



## Introduction and Contributions

Training LLMs poses challenges due to their massive scale and heterogeneous architectures. While adaptive optimizers like AdamW help address gradient variations, they still struggle with efficient and effective parameter-wise learning rate estimation, resulting in training instability, slow convergence, and poor compatibility with PEFT techniques.

- **Scaling with Gradient Grouping (SGG)** is a flexible LLM optimizer wrapper that scales adaptive learning rates with online grouping constraints rather than replace them in pre-defined groups (like Adam-mini), balancing parameter-wise dynamics and collective optimization behavior.
- Practically, SGG integrates seamlessly with existing optimizers and PEFT techniques, requiring no changes to the training pipeline or model architectures with the CPU and CPU-GPU hybrid implementations.
- SGG's consistent improvement shows the potential of scaling adaptive learning rates with group-wise constraints with LLMs (pre-training, PEFT, RL tasks) and MLLMs. SGG offers an intuitive instantiation of this scheme, while different grouping and scaling strategies are conceivable and might inspire future studies.

## Overview of SGG Optimization

### Algorithm 1 Scaling with Gradient Grouping

**Require:** Parameters  $\{\theta_l\}_{l=1}^L$ , global learning rate schedule  $\eta$ , optimizer hyperparameters  $(\beta_1, \beta_2)$ , objective  $\mathcal{L}$ , dataset  $\mathcal{D}$ . SGG hyperparameters: cluster number  $K$ , recluster interval  $T$ , scaling EMA decay  $\beta_3$ .

