

Time Series Analysis - Assignment IX

Chitra Khatri
Aniket Nighot
Hardik Kunt

February 8, 2025

What is xLstm?

xLSTM (Extended Long Short-Term Memory) is an advanced neural network architecture designed to address the limitations of traditional LSTM networks. XLSTM enhances the capability to model complex time-series data by incorporating various improvements in stability, efficiency, and the ability to capture intricate dependencies in multivariate datasets.

How XLSTM Came into Picture

The development of XLSTM was motivated by the need to overcome the inherent drawbacks of traditional LSTMs. LSTMs, although effective for sequential data, faced significant challenges in handling long-term dependencies, maintaining stable gradients during training, and modeling complex relationships in multivariate time series data. XLSTM was introduced as a solution to these challenges, integrating various enhancements to improve performance and stability.

Architecture of xLstm

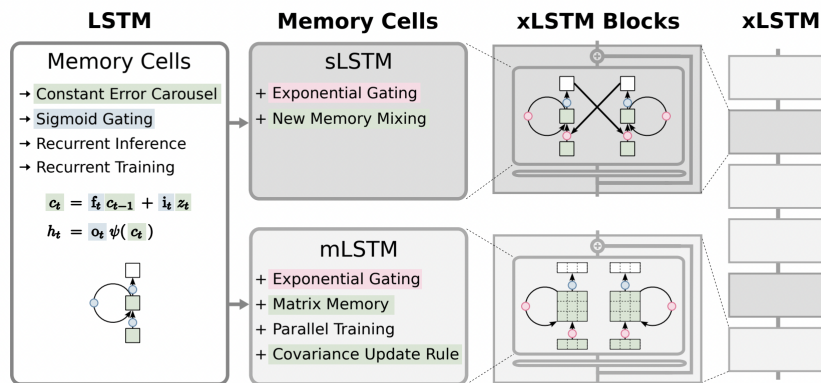


Figure 1: Architecture of XLSTM

Here's a breakdown of its components and how each part contributes to its improved functionality:

Key Changes in mLSTM and sLSTM in xLSTM Model

mLSTM

- **Matrix Memory:** Instead of storing information as single values (scalars) like in traditional LSTMs, mLSTM uses matrices (tables of numbers) to store information. This allows the model to remember and manage more complex patterns and relationships in the data.

- **Covariance Update Rule:** mLSTM updates its memory using a mathematical technique called the covariance update rule. This involves storing pairs of related information (like a key and a value) and helps the model efficiently recall stored information later.
- **Parallel Processing:** Traditional LSTMs process information step-by-step in sequence, but mLSTM removes the direct link between consecutive steps. This allows it to process multiple steps at the same time (parallel processing), which makes the model faster and more scalable.
- **Stabilization Techniques:** Because mLSTM uses exponential functions (which can produce very large numbers), it includes special techniques to keep these numbers stable and prevent errors from growing too large.
- **Pre Up-Projection:** Before doing complex calculations, mLSTM first expands the input into a higher-dimensional space. This helps the model handle and store more detailed information, improving its memory capacity.

sLSTM (Scalar LSTM)

- **Exponential Gates:** sLSTM replaces the standard sigmoid gates with exponential gates. These gates allow more flexible control over what information is stored or forgotten, making the model more adaptable and effective.
- **Normalizer State:** To keep everything balanced and stable, sLSTM introduces a normalizer state. This keeps track of the cumulative effect of inputs and forget gates over time, helping the model maintain steady information flow.
- **Stabilizer State:** Like mLSTM, sLSTM also uses stabilization techniques to manage the large numbers that can result from exponential functions. This ensures the model remains stable and performs well over long sequences.
- **Multiple Memory Cells and Heads:** sLSTM can manage several memory cells and heads, which allows it to mix different types of information within each head. This improves the model's ability to handle complex tasks but avoids mixing information between different heads to maintain clarity and stability.

xLSTM Block Composition

xLSTM blocks are constructed by integrating sLSTM and mLSTM blocks within residual structures. Two main types of blocks are used:

- **Residual Block with Post Up-Projection:** This type, similar to Transformer blocks, is used for sLSTM. It first processes the input in the original space, then projects it into a higher-dimensional space, applies non-linear activation, and projects it back.
- **Residual Block with Pre Up-Projection:** Used for mLSTM, this block projects the input to a higher-dimensional space first, performs non-linear operations there, and then maps it back to the original space. This structure maximizes the memory capacity in high-dimensional space.

Advantages of xLSTM (Extended Long Short-Term Memory)

The following are the advantages of xLSTM as mentioned in the paper:

- **Improved Storage and Revision:** xLSTM introduces exponential gating, allowing it to revise stored values efficiently. This overcomes the limitations of traditional LSTMs that struggle with revising stored information.
- **Better Memory Capacity:** The mLSTM variant includes a matrix memory, significantly improving the model’s capacity to store and retrieve complex information. This feature is particularly useful for tasks like predicting rare tokens or handling large vocabularies.
- **Parallelization:** mLSTM removes the sequential dependency in LSTM’s memory mixing, enabling parallel processing. This is critical for scaling up and speeding up computations.
- **Scalability:** The architecture supports stacking residual blocks, allowing xLSTM to be scaled to handle large and complex language modeling tasks, comparable to state-of-the-art methods like Transformers.
- **Efficient for Long Sequences:** With constant memory complexity concerning sequence length, xLSTM is efficient for processing long sequences, making it suitable for real-world applications with extensive data.
- **Enhanced Performance:** xLSTM performs better than traditional LSTMs and competes favorably with Transformers and State Space Models in tasks like language modeling and long-context processing.

Evaluation

The parity task is particularly noteworthy as Transformers and State Space Models often struggle without memory mixing or state tracking. In contrast, xLSTM, built on recurrent neural network

| | Context Sensitive | | Deterministic Context Free | | Regular | | | | | |
|--------------|-------------------|-------------------|-----------------------------|----------------|----------------|----------------|-------------------------------|----------------|----------------|----------------|
| | Bucket Sort | Missing Duplicate | Mod Arithmetic (w 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 ± 0.08 | 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 2: Comparison of Models Accuracy on Formal Language Tasks

architectures, excels in this area. The xLSTM model achieves near-perfect accuracy (as high as 1), whereas Transformer-based models like Llama face significant challenges. [1]

References

- [1] Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xlstm: Extended long short-term memory, 2024.