

# Bolstering Stochastic Gradient Descent with Model Building

Özgür Martin

Nesin Matematik Köyü  
February 10, 2023



# A "Large-Scale" Machine Learning Example

The CIFAR-10\* dataset consists of 60000 32x32 color images in 10 classes, with 6000 images per class.



bird (2)



truck (9)

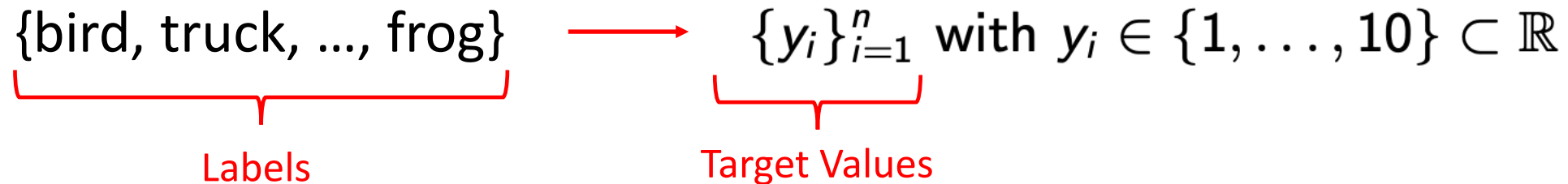
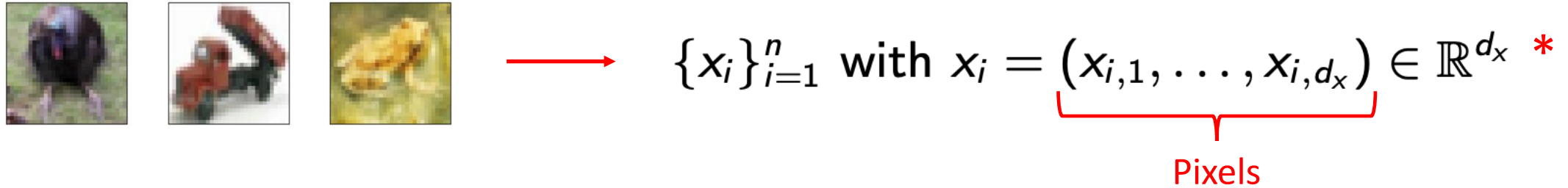


frog (6)

Image from [www.tensorflow.org](http://www.tensorflow.org)

\*Alex Krizhevsky, Learning multiple layers of features from tiny images, Tech. report, 2009.

# A "Large-Scale" Machine Learning Example



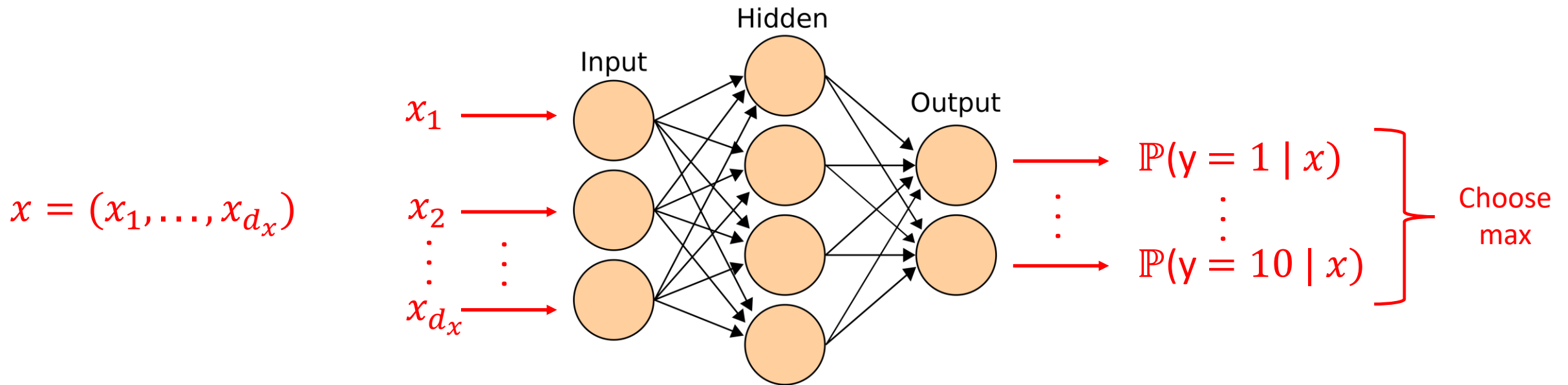
**Training points:**  $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^{d_x} \times \mathbb{R} = \mathcal{X} \times \mathcal{Y}$

Input Space      Target Space

\*  $n = 60000$ ,  $d_x = 32 \times 32 \times 3 = 3072$

# A "Large-Scale" Machine Learning Example

**Goal.** Find  $h : \mathcal{X} \rightarrow \mathcal{Y}$  s.t. for any given  $x \in \mathcal{X}$ ,  $h(x) \approx y$ , where sample  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  from a joint prob. dist. funct.  $P(x, y)$ .



An example for  $h$ : Artificial Neural Network\*

# Optimization Problem

**Goal.** Find  $h : \mathcal{X} \rightarrow \mathcal{Y}$  s.t. for any given  $x \in \mathcal{X}$ ,  $h(x) \approx y$ ,  
where sample  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  from a joint prob. dist. funct.  $P(x, y)$ .

Minimize the **expected risk**  $R(h)$ :

$$\min_h \{ R(h) := \mathbb{P}[h(x) \neq y] = \mathbb{E}[\mathbb{1}(h(x) \neq y)] \} .$$

In practice, we minimize the **empirical risk**  $R_n(h)$ :

$$\min_h \left\{ R_n(h) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i) \right\} .$$

Here,  $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$  are independently drawn.

# Optimization Problem

Thus, we solve:

$$\min_h \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i) \quad *$$

In order to make this problem easier to solve:

- Choose  $h$  from  $\mathcal{H} = \{h(\cdot; w) : w \in \mathbb{R}^d\}$
- Use a continuous (smooth) *loss* function  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$  to approximate  $\mathbb{1}$

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w), y_i)$$

\*  $\mathbb{1}[A] = 1$  if  $A$  is true, and  $\mathbb{1}[x] = 0$  otherwise.

As before, choose  $h$  from  $\mathcal{H} = \{h(\cdot; w) : w \in \mathbb{R}^d\}$ .

Choose continuous (smooth) loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ .

**Expected Risk:** 
$$\min_{w \in \mathbb{R}^d} \left\{ R(w) := \int_{\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}} \ell(h(x; w), y) dP(x, y) = \mathbb{E}[\ell(h(x; w), y)] \right\}$$

**Empirical Risk:** 
$$\min_{w \in \mathbb{R}^d} \left\{ R_n(w) := \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w), y_i) \right\}$$

# Minimizing the Empirical Risk

$$\min_{w \in \mathbb{R}^d} \left\{ F(w) := \frac{1}{n} \sum_{i=1}^n f_i(w) \right\},$$

where  $f_i(w) = \ell(h(x_i; w), y_i)$ .

So, we have a finite sum optimization problem to solve.



# Generalization: Train and Test Data

**Sample points:**  $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^{d_x} \times \mathbb{R}$

$$A \subset \{1, \dots, n\}$$

$$B = \{1, \dots, n\} / A$$

Training data

$$\{(x_i, y_i)\}_{i \in A}$$

Test data

$$\{(x_i, y_i)\}_{i \in B}$$

Solve:

$$\min_w \frac{1}{|A|} \sum_{i \in A} f_i(w)$$

Training  
Loss  
Function

Check:

$$\min_w \frac{1}{|B|} \sum_{i \in B} f_i(w)$$

Test  
Loss  
Function

# SGD Method

Given the training points  $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ , solve

$$\min_{w \in \mathbb{R}^d} \left\{ f(w) := \frac{1}{n} \sum_{i=1}^n f_i(w) \right\},$$

where  $f_i$  corresponds to the loss function evaluated at the training point  $(x_i, y_i)$ .

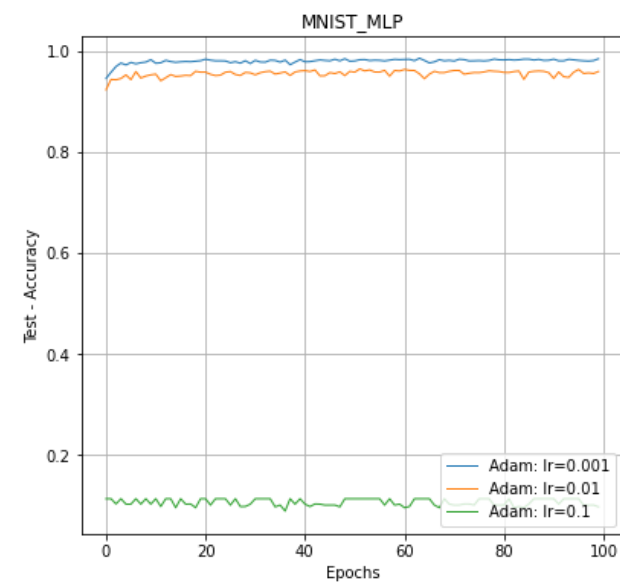
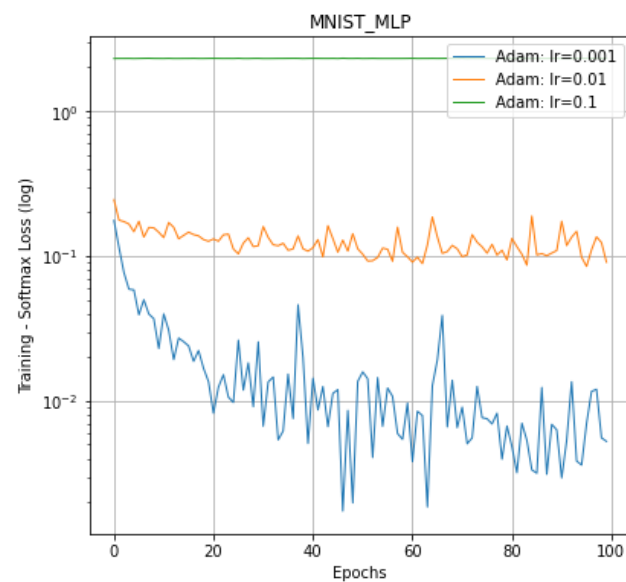
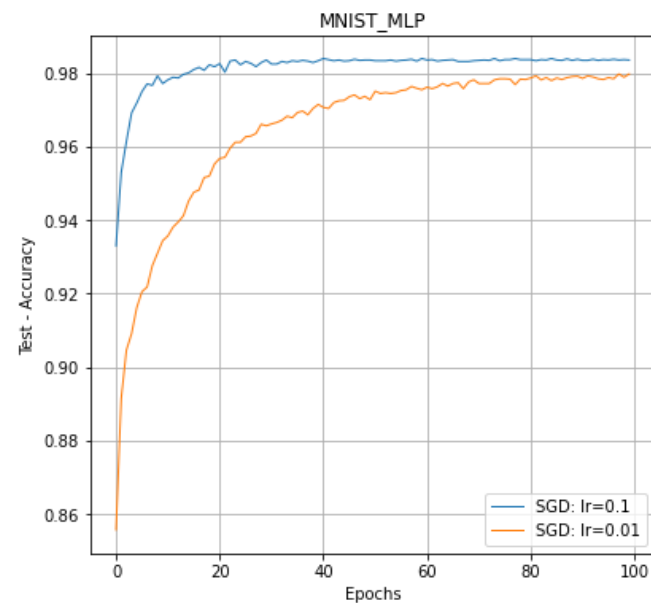
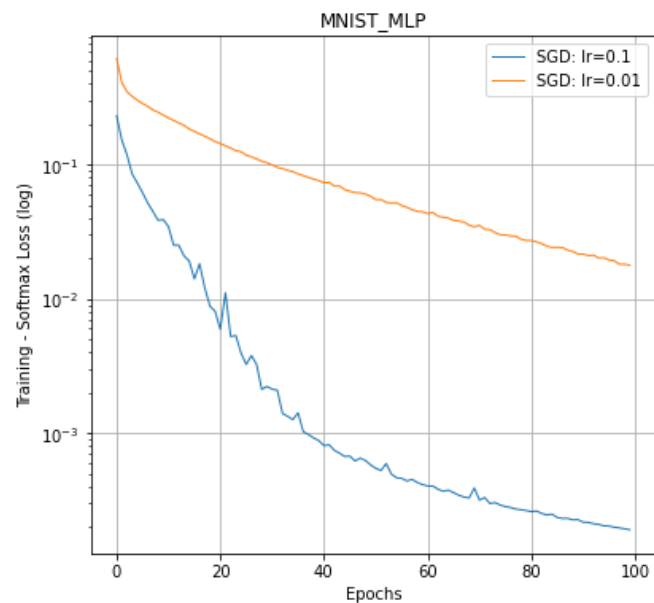
---

**Algorithm 2:** Stochastic Gradient Descent (SGD)

---

- 1 **Input:** Initial point  $w_1$ , max iteration number  $T$ , and the stepsizes  $\{\alpha_k\}_{k=1}^T$
  - 2 **for**  $k = 1, \dots, T$  **do**
  - 3     Choose  $i_k \in \{1, \dots, n\}$  randomly w.r.t. uniform distribution;
  - 4      $w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$
-

# Robustness: SGD, Adam



# Stochastic Line Search\*

**Armijo line search:** Given the maximum stepsize  $\alpha_{max}$ , find the biggest stepsize  $0 < \alpha_k \leq \alpha_{max}$  which satisfies the stochastic Armijo condition

$$f_{i_k}(w_k - \alpha_k \nabla f_{i_k}(w_k)) - f_{i_k} \leq c \alpha_k \|\nabla f_{i_k}(w_k)\|^2,$$

where  $c > 0$  is a hyper-parameter.

---

**Algorithm 3:** Stochastic Line Search (SLS)

---

```
1 Input: Initial point  $w_1$ , max iteration number  $T$ , maximum stepsize  $\alpha_{max}$ ,  $c$ 
2 for  $k = 1, \dots, T$  do
3   Choose  $i_k \in \{1, \dots, n\}$  randomly w.r.t. uniform distribution;
4   Find  $\alpha_k$  with Armijo line search;
5    $w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$ 
```

---

\*Vaswani et al., Painless Stochastic Gradient: Interpolation, Line-Search, and Convergence Rates, NeurIPS 2019.

# Stochastic Line Search

**Backtracking line search:** Start with  $\alpha_k = \alpha_{max}$  and repeatedly set  $\alpha_k = \alpha_k \cdot \gamma$  until the stochastic Armijo condition is satisfied. Here,  $0 < \gamma < 1$  is a hyper-parameter.

