

LLM Pretraining

short presentation

Andrei Semenov

MLO Group Meeting, 07.05.2025

Optimizers..., hyperparameters...

MUON IS SCALABLE FOR LLM TRAINING

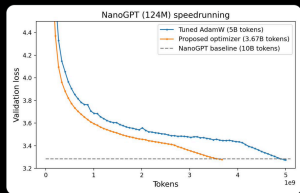
TECHNICAL REPORT

Post

Keller Jordan
@kellerjordan0

New training speed record for @karpathy's 124M-parameter NanoGPT setup: 3.28 Fineweb validation loss in 3.7B training tokens

Previous record: 5B tokens
Changelog: new optimizer
1/8



10:29 PM · Oct 4, 2024 · 252.6K Views

Dion: A Communication-Efficient Optimizer for Large Models

Kwangjun Ahn Byron Xu
Microsoft Research

BETTER, FASTER, OLDER

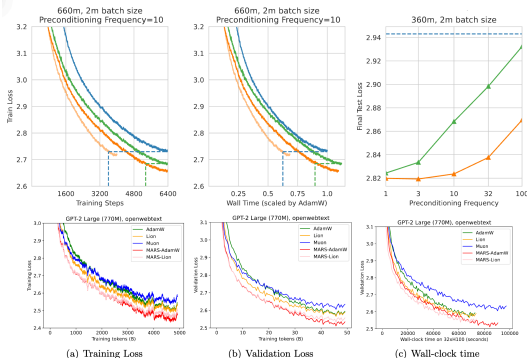


Figure 1: The training and validation loss curves, plotted against both training tokens and wall-clock time on GPT-2 large model (770M).

So we benchmarked

Algorithm 1 AdamW

```
1: Input: Initial parameter
2: Initialize:  $\mathbf{m}_0 \leftarrow \mathbf{0}, \mathbf{v}_0$ 
3: for  $t \in [T]$  do
4:    $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$ 
5:    $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$ 
6:    $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{m}_t / (1 - \beta_2^t)$ 
7:    $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t)$ 
8:    $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t \left( \frac{\hat{\mathbf{m}}_t}{\sqrt{\mathbf{v}_t}} \right)$ 
9: end for
10: Return:  $\mathbf{x}_T$ 
```

Algorithm 2 ADOPT

```
1: Input: Initial parameters  $\mathbf{x}_0, \nu$ 
2: Initialize:  $\mathbf{m}_0 \leftarrow \mathbf{0}, \mathbf{v}_0 \leftarrow \nabla \mathcal{L}(\mathbf{x}_0)$ 
3: for  $t \in [T]$  do
4:    $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$ 
5:    $c_t \leftarrow t^{1/4}$ 
6:    $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) c_t \mathbf{g}_t$ 
7:    $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{m}_t$ 
8:    $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t (\mathbf{m}_t + \lambda \mathbf{x}_t)$ 
9: end for
10: Return:  $\mathbf{x}_T$ 
```

Algorithm 3 AdEMAMix

```
1: Input: Initial parameters  $\mathbf{x}_0$ , number of iterations  $T$ , learning rate  $\gamma_t$ , weight decay  $\lambda$ ,  $\beta_1, \beta_2, \beta_3, \beta_{\text{start}}, \alpha$ ,
   beta_scheduler, alpha_scheduler, warmup parameters  $T_{\beta_3}$  and  $T_\alpha, \varepsilon$ .
2: Initialize:  $\mathbf{m}_0 \leftarrow \mathbf{0}, \mathbf{m}_t^{\text{slow}} \leftarrow \mathbf{0}, \mathbf{v}_0 \leftarrow \mathbf{0}$ 
3: for  $t \in [T]$  do
4:    $\beta_3(t) \leftarrow \text{beta\_scheduler}(t, \beta_3, \beta_{\text{start}}, T_{\beta_3}), \alpha(t) \leftarrow \text{alpha\_scheduler}(t, \alpha, T_\alpha)$   $\triangleright$  Update  $\beta_3$  and  $\alpha$  schedulers
5:    $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$ 
6:    $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$ 
7:    $\mathbf{m}_t^{\text{slow}} \leftarrow \beta_3(t) \mathbf{m}_{t-1}^{\text{slow}} + (1 - \beta_3(t)) \mathbf{g}_t$   $\triangleright$  Update slow EMA
8:    $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t$ 
9:    $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t), \hat{\mathbf{v}}_t \leftarrow \mathbf{v}_t / (1 - \beta_2^t)$ 
10:   $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t \left( \frac{\hat{\mathbf{m}}_t + \alpha(t) \mathbf{m}_t^{\text{slow}}}{\sqrt{\hat{\mathbf{v}}_t + \varepsilon}} + \lambda \mathbf{x}_t \right)$ 
11: end for
12: Return:  $\mathbf{x}_T$ 
```

Algorithm 4 Lion

Algorithm 6 Signum (our PyTorch variant)

```
1: Input: Initial parameters  $\mathbf{x}_0$ , number of iterations  $T$ ,
   learning rate  $\gamma_t$ , weight decay  $\lambda$ , momentum  $\beta$ .
2: Initialize:  $\mathbf{m}_0 \leftarrow \mathbf{0}$ 
3: for  $t \in [T]$  do
4:    $\mathbf{g}_t \leftarrow \mathcal{L}(\mathbf{x}_t, \xi_t)$ 
5:    $\mathbf{m}_t \leftarrow \beta \mathbf{m}_{t-1} + \mathbf{g}_t$ 
6:    $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t (\text{sign}(\beta \mathbf{m}_t + \mathbf{g}_t) + \lambda \mathbf{x}_t)$ 
7: end for
8: Return:  $\mathbf{x}_T$ 
```

Algorithm 10 SOAP

Algorithm 11 Sophia

Algorithm 12 SF-AdamW

Algorithm 8 Muon

Algorithm 13 Prodigy

Algorithm 14 MARS (MARS-AdamW)

