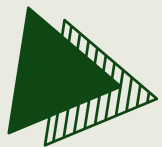# SEQUENCE ANALYSIS WITH MIN RNN
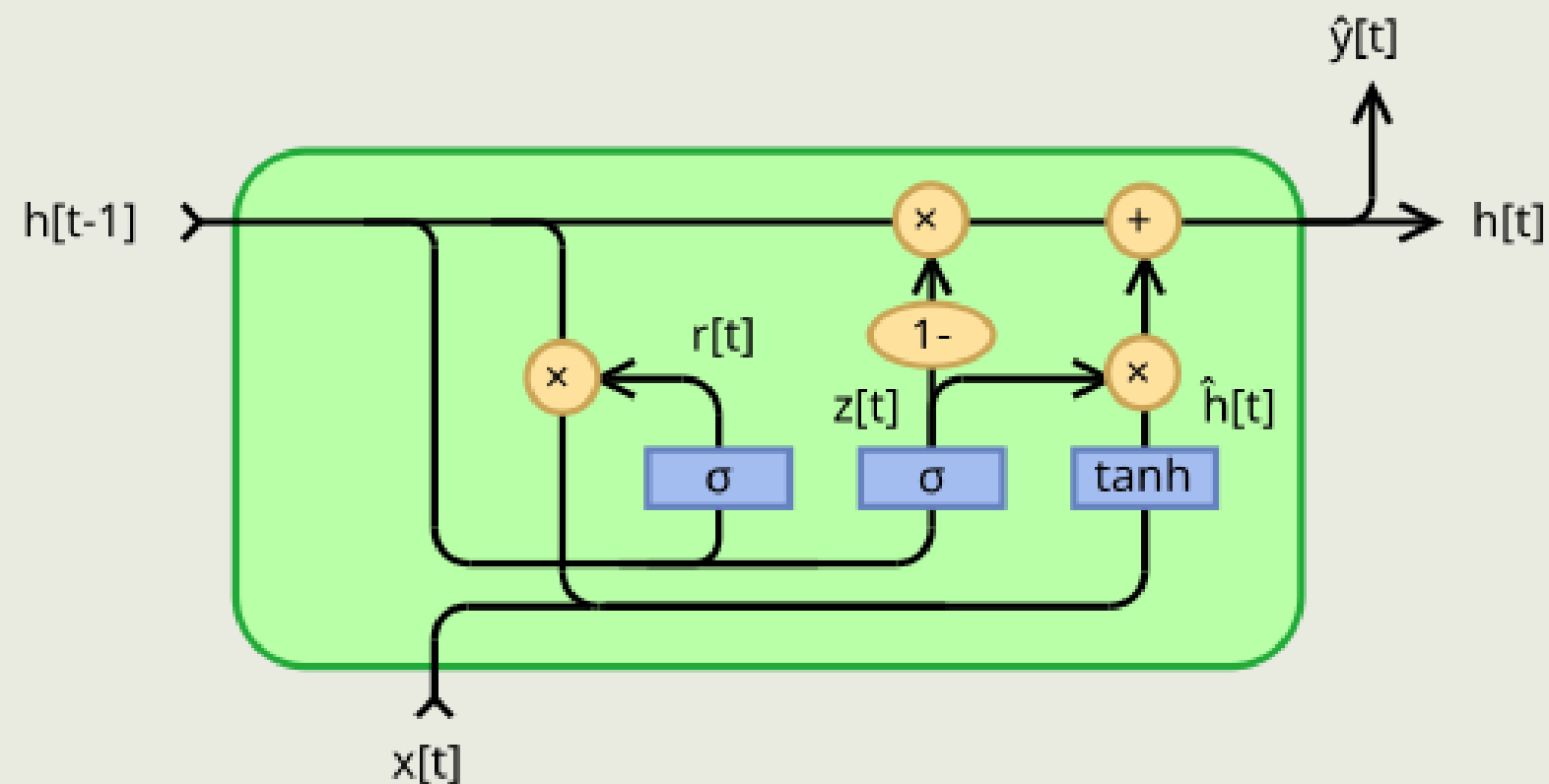
BY Mohan Sharma

