

---

# On Training Implicit Models

## (Supplementary material)

---

Anonymous Author(s)

Affiliation

Address

email

### 1 A Algorithm of Phantom Gradients

2 The following PyTorch-style [1] pseudo code describes the implementation of phantom gradients in  
3 both the unrolling form (see Alg. 1) and the Neumann form (see Alg. 2). To implement phantom  
4 gradients with TensorFlow [2], replace the *no\_grad* context manager with the *stop\_gradient*  
5 operator.

6 The unrolling-based phantom gradient is computed by the automatic differentiation engine, while  
7 the Neumann-series-based phantom gradient is given by Alg. 3. A special reminder is that, for a  
8 trained model, removing the unrolling steps in the test stage will not cause the performance decay but  
9 accelerate the inference instead. After training, even increasing the unrolling steps  $k$  to 20 in the test  
10 stage can not further improve the performance. This implies that the final results are fully obtained  
11 by the implicit model rather than the unrolling steps.

---

#### Algorithm 1 Unrolling-based phantom gradient, PyTorch-style

---

```
# solver: the solver to find  $\mathbf{h}^*$ , e.g., the Broyden solver in MDEQ.  
# func: the explicit function  $\mathcal{F}$  that defines the implicit model.  
# z: the input variables  $\mathbf{z}$  to solve  $\mathbf{h}^* = \mathcal{F}(\mathbf{h}^*, \mathbf{z})$   
# h: the solution  $\mathbf{h}^*$  of the implicit models.  
# training: a bool variable that indicates training or inference.  
# k: the unrolling step  $k$ .  
# lambda_: the damping factor  $\lambda$ .  
  
# a plain forward pass using Pytorch  
# calculate the phantom gradient by automatic differentiation  
# input: z & output: h  
def forward(z):  
    with torch.no_grad():  
        h = solver(func, z)  
  
    # define the computational graph for the backward pass.  
    # only used in the training stage  
    if training:  
        for _ in range(k):  
            h = (1 - lambda_) * h + lambda_ * func(h, z)  
  
    return h
```

---

---

**Algorithm 2** Neumann-series-based Phantom Gradient, Pytorch-style

---

```
# solver: the solver to find  $\mathbf{h}^*$ , e.g., the Broyden solver in MDEQ.
# func: the explicit function  $\mathcal{F}$  that defines the implicit model.
# grad(a, b, c): the function to compute the Jacobian-vector product  $(\partial \mathbf{a} / \partial \mathbf{b}) \mathbf{c}$ 
# z: the input variables  $\mathbf{z}$  to solve  $\mathbf{h}^* = \mathcal{F}(\mathbf{h}^*, \mathbf{z})$ 
# h: the output  $\mathbf{h}^*$  of the implicit model.
# k: the unrolling step  $k$ .
# lambda_: the damping factor  $\lambda$ .

# a plain forward pass using Pytorch
# input: z & output: h
def forward(z):
    with torch.no_grad():
        h = solver(func, z)

    return h

# phantom gradient for the backward pass
# input: dl / dh & output: dl / dz
def phantom_grad(g):
    # forward pass for automatic differentiation
    f = (1 - lambda_) * h + lambda_ * func(h, z)

    g_hat = g
    for _ in range(k-1):
        # compute Jacobian-vector product with automatic differentiation
        g_hat = g + grad(f, h, g_hat)

    # compute Jacobian-vector product to obtain dl / dz
    g_hat = grad(f, z, g_hat)
    return g_hat
```

---

---

**Algorithm 3** Neumann-series-based phantom gradient with  $\mathcal{O}(1)$  memory

---

```
1: Input  $\partial \mathcal{L} / \partial \mathbf{h}, \mathcal{F}, \mathbf{h}^*, k, \lambda$ .
2: Initialize  $\hat{\mathbf{g}} = \mathbf{g} = \partial \mathcal{L} / \partial \mathbf{h}$ ;
3:  $\mathbf{f} \leftarrow (1 - \lambda) \mathbf{h}^* + \lambda \mathcal{F}(\mathbf{h}^*, \mathbf{z})$ 
4: for  $i = 1, 2, \dots, k - 1$  do
5:    $\hat{\mathbf{g}} \leftarrow \mathbf{g} + (\partial \mathbf{f} / \partial \mathbf{h}) \hat{\mathbf{g}}$ ;  $\triangleright$  Compute Jacobian-vector product with automatic differentiation
6: end for
7:  $\hat{\mathbf{g}} \leftarrow (\partial \mathbf{f} / \partial \mathbf{z}) \hat{\mathbf{g}}$   $\triangleright$  Compute Jacobian-vector product to obtain the phantom gradient w.r.t.  $\mathbf{z}$ 
8: return  $\hat{\mathbf{g}}$ .
```

---

## 12 B Proof of Theorems

13 **Theorem 1.** Suppose the exact gradient and the phantom gradient are given by Eq. (4) and Eq. (5),  
14 respectively. Let  $\sigma_{\max}$  and  $\sigma_{\min}$  be the maximal and minimal singular value of  $\partial \mathcal{F} / \partial \boldsymbol{\theta}$ . If

$$\left\| \mathbf{A} \left( \mathbf{I} - \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \right) - \frac{\partial \mathcal{F}}{\partial \boldsymbol{\theta}} \right\| \leq \frac{\sigma_{\min}^2}{\sigma_{\max}}, \quad (\text{A-1})$$

15 then the phantom gradient provides a descent direction of the function  $\mathcal{L}$ , i.e.,

$$\left\langle \widehat{\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}}, \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\rangle \geq 0. \quad (\text{A-2})$$