---

**Algorithm 4:** Stochastic Line Search (SLS)

---

```
1 Input: Initial point  $w_1$ , max iteration number  $T$ , and maximum stepsize  $\alpha_{max}$ ,  $c$ ,  $\gamma$ ,  $N$ 
2 for  $k = 1, \dots, T$  do
3   Choose  $i_k \in \{1, \dots, n\}$  randomly w.r.t. uniform distribution;
4   Find  $\alpha_k$  with Armijo line search and backtracking;
5    $w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$ ;
6    $\alpha_{max} = \alpha_k$ ;
7   RESET  $\alpha_{max}$  if  $N \mid k$ 
```

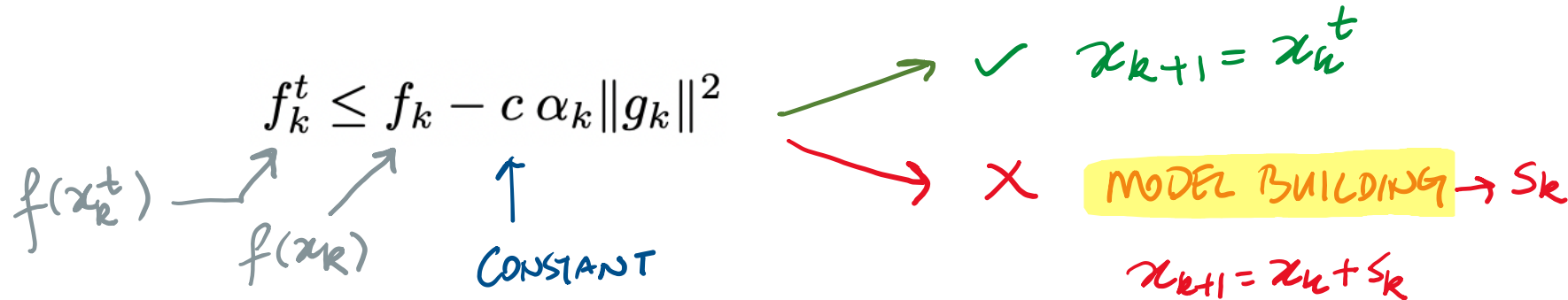
---

# Stochastic Line Search with Model Building\*

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x)$$

**TRIAL STEP**  $x_k^t = x_k + s_k^t$

**SGD**  $s_k^t = -\alpha_k g_k$



\*Birbil, Martin, Onay, and Öztoprak, Bolstering Stochastic Gradient Descent with Model Building, 2021.

$$x_k^t = x_k + s_k^t$$

$$s_k^t = -\alpha_k g_k$$

$$f_k^t \leq f_k - c \alpha_k \|g_k\|^2$$

✓

$$x_{k+1} = x_k^t$$

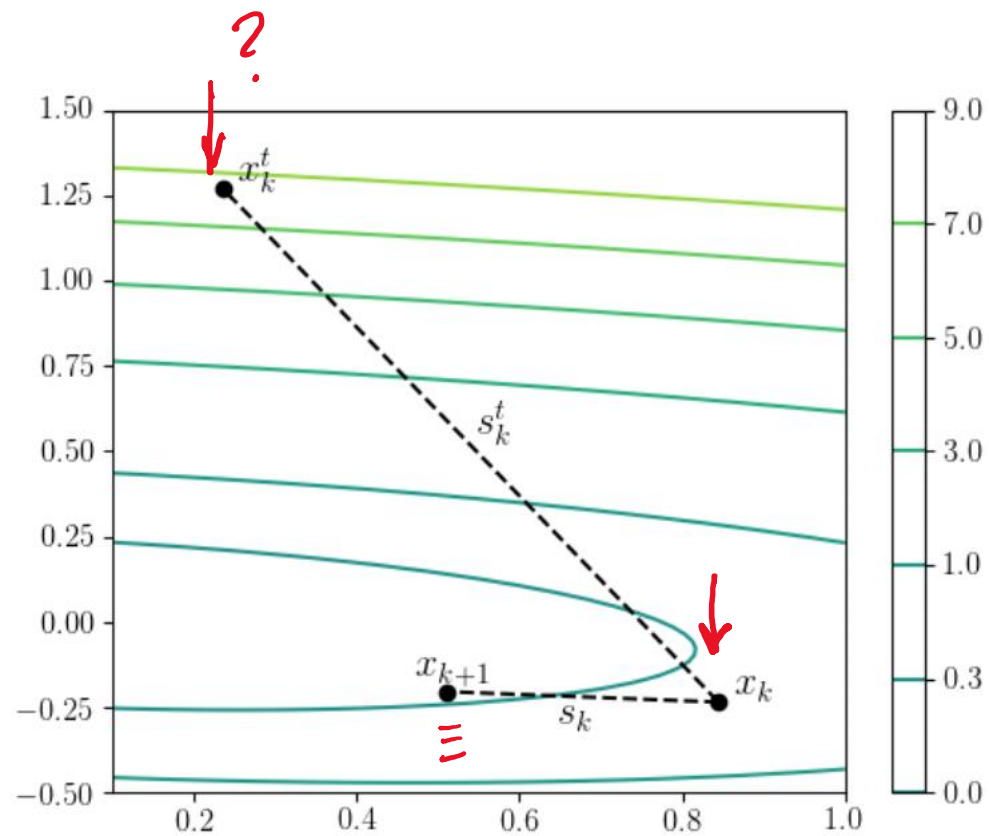
✗

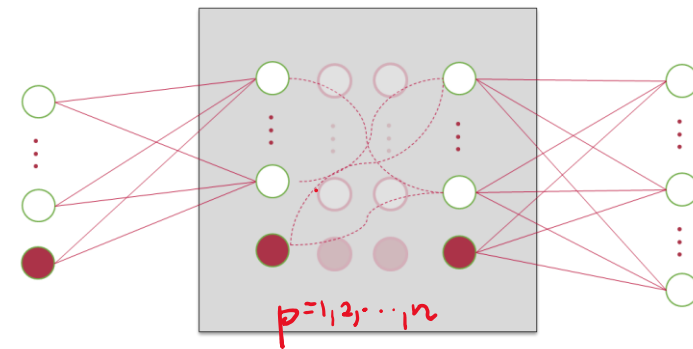
MODEL BUILDING

↓

$s_k$

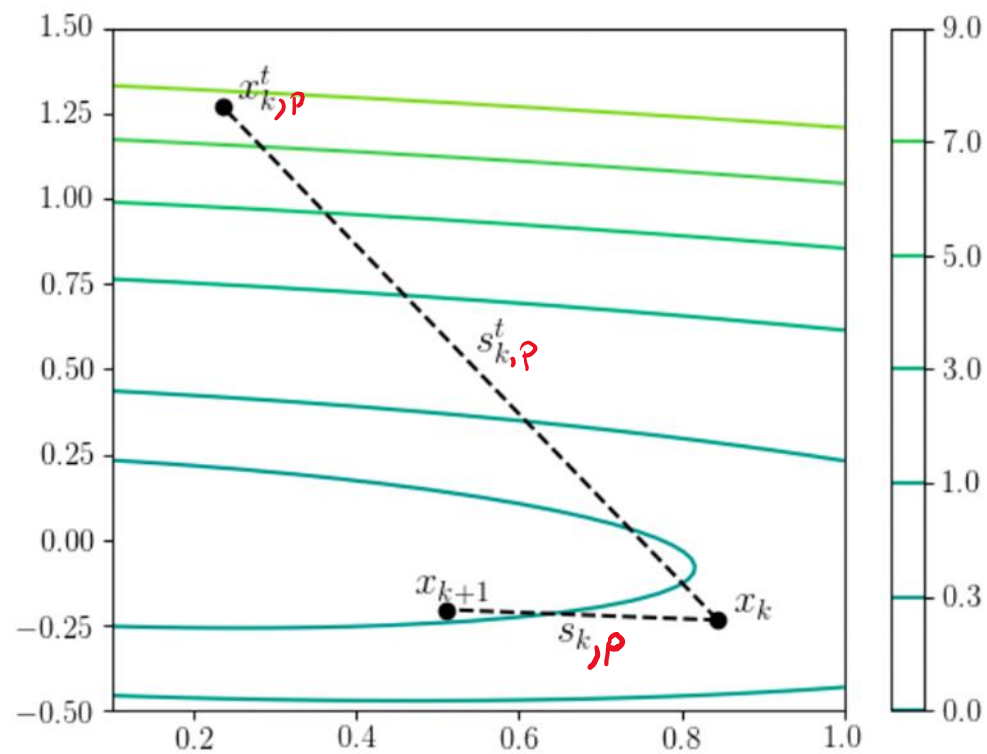
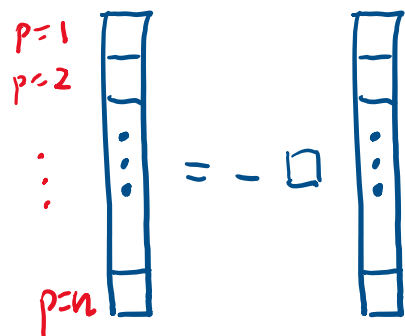
$$x_{k+1} = x_k + s_k$$





$$x_{k,p}^t = x_{k,p} + s_{k,p}^t$$

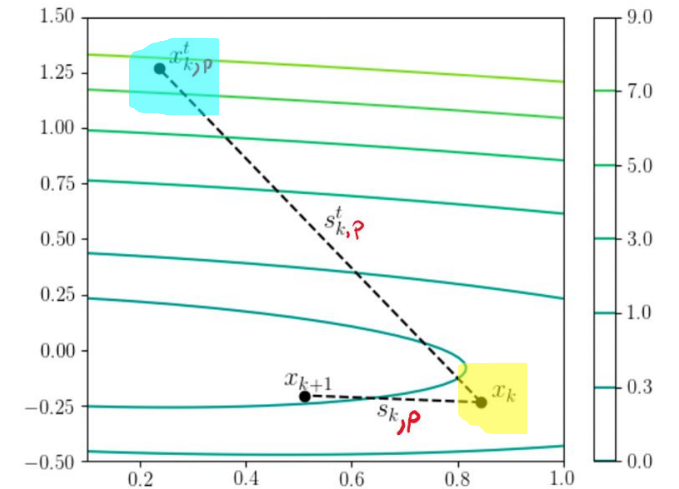
$$s_{k,p}^t = -\alpha_k g_{k,p}$$





# MODEL BUILDING

$$m_{k,p}^t(s) := \underbrace{\alpha_{k,p}^0(s)}_{\text{red}} l_{k,p}^0(s) + \underbrace{\alpha_{k,p}^t(s)}_{\text{green}} l_{k,p}^t(s - s_{k,p}^t)$$



$$l_{k,p}^0(s) := f_{k,p} + g_{k,p}^\top s$$

$$l_{k,p}^t(s - s_{k,p}^t) := f_{k,p} + (g_{k,p}^t)^\top (s - s_{k,p}^t)$$

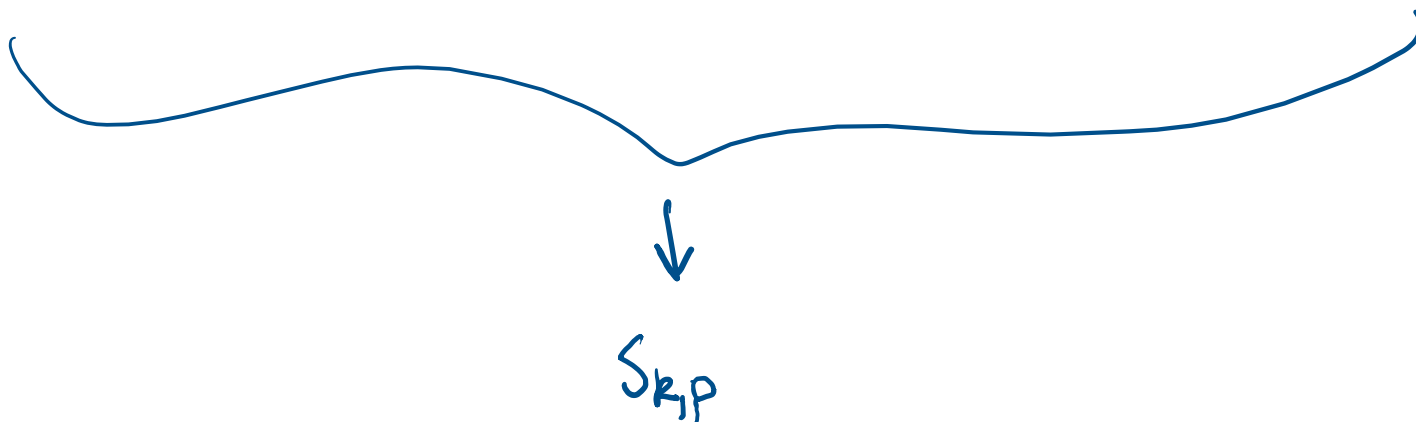
$$\underbrace{\alpha_{k,p}^0(s)}_{\text{red}} = \frac{(s - s_{k,p}^t)^\top (-s_{k,p}^t)}{(-s_{k,p}^t)^\top (-s_{k,p}^t)}$$

$$\underbrace{\alpha_{k,p}^t(s)}_{\text{green}} = \frac{s^\top s_{k,p}^t}{(s_{k,p}^t)^\top s_{k,p}^t}$$

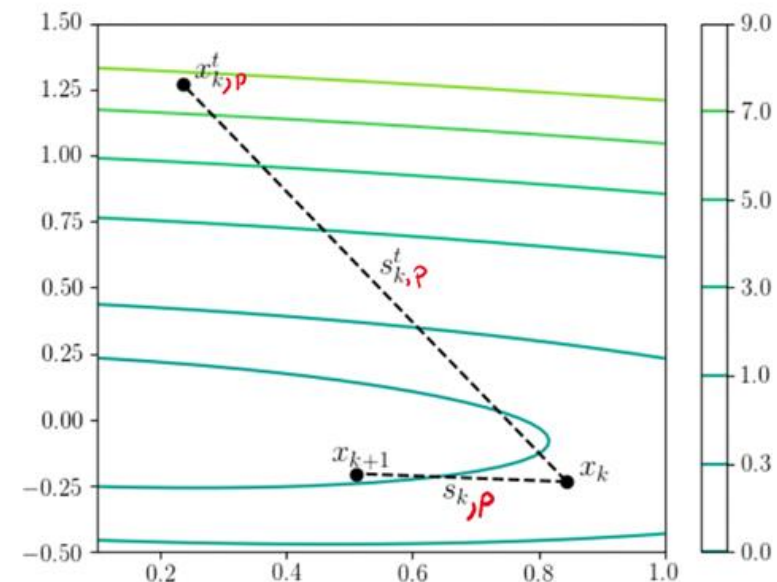
\* Figen Öztoprak and Ş İlker Birbil. An alternative globalization strategy for unconstrained optimization. *Optimization*, 67(3):377–392, 2018.

