

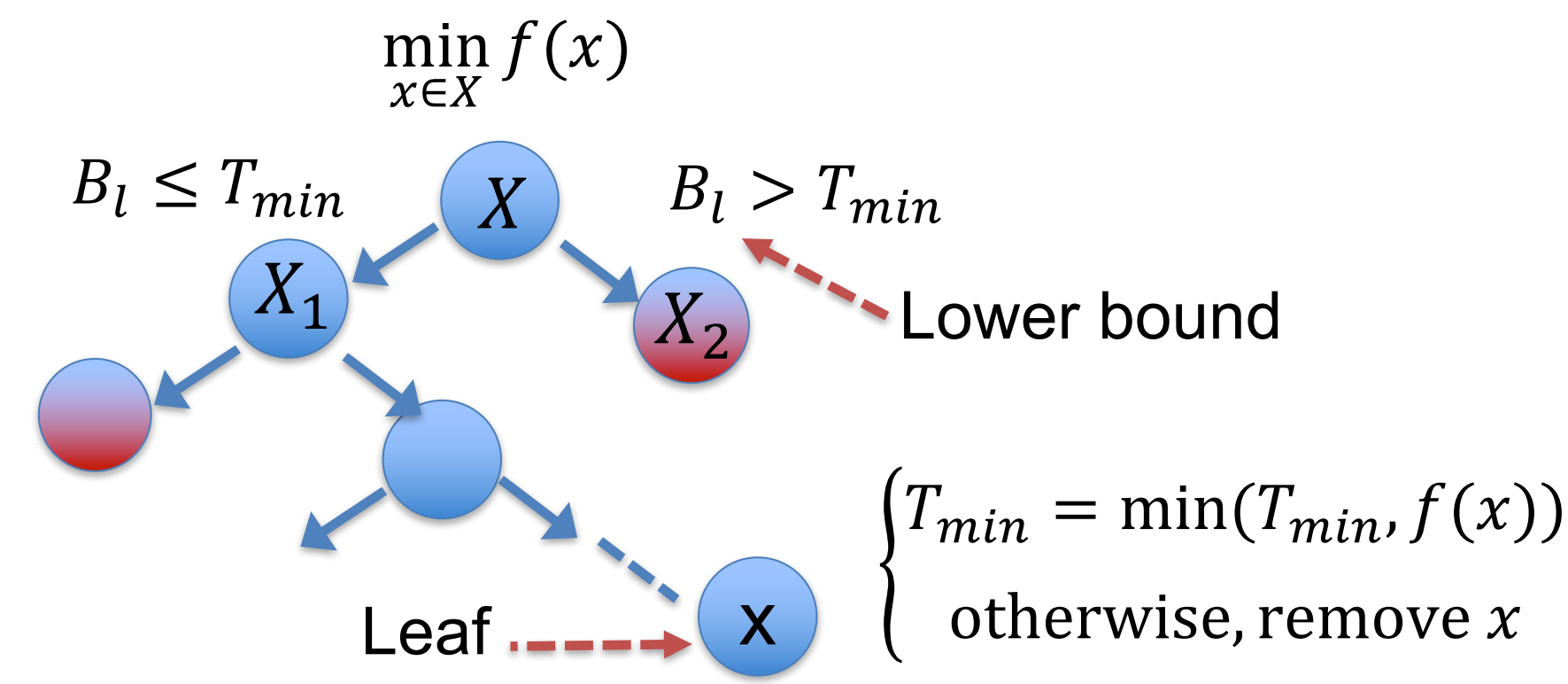
## SETTING

**Problem:** Localizing global optima in deep learning.

**Contributions:**

- 1) Explored the possibility of locating global optimality in DL from the **algorithmic** perspective;
- 2) Proposed an efficient SGD-based solver, BPGrad, with **Branch & Pruning**;
- 3) Empirically good performance of BPGrad solver on object recognition, detection, and segmentation.