Optimizers: Adam-like

Algorithm 3 ADMMAMIX

1: **Input:** Initial parameters \mathbf{x}_0 , number of iterations T , learning rate γ_t , weight decay λ , β_1 , β_2 , β_3 , β_{start} , α , β_{α} , α_{start} , α , T_{β_3} and T_{α} , ε .
2: **Initialize:** $\mathbf{m}_0 \leftarrow \mathbf{0}$, $\mathbf{m}_0^{\text{slow}} \leftarrow \mathbf{0}$, $\mathbf{v}_0 \leftarrow \mathbf{0}$
3: **for** $t \in [T]$ **do**
4: $\beta_3(t) \leftarrow \text{beta_scheduler}(t, \beta_3, \beta_{\text{start}}, T_{\beta_3})$, $\alpha(t) \leftarrow \text{alpha_scheduler}(t, \alpha, T_{\alpha})$ \triangleright Update β_3 and α schedulers
5: $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$
6: $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$
7: $\mathbf{m}_t^{\text{slow}} \leftarrow \beta_3(t) \mathbf{m}_{t-1}^{\text{slow}} + (1 - \beta_3(t)) \mathbf{g}_t$ \triangleright Update slow EMA
8: $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t$
9: $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t)$, $\hat{\mathbf{v}}_t \leftarrow \mathbf{v}_t / (1 - \beta_2^t)$
10: $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t \left(\frac{\hat{\mathbf{m}}_t + \alpha(t) \mathbf{m}_t^{\text{slow}}}{\sqrt{\hat{\mathbf{v}}_t + \varepsilon}} + \lambda \mathbf{x}_t \right)$
11: **end for**
12: **Return:** \mathbf{x}_T

Algorithm 2 ADOPT

1: **Input:** Initial parameters \mathbf{x}_0 , number of iterations T , learning rate γ_t , weight decay λ , β_1 , β_2 , ε .
2: **Initialize:** $\mathbf{m}_0 \leftarrow \mathbf{0}$, $\mathbf{v}_0 \leftarrow \nabla \mathcal{L}(\mathbf{x}_0, \xi_0) \odot \nabla \mathcal{L}(\mathbf{x}_0, \xi_0)$
3: **for** $t \in [T]$ **do**
4: $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$
5: $c_t \leftarrow t^{1/4}$ \triangleright Update clipping value schedule
6: $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \text{clamp} \left(\frac{\mathbf{g}_t}{\max(\sqrt{\mathbf{v}_{t-1}}, \varepsilon)}, c_t \right)$
7: $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t$
8: $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t (\mathbf{m}_t + \lambda \mathbf{x}_t)$ \triangleright Update without \mathbf{v}_t
9: **end for**
10: **Return:** \mathbf{x}_T

Optimizers: "Sign-based"

Algorithm 4 Lion

- 1: **Input:** Initial parameters \mathbf{x}_0 , number of iterations T , learning rate γ_t , weight decay λ , β_1 , β_2 .
 - 2: **Initialize:** $\mathbf{m}_0 \leftarrow \mathbf{0}$
 - 3: **for** $t \in [T]$ **do**
 - 4: $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$
 - 5: $\mathbf{m}_t \leftarrow \beta_2 \mathbf{m}_{t-1} + (1 - \beta_2) \mathbf{g}_t$
 - 6: $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t (\text{sign}(\beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t) + \lambda \mathbf{x}_t)$ ▷ Update EMA of \mathbf{g}_t
 - 7: **end for**
 - 8: **Return:** \mathbf{x}_T
-

Algorithm 5 Signum (basic)

- 1: **Input:** Initial parameters \mathbf{x}_0 , number of iterations T , learning rate γ_t , weight decay λ , momentum β .
 - 2: **Initialize:** $\mathbf{m}_0 \leftarrow \mathbf{0}$
 - 3: **for** $t \in [T]$ **do**
 - 4: $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$
 - 5: $\mathbf{m}_t \leftarrow \beta \mathbf{m}_{t-1} + (1 - \beta) \mathbf{g}_t$
 - 6: $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t (\text{sign}(\mathbf{m}_t) + \lambda \mathbf{x}_t)$
 - 7: **end for**
 - 8: **Return:** \mathbf{x}_T
-

Algorithm 6 Signum (our PyTorch variant)

- 1: **Input:** Initial parameters \mathbf{x}_0 , number of iterations T , learning rate γ_t , weight decay λ , momentum β .
 - 2: **Initialize:** $\mathbf{m}_0 \leftarrow \mathbf{0}$
 - 3: **for** $t \in [T]$ **do**
 - 4: $\mathbf{g}_t \leftarrow \mathcal{L}(\mathbf{x}_t, \xi_t)$
 - 5: $\mathbf{m}_t \leftarrow \beta \mathbf{m}_{t-1} + \mathbf{g}_t$
 - 6: $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t (\text{sign}(\beta \mathbf{m}_t + \mathbf{g}_t) + \lambda \mathbf{x}_t)$
 - 7: **end for**
 - 8: **Return:** \mathbf{x}_T
-

Optimizers: Using matrix information

Algorithm 7 MuonNonLD (for non-1D parameters)

```
1: Input: Initial non-1D parameters  $x_0$ , number of iterations  $T$ , learning rate  $\gamma_t$ , momentum  $\beta$ , number of Newton-Schulz iterations  $T_{NS}$ ,  $a, b, c$  coefficients.
2: Initialize:  $m_0 \leftarrow \mathbf{0}$ 
3: for  $t \in [T]$  do
4:    $g_t \leftarrow \nabla \mathcal{L}(x_t, \xi_t)$ 
5:    $m_t \leftarrow \beta m_{t-1} + g_t$ 
6:    $g_t \leftarrow \beta m_t + g_t$  ▷ Practical implementation of Nesterov momentum
7:   Set:  $w_0 \leftarrow g_t / \|g_t\|_F$ 
8:   for  $n \in [T_{NS}]$  do
9:      $w_{n+1} \leftarrow a w_n + b w_n w_n^\top + c (w_n w_n^\top)^2 w_n$  ▷ Newton-Schulz iteration
10:   end for
11:    $x_{t+1} \leftarrow x_t - \gamma_t w_{T_{NS}}$ 
12: end for
13: Return:  $x_T$ 
```

