

A Algorithm for Multi-step Selective-Scan-Update kernel

Algorithm 1 Fused multi-step State-Update kernel. This Algorithm is simplified with the only input being x, A, B, C, Δ, h in reference to the implementation of [10].

Notation $Z^* \leftarrow Z$: load Z into shared memory for computation. $Z \leftarrow Z^*$: Storing Z out to High-Bandwidth-Memory (HBM). Any variable with $*$ means that this variable is in shared memory, everything else is in HBM.

Shapes B - Batch, L - Sequence Length, H - Num heads, P - State Dimension, D - Head Dimension, G - Num groups

Input: $h_0 : (B, H, D, P)$ - initial state, $x : (B, L, H, D)$ $A : (H, D, P)$, $B : (B, L, G, P)$, $C : (B, L, G, P)$, $\Delta : (B, L, H, D)$

Output: $y : (B, L, H, D)$ - Output, $\hat{h} : (B, H, D, P)$ - updated state

```

1:  $h^* \leftarrow h_0$ 
2:  $i \leftarrow 0$ 
3: while  $i < L$  do
4:    $x^* \leftarrow x[:, i, :, :]$ 
5:    $\Delta^* \leftarrow \Delta[:, i, :, :]$ 
6:    $B^* \leftarrow B[:, i, :, :]$ 
7:    $C^* \leftarrow C[:, i, :, :]$ 
8:    $A^* \leftarrow A$ 
9:    $\bar{A}^* \leftarrow \exp(\Delta^* \times A^*)$ 
10:   $\bar{B}^* \leftarrow \Delta^* \times B^*$ 
11:   $h^* \leftarrow h^* \times \bar{A}^* + \bar{B}^* \times x^*$  ▷ State is not stored out to HBM at this point
12:   $y^* \leftarrow C^* \times h^*$ 
13:   $y[:, i, :, :] \leftarrow y^*$ 
14: end while
15:  $\hat{h} \leftarrow h^*$ 

```

B Algorithm for Activation Replay

Line 1, 3, 5, 8 are performed by the original Mamba-2 block to get the output and update the states. The added operations are the extra *Conv1d_update* and *SSM_update* at line 2 and 5 to attain the updated states using the correct cached activations. Line 4 and 7 saves the activation of the current iteration for the use of next iteration.

Algorithm 2 Abstracted Mamba-2 block with Activation Replay. The two functions are simplified abstraction of what is actually used in the implementation.

Notation $A[a : b]$: a -th position to b -th position of A in the sequence length dimension

function $Conv1d_update(x, B, C, M)$:
return (x', B', C', M') \triangleright Performs 1d convolution on x, B, C , and update conv state M

function $SSM_update(x, B, C, dt, S)$:
return (y, S') \triangleright Performs SSM update with x, B, C, dt to get output y and update SSM state S

Input: E : token embedding, S : SSM state, M : Conv state, c : number of correct tokens
 $(x_{cache}^{pre_conv}, B_{cache}^{pre_conv}, C_{cache}^{pre_conv})$: Cached activation for Conv states from last forward pass
 $(x_{cache}^{post_conv}, B_{cache}^{post_conv}, dt_{cache})$: Cached activation for SSM states from last forward pass

Output: \hat{E} : New token embedding, \hat{S} : SSM state after the i -th correct token, \hat{M} : Conv state after the i -th correct token

- 1: $x^{pre_conv}, B^{pre_conv}, C^{pre_conv}, dt \leftarrow in_projection(E)$
- 2: $(_, _, _, \hat{M}) \leftarrow Conv1d_update(x_{cache}^{pre_conv}[0 : c], B_{cache}^{pre_conv}[0 : c], C_{cache}^{pre_conv}[0 : c], M)$ \triangleright Getting updated Conv state
- 3: $(x^{post_conv}, B^{post_conv}, C^{post_conv}, _) \leftarrow Conv1d_update(x^{pre_conv}, B^{pre_conv}, C^{pre_conv}, \hat{M})$ \triangleright Actual Conv1d update
- 4: $(x_{cache}^{pre_conv}, B_{cache}^{pre_conv}, C_{cache}^{pre_conv}) \leftarrow (x^{pre_conv}, B^{pre_conv}, C^{pre_conv})$ \triangleright Saving activations for next forward pass
- 5: $(_, \hat{S}) \leftarrow SSM_update(x_{cache}^{post_conv}[0 : c], B_{cache}^{post_conv}[0 : c], dt_{cache}[0 : c], S)$ \triangleright Getting the updated state using the correctly predicted tokens
- 6: $(y, _) \leftarrow SSM_update(x^{post_conv}, B^{post_conv}, C^{post_conv}, dt, \hat{S})$ \triangleright Actual SSM update
- 7: $(x_{cache}^{post_conv}, B_{cache}^{post_conv}, dt_{cache}) \leftarrow (x^{post_conv}, B^{post_conv}, dt)$ \triangleright Saving activations for next forward pass
- 8: $\hat{E} \leftarrow out_projection(y)$