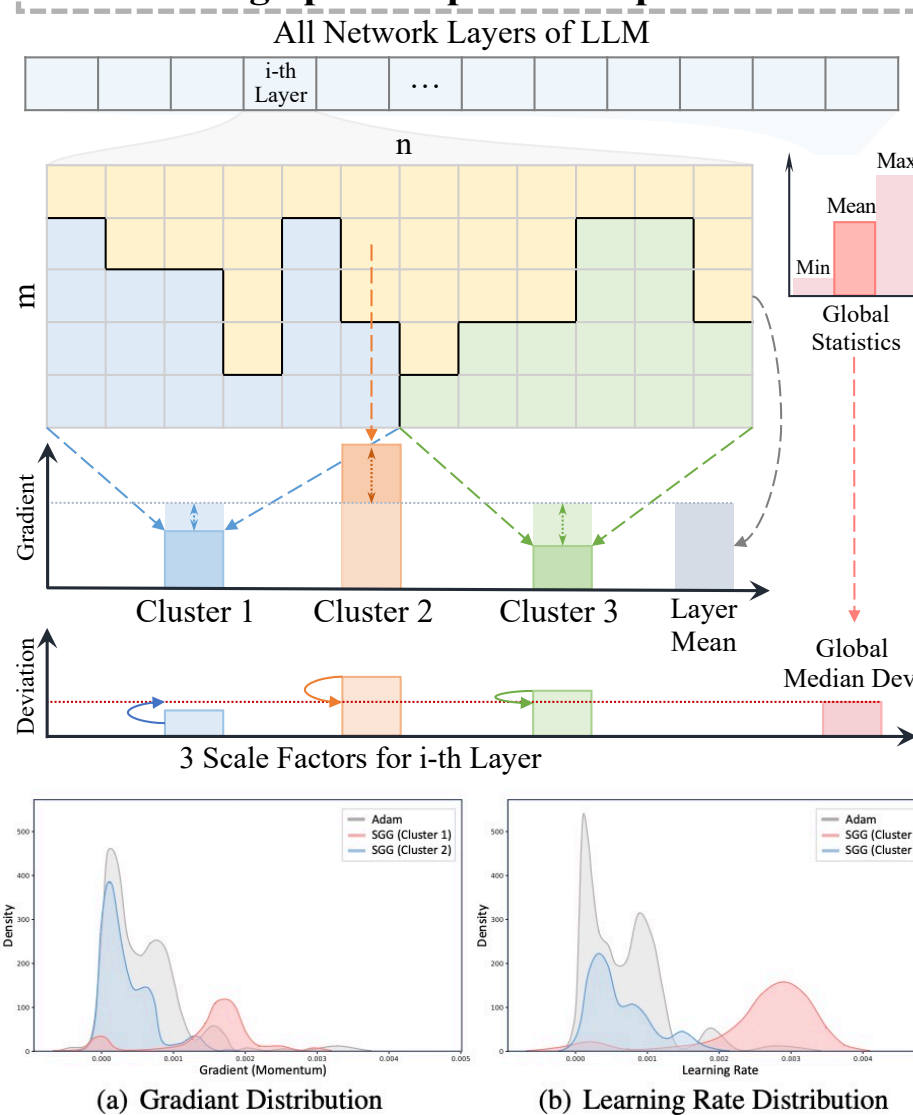
**Ensure:** Optimized model parameters  $\theta$ .

```

1: Initialize:
2: RandomInit( $\{\theta_l^0\}_{l=1}^L$ )           ▷ Model parameters
3:  $\{\alpha_l^0\}_{l=1}^L \leftarrow \eta$            ▷ Adaptive learning rates
4:  $\{C_l\}_{l=1}^L \leftarrow \mathbf{0}$            ▷ Cluster assignment
5:  $\{S_l\}_{l=1}^L \leftarrow \mathbf{1}$            ▷ Cluster scaling factor
6: for each iteration  $t = 1, 2, \dots$  do
7:    $\eta^t \leftarrow \text{LRScheduler}(\eta, t)$ 
8:   for each layer  $l = 1, 2, \dots, L$  do
9:     // — Standard Gradient-based Update Steps —
10:    Gradient Computation
11:     $g_l^t \leftarrow \nabla_{\theta_l^{t-1}} \mathcal{L}(\theta^{t-1}, \mathcal{D})$ 
12:    Momentum Estimation
13:     $m_l^t \leftarrow \text{MomentumEstimate}(g_l^t, m_l^{t-1}, \beta_1)$ 
14:    Adaptive Learning Rate Estimation
15:     $\alpha_l^t \leftarrow \text{LREstimate}(\alpha_l^{t-1}, m_l^t, \beta_2, \eta^t)$ 
16:    // — SGG Specific Steps —
17:    if  $t \bmod T == 0$  then           ▷ Re-clustering
18:      Assign Gradient Clusters
19:       $C_l^t \leftarrow \text{GradCluster}(m_l^t, K)$ 
20:      Update Cluster Scaling Factors
21:       $S_l^t \leftarrow \text{ScaleUpdate}(C_l^t, m_l^t, \beta_3)$ 
22:    end if
23:    Apply Learning Rate Scaling
24:     $\alpha_l^t \leftarrow \alpha_l^t \cdot S_l^t[C_l^t]$    ▷ Cluster-specific scaling
25:    Parameter Update
26:     $\theta_l^t \leftarrow \theta_l^{t-1} - \alpha_l^t \cdot m_l^t$ 
27:  end for
28: end for