Algorithm 8 Muon (general scheme)

```
1: Input: Initial parameters  $x_0$ , number of iterations  $T$ , Muon's parameters: learning rate  $\eta_t^M$ , momentum  $\beta$ , number of Newton-Schulz iterations  $T_{NS}$ ,  $a, b, c$  coefficients. AdamW's parameters: learning rate  $\eta_t^A$ , weight decay  $\lambda, \beta_1, \beta_2, \epsilon$ .
2: for  $t \in [T]$  do
3:   if  $x_t \in \{\text{embeds, scalar params, lm head}\}$  then
4:      $x_t^A \leftarrow x_t$ 
5:      $x_{t+1}^A \leftarrow \text{AdamW}(x_t^A, \eta_t^A, \lambda, \beta_1, \beta_2, \epsilon, T = 1)$  ▷ One iteration of AdamW
6:   else
7:      $x_t^M \leftarrow x_t$ 
8:      $x_{t+1}^M \leftarrow \text{MuonNonLD}(x_t^M, \eta_t^M, T_{NS}, \beta, a, b, c, T = 1)$  ▷ One iteration of MuonNonLD
9:   end if
10: end for
11: Return:  $x_T^A, x_T^M$ 
```

Algorithm 9 SOAFNonLD (for non-1D parameters)

```
1: Input: Initial parameters  $x_0$ , number of iterations  $T$ , learning rate  $\gamma_t$ , weight decay  $\lambda, \beta_1, \beta_2$ , preconditioning frequency  $\phi, \epsilon$ .
2: Initialize:  $m_0 \leftarrow \mathbf{0}, v_0 \leftarrow \mathbf{0}$ 
3: Initialize preconditioners:  $q_t, q_e \leftarrow \text{eigenbasis}(\nabla \mathcal{L}(x_0, \xi_0) \nabla \mathcal{L}(x_0, \xi_0)^\top)$ 
4: for  $t \in [T]$  do
5:    $g_t \leftarrow \nabla \mathcal{L}(x_t, \xi_t)$ 
6:    $g_t' \leftarrow q_t^\top g_t q_t$  ▷ Rotate  $g_t$ 
7:    $m_t \leftarrow \beta_t m_{t-1} + (1 - \beta_t) g_t'$ 
8:    $m_t' \leftarrow q_t^\top m_t q_t$  ▷ Compute Adam's statistics in rotational space
9:    $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t' \odot g_t'$ 
10:   $\gamma_t \leftarrow \gamma_t \frac{\sqrt{1 - \beta_2^t}}{1 - \beta_2}$  ▷ Optional: use bias correction
11:   $x_{t+1} \leftarrow x_t - \gamma_t \left( q_t \frac{m_t}{\sqrt{\beta_2^t + v_t}} q_t^\top + \lambda x_t \right)$  ▷ Perform update in original space
12:   $l_t \leftarrow \beta_2 l_{t-1} + (1 - \beta_2) g_t g_t^\top$ 
13:   $r_t \leftarrow \beta_2 r_{t-1} + (1 - \beta_2) g_t g_t$  ▷ Update preconditioners
14:  if  $t \equiv 1 \pmod{\phi}$  then
15:     $q_t \leftarrow \text{QR}(l_t q_t)$ 
16:     $q_e \leftarrow \text{QR}(r_t q_e)$ 
17:  end if
18: end for
19: Return:  $x_T$ 
```

Optimizers: Parameter-free effort

Algorithm 13 Prodigy

1: **Input:** Initial parameters \mathbf{x}_0 , number of iterations T , learning rate γ , weight decay λ , $\beta_1, \beta_2, \varepsilon$.
2: **Initialize:** $d_0 \leftarrow 10^{-6}$, $\gamma \leftarrow 1$, $\mathbf{m}_0 \leftarrow \mathbf{0}$, $\mathbf{v}_0 \leftarrow \mathbf{0}$, $r_0 \leftarrow 0$, $\mathbf{s}_0 \leftarrow \mathbf{0}$ ▷ Optional: use scheduler on γ
3: **for** $t \in [T]$ **do**
4: $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{x}_t, \xi_t)$
5: $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) d_t \mathbf{g}_t$
6: $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) d_t^2 \mathbf{g}_t \odot \mathbf{g}_t$
7: $\gamma_t \leftarrow \gamma \sqrt{1 - \beta_2^t} / (1 - \beta_1^t)$ ▷ Optional: use bias correction
8: $r_t \leftarrow \sqrt{\beta_2} r_{t-1} + (1 - \sqrt{\beta_2}) \gamma_t d_t^2 (\mathbf{g}_t, \mathbf{x}_0 - \mathbf{x}_t)$
9: $\mathbf{s}_t \leftarrow \sqrt{\beta_2} \mathbf{s}_{t-1} + (1 - \sqrt{\beta_2}) \gamma_t d_t^2 \mathbf{g}_t$
10: $d_{t+1} \leftarrow \max\{d_t, \frac{r_t}{\|\mathbf{s}_t\|_1}\}$
11: $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \gamma_t d_t (\mathbf{m}_t / (\sqrt{\mathbf{v}_t} + d_t \varepsilon) + \lambda \mathbf{x}_t)$
12: **end for**
13: **Return:** \mathbf{x}_T