# GRU NETWORKS
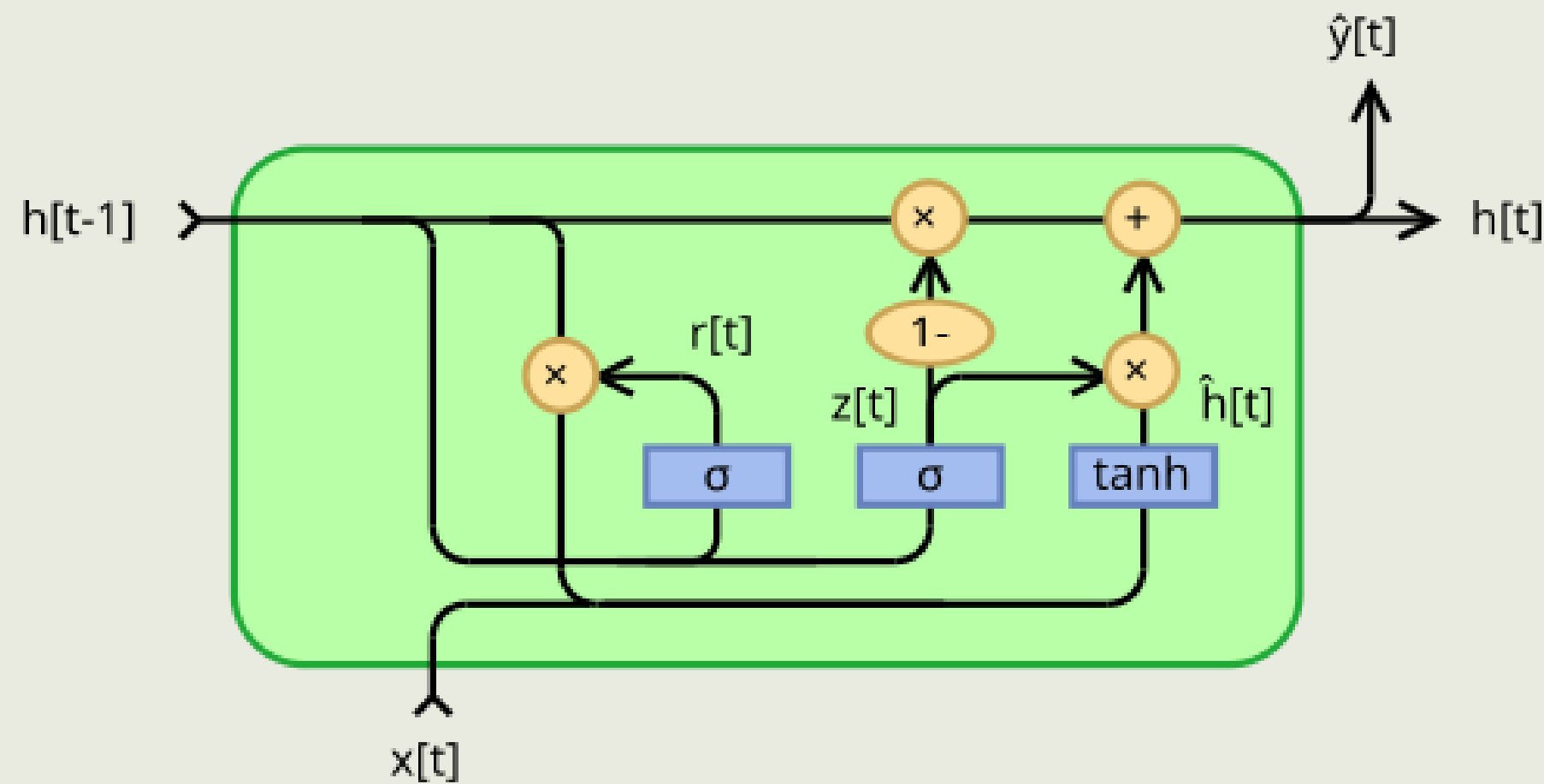
# STEP-BY-STEP GRU WALK THROUGH

$\hat{y}[t]$

$h[t-1]$

$h[t]$

$r[t]$

$z[t]$

$\hat{h}[t]$

σ   σ   tanh

$x[t]$

r_t = sigmoid(W_r * [h_t-1, x_t])

z_t = sigmoid(W_z * [h_t-1, x_t])Gate

1. The first step in our LSTM is to decide what informatioThe reset gate r and update gate z are computed using the current input x and the previous hidden state h_t-1

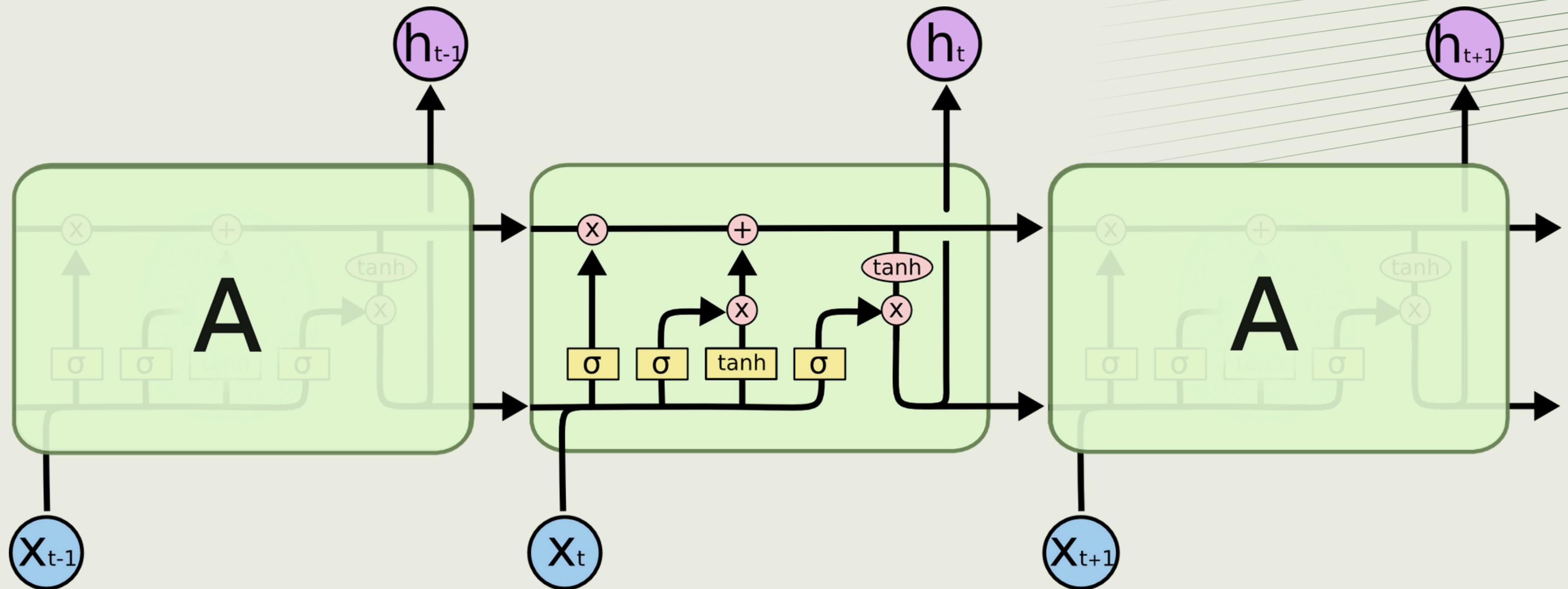n we're going to throw away from the cell state. This decision is made by a sigmoid layer called the "forget gate layer."