MINIMIZE  $m_{k,p}^t(s) := \alpha_{k,p}^0(s)l_{k,p}^0(s) + \alpha_{k,p}^t(s)l_{k,p}^t(s - s_{k,p}^t)$

SUBJECT TO  $\|s\|^2 + \|s - s_{k,p}^t\|^2 \leq \|s_{k,p}^t\|^2$



(ANALYTICAL SOLUTION)



$$s_k^{t+1} = c_g(\sigma)g_k + c_y(\sigma)y_k^t + c_s(\sigma)s_k^t,$$

where

$$c_g(\sigma) = -\frac{\|s_k^t\|^2}{2\sigma}, \quad c_y(\sigma) = -\frac{\|s_k^t\|^2}{2\sigma\theta}[-((y_k^t)^\top s_k^t + 2\sigma)((s_k^t)^\top g_k) + \|s_k^t\|^2((y_k^t)^\top g_k)],$$

$$c_s(\sigma) = -\frac{\|s_k^t\|^2}{2\sigma\theta}[-((y_k^t)^\top s_k^t + 2\sigma)((y_k^t)^\top g_k) + \|y_k^t\|^2((s_k^t)^\top g_k)],$$

with

$$\theta = ((y_k^t)^\top s_k^t + 2\sigma)^2 - \|s_k^t\|^2\|y_k^t\|^2,$$

and

$$\sigma = \frac{1}{2} \left( \|s_k^t\| \left( \|y_k^t\| + \frac{1}{\eta} \|g_k\| \right) - (y_k^t)^\top s_k^t \right).$$

---

**Algorithm 1:** SMB: Stochastic Model Building

---

```
1 Input:  $x_1 \in \mathbb{R}^n$ , stepsizes  $\{\alpha_k\}_{k=1}^T$ , mini-batch sizes  $\{m_k\}_{k=1}^T, c > 0$ , and  $\alpha_{max}$  satisfying (8)
2 for  $k = 1, \dots, T$  do
3    $f_k = f(x_k, \xi_k)$ ,  $g_k = \frac{1}{m_k} \sum_{i=1}^{m_k} g(x_k, \xi_{k,i})$ ;
4    $s_k^t = -\alpha_k g_k$ ;
5    $x_k^t = x_k + s_k^t$ ;
6    $f_k^t = f(x_k^t, \xi_k)$ ,  $g_k^t = \frac{1}{m_k} \sum_{i=1}^{m_k} g(x_k^t, \xi_{k,i})$ ;
7   if  $f_k^t \leq f_k - c \alpha_k \|g_k\|^2$  then
8      $x_{k+1} = x_k^t$ ;
9   else
10    for  $p = 1, \dots, r$  do
11       $y_{k,p} = g_{k,p}^t - g_{k,p}$ ;
12       $s_{k,p} = c_{g,p}(\delta) g_{k,p} + c_{y,p}(\delta) y_{k,p} + c_{s,p}(\delta) s_{k,p}^t$ ;
13     $x_{k+1} = x_k + s_k$  with  $s_k = (s_{k,p_1}, \dots, s_{k,p_r})$ ;
```

---

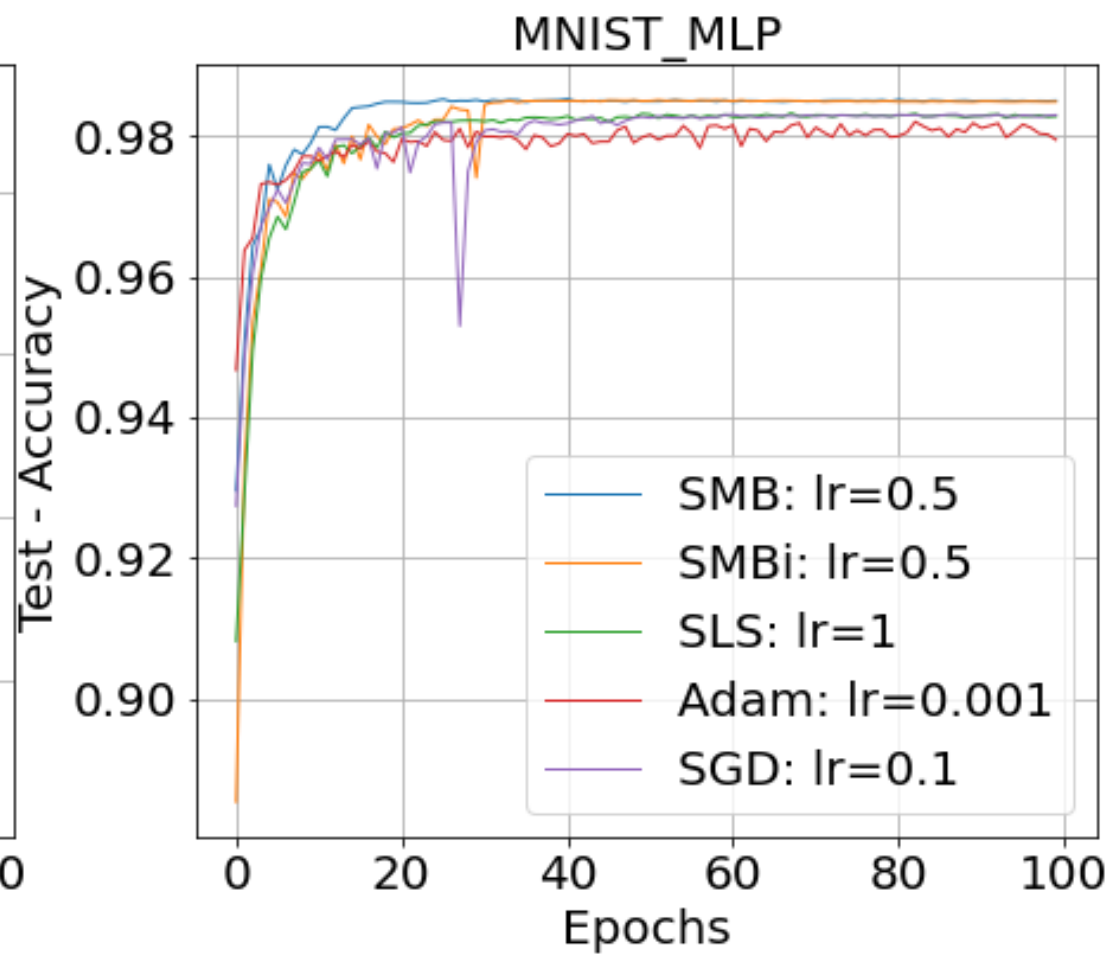
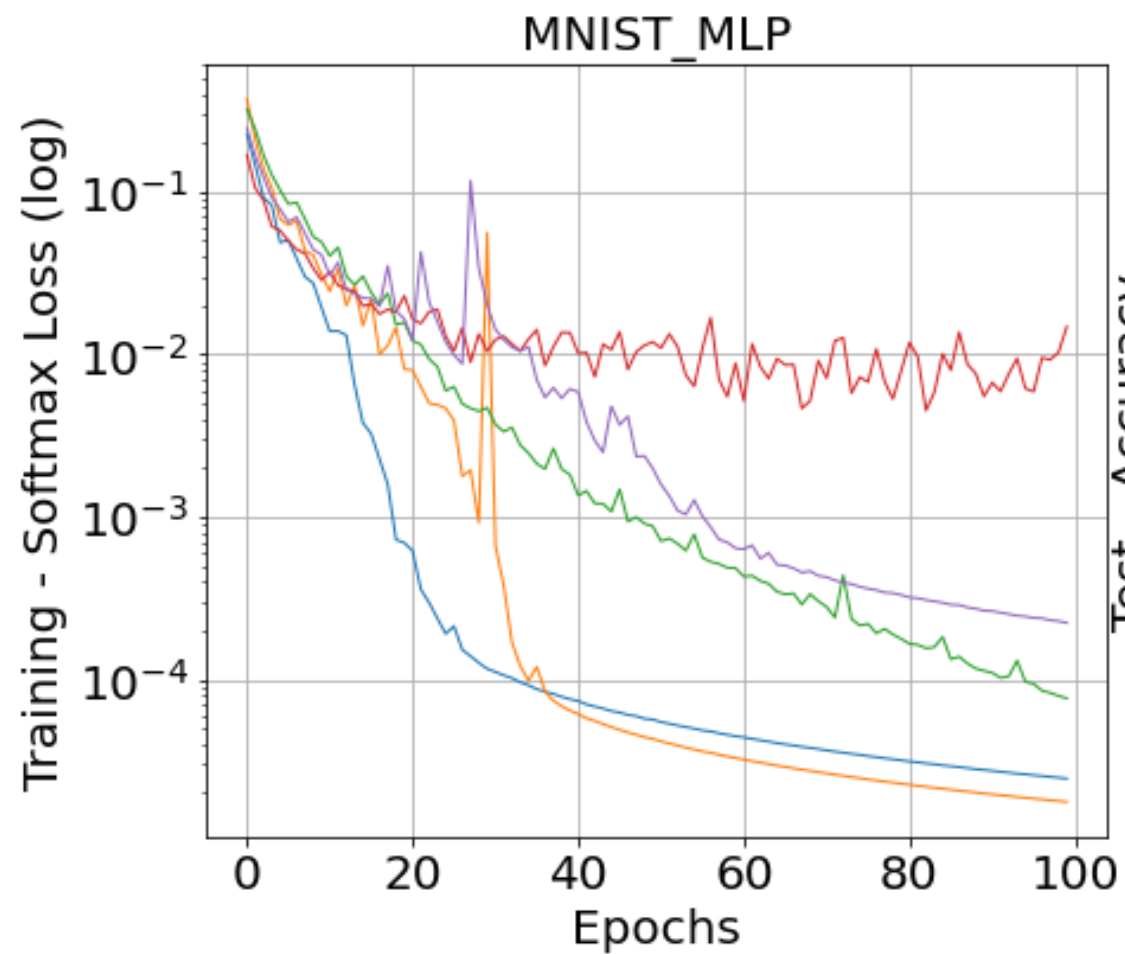
$$x_{k+1} = x_k - \alpha_k H_k g_k, \quad (s_k = -\alpha_k H_k g_k)$$

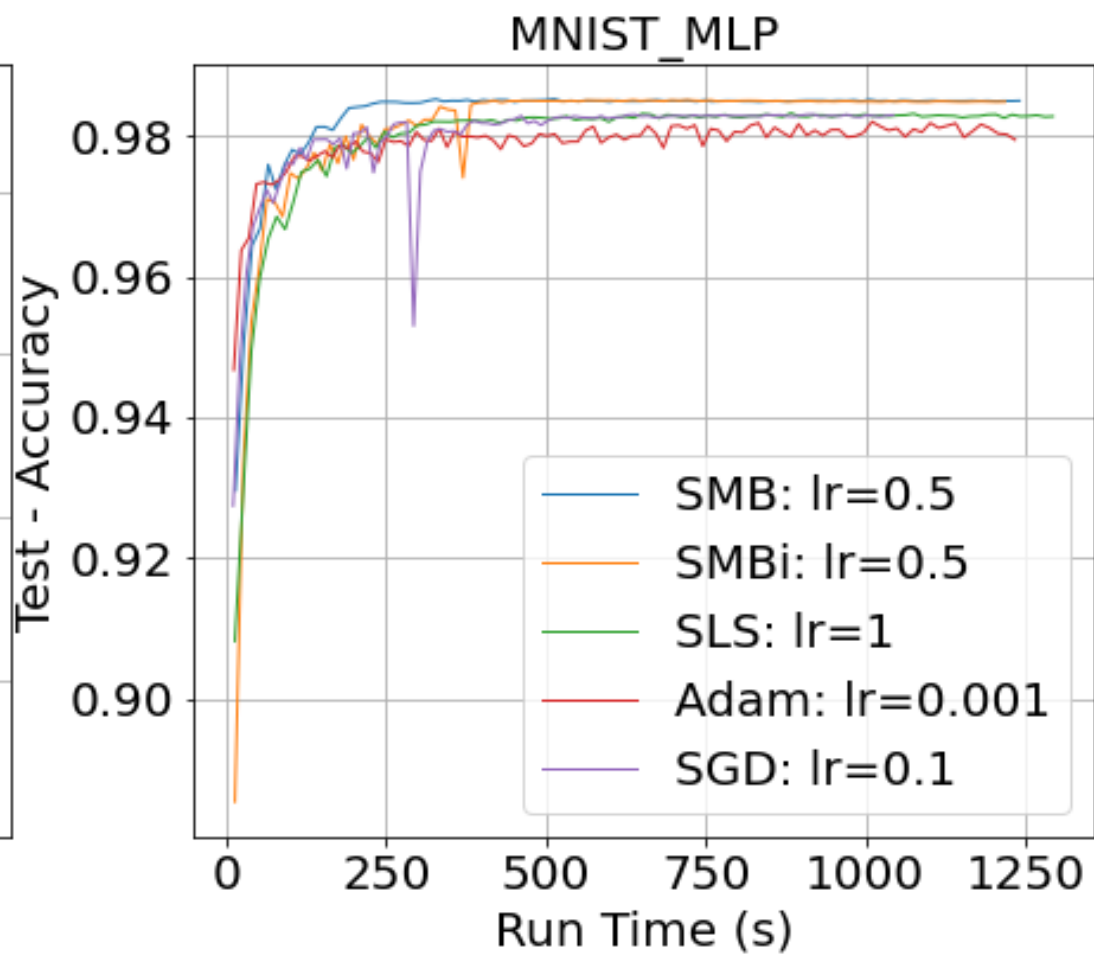
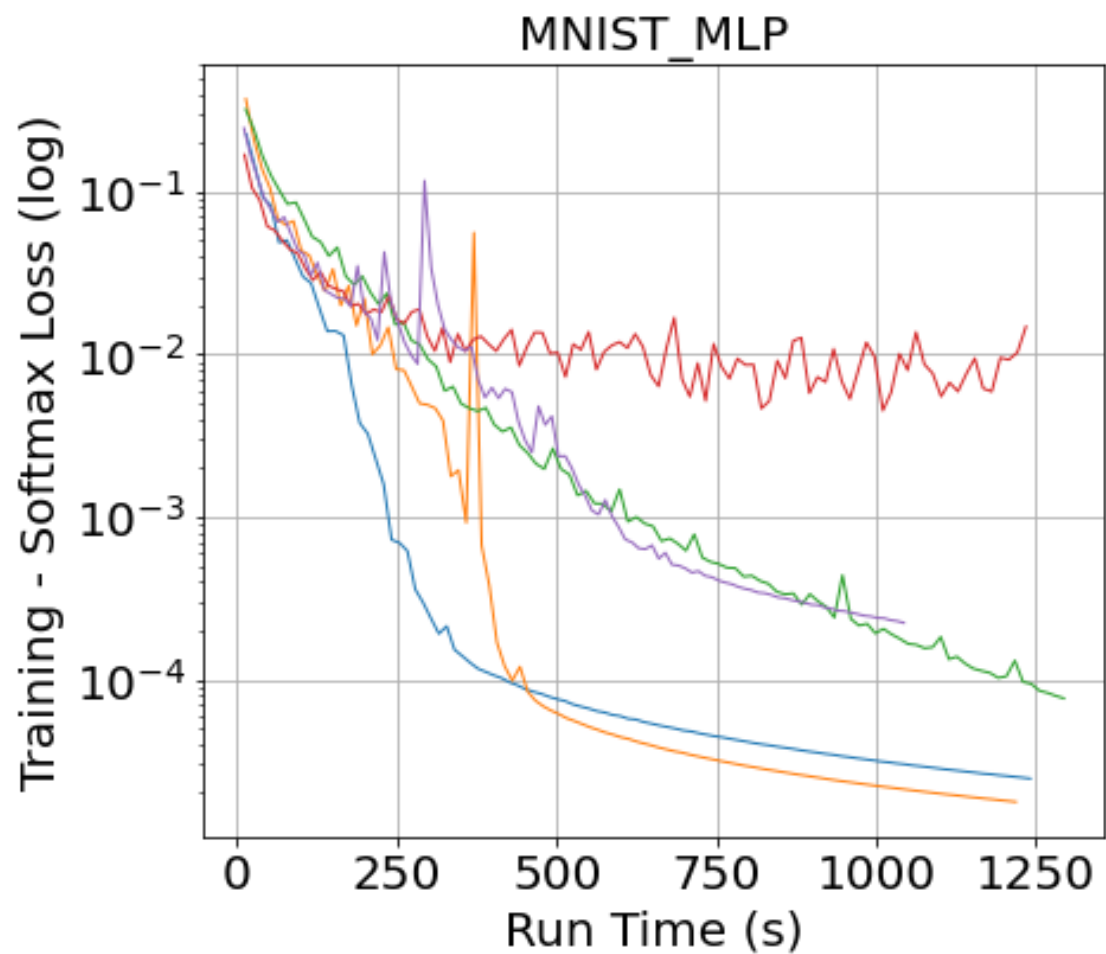
$$H_k = \begin{bmatrix} H_{k,1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & H_{k,n} \end{bmatrix}$$

