

# REPARAMETERIZING MIRROR DESCENT AS GRADIENT DESCENT

EHSAN AMID AND MANFRED K. WARMUTH

GOOGLE RESEARCH, BRAIN TEAM

{eamid, manfred}@google.com



## MIRROR DESCENT

$$\mathbf{w}_{s+1} = f^{-1}(f(\mathbf{w}_s) - h \nabla L(\mathbf{w}_s))$$

where  $f$  is (coordinate-wise) strictly monotonic link function

Gradient Descent (GD):

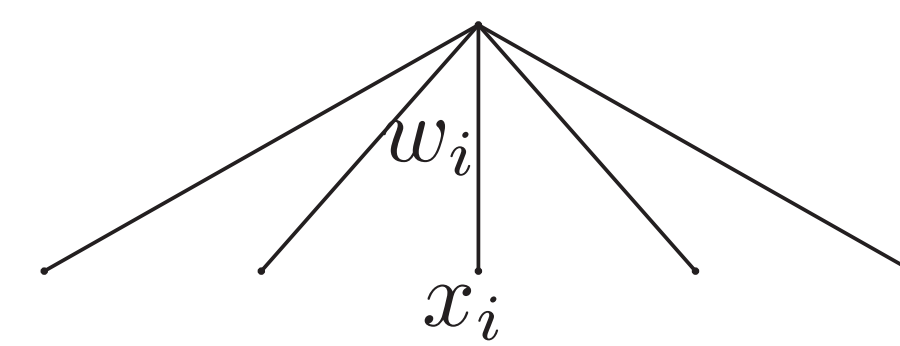
$$\mathbf{w}_{s+1} = \mathbf{w}_s - h \nabla L(\mathbf{w}_s) \quad (f(\mathbf{w}) = \mathbf{w})$$

Unnormalized Exponentiated Gradient Descent (EGU): [KW97]

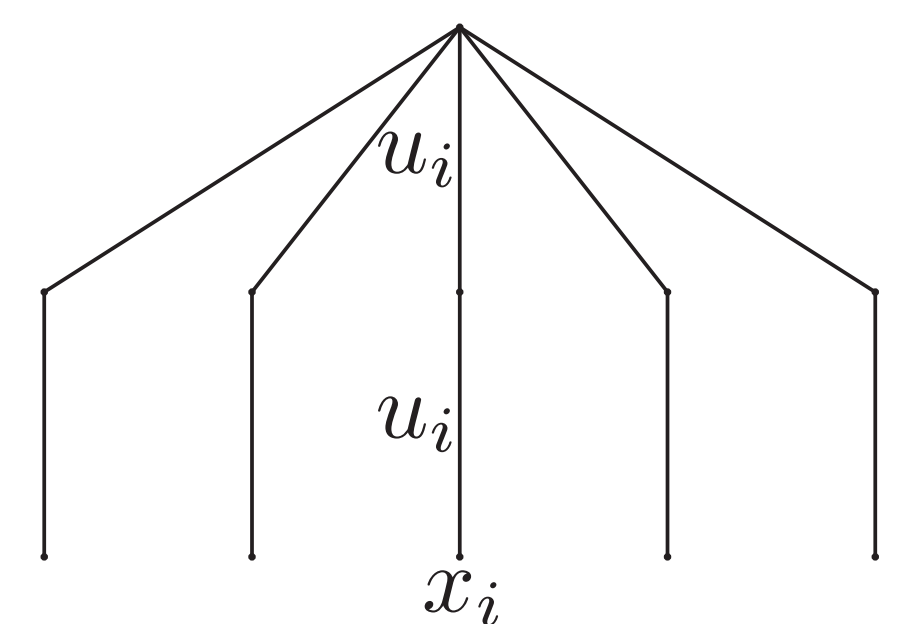
$$\mathbf{w}_{s+1} = \mathbf{w}_s \odot \exp(-h (\nabla L(\mathbf{w}_s))) \quad (f(\mathbf{w}) = \log \mathbf{w})$$