Algorithm 12 SF-AdamW

1: **Input:** Initial parameters \mathbf{x}_0 , number of iterations T , learning rate γ , weight decay λ , β_1, β_2 , warmup iterations $T_{\text{warmup}}, \varepsilon$.
2: **Initialize:** $\mathbf{z}_0 \leftarrow \mathbf{x}_0$, $\mathbf{v}_0 \leftarrow \mathbf{0}$
3: **for** $t \in [T]$ **do**
4: $\mathbf{y}_t \leftarrow (1 - \beta_1) \mathbf{z}_t + \beta_1 \mathbf{x}_t$
5: $\mathbf{g}_t \leftarrow \nabla \mathcal{L}(\mathbf{y}_t, \xi_t)$
6: $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t$
7: $\gamma_t \leftarrow \gamma \sqrt{1 - \beta_2^t} \min\{1, t / T_{\text{warmup}}\}$
8: $\mathbf{z}_{t+1} \leftarrow \mathbf{z}_t - \gamma_t (\mathbf{g}_t / (\sqrt{\mathbf{v}_t} + \varepsilon) + \lambda \mathbf{y}_t)$
9: $c_{t+1} \leftarrow \frac{\gamma_t^2}{\sum_{i=0}^t \gamma_i^2}$
10: $\mathbf{x}_{t+1} \leftarrow (1 - c_{t+1}) \mathbf{x}_t + c_{t+1} \mathbf{z}_{t+1}$
11: **end for**
12: **Return:** \mathbf{x}_T

Setup

Number of steps? batch size? scheduler? hyperparameters?

Number of steps? batch size? scheduler? hyperparameters?



Training Compute-Optimal Large Language Models

databricks

02 / Mosaic Research / How Long Should You Train Your Language Model?

How Long Should You Train Your Language Model?

Optimal Linear Decay Learning Rate Schedules and Further Refinements

Scaling Laws and Compute-Optimal Training Beyond Fixed Training Durations

The Surprising Agreement Between Convex Optimization Theory and Learning-Rate Scheduling for Large Model Training

**Beyond Chinchilla-Optimal:
Accounting for Inference in Language Model Scaling Laws**

Predictable Scale: Part I — Optimal Hyperparameter Scaling Law in Large Language Model Pretraining

November 22, 2024

How Does Critical Batch Size Scale in Pre-training? (Decoupling Data and Model Size)

Iterations & batch size

Table 2. Lengths of training for **Small batch settings** (32×512).

| # Parameters | Tokens (Iterations) | | | | | | Chinchilla Tokens |
|---------------------|----------------------------|-------------|-------------|-------------|-------------|---------------|--------------------------|
| 124M | 1B (64k) | 2.1B (128k) | 4.2B (256k) | 6.3B (384k) | 8.4B (512k) | 16.8B (1024k) | 2.5B |
| 210M | 1B (64k) | 2.1B (128k) | 4.2B (256k) | 6.3B (384k) | 8.4B (512k) | 16.8B (1024k) | 4.2B |

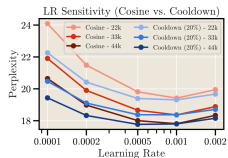
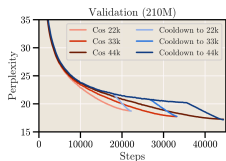
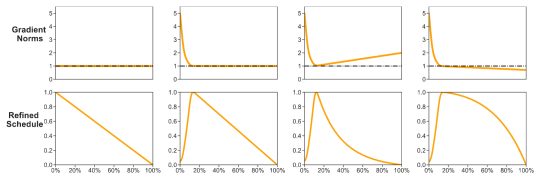
Table 3. Lengths of training for **Large batch settings** (256×512).

| # Parameters | Tokens (Iterations) | | | | | | Chinchilla Tokens |
|---------------------|----------------------------|------------|------------|------------|------------|--------------|--------------------------|
| 124M | 1B (8k) | 2.1B (16k) | 4.2B (32k) | 6.3B (48k) | 8.4B (64k) | 16.8B (128k) | 2.5B |
| 210M | 1B (8k) | 2.1B (16k) | 4.2B (32k) | 6.3B (48k) | 8.4B (64k) | 16.8B (128k) | 4.2B |

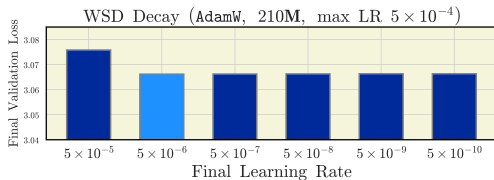
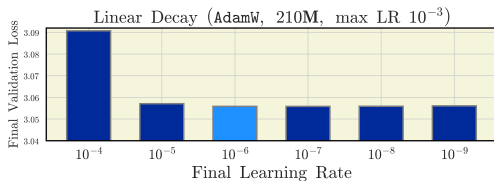
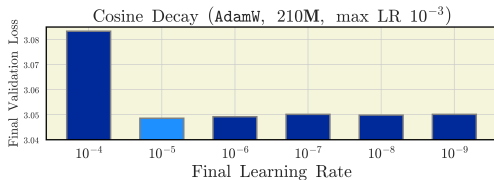
Table 4. Lengths of training for **X-Large batch settings** (1984×512).

| # Parameters | Tokens (Iterations) | | | Chinchilla Tokens |
|---------------------|----------------------------|-----------|-----------|--------------------------|
| 720M | 8B (8k) | 16B (16k) | 48B (48k) | 14.4B |

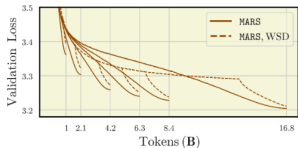
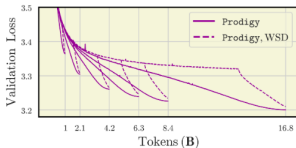
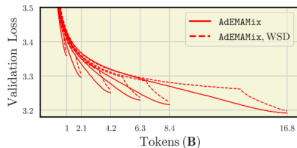
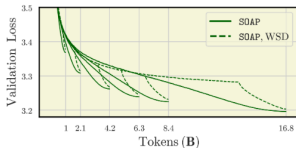
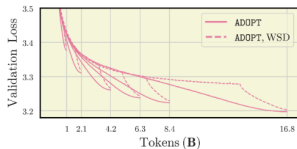
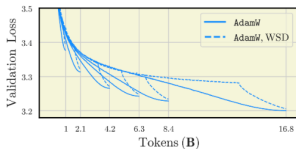
Scheduler