$$H_{k,p} = \frac{\|g_{k,p}\|^2}{\sigma_p \gamma_p} [\gamma_p I + \beta_p y_{k,p} g_{k,p}^\top + \|g_{k,p}\|^2 y_{k,p} y_{k,p}^\top + \beta_p g_{k,p} y_{k,p}^\top + \|y_{k,p}\|^2 g_{k,p} g_{k,p}^\top],$$

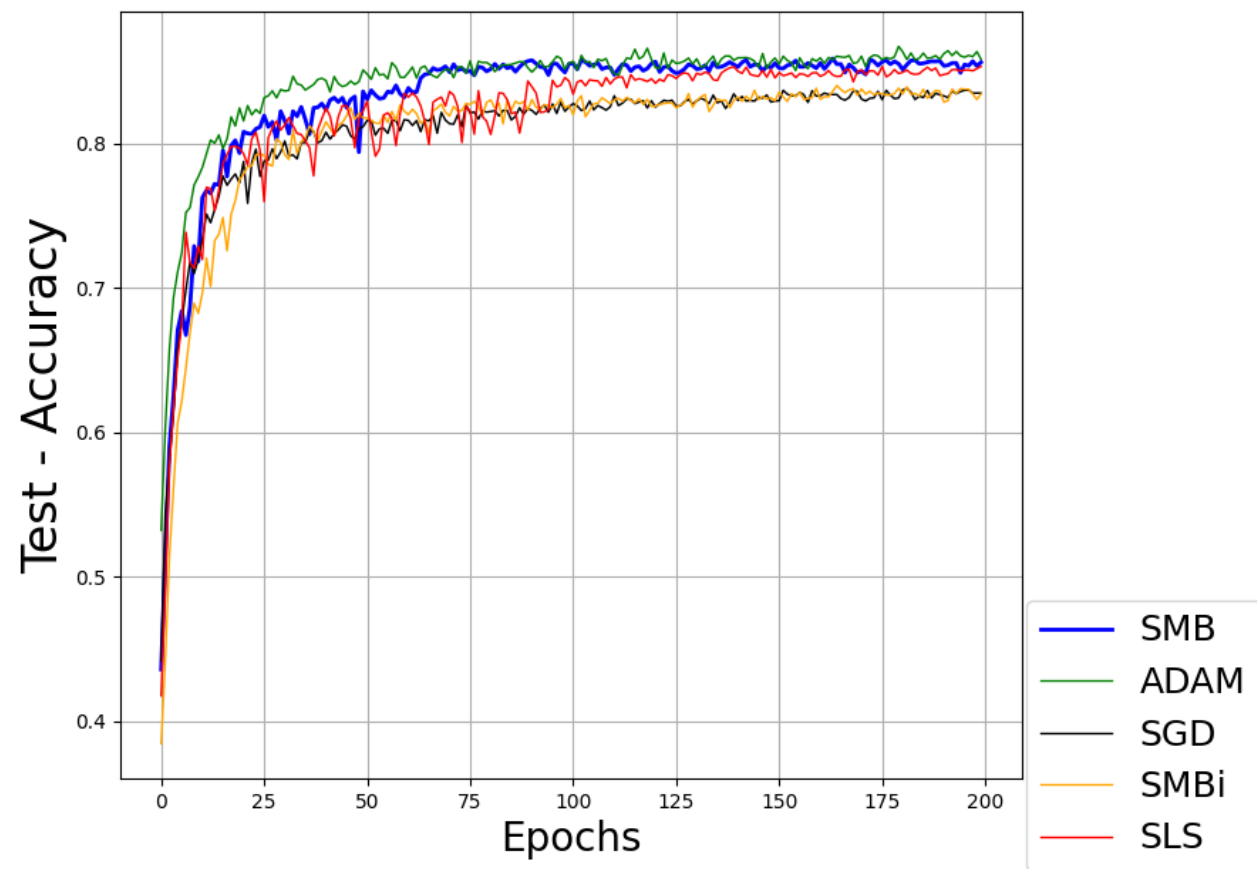
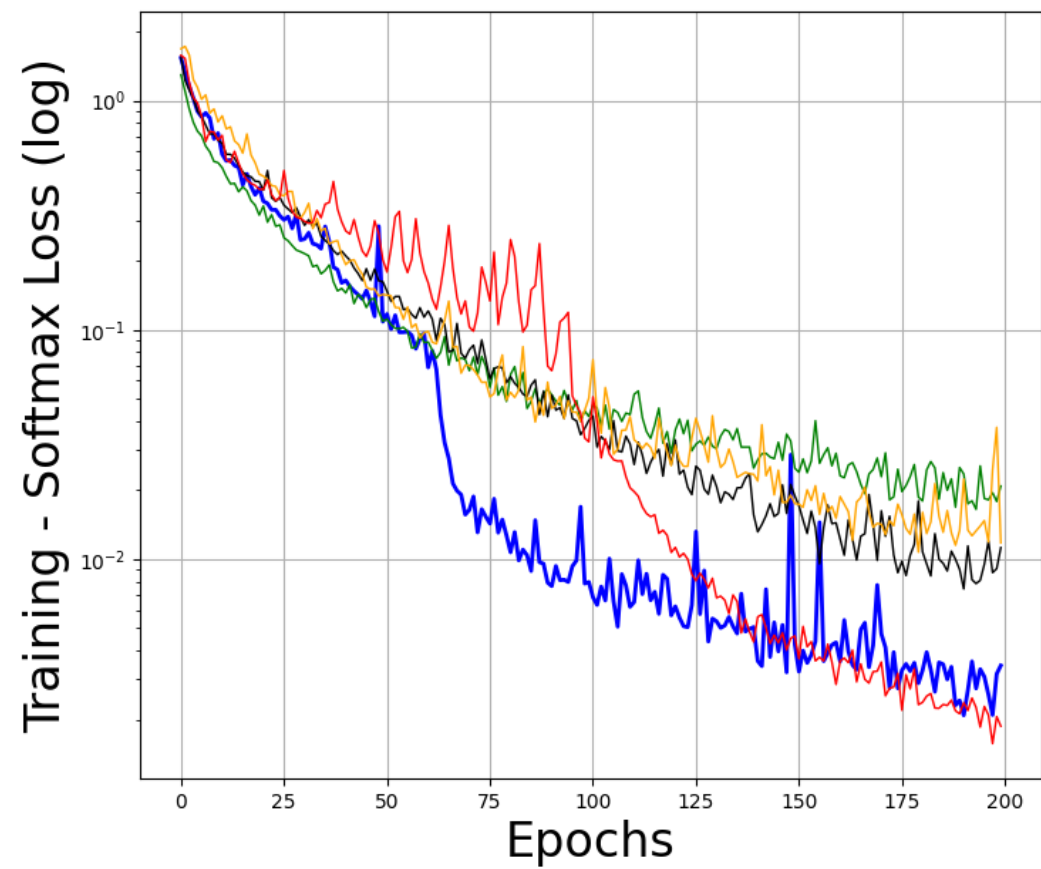
where

$$\sigma_p = \|g_{k,p}\| \|y_{k,p}\| + \frac{1}{\eta} \|g_{k,p}\|^2 + y_{k,p}^\top g_{k,p}, \quad \beta_p = \sigma_p - y_{k,p}^\top g_{k,p}, \quad \text{and} \quad \gamma_p = (\beta_p^2 - \|g_{k,p}\|^2 \|y_{k,p}\|^2).$$



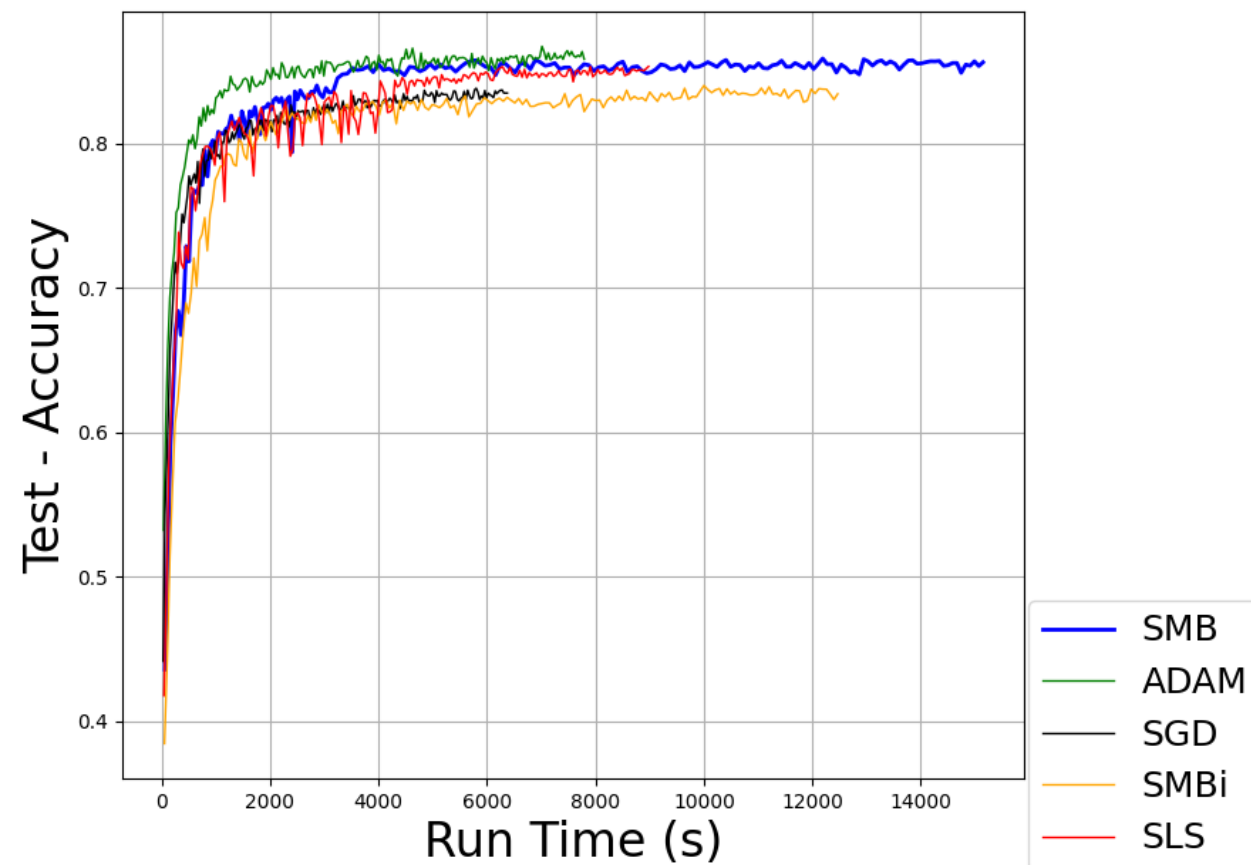
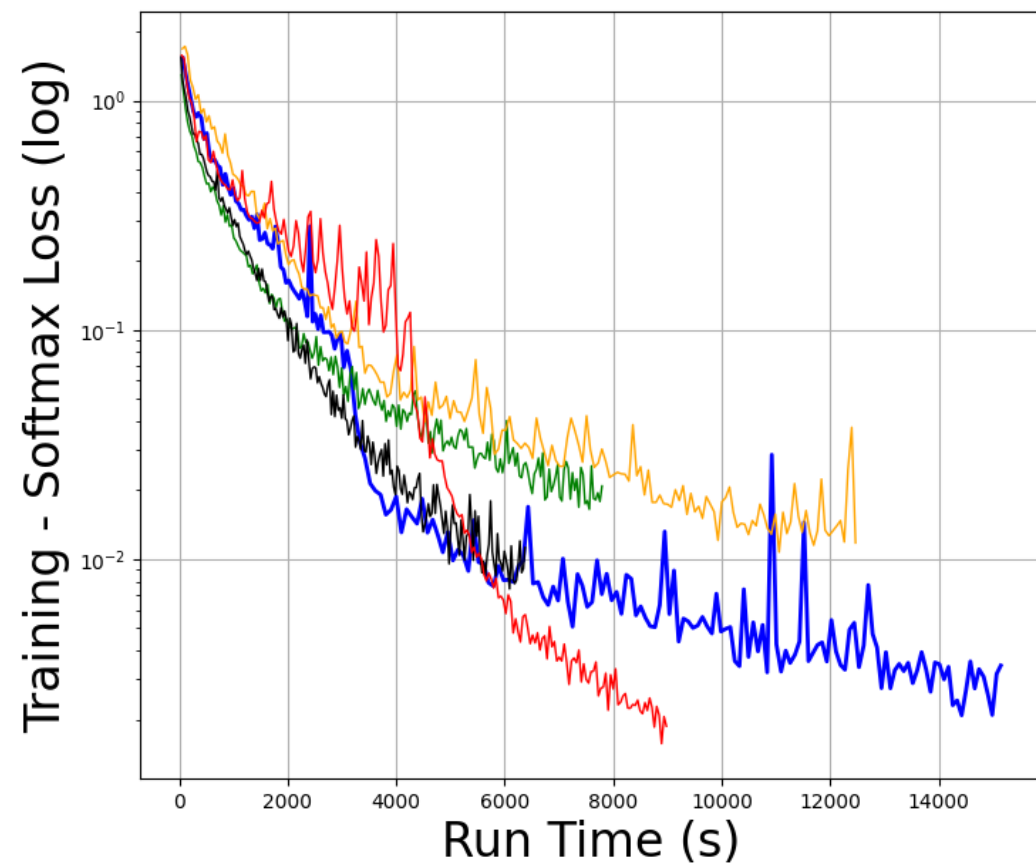


