

1 **A Supplementary Material**

2 **A.1 Backward pass for ZeCO with All-Scan communication**

3 In the backward propagation of the ZeCO algorithm, most of the process is similar to the forward
4 propagation. It is important to note the difference in notation here: $\tilde{\gamma}_{[n]}$ denotes the decay factor for
5 the reverse cumulative product. Furthermore, in the official implementation of gated linear attention,
6 $\mathbf{S}_{[n]}$ needs to be recomputed during the backward pass. However, since the global initial state has
7 already been obtained during the forward pass, there is no need for all-scan communication when
8 recomputing $\mathbf{S}_{[n]}$.

Algorithm 1 Backward pass for ZeCO with All-Scan communication

Input: $\mathbf{Q}, \mathbf{K}, \mathbf{G} \in \mathbb{R}^{L \times d_k}$, $\mathbf{V}, \mathbf{dO} \in \mathbb{R}^{L \times d_v}$, chunk size C , num_device P , device_rank $p \in \{0, 1, \dots, P-1\}$

Initialize $\mathbf{dS} = \mathbf{0} \in \mathbb{R}^{d_k \times d_v}$ on SRAM

```

1: for  $n \leftarrow N$  to 0 do
2:   Load  $\mathbf{G}_{[n]} \in \mathbb{R}^{C \times d_k}$  from HBM to SRAM
3:   Load  $\mathbf{Q}_{[n]} \in \mathbb{R}^{C \times d_k}$  from HBM to SRAM
4:   Load  $\mathbf{dO}_{[n]} \in \mathbb{R}^{C \times d_v}$  from HBM to SRAM
5:   On chip, compute  $\gamma_{[n]}, \Gamma_{[n]}$  and  $\tilde{\mathbf{Q}}_{[n]} = \mathbf{Q}_{[n]} \odot \mathbf{G}_{[n]}$ ,  $\tilde{\gamma} = \tilde{\gamma} \odot \gamma_{[n]}$ 
6:   Store  $\tilde{\gamma}$  in HBM as  $\tilde{\gamma}_{[n]}$ 
7:   On chip, compute  $\mathbf{dS} = (\gamma_{[n]}^\top \mathbf{1}) \odot \mathbf{dS} + \tilde{\mathbf{Q}}_{[n]}^\top \mathbf{dO}_{[n]}$ 
8:   Store  $\mathbf{dS}$  in HBM as  $\mathbf{dS}_{[n]}$ 
9: end for
10: In parallel do:
11: parallel stream 1:
12:  $\mathbf{dS}_{(p-1)L}, \mathbf{dS}_{pL} \leftarrow \text{All-Scan}(\mathbf{dS}_{[0]}, \tilde{\gamma}_{[0]})$ 
13: parallel stream 2:
14: Load  $\mathbf{S}_{(p-1)L}$  from HBM to SRAM
15: On chip, recompute  $\mathbf{S}_{[n]}$  with  $\mathbf{S}_{[0]} = \mathbf{S}_{(p-1)L}$ ,  $n = \{0, 1, 2, \dots, N-1\}$ 
16: Store  $\{\mathbf{S}_{[n]}, n \in \{0, 1, 2, \dots, N-1\}\}$ 
17: for  $n \leftarrow 1$  to  $N$  in parallel do
18:   Load  $\mathbf{Q}_{[n]}, \mathbf{K}_{[n]}, \mathbf{G}_{[n]}, \mathbf{V}_{[n]}, \mathbf{dO}_{[n]}$  from HBM to SRAM
19:   Load  $\gamma \in \mathbb{R}^{d_k \times d_v}$  from HBM to SRAM
20:   On chip, construct causal mask  $\mathbf{M} \in \mathbb{R}^{B \times B}$ 
21:   On chip, compute  $\Lambda_{[n]}, \Gamma_{[n]} \in \mathbb{R}^{C \times d_k}$ 