**Branch & Pruning:**



## BPGRAD ALGORITHM

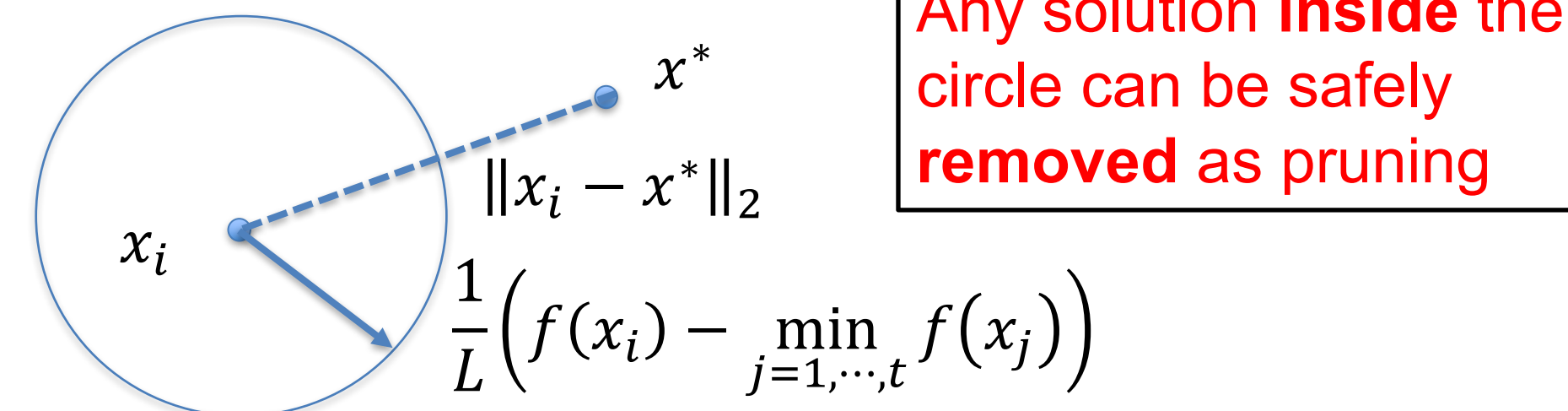
**Assumptions:**

- 1) The objective  $f$  has both lower and upper bounds;
- 2)  $f$  is differentiable in the parameter space;
- 3)  $f$  is Lipschitz continuous, or can be approximated by Lipschitz functions, with constant  $L \geq 0$ .

