

Fast and Low-Cost Genomic Foundation Models via Outlier Removal

Haozheng Luo*, Chenghao Qiu*, Maojiang Su, Zhihan Zhou,
Zoe Mehta, Guo Ye, Jerry Yao-Chieh Hu, Han Liu

Northwestern University, Computer Science

✉ <https://arxiv.org/abs/2505.00598>

ICML2025



Northwestern
University



ICML
International Conference
On Machine Learning

Problem: Transformer-based Genomic Foundation Models (GFMs) encounters outlier-inefficient in quantization and fine-tuning.

Proposal: Fast and Low-Cost Genomic Foundation Models (termed GERM) via outlier-removal architecture and continual learning.

- Serves as an outlier-free model structure to address and mitigate outliers introduced by pretrained models and low-rank adaptation,
- Retains and improves the desirable properties of GFMs in quantization and low-rank adaptation,
- All DNABERT fine-tuning tasks finish in only 5 minutes on a single NVIDIA GeForce RTX 2080 Ti GPU.
- Achieves average performance improvements of **37.98%** in finetuning and **64.34%** in quantization.

Outlier: In GFMs, tokens or activations that disproportionately influence the attention mechanism with:

- Tokens with **little or no meaningful information** receive disproportionately **high attention weights**.
- Recurring nucleotide patterns are **overemphasized** by Softmax.

GFMs: Large-scale pretrained models designed for modeling and analysing genomic sequences.

- Trained on **massive** genomic datasets
- **Classification models**: e.g., DNABERT-2, Nucleotide Transformer (NT), HyenaDNA
- **Generative models**: e.g., Evo, GenomeOcean
- Larger GFMs, especially generative models, require **substantial computational resources** for deployment and fine-tuning.

- We propose a new GFM architecture **GERM** by replacing the Softmax in the attention mechanism with Softmax_1 to achieve the **Quantization Robustness** and **Fast Low-rank Adaptation**.

$$\text{Softmax}_1(S) := \frac{\exp(S)}{1 + \sum_{i=1}^L \exp(S_i)},$$

- The original OutEffHop method requires training from scratch; we propose a **trade-off** variant, GERM-T , to achieve sub-optimal performance with small-step continual learning.

Experimental Studies: Outlier-Efficiency and Quantization Results

Compare GERM with the vanilla attention on DNABERT-2 in quantization setting.

