

# Sample-Efficient Language Modeling with Linear Attention and Lightweight Enhancements

Patrick Haller, Jonas Golde, Alan Akbik  
Humboldt Universität zu Berlin

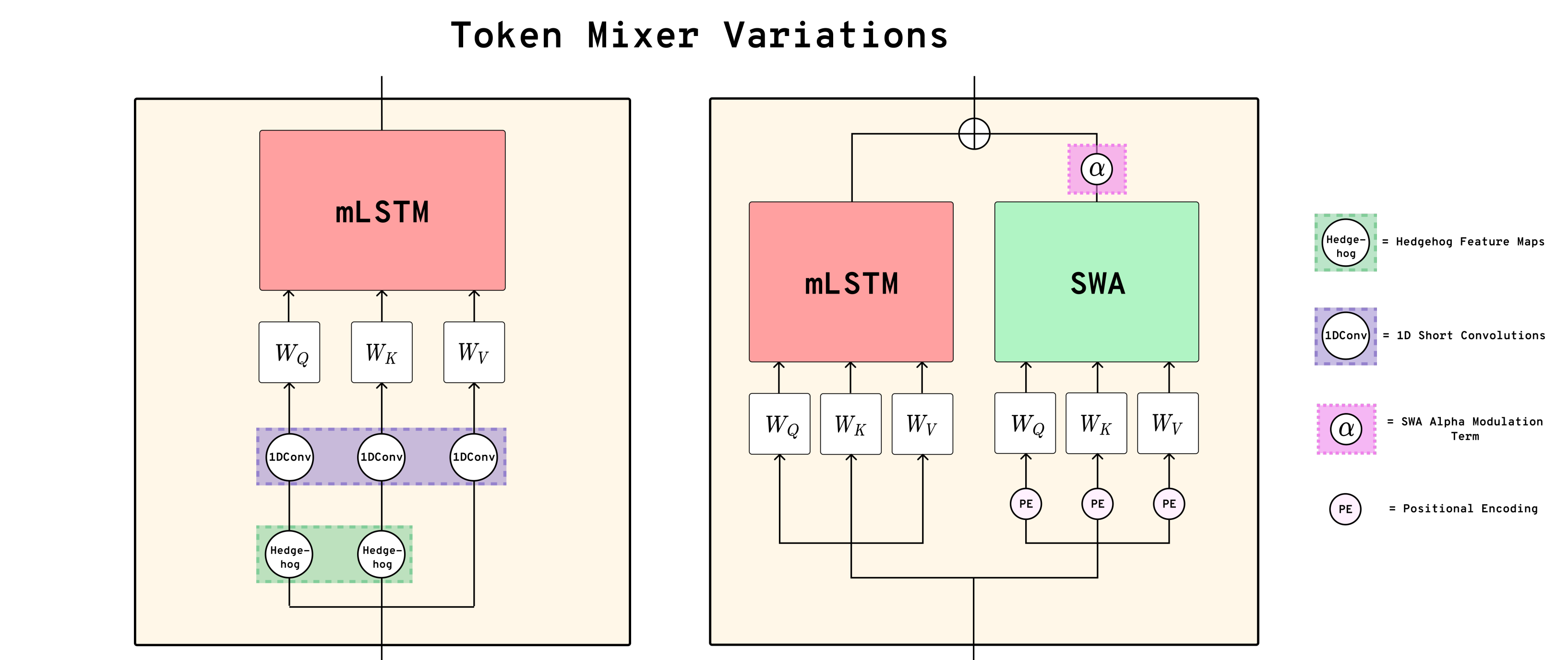


## Improving Subquadratic Language Models under a Tiny Data Budget.

### Why this matters: BabyLM Architecture

- Strict budgets: 10M or 100M words, **10 epochs**, small models.
- We study **sample-efficient** training without scaling: replace self-attention with **mLSTM** and add *lightweight* tweaks.
- Goal: Practical wins on zero-shot linguistic/educational tasks under tight compute.

| Architecture | STRICT-SMALL | STRICT |
|--------------|--------------|--------|
| Transformer  | 32.27        | 35.03  |
| mLSTM        | 35.96        | 35.42  |



$$Q, K, V = \text{conv}(Q), \text{conv}(K), \text{conv}(V) \quad Q, K = \text{hedgehog}(Q), \text{hedgehog}(Q)$$

$$h_{total} = h_{LA} + \tanh(\alpha) \cdot h_{SWA}$$

### BLaLM: Baby Linear Attention Language Model

#### Modifications and their use:

- ShortConv[1]: Local Pattern Extraction
- Sliding-Window Attention (SWA)[2]: Short-range token interaction
- DynMod: Learns when to rely on mLSTM vs. SWA
- HedgehogMaps[3]: Mimic Softmax-Attention

| MECHANISM            | STRICT-SMALL |              | STRICT |              |
|----------------------|--------------|--------------|--------|--------------|
|                      | PPL.         | AVG.         | PPL.   | AVG.         |
| BLaLM                | 20.01        | <b>37.27</b> | 7.95   | 35.08        |
| - ShortConv          | 12.37        | 36.41        | 6.48   | 34.57        |
| - SWA                | 12.08        | 36.16        | 7.38   | 35.86        |
| - SWA with Memory    | 10.08        | 34.96        | 6.67   | 37.21        |
| - SWA DynMod         | 9.44         | 36.15        | 7.76   | <b>38.82</b> |
| - SWA DynMod Bounded | 8.58         | 34.41        | 6.84   | 36.21        |
| - Hedgehog           | 6.18         | 33.58        | 6.68   | 36.65        |
| - Hedgehog + SWA     | 7.27         | 36.25        | 6.63   | 34.20        |

Layerwise learned weighting of SWA and mLSTM improves performance.

### Optimizer Choice

| OPTIMIZER | PPL              | AVG.             |
|-----------|------------------|------------------|
| AdamW     | 11.21 $\pm$ 0.11 | 35.75 $\pm$ 1.74 |
| Muon      | 7.95 $\pm$ 0.15  | 36.24 $\pm$ 1.16 |

**Muon**[4]: Scale invariant, norm-preserving updates  
→ Optimizer consistently improves perplexity and zero-shot performance.

### Takeaways

- **Linear-time token mixers (mLSTM)** are a *viable drop-in* for sample-efficient training.
- **Local attention (SWA) + Dynamic Modulation** improves downstream generalization, especially at 100M.
- **Muon** stabilizes and accelerates training in low-data regimes.