2) The candidate activation vector h_t~ is computed using the current input x and a modified version of the previous hidden state that is "reset" by the reset gate

$$h\_t\text{~} = tanh(W\_h * [r\_t * h\_t\text{-}1, x\_t])$$

3) The new hidden state h_t is computed by combining the candidate activation vector with the previous hidden state, weighted by the update gate

$$h\_t = (1 - z\_t) * h\_t\text{-}1 + z\_t * h\_t\text{~}$$

# LSTM NETWORKS

# LIMITATIONS

- Computational Cost and time is high
- Parallelization is not possible
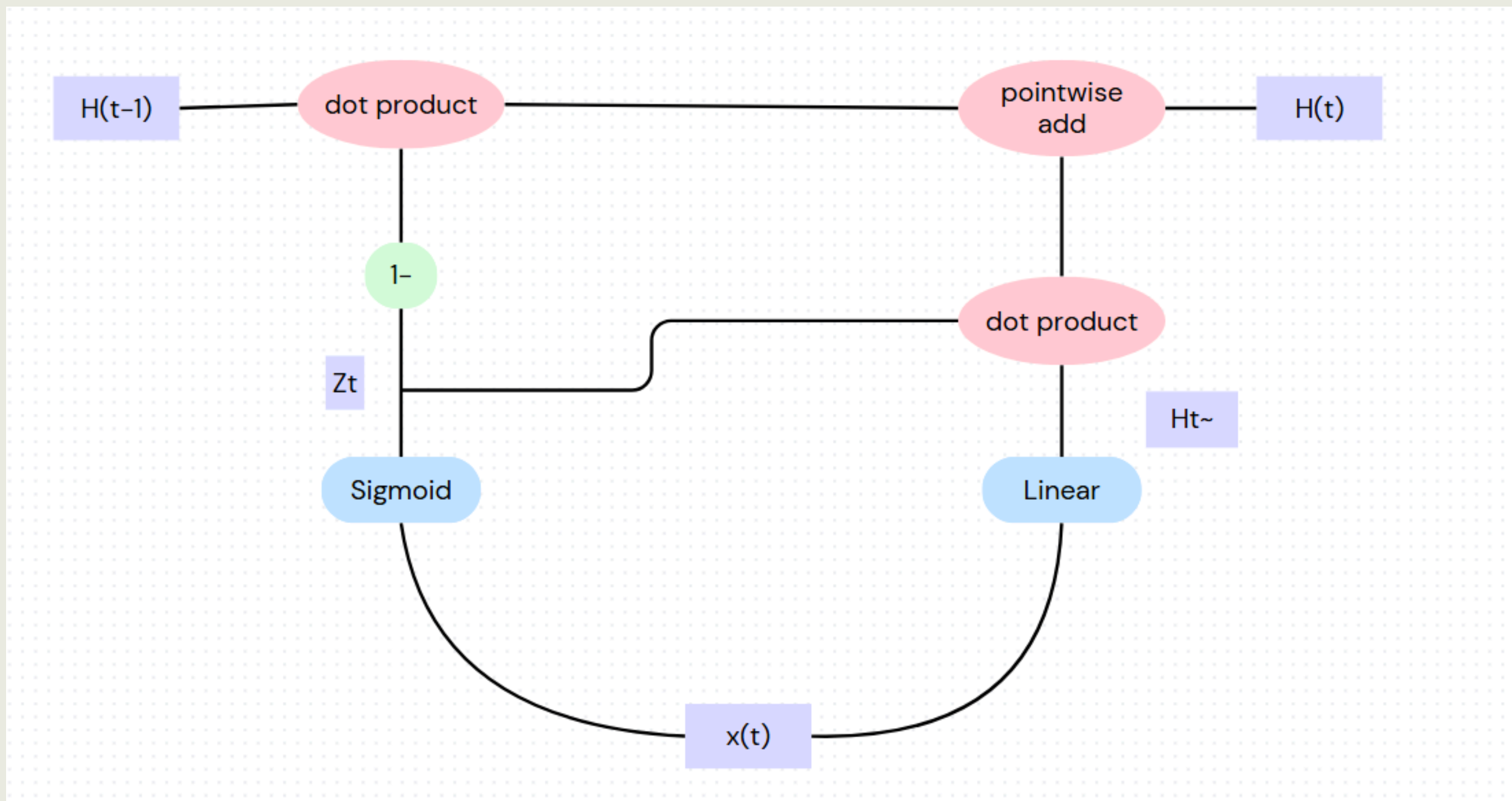- Performance on Long Sequences is not so good
- Memory Inefficiency

