

# Numerical Algorithm-driven Generative Model Design and Adaptation

Hongteng Xu

GSAI, RUC


Aug. 10, 2024



中國人民大學  
RENMIN UNIVERSITY OF CHINA


高瓴人工智能学院  
Gaoling School of Artificial Intelligence

# Demands under “Scaling Laws” of AIGC



Kaplan, Jared, et al. Scaling laws for neural language models. arXiv:2001.08361 (2020).

# Demands under “Scaling Laws” of AIGC




Kaplan, Jared, et al. Scaling laws for neural language models. arXiv:2001.08361 (2020).

- ▶ New model architectures with lower complexity
- ▶ Parameter-efficient fine-tuning (PEFT) strategies

# Parameter-efficient Fine-tuning Methods


## Partial Frozen



Sung, Yi-Lin, et al. "Training neural networks with fixed sparse masks."  
NeurIPS 2021.


# Parameter-efficient Fine-tuning Methods

## Partial Frozen



Sung, Yi-Lin, et al. "Training neural networks with fixed sparse masks." NeurIPS 2021.


## Soft prompt fine-tuning



Li, Xiang Lisa, and Percy Liang. "Prefix-Tuning: Optimizing Continuous Prompts for Generation." ACL 2021.

# LoRA: The Mainstream Adapter-based Strategy


## Adapter-based fine-tuning



# LoRA: The Mainstream Adapter-based Strategy

## Adapter-based fine-tuning

Hu, Edward J., et al. "LoRA: Low-Rank Adaptation of Large Language Models." ICLR 2022.




$$W + A \times B$$

# LoRA: The Mainstream Adapter-based Strategy

## Adapter-based fine-tuning



Hu, Edward J., et al. "LoRA: Low-Rank Adaptation of Large Language Models." ICLR 2022.



$$W + A \times B$$


- ▶ **Structure Adjustment:** adjust matrices' ranks
- ▶ **Initialization Improvement:** mainly based on the SVD of weight matrices
- ▶ **Parameter Quantization:** lower bits, sparser matrices

# OFT: Another Potential Solution




Qiu, Zeju, et al. "Controlling text-to-image diffusion by orthogonal finetuning." NeurIPS 2023.

# OFT: Another Potential Solution




Qiu, Zeju, et al. "Controlling text-to-image diffusion by orthogonal finetuning." NeurIPS 2023.



Angular information matters  
→ Orthogonal fine-tuning (OFT)


# Is There A Bridge between LoRA and OFT?



# Householder Reflection: A Simple Orthogonal Transform



Alston Householder




$$H = I - 2uu^\top, \quad u \in \mathbb{S}^{d-1} \quad (1)$$

## Householder Reflection Adaptation (HRA)

- Implement OFT by a chain of Householder reflections


$$z = W \underbrace{\left( \prod_{i=1}^r H_i \right)}_{H^{(r)}} x = W \left( \prod_{i=1}^r (\mathbf{I} - 2\mathbf{u}_i \mathbf{u}_i^\top) \right) x, \text{ with } \{\mathbf{u}_i \in \mathbb{S}^{d-1}\}_{i=1}^r. \quad (2)$$



## Householder Reflection Adaptation (HRA)

- Implement OFT by a chain of Householder reflections




$$z = \mathbf{W} \underbrace{\left( \prod_{i=1}^r \mathbf{H}_i \right)}_{\mathbf{H}^{(r)}} \mathbf{x} = \mathbf{W} \left( \prod_{i=1}^r (\mathbf{I} - 2\mathbf{u}_i \mathbf{u}_i^\top) \right) \mathbf{x}, \text{ with } \{\mathbf{u}_i \in \mathbb{S}^{d-1}\}_{i=1}^r. \quad (2)$$