	Model	#Bits	Quantization Method	MCC (\uparrow)	Delta MCC (\downarrow)	Avg Performance Drop (\downarrow)	Avg. Kurtosis (\downarrow)	Max inf. norm (\downarrow)
	Official	16W/16A	-	66.11	-	-	<u>39.68</u>	53.61
DNABERT-2	16W/16A	-	-	59.11	7.00	-	270.90	61.64
	8W/8A	-	-	33.60 \pm 0.41	32.51	43.81%		
	8W/8A	-	-	36.51 \pm 0.02	45.37	38.63%		
	6W/6A	-	SmoothQuant	20.74 \pm 0.04	45.37	66.18%		
	4W/4A	-	SmoothQuant	-1.03 \pm 0.06	67.06	101.24%		
	8W/8A	-	Outlier	25.26 \pm 0.02	40.85	57.60%		
	6W/6A	-	Outlier	27.84 \pm 0.28	38.27	52.71%		
	8W/8A	-	OmniQuant	49.92 \pm 0.05	16.19	15.76%		
GERM	6W/6A	-	OmniQuant	48.47 \pm 0.14	17.64	18.61%	21.29	10.62
	4W/4A	-	OmniQuant	2.94 \pm 0.19	63.17	94.78%		
	16W/16A	-	-	59.73	6.38	-		
	8W/8A	-	-	57.30 \pm 0.08	8.81	<u>3.77%</u>		
	8W/8A	-	-	56.65 \pm 0.15	9.46	<u>4.82%</u>		
	6W/6A	-	SmoothQuant	56.48 \pm 0.07	9.63	<u>5.45%</u>		
	4W/4A	-	SmoothQuant	20.05 \pm 0.00	46.06	<u>69.44%</u>		
	8W/8A	-	Outlier	45.87 \pm 0.08	20.24	<u>25.23%</u>		
GERM-T	6W/6A	-	Outlier	40.57 \pm 0.56	25.54	<u>36.27%</u>	251.40	28.49
	8W/8A	-	OmniQuant	55.99 \pm 0.09	10.12	5.95%		
	6W/6A	-	OmniQuant	55.70 \pm 0.03	10.41	<u>6.41%</u>		
	4W/4A	-	OmniQuant	49.42 \pm 0.00	16.69	<u>17.17%</u>		
	16W/16A	-	-	59.30	6.81	-		
	8W/8A	-	-	38.38 \pm 0.15	27.73	<u>35.27%</u>		
	8W/8A	-	-	57.52 \pm 0.00	8.59	<u>3.01%</u>		
	6W/6A	-	SmoothQuant	30.34 \pm 0.04	35.77	<u>48.83%</u>		
GERM-T	4W/4A	-	SmoothQuant	0.22 \pm 0.00	65.89	<u>99.63%</u>	251.40	28.49
	8W/8A	-	Outlier	42.57 \pm 0.05	23.54	<u>28.31%</u>		
	6W/6A	-	Outlier	46.02 \pm 0.06	20.06	<u>22.34%</u>		
	8W/8A	-	OmniQuant	56.80 \pm 0.12	9.31	<u>4.21%</u>		
GERM-T	6W/6A	-	OmniQuant	55.41 \pm 0.00	10.71	<u>6.57%</u>	251.40	28.49
	4W/4A	-	OmniQuant	3.86 \pm 0.00	62.25	<u>93.49%</u>		

Results: GERM achieves an average performance improvement of **64.34%** in PTQ experiments. Similarly, GERM-T shows an average performance improvement of **43.04%** over the same baseline.

Compare GERM with the vanilla attention on NT 2.5B in quantization setting.

Model	#Bits	Quantization Method	MCC	Delta MCC	Average Performance Drop
NT-2.5B-multi	16W/16A	-	56.98	-	-
	6W/6A	-	18.52	38.46	67.50%
	4W/4A	-	1.39	55.59	97.56%
	6W/6A	Outlier	50.23	6.75	11.85%
	4W/4A		40.74	16.24	28.50%
	6W/6A	SmoothQuant	47.23	9.75	17.11%
	4W/4A		35.16	21.82	38.29%
	6W/6A	OmniQuant	49.55	7.43	13.04%
	4W/4A		43.63	13.35	23.43%
GERM (NT-2.5B-multi)	16W/16A	-	57.16	-0.18	-
	6W/6A	-	45.96	11.2	19.59%
	4W/4A	-	42.48	14.68	25.68%
	6W/6A	Outlier	52.24	4.92	<u>8.61%</u>
	4W/4A		49.00	8.16	14.28%
	6W/6A	SmoothQuant	51.95	5.21	<u>9.11%</u>
	4W/4A		48.15	31.09	<u>15.76%</u>
	6W/6A	OmniQuant	52.55	4.61	<u>8.07%</u>
	4W/4A		49.26	7.90	13.82%
GERM-T (NT-2.5B-multi)	16W/16A	-	56.82	0.16	-
	6W/6A	-	32.58	24.24	<u>42.66%</u>
	4W/4A	-	10.49	46.33	<u>81.54%</u>
	6W/6A	Outlier	52.14	4.68	8.24%
	4W/4A		46.24	10.58	<u>18.62%</u>
	6W/6A	SmoothQuant	51.61	5.21	<u>9.17%</u>
	4W/4A		48.12	8.70	15.31%
	6W/6A	OmniQuant	52.43	4.39	7.73%
	4W/4A		47.28	9.54	<u>16.79%</u>

Results: GERM achieves an average performance improvement of **50.83%** in PTQ experiments. Similarly, GERM-T shows an average performance improvement of **36.73%** over the same baseline.

