

Batch-wise Convergent Pre-training: Step-by-Step Learning Inspired by Child Language Development

○Ko Yoshida¹, Daiki Shiono¹, Kai Sato¹, Toko Miura¹, Momoka Furuhashi¹, Jun Suzuki^{1,2,3}
¹Tohoku University, ²RIKEN, ³NII LLMC

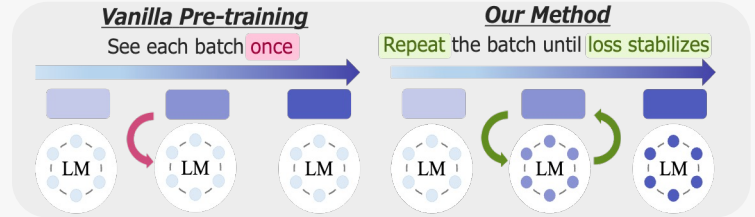
Overview

Concept: What if a LM could gradually accumulate knowledge through repetitions within limited contexts – like children?

Method :

Learn each batch repeatedly with regularization to reduce forgetting

Method Overview



Method: Batch-wise Convergent Pre-training

Step1 ➡ Get \mathcal{L}_{CE} for batch \mathcal{X}_t

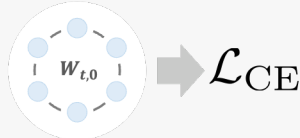
Let's Say:

\mathcal{X}_t : Target Batch

$\mathbf{W}_{t,0}$: Learnable Params^{*1}

α : Adaptive learning strength^{*2}

Calculate \mathcal{L}_{CE}



- *1** • Start with $\mathbf{W}_{t,0} = \mathbf{W}_{t-1,n}$
• Update iteratively to $\mathbf{W}_{t,k}$

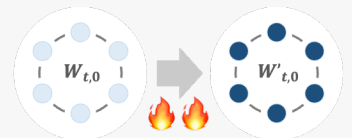
- *2** • Start with $\alpha = 1$
• α grows when $\mathcal{L}_{CE} > \mathcal{L}_{CE_{goal}}$

Step2 ➡ Parameter Update with α

- ①How far from $\mathcal{L}_{CE_{goal}}$ ②Accelerate with α^{*4}

$$\alpha' \leftarrow f(\mathcal{L}_{CE}, \alpha)^{*3}$$

Boost learning for this batch! 🔥



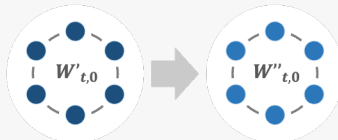
- *3** • What is $f(\mathcal{L}_{CE}, \alpha)$?
 $\alpha' \leftarrow \alpha + \eta(\mathcal{L}_{CE} - \mathcal{L}_{CE_{goal}})$
• η is learning rate for α update

- *4** • Use AdamW
 $\mathbf{W}' \leftarrow \text{AdamW}(\alpha' \mathcal{L}_{CE})$

Step3 ➡ Pull back \mathbf{W}' toward $\mathbf{W}_{t,0}$

- ①Regularization term ②Regularize \mathbf{W}'^{*6}

$$\nabla R = \frac{C}{p} \|\mathbf{W}' - \mathbf{W}_{t,0}\|_p^{p*5}$$



- *5** • C controls the regularization strength
• L_1 when $p = 1$, L_2 when $p = 2$

- *6** • Regularization
 $\mathbf{W}'' = \mathbf{W}' - \nabla R$

Step4 ➡ Convergence Check

Has the model learned \mathcal{X}_t sufficiently?

→ Check α' , Not \mathcal{L}_{CE}

Check List

- ☐ Is α decreasing continuously?
- ☐ $\alpha < \text{threshold}$?

✓ Next Batch \mathcal{X}_{t+1}

$\mathbf{W}_{t+1,0} \leftarrow \mathbf{W}_{t,0}$

✗ Repeat Batch \mathcal{X}_t

$\mathbf{W}_{t,1} \leftarrow \mathbf{W}_{t,0}$

Experiments/Analysis

Comparison with Official Baselines

Setup:

- We compare our model with **GPT-2**, **GPT-BERT**
- Model size is 117M, Qwen2.5 architecture
- A curriculum based on our original difficulty score

Result: No significant improvement

Model	BLiMP↑	BLiMP-S↑	WUG-ADJ↑	Text-Avg.↑
GPT-BERT	80.5	73.0	41.2	70.9
GPT-2	74.9	63.3	50.2	54.7
Ours(p=1)	49.2	50.4	57.5	32.8
Ours(p=2)	52.2	50.2	57.1	32.5

Next TODO : Isolated batch learning breaks the distributional assumption
→ Batch Design? or Architecture Design?

Training Orders & Repetition Strategies

Setup:

1. **Random**: 10 epochs with the random data
2. **Curriculum**: 10 epochs with the curriculum data
3. **Curriculum-Repeat**: the curriculum data + repeat each batch 10 times in a row
4. **Our Proposed Method**

Result: Comparable but not better

Analysis:



Our Proposed Method performs slightly better in Text-Avg, by only 2 points



Random performed best on BLiMP

Curriculum performed best on WUG-ADJ

Curriculum-Repeat performed best on Entity