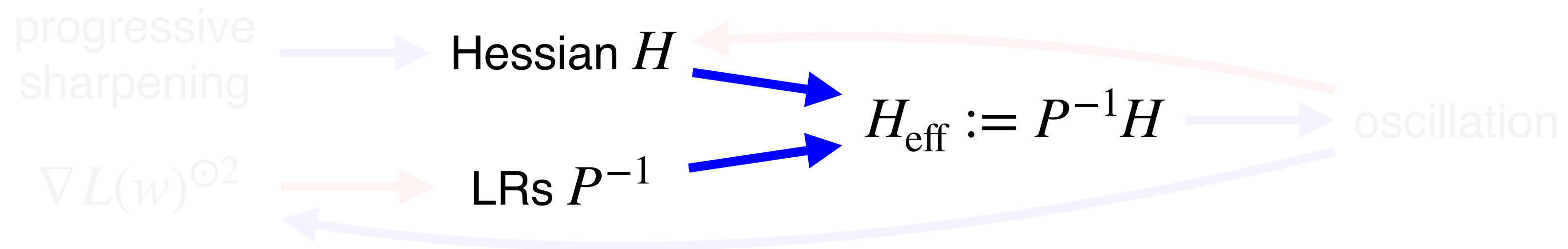


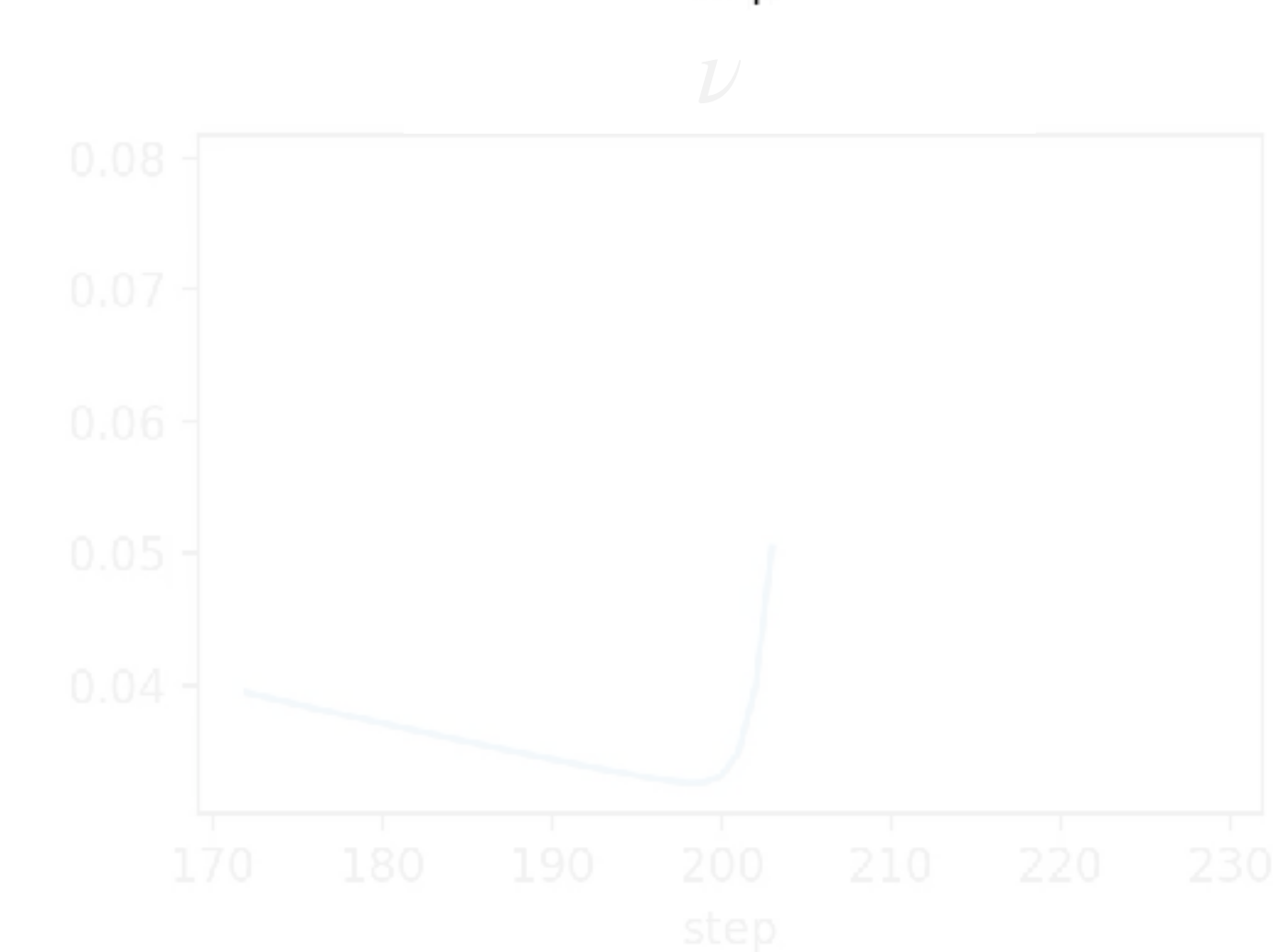
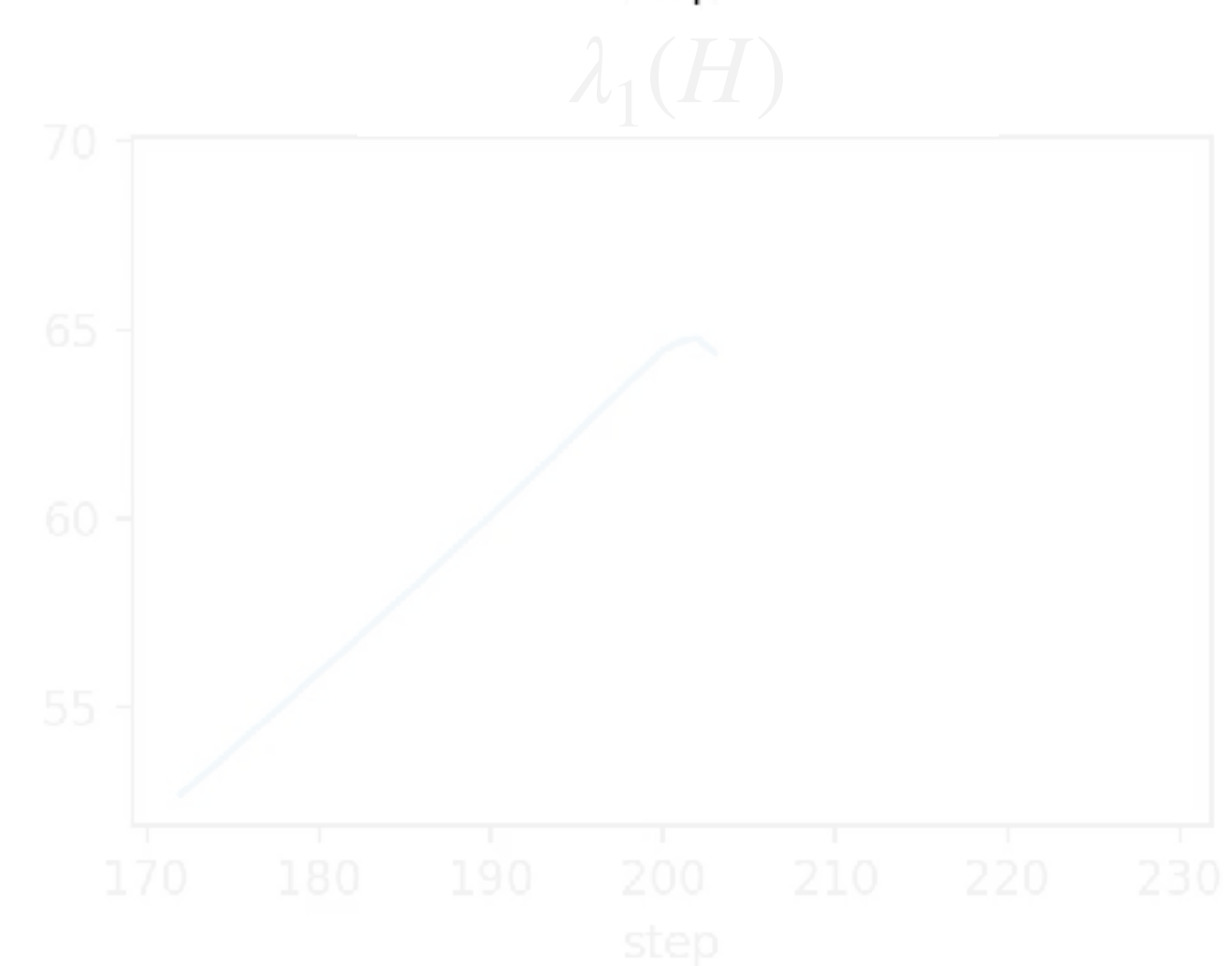
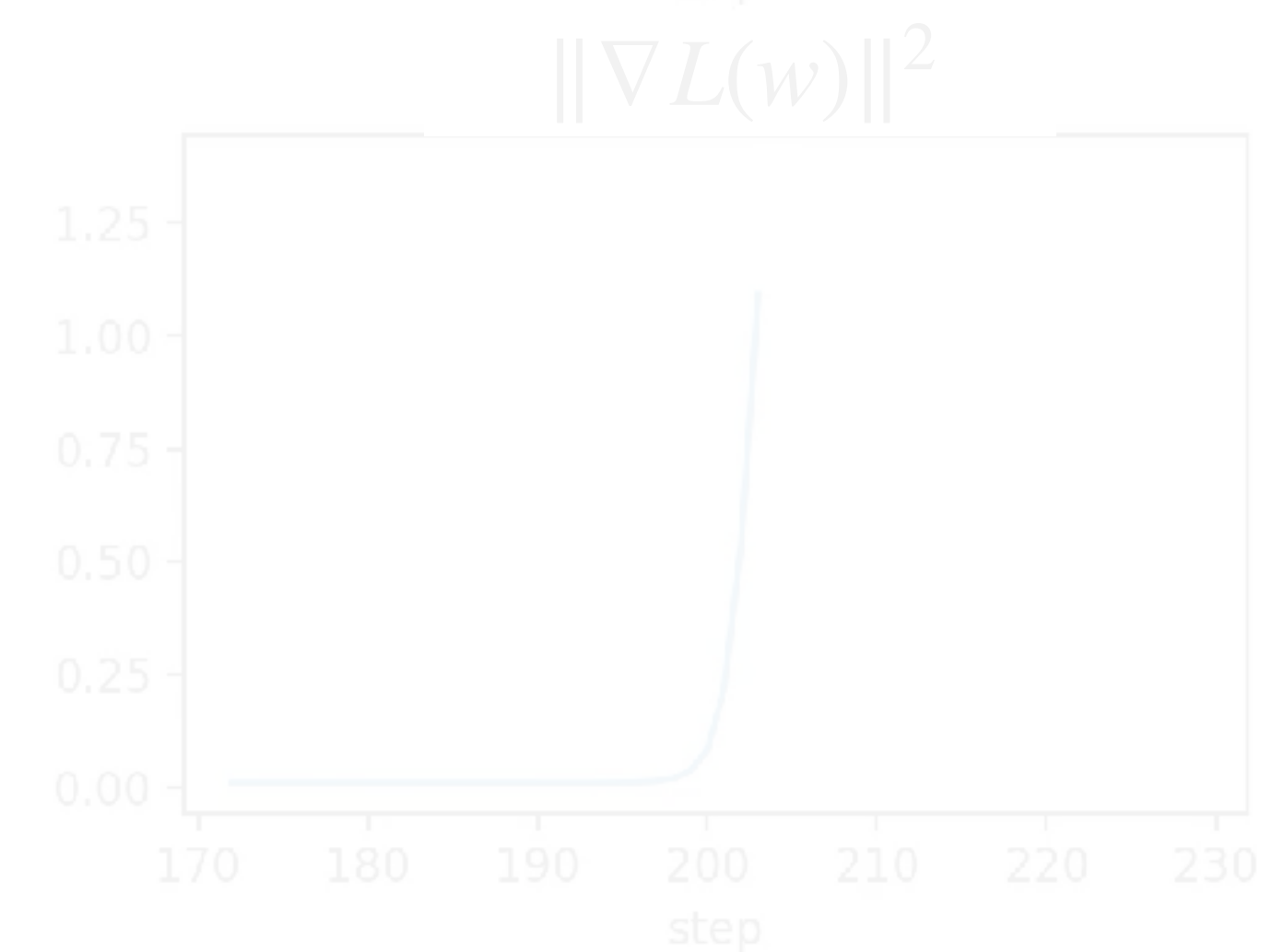
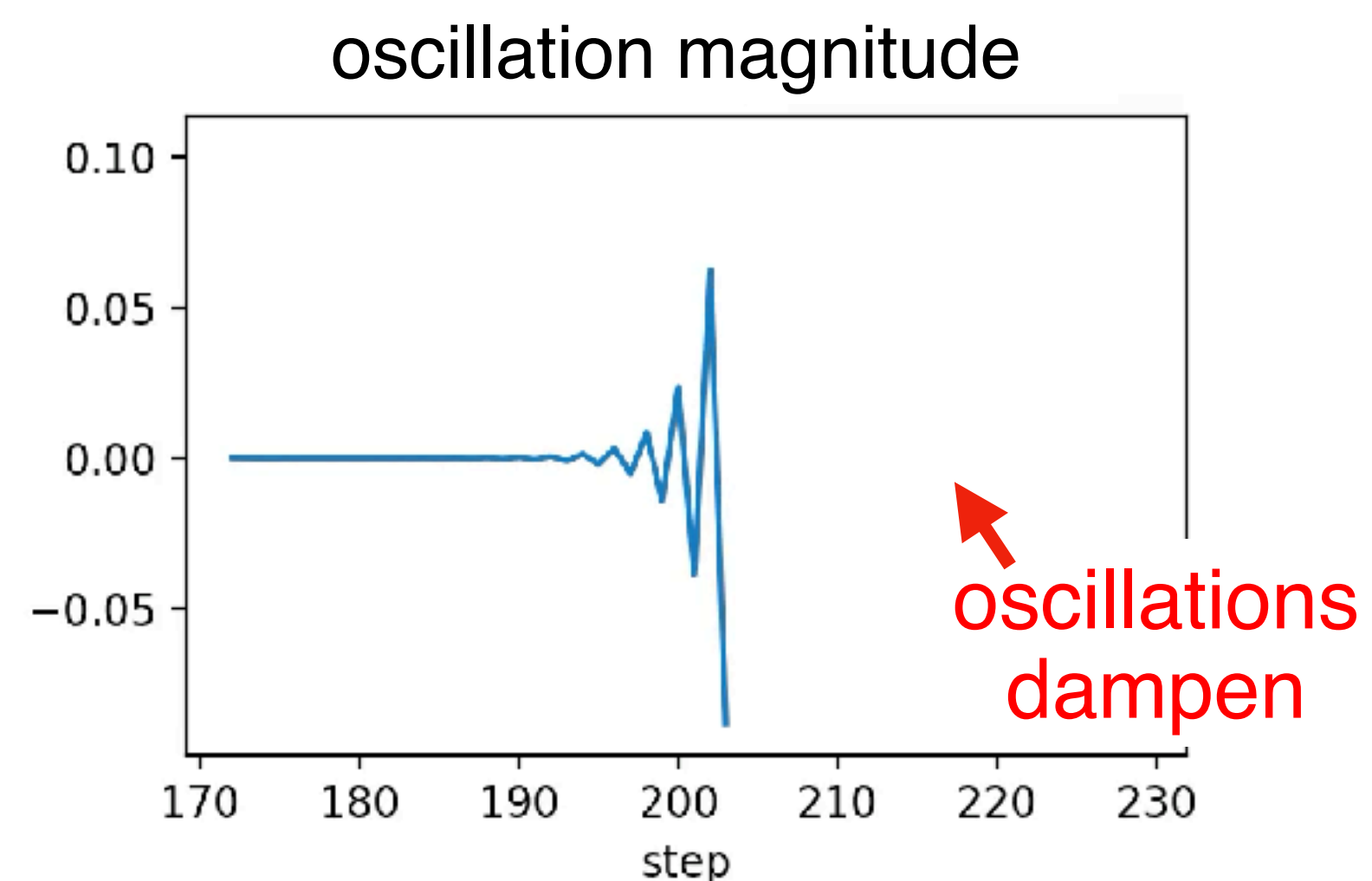
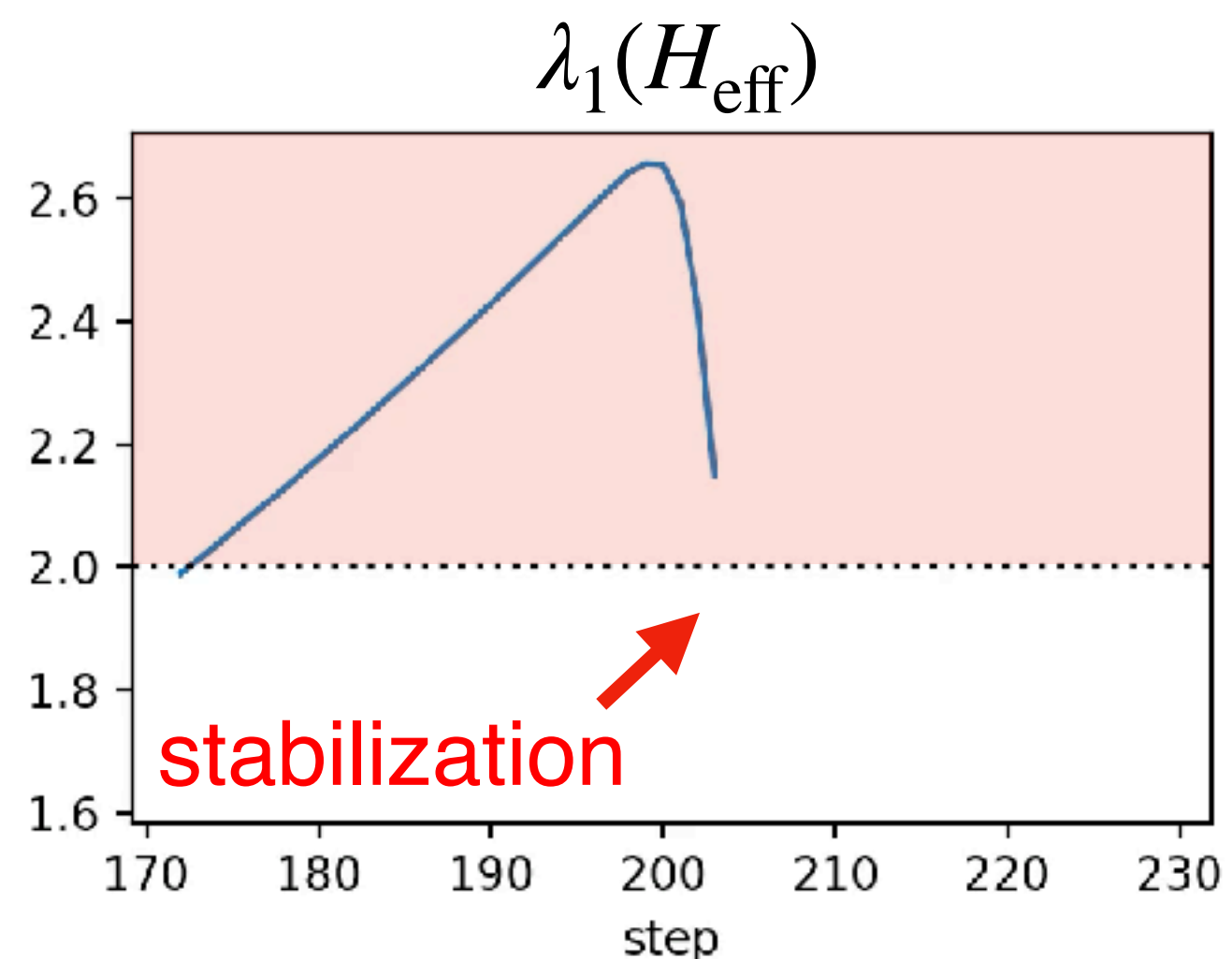
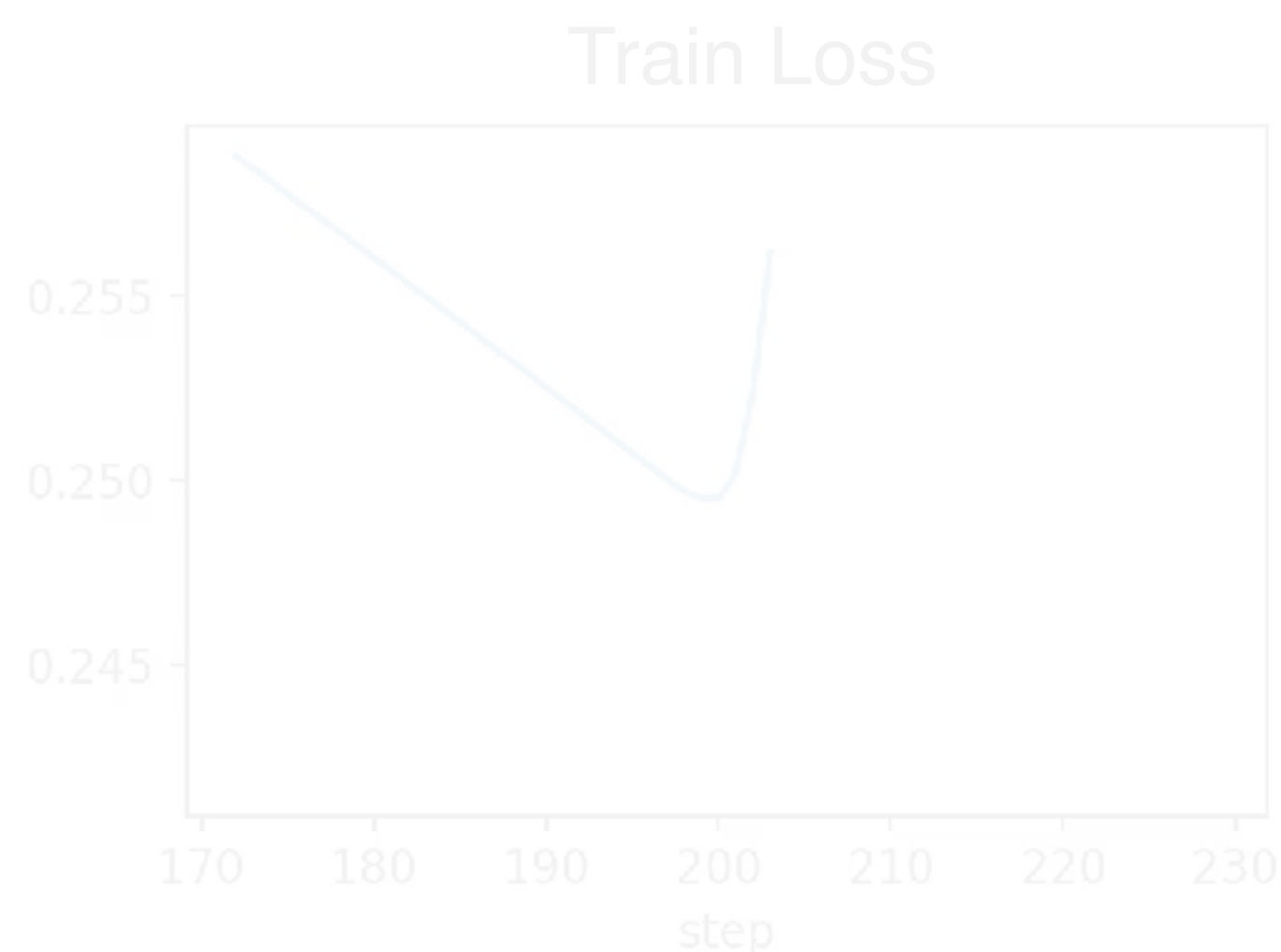


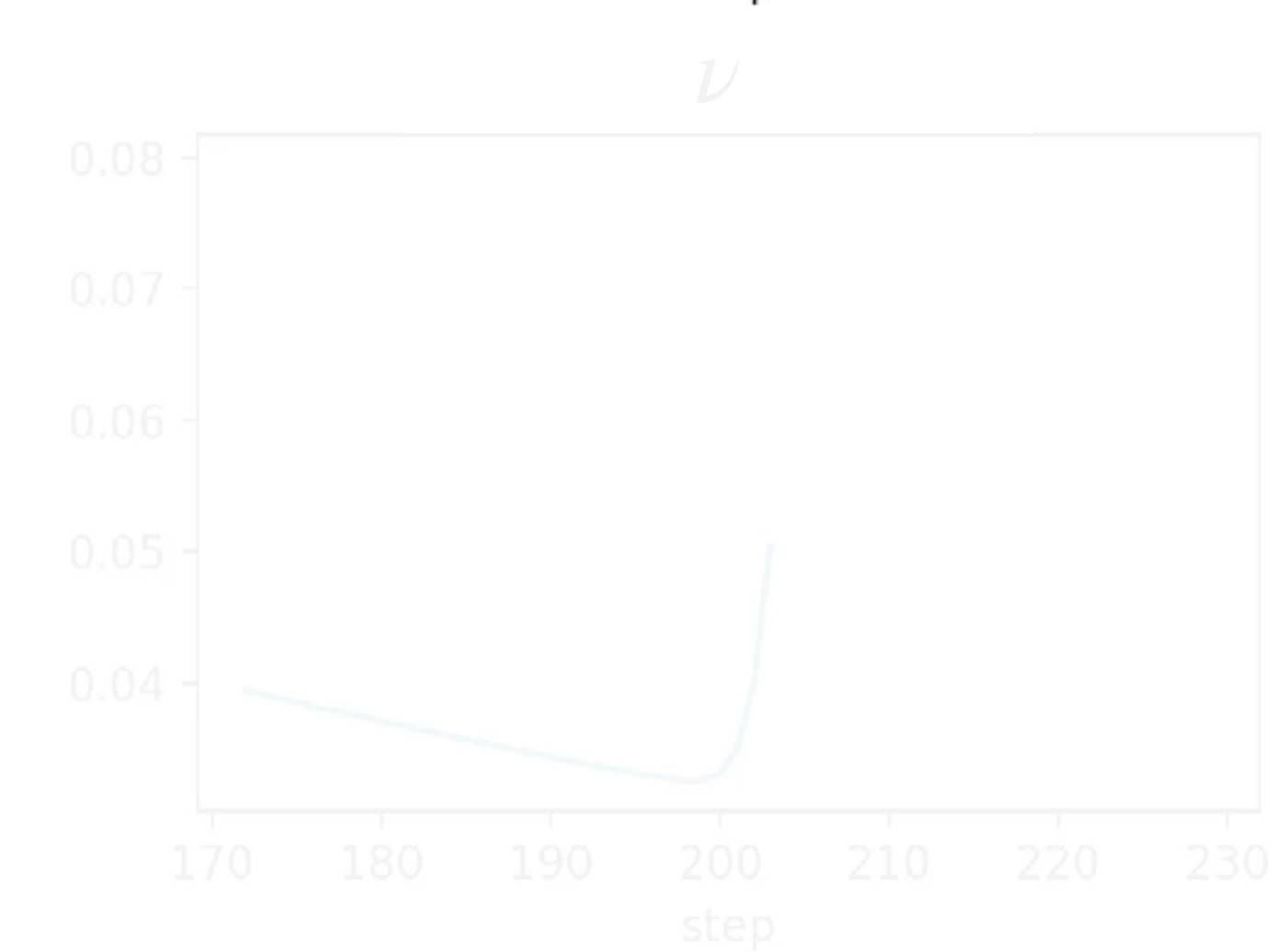
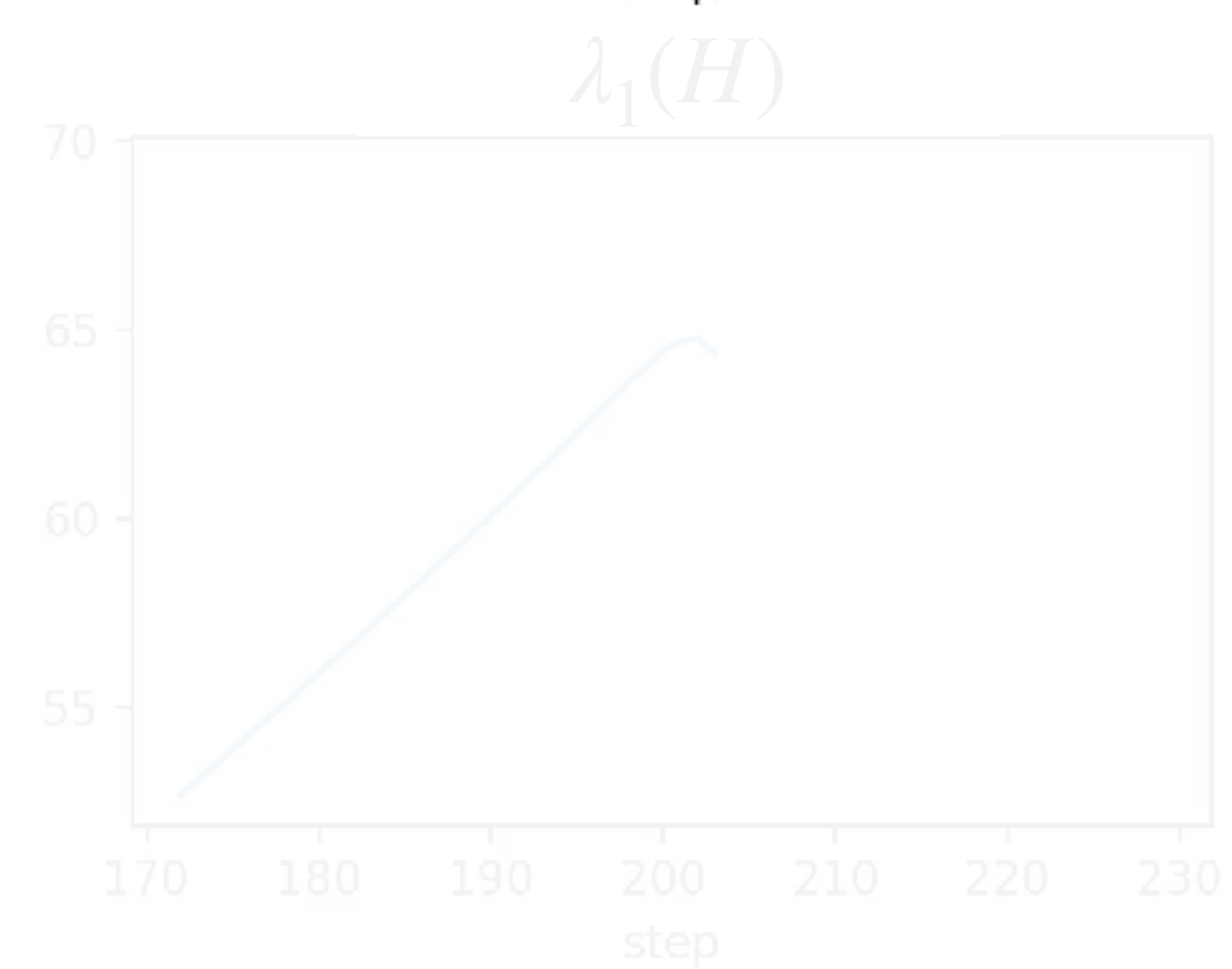
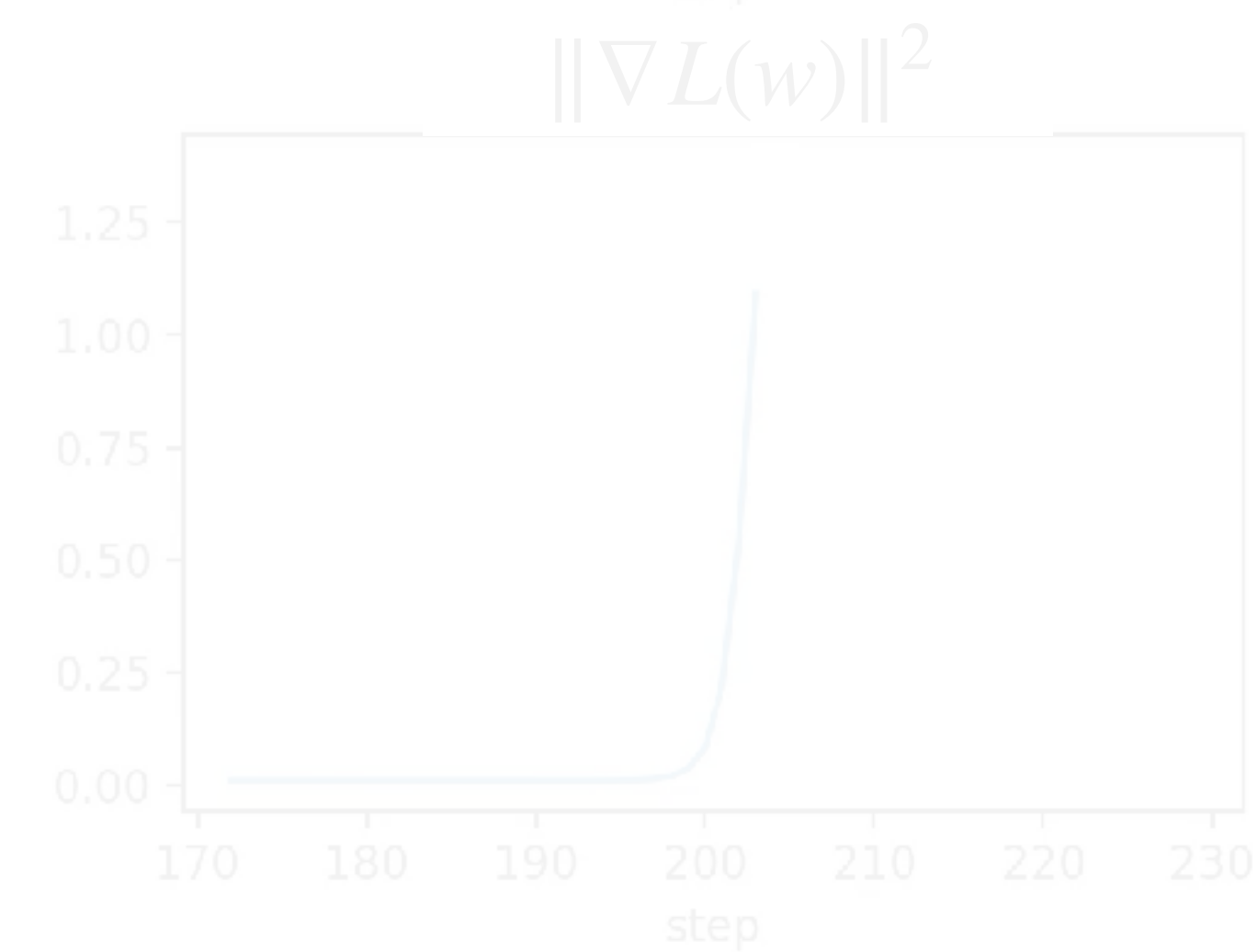
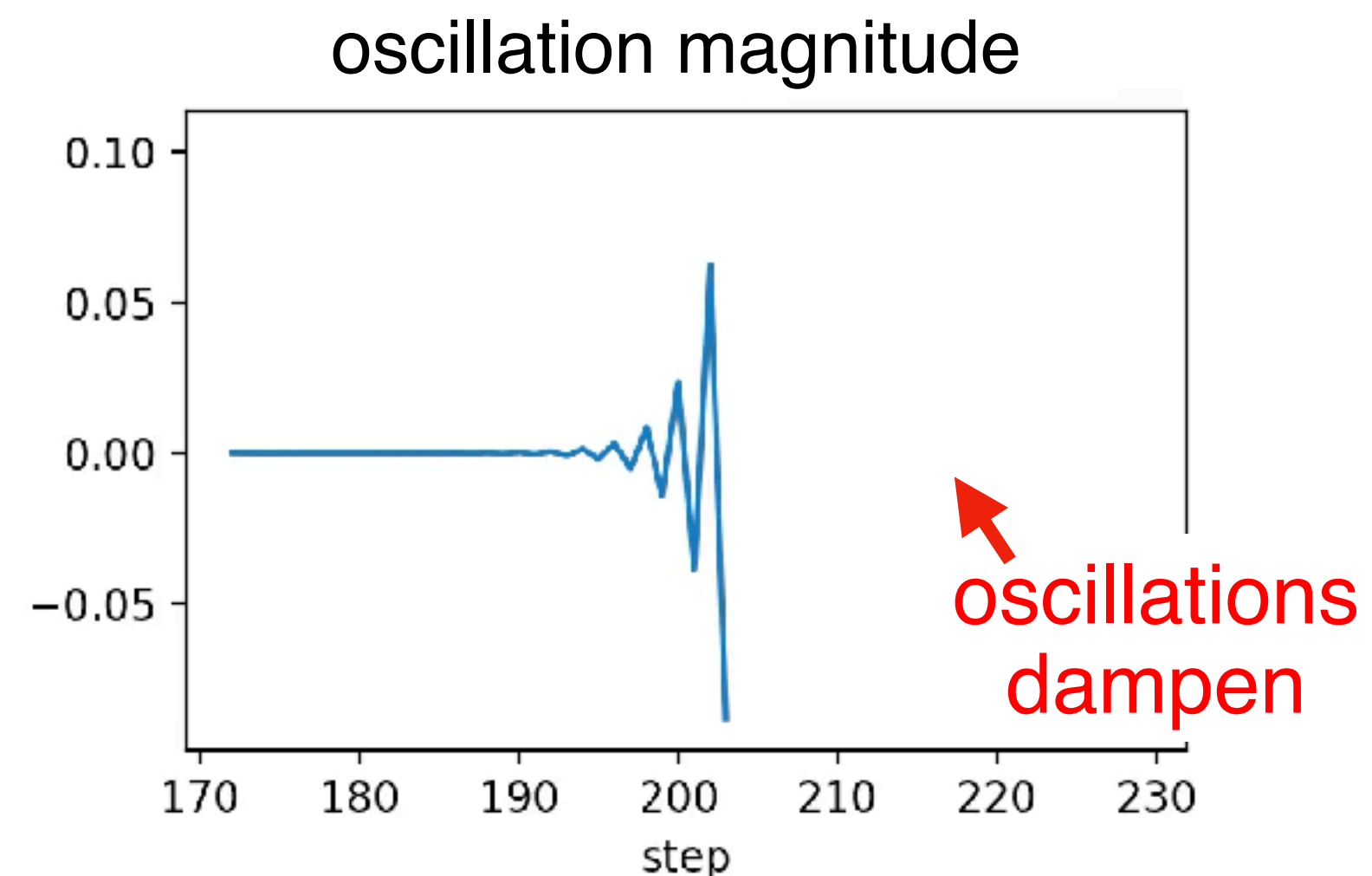
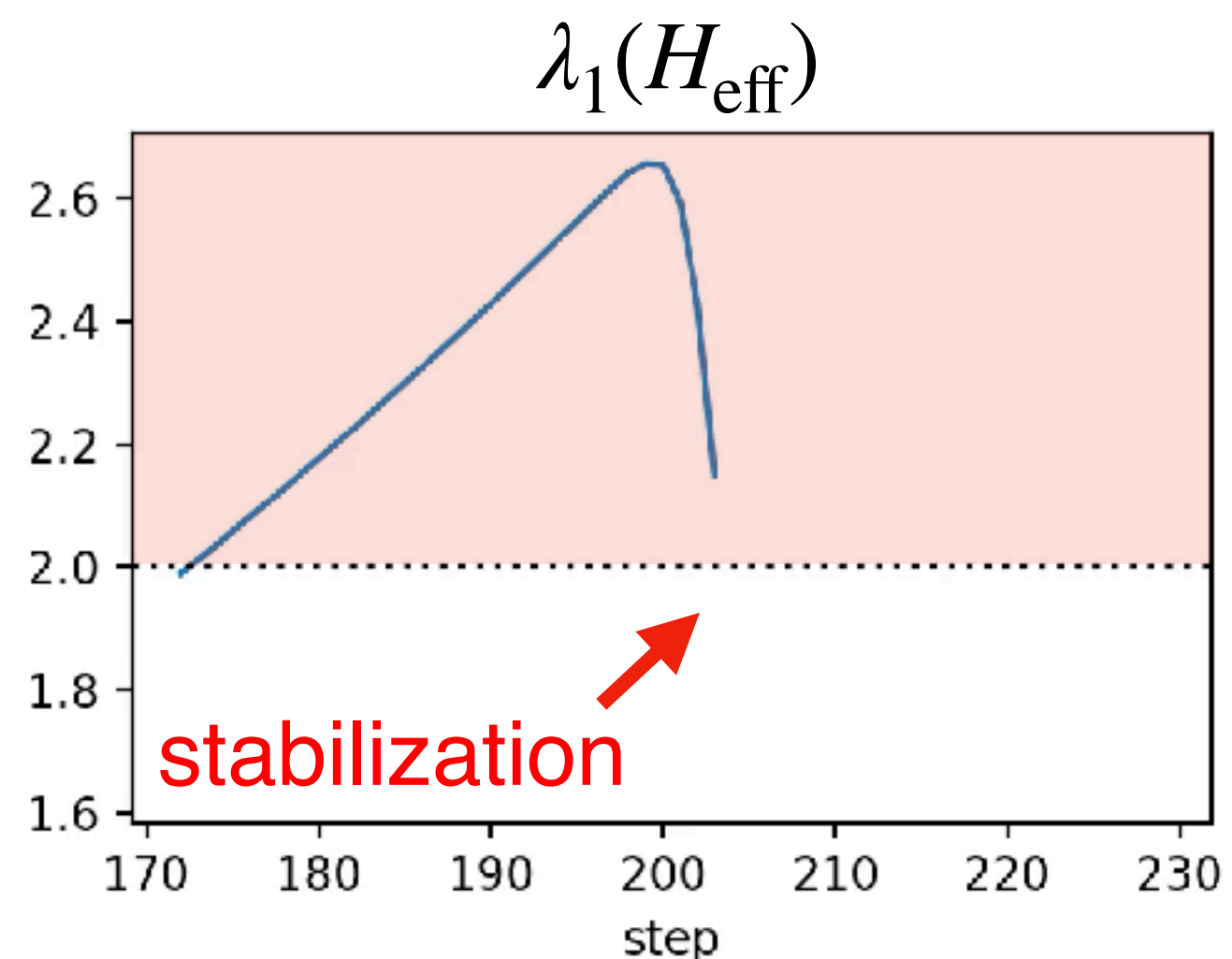
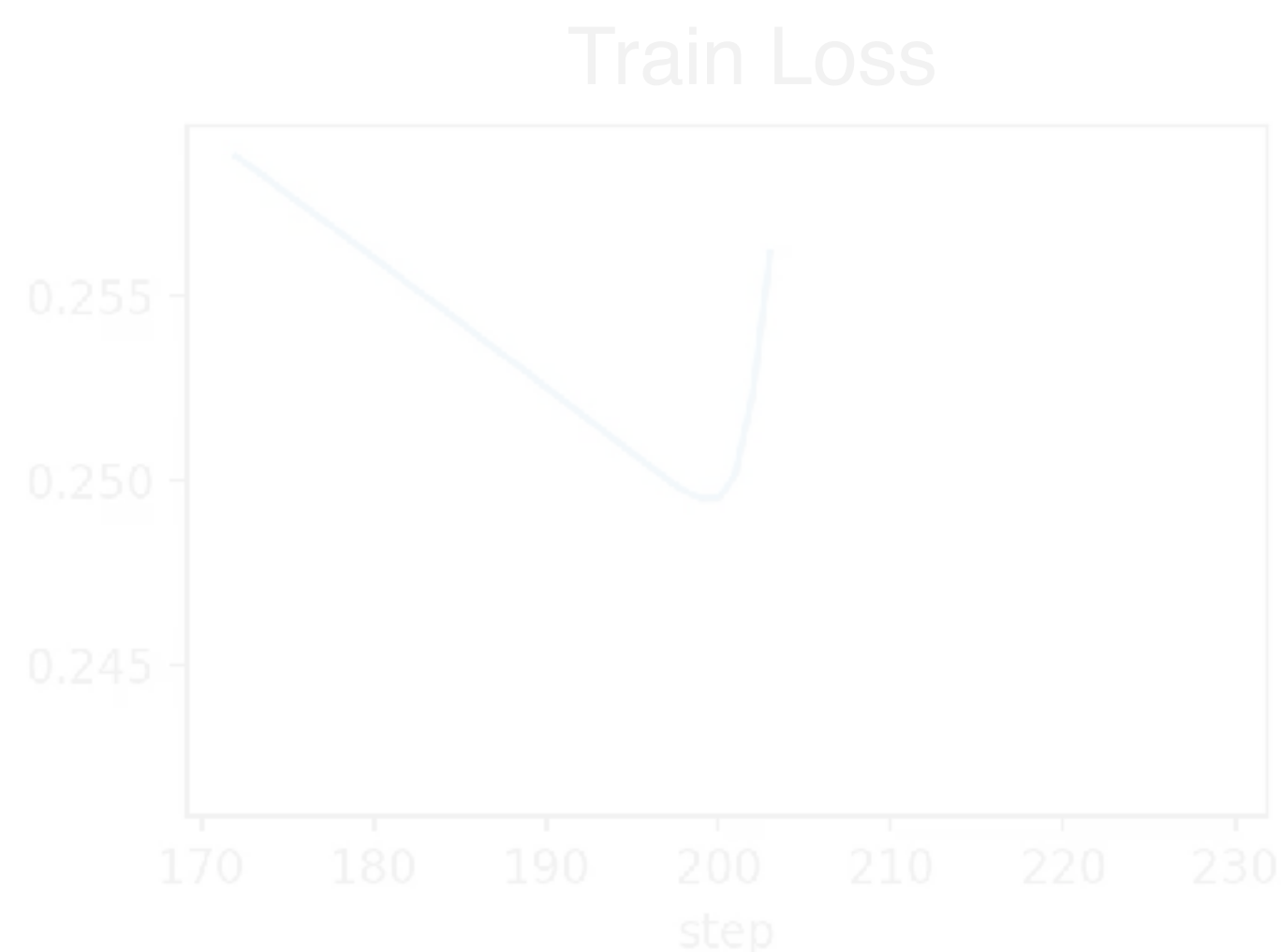