**Lipschitz Continuity:**

$$|f(x_1) - f(x_2)| \leq L \|x_1 - x_2\|_2, \forall x_1, x_2 \in X$$

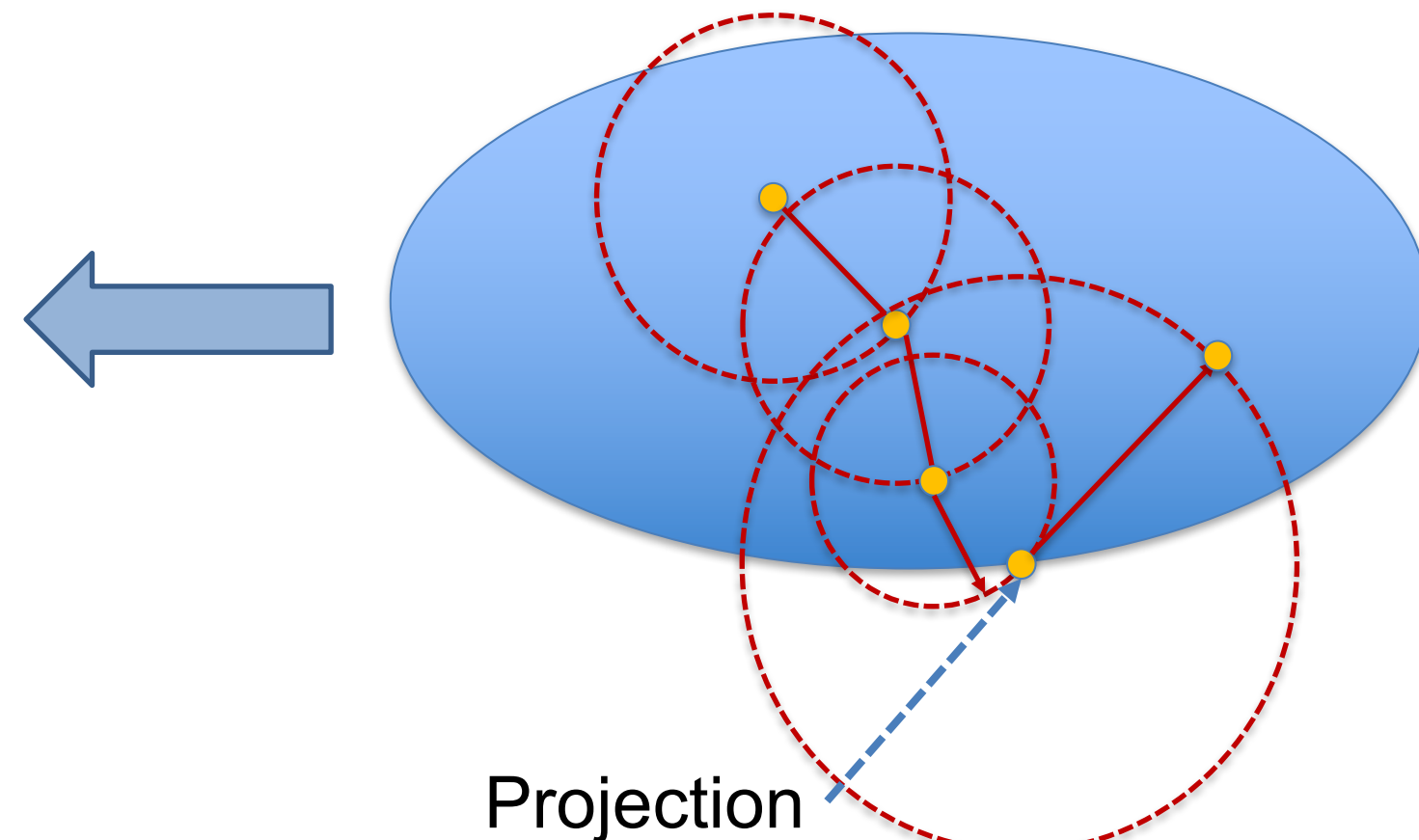
$$f(x_i) - L \|x_i - x^*\|_2 \leq f^* \leq \min_{j=1, \dots, t} f(x_j), \forall i \in [t]$$



**BPGrad Algorithm for Lipschitz functions:**

$$\min_{x_{t+1} \in X, \eta_t \geq 0} \|x_{t+1} - (x_t - \eta_t \nabla \tilde{f}(x_t))\|_2^2 + \gamma \eta_t^2, (5)$$

$$s. t. \max_{i=1, \dots, t} \{f(x_i) - L \|x_i - x_{t+1}\|_2\} \leq \rho \min_{i=1, \dots, t} f(x_i), (4)$$



Projection



- Global optimality
- Tightness bound guarantee
- Provable convergence



- Tracking of infeasible solutions
- Exponential samples

## BPGRAD FOR DEEP LEARNING

**Approximate DL solver based on BPGrad:**

- A1. Minimizing distortion is more important than minimizing step sizes, i.e.  $\gamma \ll 1$ ;
- A2.  $X$  is sufficiently large where  $\exists \eta_t \geq 0$  so that  $x_{t+1} = x_t - \eta_t \nabla \tilde{f}(x_t) \in X \setminus X_{\mathbb{R}}(t)$  always holds;
- A3.  $\eta_t \geq 0$  is always sufficiently small for local update;
- A4.  $x_{t+1}$  can be sampled only based on  $x_t$  and  $\nabla \tilde{f}(x_t)$ .