16

17 *Proof.* Denote  $\mathbf{J} = \partial\mathcal{F}/\partial\boldsymbol{\theta}$ ,  $\mathbf{v} = \partial\mathcal{L}/\partial\mathbf{h}$ , and  $\mathbf{u} = (\mathbf{I} - \partial\mathcal{F}/\partial\mathbf{h})^{-1} \mathbf{v}$ . Let

$$\mathbf{E} = \mathbf{A} \left( \mathbf{I} - \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \right) - \frac{\partial\mathcal{F}}{\partial\boldsymbol{\theta}}, \quad (\text{A-3})$$

18 and we have  $\|\mathbf{E}\| \leq \sigma_{\min}^2/\sigma_{\max}$ . Then,

$$\begin{aligned} \left\langle \frac{\partial\mathcal{L}}{\partial\boldsymbol{\theta}}, \frac{\partial\mathcal{L}}{\partial\boldsymbol{\theta}} \right\rangle &= \mathbf{v}^\top \mathbf{A}^\top \mathbf{J} \left( \mathbf{I} - \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \right)^{-1} \mathbf{v} = \mathbf{u}^\top \left( \mathbf{I} - \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \right)^\top \mathbf{A}^\top \mathbf{J} \mathbf{u} = \mathbf{u}^\top (\mathbf{J} + \mathbf{E})^\top \mathbf{J} \mathbf{u} \\ &\geq \|\mathbf{J}\mathbf{u}\|^2 - \|\mathbf{E}\| \|\mathbf{J}\| \|\mathbf{u}\|^2 \geq (\sigma_{\min}^2 - \sigma_{\max} \|\mathbf{E}\|) \|\mathbf{u}\|^2 \geq 0, \end{aligned} \quad (\text{A-4})$$

19 which concludes the proof.  $\square$

20 *Proof of Remark 1.* Suppose  $\mathbf{A} = (\partial\mathcal{F}/\partial\boldsymbol{\theta}) \mathbf{D}$  and the condition in Eq. (8). Then,

$$\left\| \mathbf{A} \left( \mathbf{I} - \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \right) - \frac{\partial\mathcal{F}}{\partial\boldsymbol{\theta}} \right\| \leq \left\| \frac{\partial\mathcal{F}}{\partial\boldsymbol{\theta}} \right\| \left\| \mathbf{D} \left( \mathbf{I} - \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \right) - \mathbf{I} \right\| \leq \sigma_{\max} \cdot \frac{1}{\kappa^2} = \frac{\sigma_{\min}^2}{\sigma_{\max}}, \quad (\text{A-5})$$

21 indicating the condition in Eq. (6) is satisfied.  $\square$

22 **Theorem 2.** Suppose the matrix  $\partial\mathcal{F}/\partial\mathbf{h}$  is a contractive mapping. Then,

- 23 (i) the Neumann series in (15) converges to the Jacobian-inverse  $(\mathbf{I} - \partial\mathcal{F}/\partial\mathbf{h})^{-1}$ ; and  
 24 (ii) if the function  $\mathcal{F}$  is continuously differentiable w.r.t. both  $\mathbf{h}$  and  $\boldsymbol{\theta}$ , the sequence in Eq. (14)  
 25 converges to the exact Jacobian  $\partial\mathbf{h}^*/\partial\boldsymbol{\theta}$  as  $T \rightarrow \infty$ , i.e.,

$$\lim_{T \rightarrow \infty} \frac{\partial\mathbf{h}_T}{\partial\boldsymbol{\theta}} = \frac{\partial\mathcal{F}}{\partial\boldsymbol{\theta}} \bigg|_{\mathbf{h}^*} \left( \mathbf{I} - \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \bigg|_{\mathbf{h}^*} \right)^{-1}. \quad (\text{A-6})$$

26 *Proof.* (i) Since  $\|\partial\mathcal{F}/\partial\mathbf{h}\| < 1$ ,

$$\|\mathbf{B}\| \leq \lambda \left\| \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \right\| + (1 - \lambda) \|\mathbf{I}\| < 1. \quad (\text{A-7})$$

27 Let  $\mathbf{B}_k = \sum_{t=0}^{k-1} \mathbf{B}^t$ , and for each  $p \in \mathbb{N}_+$ , we have

$$\|\mathbf{B}_{k+p} - \mathbf{B}_k\| = \left\| \sum_{t=k}^{k+p-1} \mathbf{B}^t \right\| \leq \|\mathbf{B}\|^k \left\| \sum_{t=0}^{p-1} \mathbf{B}^t \right\| \leq \|\mathbf{B}\|^k \sum_{t=0}^{p-1} \|\mathbf{B}\|^t < \frac{\|\mathbf{B}\|^k}{1 - \|\mathbf{B}\|}. \quad (\text{A-8})$$

28 By the Cauchy's convergence test, the sequence  $\{\mathbf{B}_k\}$  is convergent. Since

$$(\mathbf{I} - \mathbf{B})\mathbf{B}_k = \mathbf{I} - \mathbf{B}^k \rightarrow \mathbf{I}, \quad \text{as } k \rightarrow \infty, \quad (\text{A-9})$$

29 it follows that  $\mathbf{B}_k \rightarrow (\mathbf{I} - \mathbf{B})^{-1}$ , as  $k \rightarrow \infty$ . Therefore,

$$\lambda \sum_{t=0}^{\infty} \mathbf{B}^t = \lambda (\mathbf{I} - \mathbf{B})^{-1} = \left( \mathbf{I} - \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \right)^{-1}. \quad (\text{A-10})$$

30 (ii) Let  $\mathcal{H}(\mathbf{h}, \mathbf{z}) = \lambda\mathcal{F}(\mathbf{h}, \mathbf{z}) + (1 - \lambda)\mathbf{h}$ , and

$$\frac{\partial\mathcal{H}}{\partial\mathbf{h}} = \lambda \frac{\partial\mathcal{F}}{\partial\mathbf{h}} + (1 - \lambda)\mathbf{I}. \quad (\text{A-11})$$

31 Similar to Eq. (A-7),  $\partial\mathcal{H}/\partial\mathbf{h}$  is also a contractive mapping. By the Banach Fixed Point Theorem [3],  
 32 the sequence  $\{\mathbf{h}_t\}$  converges to an exact fixed point  $\mathbf{h}^*$  of  $\mathcal{H}$ , which is also a fixed point of  $\mathcal{F}$ .