# MIN GRU



GRU

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$
$$z_t = \sigma(\text{Linear}_{d_h}([x_t, h_{t-1}]))$$
$$r_t = \sigma(\text{Linear}_{d_h}([x_t, h_{t-1}]))$$
$$\tilde{h}_t = \tanh(\text{Linear}_{d_h}([x_t, r_t \odot h_{t-1}]))$$

$\Rightarrow$

minGRU

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$
$$z_t = \sigma(\text{Linear}_{d_h}(x_t))$$
$$\tilde{h}_t = \text{Linear}_{d_h}(x_t)$$

# MIN LSTM

To remove the previous dependencies by decomposing to
Model to bare Bones



- Remove cell state
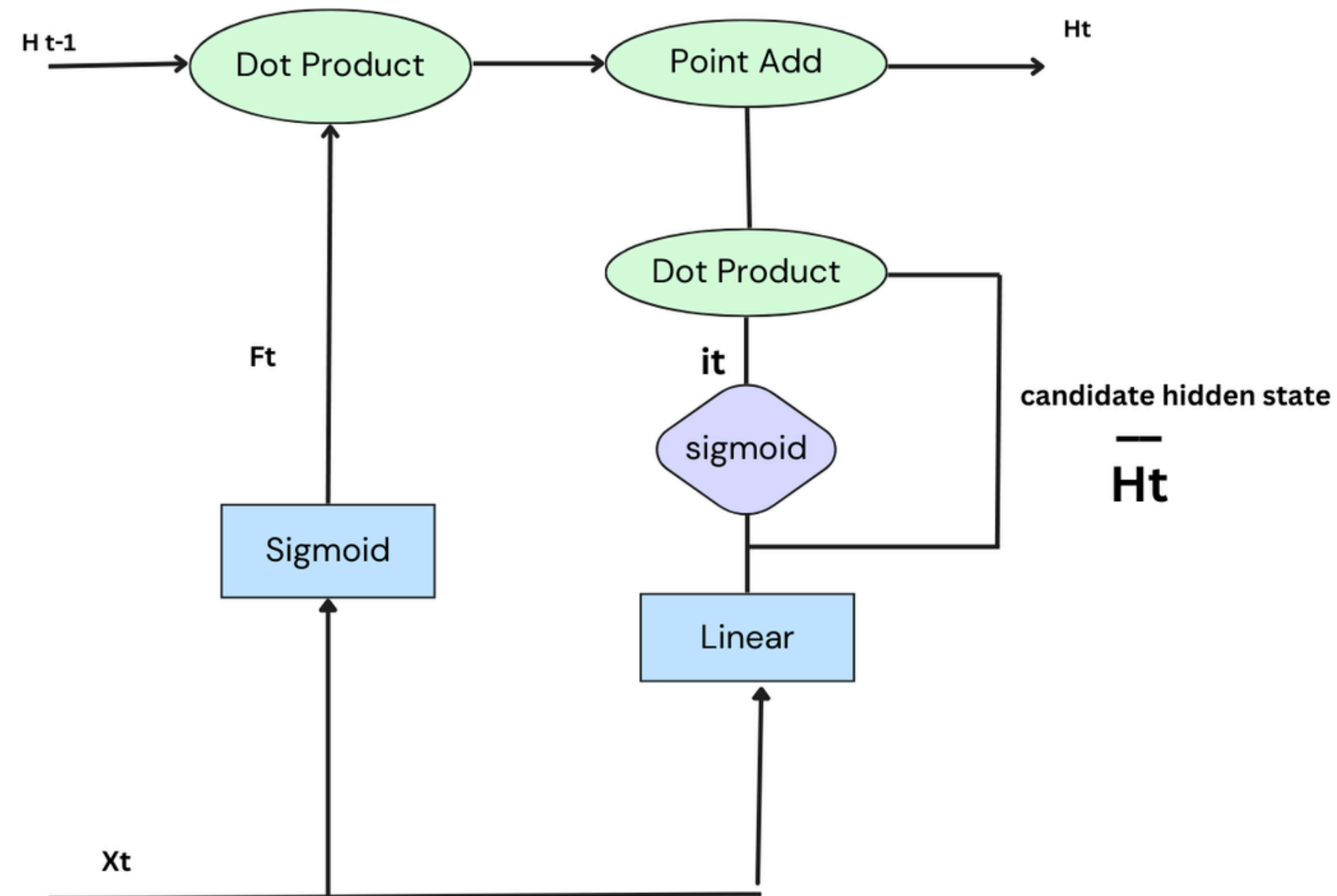- Remove previous dependencies on previous hidden and cell state