**Algorithm 2** BPGrad based Solver for Deep Learning

**Input** : number of evaluations  $n$  repeating  $N$  times at most, objective function  $f$  with Lipschitz constant  $L \geq 0$ , momentum  $0 \leq \mu \leq 1$

**Output** : minimizer  $x^*$

$t \leftarrow 1, v_1 \leftarrow 0$ , and randomly initialize  $x_1$ ;

**for**  $m \leftarrow 1$  **to**  $N$  **do**

$\rho \leftarrow 1 - \frac{1}{m}$ ;

**while**  $t < mn$  **do**

$v_{t+1} \leftarrow \mu v_t - \frac{f(x_t) - \rho \min_{i=1, \dots, t} f(x_i)}{L} \cdot \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|_2}$ ;

$x_{t+1} \leftarrow x_t + v_{t+1}$ ;

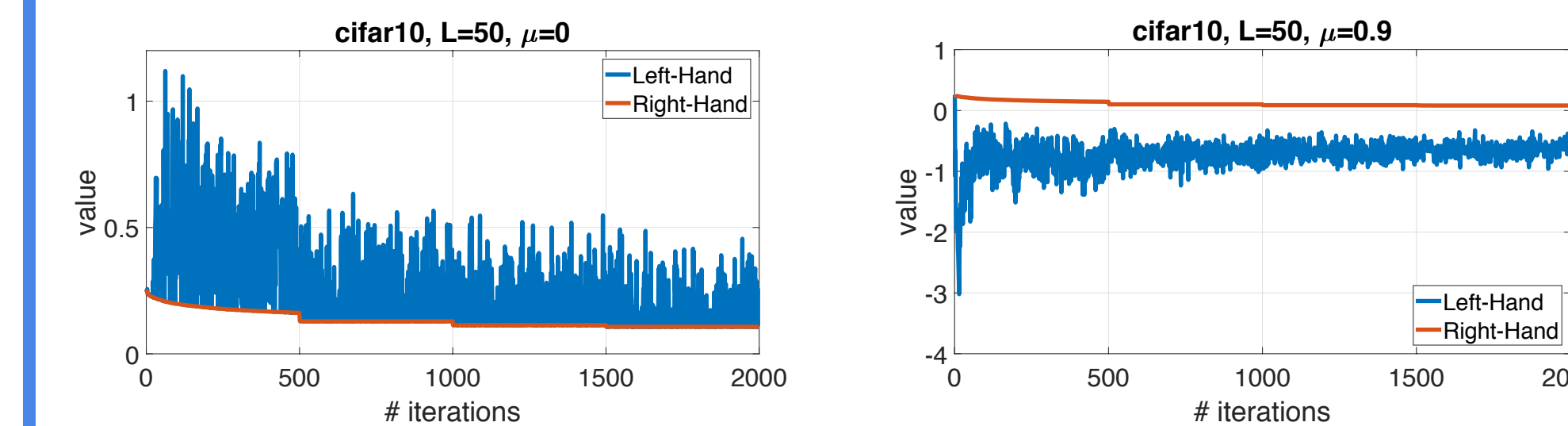
$t \leftarrow t + 1$ ;