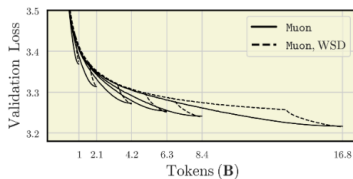
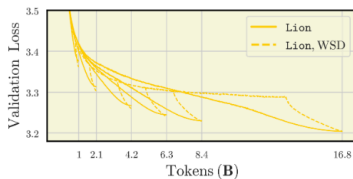
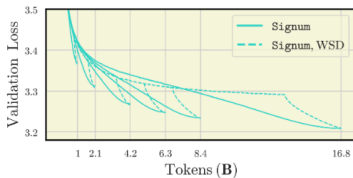
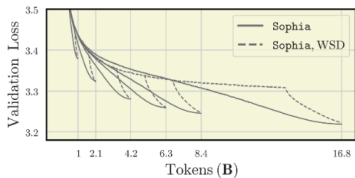
Learning rate decay: down to zero



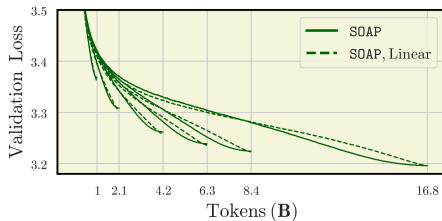
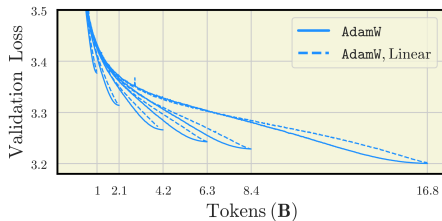
Benchmarking schedulers for different optimizers



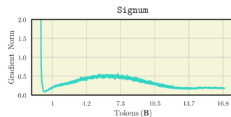
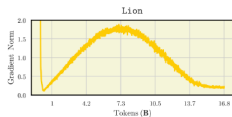
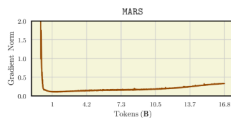
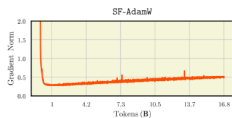
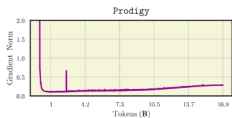
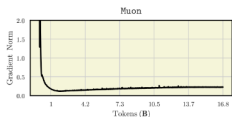
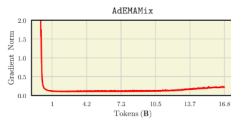
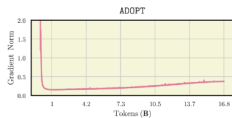
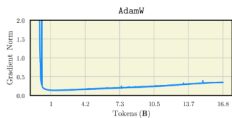
Benchmarking schedulers for different optimizers



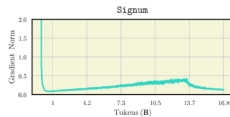
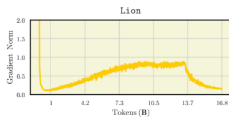
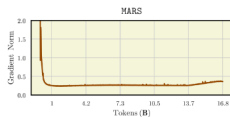
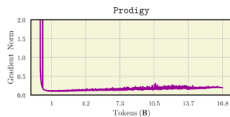
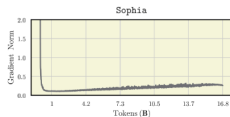
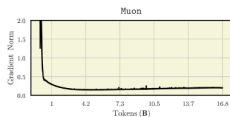
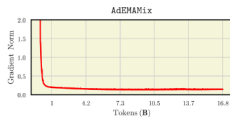
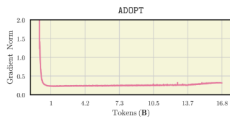
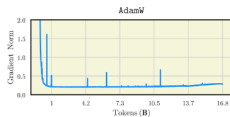
Benchmarking schedulers for different optimizers



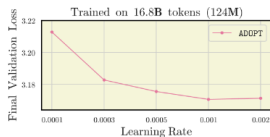
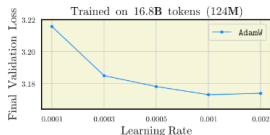
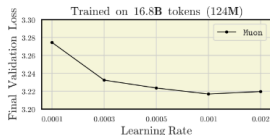
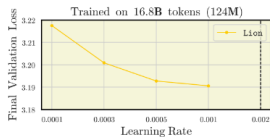
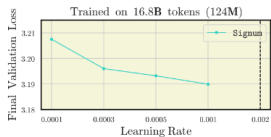
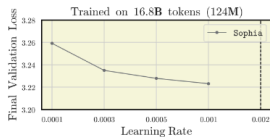
Scheduler: gradient norms perspective