HOW TO COMPUTE THIS IN PARALLEL??
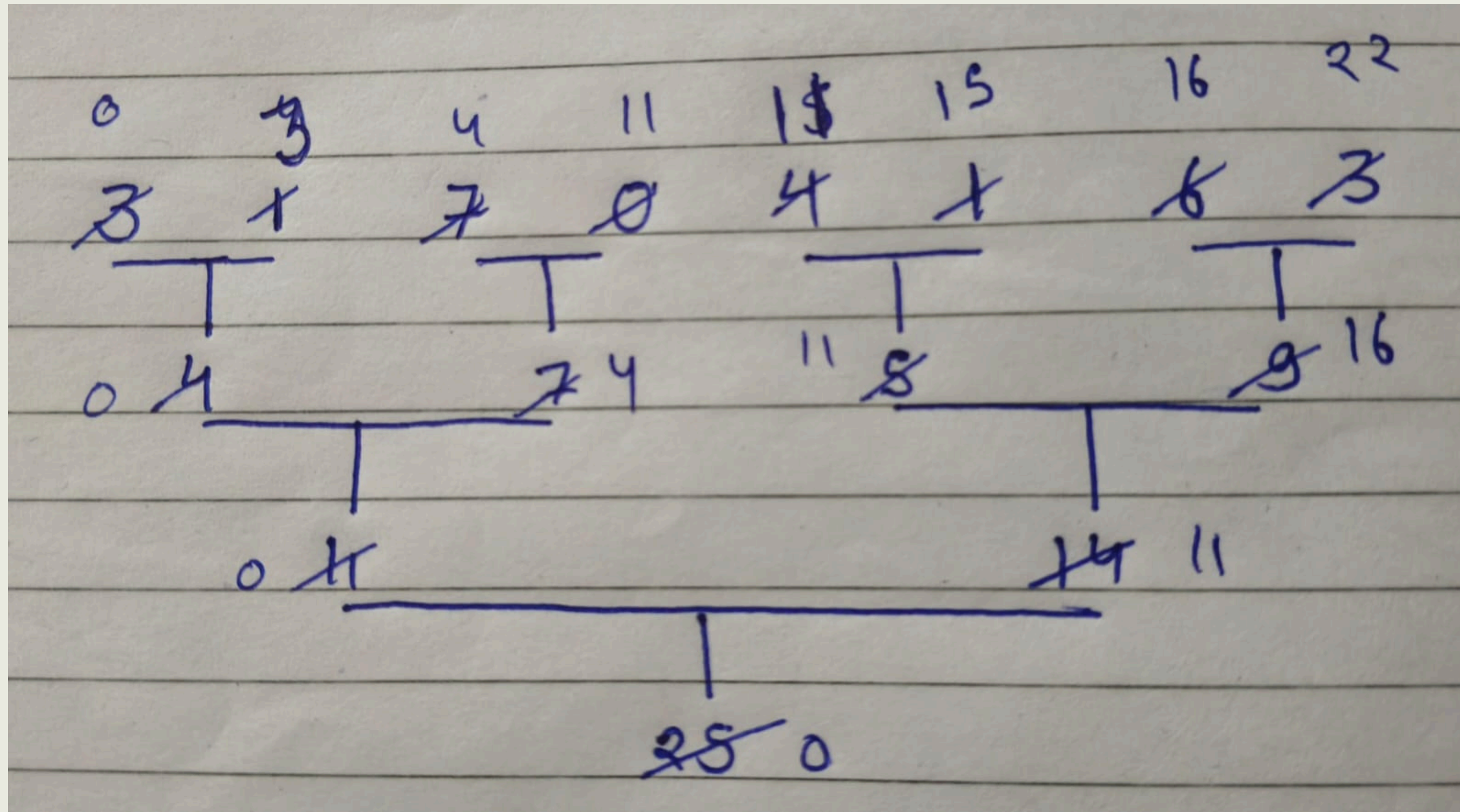
# PARALLEL SCAN ALGORITHM

1. UP-SWEEP PHASE (REDUCE)

BUILD A BINARY TREE FROM THE ARRAY BOTTOM-UP BY SUMMING PAIRS OF ELEMENTS:
- AT EACH LEVEL, EVERY PROCESSOR WORKS ON 2 ELEMENTS AND STORES THE SUM IN THE PARENT NODE.
- THE LAST ELEMENT WILL HAVE THE TOTAL SUM OF THE ARRAY.

2. DOWN-SWEEP PHASE
- SET THE ROOT NODE TO 0.
- TRAVERSE THE TREE BACK DOWN:
  - SWAP VALUES AND COMPUTE PREFIX SUMS USING VALUES FROM THE UP-SWEEP PHASE.