**end** **if**  $\min_{i=1, \dots, t} f(x_i) \leq \frac{\epsilon}{1-\rho}$  **holds then** Break ;

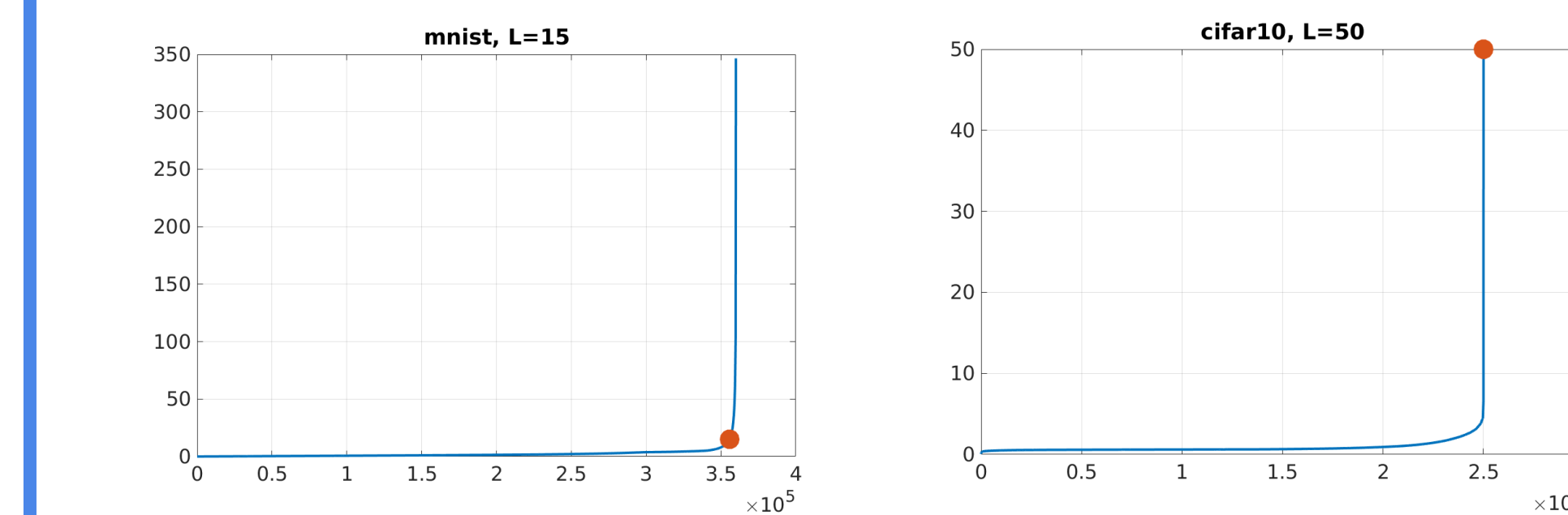
**end**

**return**  $x^* = x_{i^*}$  where  $i^* \in \arg \min_{i=1, \dots, n} f(x_i)$ ;

**Momentum helps evolve toward global optima**



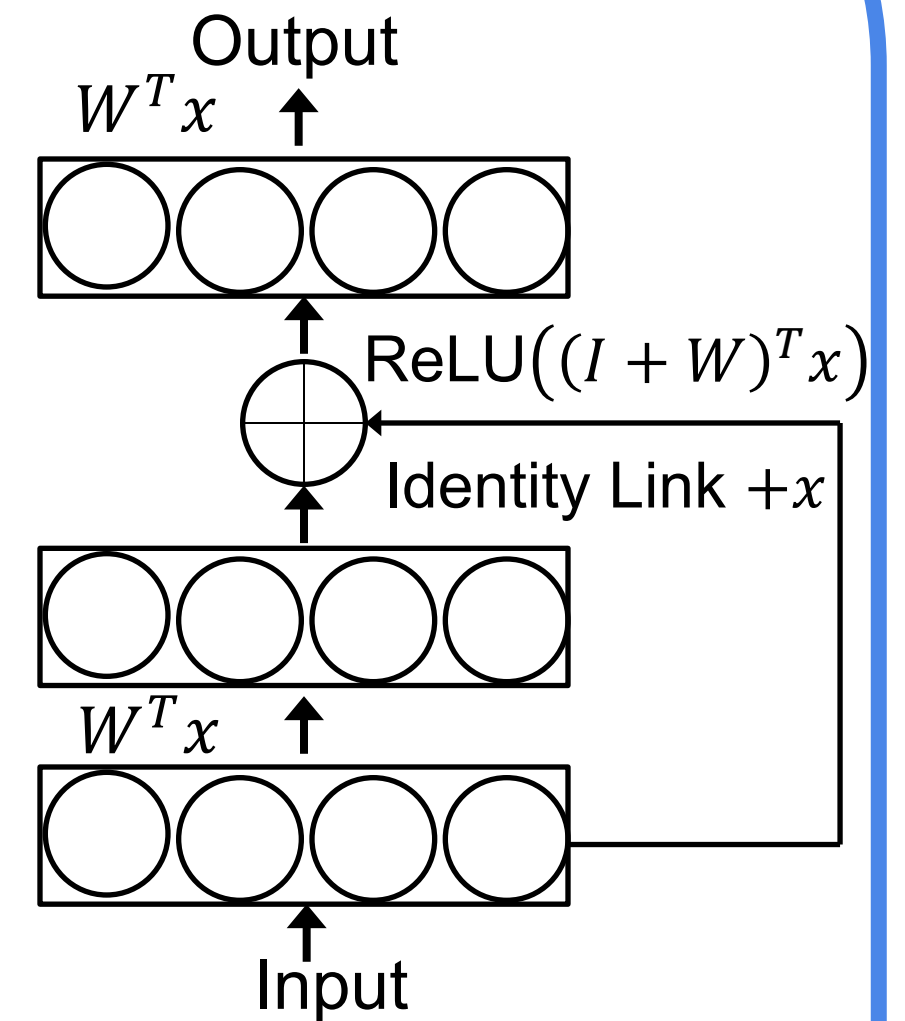
**Estimation of Lipschitz Constant  $L$ :**