- Implement  $\mathbf{W}\mathbf{H}^{(r)}\mathbf{x}$  with low complexity ( $\mathcal{O}(d(r + d_{\text{out}}))$ ) for  $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d}$

$$1) \mathbf{x}^{(j+1)} = \mathbf{x}^{(j)} - 2\langle \mathbf{u}_{r-j}, \mathbf{x}^{(j)} \rangle \mathbf{u}_{r-j}, \text{ for } j = 0, \dots, r-1. \quad 2) \mathbf{z} = \mathbf{W}\mathbf{x}^{(r)}. \quad (3)$$




## Comparisons with Existing OFTs

| Method         | OFT   | BOFT   | Our HRA   |
|----------------|---|--|---|
| Implementation | $\mathbf{R}^{(b)} = \text{diag}(\{\mathbf{R}_i\}_{i=1}^{d/b})$                    | $\mathbf{B}^{(m,b)} = \prod_{i=1}^m \mathbf{B}_i^{(b)}$                            | $\mathbf{H}^{(r)} = \prod_{i=1}^r \mathbf{I} - 2\mathbf{u}_i\mathbf{u}_i^\top$      |
| Illustration   |  |  |  |




# Comparisons with Existing OFTs

| Method   | OFT  | BOFT   | Our HRA  |
|--|--|--|--|
| Implementation   | $\mathbf{R}^{(b)} = \text{diag}(\{\mathbf{R}_i\}_{i=1}^{d/b})$ | $\mathbf{B}^{(m,b)} = \prod_{i=1}^m \mathbf{B}_i^{(b)}$  | $\mathbf{H}^{(r)} = \prod_{i=1}^r \mathbf{I} - 2\mathbf{u}_i\mathbf{u}_i^\top$ |
| Illustration   |  |  |  |
| #Parameters  | $\frac{d(b-1)}{2} \sim db$                                     | $\frac{dm(b-1)}{2} \sim dm b$  | $rd$   |
| Complexity   | $\mathcal{O}(d(b^2 + b + d_{\text{out}}))$                     | $\mathcal{O}(d((b^2 + b)m + d_{\text{out}})) \sim \mathcal{O}(d((b^2 + d)m + d_{\text{out}}))$ | $\mathcal{O}(d(r + d_{\text{out}}))$   |
| Cayley Transform: $\mathbf{R} = (\mathbf{I} - \mathbf{A})(\mathbf{I} + \mathbf{A})^{-1}$ |  |  |  |

# Comparisons with Existing OFTs

| Method   | OFT   | BOFT   | Our HRA   |
|--|---|--|---|
| Implementation   | $\mathbf{R}^{(b)} = \text{diag}(\{\mathbf{R}_i\}_{i=1}^{d/b})$                    | $\mathbf{B}^{(m,b)} = \prod_{i=1}^m \mathbf{B}_i^{(b)}$  | $\mathbf{H}^{(r)} = \prod_{i=1}^r \mathbf{I} - 2\mathbf{u}_i\mathbf{u}_i^\top$      |
| Illustration   |  |              |  |
| #Parameters  | $\frac{d(b-1)}{2} \sim db$  | $\frac{dm(b-1)}{2} \sim dmb$   | $rd$  |
| Complexity   | $\mathcal{O}(d(b^2 + b + d_{\text{out}}))$  | $\mathcal{O}(d((b^2 + b)m + d_{\text{out}})) \sim \mathcal{O}(d((b^2 + d)m + d_{\text{out}}))$ | $\mathcal{O}(d(r + d_{\text{out}}))$  |
| Cayley Transform: $\mathbf{R} = (\mathbf{I} - \mathbf{A})(\mathbf{I} + \mathbf{A})^{-1}$ |   |  |   |
| Cover $\mathbb{O}_{d \times d}$  | $b = d$   | $\{\mathbf{B}_i^{(m=\log d, b=2)}\}_{i=1}^{d-1}$   | $\{\mathbf{u}_i\}_{i=1}^{d-1}$  |