(with  $w_i \geq 0$ )

## SURPRISES



When linear neuron is trained with GD:  
lower bound for linear decrease of av. loss [WV05]



Reparameterize weights  $w_i$  by  $u_i^2$  [Akin79, GWBNS17]

Continuous GD on  $u_i$  simulates continuous-time EGU on  $w_i$

Discretizations learn Hadamard with Backprop with essentially  $\mathcal{O}(\log n)$  examples

Experimentally indistinguishable from discretized EGU

## MAIN CASE: EGU AS GD

Link

$$f(\mathbf{w}) = \log(\mathbf{w})$$

Reparameterization

$$\mathbf{w} = q(\mathbf{u}) := 1/4 \mathbf{u} \odot \mathbf{u}$$

$$(\mathbf{J}_f(\mathbf{w}))^{-1} = (\text{diag}(\mathbf{w})^{-1})^{-1} = \text{diag}(\mathbf{w})$$

$$\mathbf{J}_q(\mathbf{u})(\mathbf{J}_q(\mathbf{u}))^\top = 1/2 \text{diag}(\mathbf{u}) (1/2 \text{diag}(\mathbf{u}))^\top = \text{diag}(\mathbf{w})$$

Therefore

$$\dot{\log}(\mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{equals} \quad \dot{\mathbf{u}} = \underbrace{-\nabla L \circ q}_{\nabla_{\mathbf{u}} L(1/4 \mathbf{u} \odot \mathbf{u})}(\mathbf{u})$$