## EXPERIMENTS

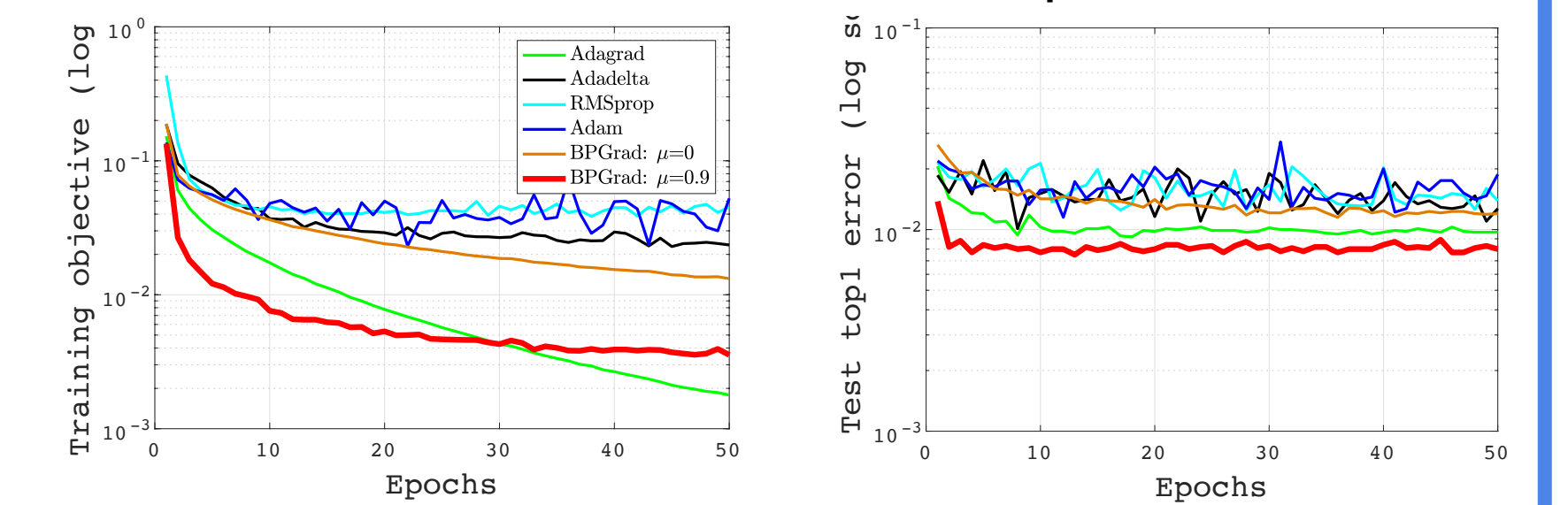
**Global Convergence of BPGrad Solver:**