Experimental Studies: Outlier-Efficiency and Low-rank Adaptation Results

Compare GERM with the vanilla attention on DNABERT-2 in low-rank adaptation setting.

Models	Low-Rank Adaptation Method	MCC (\uparrow)	Delta MCC different (\downarrow)	Avg Performance Drop (\downarrow)	Avg. kurtosis(\downarrow)	Max inf. norm(\downarrow)
DNA BERT-2	Full	59.11	7.00	-	270.90	61.41
	LoRA	50.91 \pm 1.67	15.2	13.87%	-	219.20
	QLoRA	50.65 \pm 0.13	15.46	14.31%	292.85	<u>53.91</u>
	LoftQ	50.76 \pm 0.06	15.31	14.05%	299.18	54.18
GERM	Full	59.73	6.38	-	21.29	10.62
	LoRA	57.27 \pm 0.70	8.84	4.12%	-	19.41
	QLoRA	53.16 \pm 0.21	12.95	10.99%	34.29	27.27
	LoftQ	53.11 \pm 0.08	13.00	11.08%	33.02	27.41
GERM-T	Full	59.30	6.81	-	<u>251.40</u>	<u>28.49</u>
	LoRA	55.60 \pm 0.28	10.51	<u>6.23%</u>	-	<u>140.86</u>
	QLoRA	51.05 \pm 0.07	15.06	<u>13.90%</u>	<u>287.95</u>	53.92
	LoftQ	51.20 \pm 0.13	14.91	<u>13.65%</u>	<u>286.16</u>	<u>53.35</u>

Results: GERM achieves an average performance improvement of **37.98%** in low-rank adaptation compared to DNABERT-2 model. Similarly, GERM-T shows an average performance improvement of **20.01%** over the same baseline.

Experimental Studies: Outlier-Efficiency and Low-rank Adaptation Results2

Compare GERM with the vanilla attention on NT 2.5B in low-rank adaptation setting.

Model	Fine-Tuning Method	MCC	Delta MCC	Average Performance Drop
NT-2.5B-multi	Full	56.98	-	-
	LoRA	53.50	3.48	6.11%
	QLoRA	52.29	4.69	8.19%
	LoftQ	52.89	4.09	7.17%
GERM (NT-2.5B-multi)	Full	57.16	-0.18	-
	LoRA	55.98	1.18	2.06%
	QLoRA	55.52	1.64	2.87%
	LoftQ	55.80	1.36	2.38%
GERM-T (NT-2.5B-multi)	Full	56.82	0.16	-
	LoRA	55.24	1.58	<u>2.78%</u>
	QLoRA	53.32	3.50	<u>6.16%</u>
	LoftQ	53.74	3.08	<u>5.42%</u>

Results: GERM achieves an average performance improvement of **66.02%** in low-rank adaptation. Similarly, GERM-T shows an average performance improvement of **34.56%** over the same baseline.

Experimental Studies: Outlier-Efficiency on Various Continual Learning Steps

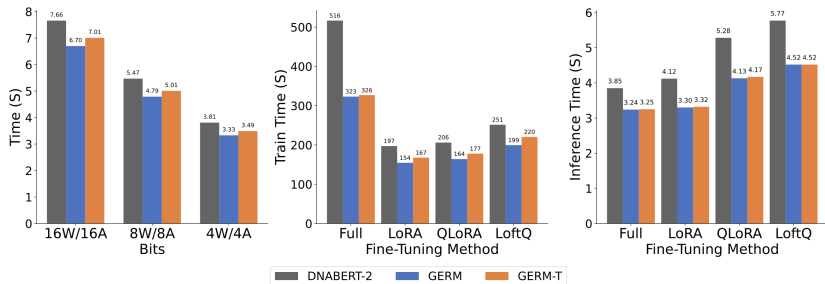
Compare GERM-T with the vanilla attention on various continual learning steps.