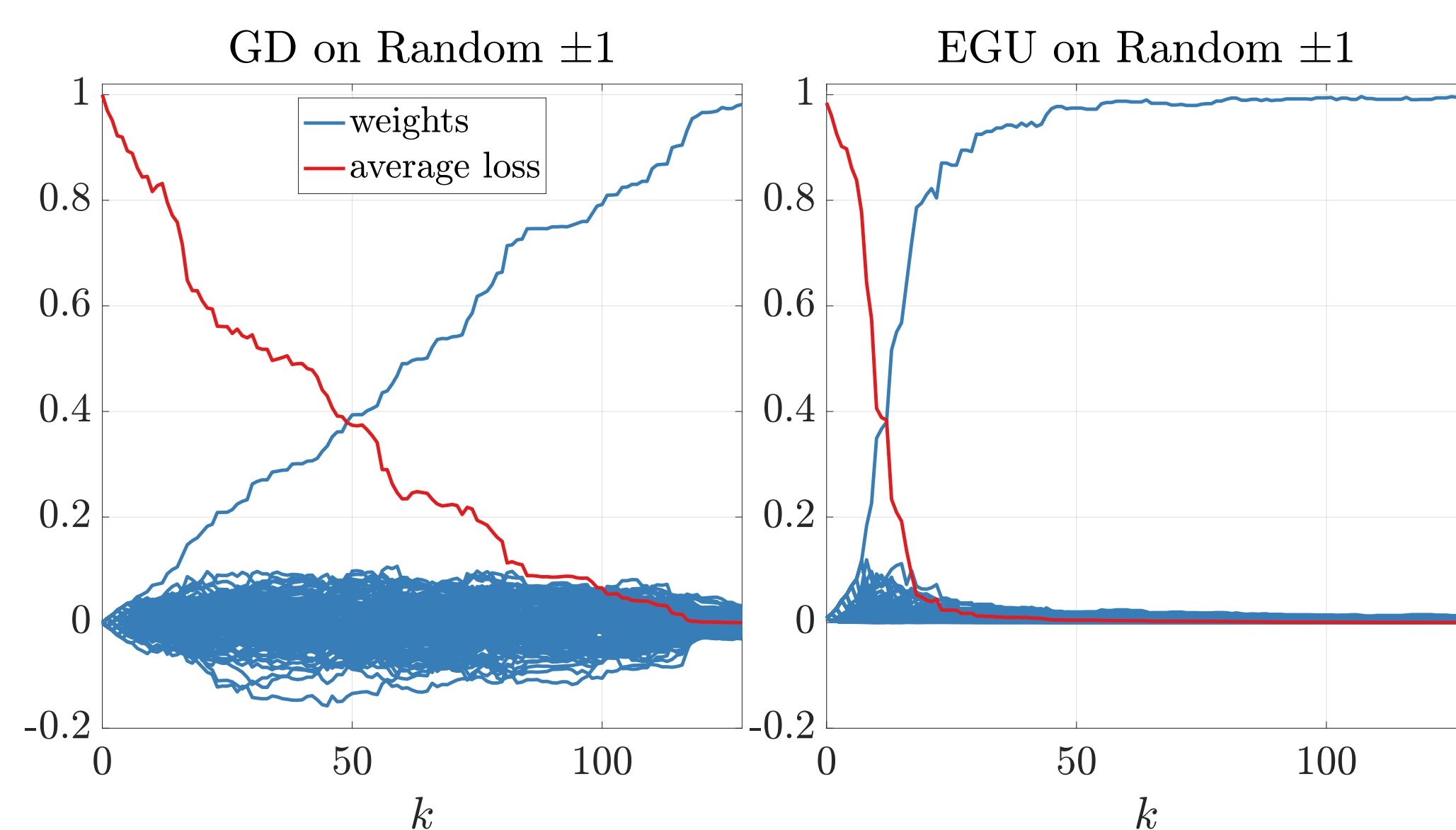
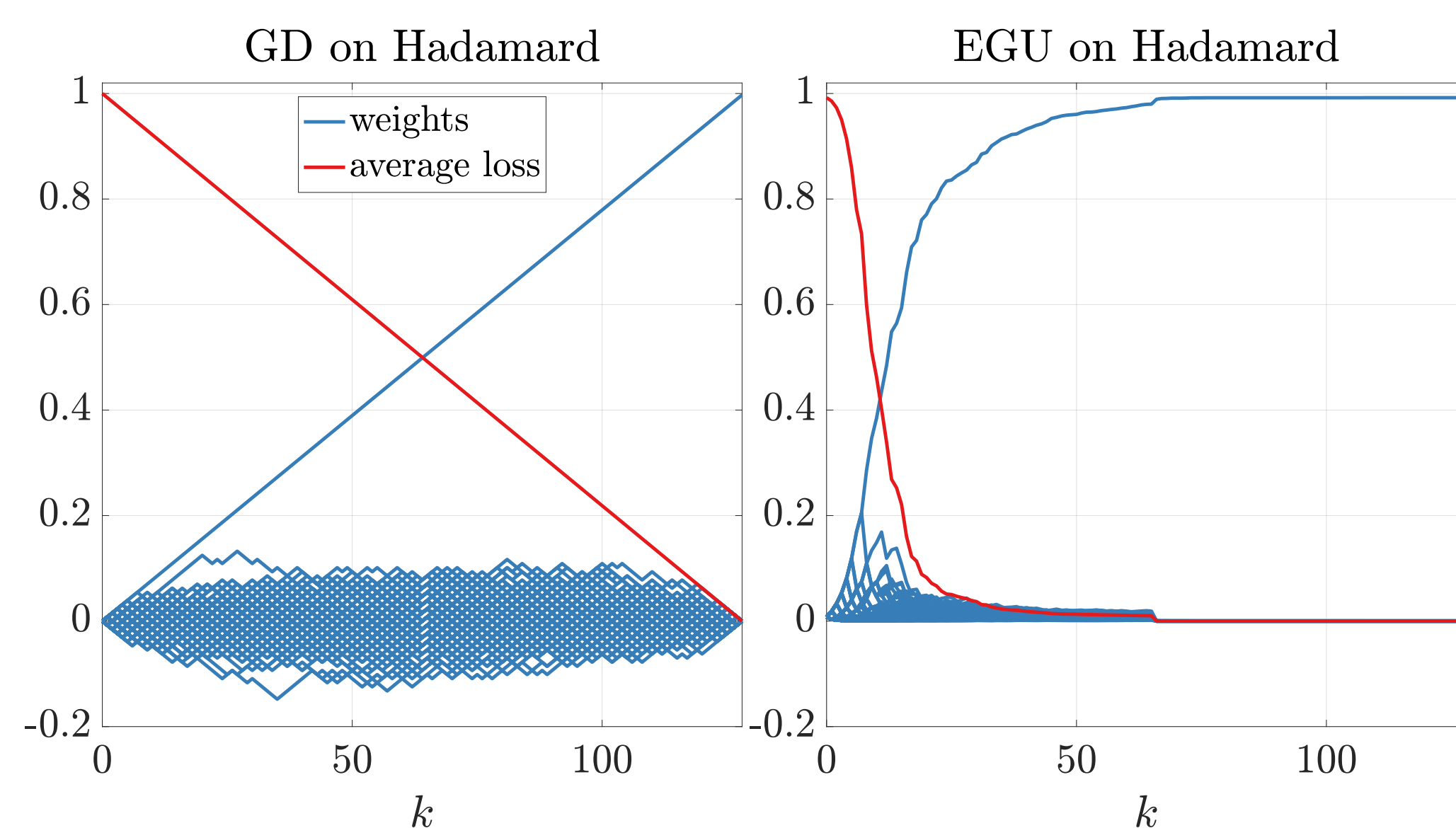
## MAJOR DIFFERENCES