33 Denote

$$\mathbf{U}_t = \frac{\partial\mathcal{F}}{\partial\boldsymbol{\theta}} \bigg|_{\mathbf{h}_t}, \quad \mathbf{V}_t = \lambda \frac{\partial\mathcal{F}}{\partial\mathbf{h}} \bigg|_{\mathbf{h}_t} + (1 - \lambda)\mathbf{I}. \quad (\text{A-12})$$

34 Since the function  $\mathcal{F}$  is continuously differentiable w.r.t. both  $\mathbf{h}$  and  $\boldsymbol{\theta}$ , we have

$$\lim_{t \rightarrow \infty} \mathbf{U}_t = \frac{\partial \mathcal{F}}{\partial \boldsymbol{\theta}} \Big|_{\mathbf{h}^*} = \mathbf{U}_\infty, \quad \lim_{t \rightarrow \infty} \mathbf{V}_t = \lambda \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \Big|_{\mathbf{h}^*} + (1 - \lambda) \mathbf{I} = \mathbf{V}_\infty. \quad (\text{A-13})$$

35 According to the conclusion in (i), we have

$$\frac{\partial \mathcal{F}}{\partial \boldsymbol{\theta}} \Big|_{\mathbf{h}^*} \left( \mathbf{I} - \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \Big|_{\mathbf{h}^*} \right)^{-1} = \lambda \mathbf{U}_\infty \sum_{t=0}^{\infty} \mathbf{V}_\infty^t. \quad (\text{A-14})$$

36 Comparing Eq. (14) with Eq. (A-14), we have

$$\begin{aligned} & \left\| \frac{\partial \mathbf{h}_T}{\partial \boldsymbol{\theta}} - \frac{\partial \mathcal{F}}{\partial \boldsymbol{\theta}} \Big|_{\mathbf{h}^*} \left( \mathbf{I} - \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \Big|_{\mathbf{h}^*} \right)^{-1} \right\| = \lambda \left\| \sum_{t=0}^{T-1} \mathbf{U}_t \prod_{s=t+1}^{T-1} \mathbf{V}_s - \mathbf{U}_\infty \sum_{t=0}^{\infty} \mathbf{V}_\infty^t \right\| \\ & \leq \lambda \left( \underbrace{\left\| \sum_{t=0}^{T-2} \mathbf{U}_t \left( \prod_{s=t+1}^{T-1} \mathbf{V}_s - \mathbf{V}_\infty^{T-t-1} \right) \right\|}_{\Delta_1} + \underbrace{\left\| \sum_{t=0}^{T-1} (\mathbf{U}_t - \mathbf{U}_\infty) \mathbf{V}_\infty^{T-t-1} \right\|}_{\Delta_2} + \underbrace{\left\| \mathbf{U}_\infty \sum_{t=T}^{\infty} \mathbf{V}_\infty^t \right\|}_{\Delta_3} \right). \end{aligned} \quad (\text{A-15})$$

37 In the following context, we prove Eq. (A-6) by showing that  $\Delta_1$ ,  $\Delta_2$ , and  $\Delta_3$  can be arbitrarily  
38 small when  $T$  is sufficiently large.

39 **Preparations.** For any  $\epsilon > 0$ , since  $\mathbf{U}_t \rightarrow \mathbf{U}_\infty$  and  $\mathbf{V}_t \rightarrow \mathbf{V}_\infty$  as  $t \rightarrow \infty$ , there exists  $N \in \mathbb{N}_+$  s.t.

$$\|\mathbf{U}_t - \mathbf{U}_\infty\| < \epsilon, \quad \|\mathbf{V}_t - \mathbf{V}_\infty\| < \epsilon, \quad \forall t > N. \quad (\text{A-16})$$

40 Since  $\partial \mathcal{H} / \partial \mathbf{h}$  is a contractive mapping, there exists  $\gamma \in (0, 1)$  s.t.

$$\|\mathbf{V}_t\| \leq \gamma, \quad \|\mathbf{V}_\infty\| \leq \gamma. \quad (\text{A-17})$$

41 Besides, since  $\partial \mathcal{F} / \partial \boldsymbol{\theta}$  is a continuous function and  $\{\mathbf{h}_t\}$  is a convergent sequence, it follows that  
42  $\{\mathbf{h}_t\}$  forms a compact set and that  $\partial \mathcal{F} / \partial \boldsymbol{\theta}$  is bounded on  $\{\mathbf{h}_t\}$ . Therefore, there exists  $M > 0$ , s.t.