# Comparisons with Existing OFTs

| Method   | OFT   | BOFT   | Our HRA   |
|--|---|--|---|
| Implementation   | $\mathbf{R}^{(b)} = \text{diag}(\{\mathbf{R}_i\}_{i=1}^{d/b})$                    | $\mathbf{B}^{(m,b)} = \prod_{i=1}^m \mathbf{B}_i^{(b)}$  | $\mathbf{H}^{(r)} = \prod_{i=1}^r \mathbf{I} - 2\mathbf{u}_i\mathbf{u}_i^\top$      |
| Illustration   |  |              |  |
| #Parameters  | $\frac{d(b-1)}{2} \sim db$  | $\frac{dm(b-1)}{2} \sim dmb$   | $rd$  |
| Complexity   | $\mathcal{O}(d(b^2 + b + d_{\text{out}}))$  | $\mathcal{O}(d((b^2 + b)m + d_{\text{out}})) \sim \mathcal{O}(d((b^2 + d)m + d_{\text{out}}))$ | $\mathcal{O}(d(r + d_{\text{out}}))$  |
| Cayley Transform: $\mathbf{R} = (\mathbf{I} - \mathbf{A})(\mathbf{I} + \mathbf{A})^{-1}$ |   |  |   |
| Cover $\mathbb{O}_{d \times d}$  | $b = d$   | $\{\mathbf{B}_i^{(m=\log d, b=2)}\}_{i=1}^{d-1}$   | $\{\mathbf{u}_i\}_{i=1}^{d-1}$  |

Have potentials to be more efficient in practice.

## Connections to LoRA: HRA is An Adaptive LoRA

- Reformulation of the HR chain:

$$\mathbf{H}^{(r)} = \prod_{i=1}^r (\mathbf{I} - 2\mathbf{u}_i \mathbf{u}_i^\top) = \mathbf{I} + \mathbf{U}_r \boldsymbol{\Gamma}_r \mathbf{U}_r^\top, \quad (4)$$


- $\boldsymbol{\Gamma}_r = [\gamma_{ij}] \in \mathbb{R}^{r \times r}$  is a upper-triangular matrix, and its upper-triangular element is

$$\gamma_{ij} = \begin{cases} -2 & i = j \\ (-2)^{j-i+1} \prod_{i=1}^{j-1} \langle \mathbf{u}_i, \mathbf{u}_{i+1} \rangle & i < j. \end{cases} \quad (5)$$

- HRA is equivalent to an adaptive LoRA, making  $\text{Range}(\mathbf{W})$  unchanged.

$$\mathbf{W} \mathbf{H}^{(r)} = \mathbf{W} + \underbrace{\mathbf{W} \mathbf{U}_r \boldsymbol{\Gamma}_r \mathbf{U}_r^\top}_{\mathbf{A}_{\mathbf{W}, \mathbf{U}}}. \quad (6)$$

# Orthogonality: The Key of Balancing Expressiveness and Regularity



$$\min_{\{\mathbf{U}_r^{(l)}\}_{l=1}^L} \text{Loss}(\mathcal{D}; \{\mathbf{U}_r^{(l)}\}_{l=1}^L) + \lambda \underbrace{\sum_{l=1}^L \|\mathbf{I}_r - (\mathbf{U}_r^{(l)})^\top \mathbf{U}_r^{(l)}\|_F^2}_{\text{Orthogonal regularizer}}, \quad (7)$$

- $\lambda \in [0, \infty)$ : Normalization
- $\lambda = \infty$ : (Modified) Gram-Schmidt Orthogonalization

## Experiments: NLP Tasks

Table: Results (%) of various methods on GLUE development set.

