

(a)

initial hidden state : $\mathbf{h0} = \mathbf{0}$

initial cell state: $\mathbf{C0} = \mathbf{0}$

Step	Forget Gate	Input Gate	Cell State	Output Gate
1	$f1 = \sigma(Wf * [h0, x1] + bf)$	$i1 = \sigma(Wi * [h0, x1] + bi)$	$C1 = f1 * C0 + i1 * \tanh(Wc * [h0, x1] + bc)$	$o1 = \sigma(Wo * [h0, x1] + bo)$
2	$f2 = \sigma(Wf * [h1, x2] + bf)$	$i2 = \sigma(Wi * [h1, x2] + bi)$	$C2 = f2 * C1 + i2 * \tanh(Wc * [h1, x2] + bc)$	$o2 = \sigma(Wo * [h1, x2] + bo)$
3	$f3 = \sigma(Wf * [h2, x3] + bf)$	$i3 = \sigma(Wi * [h2, x3] + bi)$	$C3 = f3 * C2 + i3 * \tanh(Wc * [h2, x3] + bc)$	$o3 = \sigma(Wo * [h2, x3] + bo)$
4	$f4 = \sigma(Wf * [h3, x4] + bf)$	$i4 = \sigma(Wi * [h3, x4] + bi)$	$C4 = f4 * C3 + i4 * \tanh(Wc * [h3, x4] + bc)$	$o4 = \sigma(Wo * [h3, x4] + bo)$
5	$f5 = \sigma(Wf * [h4, x5] + bf)$	$i5 = \sigma(Wi * [h4, x5] + bi)$	$C5 = f5 * C4 + i5 * \tanh(Wc * [h4, x5] + bc)$	$o5 = \sigma(Wo * [h4, x5] + bo)$

(b)

LSTMs address the vanishing gradient problem by:

- Introducing a cell state that can store information over long sequences.
- Utilizing gates to selectively control the flow of information.
- Employing additive operations in the cell state update, which helps to preserve the gradient signal.

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

When calculating the gradient of C_t with respect to C_{t-1} (which is essential for backpropagation), the gradient of C_{t-1} will be directly added to the overall gradient, rather than being multiplied. This additive nature helps to preserve the gradient signal and prevent it from vanishing, even over long sequences.

(c)

- The dimension of each weight matrix: $W_x = 100 \times 100$, $W_h = 64 \times 64$
- The total number of parameters : 64 biases for each of 4 gates: $64 \times 4 = 256$, and 2 matrices: $100 \times 64 + 64 \times 64 = 6400 + 4096$. Total: **10752**
- The memory requirements during training: assuming number is a float type = 4byte.

Then: memory for parameters: $10752 \times 4 = \mathbf{43008B}$

For each sample in a batch + gradients: hidden state + cell state + input: $64B + 64B + 100B + 64B + 64B + 100B = 396B$

For entire batch: $396B \times 32 \text{ samples} = \mathbf{12672B}$

Total memory: $\mathbf{43008B + 12672B = 55680B = 54.37K}$