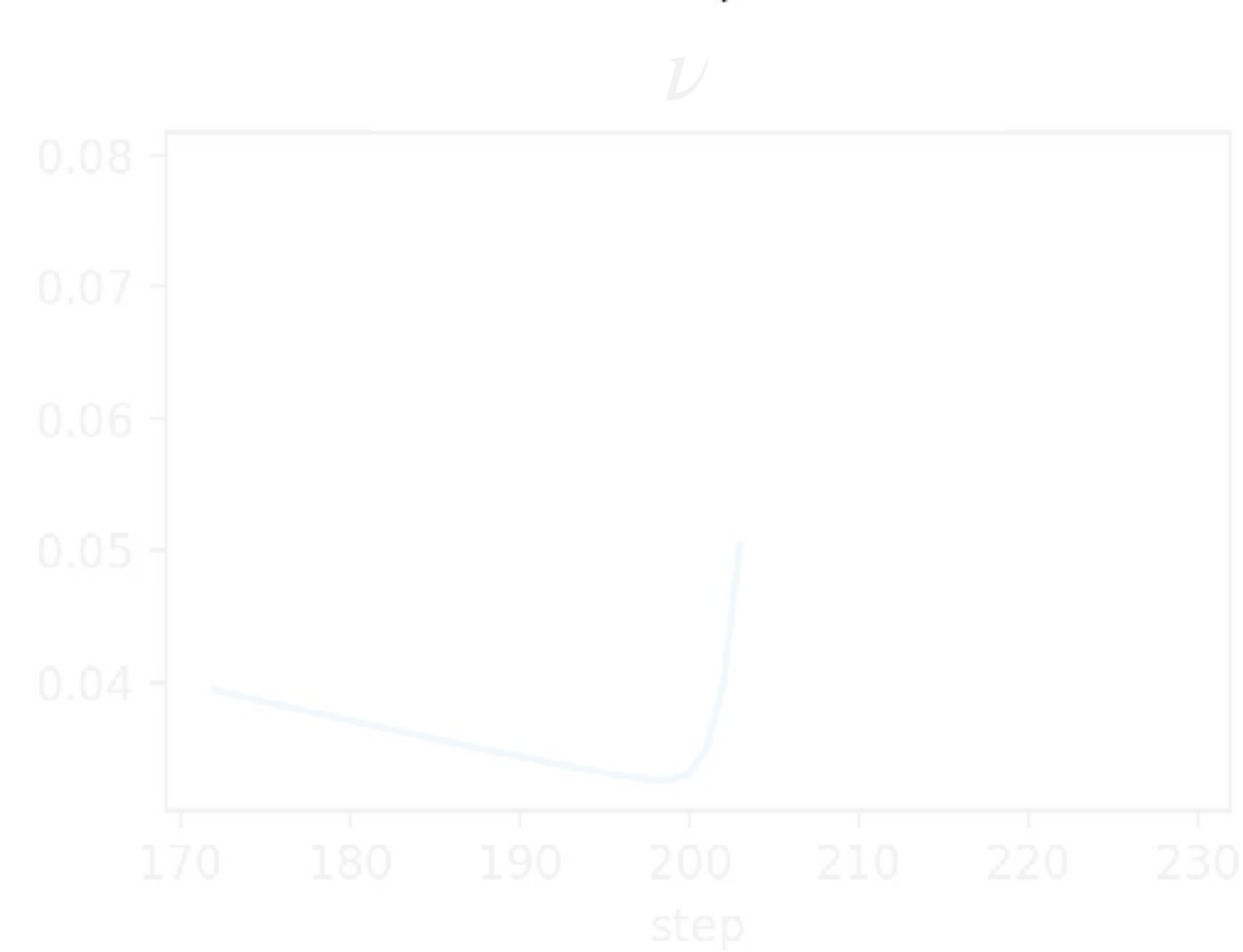
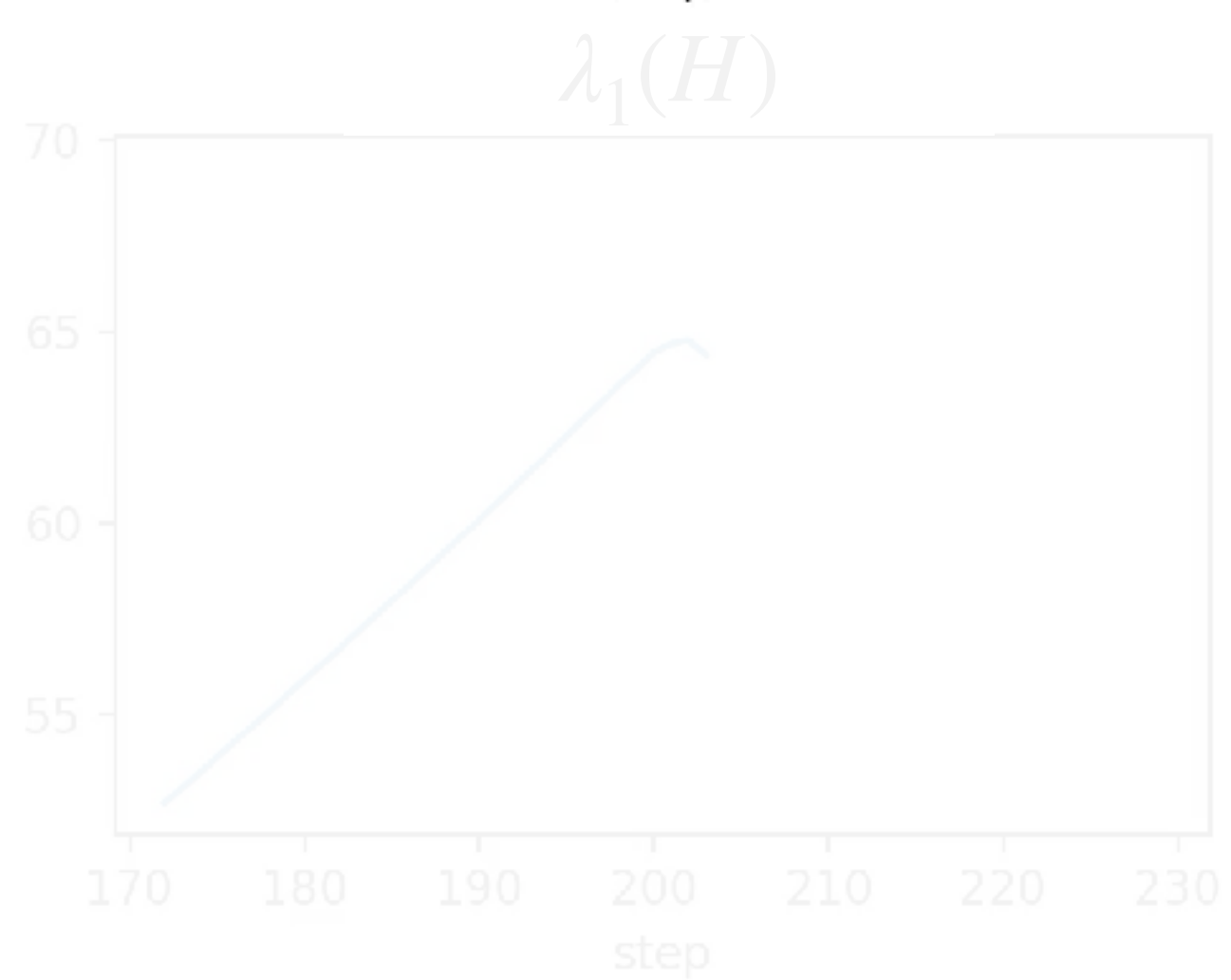
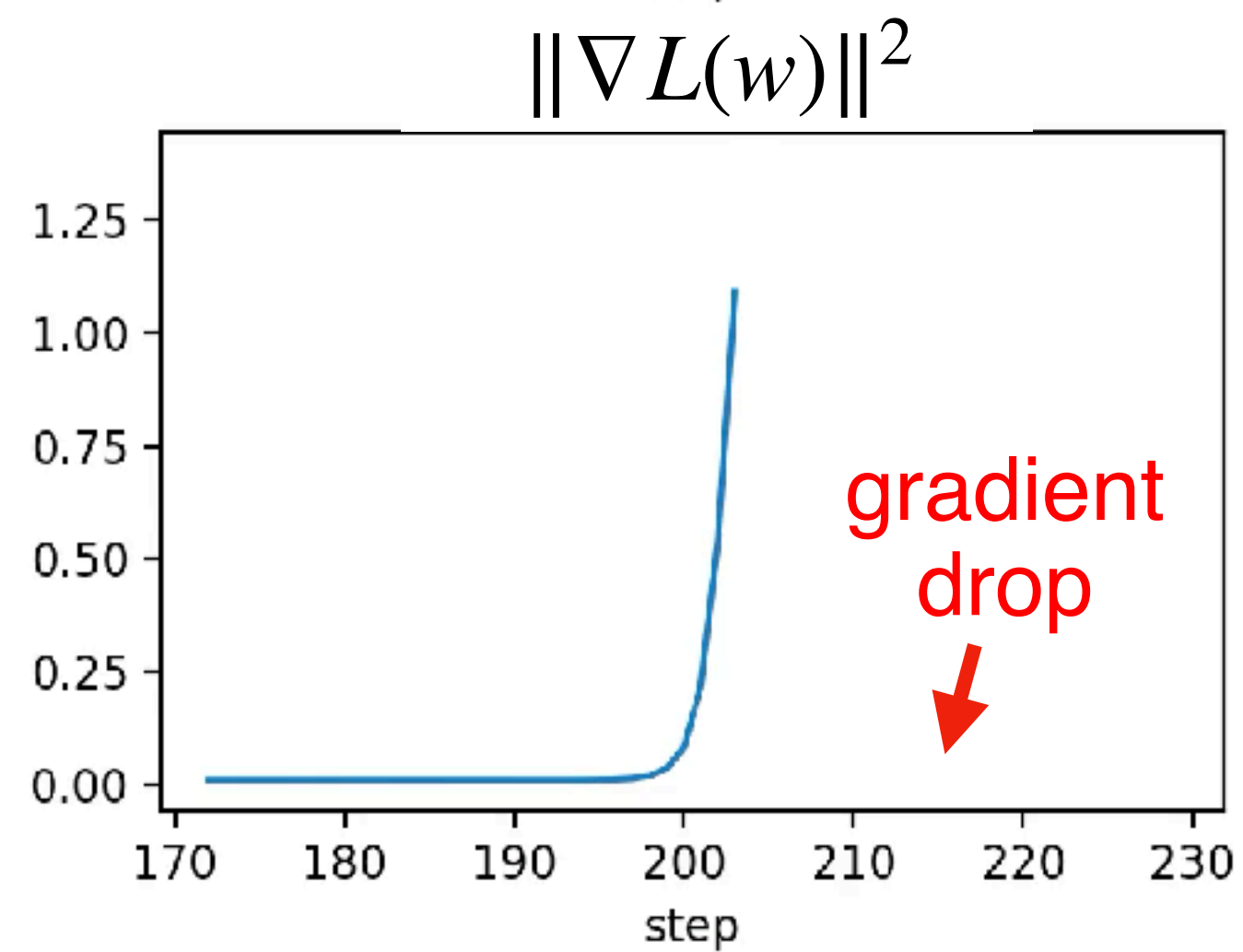
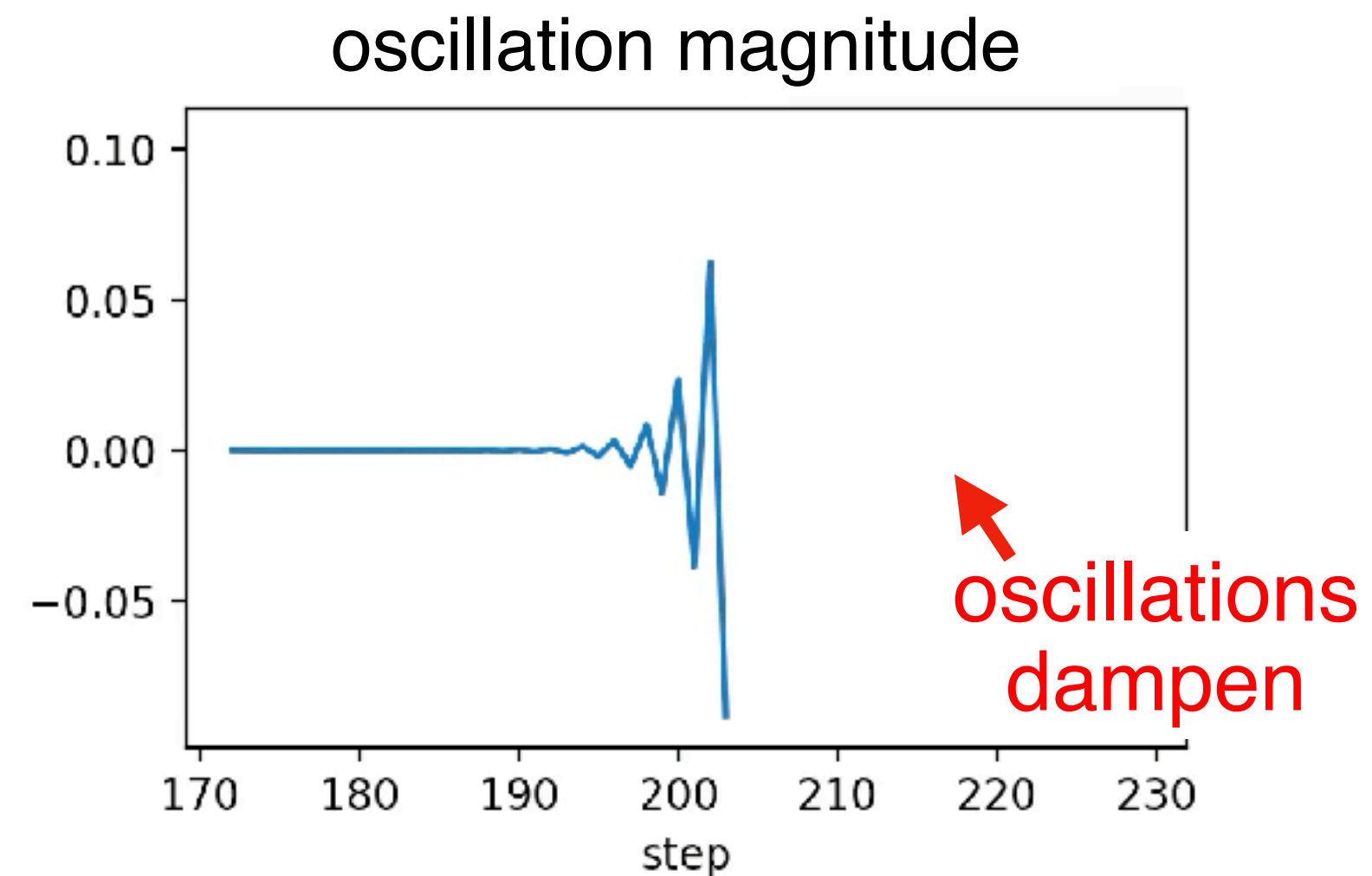
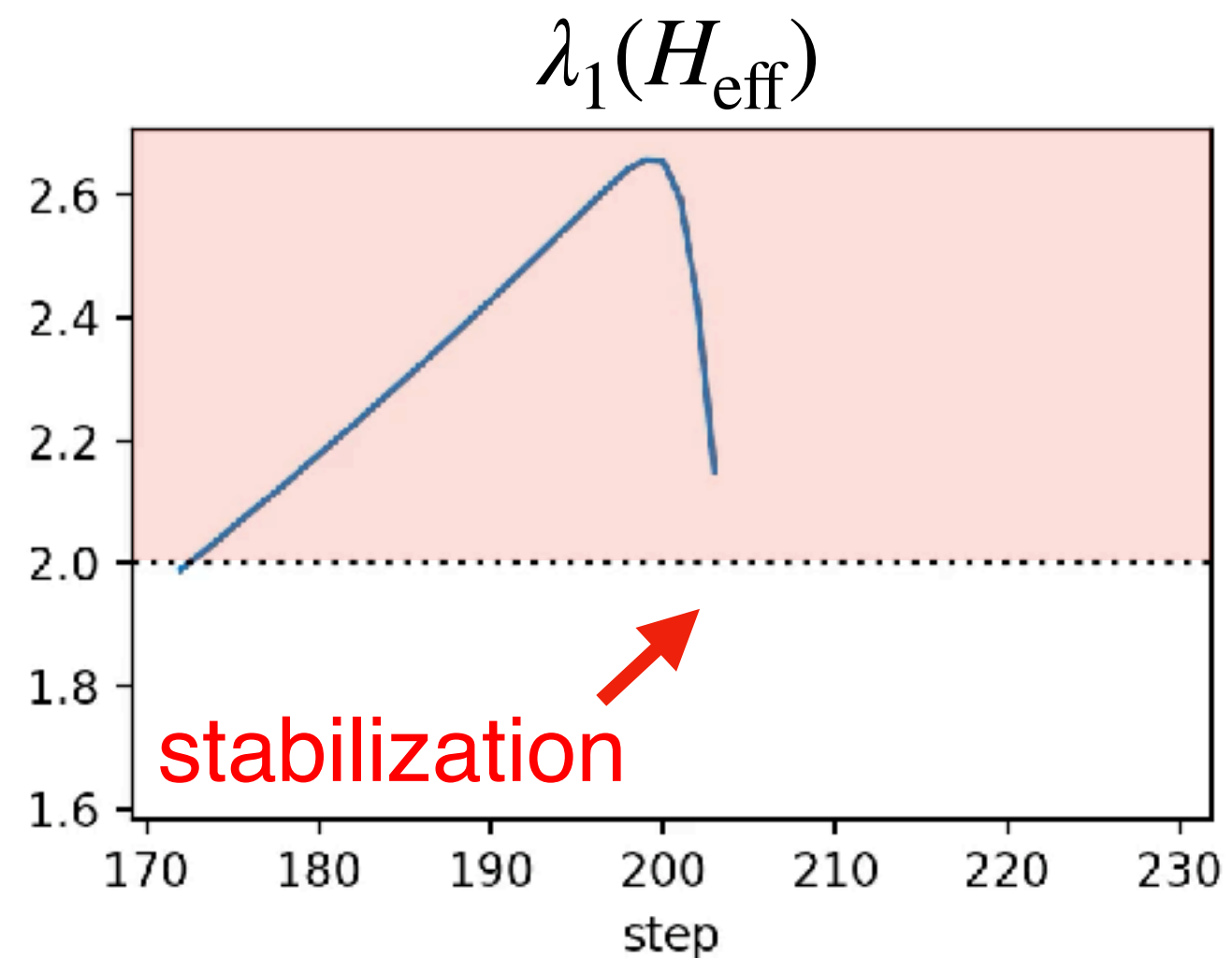
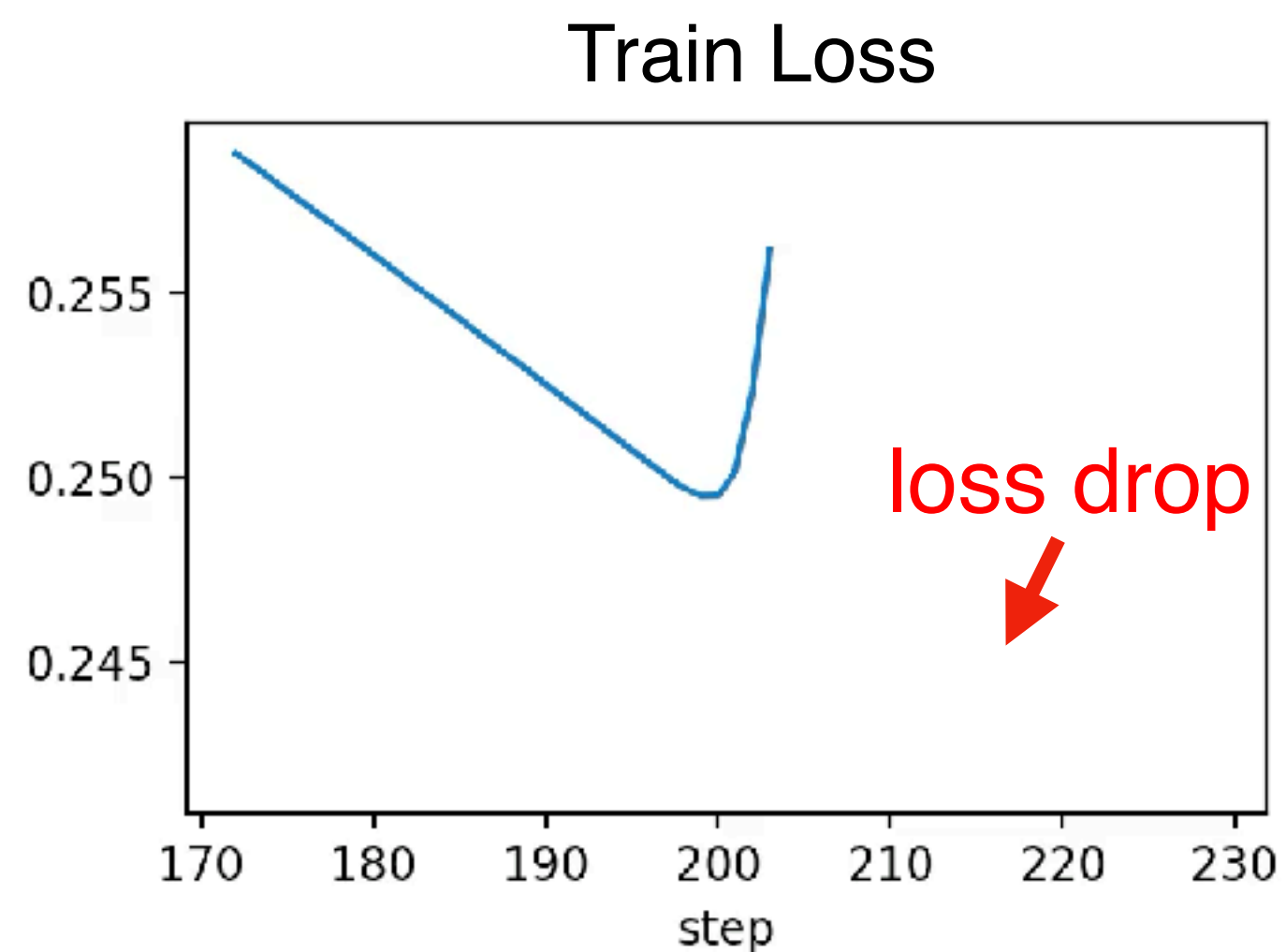
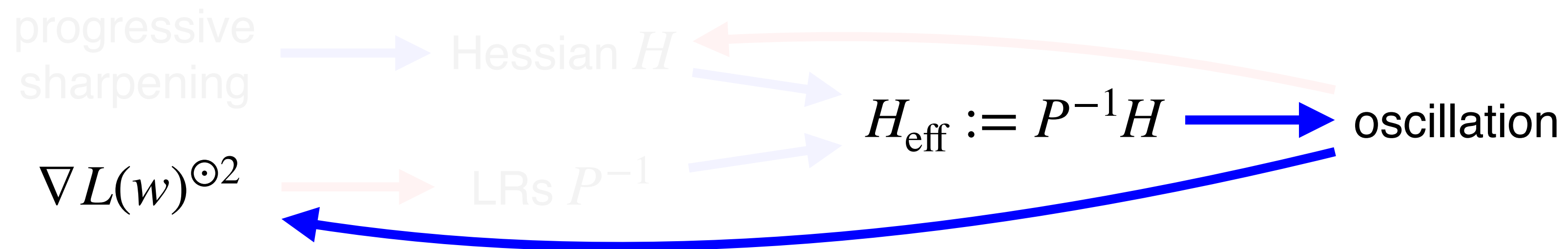