Method	Fine-Tuning Method	MCC (\uparrow)	Avg Performance Drop (\downarrow)
DNABERT-2	Full	59.11	-
GERM	Full	59.73	-
Out20k	Full	59.21	-
GERM-T	Full	59.30	-
Out100k	Full	60.56	-
DNABERT-2	LoRA	50.91	13.87%
GERM	LoRA	56.78	4.94%
Out20k	LoRA	54.75	7.53%
GERM-T	LoRA	55.60	<u>6.24%</u>
Out100k	LoRA	56.61	6.52%
DNABERT-2	QLoRA	50.65	14.31%
GERM	QLoRA	53.16	11.00%
Out20k	QLoRA	50.61	14.52%
GERM-T	QLoRA	51.05	<u>13.91%</u>
Out100k	QLoRA	51.24	15.39%
DNABERT-2	LoftQ	50.76	14.13%
GERM	LoftQ	53.11	11.08%
Out20k	LoftQ	50.94	13.97%
GERM-T	LoftQ	51.20	<u>13.66%</u>
Out100k	LoftQ	50.77	16.17%

Results: Our method outperforms the vanilla approach across all test sets. Also, we observe that GERM-T exhibits the most optimal performance drop during quantization and low-rank adaptation compared to other continual learning steps.

Experimental Studies: Comparison of Performance in Resource-Constrained Environments

Compare GERM with the vanilla attention on DNABERT-2 in resource-constrained setting.



Results: Both GERM and GERM-T achieve shorter full-rank fine-tuning times per epoch compared to DNABERT-2. Additionally, the model quantization latency for both GERM and GERM-T is lower than that of DNABERT-2, while delivering superior quantization performance.

Compare GERM with the vanilla attention on DNABERT-2 in CPU-only environments.

Method	Fine-Tuning Method	MCC (\uparrow)	Time (sec.)	
			Train	Inference
DNABERT-2	LoRA	50.91	808.23	29.66
GERM	LoRA	57.27	618.68	23.10
GERM-T	LoRA	<u>55.60</u>	<u>674.40</u>	<u>23.57</u>
DNABERT-2	QLoRA	50.65	516.04	63.17
GERM	QLoRA	53.16	358.34	45.28
GERM-T	QLoRA	<u>51.50</u>	<u>418.13</u>	<u>46.91</u>

Results: Both GERM and GERM-T achieve shorter fine-tuning times per epoch compared to DNABERT-2, with the only exception being QLoRA when deployed, where the time is slightly longer.

- Fast and Low-Cost Genomic Foundation Models
 - Manages outliers in transformer-based GFM.
 - Remove outlier in model pretraining and fine tuning period.
- Theoretical Enhancements
 - Provide expressive guarantee of low-rank adaption.
- Small-Step Continual Learning
 - Leverages continual learning to address the training-from-scratch limitation in [Hu et al., 2024].
 - Achieves sub-optimal yet effective performance.
- Empirical Performance of GERM
 - Achieves 92.14% lower average kurtosis and 82.77% lower maximum infinity norm $|\mathbf{x}|_\infty$, enabling robust quantization and fast low-rank adaptation.
 - Improves fine-tuning performance by 37.98% and quantization performance by 64.34% over the baseline.

Thank You!

Haozheng Luo*, Chenghao Qiu*, Maojiang Su, Zhihan Zhou, Zoe Mehta, Guo Ye, Jerry Yao-Chieh Hu, Han Liu

- ✉ hluo@u.northwestern.edu
- ✉ q1320460765@tju.edu.cn
- ✉ maojiangsu2030@u.northwestern.edu
- ✉ zhihanzhou2020@u.northwestern.edu
- ✉ zoe.mehta@vhhsougars.org
- ✉ guoye2018@u.northwestern.edu
- ✉ jhu@u.northwestern.edu
- ✉ hanliu@northwestern.edu
- 🏠 <https://github.com/MAGICS-LAB/GERM>