## CIFAR10-DENSENET10

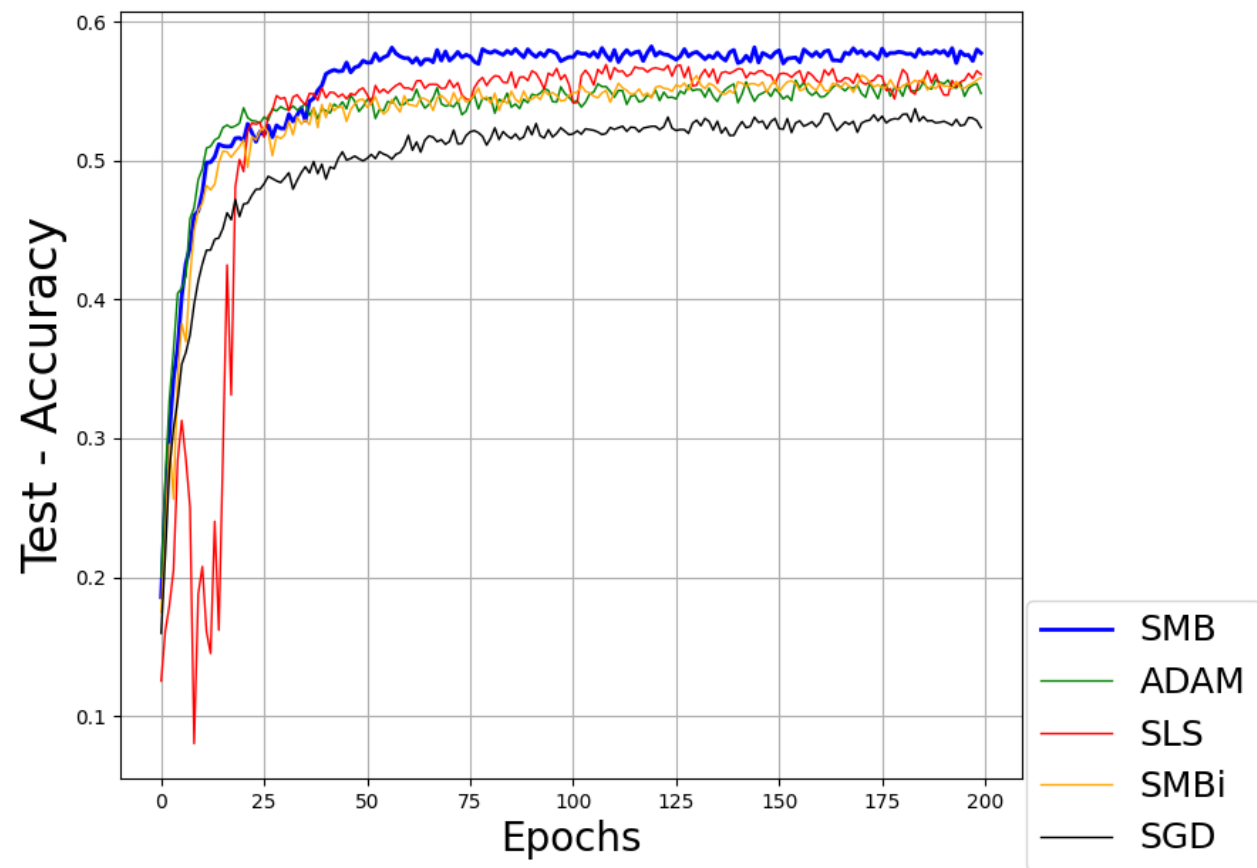
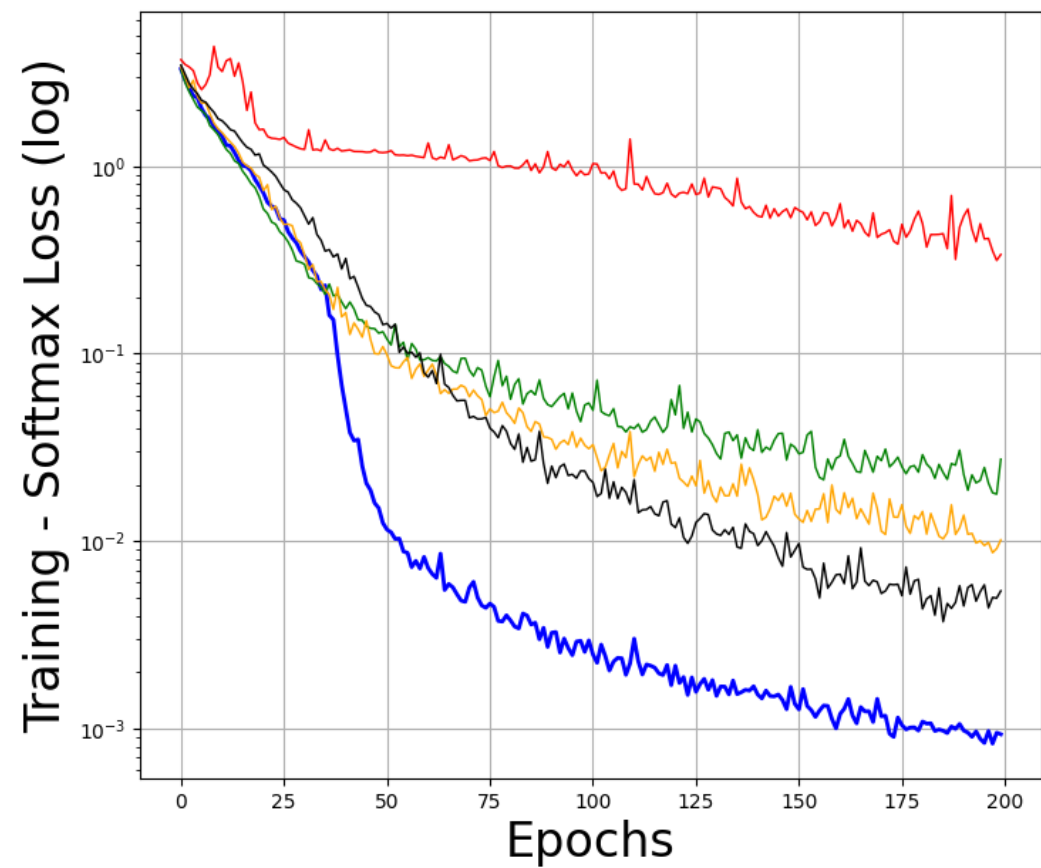




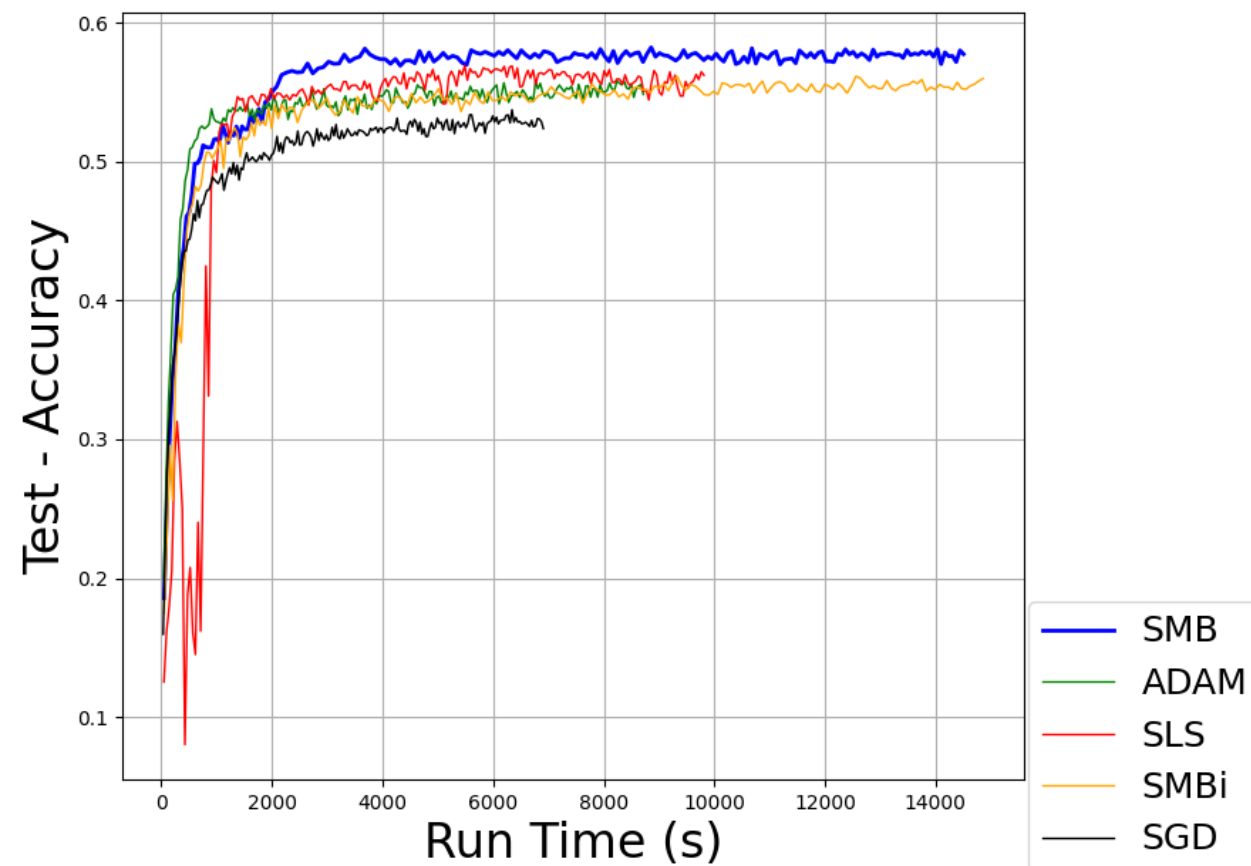
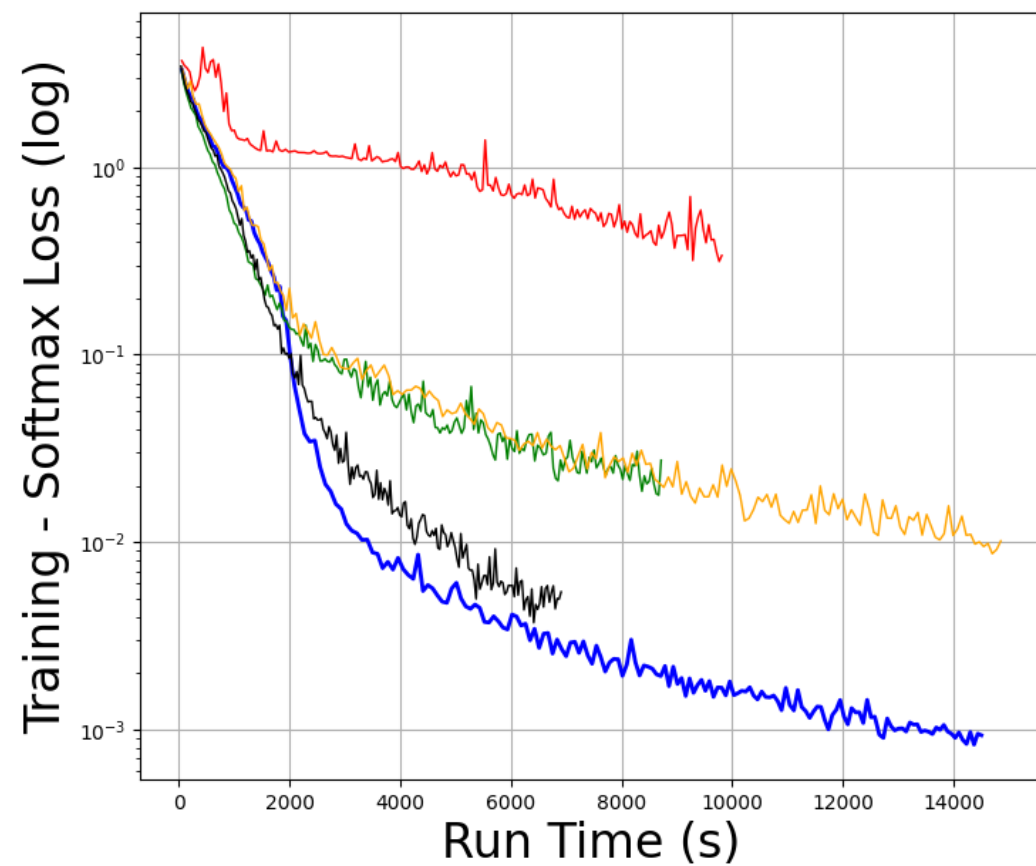
## CIFAR10-DENSENET10