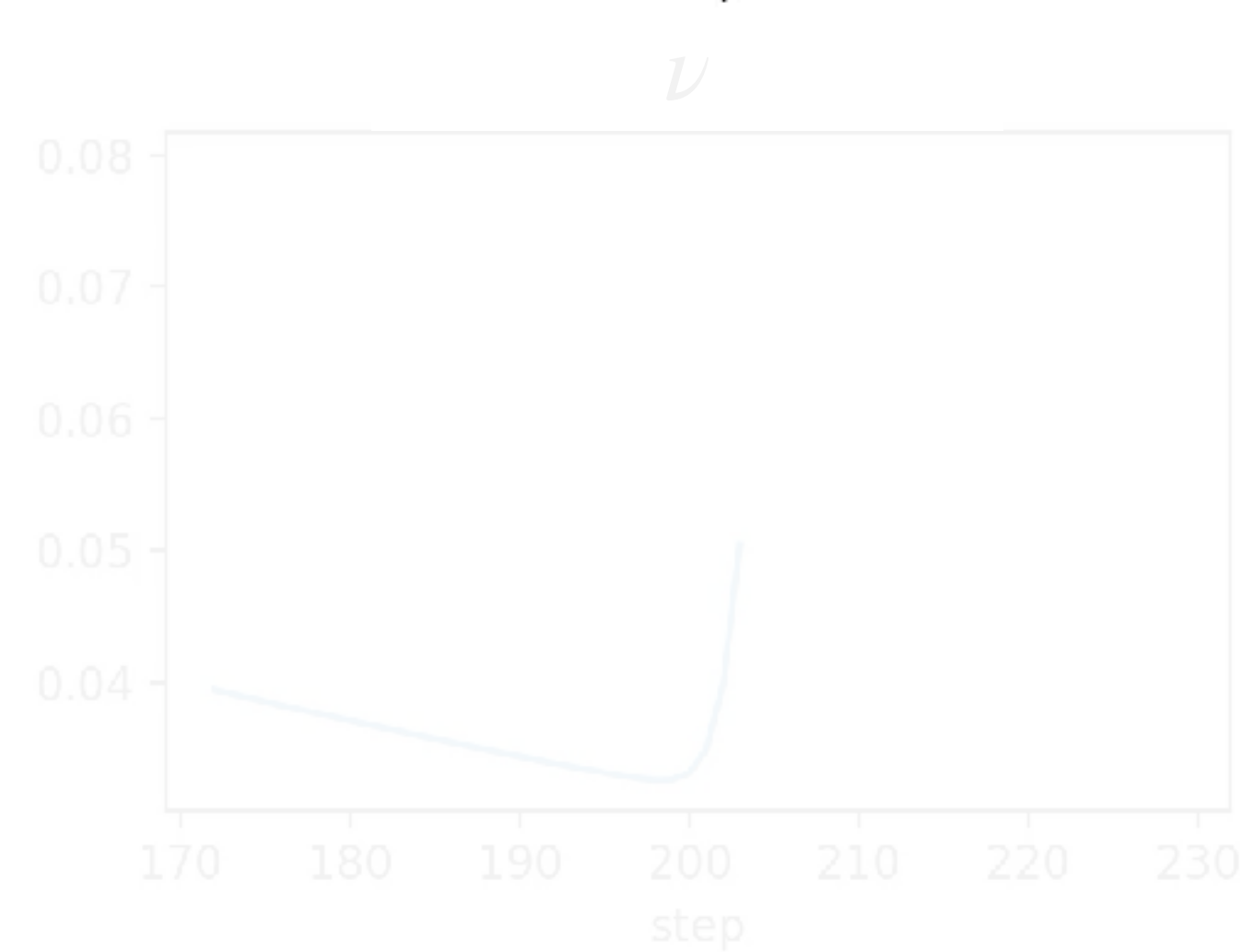
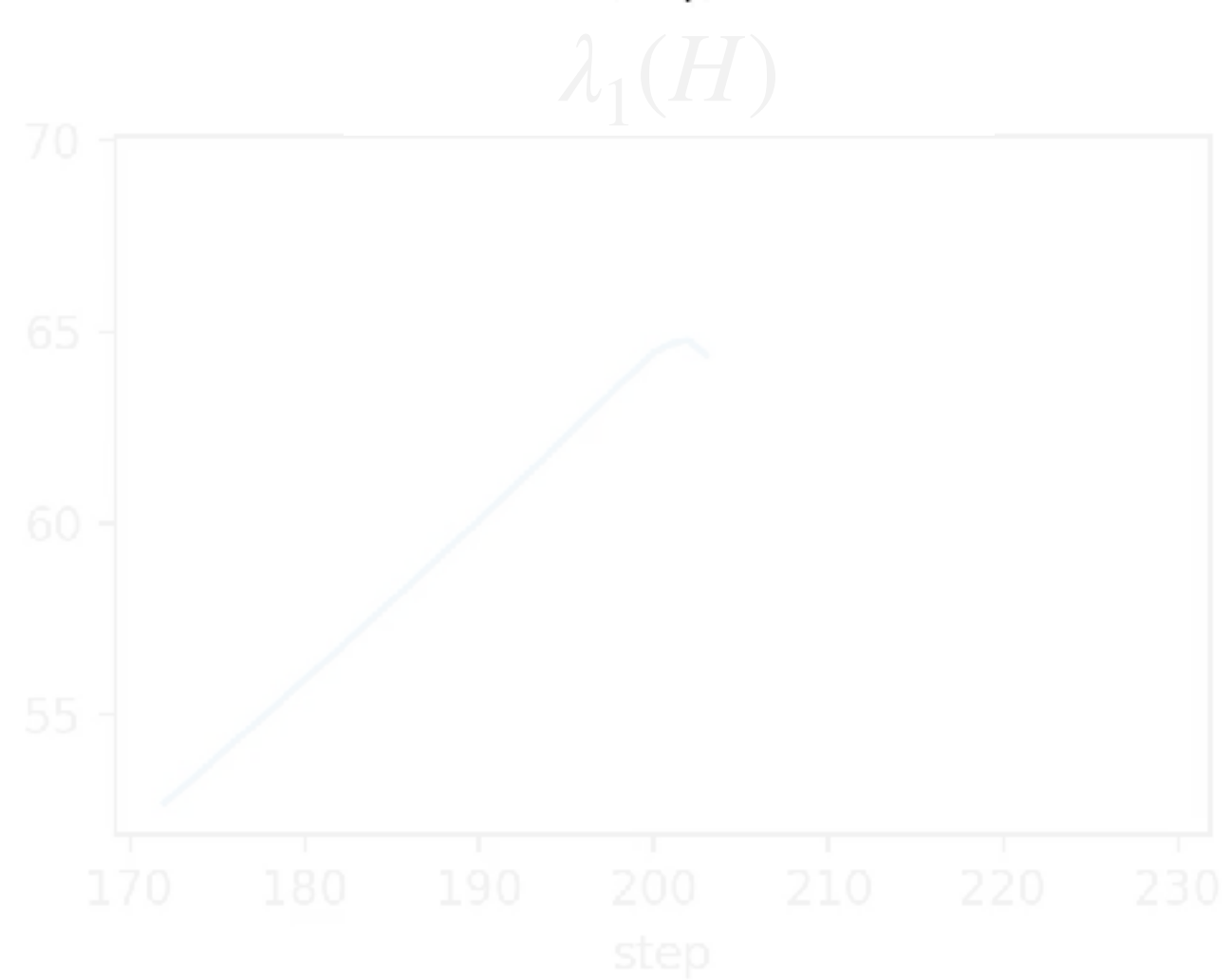
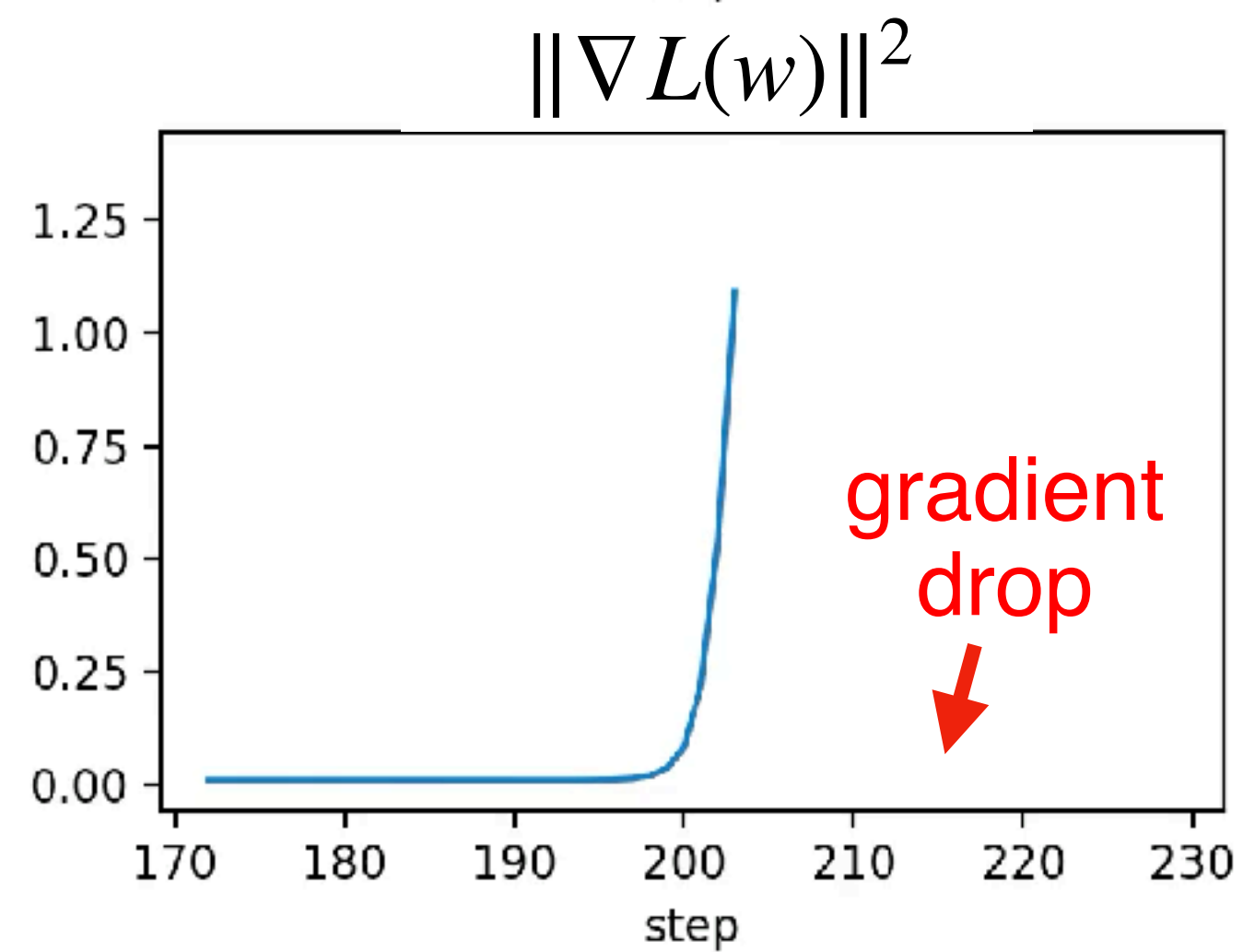
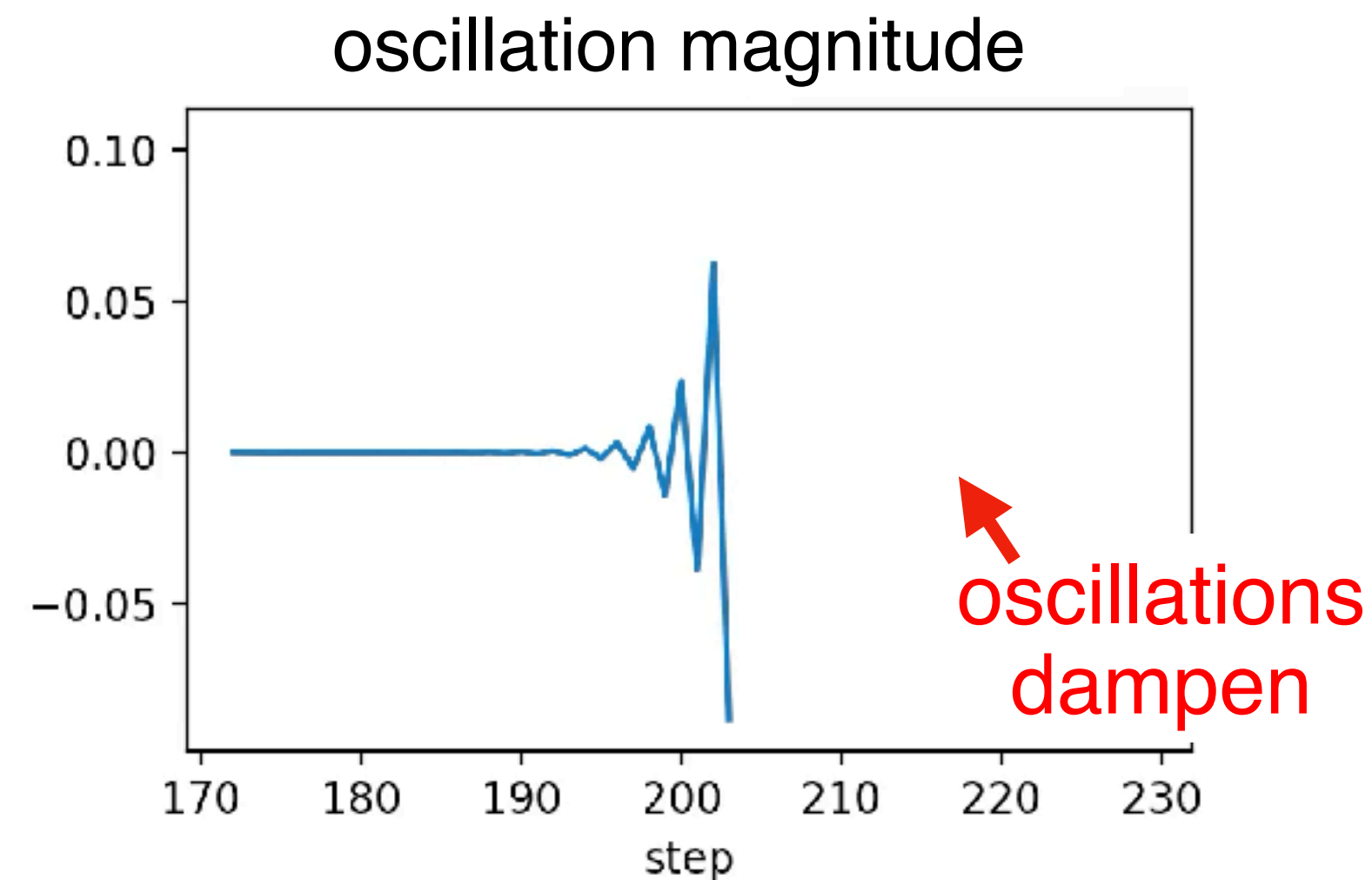
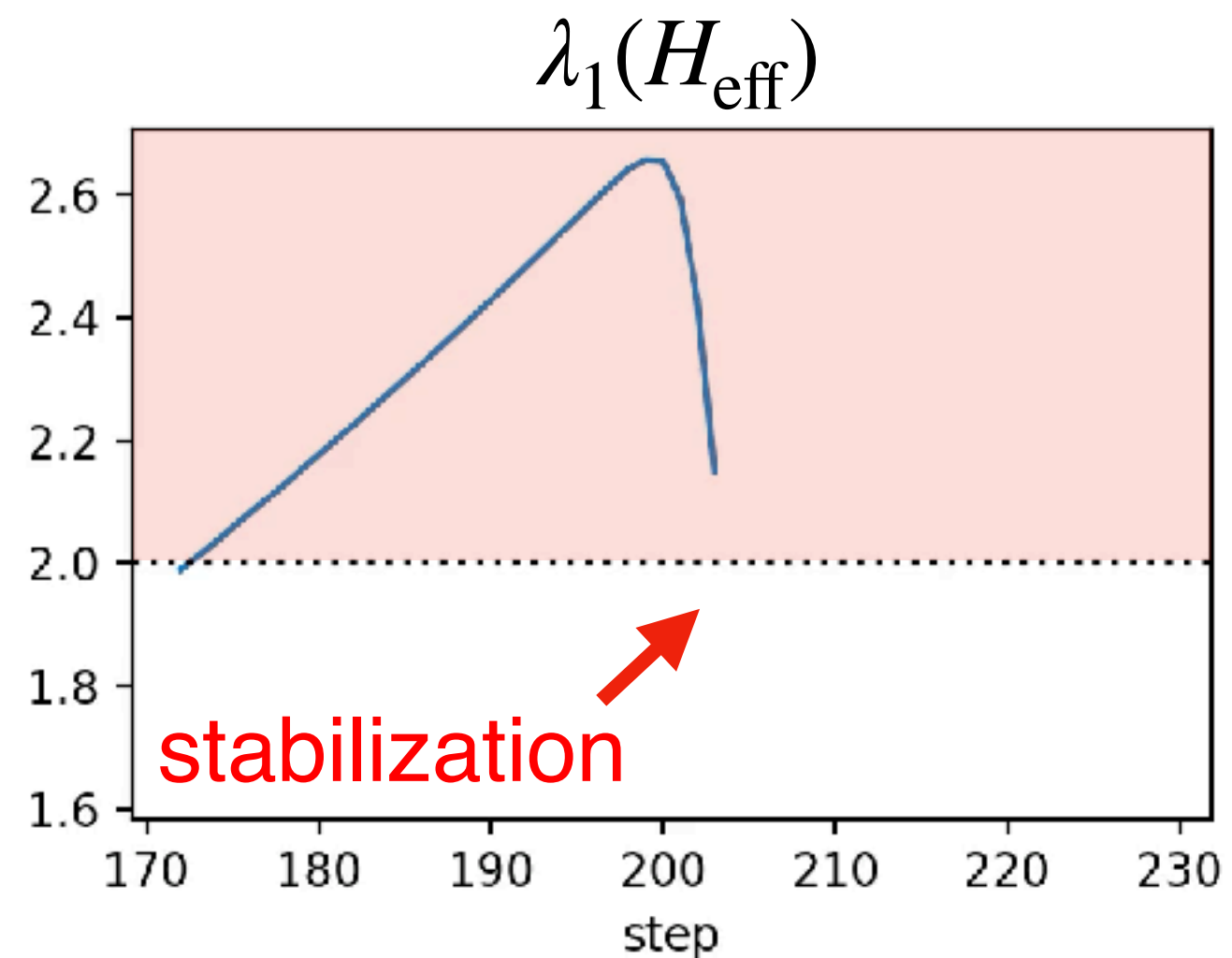
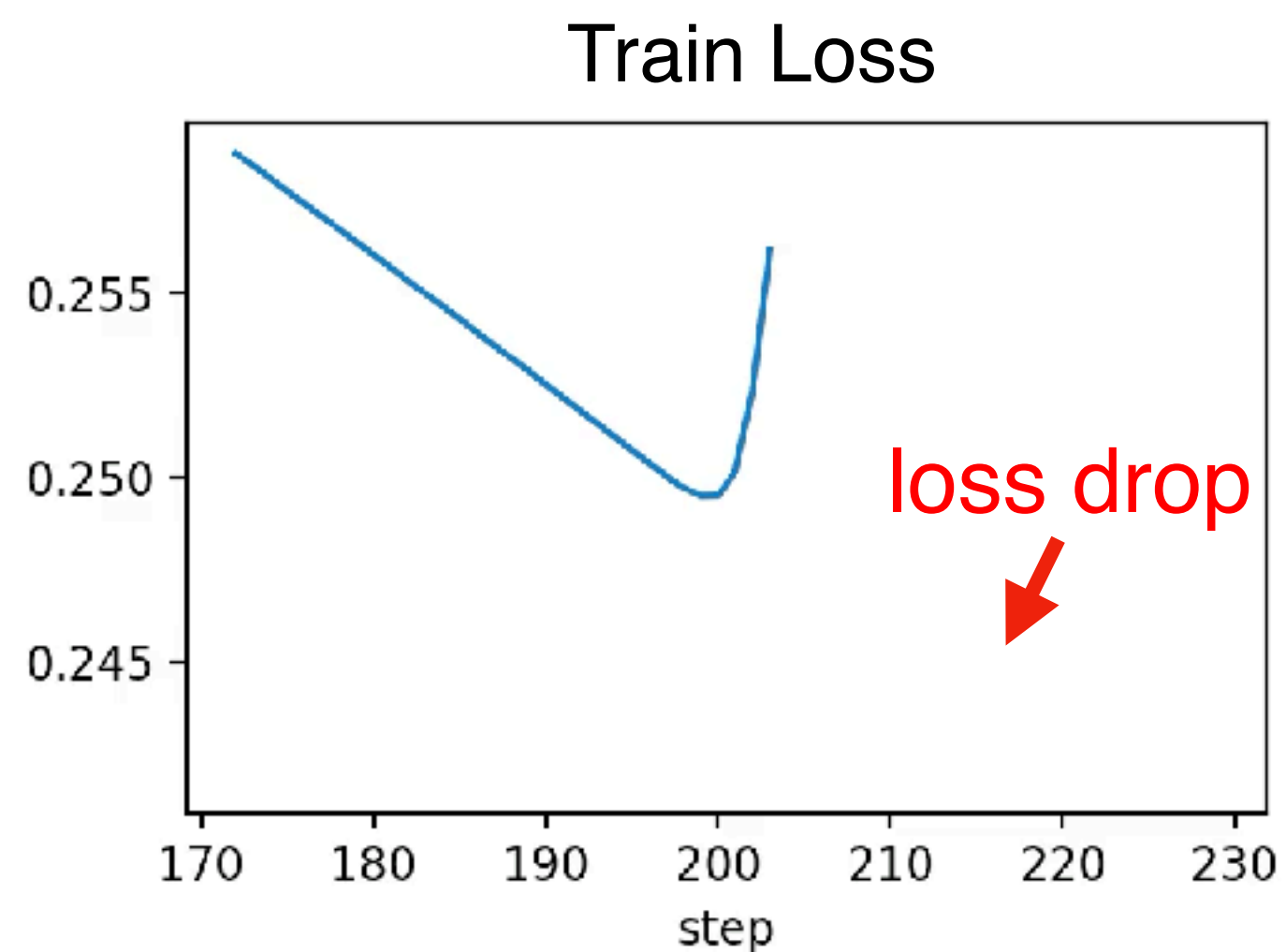


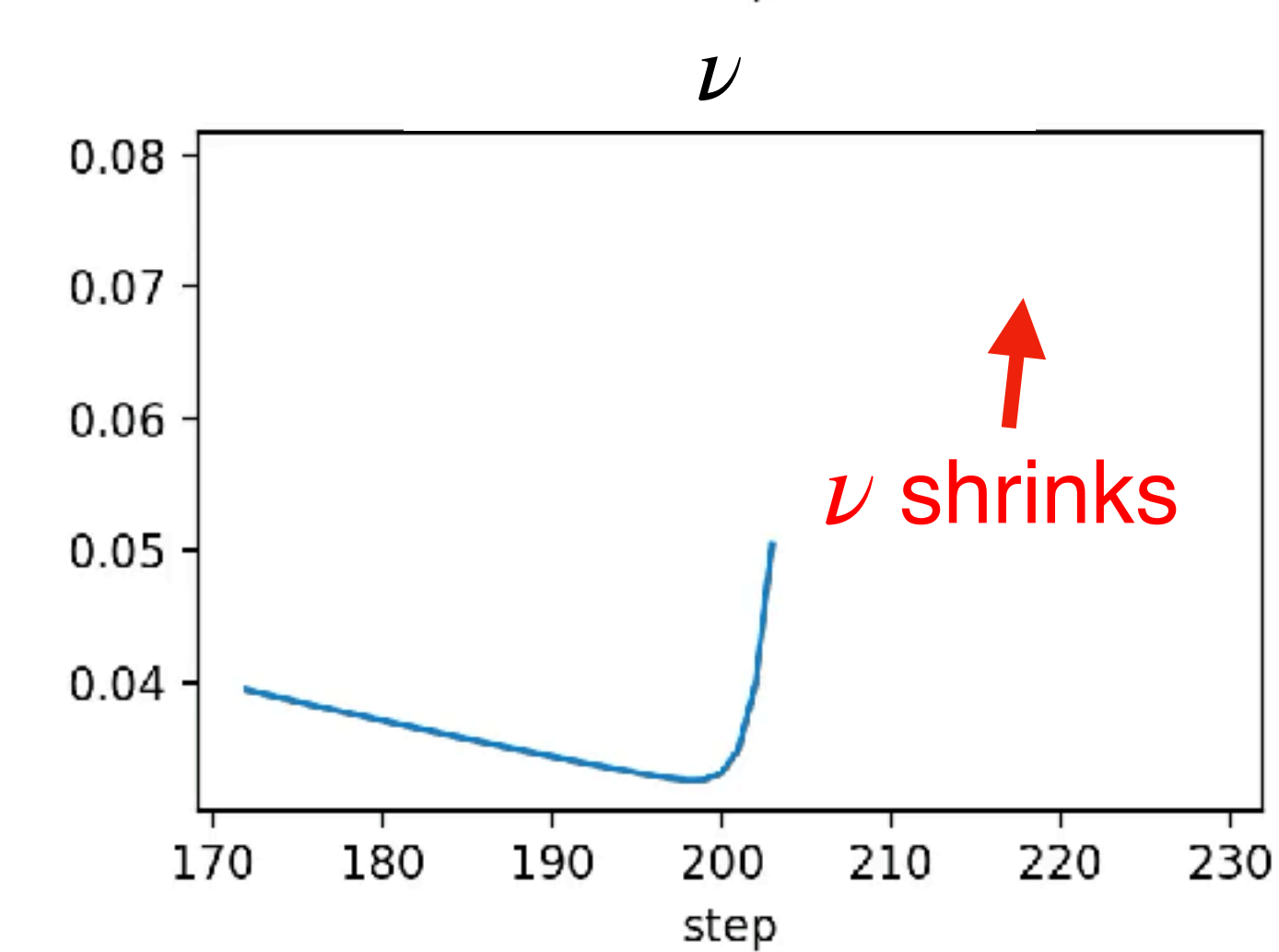
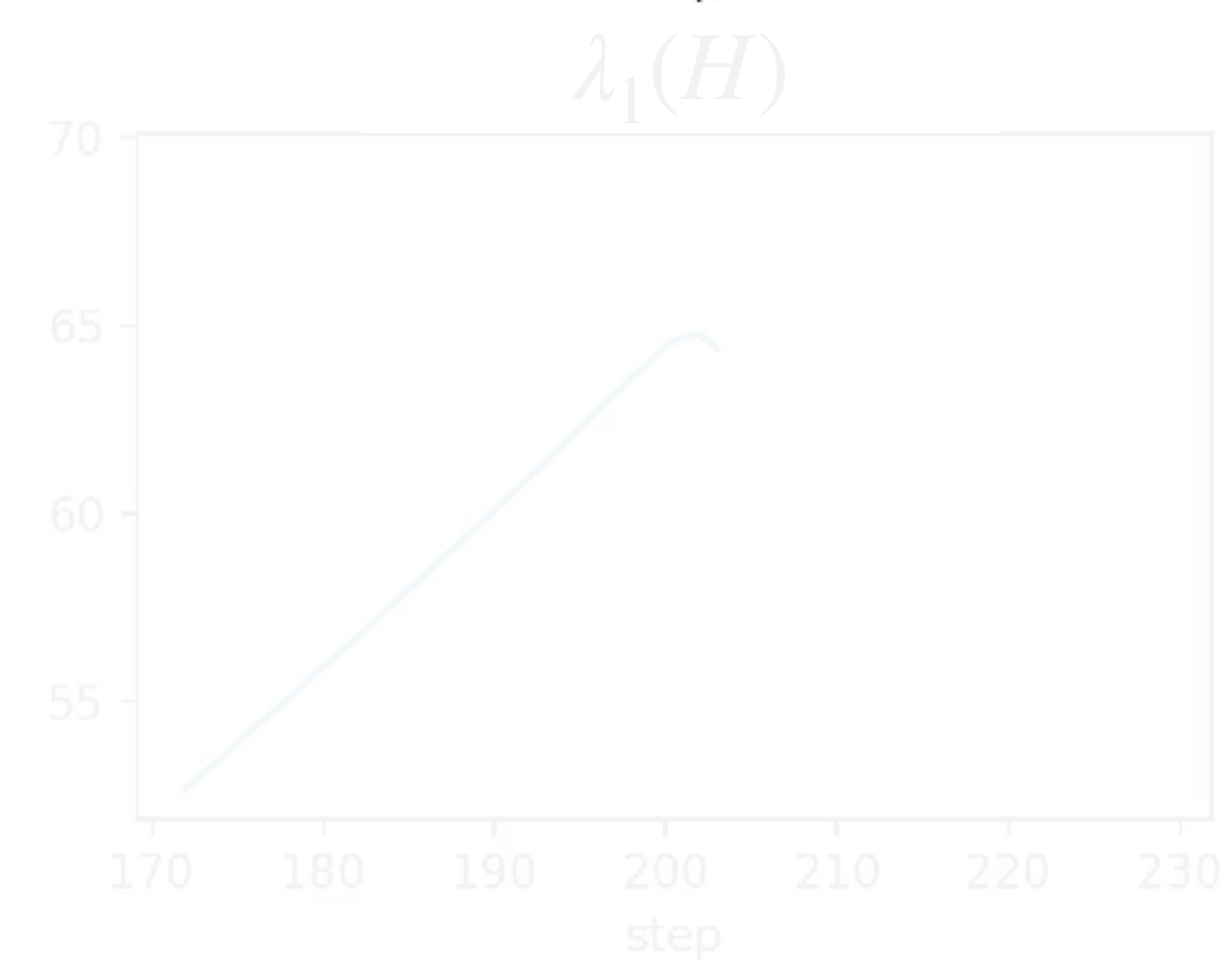
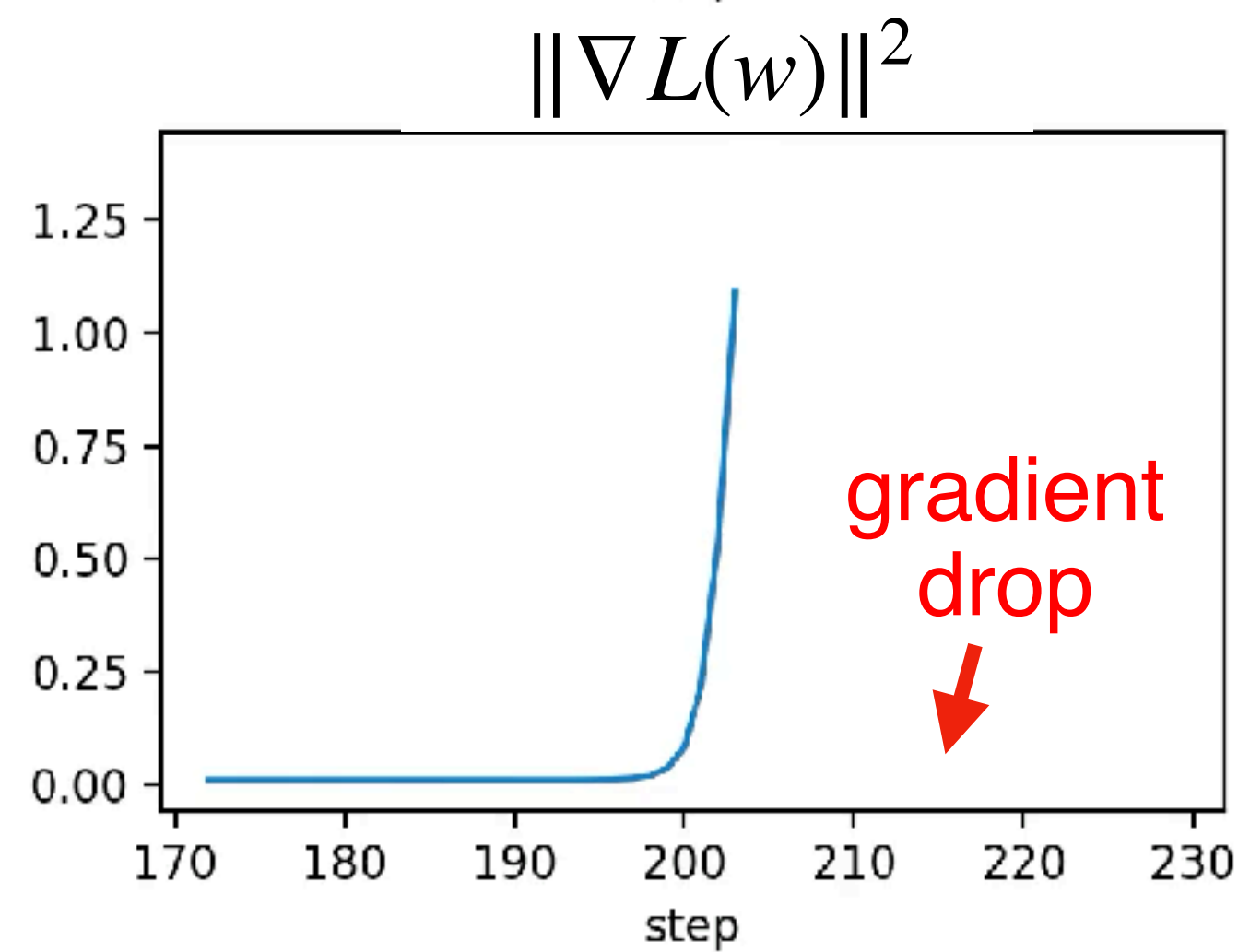
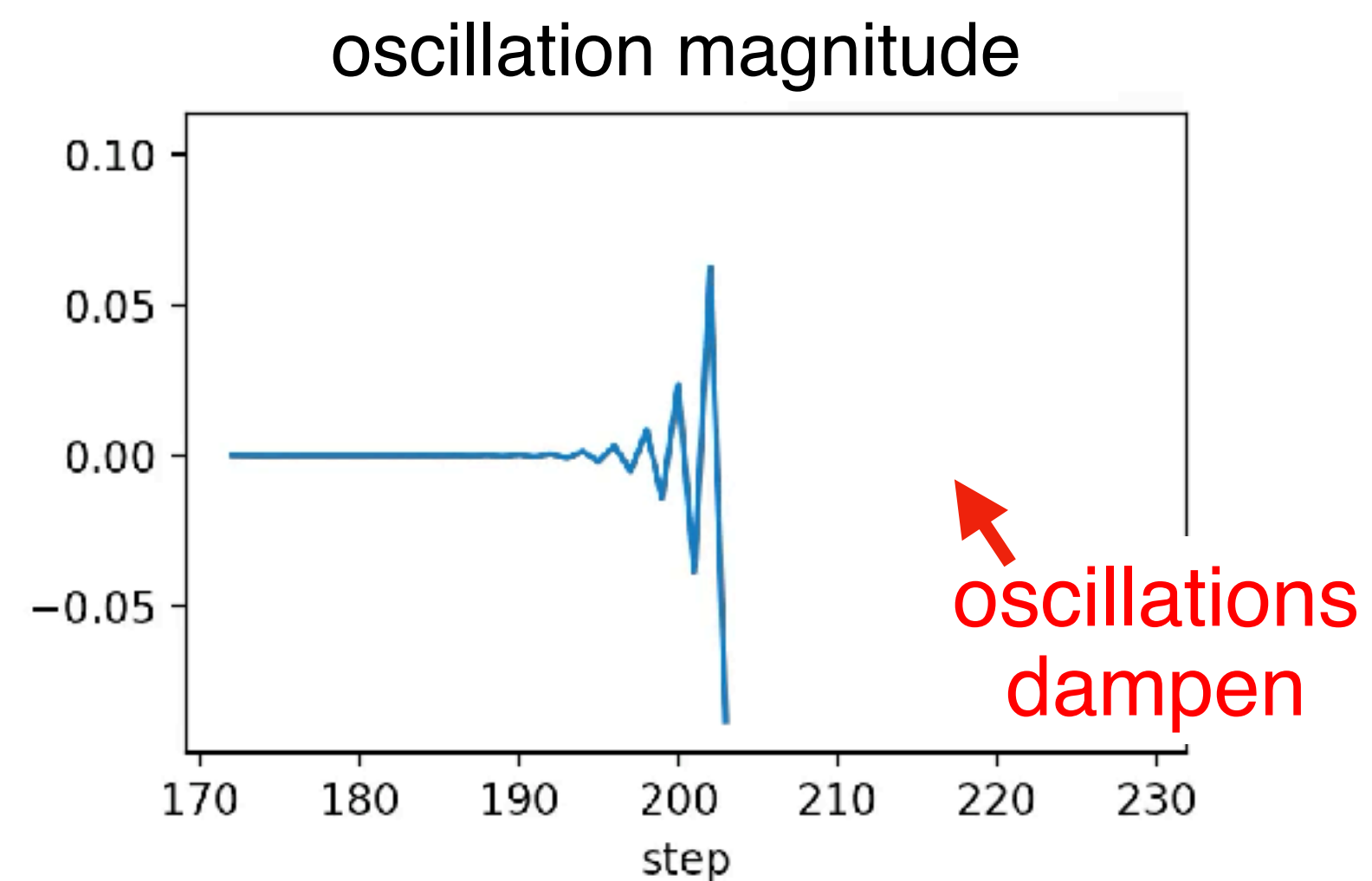
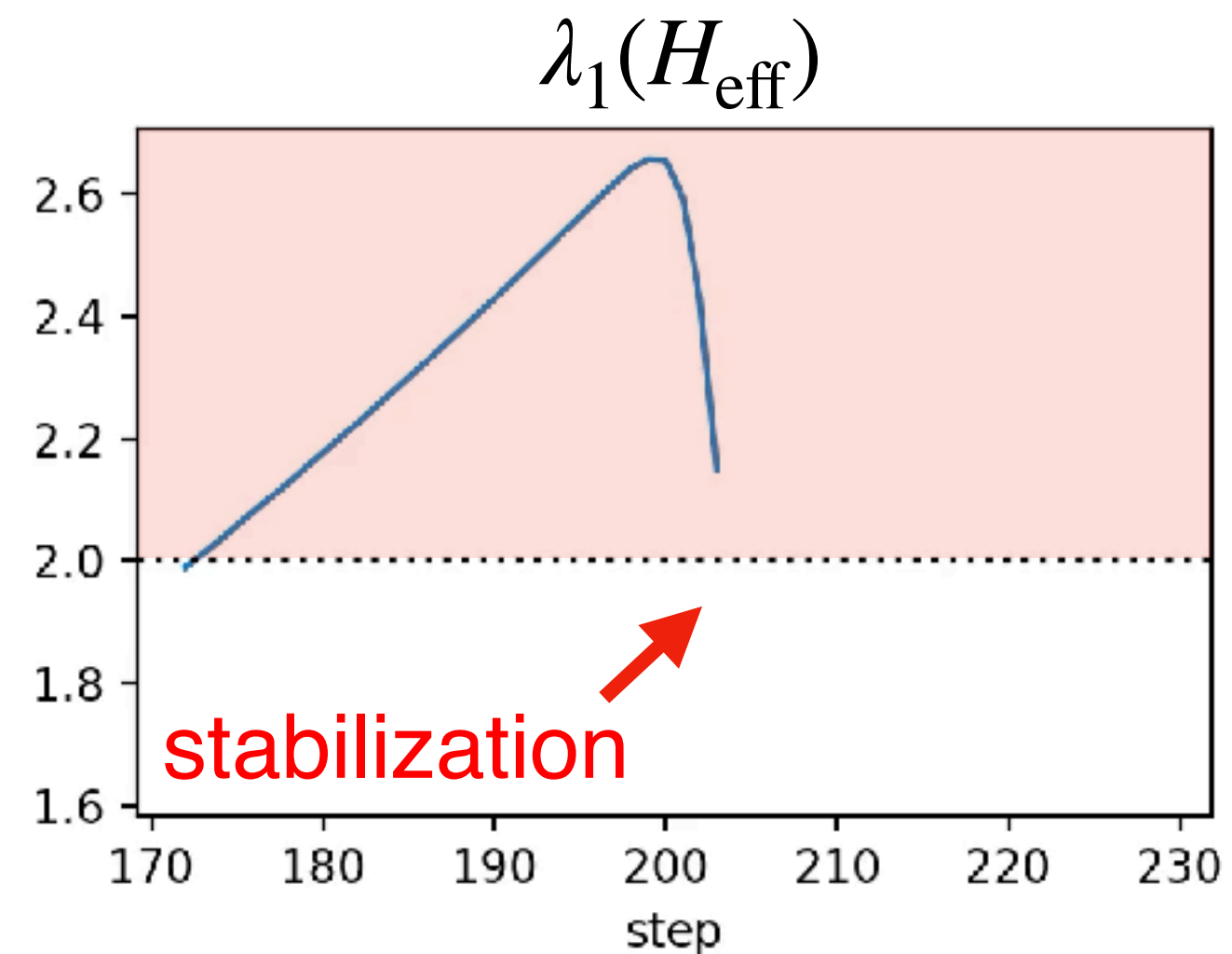
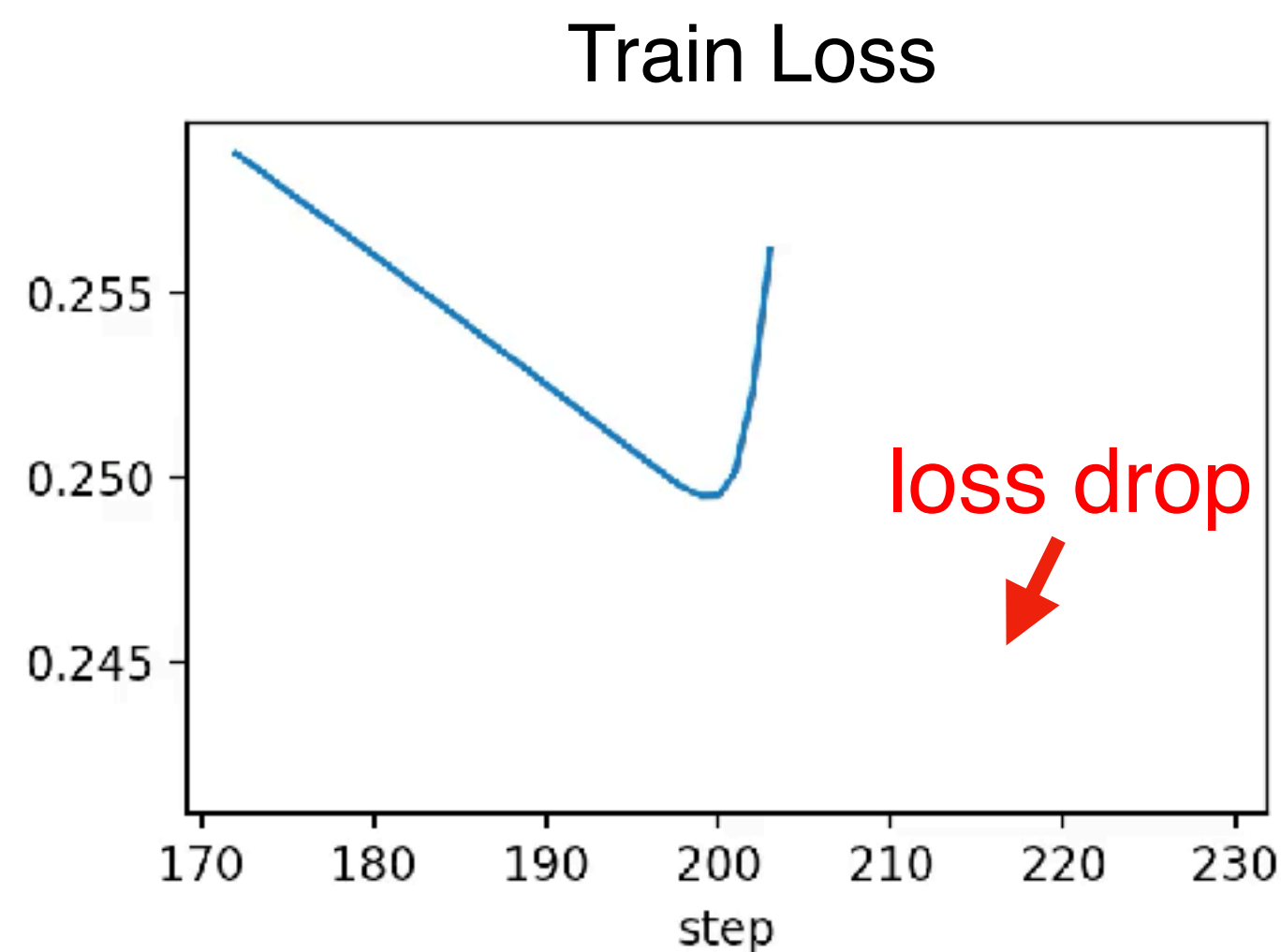