| Method                              | #Param | MNLI         | SST-2        | CoLA         | QQP          | QNLI         | RTE          | MRPC         | STS-B        | All          |
|-------------------------------------|--------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Full Fine-tune                      | 184M   | 89.90        | 95.63        | 69.19        | <b>92.40</b> | 94.03        | 83.75        | 89.46        | 91.60        | 88.25        |
| BitFit                              | 0.10M  | 89.37        | 94.84        | 66.96        | 88.41        | 92.24        | 78.70        | 87.75        | 91.35        | 86.20        |
| H-Adapter                           | 1.22M  | 90.13        | 95.53        | 68.64        | 91.91        | 94.11        | 84.48        | 89.95        | 91.48        | 88.28        |
| P-Adapter                           | 1.18M  | 90.33        | 95.61        | 68.77        | 92.04        | 94.29        | 85.20        | 89.46        | 91.54        | 88.41        |
| LoRA <sub>r=8</sub>                 | 1.33M  | 90.65        | 94.95        | 69.82        | 91.99        | 93.87        | 85.20        | 89.95        | 91.60        | 88.50        |
| AdaLoRA                             | 1.27M  | <b>90.76</b> | 96.10        | 71.45        | <u>92.23</u> | <u>94.55</u> | 88.09        | 90.69        | 91.84        | 89.46        |
| OFT <sub>b=16</sub>                 | 0.79M  | 90.33        | 96.33        | <b>73.91</b> | 92.10        | 94.07        | 87.36        | 92.16        | <u>91.91</u> | 89.77        |
| BOFT <sub>b=8</sub> <sup>m=2</sup>  | 0.75M  | 90.25        | 96.44        | 72.95        | 92.10        | 94.23        | <u>88.81</u> | 92.40        | <b>91.92</b> | 89.89        |
| HRA <sub>r=8, \lambda=0</sub>       | 0.66M  | <u>90.70</u> | <u>96.45</u> | <u>73.70</u> | 91.29        | <b>94.66</b> | 88.45        | <u>93.69</u> | 91.86        | <b>90.10</b> |
| HRA <sub>r=8, \lambda=10^{-5}</sub> | 0.66M  | 90.43        | <b>96.79</b> | 71.91        | 91.02        | 94.44        | <b>89.53</b> | <b>94.10</b> | 91.74        | <u>90.00</u> |
| HRA <sub>r=8, \lambda=\infty</sub>  | 0.66M  | 90.52        | 95.87        | 70.71        | 90.71        | 94.12        | 87.00        | 92.59        | 91.54        | 89.13        |

# Experiments: Controllable Text-to-Image Generation




Figure: Qualitative results on subject-driven generation.

## Experiments: Controllable Text-to-Image Generation



Prompt: A baseball game being played.



Prompt: A man smiling for the camera.



Prompt: A tree stump.


# Experiments: Controllable Text-to-Image Generation

Table: Results of various methods on subject-driven generation and controllable generation.