- Numerical test (Li&Yuan, NIPS'17)
  - \* Two-layer neural network
  - \* 10,302 parameters
  - \* On MNIST with 20 epochs and batch size of 200
- SGD vs. BPGrad
  - \* Momentum coefficient 0.9
  - \* Best performance
  - \* Euclidean distance between the solutions is 0.6

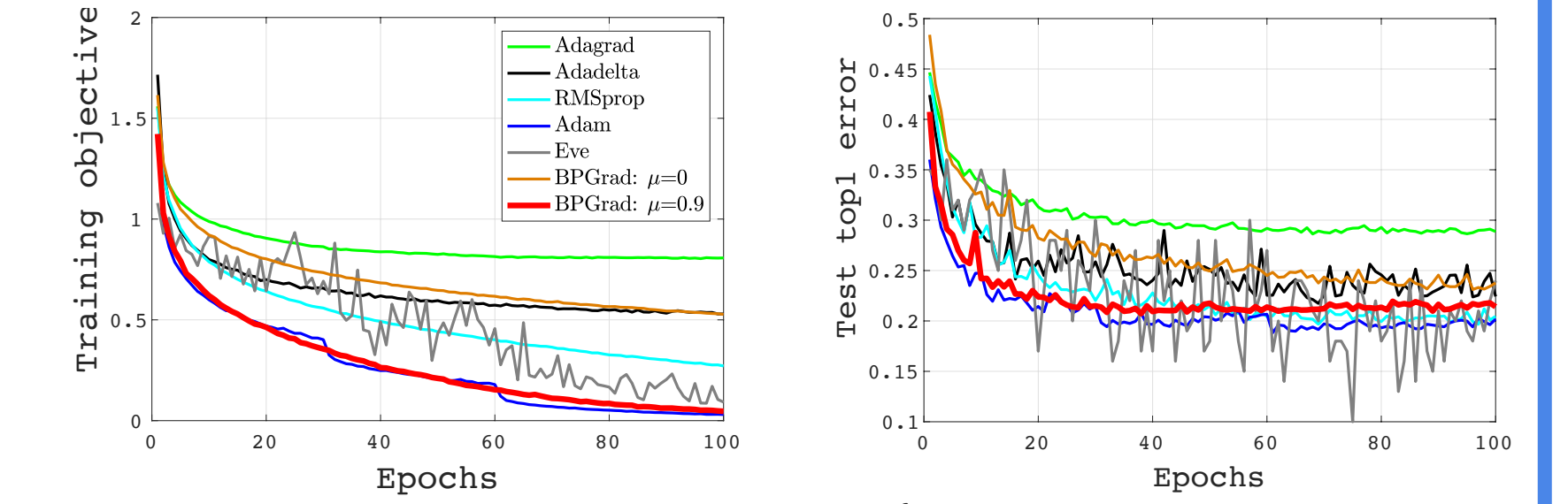


**Object Recognition:**

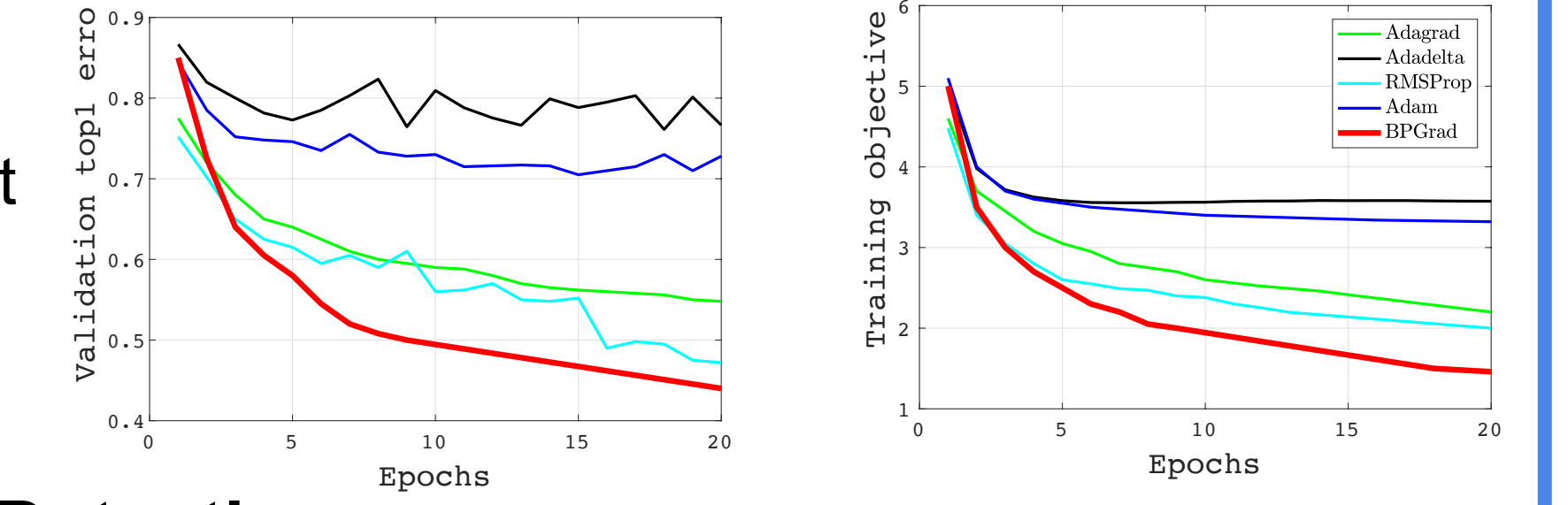
**MNIST**



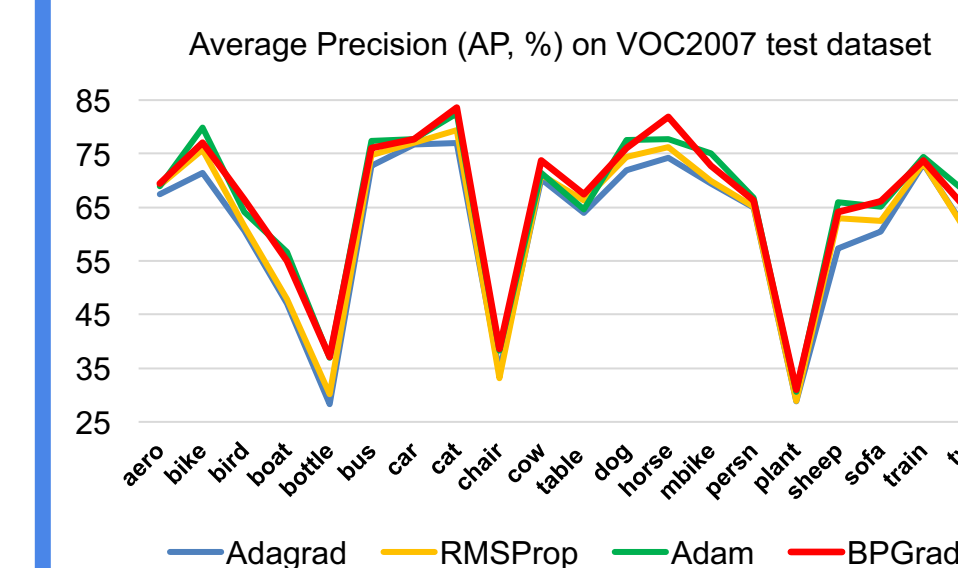
**CIFAR10**



**ImageNet**



**Object Detection:**



**Segmentation:**

	mean IU	pixel accuracy	mean accuracy	average
Adagrad	60.8	89.5	77.4	75.9
Adadelta	46.6	86.0	54.4	62.3
RMSProp	60.5	<b>90.2</b>	71.0	73.9
Adam	50.9	87.2	66.4	68.2
BPGrad	<b>62.4</b>	89.8	<b>79.6</b>	<b>77.3</b>

Table 3. Numerical comparison on semantic segmentation performance (%) using VOC2011 test dataset at the 50-th epoch.

(\* denotes equal contributions)