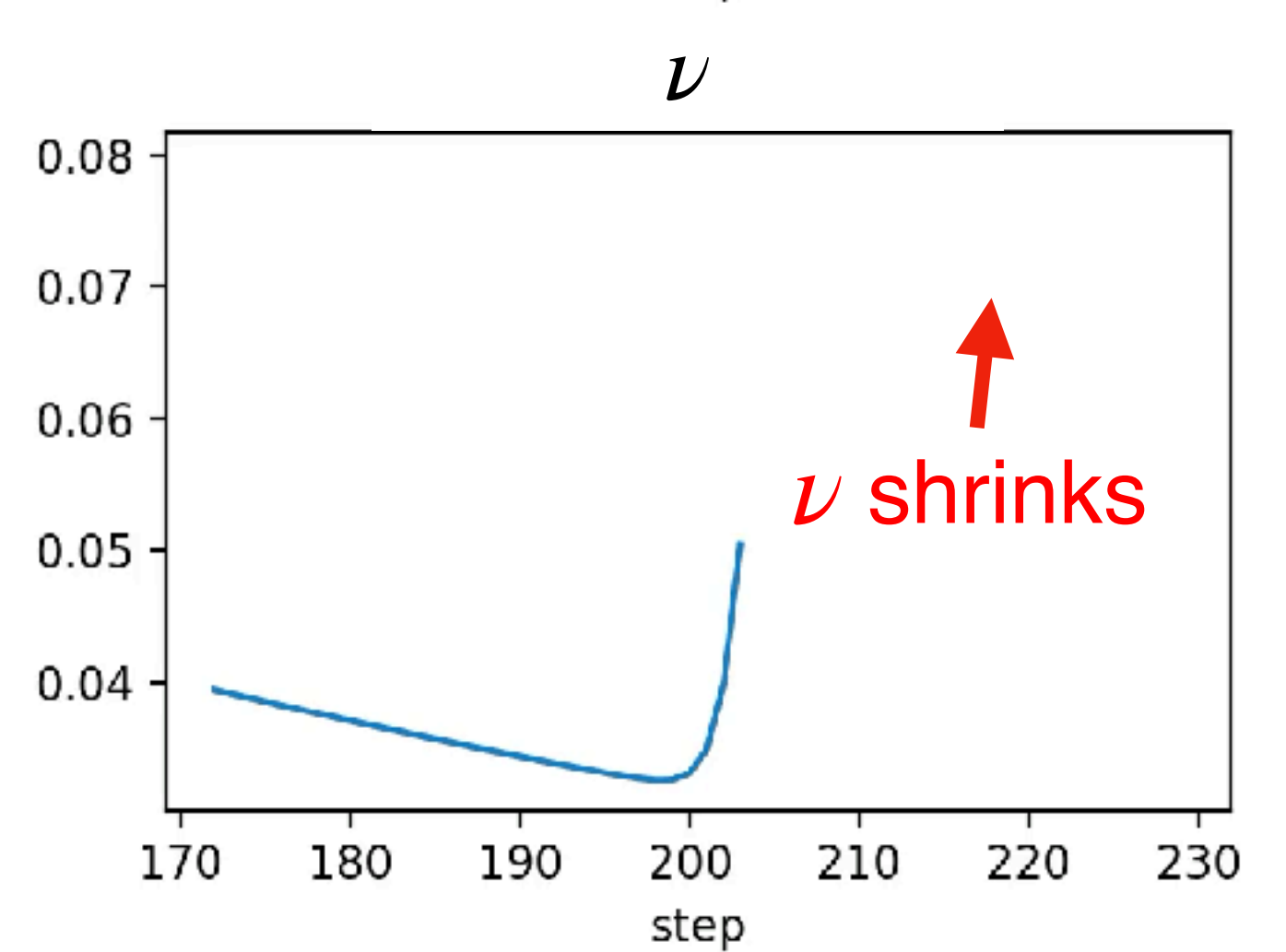
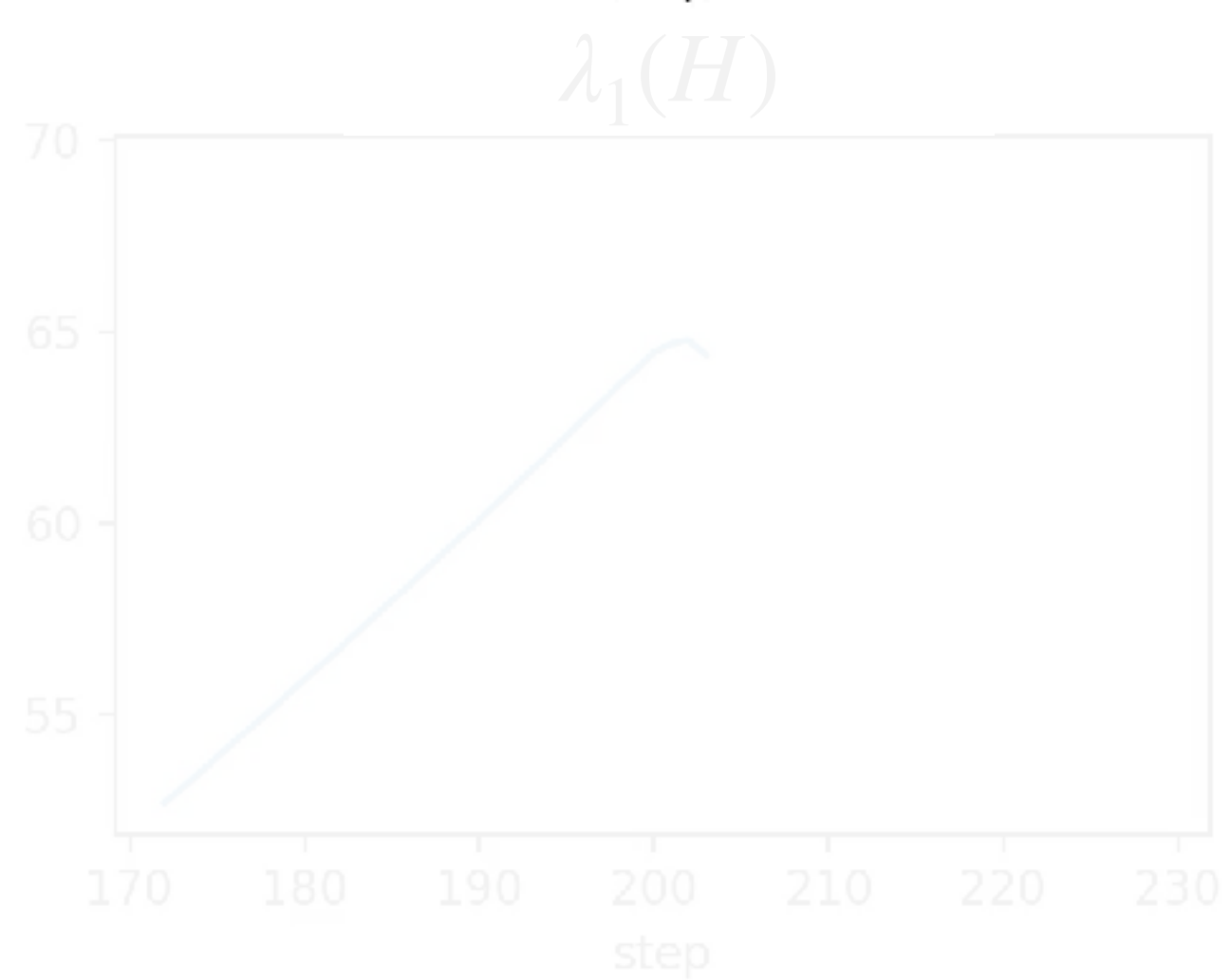
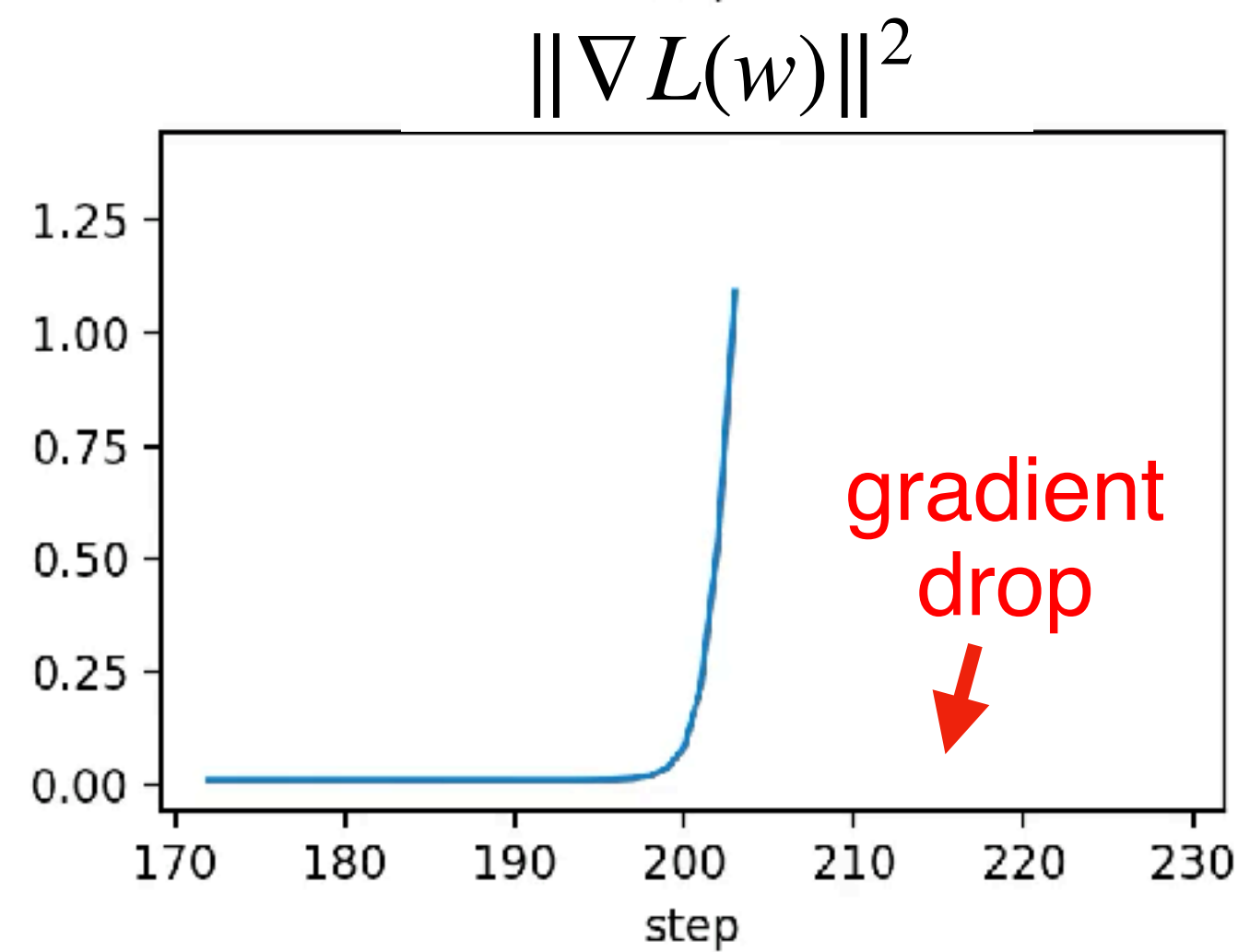
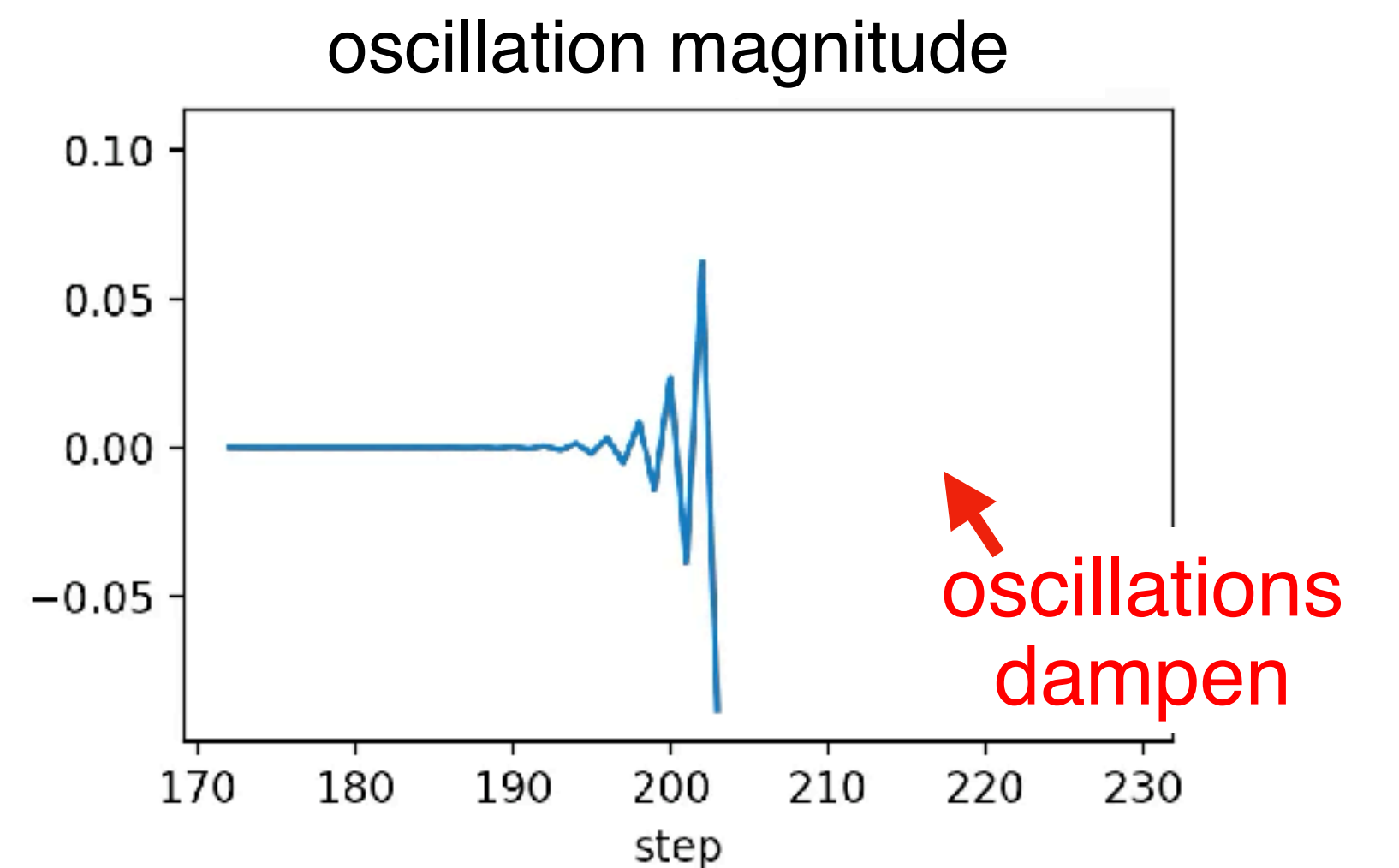
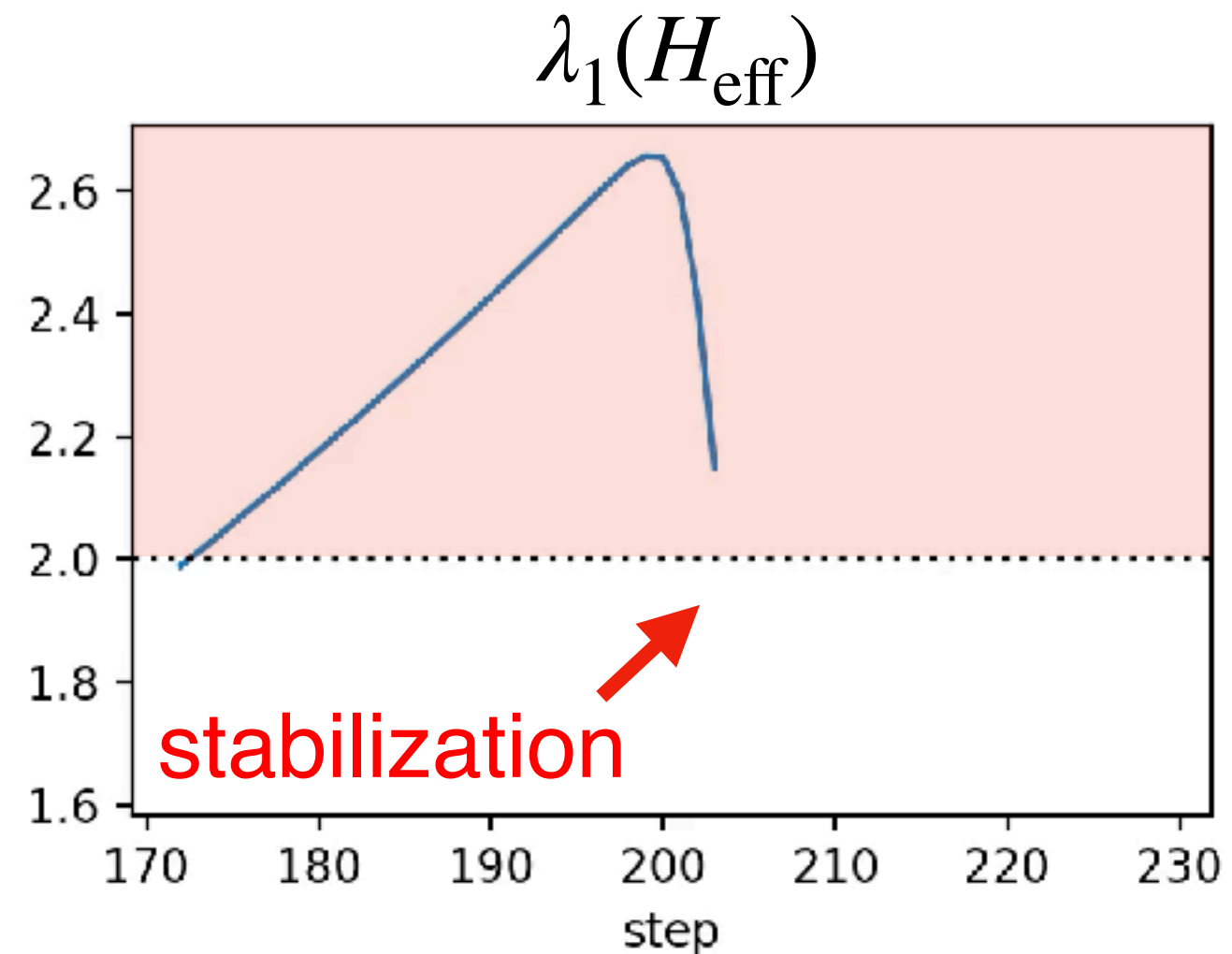
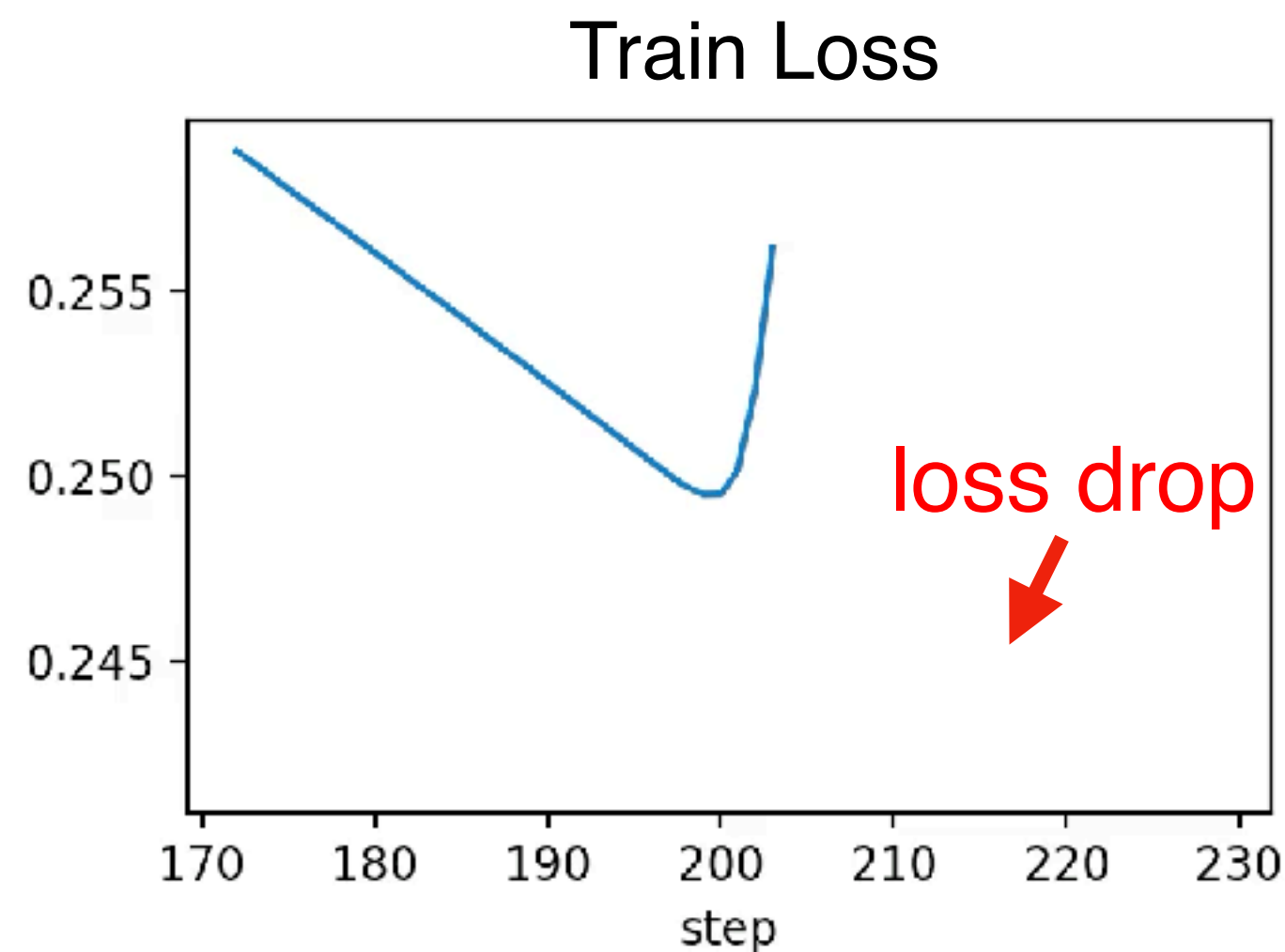


