

## xLSTM

Extended Long short term memory by Sepp Hochreiter.



KingGongzilla · 14h ago

I just recently attended a lecture by Sepp Hochreiter where he compares Mamba to his new (yet to be released) xLSTM. Apparently xLSTM outperforms Mamba and transformers. Hope its true 🤞

(-)    14    Reply       ...

## What is RNN?

A Recurrent Neural Network is a type of neural network that is designed for sequential data. RNN have connections that form directed cycles allowing them to maintain a memory of previous input & utilize this information when processing new inputs.

## What is LSTM?

LSTM is a special type of RNN designed to overcome the vanishing gradient problem and effectively capture long-term dependencies in sequential data.

memory }  
cell      LSTM contain memory cells that can store  
            information for long period.

Gating Mechanisms ↴ LSTM uses three main gates to control the flow of information into and out of the memory cells.

1. Forget Gate:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

This gate decides what fraction of the previous cell state to forget.

2. Input Gate:

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

This gate determines which new information to add to the cell state.

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

This is the candidate cell state, which contains the new information.

3. Update Cell State:

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

The cell state is updated by combining the old cell state (partially forgotten) and the new candidate cell state.

4. Output Gate:

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

This gate determines the output of the LSTM cell.

$$h_t = o_t \cdot \tanh(C_t)$$

The hidden state is updated based on the cell state and the output gate.

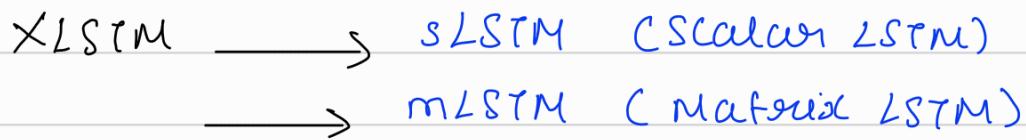
What is GRU?

The Gated Recurrent Unit, are similar to LSTM network but a simpler structure. This simplicity can lead to faster training times & fewer computational resources while still maintaining performances comparable to LSTM.

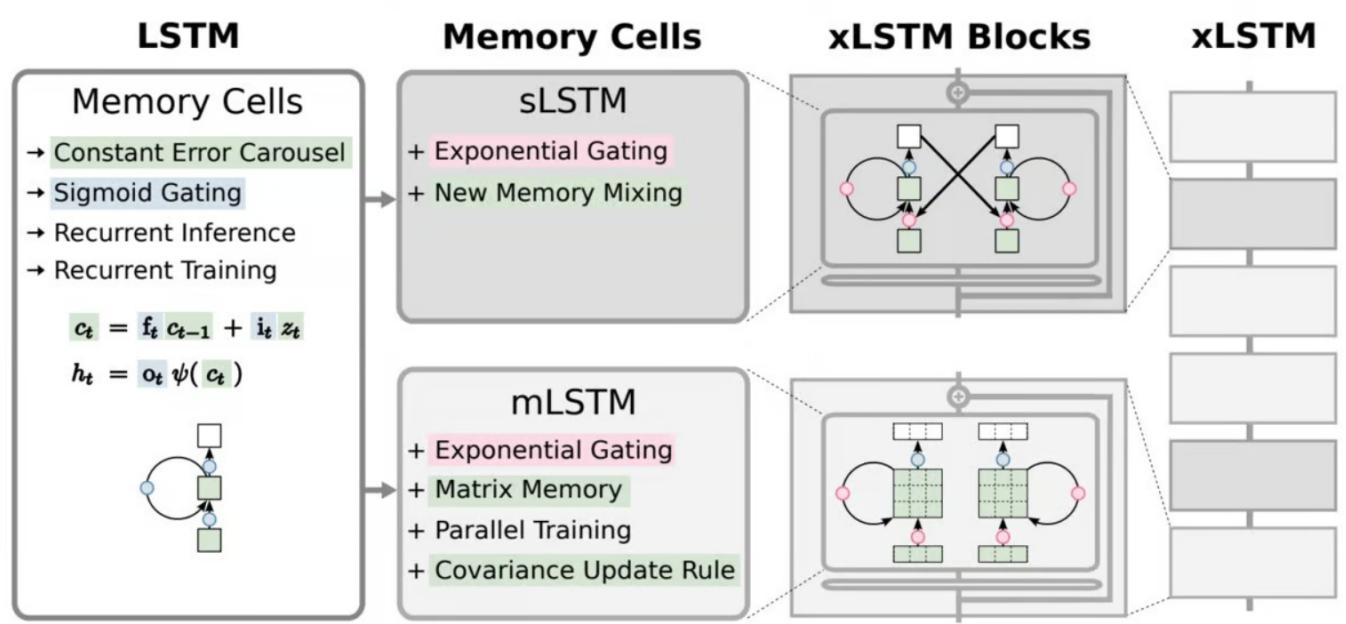
\* NO separate cell state.