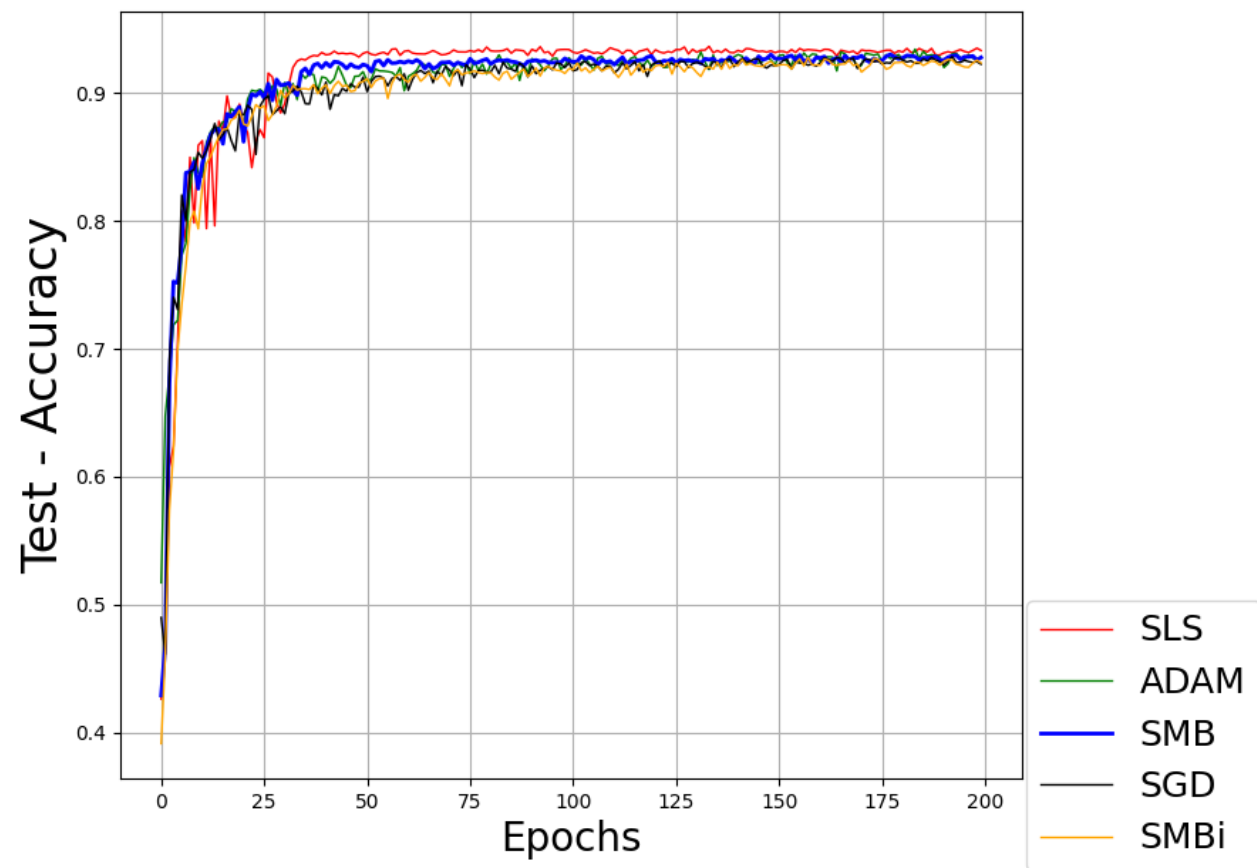
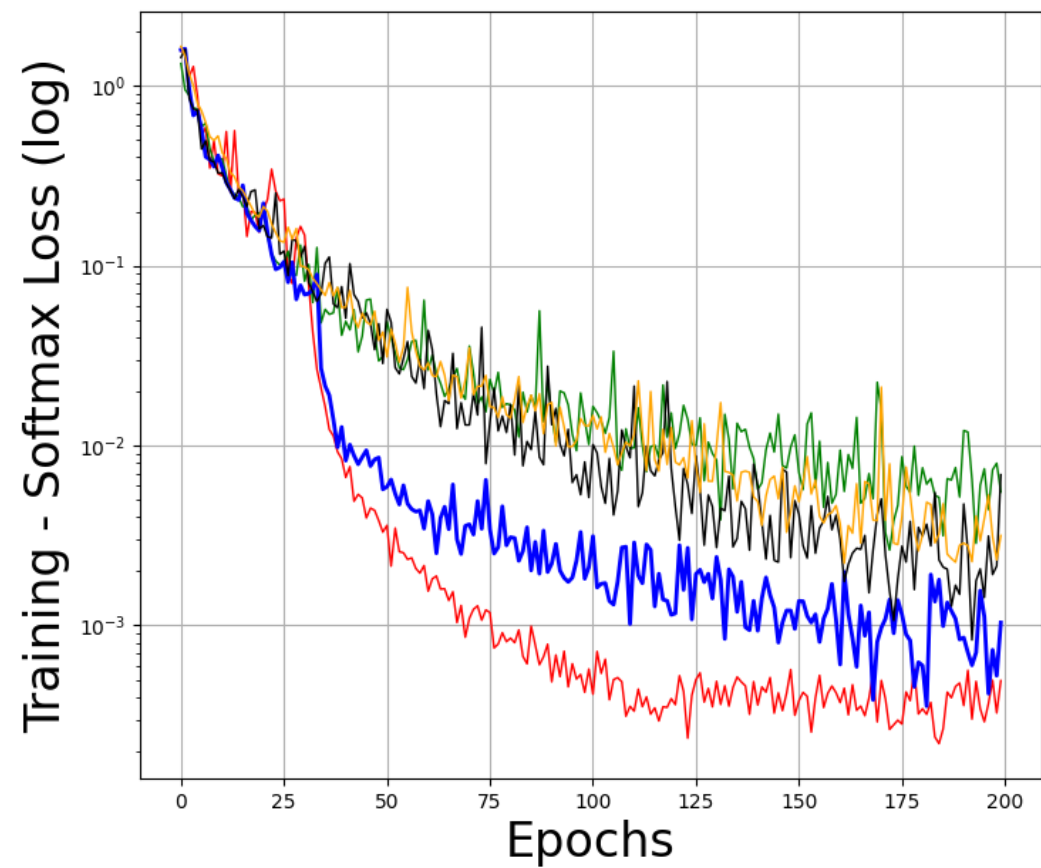
## CIFAR100-DENSENET100



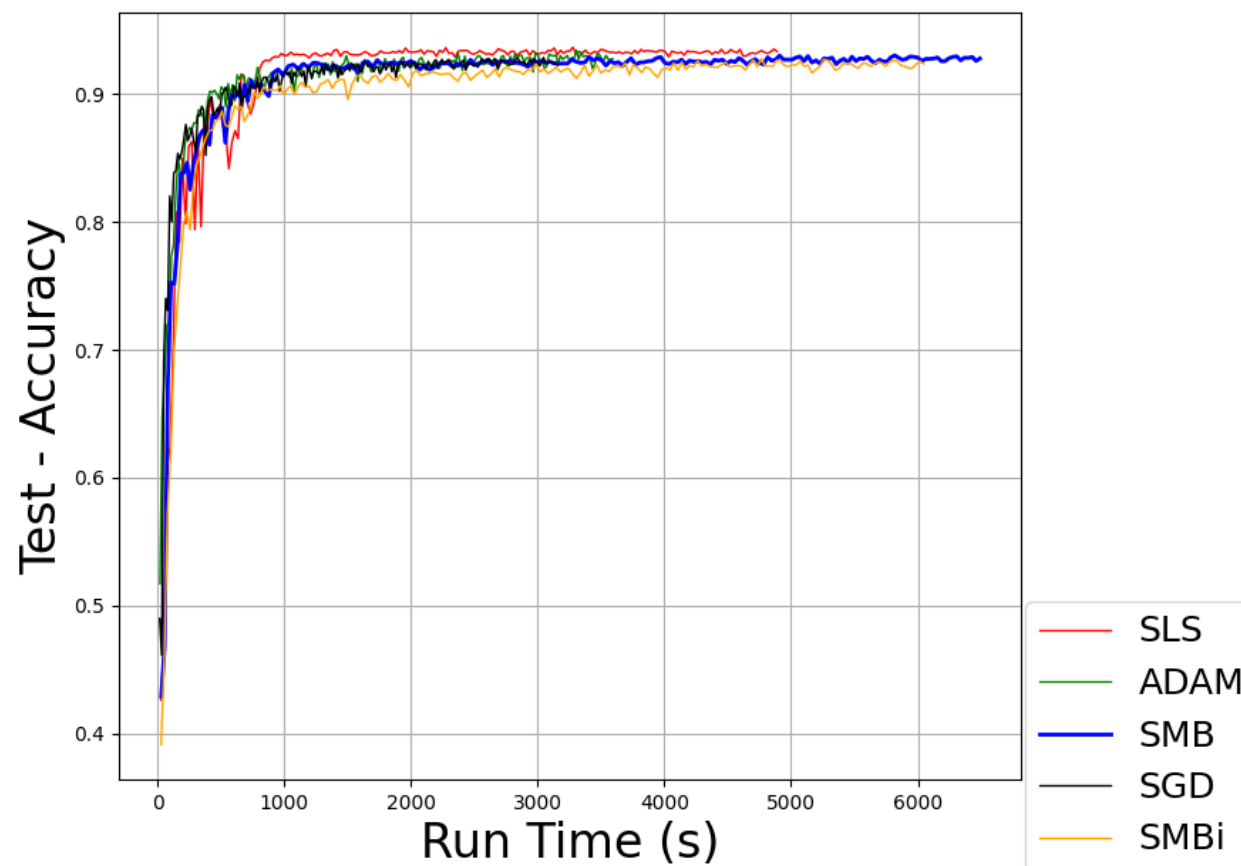
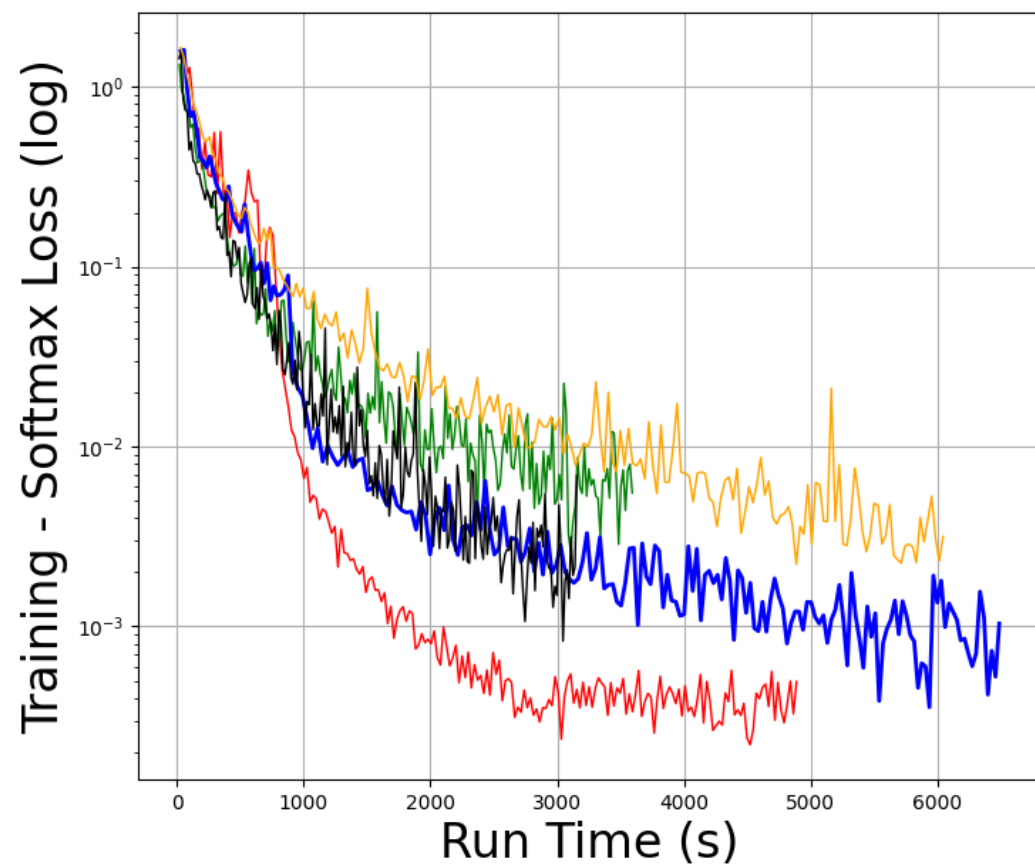
## CIFAR100-DENSENET100