| Method                     | #Param<br>(M) | Subject-driven generation |              |               |              | #Param<br>(M) | Controllable generation |              |              |              |              |             |  |
|----------------------------|---------------|---------------------------|--------------|---------------|--------------|---------------|-------------------------|--------------|--------------|--------------|--------------|-------------|--|
|                            |               | Image fidelity            |              | Text fidelity |              |               | C2I                     |              | S2I          |              | L2F          |             |  |
|                            |               | DINO↑                     | CLIP-I↑      | CLIP-T↑       | LPIPS↑       |               | IoU↑                    | F1↑          | mIoU↑        | mAcc↑        | aAcc↑        | Error↓      |  |
| Real Images                | -             | 0.764                     | 0.890        | -             | 0.562        | -             | -                       | -            | -            | -            | -            | -           |  |
| DreamBooth                 | 859.52        | 0.614                     | 0.778        | 0.239         | 0.737        | 859.52        | 0.049                   | 0.093        | 7.72         | 14.40        | 33.61        | 146.19      |  |
| ControlNet                 | -             | -                         | -            | -             | -            | 361.30        | 0.189                   | 0.317        | 20.88        | 30.91        | 61.42        | 7.61        |  |
| T2I-Adapter                | -             | -                         | -            | -             | -            | 77.00         | 0.078                   | 0.143        | 16.38        | 26.31        | 51.63        | 23.75       |  |
| LoRA                       | 0.8           | 0.613                     | 0.765        | 0.237         | 0.744        | 1.25          | 0.168                   | 0.286        | 22.98        | 35.52        | 58.03        | 7.68        |  |
| COFT $b=4$                 | 23.3          | 0.630                     | 0.783        | 0.235         | 0.744        | 26.40         | 0.195                   | 0.325        | 26.92        | 40.08        | 62.96        | 6.92        |  |
| OFT $b=4$                  | 23.3          | 0.632                     | 0.785        | 0.237         | 0.746        | 26.40         | 0.193                   | 0.323        | 27.06        | 40.09        | 62.42        | 7.07        |  |
| BOFT $m=4$<br>$r=8$        | -             | -                         | -            | -             | -            | 20.76         | -                       | -            | 28.83        | 41.24        | 67.74        | 5.67        |  |
| HRA $r=8, \lambda=0$       | 0.69          | <b>0.670</b>              | <b>0.803</b> | 0.238         | 0.758        | 0.89          | <b>0.213</b>            | <b>0.350</b> | <b>29.45</b> | <b>42.02</b> | 66.83        | <u>5.56</u> |  |
| HRA $r=8, \lambda=10^{-3}$ | 0.69          | <u>0.661</u>              | <u>0.799</u> | <u>0.255</u>  | <u>0.760</u> | 0.89          | <u>0.205</u>            | <u>0.339</u> | <u>29.27</u> | <u>40.89</u> | <b>67.86</b> | <b>5.46</b> |  |
| HRA $r=8, \lambda=\infty$  | 0.69          | 0.651                     | 0.794        | <b>0.274</b>  | <b>0.778</b> | 0.89          | 0.201                   | 0.334        | 28.15        | 40.22        | 64.95        | 11.11       |  |


# Experiments: Mathematical Reasoning

| Method                      | Param. Ratio | GSM8K       | MATH       |
|-----------------------------|--------------|-------------|------------|
| LLaMA2-7B                   | -            | 14.6        | 2.5        |
| LoRA $r=32$                 | 0.25%        | 50.2        | 7.8        |
| OFT $b=16$                  | 0.13%        | 50.1        | 8.4        |
| BOFT $^{m=2}_{b=8}$         | 0.12%        | 50.6        | 8.6        |
| PiSSA                       | 4.75%        | 53.1        | 7.4        |
| HRA $r=8, \lambda=0$        | 0.03%        | 47.1        | 6.6        |
| HRA $r=16, \lambda=0$       | 0.06%        | 52.1        | 8.1        |
| HRA $r=32, \lambda=0$       | 0.12%        | <u>55.8</u> | 9.0        |
| HRA $r=32, \lambda=\infty$  | 0.12%        | 52.8        | <u>9.2</u> |
| HRA $r=32, \lambda=10^{-4}$ | 0.12%        | <b>56.3</b> | <b>9.3</b> |



# Experiments: Mathematical Reasoning

| Method                      | Param. Ratio | GSM8K       | MATH       |
|-----------------------------|--------------|-------------|------------|
| LLaMA2-7B                   | -            | 14.6        | 2.5        |
| LoRA $r=32$                 | 0.25%        | 50.2        | 7.8        |
| OFT $b=16$                  | 0.13%        | 50.1        | 8.4        |
| BOFT $^{m=2}_{b=8}$         | 0.12%        | 50.6        | 8.6        |
| PiSSA                       | 4.75%        | 53.1        | 7.4        |
| HRA $r=8, \lambda=0$        | 0.03%        | 47.1        | 6.6        |
| HRA $r=16, \lambda=0$       | 0.06%        | 52.1        | 8.1        |
| HRA $r=32, \lambda=0$       | 0.12%        | 55.8        | 9.0        |
| HRA $r=32, \lambda=\infty$  | 0.12%        | 52.8        | 9.2        |
| HRA $r=32, \lambda=10^{-4}$ | 0.12%        | <b>56.3</b> | <b>9.3</b> |