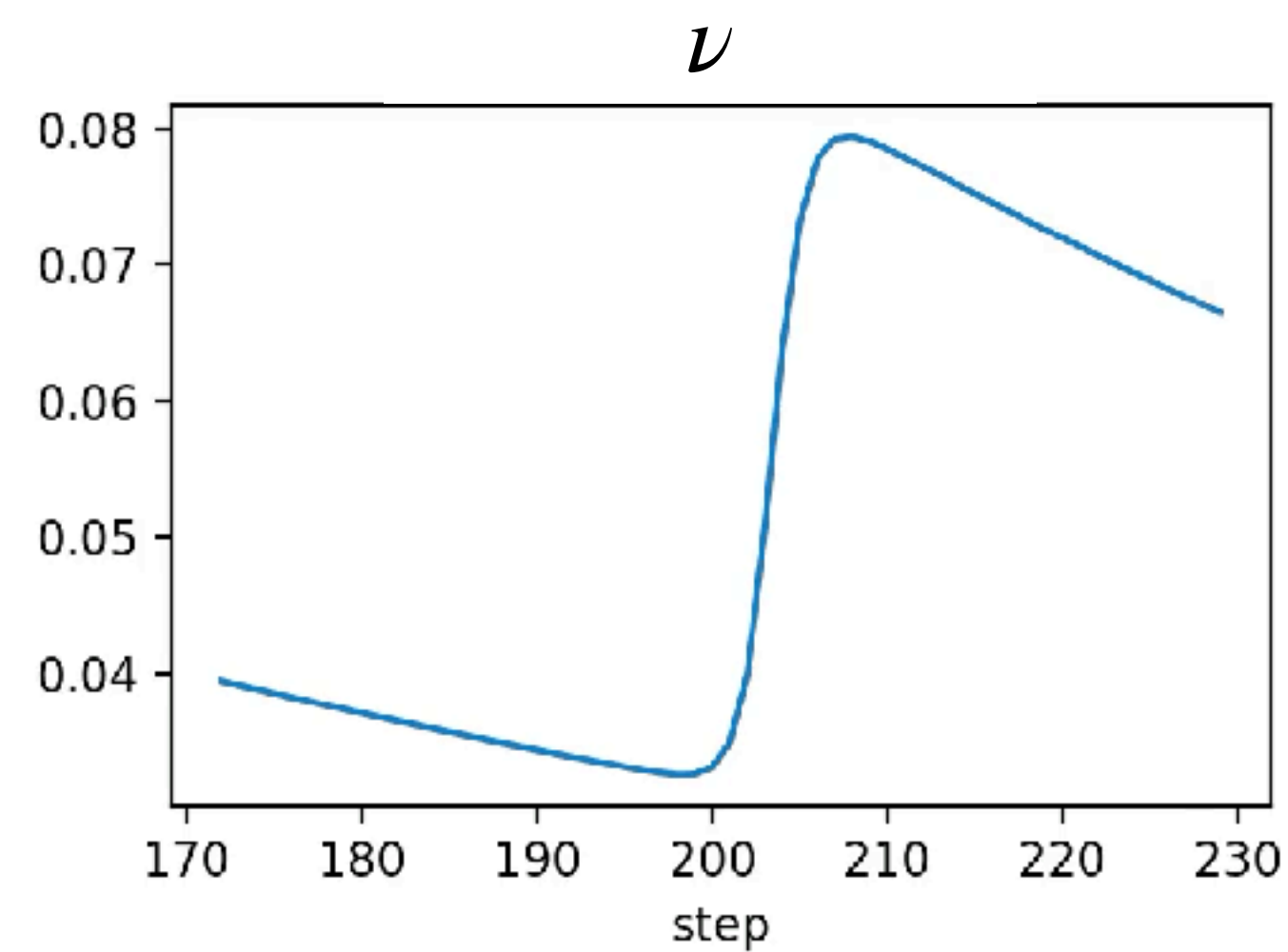
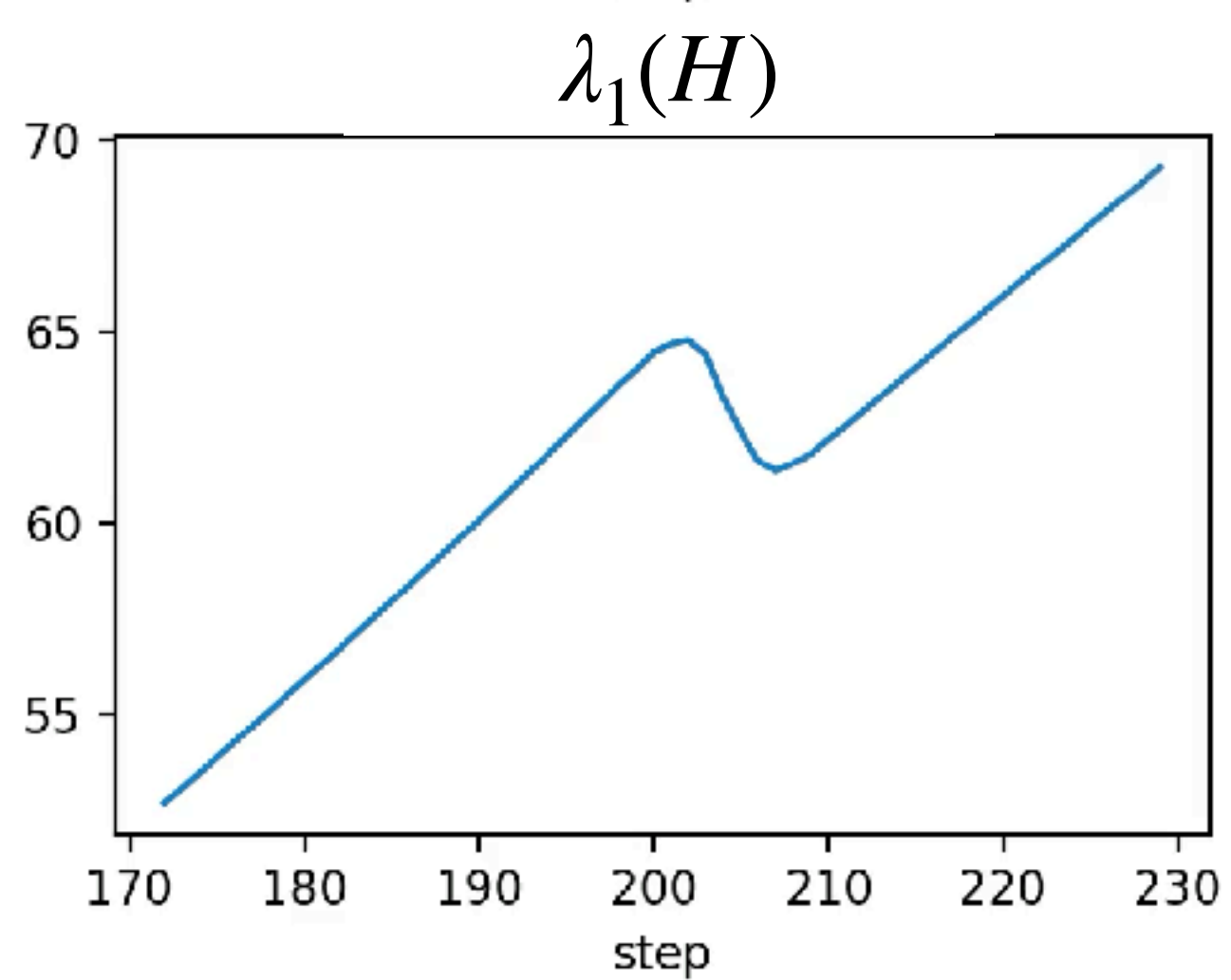
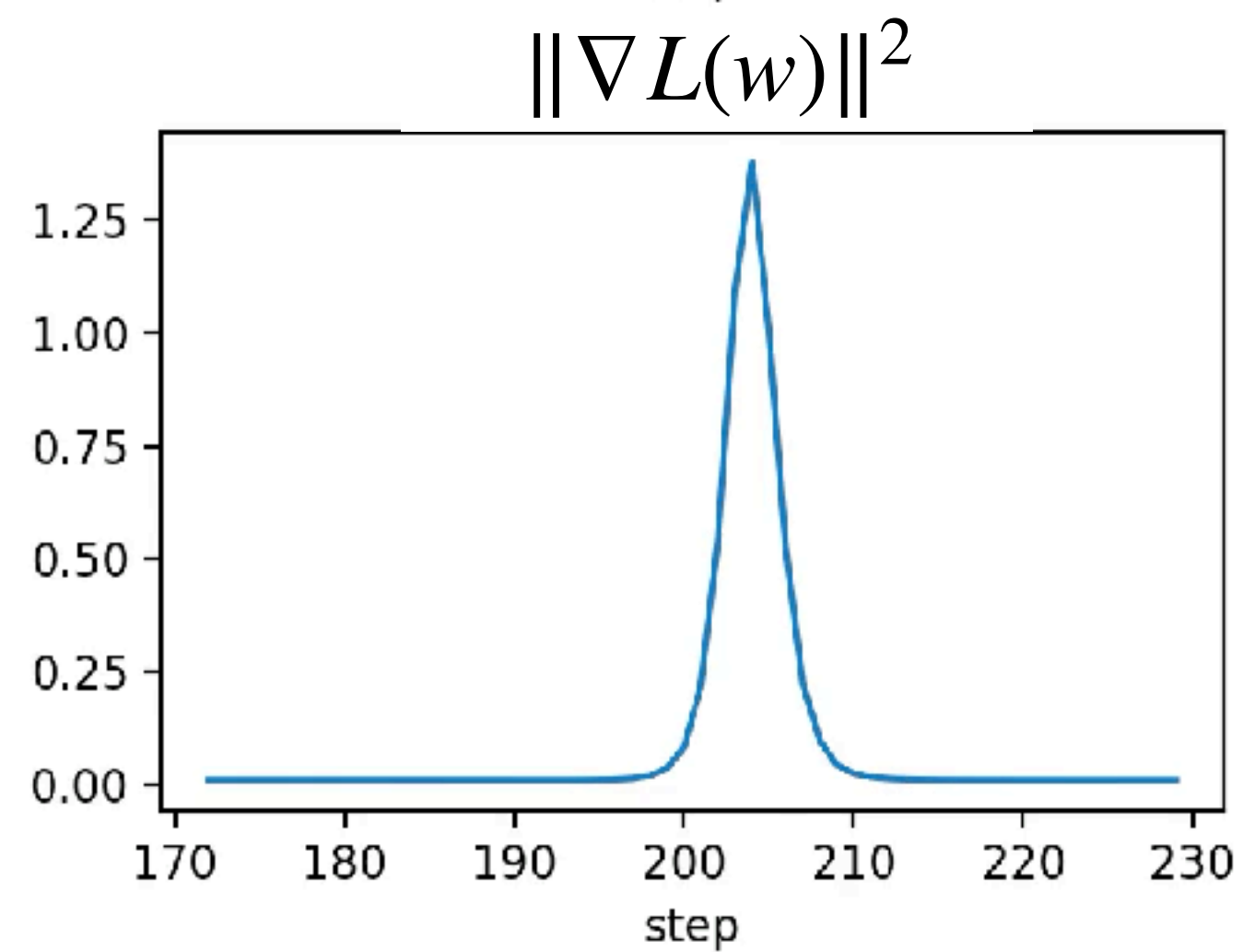
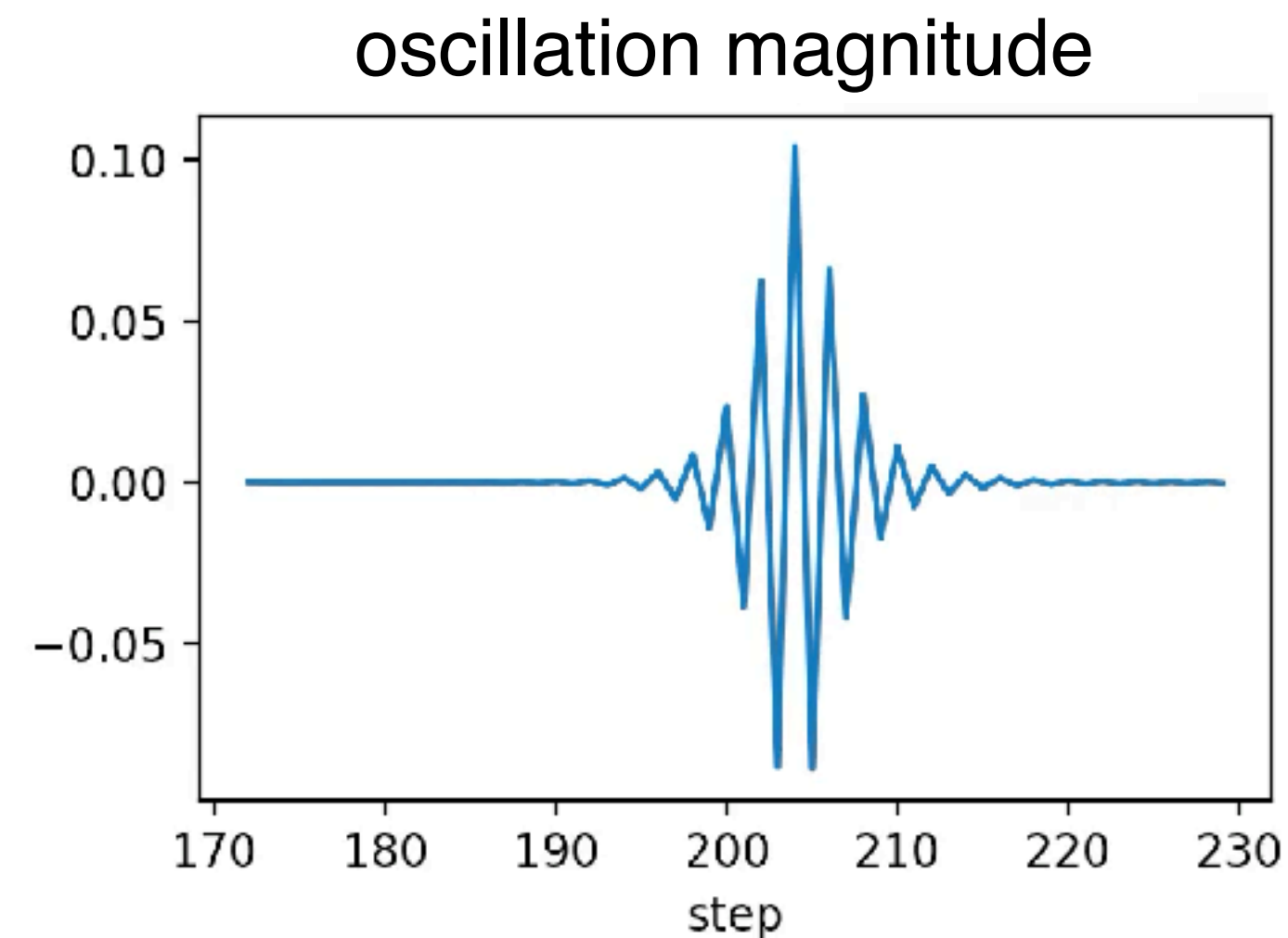
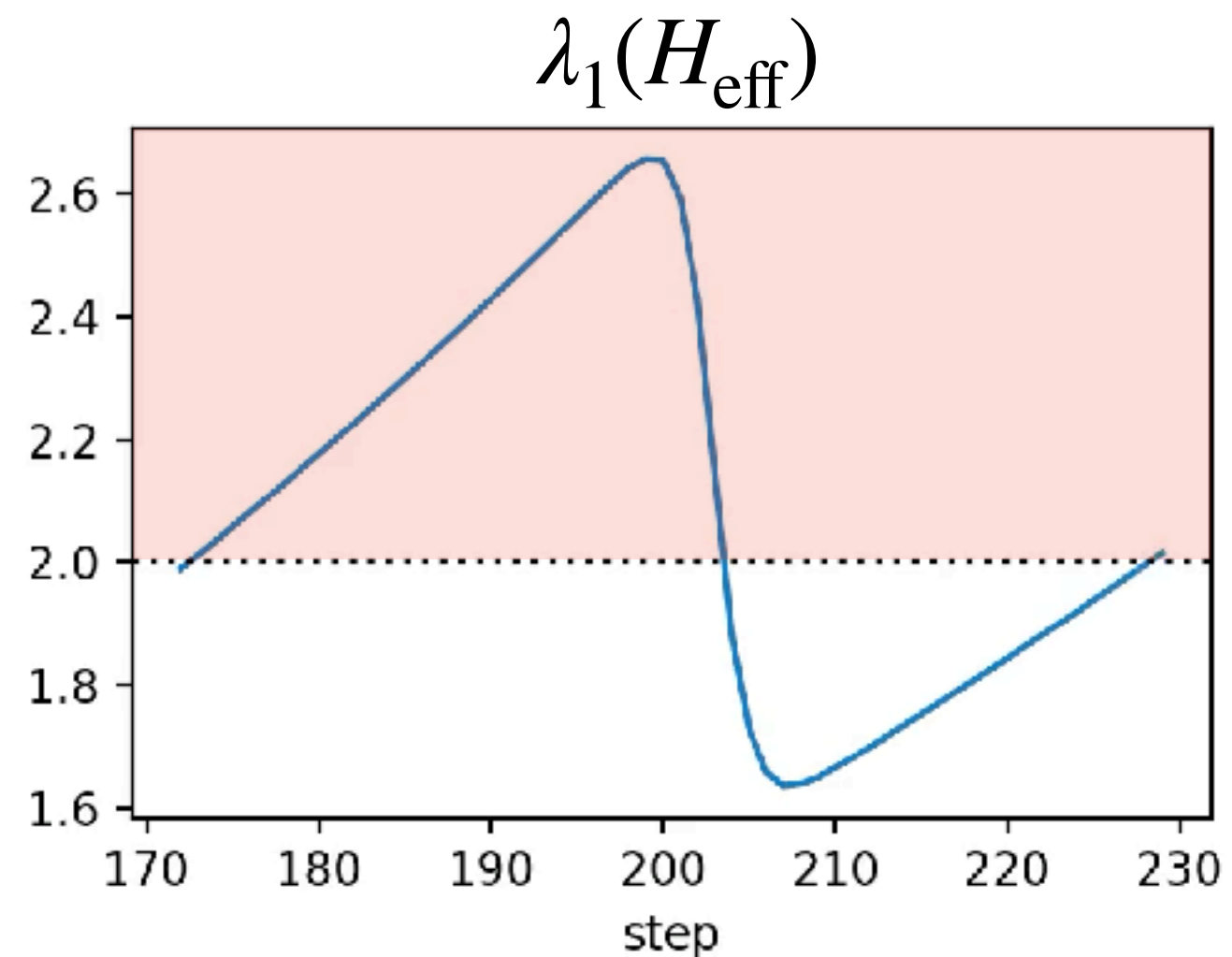
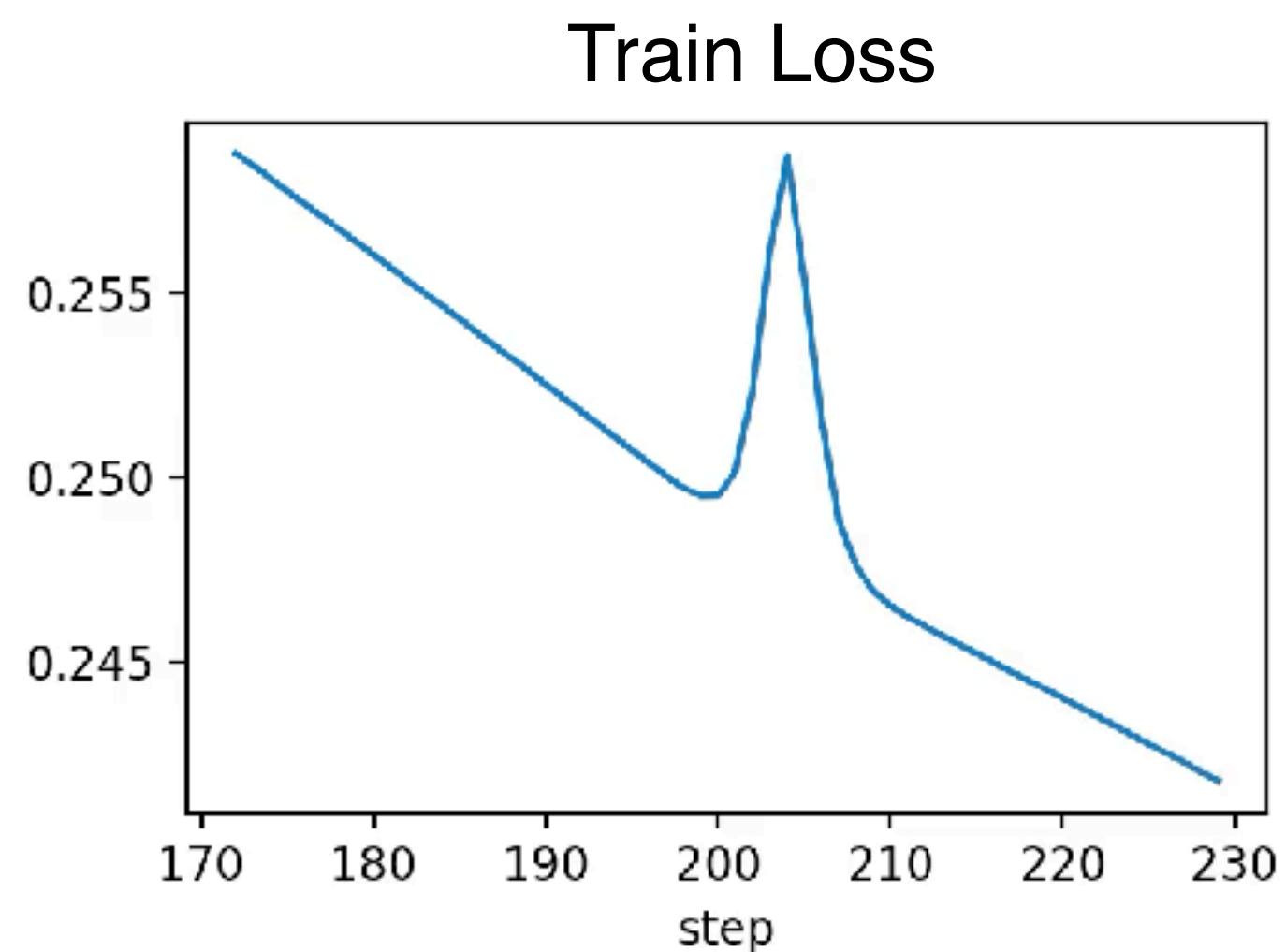
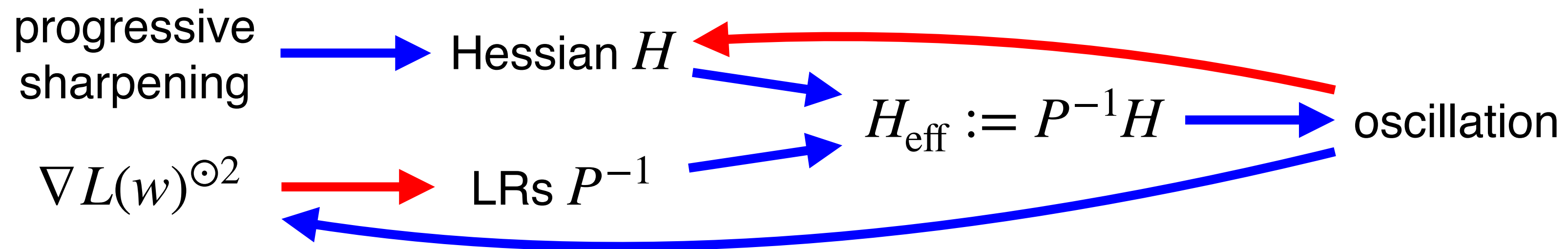


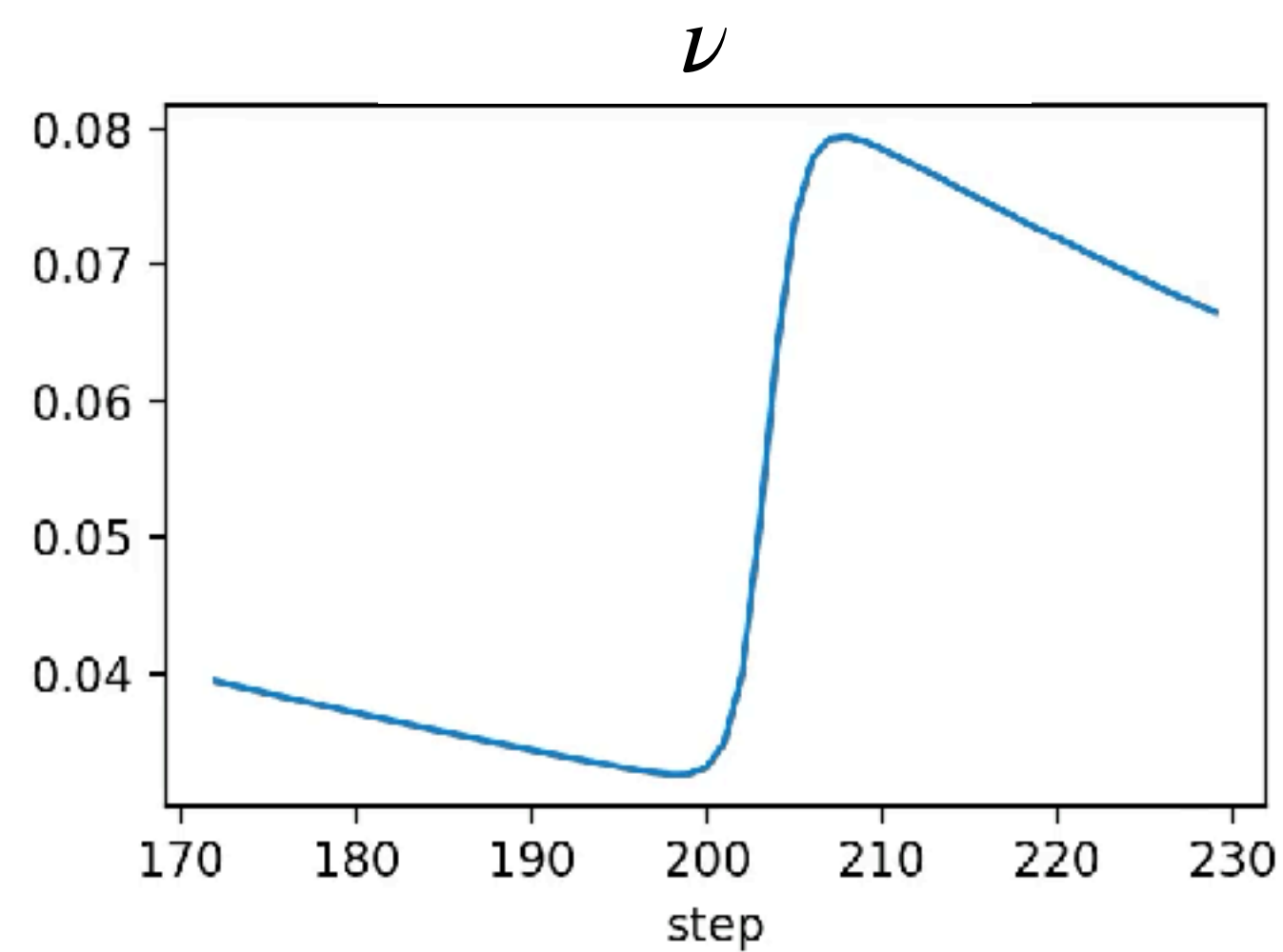
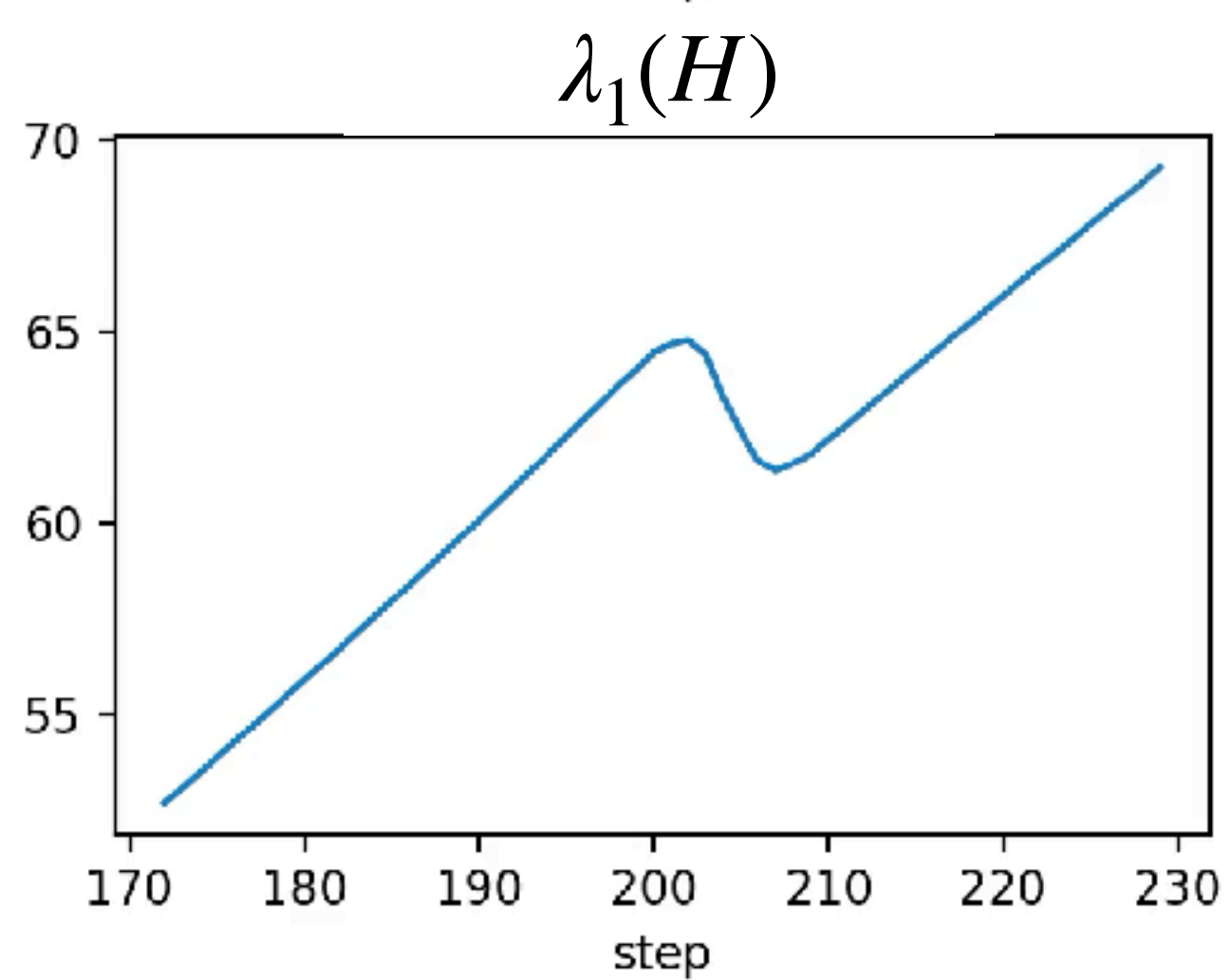
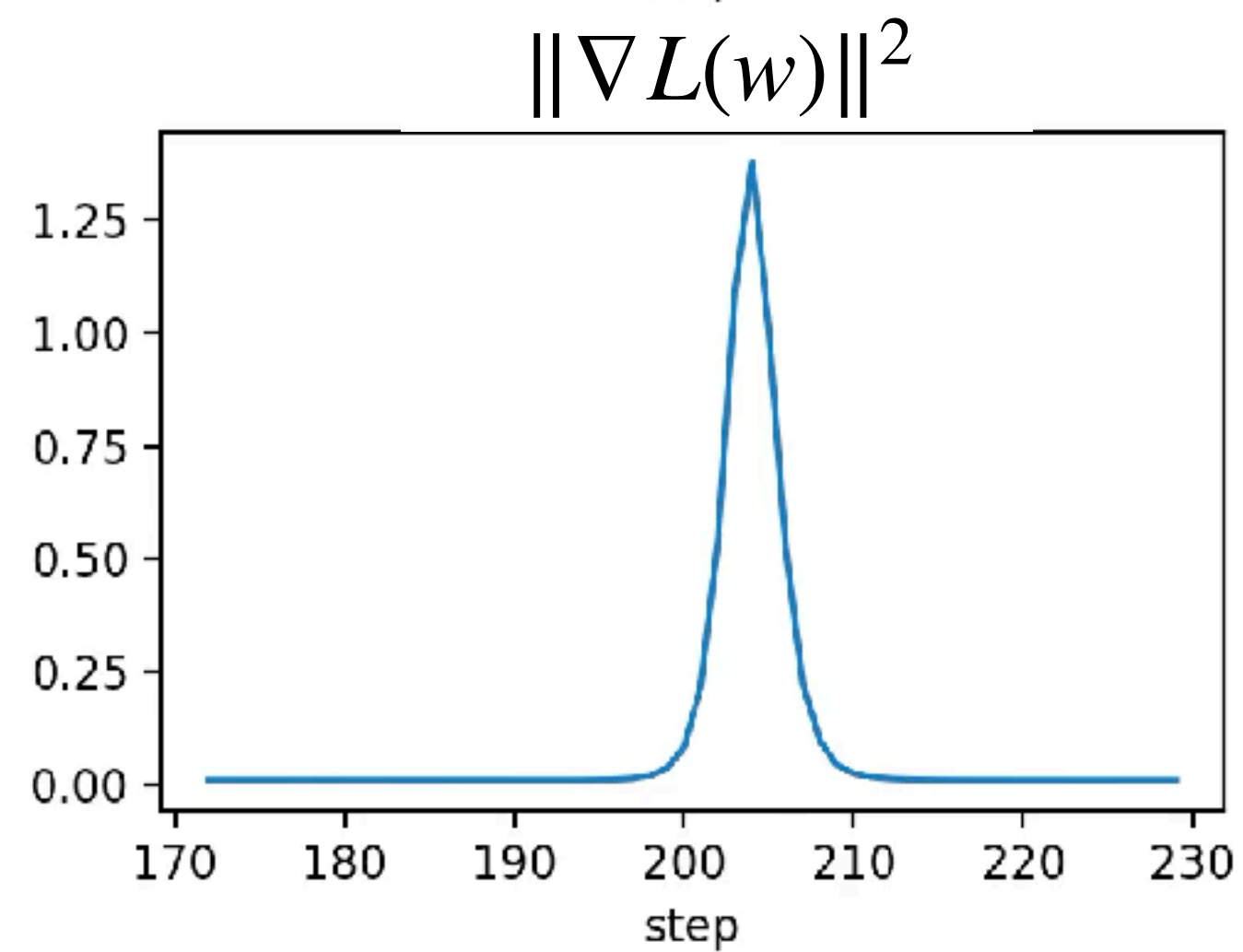
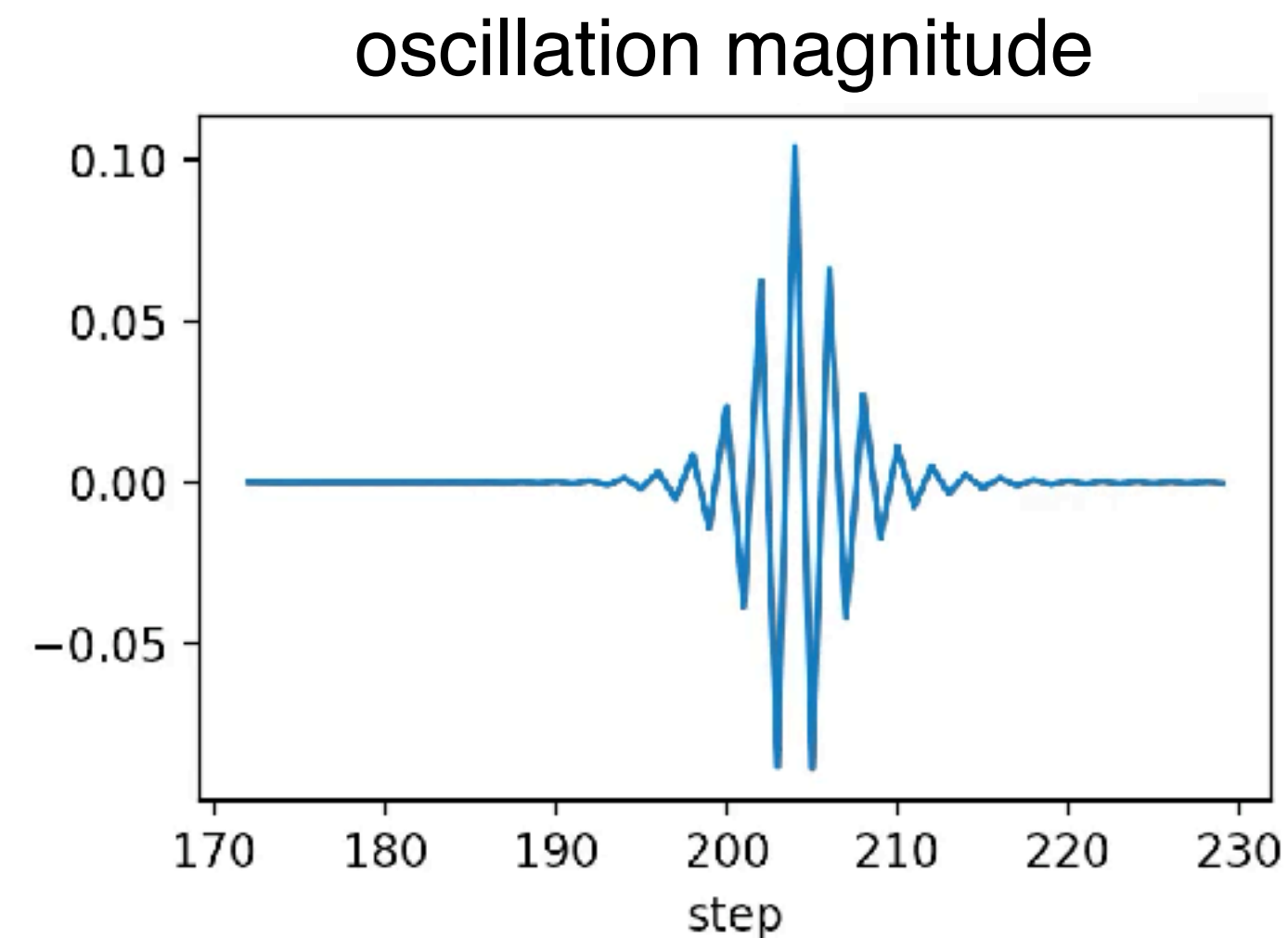
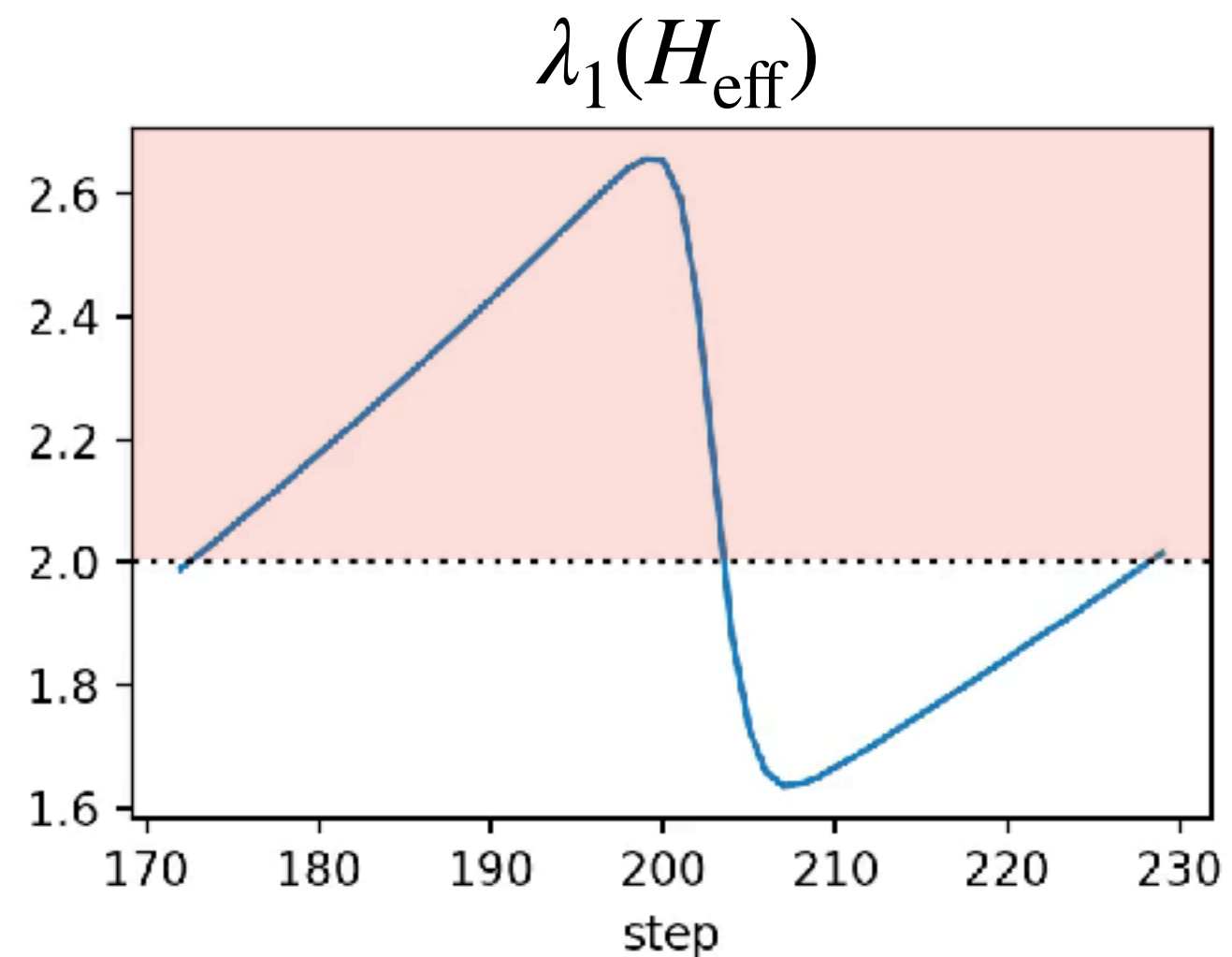
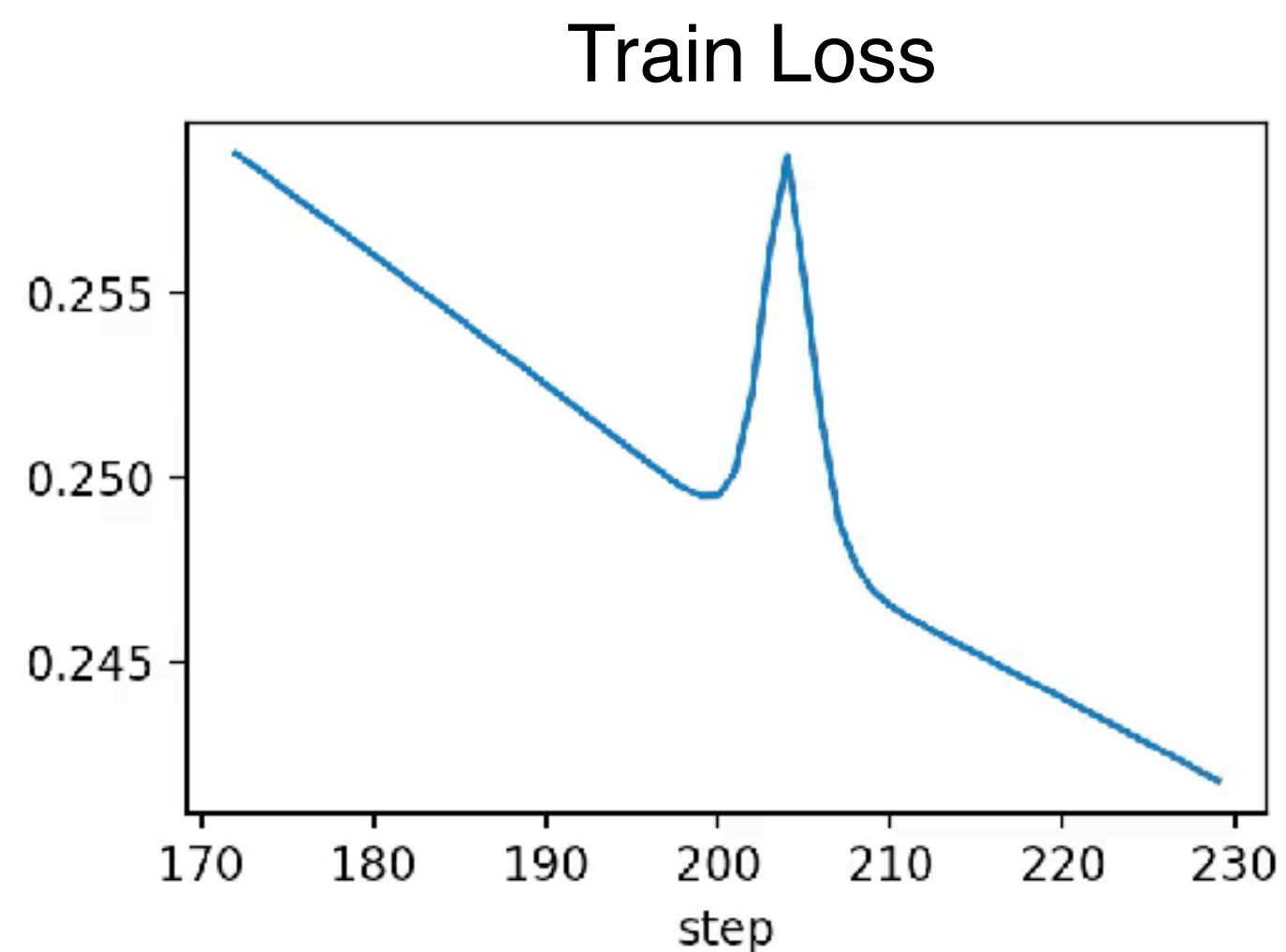
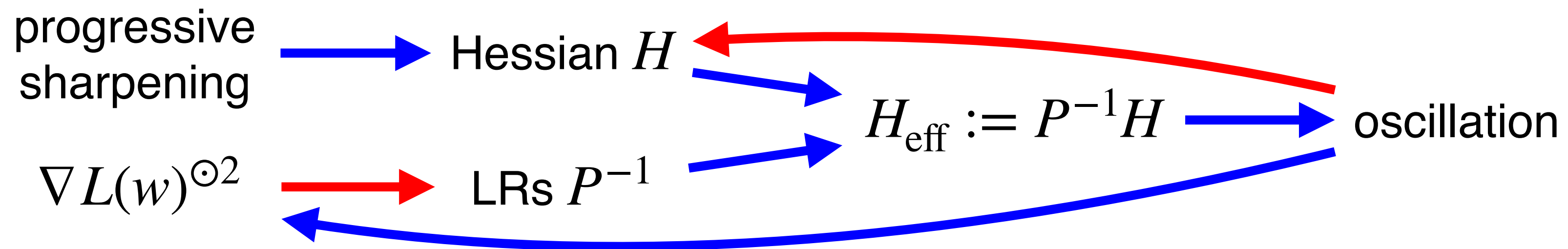


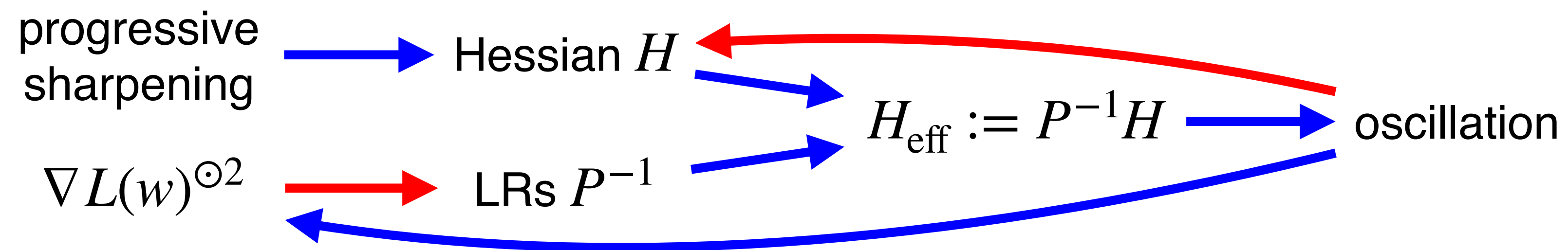




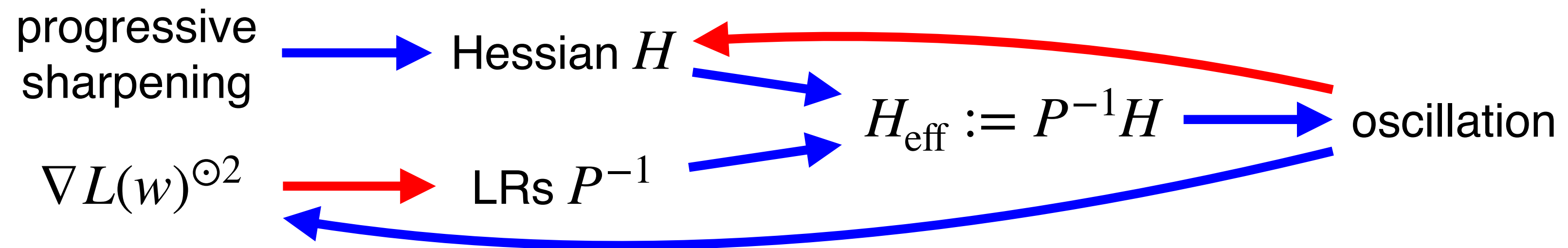






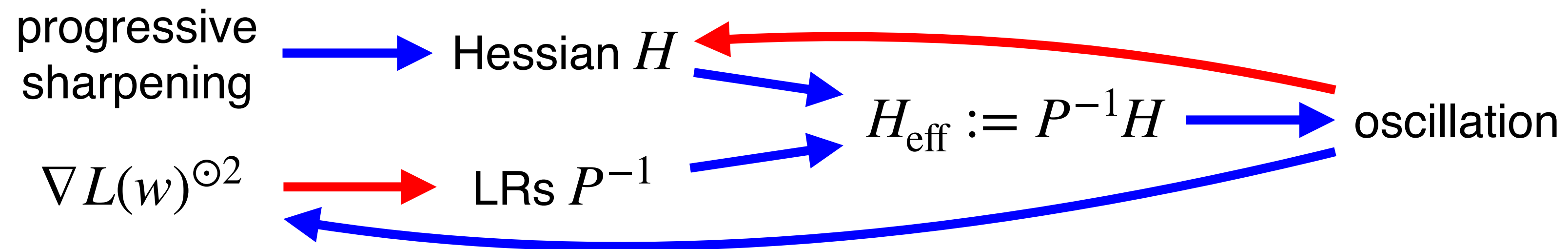


How can we analyze this system?