0    3    4    11    1$    15    16    22

3    1    7    0    4    1    6    3

0 4      7 4     11 8     8 16

0 11              14 11

25 0

**INPUT - 3 1 7 0 4 1 6 3**

**INPUT - 0 3 4 11 11 12 16 22**

(exclusive prefix sum)

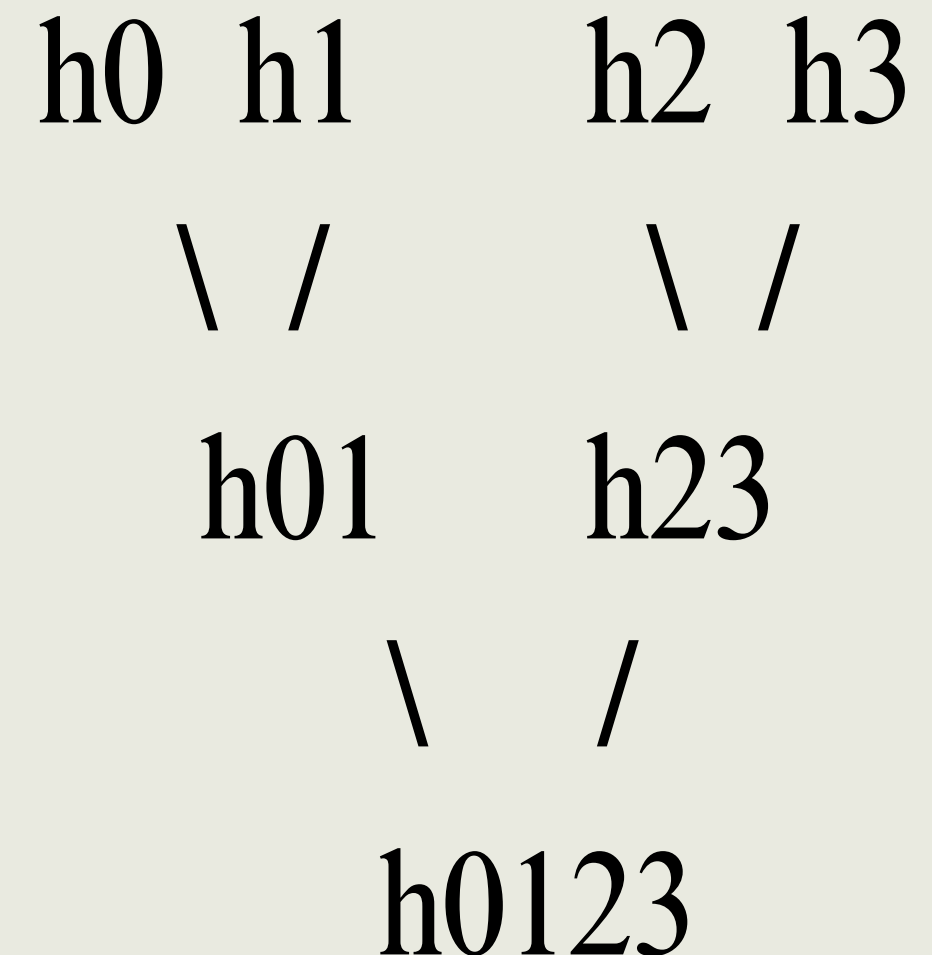# How Is the parallel sum algo is applied to min GRU and LSTM?