$$\|\mathbf{U}_t\| \leq M, \quad t = 0, 1, 2, \dots \quad (\text{A-18})$$

43 Taking  $t \rightarrow \infty$ , we have  $\|\mathbf{U}_\infty\| \leq M$ .

44 **For  $\Delta_1$ .** For  $t > N$ , consider

$$\begin{aligned} \left\| \mathbf{U}_t \left( \prod_{s=t+1}^{T-1} \mathbf{V}_s - \mathbf{V}_\infty^{T-t-1} \right) \right\| & \leq \|\mathbf{U}_t\| \sum_{s=t+1}^{T-1} \left\| \mathbf{V}_{t+1} \mathbf{V}_{t+2} \cdots \mathbf{V}_s \mathbf{V}_\infty^{T-s-1} - \mathbf{V}_{t+1} \mathbf{V}_{t+2} \cdots \mathbf{V}_{s-1} \mathbf{V}_\infty^{T-s} \right\| \\ & \leq \|\mathbf{U}_t\| \sum_{s=t+1}^{T-1} \|\mathbf{V}_{t+1}\| \|\mathbf{V}_{t+2}\| \cdots \|\mathbf{V}_{s-1}\| \|\mathbf{V}_s - \mathbf{V}_\infty\| \|\mathbf{V}_\infty\|^{T-s-1} \\ & \leq M(T-t-1)\gamma^{T-t-2}\epsilon, \end{aligned} \quad (\text{A-19})$$

45 and for  $t \leq N$ , we simply have

$$\left\| \mathbf{U}_t \left( \prod_{s=t+1}^{T-1} \mathbf{V}_s - \mathbf{V}_\infty^{T-t-1} \right) \right\| \leq \|\mathbf{U}_t\| \left( \prod_{s=t+1}^{T-1} \|\mathbf{V}_s\| + \|\mathbf{V}_\infty\|^{T-t-1} \right) \leq 2M\gamma^{T-t-1}. \quad (\text{A-20})$$

46 Therefore, when  $T > N + 2$ ,  $\Delta_1$  can be bounded as follows:

$$\begin{aligned} \Delta_1 & \leq \left( \sum_{t=0}^N + \sum_{t=N+1}^{T-2} \right) \left\| \mathbf{U}_t \left( \prod_{s=t+1}^{T-1} \mathbf{V}_s - \mathbf{V}_\infty^{T-t-1} \right) \right\| \\ & \leq 2M \sum_{t=0}^N \gamma^{T-t-1} + M\epsilon \sum_{t=N+1}^{T-2} (T-t-1)\gamma^{T-t-2} \\ & \leq 2M\gamma^{T-N-1} \frac{1-\gamma^{N+1}}{1-\gamma} + \left( \frac{1-\gamma^{T-N-2}}{(1-\gamma)^2} - \frac{(T-N-2)\gamma^{T-N-2}}{1-\gamma} \right) M\epsilon \\ & \leq \frac{2M}{1-\gamma} \gamma^{T-N-1} + \frac{M}{(1-\gamma)^2} \epsilon. \end{aligned} \quad (\text{A-21})$$

47 Since  $M/(1-\gamma)^2$  is a constant and  $\gamma^{T-N-1} \rightarrow 0$  as  $T \rightarrow \infty$ ,  $\Delta_1$  can be arbitrarily small for a  
 48 sufficiently large  $T$ .

49 **For  $\Delta_2$ .** Consider

$$\|(U_t - U_\infty) \mathbf{V}_\infty^{T-t-1}\| \leq \|U_t - U_\infty\| \|\mathbf{V}_\infty\|^{T-t-1} \leq \begin{cases} \gamma^{T-t-1} \epsilon, & \text{when } t \geq N; \\ 2M\gamma^{T-t-1} & \text{when } t < N. \end{cases} \quad (\text{A-22})$$

50 Therefore, when  $T > N + 2$ ,  $\Delta_2$  can be bounded as follows:

$$\begin{aligned} \Delta_2 &\leq \left( \sum_{t=0}^N + \sum_{t=N+1}^{T-1} \right) \|(U_t - U_\infty) \mathbf{V}_\infty^{T-t-1}\| \\ &\leq 2M \sum_{t=0}^N \gamma^{T-t-1} + \epsilon \sum_{t=N+1}^{T-1} \gamma^{T-t-1} \\ &\leq \frac{2M}{1-\gamma} \gamma^{T-N-1} + \frac{\epsilon}{1-\gamma}. \end{aligned} \quad (\text{A-23})$$

51 Since  $1/(1-\gamma)$  is a constant and  $\gamma^{T-N-1} \rightarrow 0$  as  $T \rightarrow \infty$ ,  $\Delta_2$  can be arbitrarily small for a  
 52 sufficiently large  $T$ .

53 **For  $\Delta_3$ .** As  $t \rightarrow \infty$ , we have

$$\left\| U_\infty \sum_{t=T}^{\infty} \mathbf{V}_\infty^t \right\| \leq \|U_\infty\| \|\mathbf{V}_\infty\|^T \left\| (\mathbf{I} - \mathbf{V}_\infty)^{-1} \right\| \leq M \cdot \gamma^T \cdot \frac{1}{1-\gamma} \rightarrow 0. \quad (\text{A-24})$$

54 As a result, the conclusion in Eq. (A-6) is proved.  $\square$

55 **Theorem 3.** Suppose the loss function  $\mathcal{R}$  in Eq. (3) is  $\ell$ -smooth, lower-bounded, and has bounded  
 56 gradient almost surely in the training process. Besides, assume the gradient in Eq. (4) is an  
 57 unbiased estimator of  $\nabla \mathcal{R}(\theta)$  with a bounded covariance. If the phantom gradient in Eq. (5) is an  
 58  $\epsilon$ -approximation to the gradient in Eq. (4), i.e.,

$$\left\| \widehat{\frac{\partial \mathcal{L}}{\partial \theta}} - \frac{\partial \mathcal{L}}{\partial \theta} \right\| \leq \epsilon, \quad \text{almost surely}, \quad (\text{A-25})$$