\* Simplified Gating Mechanism (Reset & update gate)

NOW, it's back



sLSTM & mLSTM is optionally joined together to create  $\times$ LSTM block.



Due to nature of LSTM, the problem can't be solved using one architecture

In LSTM, each hidden block relies on the output of previous block (sequential manner) it's called memory leaking which is inefficient.

While, transformer quickly get the entire output with a matrix multiplication.

LSTM is not a Scalable to sequential processing.

One of the key idea of LSTM is Gating, which was rediscovered & reinterpreted in Gated state space model and memory.

All Gated activation functions are sigmoid.

$$\sigma(x) = 1 / (1 + e^{-x})$$

earlier, LSTM incorporated constant error carousel along with structured gating mechanisms for vanishing gradient problem in RNN.

### Limitation of LSTM

- \* Inability to revise storage decision, Once information is stored, traditional LSTMs have limited ability to revise or overwrite it based on new inputs.
- \* Limited storage capacity, The scalar nature of LSTM memory cells restricts the volume and complexity of information that can be maintained.
- \* Lack of parallelizability, The inherent sequential processing between time steps prevents LSTMs from leveraging modern parallel computing architecture effectively.

# Working XLSTM

\* Exponential Gating, XLSTM introduces exponential gating with appropriate normalization and stabilization techniques which enhances the gating mechanism to address inability of revising storage deletion.

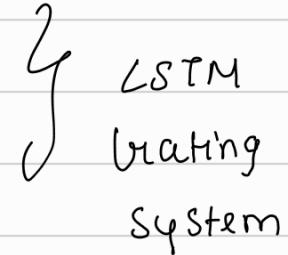
The gates (input & forget) may use exponential activations, adjusted by a stabilizing factors to maintain numerical stability.

$$\begin{aligned}
 c_t &= f_t c_{t-1} + i_t z_t \\
 h_t &= o_t \tilde{h}_t, \\
 z_t &= \varphi(\tilde{z}_t), \\
 i_t &= \sigma(\tilde{i}_t), \\
 f_t &= \sigma(\tilde{f}_t), \\
 o_t &= \sigma(\tilde{o}_t),
 \end{aligned}$$

Sigmoid

$$\begin{aligned}
 \text{cell state} &\quad (2) \\
 \tilde{h}_t &= \psi(c_t) \\
 \tilde{z}_t &= w_z^\top x_t + r_z h_{t-1} + b_z \\
 \tilde{i}_t &= w_i^\top x_t + r_i h_{t-1} + b_i \\
 \tilde{f}_t &= w_f^\top x_t + r_f h_{t-1} + b_f \\
 \tilde{o}_t &= w_o^\top x_t + r_o h_{t-1} + b_o
 \end{aligned}$$

$$\begin{aligned}
 \text{hidden state} &\quad (3) \\
 \text{cell input} &\quad (4) \\
 \text{input gate} &\quad (5) \\
 \text{forget gate} &\quad (6) \\
 \text{output gate} &\quad (7)
 \end{aligned}$$



$$C_t = f_t C_{t-1} + i_t v_t k_t^\top$$

$$n_t = f_t n_{t-1} + i_t k_t$$

$$h_t = o_t \odot \tilde{h}_t, \quad \tilde{h}_t = C_t q_t / \max \left\{ \| n_t^\top q_t \|, 1 \right\}$$

$$q_t = W_q x_t + b_q$$

$$k_t = \frac{1}{\sqrt{d}} W_k x_t + b_k$$

$$v_t = W_v x_t + b_v$$

$$i_t = \exp(\tilde{i}_t), \quad \tilde{i}_t = w_i^\top x_t + b_i$$

$$f_t = \sigma(\tilde{f}_t) \text{ OR } \exp(\tilde{f}_t), \quad \tilde{f}_t = w_f^\top x_t + b_f$$