How can we analyze this system?

- ▶ The fine-grained dynamics are super complicated... 🤔



How can we analyze this system?

- ▶ The fine-grained dynamics are super complicated... 🤔
- ▶ Luckily, the central flows time-averaging argument generalizes easily!

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(w) \right]$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w} + \delta) \right]$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \nabla^2 L(\bar{w}) \delta + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \delta \delta^T \rangle \right] + \mathcal{O}(\delta^3)$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \nabla^2 L(\bar{w}) \overset{0}{\delta} + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \delta \delta^T \rangle \right] + \mathcal{O}(\delta^3)$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\underbrace{\nabla L(\bar{w})}_{\text{gradient}} + \underbrace{\frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle}_{\text{curvature penalty}} \right]$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\underbrace{\nabla L(\bar{w})}_{\text{gradient}} + \underbrace{\frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle}_{\text{curvature penalty}} \right]$$

- ▶ We assume $\nu(t)$ follows the average **squared** gradient:

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\mathbb{E} \left[\nabla L(w)^{\odot 2} \right] - \nu \right]$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\underbrace{\nabla L(\bar{w})}_{\text{gradient}} + \underbrace{\frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle}_{\text{curvature penalty}} \right]$$

- ▶ We assume $\nu(t)$ follows the average **squared** gradient:

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\mathbb{E} \left[\nabla L(\bar{w} + \delta)^{\odot 2} \right] - \nu \right]$$

Sketch: Deriving the RMSProp Central Flow

- ▶ Pretend w_t oscillates around a **central flow** $\bar{w}(t)$ with covariance $\Sigma(t)$:

$$w_t = \bar{w}(t) + \delta_t \quad \text{where} \quad \mathbb{E}[\delta_t \delta_t^T] = \Sigma(t)$$

- ▶ We again assume $\bar{w}(t)$ follows the average gradient:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\underbrace{\nabla L(\bar{w})}_{\text{gradient}} + \underbrace{\frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle}_{\text{curvature penalty}} \right]$$

- ▶ We assume $\nu(t)$ follows the average **squared** gradient:

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w}) \Sigma H(\bar{w})] - \nu \right]$$

Sketch: Deriving the RMSProp Central Flow

- ▶ If we knew $\Sigma(t)$, then \bar{w} , ν would evolve by:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

Sketch: Deriving the RMSProp Central Flow

- ▶ If we knew $\Sigma(t)$, then \bar{w} , ν would evolve by:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

- ▶ We show there is only one choice of $\Sigma(t)$ which satisfies the constraints:
 - 1. Positivity:** As a covariance matrix, $\Sigma(t) \succeq 0$

Sketch: Deriving the RMSProp Central Flow

- ▶ If we knew $\Sigma(t)$, then \bar{w} , ν would evolve by:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

- ▶ We show there is only one choice of $\Sigma(t)$ which satisfies the constraints:
 1. **Positivity:** As a covariance matrix, $\Sigma(t) \succeq 0$
 2. **Stability:** $H(\bar{w}) \preceq 2P(\nu)$ where $P(\nu) = \text{diag}(\sqrt{\nu}/\eta)$

Sketch: Deriving the RMSProp Central Flow

- ▶ If we knew $\Sigma(t)$, then \bar{w} , ν would evolve by:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

- ▶ We show there is only one choice of $\Sigma(t)$ which satisfies the constraints:
 1. **Positivity:** As a covariance matrix, $\Sigma(t) \succeq 0$
 2. **Stability:** $H(\bar{w}) \preceq 2P(\nu)$ where $P(\nu) = \text{diag}(\sqrt{\nu}/\eta)$
 3. **Complementarity:** $\text{span } \Sigma(t) \subseteq \text{span}[\text{unstable evecs}]$

Sketch: Deriving the RMSProp Central Flow

- ▶ If we knew $\Sigma(t)$, then \bar{w} , ν would evolve by:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

- ▶ We show there is only one choice of $\Sigma(t)$ which satisfies the constraints:

$$\underbrace{0 \preceq 2P(\nu) - H(\bar{w})}_{\text{Stability}} \perp \underbrace{\Sigma \succeq 0}_{\text{Positivity}} \quad \text{for all } t$$

Complementarity

Sketch: Deriving the RMSProp Central Flow

- ▶ If we knew $\Sigma(t)$, then \bar{w} , ν would evolve by:

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

- ▶ We show there is only one choice of $\Sigma(t)$ which satisfies the constraints:

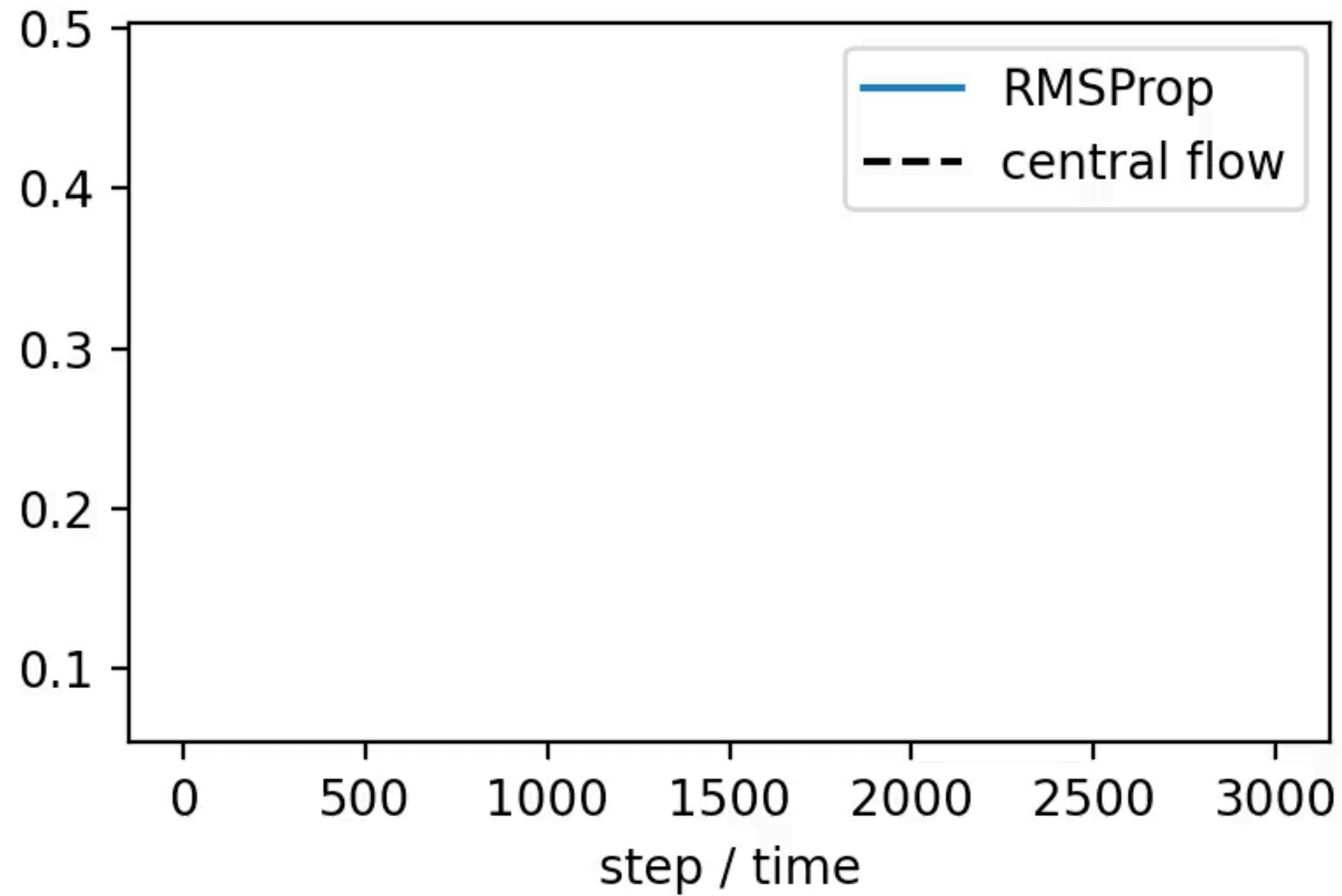
$$\underbrace{0 \leq 2P(\nu) - H(\bar{w})}_{\text{Stability}} \perp \underbrace{\Sigma \geq 0}_{\text{Positivity}} \quad \text{for all } t$$

Complementarity

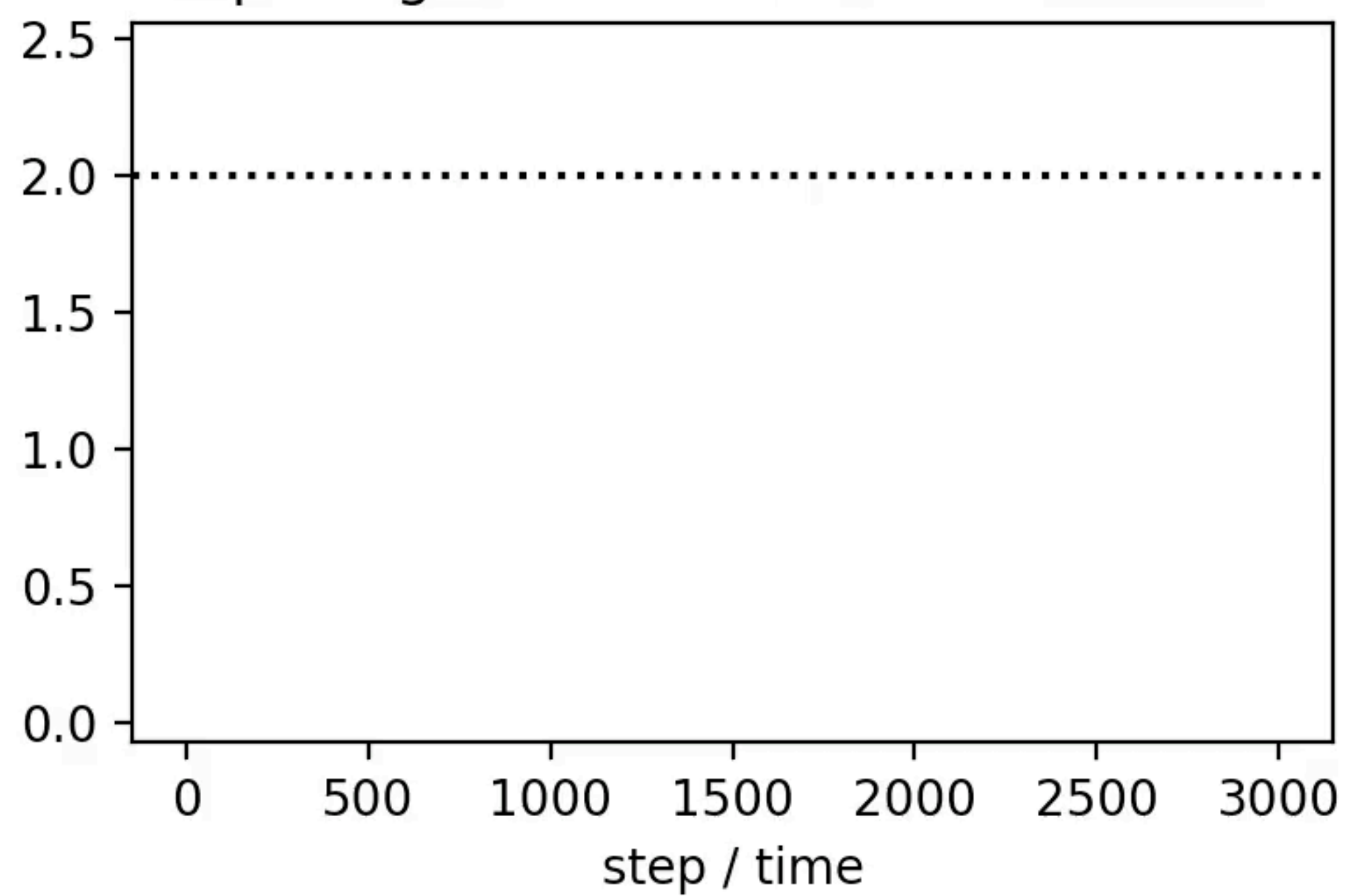
- ▶ This defines a **differential complementarity problem** (DCP) for (\bar{w}, ν, Σ) which uniquely determines the central flow $\{(\bar{w}(t), \nu(t), \Sigma(t))\}_{t \geq 0}$ and can be efficiently simulated

Validating The RMSProp Central Flow

train loss

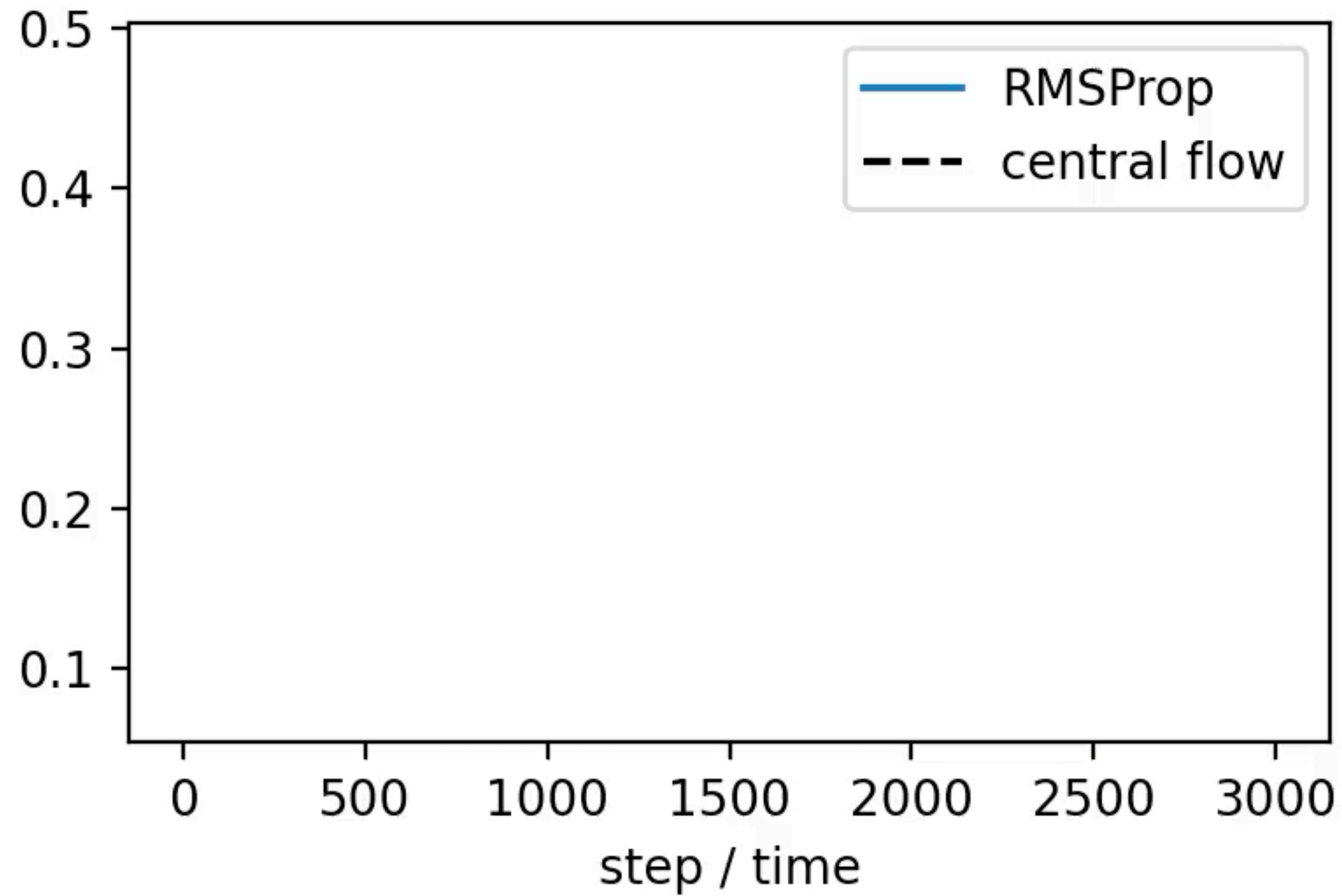


top 4 eigenvalues of effective Hessian

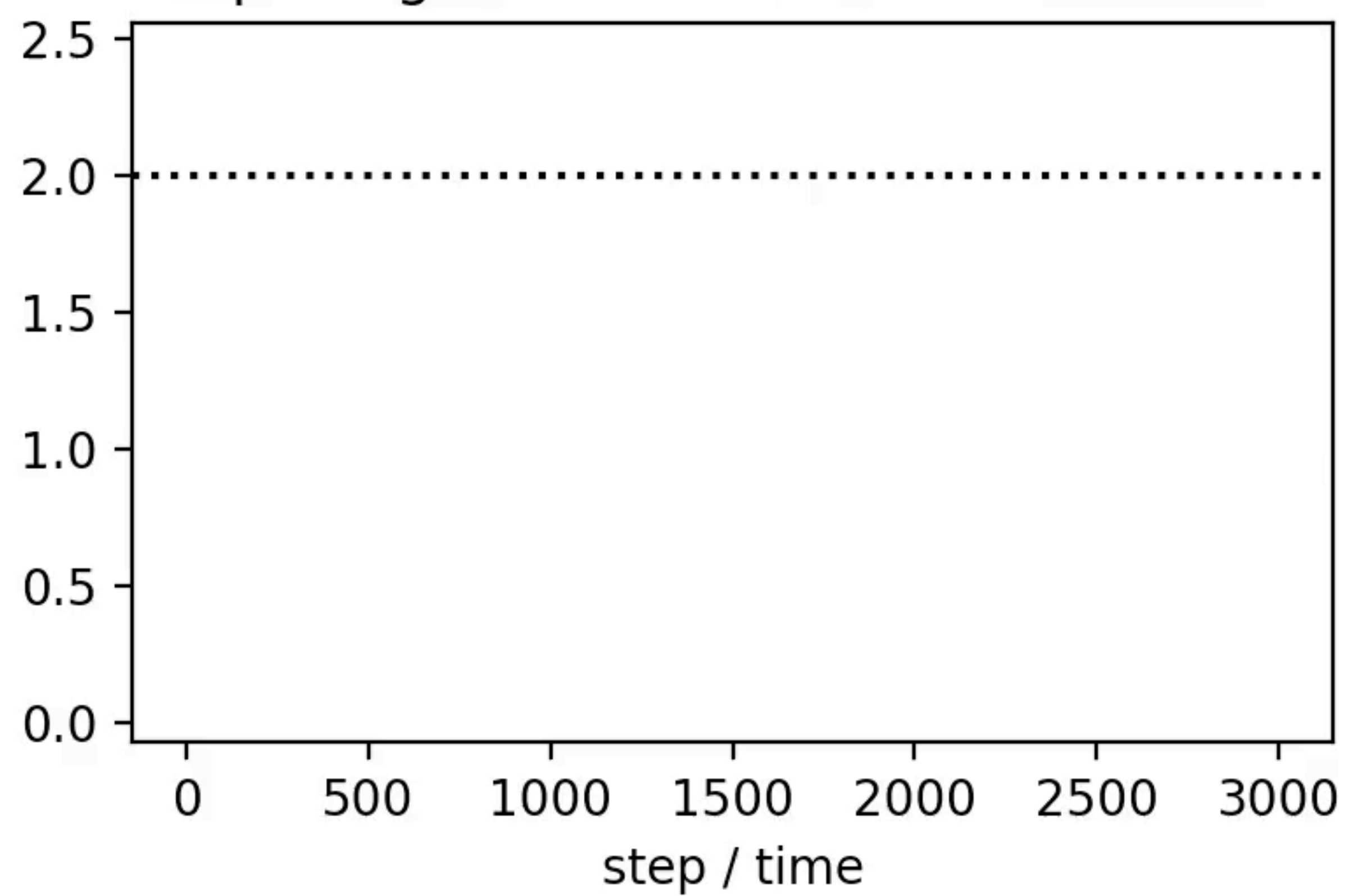


Validating The RMSProp Central Flow

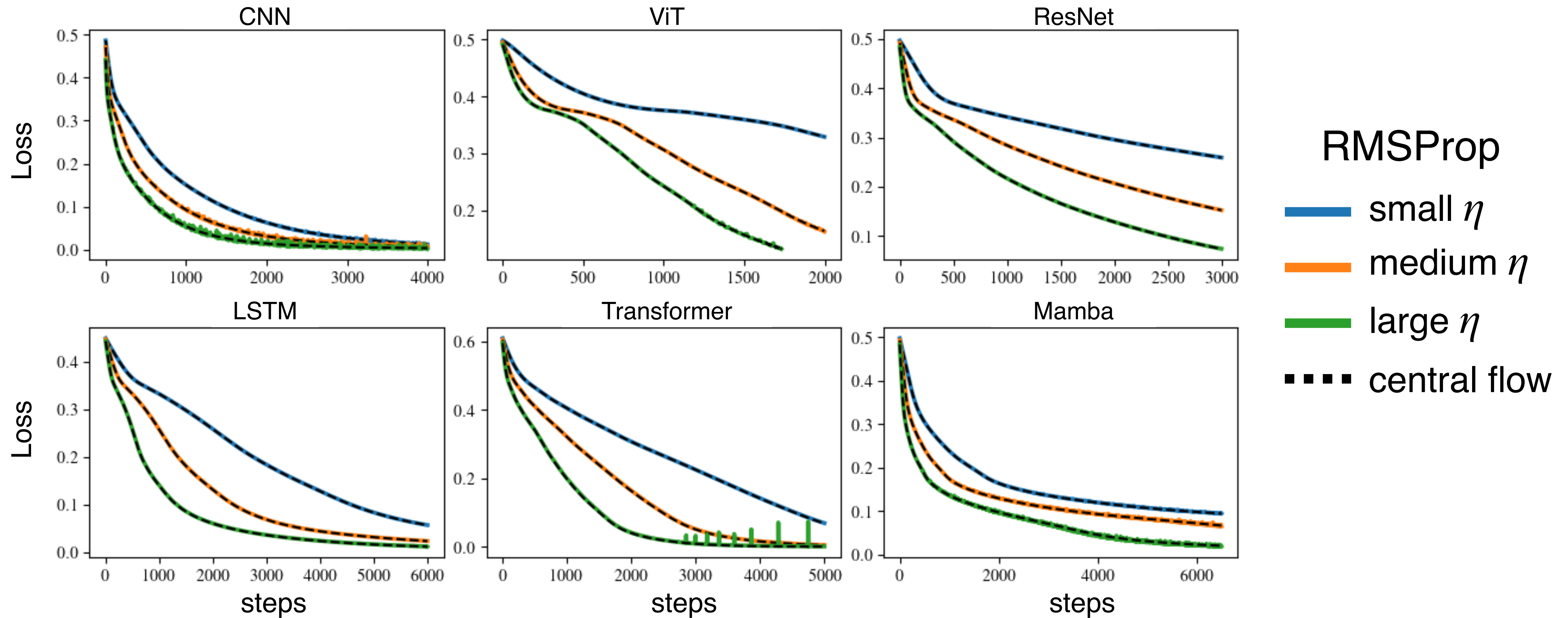
train loss



top 4 eigenvalues of effective Hessian



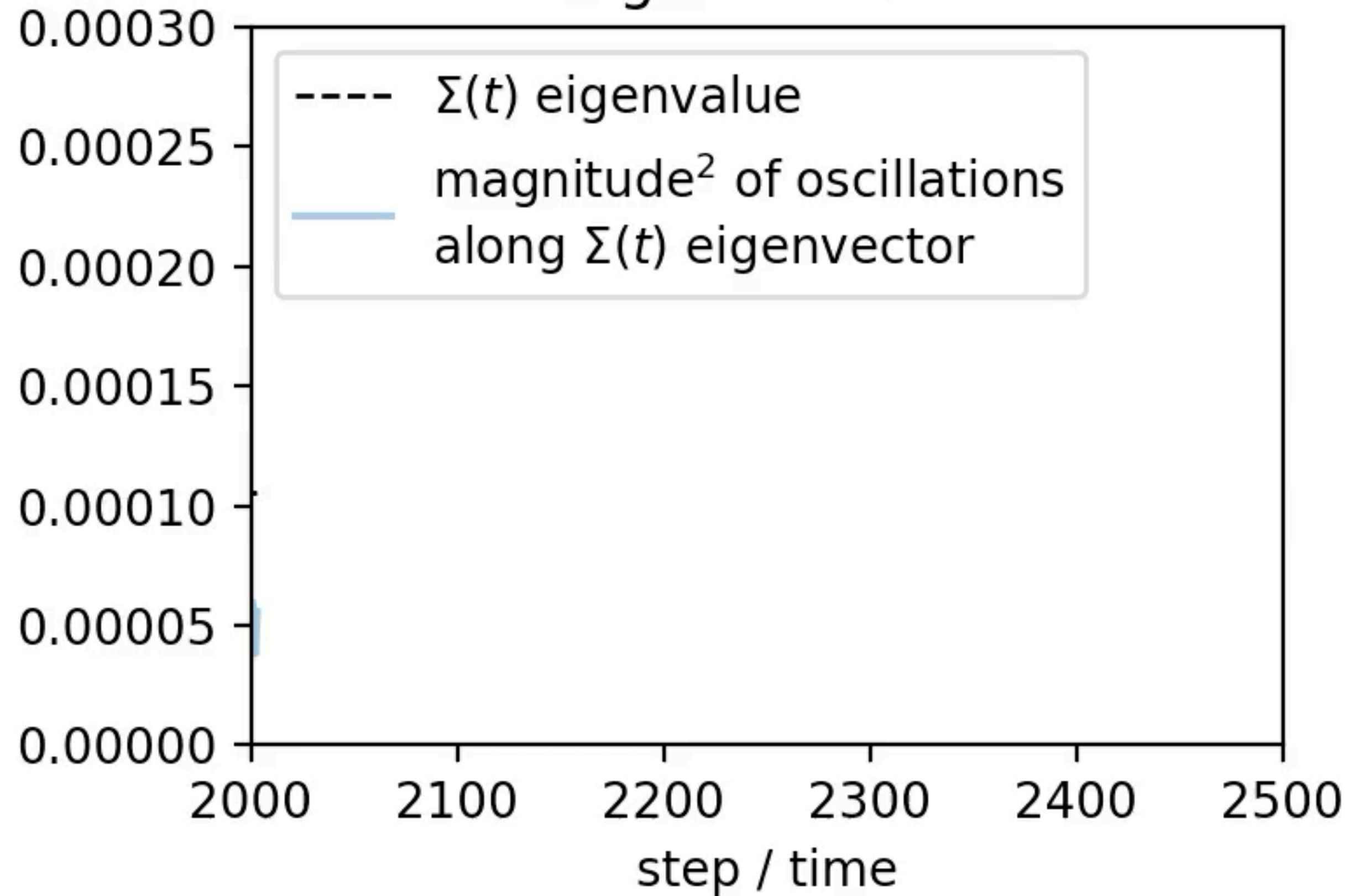
The RMSProp Central Flow Holds Across Architectures



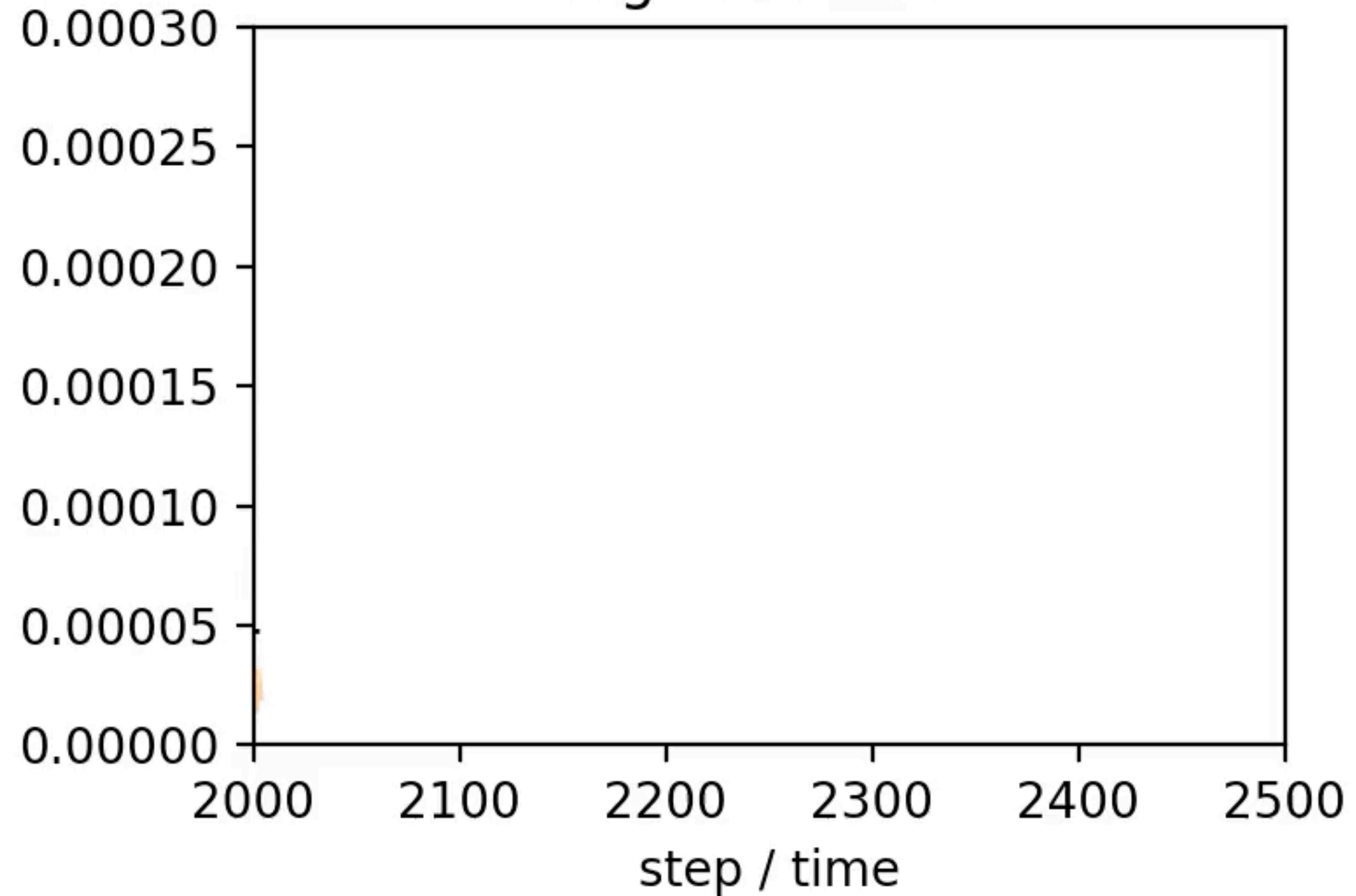
Strong agreement across architectures!