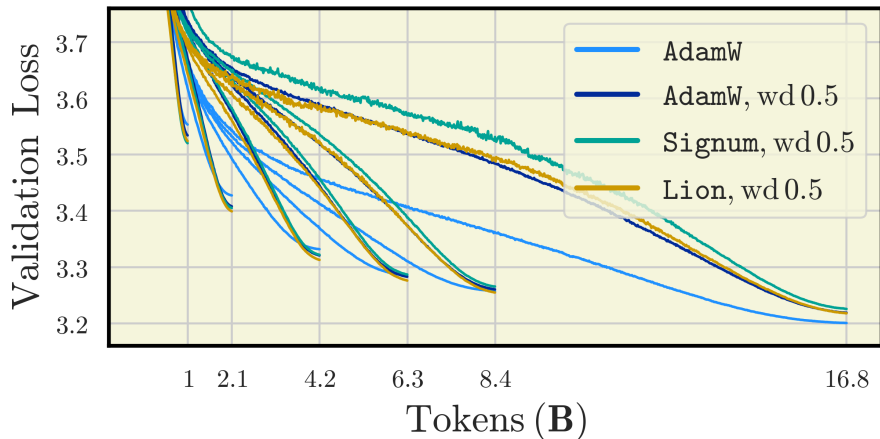
Scheduler: gradient norms perspective



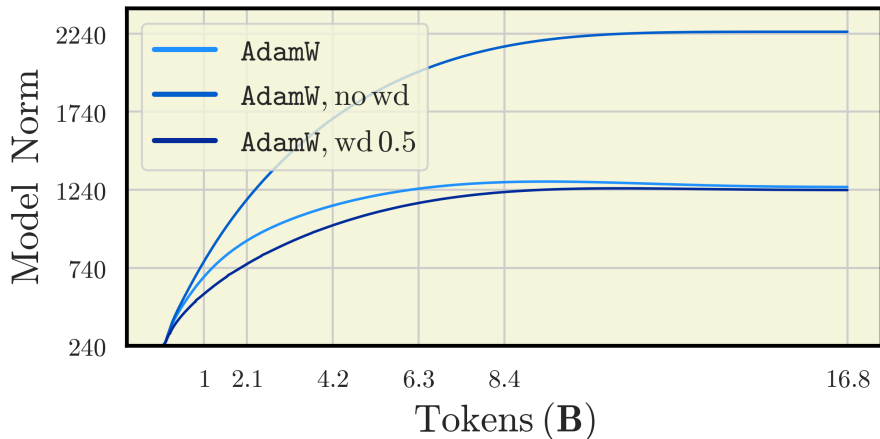
Hyperparameters: learning rate



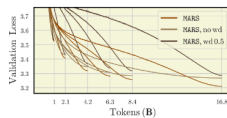
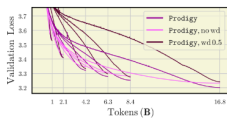
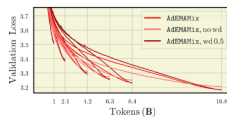
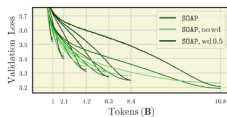
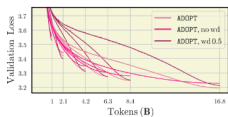
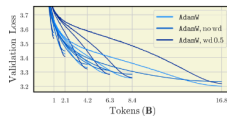
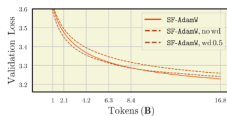
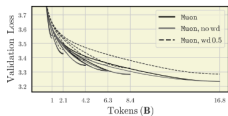
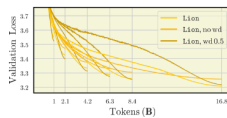
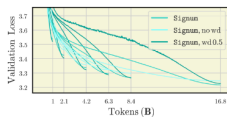
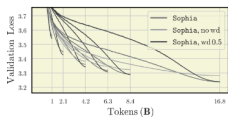
Hyperparameters: weight decay



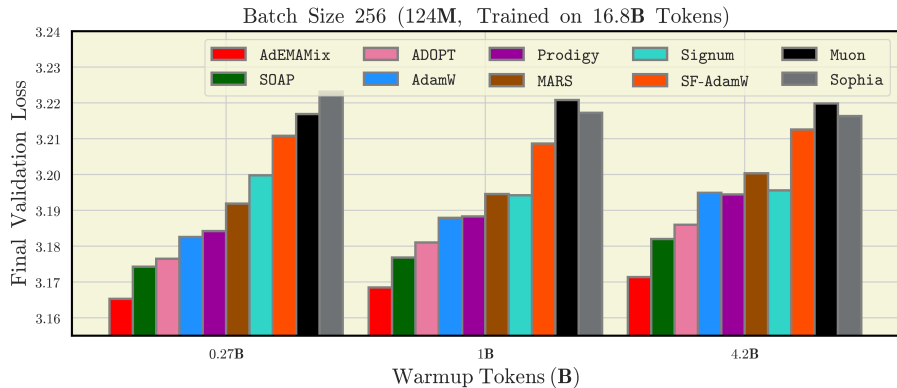
Hyperparameters: weight decay



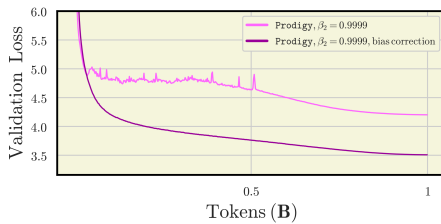
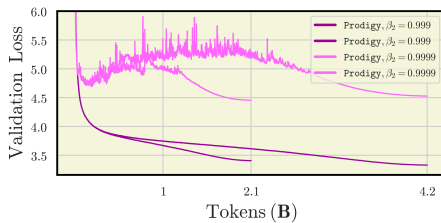
Hyperparameters: weight decay