59 then using Eq. (5) as a stochastic first-order oracle with a step size of  $\eta_\tau = O(1/\sqrt{\tau})$  to update  $\theta$   
 60 with gradient descent, it follows after  $T$  iterations that

$$\mathbb{E} \left[ \frac{\sum_{\tau=1}^T \eta_\tau \|\nabla \mathcal{R}(\theta_\tau)\|^2}{\sum_{\tau=1}^T \eta_\tau} \right] \leq O \left( \epsilon + \frac{\log T}{\sqrt{T}} \right). \quad (\text{A-26})$$

61

62 *Proof.* Let  $\widehat{\frac{\partial \mathcal{L}_\tau}{\partial \theta}}$  be the phantom gradient at the  $\tau^{\text{th}}$  iteration. By  $\ell$ -smoothness of  $\mathcal{R}$ , we have

$$\begin{aligned} \mathcal{R}(\theta_{\tau+1}) &\leq \mathcal{R}(\theta_\tau) + \langle \nabla \mathcal{R}(\theta_\tau), \theta_{\tau+1} - \theta_\tau \rangle + \frac{\ell}{2} \|\theta_{\tau+1} - \theta_\tau\|^2 \\ &= \mathcal{R}(\theta_\tau) - \eta_\tau \left\langle \nabla \mathcal{R}(\theta_\tau), \widehat{\frac{\partial \mathcal{L}_\tau}{\partial \theta}} \right\rangle + \frac{\ell \eta_\tau^2}{2} \left\| \widehat{\frac{\partial \mathcal{L}_\tau}{\partial \theta}} \right\|^2. \end{aligned} \quad (\text{A-27})$$

63 Let

$$\mathbf{e}_\tau = \frac{\partial \mathcal{L}}{\partial \theta} \Big|_{\theta=\theta_\tau} - \widehat{\frac{\partial \mathcal{L}_\tau}{\partial \theta}} \quad (\text{A-28})$$

64 be the approximation error at the  $\tau^{\text{th}}$  iteration. Taking expectation w.r.t. the first  $\tau$  iterations, we have

$$\mathbb{E}_{1 \sim \tau} [\mathcal{R}(\theta_{\tau+1})] = \mathbb{E}_{1 \sim \tau-1} [\mathbb{E}_\tau [\mathcal{R}(\theta_{\tau+1}) \mid 1 \sim \tau-1]] = \mathbb{E}_{1 \sim \tau-1} [\mathbb{E}_\tau [\mathcal{R}(\theta_{\tau+1}) \mid \theta_\tau]], \quad (\text{A-29})$$

65 where the first equality comes from the law of total expectation, while the second from the fact  
 66 that the stochasticity of the first  $\tau-1$  steps is totally captured by the value  $\theta_\tau$ . Consider the inner

67 expectation in Eq. (A-29), and we omit the condition on  $\boldsymbol{\theta}_\tau$  when no ambiguity is made. Note that in  
 68 the following derivation, all expectations and variances are conditioned on  $\boldsymbol{\theta}_\tau$ :

$$\begin{aligned}\mathbb{E}_\tau [\mathcal{R}(\boldsymbol{\theta}_{\tau+1})] &\leq \mathbb{E}_\tau \left[ \mathcal{R}(\boldsymbol{\theta}_\tau) - \eta_\tau \left\langle \nabla \mathcal{R}(\boldsymbol{\theta}_\tau), \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right\rangle + \frac{\ell \eta_\tau^2}{2} \left\| \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right\|^2 \right] \\ &= \mathcal{R}(\boldsymbol{\theta}_\tau) - \eta_\tau \left\langle \nabla \mathcal{R}(\boldsymbol{\theta}_\tau), \mathbb{E}_\tau \left[ \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right] \right\rangle + \frac{\ell \eta_\tau^2}{2} \mathbb{E}_\tau \left[ \left\| \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right\|^2 \right],\end{aligned}\quad (\text{A-30})$$

69 where

$$\mathbb{E}_\tau \left[ \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right] = \mathbb{E}_\tau \left[ \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_\tau} - \mathbf{e}_\tau \right] = \nabla \mathcal{R}(\boldsymbol{\theta}_\tau) - \mathbb{E}_\tau [\mathbf{e}_\tau], \quad (\text{A-31})$$

70 and

$$\mathbb{E}_\tau \left[ \left\| \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right\|^2 \right] = \left\| \mathbb{E}_\tau \left[ \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right] \right\|^2 + \text{tr} \left( \text{Cov}_\tau \left( \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right) \right). \quad (\text{A-32})$$

71 Suppose  $\|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\| \leq G$  almost surely, and then we have

$$\left\| \mathbb{E}_\tau \left[ \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right] \right\|^2 = \|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau) - \mathbb{E}_\tau [\mathbf{e}_\tau]\|^2 \leq (G + \epsilon)^2. \quad (\text{A-33})$$

72 Moreover, by the properties of covariance,

$$\begin{aligned}\text{tr} \left( \text{Cov}_\tau \left( \frac{\widehat{\partial \mathcal{L}_\tau}}{\partial \boldsymbol{\theta}} \right) \right) &= \text{tr} \left( \text{Cov}_\tau \left( \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_\tau} - \mathbf{e}_\tau \right) \right) \\ &= \text{tr} \left( \text{Cov}_\tau \left( \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_\tau} \right) \right) + \text{tr} (\text{Cov}_\tau (\mathbf{e}_\tau)) - 2 \text{tr} \left( \text{Cov}_\tau \left( \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_\tau}, \mathbf{e}_\tau \right) \right) \\ &\leq 2 \text{tr} \left( \text{Cov}_\tau \left( \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_\tau} \right) \right) + 2 \text{tr} (\text{Cov}_\tau (\mathbf{e}_\tau)),\end{aligned}\quad (\text{A-34})$$