The RMSProp Central Flow Predicts Σ

eigenvector 1

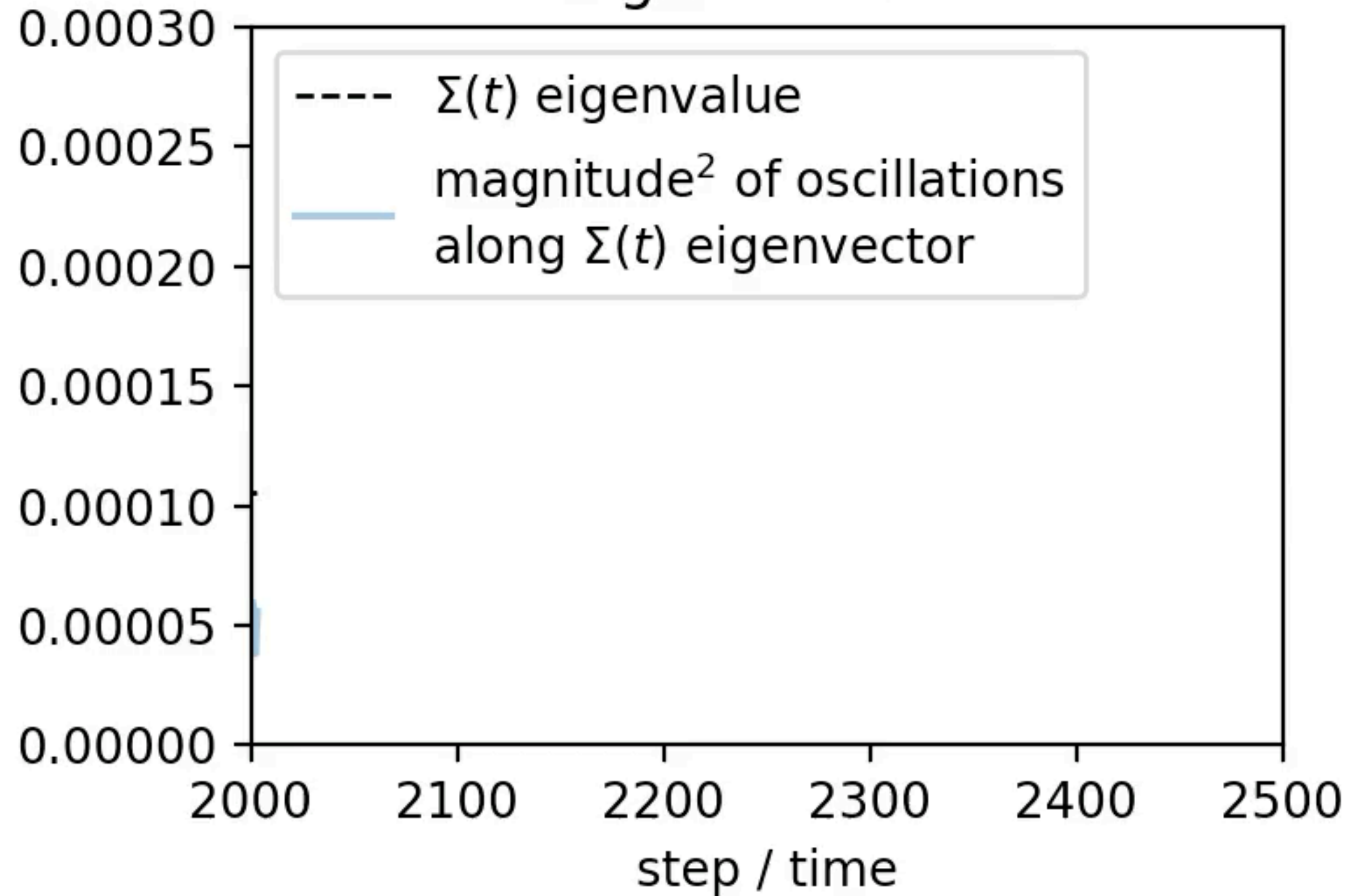


eigenvector 2

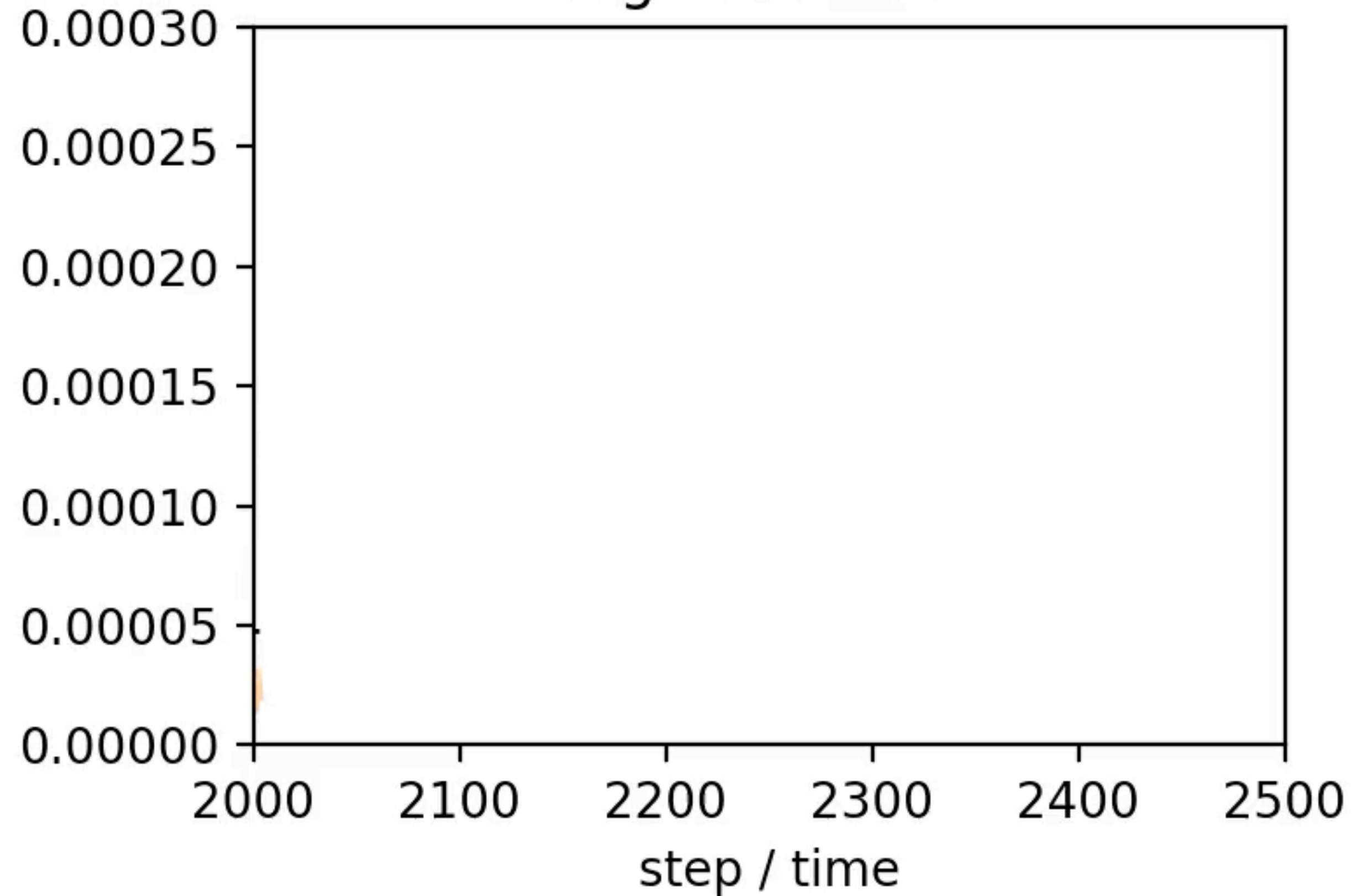


The RMSProp Central Flow Predicts Σ

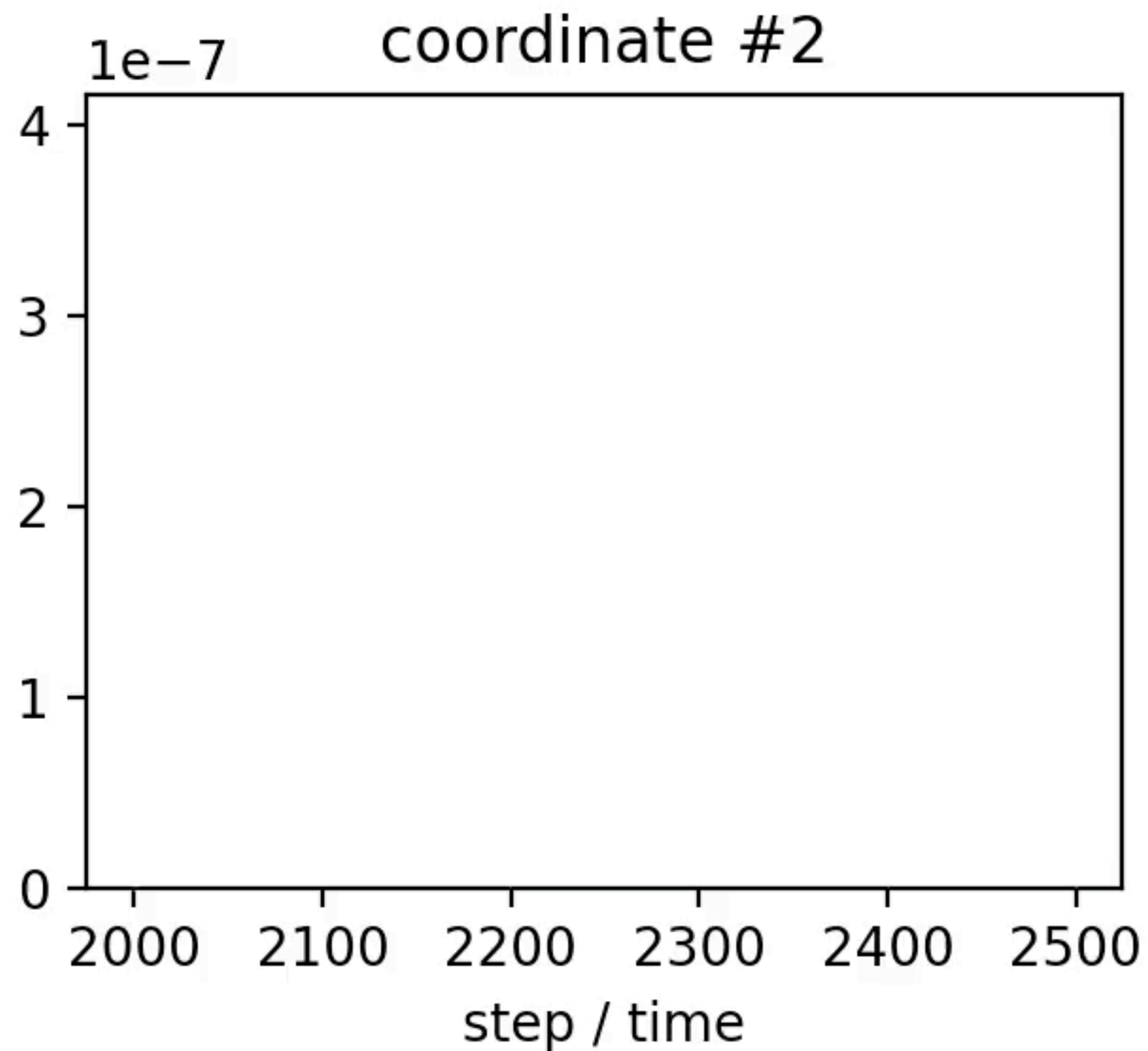
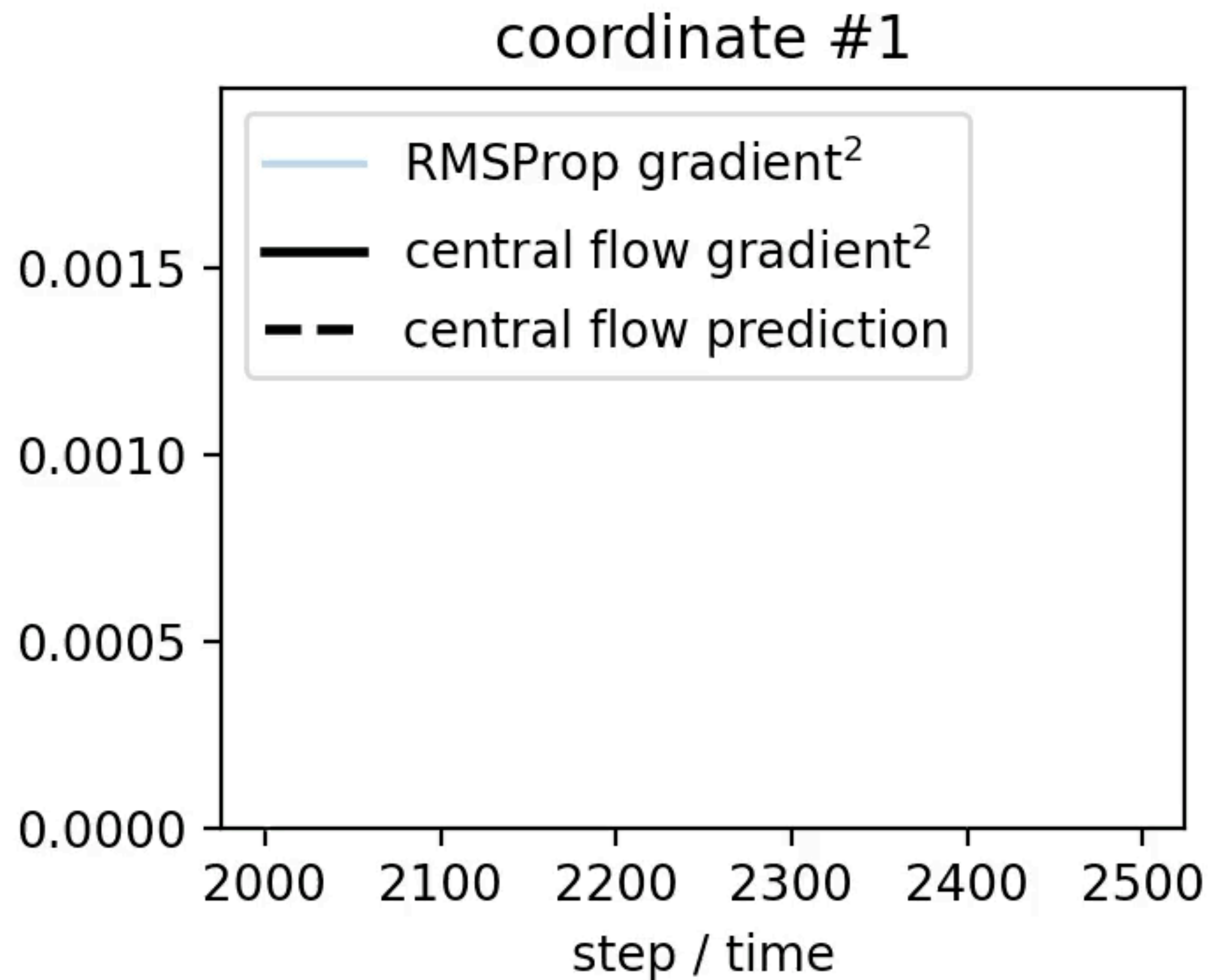
eigenvector 1



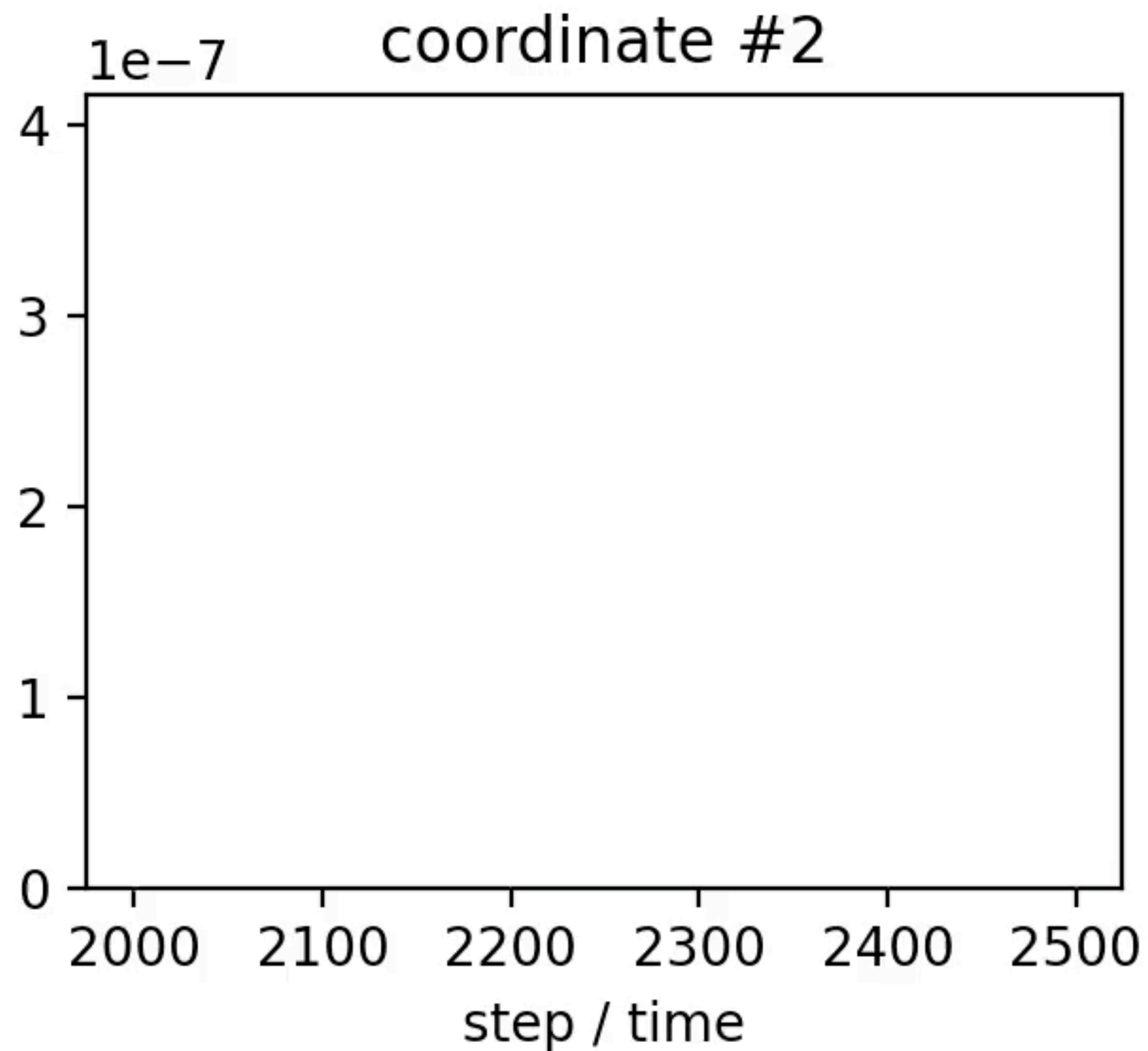
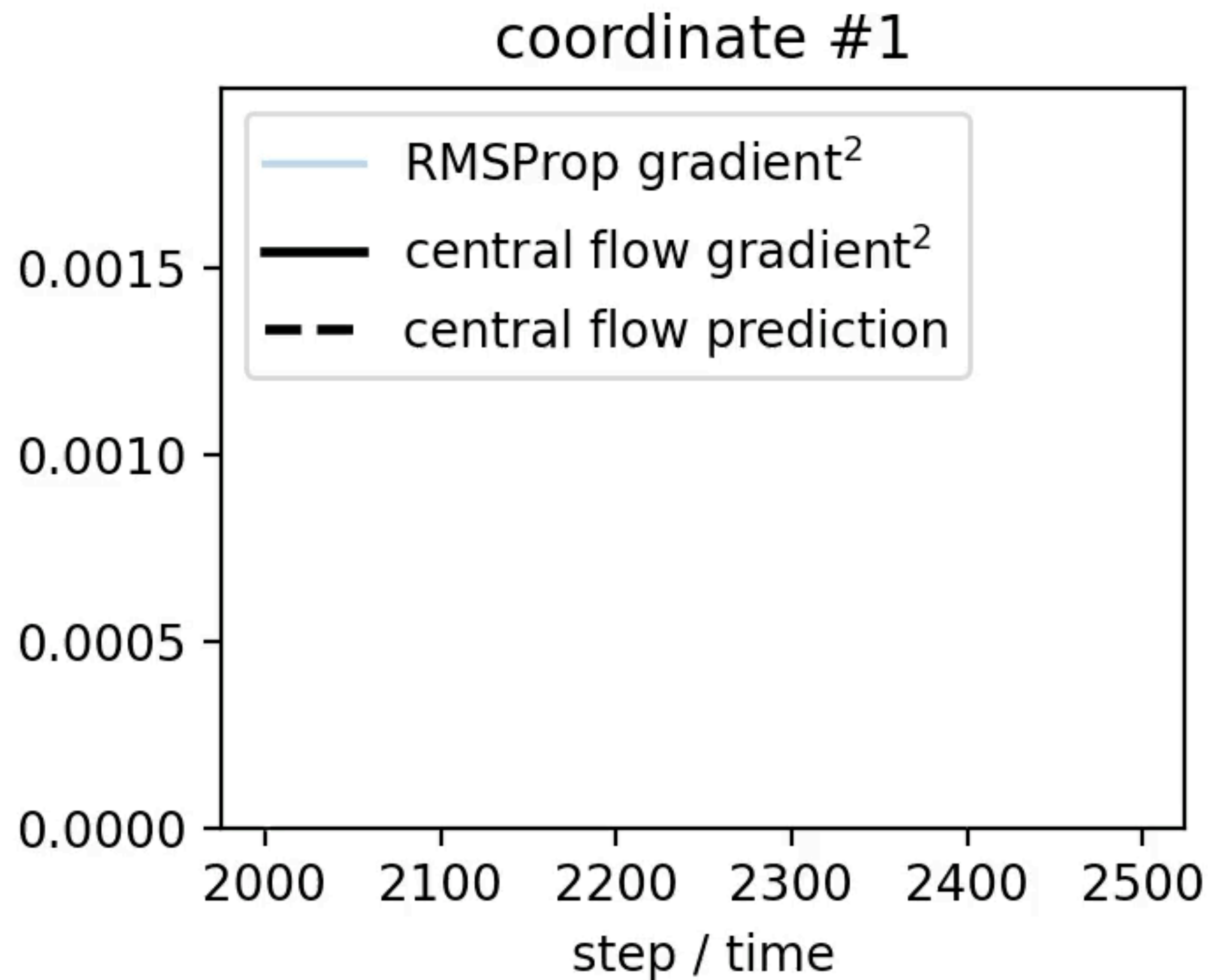
eigenvector 2



The RMSProp Central Flow Predicts Gradient Norms



The RMSProp Central Flow Predicts Gradient Norms



Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -\frac{\eta}{\sqrt{\nu}} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma \rangle \right]$$

$$\frac{d\nu}{dt} = \frac{1 - \beta_2}{\beta_2} \left[\nabla L(\bar{w})^{\odot 2} + \text{diag}[H(\bar{w})\Sigma H(\bar{w})] - \nu \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$
- ▶ Gives a stationary preconditioner $P^*(\bar{w}) := \text{diag}[\nu^*(\bar{w})^{1/2}/\eta]$

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma^*(\bar{w}) \rangle \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$
- ▶ Gives a stationary preconditioner $P^*(\bar{w}) := \text{diag}[\nu^*(\bar{w})^{1/2}/\eta]$

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma^*(\bar{w}) \rangle \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$
- ▶ Gives a stationary preconditioner $P^*(\bar{w}) := \text{diag}[\nu^*(\bar{w})^{1/2}/\eta]$

Lemma: The stationary preconditioner $P^*(\bar{w})$ solves the following SDP:

$$P^*(w) = \underset{P \text{ diagonal}}{\text{argmin}} \quad \text{tr}(P) + \frac{1}{\eta^2} \|\nabla L(w)\|_{P^{-1}} \quad \text{such that} \quad H \preceq 2P$$

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma^*(\bar{w}) \rangle \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$
- ▶ Gives a stationary preconditioner $P^*(\bar{w}) := \text{diag}[\nu^*(\bar{w})^{1/2}/\eta]$

Lemma: The stationary preconditioner $P^*(\bar{w})$ solves the following SDP:

$$P^*(w) = \underset{P \text{ diagonal}}{\text{argmin}} \quad \text{tr}(P) + \frac{1}{\eta^2} \|\nabla L(w)\|_{P^{-1}} \quad \text{such that} \quad H \preceq 2P$$

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma^*(\bar{w}) \rangle \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$
- ▶ Gives a stationary preconditioner $P^*(\bar{w}) := \text{diag}[\nu^*(\bar{w})^{1/2}/\eta]$

Lemma: The stationary preconditioner $P^*(\bar{w})$ solves the following SDP:

$$P^*(w) = \underset{P \text{ diagonal}}{\text{argmin}} \quad \text{tr}(P) + \frac{1}{\eta^2} \|\nabla L(w)\|_{P^{-1}} \quad \text{such that} \quad H \preceq 2P$$

local stability

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma^*(\bar{w}) \rangle \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$
- ▶ Gives a stationary preconditioner $P^*(\bar{w}) := \text{diag}[\nu^*(\bar{w})^{1/2}/\eta]$

Lemma: The stationary preconditioner $P^*(\bar{w})$ solves the following SDP:

$$P^*(w) = \underset{P \text{ diagonal}}{\text{argmin}} \text{tr}(P) + \frac{1}{\eta^2} \|\nabla L(w)\|_{P^{-1}} \quad \text{such that} \quad H \preceq 2P$$

max *harmonic mean* of the learning rates

local stability

Interpreting the Flow: The Stationary Preconditioner

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \mathbb{E} \left[\nabla L(\bar{w}) + \frac{1}{2} \nabla \langle \nabla^2 L(\bar{w}), \Sigma^*(\bar{w}) \rangle \right]$$

- ▶ Freeze \bar{w} and evolve (ν, Σ) until they reach stationary values $\nu^*(\bar{w}), \Sigma^*(\bar{w})$
- ▶ Gives a stationary preconditioner $P^*(\bar{w}) := \text{diag}[\nu^*(\bar{w})^{1/2}/\eta]$

Lemma: The stationary preconditioner $P^*(\bar{w})$ solves the following SDP:

$$P^*(w) = \underset{P \text{ diagonal}}{\text{argmin}} \quad \text{tr}(P) + \frac{1}{\eta^2} \|\nabla L(w)\|_{P^{-1}} \quad \text{such that} \quad H \preceq 2P$$

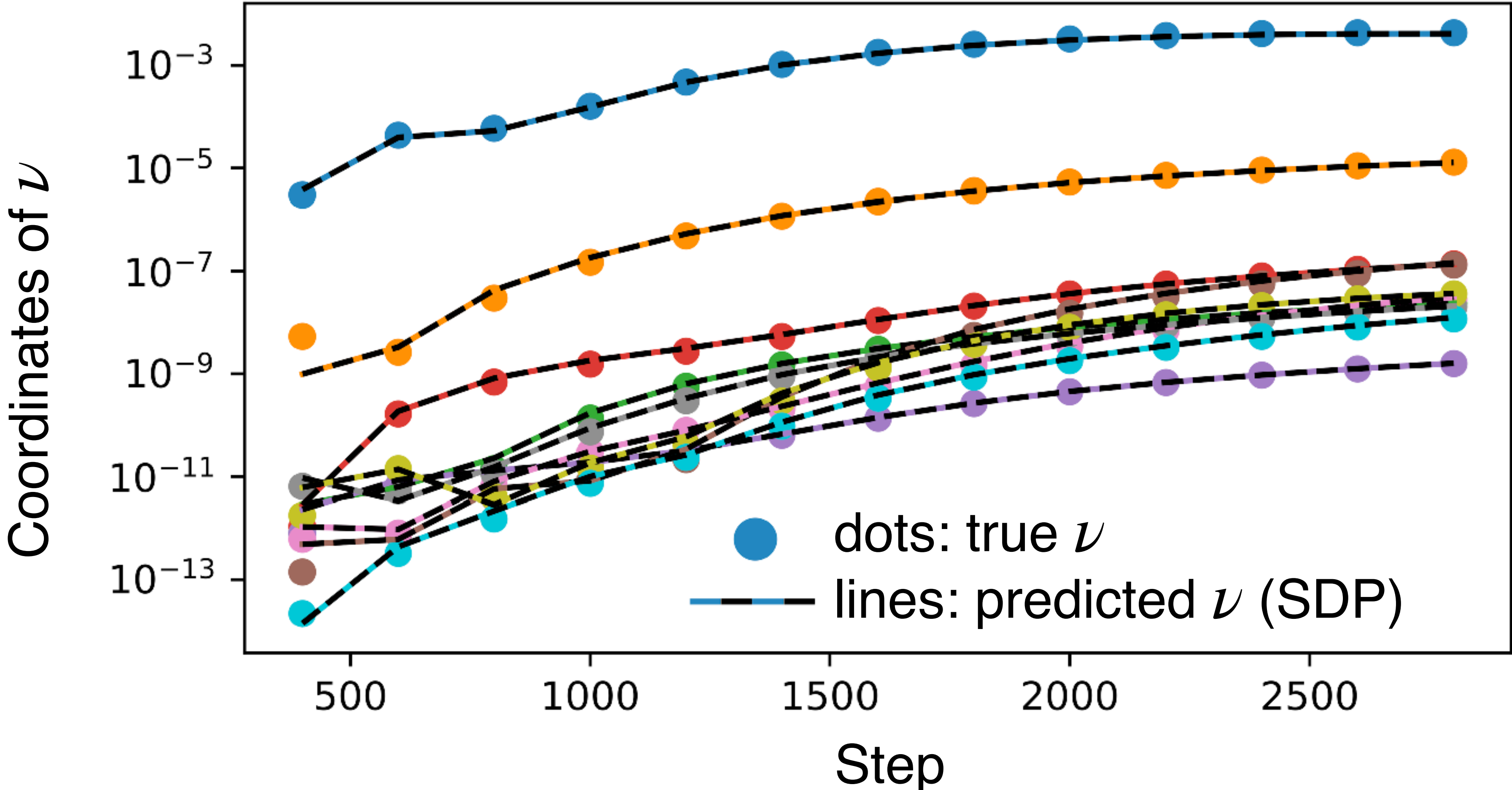
max *harmonic mean* of the learning rates

optimization speed

local stability

Interpreting the Flow: The Stationary Preconditioner

True and Predicted Coordinates of ν



The Stationary RMSProp Central Flow*

$$\frac{d\bar{w}}{dt} = - P^{\star}(\bar{w})^{-1} \left[\nabla L(\bar{w}) + \frac{\eta^2}{4} \nabla \text{tr} P^{\star}(\bar{w}) \right]$$

stationary preconditioner



* this is a slightly simplified form which gives clearer intuition

The Stationary RMSProp Central Flow*

$$\frac{d\bar{w}}{dt} = - P^{\star}(\bar{w})^{-1} \left[\nabla L(\bar{w}) + \frac{\eta^2}{4} \nabla \text{tr} P^{\star}(\bar{w}) \right]$$

stationary preconditioner

curvature penalty

* this is a slightly simplified form which gives clearer intuition

The Stationary RMSProp Central Flow*

$$\frac{d\bar{w}}{dt} = -P^\star(\bar{w})^{-1} \left[\nabla L(\bar{w}) + \frac{\eta^2}{4} \nabla \text{tr} P^\star(\bar{w}) \right]$$

stationary preconditioner

curvature penalty

Acceleration via Regularization:

Use the largest step sizes you can (**adapt**), but avoid regions of the loss landscape where you are forced to take smaller steps (**regularize**).

* this is a slightly simplified form which gives clearer intuition

Acceleration via Regularization:

Use the largest step sizes you can (**adapt**), but avoid regions of the loss landscape where you are forced to take smaller steps (**regularize**).

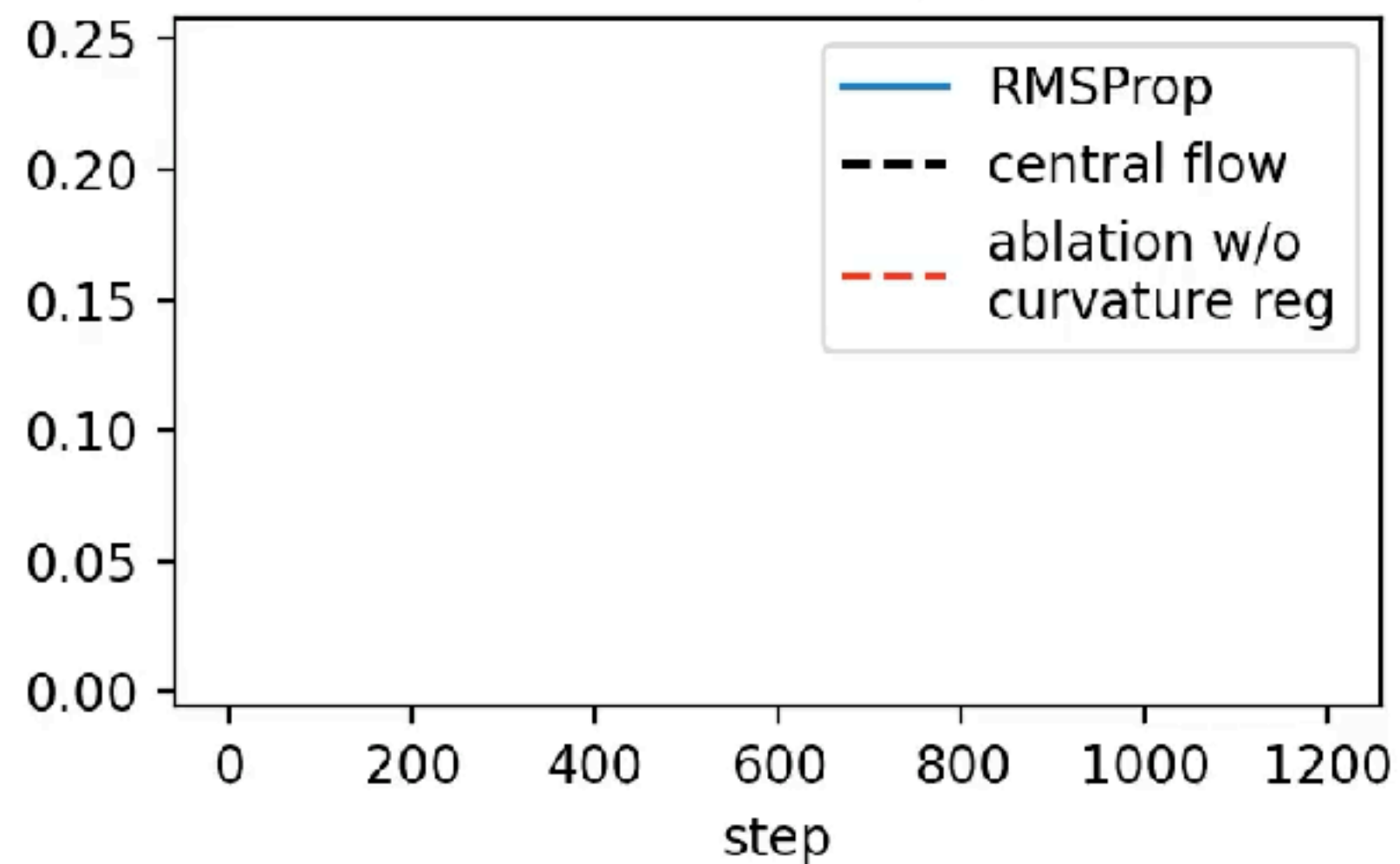


Acceleration via Regularization:

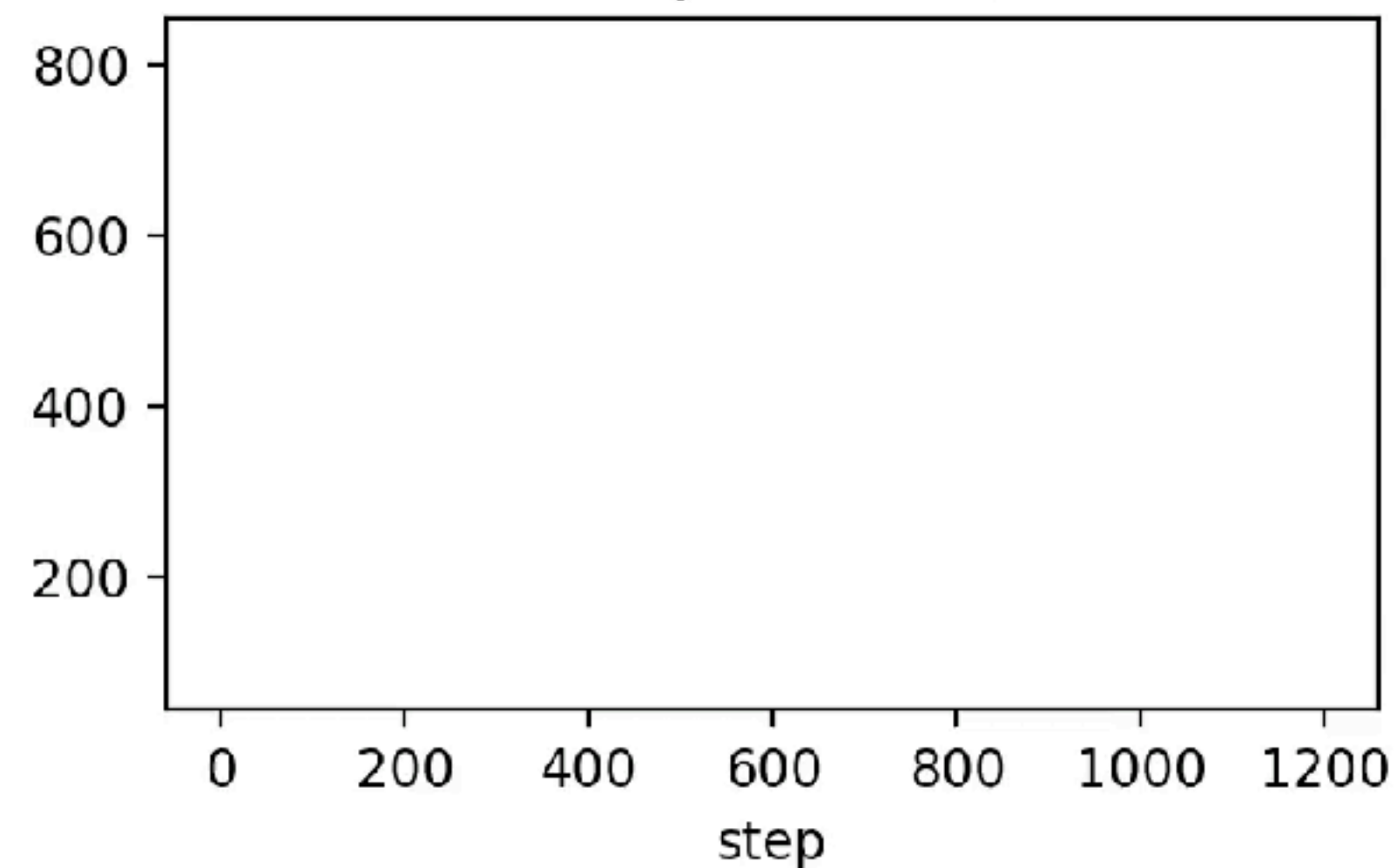
Use the largest step sizes you can (**adapt**), but avoid regions of the loss landscape where you are forced to take smaller steps (**regularize**).