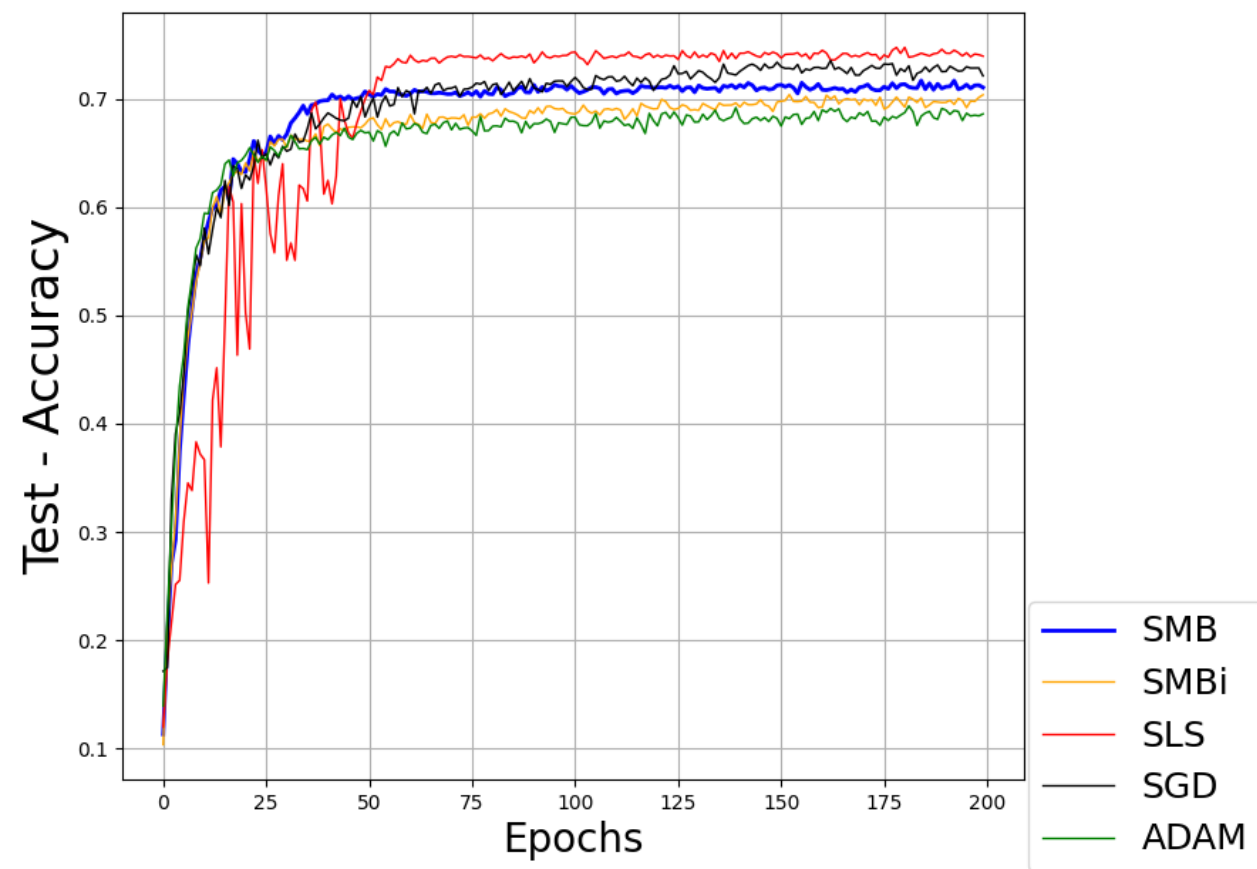
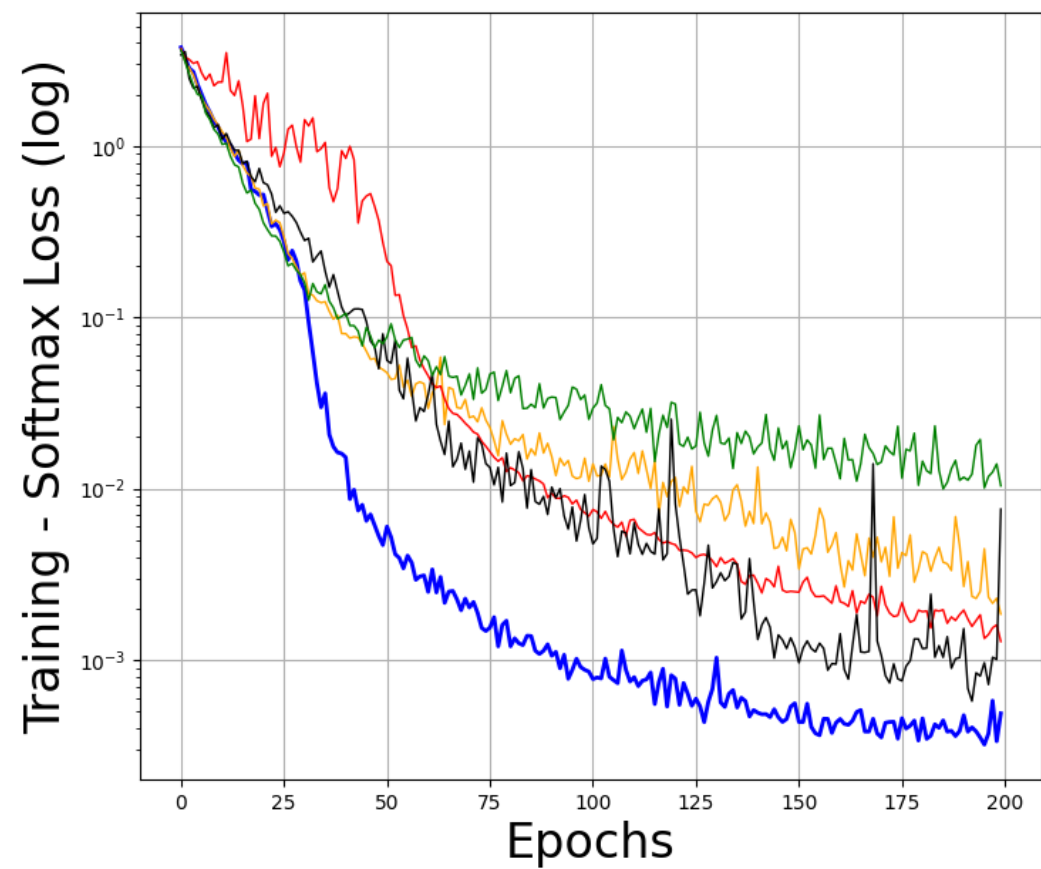
## CIFAR10-RESNET34\_10



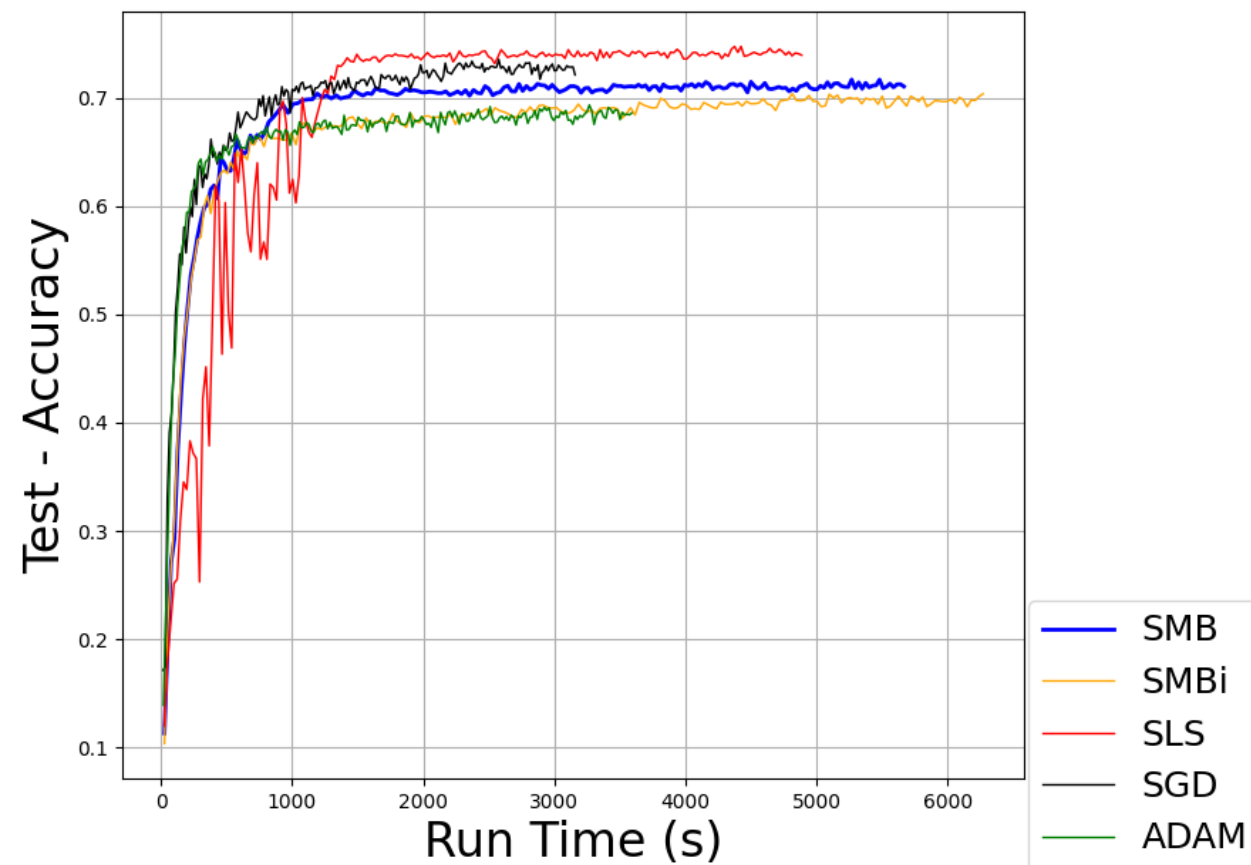
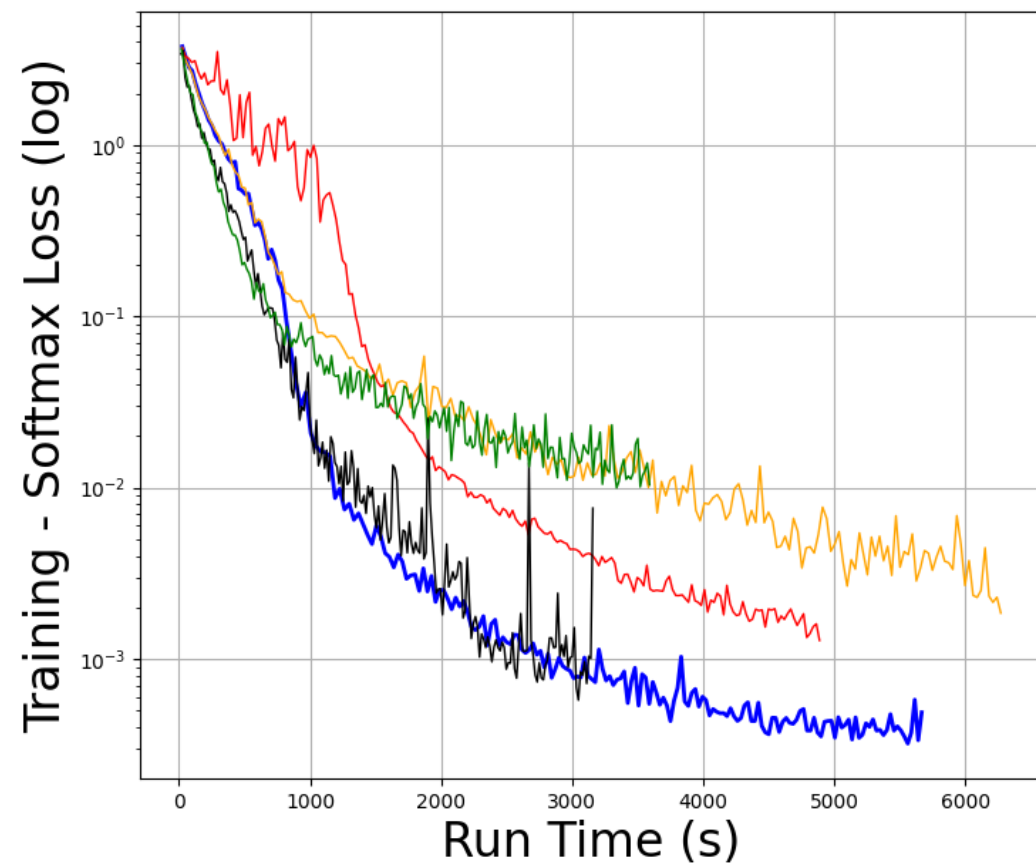
## CIFAR10-RESNET34\_10

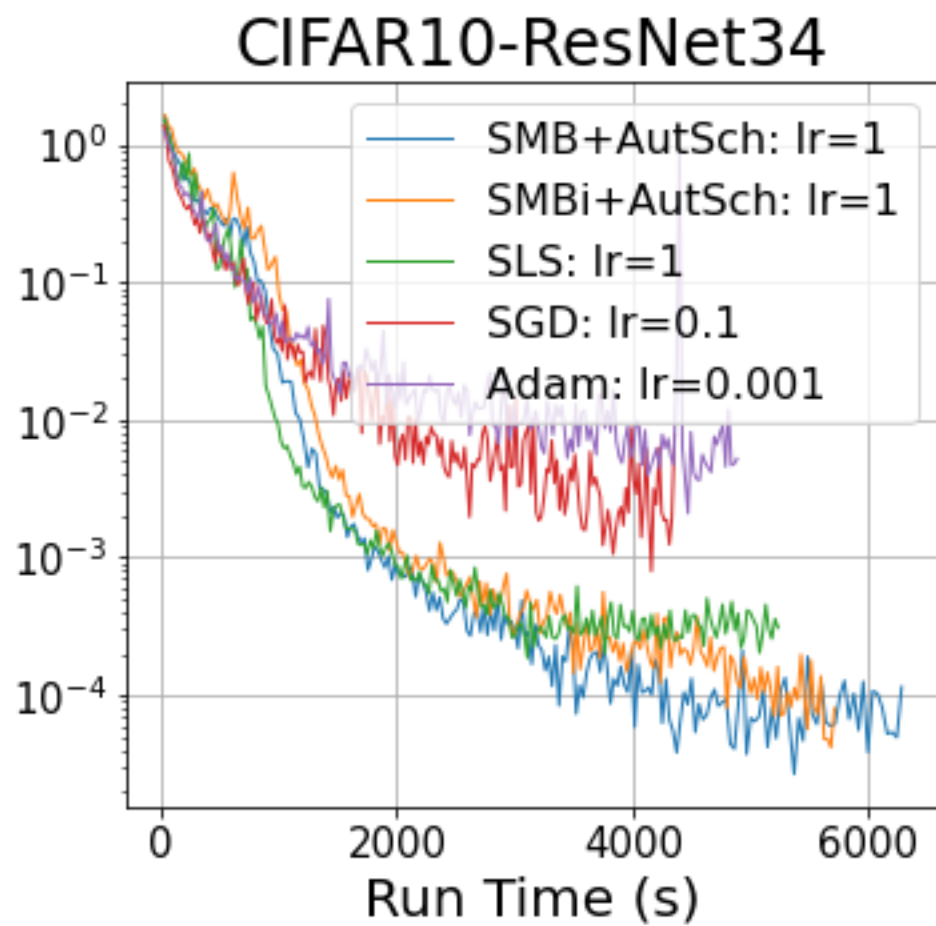
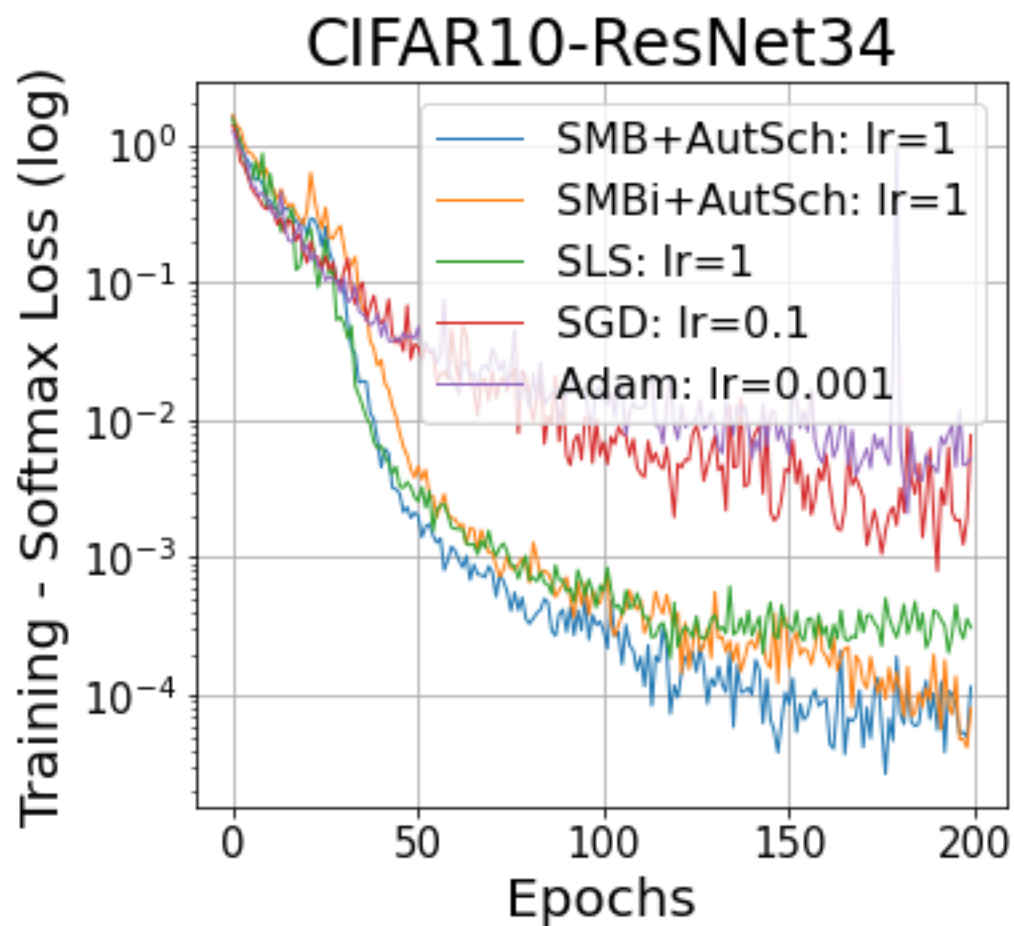