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$
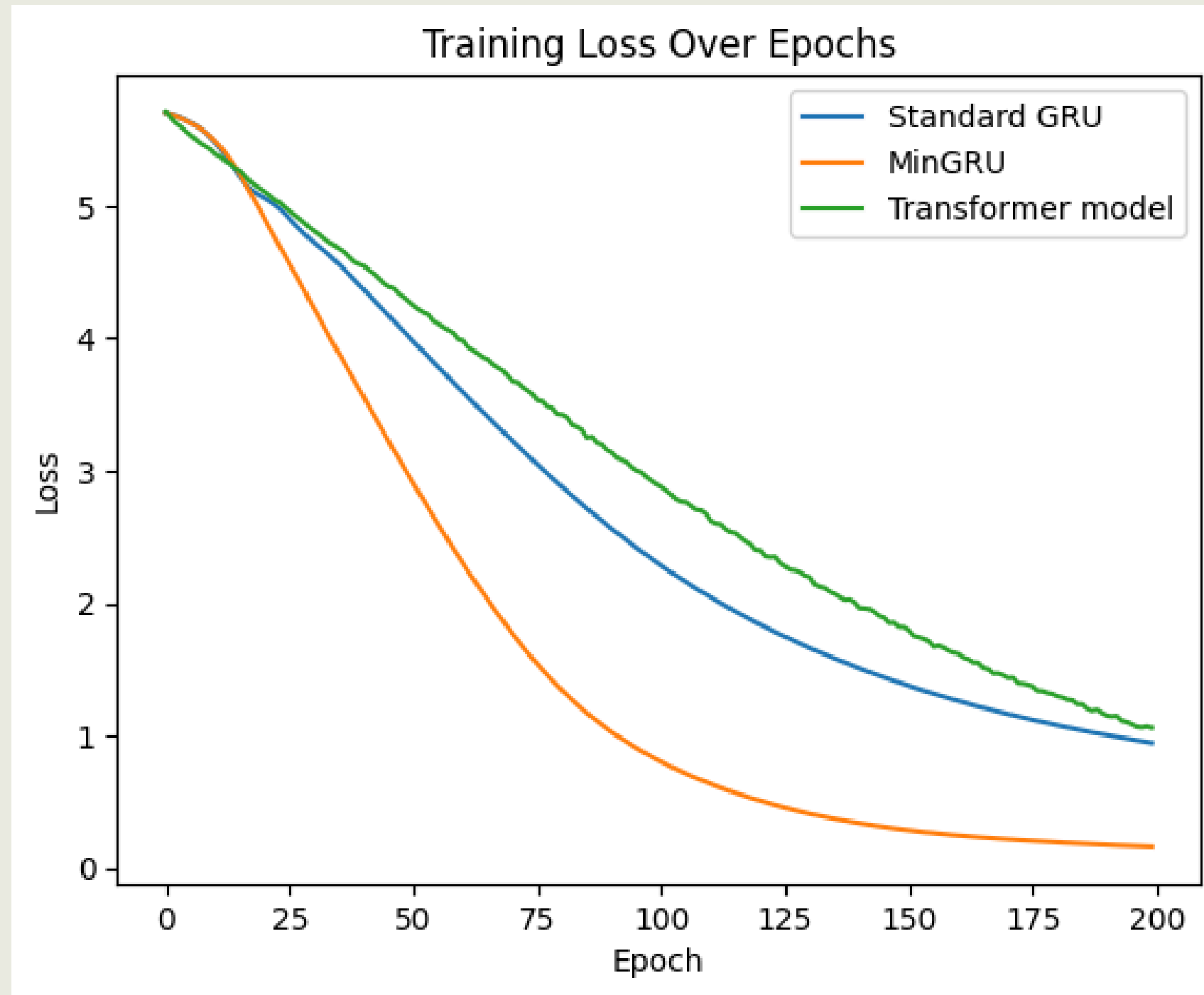
This is similar to Parallel scan