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \left[\nabla L(\bar{w}) + \frac{\eta^2}{4} \nabla \text{tr} P^*(\bar{w}) \right]$$

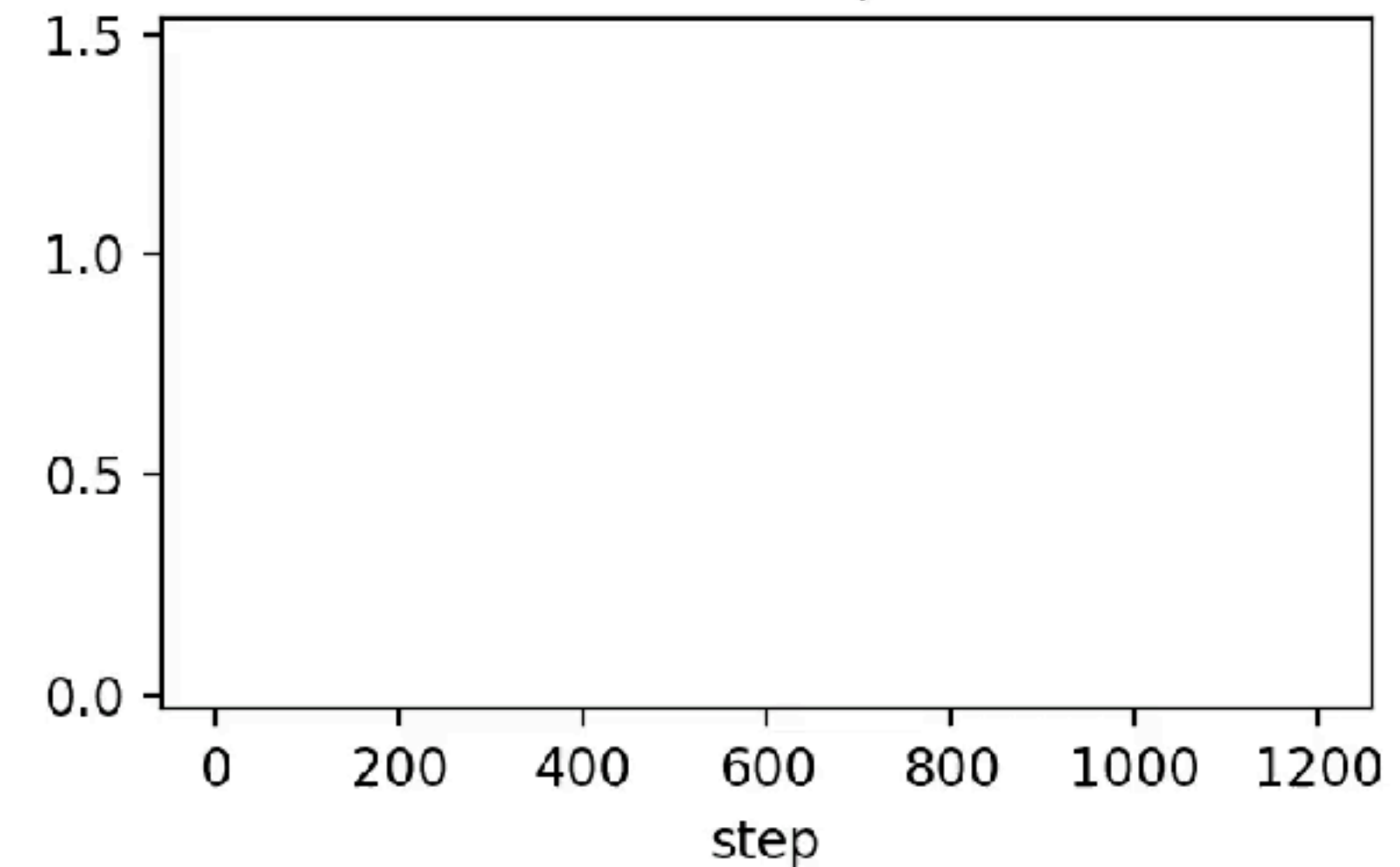
train loss $L(w)$



sharpness $S(w)$



harmonic mean of effective step sizes

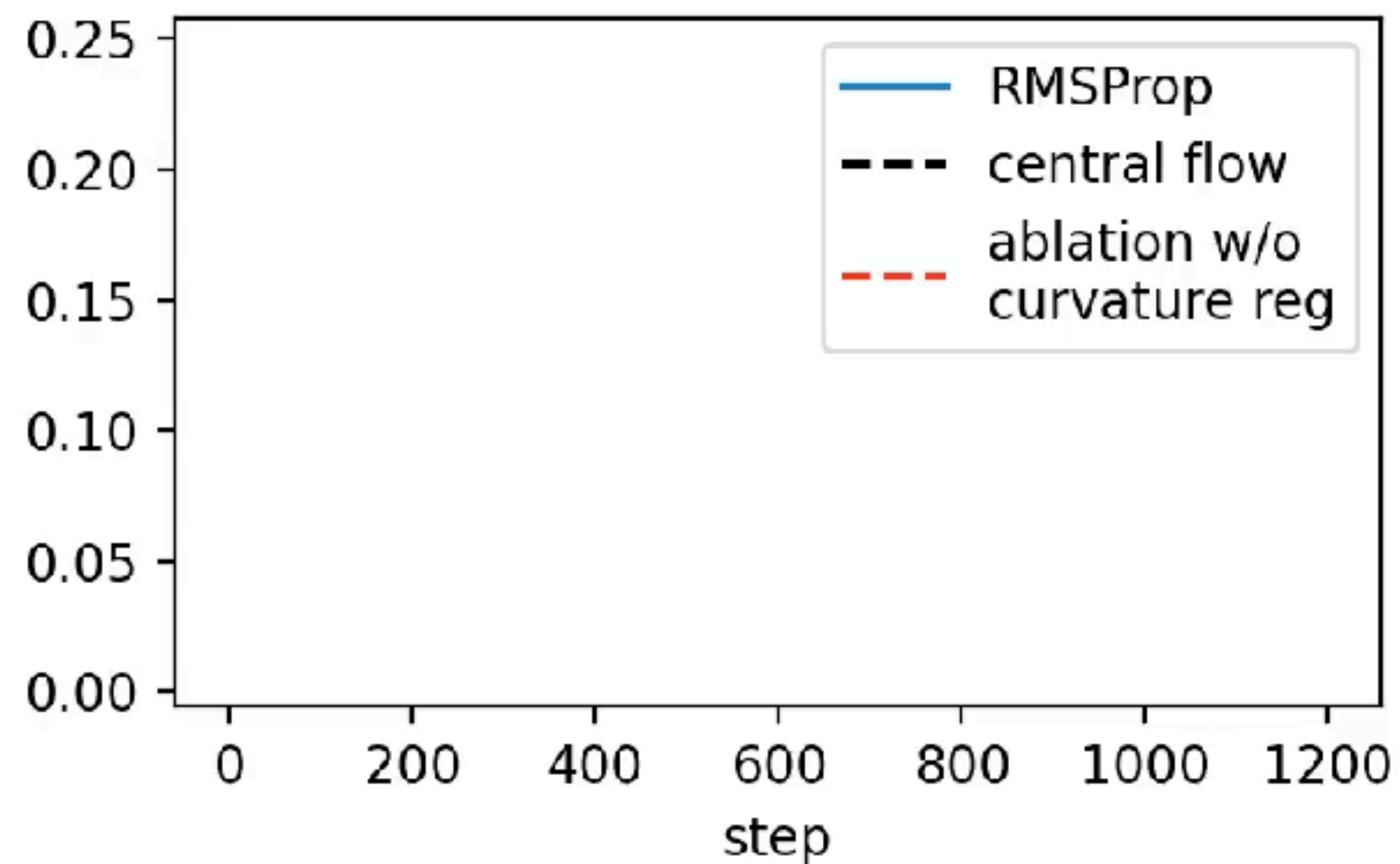


Acceleration via Regularization:

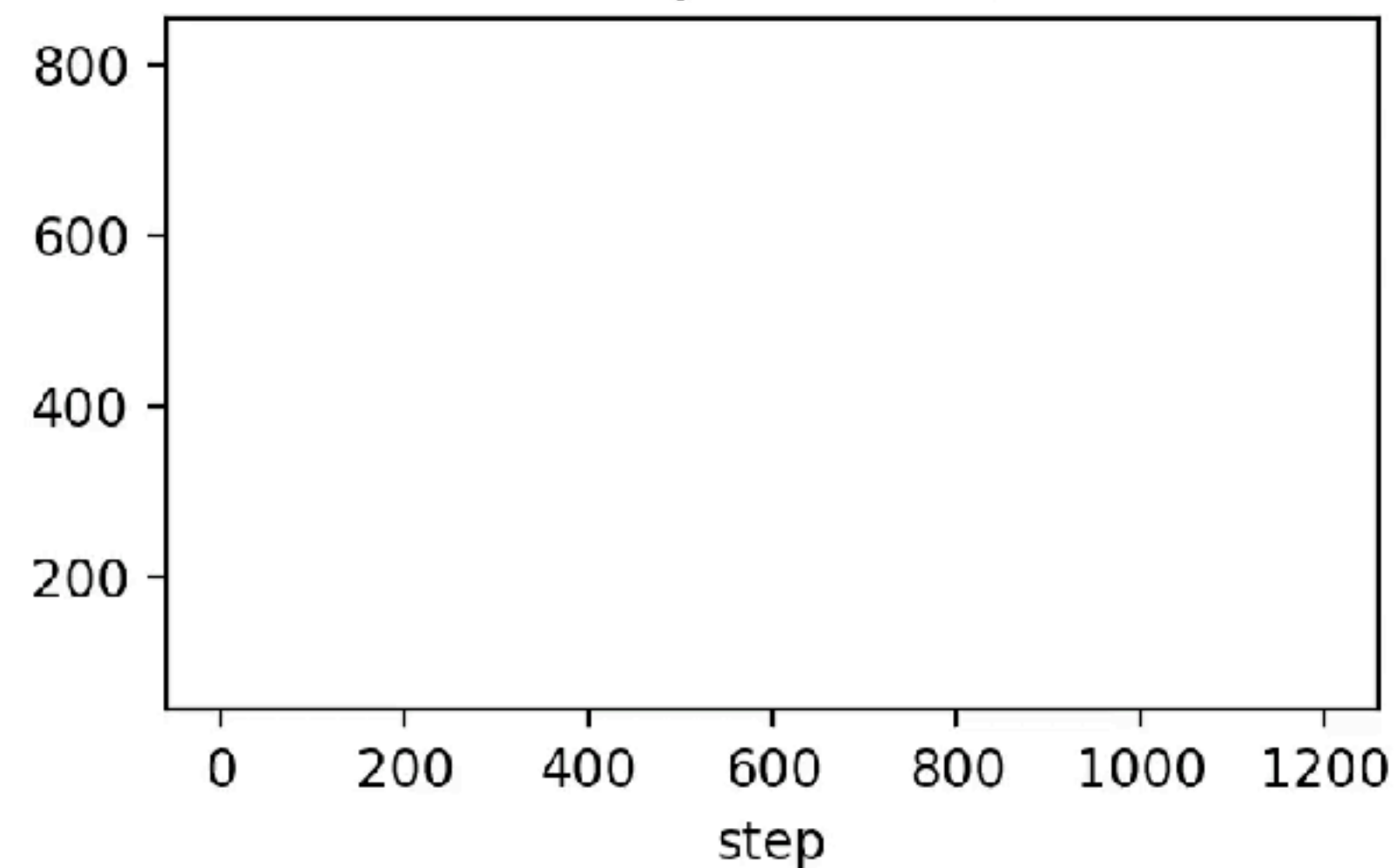
Use the largest step sizes you can (**adapt**), but avoid regions of the loss landscape where you are forced to take smaller steps (**regularize**).

$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \left[\nabla L(\bar{w}) + \frac{\eta^2}{4} \nabla \text{tr} P^*(\bar{w}) \right]$$

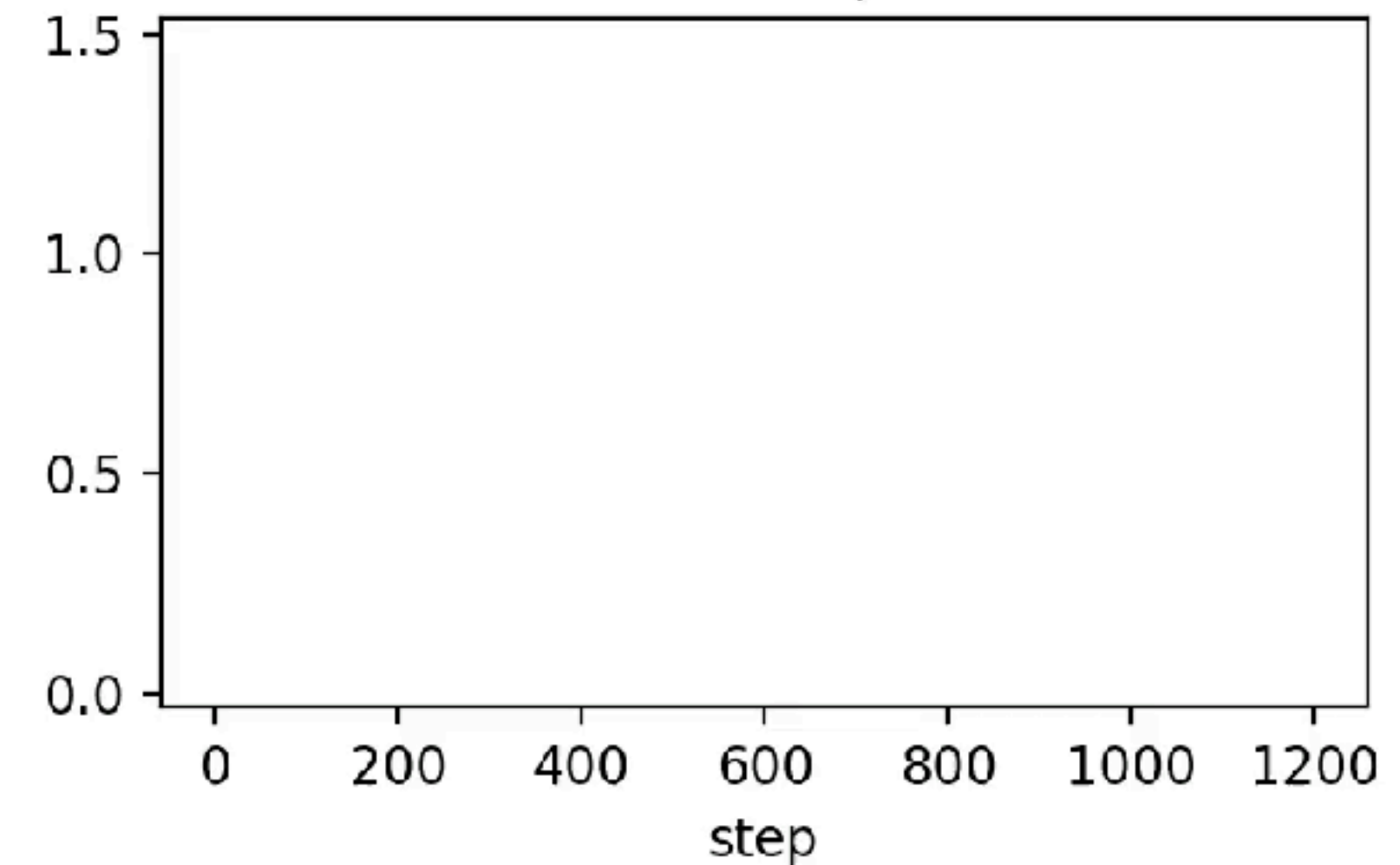
train loss $L(w)$



sharpness $S(w)$



harmonic mean of effective step sizes



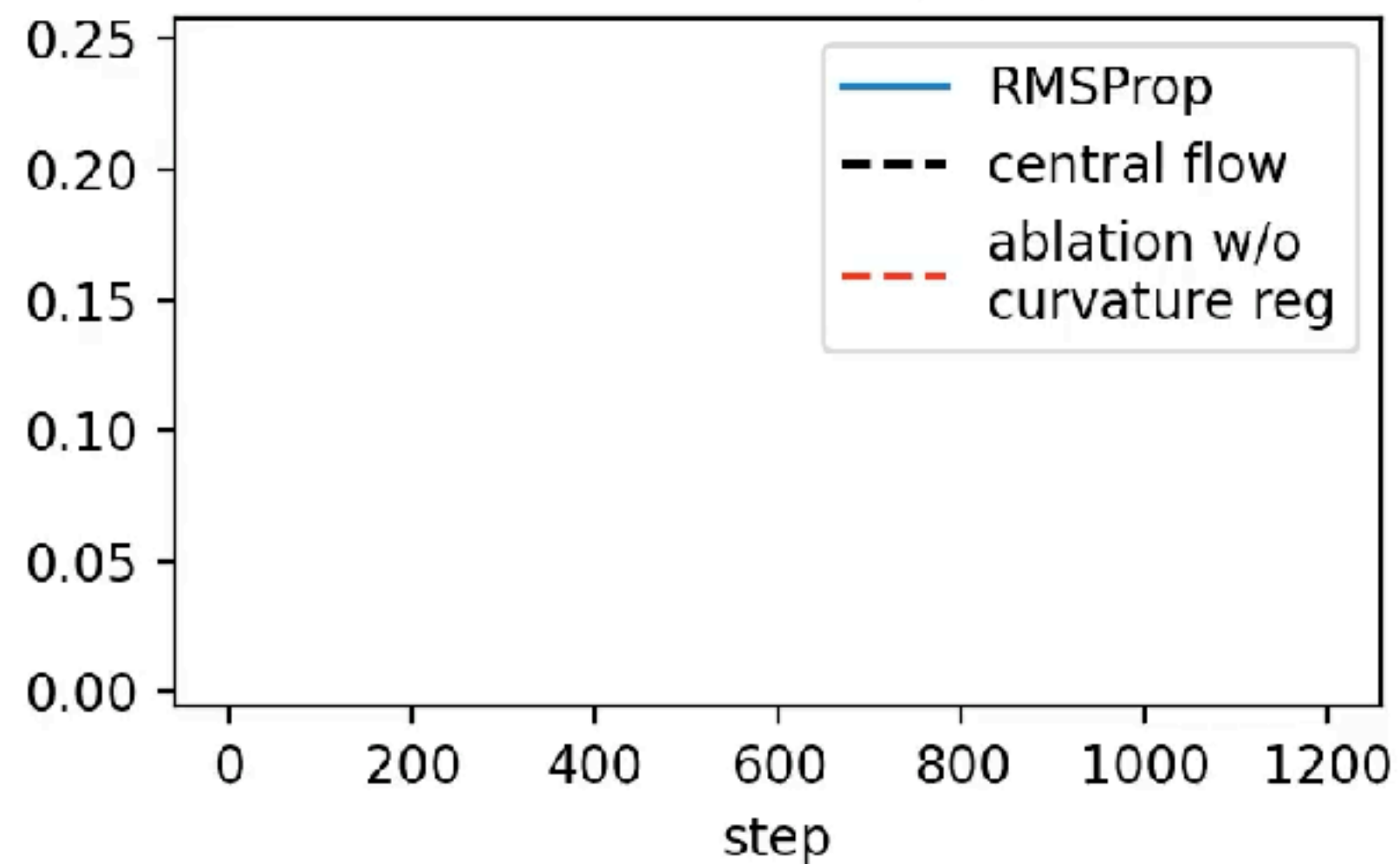
Acceleration via Regularization:

Use the largest step sizes you can (**adapt**), but avoid regions of the loss landscape where you are forced to take smaller steps (**regularize**).

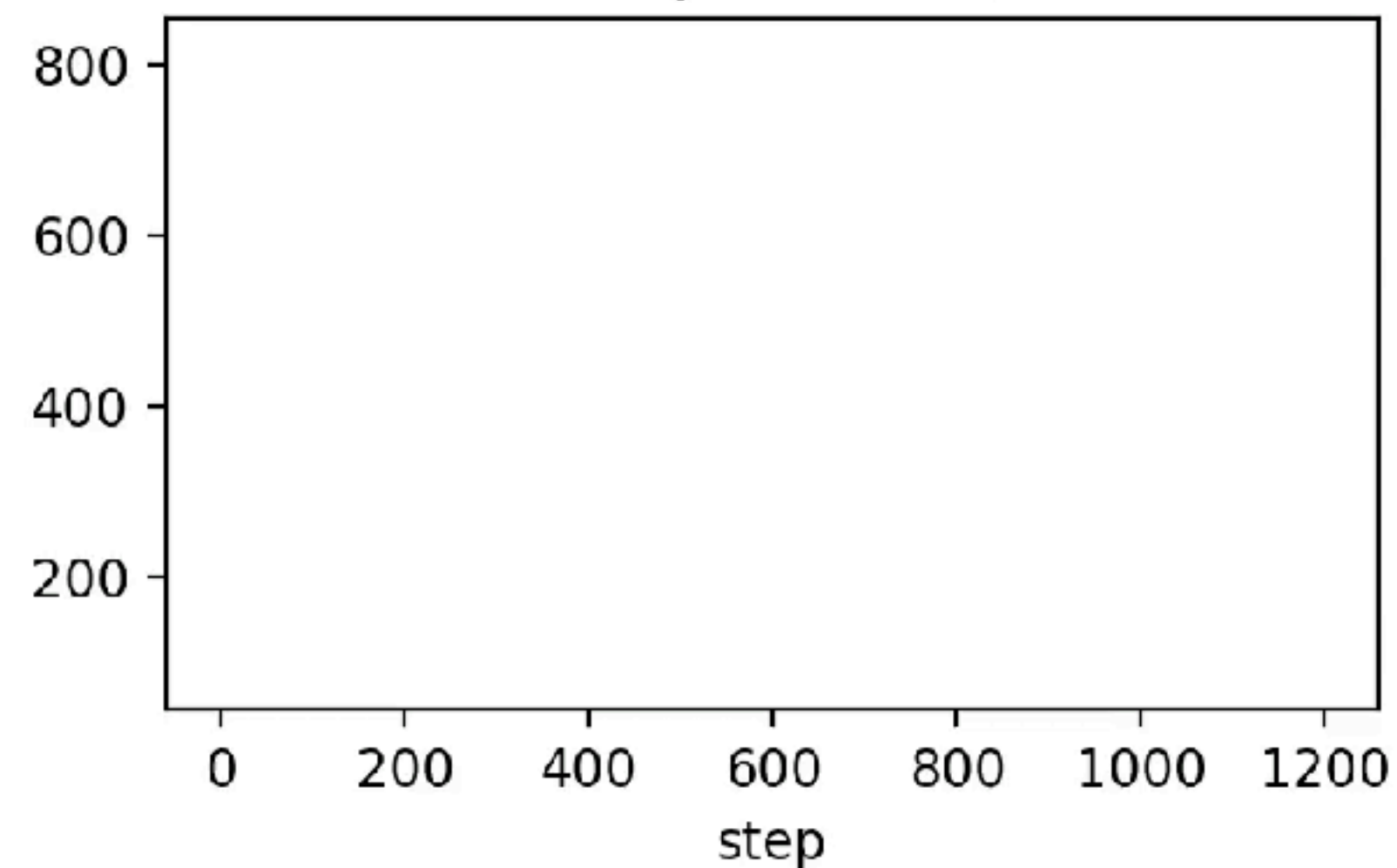
$$\frac{d\bar{w}}{dt} = -P^*(\bar{w})^{-1} \left[\nabla L(\bar{w}) + \frac{\eta^2}{4} \nabla \text{tr} P^*(\bar{w}) \right]$$

controlling this reg. is the main role of η !

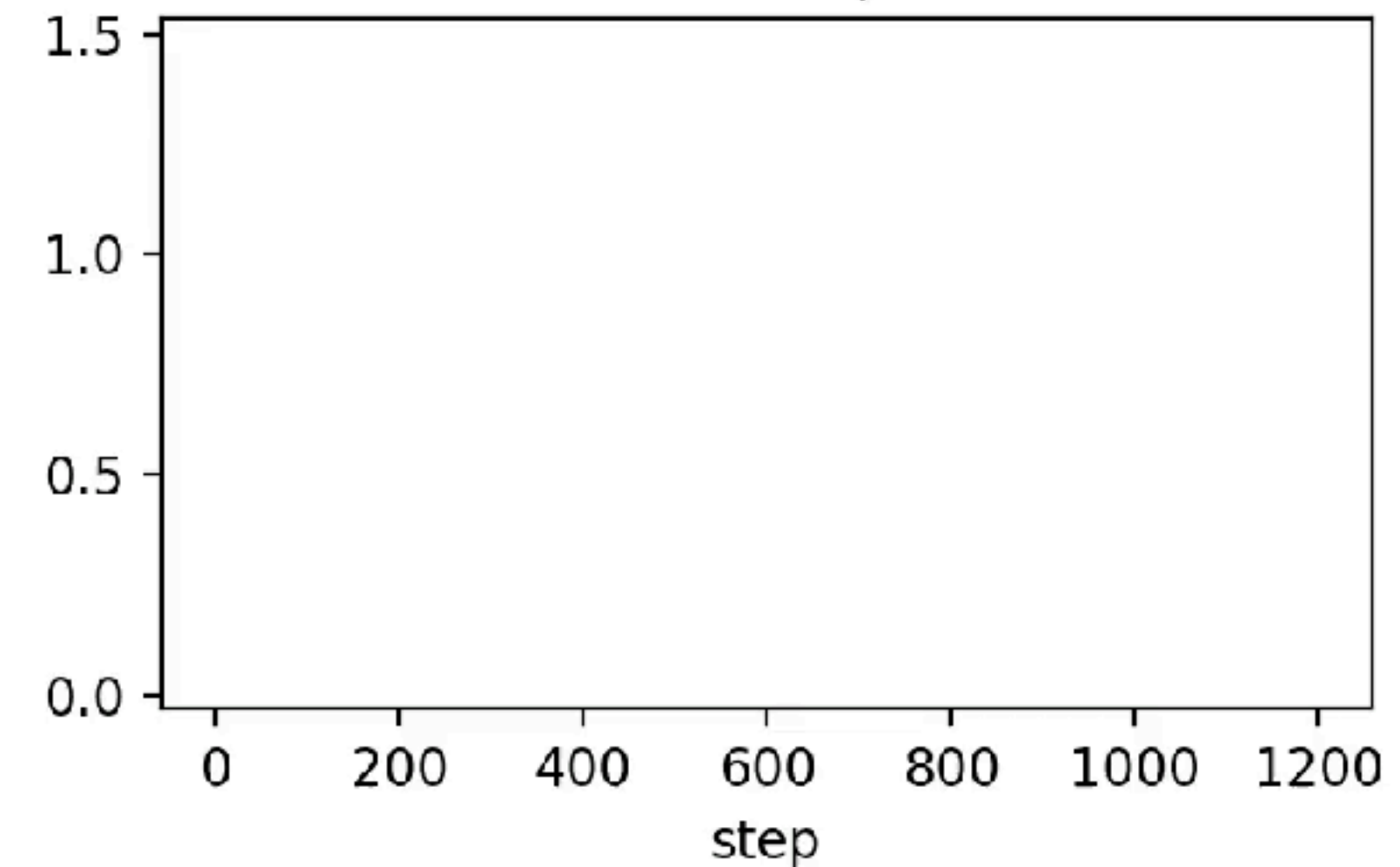
train loss $L(w)$



sharpness $S(w)$



harmonic mean of effective step sizes



Takeaways

Takeaways

- ▶ We derive *central flows*, which model time-averaged trajectories of oscillatory optimizers

Takeaways

- ▶ We derive ***central flows***, which model time-averaged trajectories of oscillatory optimizers
- ▶ We empirically verify these flows predict long-term optimization trajectories

Takeaways

- ▶ We derive ***central flows***, which model time-averaged trajectories of oscillatory optimizers
- ▶ We empirically verify these flows predict long-term optimization trajectories
- ▶ By interpreting these flows, we can easily read off an optimizer's behavior

Takeaways

- ▶ We derive *central flows*, which model time-averaged trajectories of oscillatory optimizers
- ▶ We empirically verify these flows predict long-term optimization trajectories
- ▶ By interpreting these flows, we can easily read off an optimizer's behavior

-
- ▶ **First order** optimizers implicitly implement sophisticated **second-order** strategies

Takeaways

- ▶ We derive *central flows*, which model time-averaged trajectories of oscillatory optimizers
 - ▶ We empirically verify these flows predict long-term optimization trajectories
 - ▶ By interpreting these flows, we can easily read off an optimizer's behavior
-
- ▶ **First order** optimizers implicitly implement sophisticated **second-order** strategies
 - ▶ **Acceleration via regularization:** good optimizers not only adapt to the loss landscape, but also avoid high curvature regions where they are forced to take smaller steps

Takeaways

- ▶ We derive *central flows*, which model time-averaged trajectories of oscillatory optimizers
 - ▶ We empirically verify these flows predict long-term optimization trajectories
 - ▶ By interpreting these flows, we can easily read off an optimizer's behavior
-
- ▶ **First order** optimizers implicitly implement sophisticated **second-order** strategies
 - ▶ **Acceleration via regularization:** good optimizers not only adapt to the loss landscape, but also avoid high curvature regions where they are forced to take smaller steps
 - ▶ RMSProp implicitly solves an SDP related to max-cut to determine its preconditioner

Promising Future Work

Promising Future Work

Rigor:

- ▶ Under what assumptions is the central flow approximation provably correct?
- ▶ What do “reasonable” assumptions for deep learning optimization look like?

Promising Future Work

Rigor:

- ▶ Under what assumptions is the central flow approximation provably correct?
- ▶ What do “reasonable” assumptions for deep learning optimization look like?

Noise:

- ▶ Can similar arguments be used to derive a central flow for SGD?
- ▶ When does an SGD central flow exist?

Promising Future Work

Rigor:

- ▶ Under what assumptions is the central flow approximation provably correct?
- ▶ What do “reasonable” assumptions for deep learning optimization look like?

Noise:

- ▶ Can similar arguments be used to derive a central flow for SGD?
- ▶ When does an SGD central flow exist?

Applications:

- ▶ Can these ideas be used to design new optimizers?
- ▶ How can we design architectures that are easy to optimize?

Thanks!

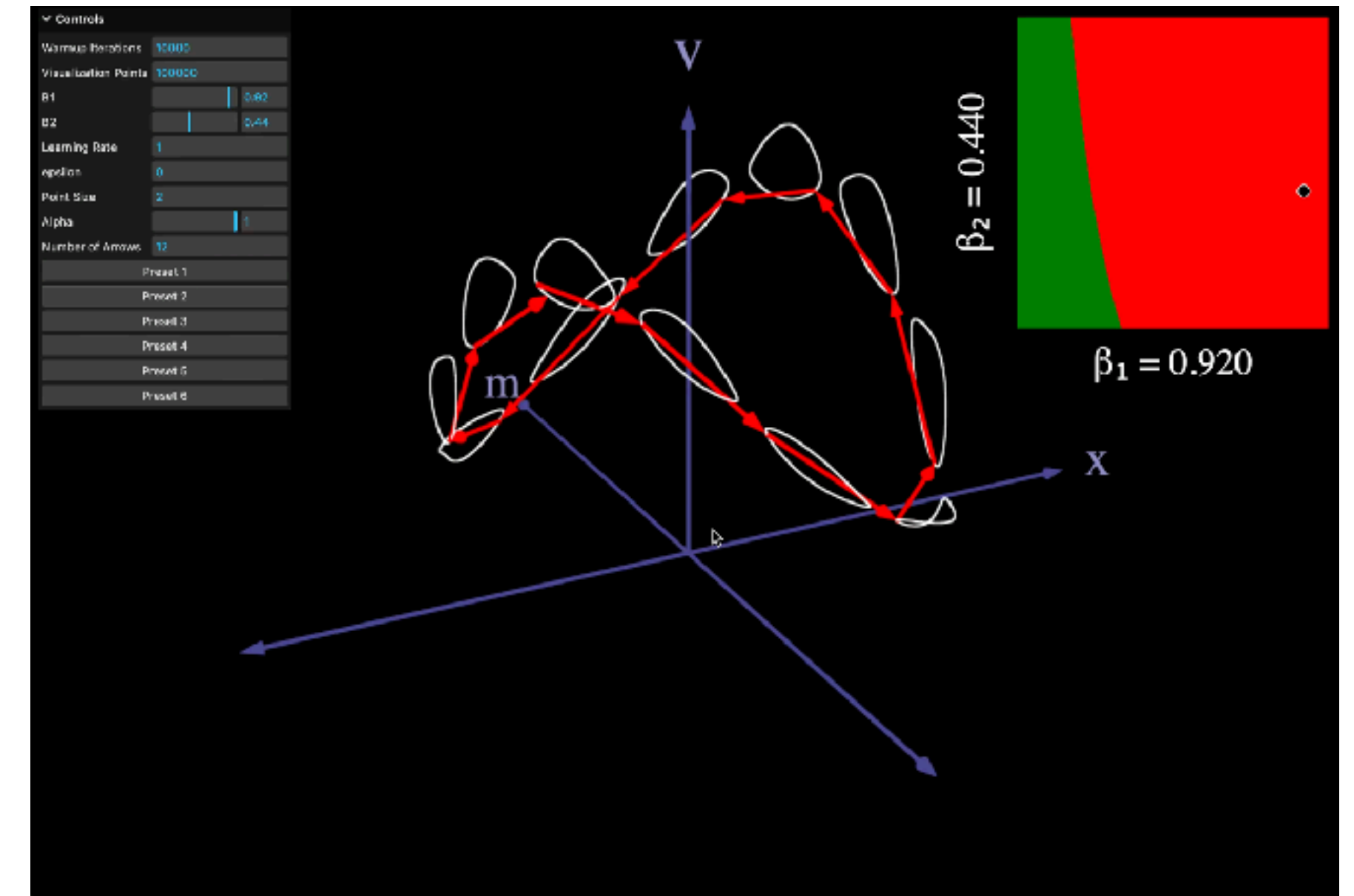


Jeremy Cohen
Flatiron Institute



<https://centralflows.github.io>

Animated Blog Post!



<https://alex-damian.github.io/adam/>

Bonus: crazy limit cycles in Adam

Thanks!

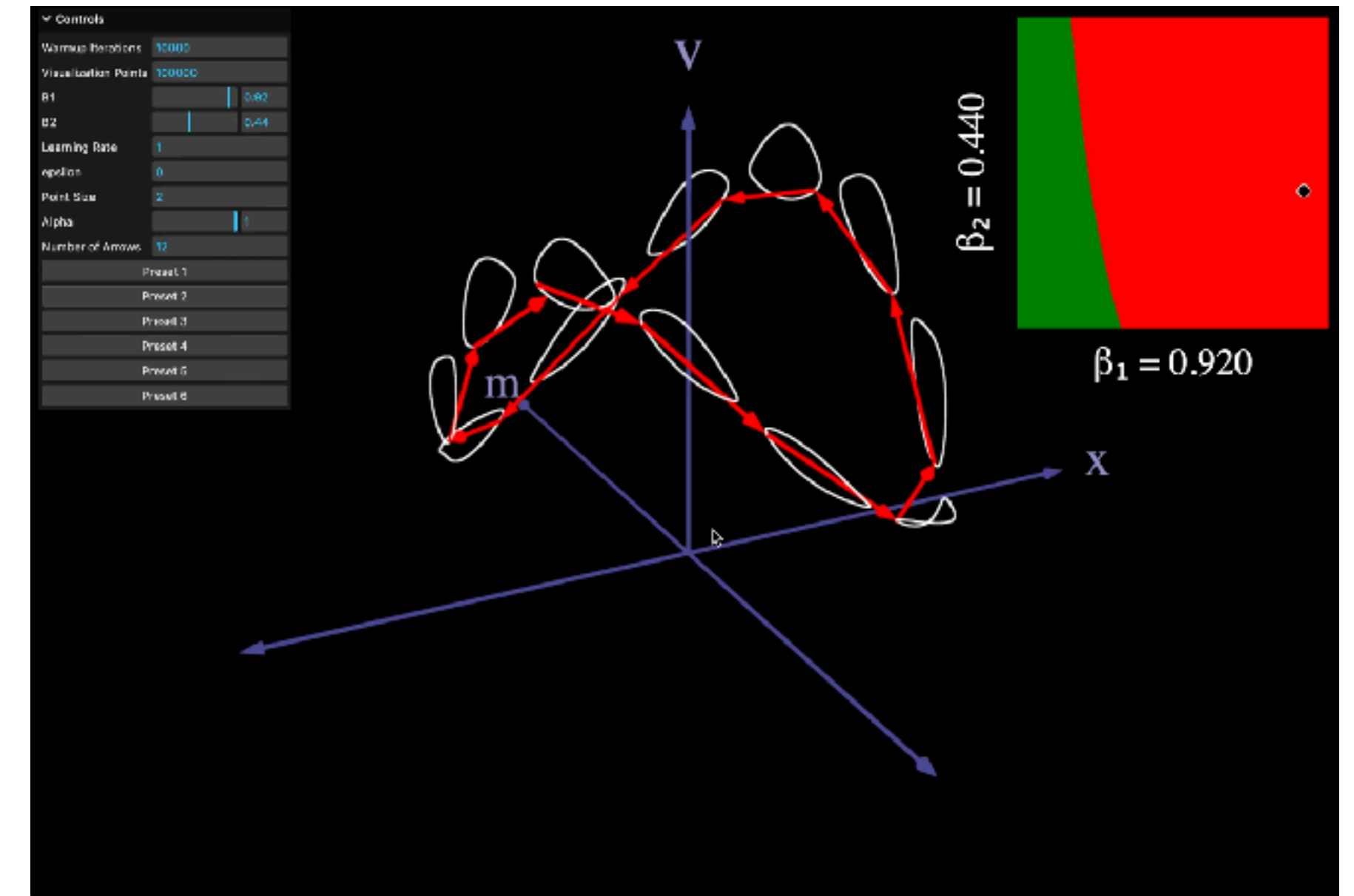


Jeremy Cohen
Flatiron Institute



<https://centralflows.github.io>

Animated Blog Post!



<https://alex-damian.github.io/adam/>

Bonus: crazy limit cycles in Adam