$$o_t = \sigma(\tilde{o}_t), \quad \tilde{o}_t = W_o x_t + b_o$$

$$\text{cell state (19)}$$

$$\text{normalizer state (20)}$$

$$\text{hidden state (21)}$$

$$\text{query input (22)}$$

$$\text{key input (23)}$$

$$\text{value input (24)}$$

$$\text{input gate (25)}$$

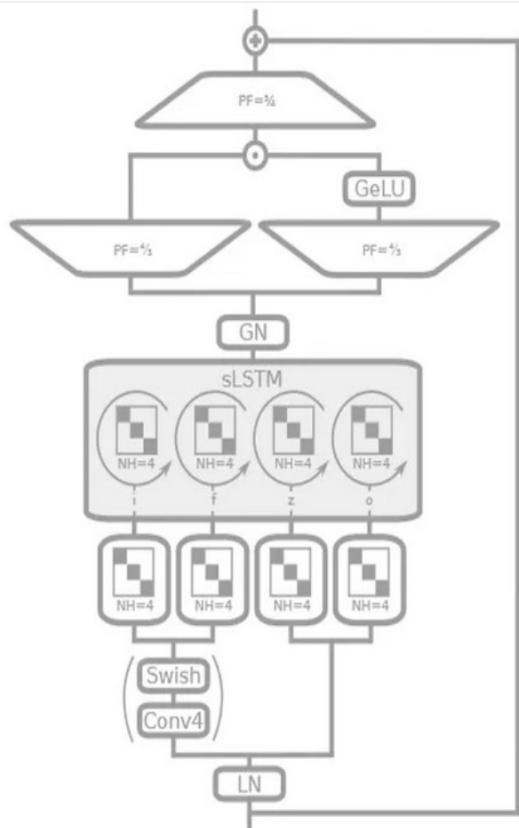
$$\text{forget gate (26)}$$

$$\text{output gate (27)}$$

Exponential function.

\* Modified memory structure,

(i) SLSTM : Scalar LSTM with Exponential Gating & Memory Mixing



Enhanced version of LSTM with scalar or sequences level updates.

includes improvement to the gating mechanisms (such as exponential gating) and optimization of memory structure.

SLSTM, utilizes the post up projections.

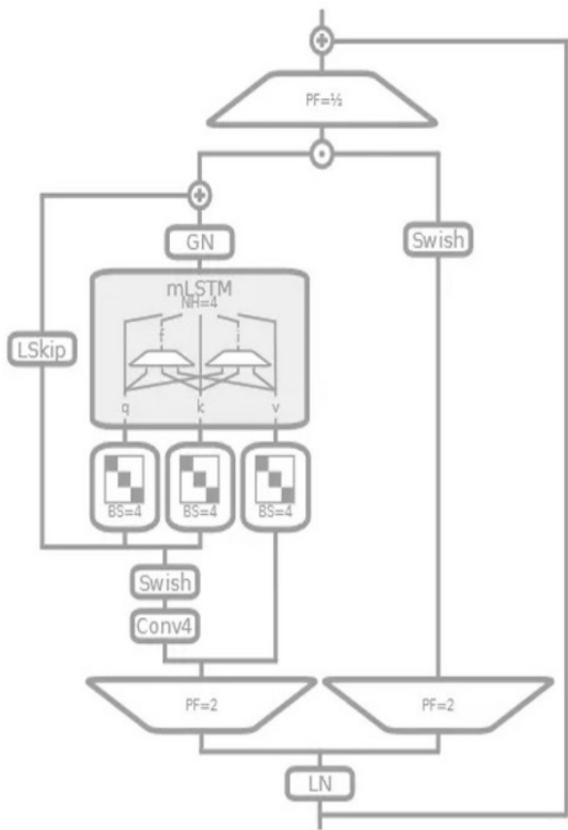
\* Exponential Gating,

SLSTM incorporates exponential activation functions for the Input & forget gates. Providing more flexible control over information flow.

\* Normalization & Stabilization, To prevent numerical instabilities, SLSTM introduces a normalizer state that tracks the product of input & future forget gates.

\* Memory Mixing, SLSTM supports multiple memory cell and allows for memory mixing via recurrent connections. Enhances ability to extract complex pattern.

### ciii MLSTM: Matrix LSTM with Enhanced Storage Capacities



LSTM variant featuring matrix memory, which allows it to process and store more information in parallel.

It uses matrices instead of scalars to store the LSTM cell states.

MLSTM utilizes prep-up projections.