Hyperparameters: warmup

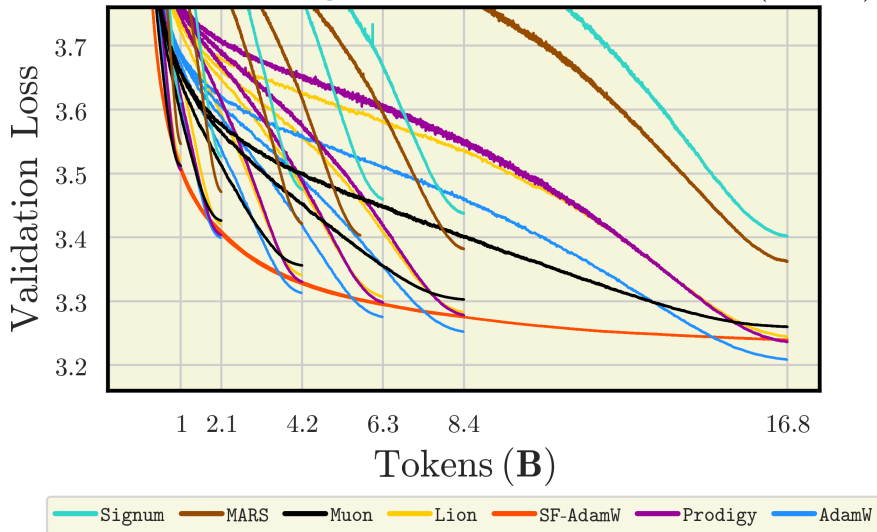


Hyperparameters: betas



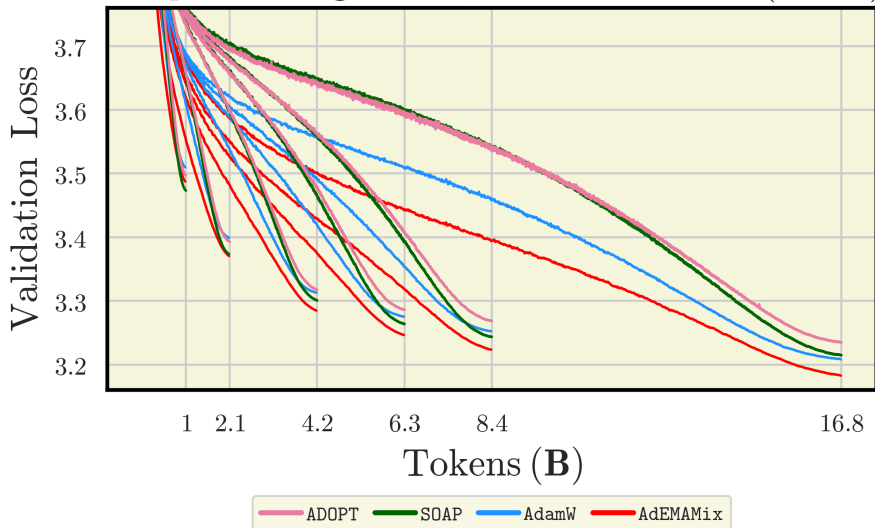
Benchmarking at small scale

Underperforming AdamW, Batch Size 32 (124M)



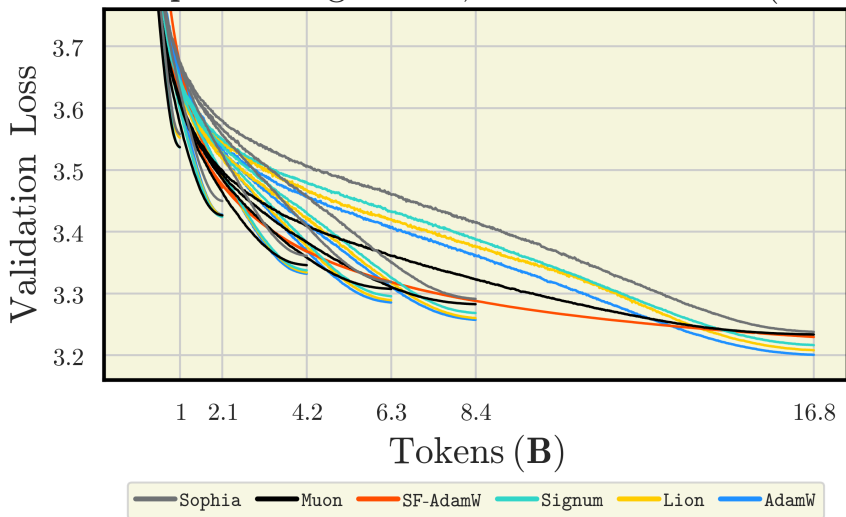
Benchmarking at small scale

Outperforming AdamW, Batch Size 32 (124M)



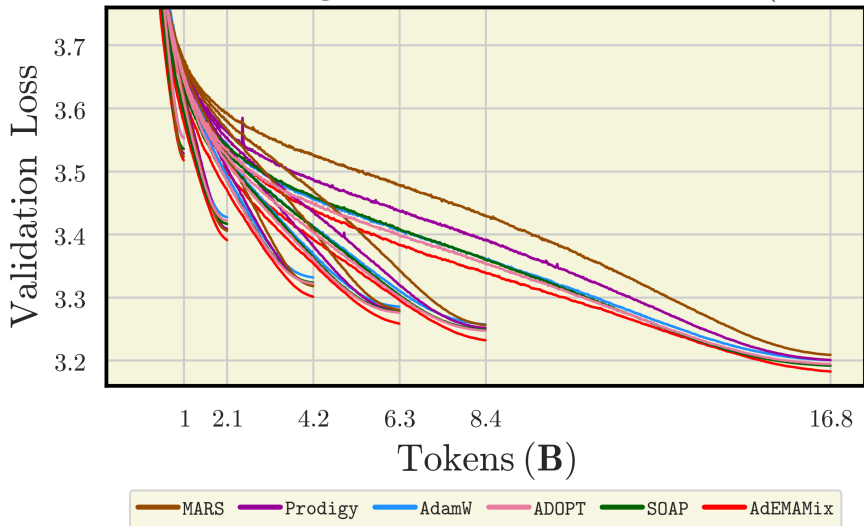
Benchmarking at small scale

Underperforming AdamW, Batch Size 256 (124M)



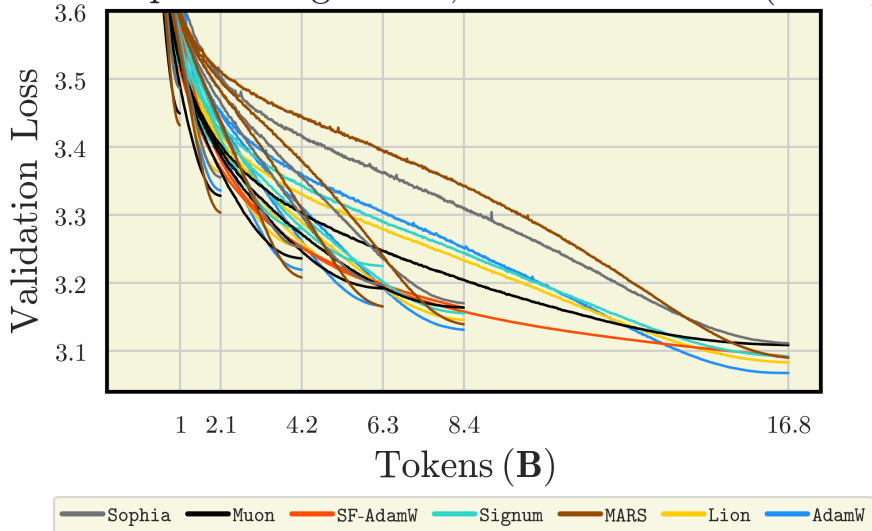
Benchmarking at small scale

Outperforming AdamW, Batch Size 256 (124M)