22:   On chip, compute  $\tilde{\mathbf{Q}}_{[n]} = \mathbf{Q}_{[n]} \odot \Lambda_{[n]}$ ,  $\tilde{\mathbf{K}}_{[n]} = \mathbf{K}_{[n]} \odot \Gamma_{[n]}$ 
23:   On chip, compute  $\mathbf{P}_{[n]} = (\tilde{\mathbf{Q}}_{[n]} \tilde{\mathbf{K}}_{[n]}^\top) \odot \mathbf{M} \in \mathbb{R}^{C \times C}$ 
24:   On chip, compute  $\mathbf{dP}_{[n]} = (\mathbf{dO}_{[n]} \mathbf{V}_{[n]}^\top) \odot \mathbf{M}$ 
25:   On chip, compute  $\mathbf{d}\tilde{\mathbf{K}}_{[n]} = \tilde{\mathbf{Q}}_{[n]}^\top \mathbf{dP}_{[n]}$ 
26:   On chip, compute  $\mathbf{d}\mathbf{K}_{[n]} = \mathbf{d}\tilde{\mathbf{K}}_{[n]} / \Lambda_{[n]}$ 
27:   On chip, compute  $\mathbf{d}\tilde{\mathbf{Q}}_{[n]} = \mathbf{dP}_{[n]} \tilde{\mathbf{K}}_{[n]}$ 
28:   On chip, compute  $\mathbf{dQ}_{[n]} = \mathbf{d}\tilde{\mathbf{Q}}_{[n]} \odot \Lambda_{[n]}$ 
29:   Store  $\mathbf{P}_{[n]}, \mathbf{dQ}_{[n]}, \mathbf{dK}_{[n]}$  in HBM.
30: end for
31: stream barrier
32: for  $n \leftarrow 1$  to  $N$  in parallel do
33:   Load  $\mathbf{P}_{[n]}, \mathbf{dQ}_{[n]}, \mathbf{dK}_{[n]}, \mathbf{dO}_{[n]}, \mathbf{Q}_{[n]}, \mathbf{K}_{[n]}, \mathbf{G}_{[n]}, \tilde{\gamma}_{[n-1]}, \mathbf{dS}_{pL}, \mathbf{S}_{[n-1]}$  from HBM to SRAM
34:   On chip, compute  $\Lambda_{[n]}, \Gamma_{[n]} \in \mathbb{R}^{C \times d_k}$ 
35:   On chip, compute  $\tilde{\mathbf{K}}_{[n]} = \mathbf{K}_{[n]} \odot \Gamma_{[n]}$ 
36:   On chip, compute  $\mathbf{d}\tilde{\mathbf{K}}_{[n]} = \mathbf{V}_{[n]} (\mathbf{dS}_{[n-1]}^\top + (\tilde{\gamma}_{[n-1]}^\top \mathbf{1}) \odot \mathbf{dS}_{pL}^\top)$ 
37:   On chip, compute  $\mathbf{d}\mathbf{K}_{[n]} = \mathbf{d}\mathbf{K}_{[n]} + \mathbf{d}\tilde{\mathbf{K}}_{[n]} \odot \Gamma_{[n]}$ 
38:   On chip, compute  $\mathbf{d}\tilde{\mathbf{Q}}_{[n]} = \mathbf{dO}_{[n]} \mathbf{S}_{[n-1]}^\top$ 
39:   On chip, compute  $\mathbf{dQ}_{[n]} = \mathbf{dQ}_{[n]} + \mathbf{d}\tilde{\mathbf{Q}}_{[n]} \odot \Lambda_{[n]}$ 
40:   On chip, compute  $\mathbf{dV}_{[n]} = \mathbf{P}_{[n]}^\top \mathbf{dO}_{[n]} + \tilde{\mathbf{K}}_{[n]} (\mathbf{dS}_{[n-1]}^\top + (\tilde{\gamma}_{[n-1]}^\top \mathbf{1}) \odot \mathbf{dS}_{pL}^\top)$ 
41:   Store  $\mathbf{dK}_{[n]}, \mathbf{dV}_{[n]}$  in HBM
42: end for
43: Let  $\mathbf{dQ} = \{\mathbf{dQ}_{[1]}, \dots, \mathbf{dQ}_{[N]}\}$ ,  $\mathbf{dK} = \{\mathbf{dK}_{[1]}, \dots, \mathbf{dK}_{[N]}\}$ ,  $\mathbf{dV} = \{\mathbf{dV}_{[1]}, \dots, \mathbf{dV}_{[N]}\}$ 
44: Compute  $\mathbf{dA} = \mathbf{Q} \odot \mathbf{dQ} - \mathbf{K} \odot \mathbf{dK}$ ,  $\mathbf{dG} = \text{revcum}(\mathbf{dA})$ 
45: return  $\mathbf{dQ}, \mathbf{dK}, \mathbf{dV}, \mathbf{dG}$ 

```