As

$$h_t = a_t \cdot h_{t-1} + b_t$$

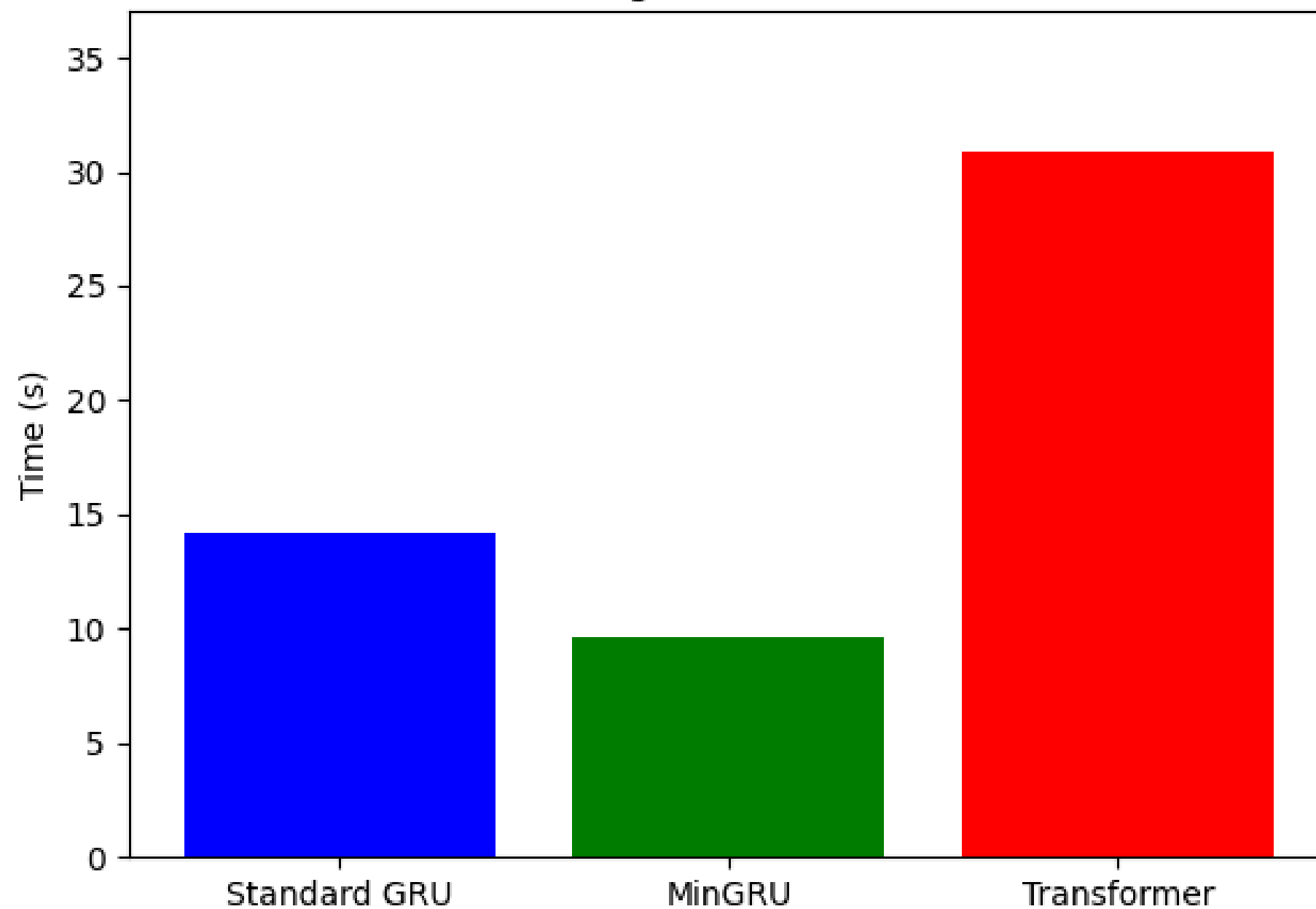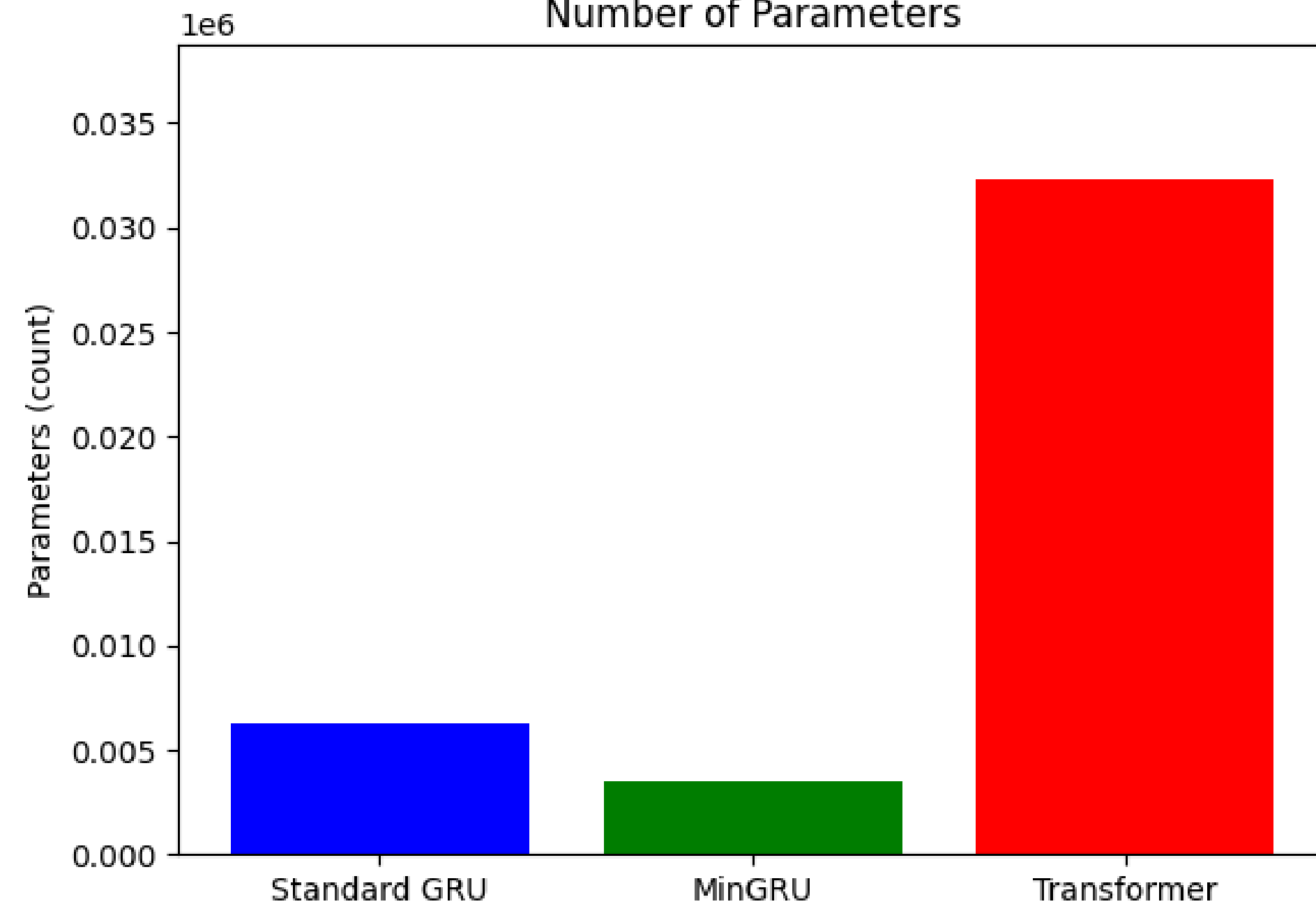where a & b are independent on h(t-1)

```
h0  h1      h2  h3
 \ /         \ /
 h01        h23
    \      /
      h0123
```

# TRAINING LOSS



## Training Loss Over Epochs

# GRU & MIN GRU &TRANSFORMER MODEL

Model Comparison: Training Time and Parameters

*Thank you.*