GD: backprop, kernel methods

EG: Winnow, expert algorithms, Boosting, Bayes

Setup: 128x128 Hadamard (top) and random (bottom)  $\pm 1$  matrix

Rows are instances, labels are the first column, square loss



x-axis:  $k = 1..128$

y-axis: all 128 weights Loss when trained on examples 1..k

Upshot: After half examples, GD has average loss  $1/2$

EG family converges in  $\mathcal{O}(\log(n))$  many examples

## FOCUS ON CONTINUOUS MD

$$\dot{f}(\mathbf{w}) = -\nabla L(\mathbf{w})$$

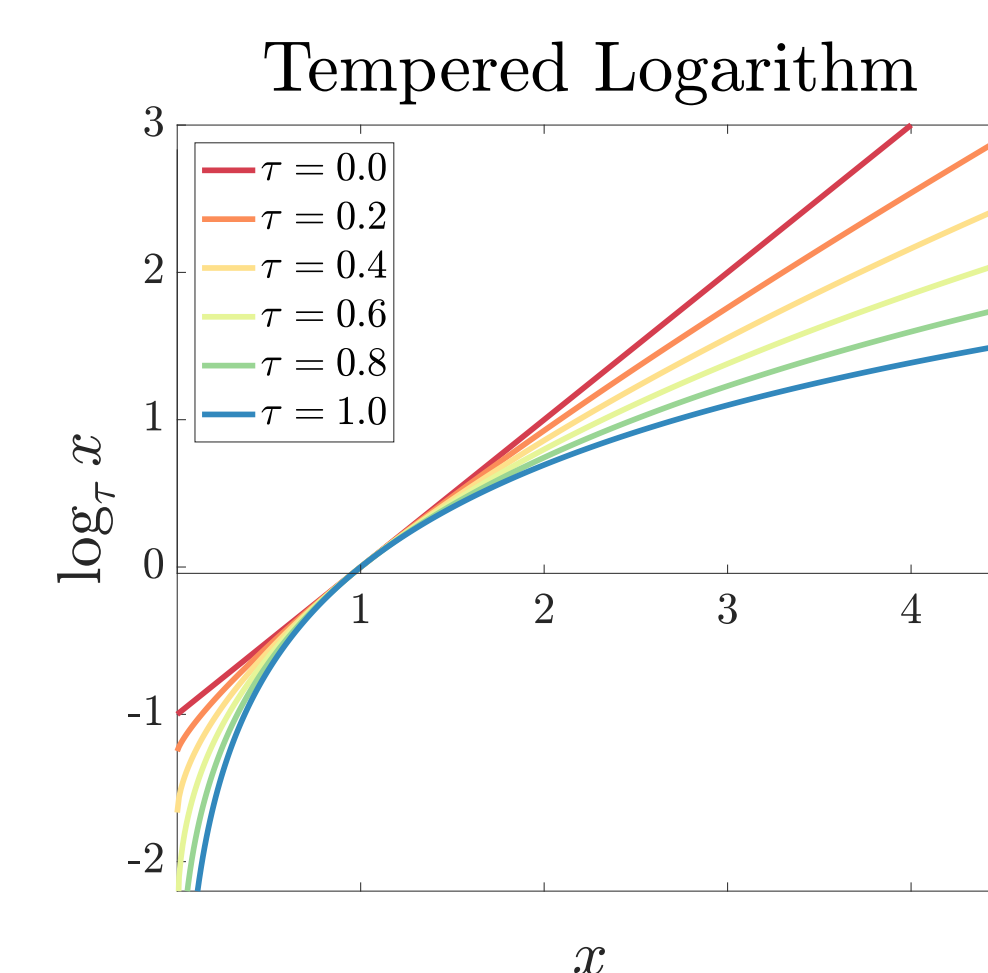
Main examples:

GD ( $f(\mathbf{w}) = \mathbf{w}$ ) and EGU ( $f(\mathbf{w}) = \log(\mathbf{w})$ )

Between  $f(\mathbf{w}) = \mathbf{w}$  &  $f(\mathbf{w}) = \log \mathbf{w}$ :

$$\log_\tau(\mathbf{w}) := \frac{1}{1-\tau} (\mathbf{w}^{1-\tau} - 1)$$

(for  $\tau \in [0, 1]$ )



## BURG AS GD

$$(-1 \odot \mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{equals} \quad \dot{\mathbf{u}} = \underbrace{-\nabla L \circ q}_{\nabla_{\mathbf{u}} L(\exp(\mathbf{u}))}(\mathbf{u})$$

## TEMPERED $\log_\tau$ AS GD

$$\dot{\log}_\tau(\mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{equals} \quad \dot{\mathbf{u}} = \underbrace{-\nabla L \circ q}_{\nabla_{\mathbf{u}} L\left(\left(\frac{2-\tau}{2}\right)^{\frac{2}{2-\tau}} \mathbf{u}^{\frac{2}{2-\tau}}\right)}(\mathbf{u})$$

## CONCLUSION

- World of continuous updates more succinct

$$\text{Euler discr.: } \frac{f(\mathbf{w}(t+h)) - f(\mathbf{w}(t))}{h} = -\nabla L(\mathbf{w}(t))$$

$$\iff \mathbf{w}(t+h) = f^{-1}(f(\mathbf{w}(t)) - h \nabla L(\mathbf{w}(t)))$$

- Under what conditions does the discrete MD track continuous MD
- Discretization of reparameterized EGU as GD tracks discrete EGU well enough so that same regret bounds hold [AW20]
- Discretization of reparameterized EGU as GD sample efficiently learns Hadamard problem

## REPARAMETERIZATION

Main theorem: For the reparameterization function  $\mathbf{w} = q(\mathbf{u})$  with the property that  $\text{range}(q) = \text{dom}(f)$ , the two updates

$$\dot{f}(\mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{and} \quad \dot{\mathbf{u}} = -\nabla L \circ q(\mathbf{u}),$$

coincide if that  $\mathbf{w}(0) = q(\mathbf{u}(0))$ ,  $\text{range}(q) \subseteq \text{dom}(F)$ , and we have

$$(\mathbf{J}_f(\mathbf{w}))^{-1} = \mathbf{J}_q(\mathbf{u}) (\mathbf{J}_q(\mathbf{u}))^\top$$