

# Orthonormality changes to Simplicits

Rishit Dagli

## Contents

<b>1</b>	<b>Motivation</b>	<b>2</b>
<b>2</b>	<b>What do I do?</b>	<b>2</b>
2.1	Goals . . . . .	2
2.2	Main Idea . . . . .	2
2.3	Algorithm . . . . .	2
<b>3</b>	<b>Implementation Details</b>	<b>4</b>
3.1	Parameter State Management . . . . .	4
3.2	Distributed Processing . . . . .	4
3.3	Memory Considerations . . . . .	4

# 1 Motivation

When minimizing elastic energy functions, models often collapse to trivial solutions where all handles encode constant weights, effectively limiting the model to rigid motions.

$$\mathcal{L}_{\text{total}} = \lambda_1 \cdot \frac{\mathcal{E}_{\text{elastic}}}{N} + \lambda_2 \|\mathbf{W}^T \mathbf{W} - \mathbf{I}\|_F^2$$

where  $\mathcal{E}_{\text{elastic}}$  is the elastic energy,  $N$  is the number of points, the second term enforces orthonormality of the outputs  $\mathbf{W}$ , and  $\lambda_1, \lambda_2$  are Lagrangian multipliers.

Originally, in Simplicits  $\lambda_1 = 10^{-1}$  and  $\lambda_2 = 10^6$ , which makes this way of enforcing these constraints seem like a hack to me. Thus, I was thinking of better ways to handle this.

## 2 What do I do?

### 2.1 Goals

Before stating the Algorithm, let us state our goals. We want to (1) at least get as good quality as the paper (we get better quality) and (2) run in times comparable to those listed in the supplementary (we run faster from my experiments!).

### 2.2 Main Idea

The primary idea is for 2D parameters, we apply a new recipe based on the quintic Newton-Schulz iteration to approximate orthogonalization:

$$\begin{aligned} \mathbf{X}_0 &= \frac{\mathbf{X}}{\|\mathbf{X}\| + \epsilon} \\ \mathbf{X}_{k+1} &= a\mathbf{X}_k + (b\mathbf{A}_k + c\mathbf{A}_k^2)\mathbf{X}_k \\ \text{where } \mathbf{A}_k &= \mathbf{X}_k \mathbf{X}_k^T \end{aligned}$$

with coefficients:

$$a = 3.4445, \quad b = -4.7750, \quad c = 2.0315$$

This iteration approximates the polar factor in the polar decomposition, effectively orthogonalizing the matrix while preserving its essential structure.

### 2.3 Algorithm

I present the algorithm for training in Algorithm 2. Alongside using this algorithm, I also remove the orthogonalization term from the loss completely.

---

**Algorithm 1** A Modified Newton-Schulz Iteration.

---

**Require:** Matrix  $X$ , Steps  $K$ 

```
1: function MODIFIEDNS( $X$ ,  $K$ )
2:   Convert  $X$  to bfloat16                                ▷ Reduce memory usage and improve speed
3:    $X \leftarrow \frac{X}{(\|X\|+\epsilon)}$ 
4:   if rows < columns then                                ▷ Work with tall matrices if they are not square
5:      $X \leftarrow X^T$ 
6:   end if
7:   for  $k = 1$  to  $K$  do
8:      $A \leftarrow XX^T$                                        ▷ Gram matrix
9:      $B \leftarrow bA + cA^2$                                    ▷ Quintic parameters
10:     $X \leftarrow aX + BX$                                        ▷ Quintic iteration
11:    if  $X$  contains NaN or Inf then
12:       $X \leftarrow \frac{X}{\|X\|}$ 
13:    end if
14:  end for
15:  if transposed in line 5 then
16:     $X \leftarrow X^T$ 
17:  end if
18:  Convert  $X$  to float32
19:  return  $X$ 
20: end function
```

---

---

**Algorithm 2** Algorithm for Training, uses Algorithm 1.

---

**Require:** Parameters  $\theta$ , Learning rate  $\alpha$ , Momentum coefficients  $\beta_1, \beta_2$ 

```
1: for all parameter  $p$  in  $\theta$  do                                ▷ Adam Updates
2:   Compute gradient  $g$ 
3:    $g \leftarrow \text{clip}(g, \text{max\_norm} = 1.0)$ 
4:    $m \leftarrow \beta_1 m + (1 - \beta_1)g$ 
5:    $v \leftarrow \beta_2 v + (1 - \beta_2)g^2$ 
6:    $\hat{m} \leftarrow m / (1 - \beta_1^t)$ 
7:    $\hat{v} \leftarrow v / (1 - \beta_2^t)$ 
8:    $p \leftarrow p - \alpha \hat{m} / \sqrt{\hat{v} + \epsilon}$ 
9: end for
10:
11: for all 2D parameters do
12:   Partition parameters across GPUs
13:   for each GPU  $i$  do
14:     for all assigned parameter  $p$  do
15:        $p \leftarrow \text{ModifiedNS}(p)$ 
16:     end for
17:   end for
18:   Synchronize updates via all-reduce
19:   Apply synchronized updates
20: end for
```

---

## 3 Implementation Details

### 3.1 Parameter State Management

Each parameter maintains the following state:

1. Step count ( $t$ )
2. Exponential average of gradients ( $\mathbf{m}$ )
3. Exponential average of squared gradients ( $\mathbf{v}$ )

### 3.2 Distributed Processing

For  $N$  GPUs:

1. Each GPU processes approximately  $M/N$  parameters (where  $M$  is total 2D parameters)
2. Updates are accumulated in a flat buffer
3. All-reduce operation synchronizes updates across devices
4. Updates are redistributed to original parameter shapes

### 3.3 Memory Considerations

The optimizer requires the following additional memory per parameter:

- First moment buffer: Same size as parameter
- Second moment buffer: Same size as parameter
- Temporary buffers for Newton-Schulz:  $\approx 3\times$  parameter size for 2D parameters