73 where the last inequality comes from

$$\begin{aligned}|\text{tr} (\text{Cov} (\mathbf{a}, \mathbf{b}))| &\leq \sum_i |\text{Cov} (a_i, b_i)| \leq \sum_i \sqrt{\text{Var} (a_i) \text{Var} (b_i)} \leq \sum_i \frac{\text{Var} (a_i) + \text{Var} (b_i)}{2} \\ &= \frac{1}{2} (\text{tr} (\text{Cov} (\mathbf{a})) + \text{tr} (\text{Cov} (\mathbf{b}))).\end{aligned}\quad (\text{A-35})$$

74 By the Popoviciu's inequality on variances [4], the second term in (A-34) can be bounded by  $d_\theta \epsilon^2$ ,  
 75 i.e.,

$$\text{tr} (\text{Cov}_\tau (\mathbf{e}_\tau)) \leq d_\theta \epsilon^2, \quad (\text{A-36})$$

76 where  $d_\theta$  denotes the dimension of  $\boldsymbol{\theta}$ . Finally, since the gradient estimator  $\partial \mathcal{L} / \partial \boldsymbol{\theta}$  has a bounded  
 77 covariance, there exists  $M > 0$ , s.t.

$$\text{tr} \left( \text{Cov}_\tau \left( \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_\tau} \right) \right) \leq M, \quad \text{almost surely.} \quad (\text{A-37})$$

78 Combining (A-30), (A-31), (A-32), (A-33), (A-34), (A-37), we have

$$\begin{aligned}\mathbb{E}_\tau [\mathcal{R}(\boldsymbol{\theta}_{\tau+1})] &\leq \mathcal{R}(\boldsymbol{\theta}_\tau) - \eta_\tau \|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\|^2 + \eta_\tau \langle \nabla \mathcal{R}(\boldsymbol{\theta}_\tau), \mathbb{E}_\tau [\mathbf{e}_\tau] \rangle + K \eta_\tau^2, \\ &\leq \mathcal{R}(\boldsymbol{\theta}_\tau) - \eta_\tau \|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\|^2 + \eta_\tau \|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\| \|\mathbb{E}_\tau [\mathbf{e}_\tau]\| + K \eta_\tau^2 \\ &\leq \mathcal{R}(\boldsymbol{\theta}_\tau) - \eta_\tau \|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\|^2 + \eta_\tau G \epsilon + K \eta_\tau^2,\end{aligned}\quad (\text{A-38})$$

where  $K = \ell((G + \epsilon)^2 + 2M + 2d_\theta \epsilon^2)/2$  is a constant. Substitute (A-38) into Eq. (A-29), and it becomes

$$\mathbb{E}_{1 \sim \tau} [\mathcal{R}(\boldsymbol{\theta}_{\tau+1})] \leq \mathbb{E}_{1 \sim \tau-1} [\mathcal{R}(\boldsymbol{\theta}_\tau)] - \eta_\tau \mathbb{E}_{1 \sim \tau-1} [\|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\|^2] + \eta_\tau G\epsilon + K\eta_\tau^2. \quad (\text{A-39})$$

By taking a summation over the first  $T$  steps,

$$\begin{aligned} \mathbb{E}_{1 \sim T} \left[ \sum_{\tau=1}^T \eta_\tau \|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\|^2 \right] &\leq \mathcal{R}(\boldsymbol{\theta}_1) - \mathbb{E}_{1 \sim T} [\mathcal{R}(\boldsymbol{\theta}_{T+1})] + G\epsilon \sum_{\tau=1}^T \eta_\tau + K \sum_{\tau=1}^T \eta_\tau^2 \\ &\leq \mathcal{R}(\boldsymbol{\theta}_1) - m + G\epsilon \sum_{\tau=1}^T \eta_\tau + K \sum_{\tau=1}^T \eta_\tau^2, \end{aligned} \quad (\text{A-40})$$

where  $m = \inf_{\boldsymbol{\theta}} \mathcal{R}(\boldsymbol{\theta})$  since  $\mathcal{R}$  is lower-bounded. Dividing a factor of  $\sum_{\tau=1}^T \eta_\tau$ , we have

$$\mathbb{E}_{1 \sim T} \left[ \frac{\sum_{\tau=1}^T \eta_\tau \|\nabla \mathcal{R}(\boldsymbol{\theta}_\tau)\|^2}{\sum_{\tau=1}^T \eta_\tau} \right] \leq G\epsilon + \frac{\mathcal{R}(\boldsymbol{\theta}_1) - m}{\sum_{\tau=1}^T \eta_\tau} + K \frac{\sum_{\tau=1}^T \eta_\tau^2}{\sum_{\tau=1}^T \eta_\tau}. \quad (\text{A-41})$$

Since  $\eta_\tau = O(1/\sqrt{T})$ , it follows that

$$\sum_{\tau=1}^T \eta_\tau = O(\sqrt{T}), \quad \frac{\sum_{\tau=1}^T \eta_\tau^2}{\sum_{\tau=1}^T \eta_\tau} = O\left(\frac{\log T}{\sqrt{T}}\right). \quad (\text{A-42})$$

Combining (A-41) and Eq. (A-42) concludes the proof.  $\square$

**Remark 4.** The assumption that  $\mathcal{R}$  has almost-surely bounded gradient at  $\{\boldsymbol{\theta}_\tau\}_{\tau=0}^T$  is reasonable. Because of the norm-based regularization, *e.g.*, weight decay, we can assume  $\boldsymbol{\theta}$  is almost surely optimized within a compact set in the parameter space. If we further assume  $\mathcal{R}$  is continuously differentiable, the almost-sure boundedness of  $\|\nabla \mathcal{R}\|$  within the compact set follows its continuity.