```

### SGG with online grouping and group-specific LR scaling upon adaptive LR optimizer.



## Comparison Experiments

**LLM Comparison:** C4 Pre-training with diverse LLaMA sizes and popular LLM optimizers.

Method	Venue	60M	130M	350M	1B
<i>Pre-training with Full-Rank Optimizers</i>					
Adam <sup>†</sup>	ICLR'15	34.06	25.08	18.80	15.56
NAdam	ICLR'18	35.86	28.88	19.24	15.78
RAdam	ICLR'20	30.43	25.17	19.13	15.65
LAMB	ICLR'20	33.04	24.37	18.26	15.84
Adan	TPAMI'23	32.01	23.14	17.32	14.70
<b>Adam+SGG</b>	<b>Ours</b>	<b>30.31</b>	<b>22.18</b>	<b>17.28</b>	<b>14.30</b>
$\Delta$ Gains		<b>-3.75</b>	<b>-2.90</b>	<b>-1.52</b>	<b>-1.26</b>
<i>Pre-training with Memory-efficient Optimizers</i>					
Adam-mini <sup>†</sup>	ICLR'25	34.10	24.85	19.05	16.07
Adafactor <sup>†</sup>	ICML'18	32.57	23.98	17.74	15.19
Low-Rank <sup>†</sup>	arXiv'22	78.18	45.51	37.41	34.53
CAME	ACL'23	31.37	23.38	17.45	14.68
<b>CAME+SGG</b>	<b>Ours</b>	<b>30.15</b>	<b>22.91</b>	<b>17.09</b>	<b>14.35</b>
$\Delta$ Gains		<b>-1.22</b>	<b>-0.46</b>	<b>-0.36</b>	<b>-0.33</b>
APOLLO <sup>†</sup>	MLSys'25	31.55	22.94	16.85	14.20
<b>APOLLO+SGG</b>	<b>Ours</b>	<b>30.18</b>	<b>22.52</b>	<b>16.54</b>	<b>13.95</b>
$\Delta$ Gains		<b>-1.37</b>	<b>-0.42</b>	<b>-0.31</b>	<b>-0.25</b>
<i>Low-Rank Pre-training</i>					
LoRA <sup>†</sup>	ICLR'22	34.99	33.92	25.58	19.21
ReLoRA <sup>†</sup>	ICLR'23	37.04	29.37	29.08	18.33
GaLore <sup>†</sup>	ICML'24	34.88	25.36	18.95	15.64
<b>LoRA+SGG</b>	<b>Ours</b>	<b>30.62</b>	<b>23.62</b>	<b>17.86</b>	<b>14.73</b>
$\Delta$ Gains		<b>-4.37</b>	<b>-10.30</b>	<b>-7.72</b>	<b>-4.48</b>
Training Tokens		1.1B	2.2B	6.4B	13.1B

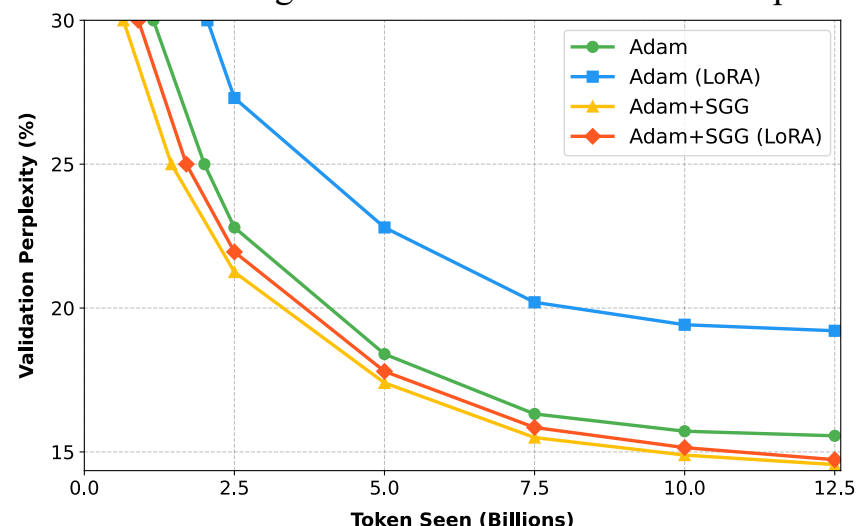
Considering a neural net layer  $W \in \mathbb{R}^{m \times n}$  ( $m \leq n$ ) with LoRA's rank  $r \leq n$  and SGG's clusters  $K \ll m$ .