- \* **Matrix memory**, MLSTM utilizes a matrix memory instead of a scalar memory cell, increasing storage capacity and enabling more efficient information retrieval.
- \* **Covariance update rule**; MLSTM employs a covariance update rule to store and retrieve key-value pairs efficiently.
- \* **parallelizability**, By eliminating memory mixing, MLSTM achieves full parallelizability allowing for efficient computation on modern hardware accelerators.

These variants SLSTM & MLSLM can be integrated into residual block architectures forming xLSTM Blocks.

	Context Sensitive		Deterministic Context Free		Regular				Majority	
	Bucket Sort	Missing Duplicate	Mod Arithmetic (w/o Brackets)	Solve Equation	Cycle Nav	Even Pairs	Mod Arithmetic (w/o Brackets)	Parity	Majority	Majority Count
Llama	0.92 ± 0.02	0.08 ± 0.0	0.02 ± 0.0	0.02 ± 0.0	0.04 ± 0.01	1.0 ± 0.0	0.03 ± 0.0	0.03 ± 0.01	0.37 ± 0.01	0.13 ± 0.0
Mamba	0.69 ± 0.0	0.15 ± 0.0	0.04 ± 0.01	0.05 ± 0.02	0.86 ± 0.04	1.0 ± 0.0	0.05 ± 0.02	0.13 ± 0.02	0.69 ± 0.01	0.45 ± 0.03
Retention	0.13 ± 0.01	0.03 ± 0.0	0.03 ± 0.0	0.03 ± 0.0	0.05 ± 0.01	0.51 ± 0.07	0.04 ± 0.0	0.05 ± 0.01	0.36 ± 0.0	0.12 ± 0.01
Hyena	0.3 ± 0.02	0.06 ± 0.02	0.05 ± 0.0	0.02 ± 0.0	0.06 ± 0.01	0.93 ± 0.07	0.04 ± 0.0	0.04 ± 0.0	0.36 ± 0.01	0.18 ± 0.02
RWKV-4	0.54 ± 0.0	0.21 ± 0.01	0.06 ± 0.0	0.07 ± 0.0	0.13 ± 0.0	1.0 ± 0.0	0.07 ± 0.0	0.06 ± 0.0	0.63 ± 0.0	0.13 ± 0.0
RWKV-5	0.49 ± 0.04	0.15 ± 0.01	0.08 ± 0.0	0.08 ± 0.0	0.26 ± 0.05	1.0 ± 0.0	0.15 ± 0.02	0.06 ± 0.03	0.73 ± 0.01	0.34 ± 0.03
RWKV-6	0.96 ± 0.0	0.23 ± 0.06	0.09 ± 0.01	0.09 ± 0.02	0.31 ± 0.14	1.0 ± 0.0	0.16 ± 0.0	0.22 ± 0.12	0.76 ± 0.01	0.24 ± 0.01
LSTM (Block)	0.99 ± 0.0	0.15 ± 0.0	0.76 ± 0.0	0.5 ± 0.05	0.97 ± 0.03	1.0 ± 0.0	0.91 ± 0.09	1.0 ± 0.0	0.58 ± 0.02	0.27 ± 0.0
LSTM	0.94 ± 0.01	0.2 ± 0.0	0.72 ± 0.04	0.38 ± 0.05	0.93 ± 0.07	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.82 ± 0.02	0.33 ± 0.0
xLSTM[0:1]	0.84 ± 0.06	0.23 ± 0.01	0.57 ± 0.09	0.55 ± 0.09	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.75 ± 0.02	0.22 ± 0.0
xLSTM[1:0]	0.97 ± 0.0	0.33 ± 0.22	0.03 ± 0.0	0.03 ± 0.01	0.86 ± 0.01	1.0 ± 0.0	0.04 ± 0.0	0.04 ± 0.01	0.74 ± 0.01	0.46 ± 0.0
xLSTM[1:1]	0.7 ± 0.21	0.2 ± 0.01	0.15 ± 0.06	0.24 ± 0.04	0.8 ± 0.03	1.0 ± 0.0	0.6 ± 0.4	1.0 ± 0.0	0.64 ± 0.04	0.5 ± 0.0

Figure 4: Test of xLSTM's exponential gating with memory mixing. Results are given by the scaled accuracy of different models at solving formal language tasks, of which some require state tracking. The different tasks are grouped by the Chomsky hierarchy.

Comparative Study

## Application of xLSTM

1. Language Modelling.
2. Time Series Analysis
3. Machine Translation.
4. Speech Recognition & Generation.