**Remark 5.** We justify the assumption that the gradient in Eq. (4) has a bounded covariance. For the SGD algorithm, the stochasticity of Eq. (4) comes from the random sampling of the training example (or the training mini-batch) from the dataset. Since there are finite samples in the training set, the covariance of Eq. (4) remains finite. Moreover, as Theorem 2 only considers a finite training schedule, *i.e.*,  $T$  steps, the possible combination of the selected sample at each step is still finite (even though its number grows exponentially). Therefore, it is reasonable to assume the gradient in Eq. (4) has a bounded covariance.

## C Implementation Details

### C.1 Synthetic Settings

For the synthetic setting, the following model is used:

$$\mathbf{h}^* = \mathcal{F}(\mathbf{h}^* + \mathbf{u}) \quad (\text{A-43})$$

where  $\mathcal{F}$  is an 1-layer convolutional network with spectral normalization [5], and  $\mathbf{u}, \mathbf{h}^* \in \mathbb{R}^{B \times C \times N}$ . The loss  $\mathcal{L}$  is given by the mean squared error between  $\mathbf{h}^*$  and  $\mathbf{y}$ . We choose  $C = 128$ ,  $N = 32$ , and randomly sample 50000 data pairs  $(\mathbf{u}, \mathbf{y})$  to compute the gradient  $\partial \mathcal{L} / \partial \mathbf{u}$ .

We generate a symmetric weight matrix for the network and constrain the Lipschitz constant  $L_h$  to a given level with spectral normalization. For the visualization in the main paper, we adopt  $L_h = 0.9$ . For the additional visualization on the stability of the solver in Fig. 1, we choose  $L_h$  from  $\{0.9, 0.99, 0.999, 0.9999\}$ .

To solve  $\mathbf{h}^*$ , we employ the fixed-point iteration as the solver. For the synthetic setting shown in the main paper, we use 100 fixed-point iterations to obtain  $\mathbf{h}^*$  that satisfies the relative error  $\|\mathbf{h}^* - \mathcal{F}(\mathbf{h}^*, \mathbf{u})\| / \|\mathbf{h}^*\| \leq 10^{-5}$ . For the visualization in Fig. 1, we also apply 100 fixed-point iterations for each  $L_h$ .

## 110 C.2 Ablation Settings

111 For the ablation settings, we use the original MDEQ-Tiny [6] model (170K parameters) on the CIFAR-  
 112 10 [7] classification benchmark without any architecture modification. Therefore, the performance  
 113 gain upon state-of-the-art methods is due to the improved training efficiency thanks to the proposed  
 114 phantom gradient.

115 The experiments are conducted without data augmentation as in [6]. The training schedule, batch  
 116 size, cosine learning rate annealing strategy, and other hyperparameters are kept unchanged for all  
 117 ablation experiments. We also follow the official training protocol of MDEQ<sup>1</sup> to reproduce its results.