Category	Method	Adaptive LR	Basic State	Extra State	Low-Rank Plugin	Extra Branch	C4↓	GPU Memory
Classical Opt.	SGD	×	Weight & Grad.	×	×	×	—	$2mn$
Adaptive LR Opt.	Adam	Param-wise $mn$	Weight & Grad.	2 <sup>nd</sup> -Moment $mn$	×	×	23.36	$3mn$
Efficient Opt.	CAME	Param-wise $mn$	Weight & Grad.	NMF $2(m+n)$	NMF	×	<b>-1.64</b>	$2mn+2(m+n)$
					LoRA	✓	<b>+5.06</b>	<b><math>+3r(m+n)</math></b>
PEFT	LoRA	×	Full-rank Grad.	×	LoRA	✓	<b>-1.99</b>	$+0$
Opt. Wrapper	<b>SGG</b>	Group-wise $K$	Base Opt.	Indices $(mn+K)$	Clustering	✓		

**MLLM Comparison:** Top-1 accuracy (%) with LLaVA variants is reported, where MMB and MMB<sup>CN</sup> denote MMBench and MMBench (Chinese).

Optimizer	Image Question Answering				Benchmarks			Avg
	GQA	VizWiz	SciVQA	VQA <sup>T</sup>	MMB	MMB <sup>CN</sup>	POPE	
LLaVA-v1.5 Full-Rank SFT								
AdamW	62.0	50.0	66.8	<b>58.2</b>	64.3	58.3	85.9	63.6
Adafactor	62.7	48.2	70.7	57.1	66.1	60.4	86.0	64.5
LAMB	43.8	53.3	61.5	43.4	43.2	41.8	81.2	52.6
AdamW+SGG	62.4	50.2	69.8	57.4	65.9	60.1	<b>86.3</b>	64.6
$\Delta$ Gains	<b>+0.4</b>	<b>+0.2</b>	<b>+3.0</b>	<b>-0.8</b>	<b>+1.6</b>	<b>+1.8</b>	<b>+0.4</b>	<b>+1.0</b>
Adafactor+SGG	<b>62.8</b>	50.6	<b>71.6</b>	57.3	<b>66.3</b>	<b>60.8</b>	86.0	<b>65.1</b>
$\Delta$ Gains	<b>+0.1</b>	<b>+2.4</b>	<b>+0.9</b>	<b>+0.2</b>	<b>+0.2</b>	<b>+0.4</b>	<b>+0.0</b>	<b>+0.6</b>
LAMB+SGG	44.0	<b>53.3</b>	61.8	43.5	43.3	41.9	81.3	52.7
$\Delta$ Gains	<b>+0.2</b>	<b>+0.0</b>	<b>+0.3</b>	<b>+0.1</b>	<b>+0.1</b>	<b>+0.1</b>	<b>+0.1</b>	<b>+0.1</b>
LLaVA-v1.5 Low-Rank SFT (AdamW)								
LoRA	63.0	47.8	68.4	58.2	66.1	58.9	86.4	64.1
LoRA+SGG	<b>63.4</b>	<b>51.0</b>	<b>70.1</b>	<b>58.6</b>	<b>66.7</b>	<b>59.4</b>	<b>86.6</b>	<b>65.1</b>
$\Delta$ Gains	<b>+0.4</b>	<b>+2.2</b>	<b>+1.5</b>	<b>+0.4</b>	<b>+0.6</b>	<b>+0.5</b>	<b>+0.2</b>	<b>+1.0</b>
LLaVA-v1.5 8-bit Low-Rank SFT (AdamW)								
Q-LoRA	54.3	50.7	66.4	52.5	56.0	49.8	82.9	58.9
Q-LoRA+SGG	<b>55.1</b>	<b>51.3</b>	<b>66.7</b>	<b>53.0</b>	<b>56.1</b>	<b>51.0</b>	<b>83.4</b>	<b>59.5</b>
$\Delta$ Gains	<b>+0.8</b>	<b>+0.6</b>	<b>+0.3</b>	<b>+0.5</b>	<b>+0.1</b>	<b>+0.2</b>	<b>+0.5</b>	<b>+0.6</b>

C4 Pre-training with Full-rank and Low-rank Opt.



Alpaca FT with LR and batch size scaling-up