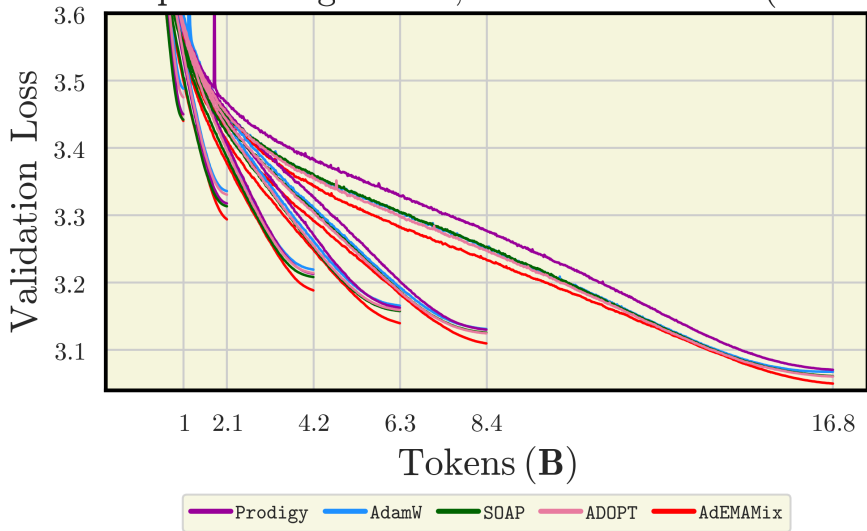
Benchmarking at small scale

Underperforming AdamW, Batch Size 256 (210M)

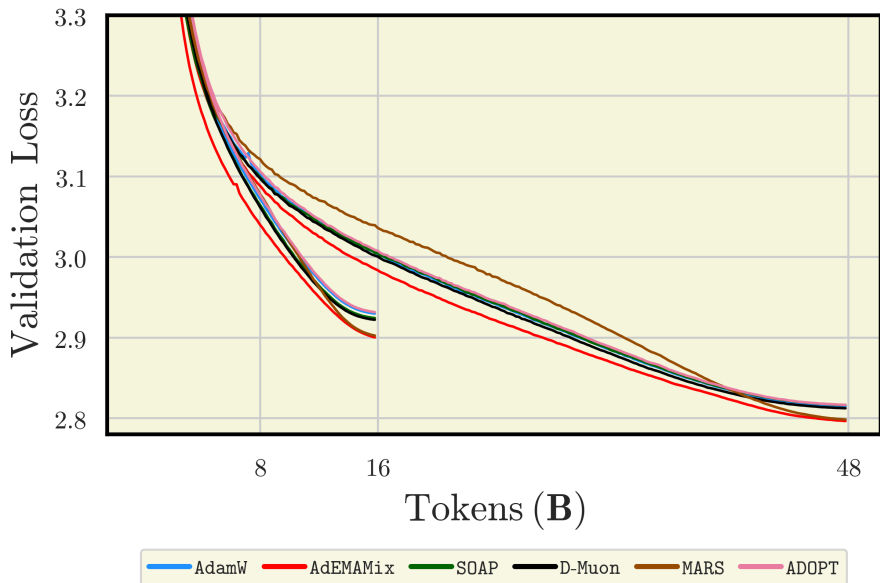


Benchmarking at small scale

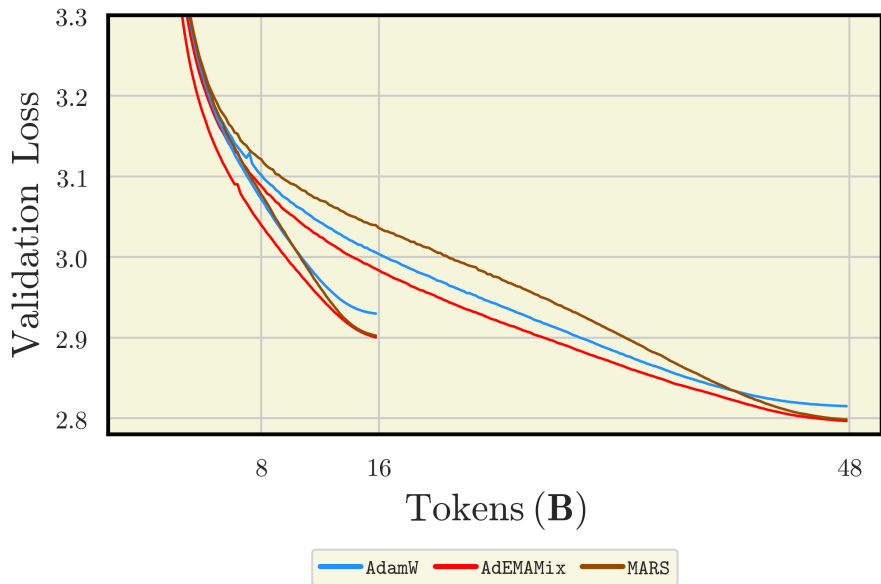
Outperforming AdamW, Batch Size 256 (210M)



Benchmarking: Scaling Up

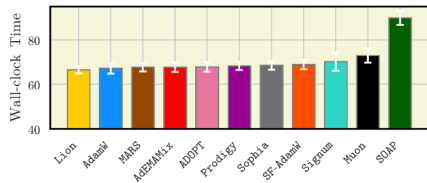


Benchmarking: Scaling Up

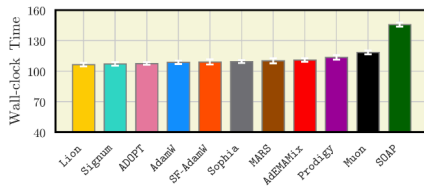


Benchmarking: Wall-Time

124M



210M



The End?

Thanks!