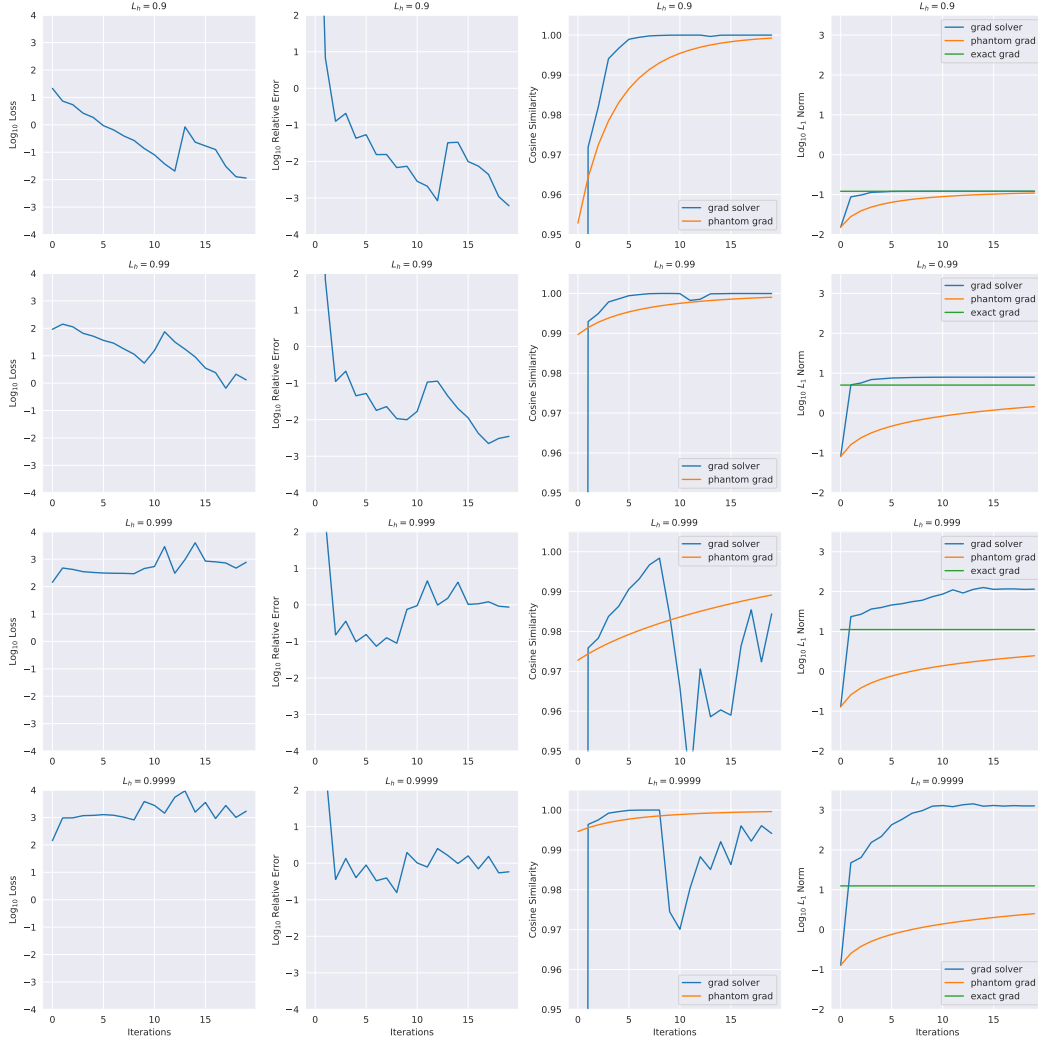


Figure 1: Visualization of gradient solver under different  $L_h$ .

118 For the training protocol without pretraining, we substitute the unrolled pretraining stage by implicit  
 119 differentiation. For the training protocol without Dropout, we remove the variational Dropout from  
 120 the model. We also try the SGD optimizer with a regular hyperparameter setting of learning rate 0.1  
 121 and weight decay 0.0001.

122 We train the MDEQ model with two types of phantom gradient using the SGD optimizer with a  
 123 learning rate of 0.1, a weight decay of 0.0001, and other hyperparameters unchanged from the original  
 124 setting. The model is trained without shallow-layer pretraining, suggesting an  $\mathcal{O}(k)$  or  $\mathcal{O}(1)$  peak

<sup>1</sup>Code available at <https://github.com/locuslab/mdeq>.



memory usage in the unrolling form or the Neumann form, as presented in the main paper. In both cases, the damped fixed-point iteration starts at the solution  $\hat{\mathbf{h}}^*$  obtained by the Broyden’s method.

### C.3 Large-Scale Experiments

For large-scale experiments, we adopt the MDEQ model (10M parameters) and MDEQ-Small model (18M parameters) on the CIFAR-10 [7] and ImageNet [8] benchmarks, respectively. To train MDEQ on CIFAR-10, we employ the unrolling-based phantom gradient with  $\lambda = 0.5$  and  $k = 5$ , *i.e.*,  $\mathbf{A}_{5,0.5}$ . Besides, we use the SGD optimizer with a learning rate of 0.1 and a weight decay 0.0001, and keep other experimental setting unchanged, including the number of training epochs, batch size, the learning rate annealing strategy, and *etc.* On the ImageNet dataset, we follow the practice of [6] to pretrain the model for the same number of epochs. Afterwards, the unrolling-based phantom gradient (*i.e.*,  $\mathbf{A}_{5,0.5}$ ) is used to train the model for the rest training schedule.

## D Additional Analysis on the Gradient Solver

To illustrate the vulnerability the gradient solver for implicit differentiation in the ill-conditioned cases, we provide the optimization dynamics in Fig. 1 and its comparison with the phantom gradient in the synthetic setting. We plot the optimization objective  $\|(I - \partial\mathcal{F}/\partial\mathbf{h})\mathbf{g} - \partial\mathcal{L}/\partial\mathbf{h}\|$ , the relative error  $\|(I - \partial\mathcal{F}/\partial\mathbf{h})\mathbf{g} - \partial\mathcal{L}/\partial\mathbf{h}\|/\|\mathbf{g}\|$ , the cosine similarity between the solved gradient  $\mathbf{g}$  and the exact gradient, and the  $L_1$  norm of the gradient. Here, in the optimization-based context,  $\mathbf{g}$  is the solution of the backward linear system solved by the Broyden solver.

Fig. 1 shows that the gradient solver may diverge in ill-conditioned situations. It can be seen that the phantom gradient demonstrates much better stability especially in the extremely ill-conditioned cases. On the contrary, the optimization step does not necessarily lead to the solved gradient being more aligned to the exact gradient, as indicated by the cosine similarity, and its angle may oscillate in the great range. Besides, the norm of the solved gradient also tends to explode in the optimization process, while the phantom gradient maintains a moderate norm throughout.

## References

- [1] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019. 1
- [2] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: A system for large-scale machine learning. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation*, page 265–283, 2016. 1
- [3] Stefan Banach. Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales. *Fund. math*, 3(1):133–181, 1922. 3
- [4] Tiberiu Popoviciu. Sur les équations algébriques ayant toutes leurs racines réelles. *Mathematica*, 9:129–145, 1935. 6
- [5] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral Normalization for Generative Adversarial Networks. In *International Conference on Learning Representations (ICLR)*, 2018. 7
- [6] Shaojie Bai, Vladlen Koltun, and J. Zico Kolter. Multiscale Deep Equilibrium Models. In *Neural Information Processing Systems (NeurIPS)*, pages 5238–5250, 2020. 8, 9
- [7] Alex Krizhevsky, Geoffrey Hinton, et al. Learning Multiple Layers of Features from Tiny Images. Technical report, Citeseer, 2009. 8, 9
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009. 9