## CIFAR100-RESNET34\_100

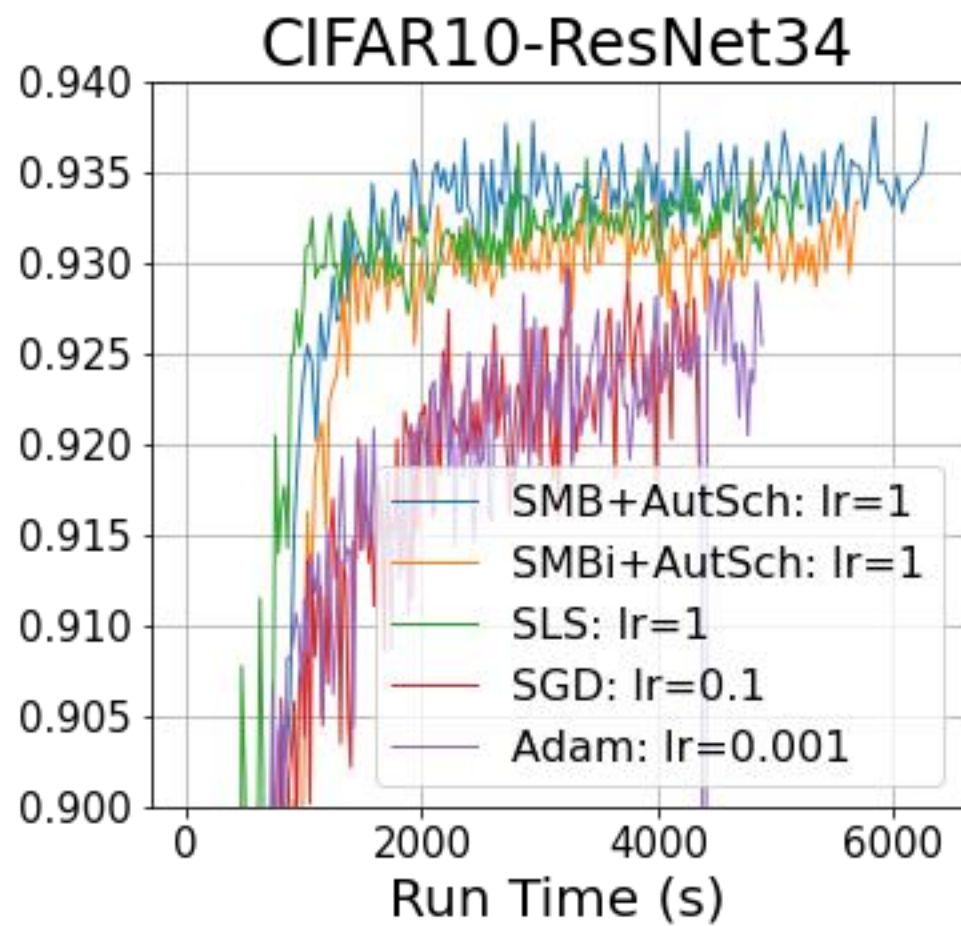
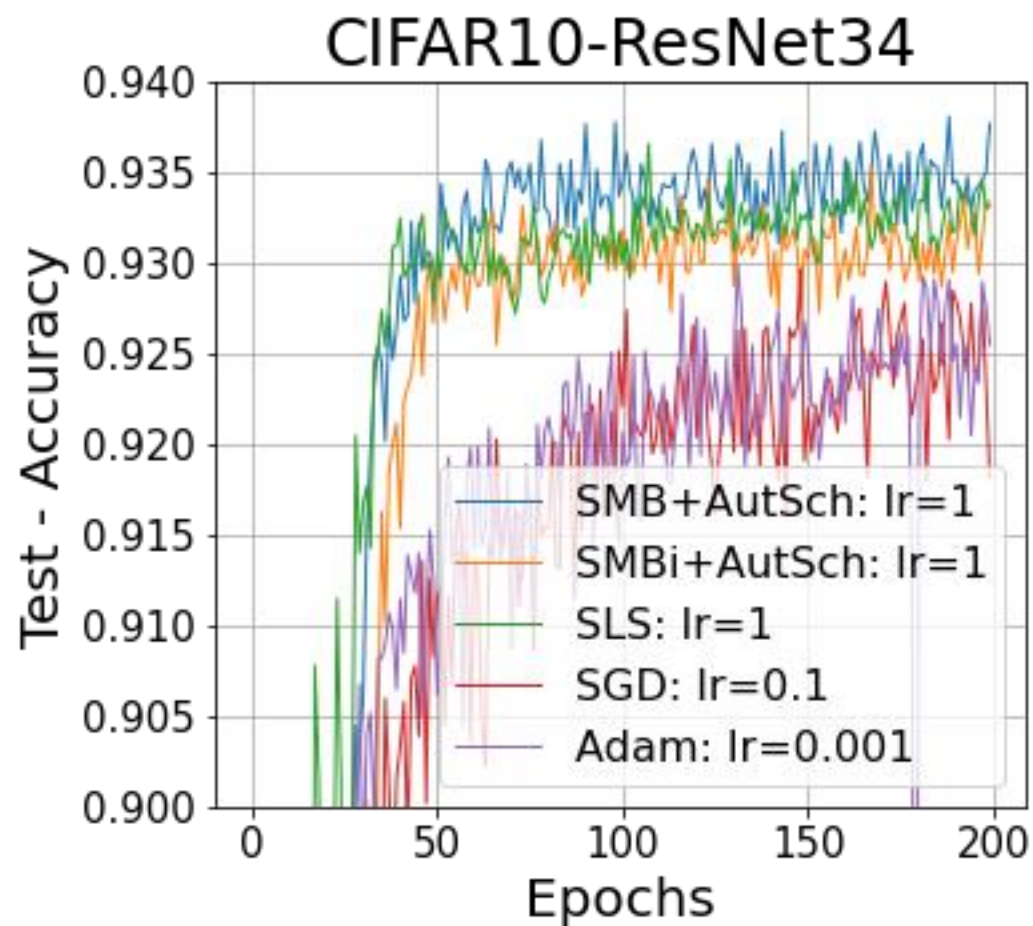


## CIFAR100-RESNET34\_100



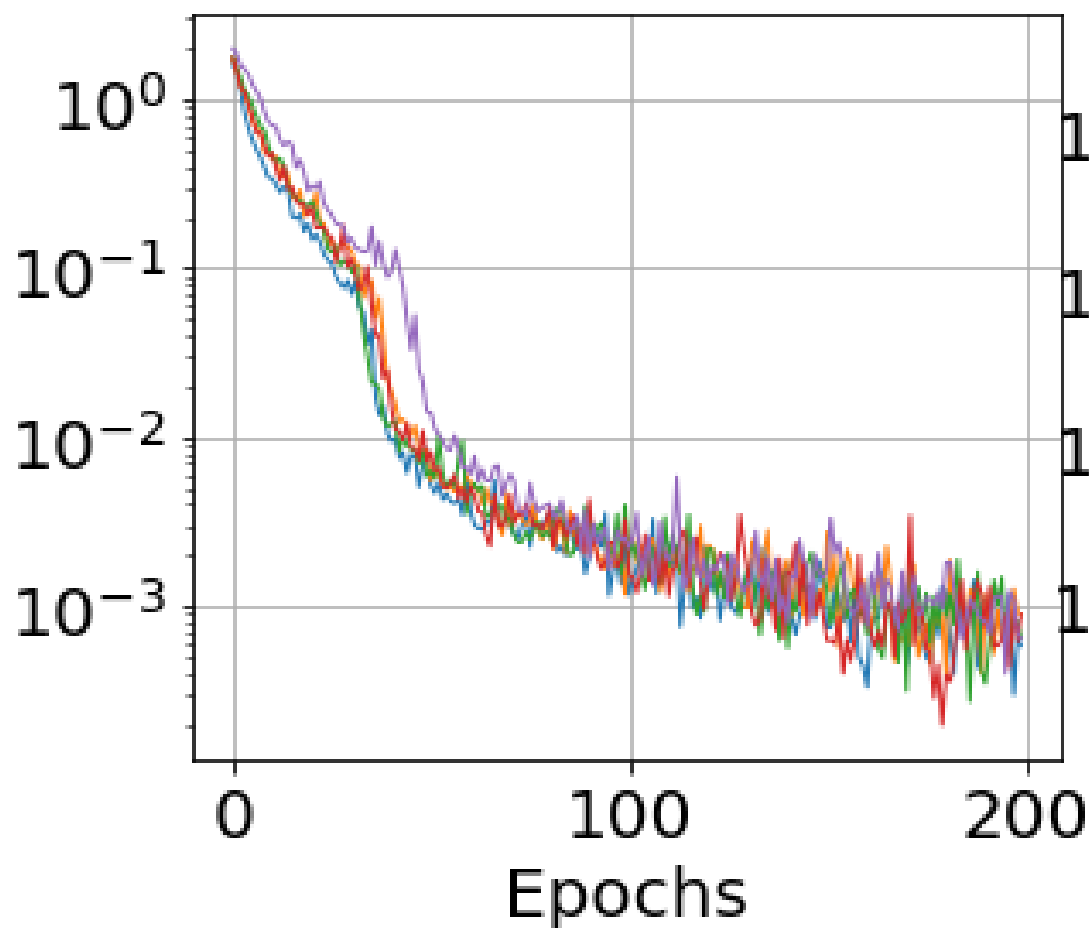




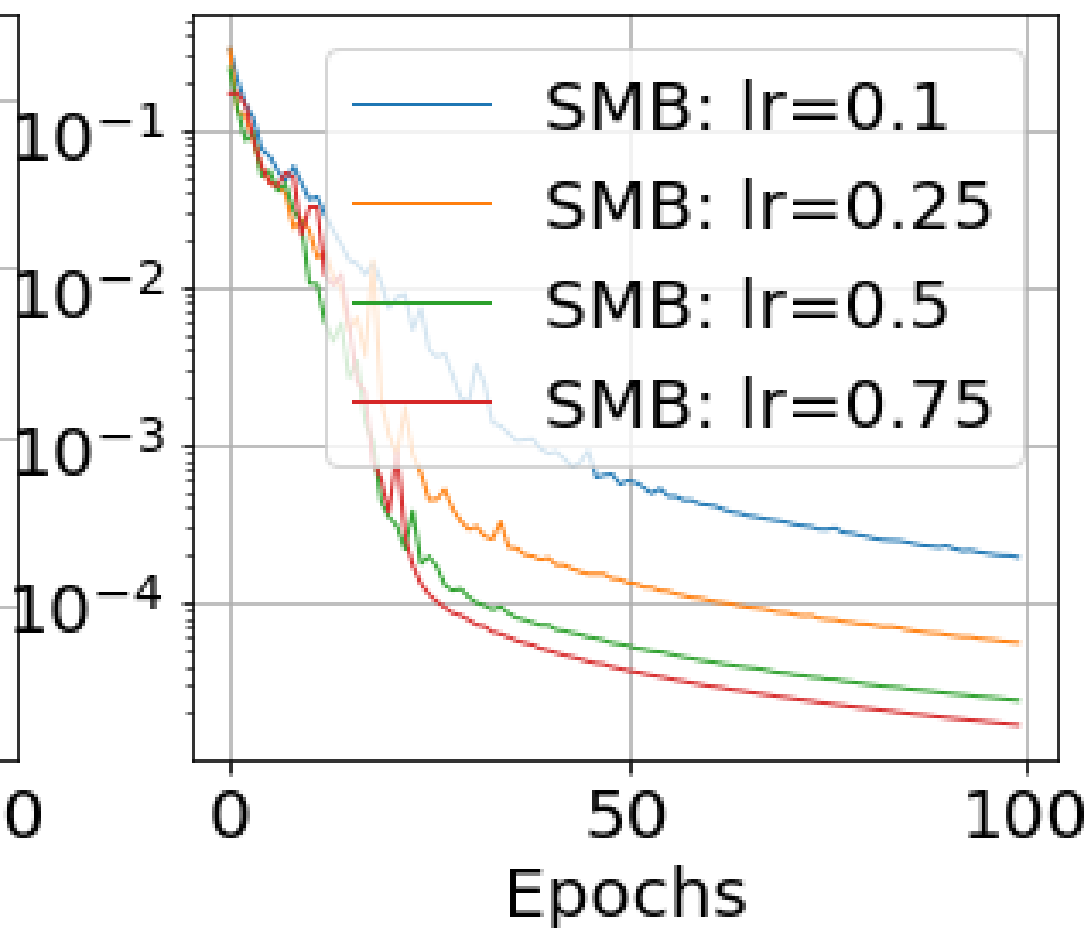


Training - Softmax Loss (log)

CIFAR10-ResNet34



MNIST-MLP



Paper: <https://arxiv.org/abs/2111.07058>

Codes: <https://github.com/sibirbil/SMB>

Thanks!