| Method                      | Param. Ratio | ARC   | HellaSwag | MMLU  | Winogrande | HumanEval |
|-----------------------------|--------------|-------|-----------|-------|------------|-----------|
| LLaMA2-7B                   | -            | 49.74 | 58.90     | 45.92 | 74.11      | 12.80     |
| LoRA $r=16$                 | 0.12%        | 48.81 | 56.89     | 40.60 | 71.27      | 11.59     |
| HRA $r=32, \lambda=10^{-4}$ | 0.12%        | 49.57 | 57.72     | 41.20 | 73.32      | 13.41     |

## Discussion: From Model Adaptation to Model Design



- ▶ Householder reflection  $\mathbf{H}\mathbf{x} = \mathbf{x} - 2\mathbf{u}\mathbf{u}^\top\mathbf{x}$  can work as a (restricted) ResNet.
- ▶  $|\mathbf{x}|$  can be treated as an adaptive Householder reflection transform.

## Discussion: From Model Adaptation to Model Design


- ▶ Householder reflection  $\mathbf{H}\mathbf{x} = \mathbf{x} - 2\mathbf{u}\mathbf{u}^\top\mathbf{x}$  can work as a (restricted) ResNet.
- ▶  $|\mathbf{x}|$  can be treated as an adaptive Householder reflection transform.
- ▶ **Stacking HRs leads to a Lipschitz-1 Network.**
  - ▶ Support deeper networks without other tricks (LayerNorm, ResNet connection, ...)
  - ▶ Robust to attack because of its Lipschitz-1 property.

## Discussion: From Model Adaptation to Model Design

- ▶ Householder reflection  $Hx = x - 2uu^\top x$  can work as a (restricted) ResNet.
- ▶  $|x|$  can be treated as an adaptive Householder reflection transform.
- ▶ **Stacking HRs leads to a Lipschitz-1 Network.**
  - ▶ Support deeper networks without other tricks (LayerNorm, ResNet connection, ...)
  - ▶ Robust to attack because of its Lipschitz-1 property.



# Roles of “Traditional” Algorithms in The Era of AIGC




## Discrete Algorithms (Computer Science)

- ▶ Sorting, Searching, Discrete Optimization, ...

## Continuous Algorithms (Scientific Computing)

- ▶ SVD, Eigenproblem Solver, Equation Solver, ...

# Roles of “Traditional” Algorithms in The Era of AIGC



## Discrete Algorithms (Computer Science)

- ▶ Sorting, Searching, Discrete Optimization, ...

## Continuous Algorithms (Scientific Computing)

- ▶ SVD, Eigenproblem Solver, Equation Solver, ...

- ▶ Ironically, everyone admits the significance of algorithms, while few of them fully exploit the power of algorithms in their research:(
- ▶ “Traditional” algorithm matters!

# Summary

## Householder Reflection Adaptation

- ▶ Bridge the gap between low-rank and orthogonal adaptation
- ▶ An typical and effective application of algorithm to model adaptation
- ▶ Have potentials to model design

### Resources

- ▶ Paper: <https://arxiv.org/pdf/2405.17484>
- ▶ Code: <https://github.com/DaShenZi721/HRA>
- ▶ Email: hongtengxu@ruc.edu.cn

# Summary

## Householder Reflection Adaptation

- ▶ Bridge the gap between low-rank and orthogonal adaptation
- ▶ An typical and effective application of algorithm to model adaptation
- ▶ Have potentials to model design

## Resources

- ▶ Paper: <https://arxiv.org/pdf/2405.17484>
- ▶ Code: <https://github.com/DaShenZi721/HRA>
- ▶ Email: hongtengxu@ruc.edu.cn

Thanks! Q & A