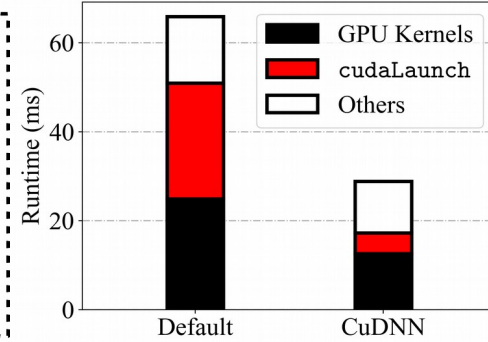
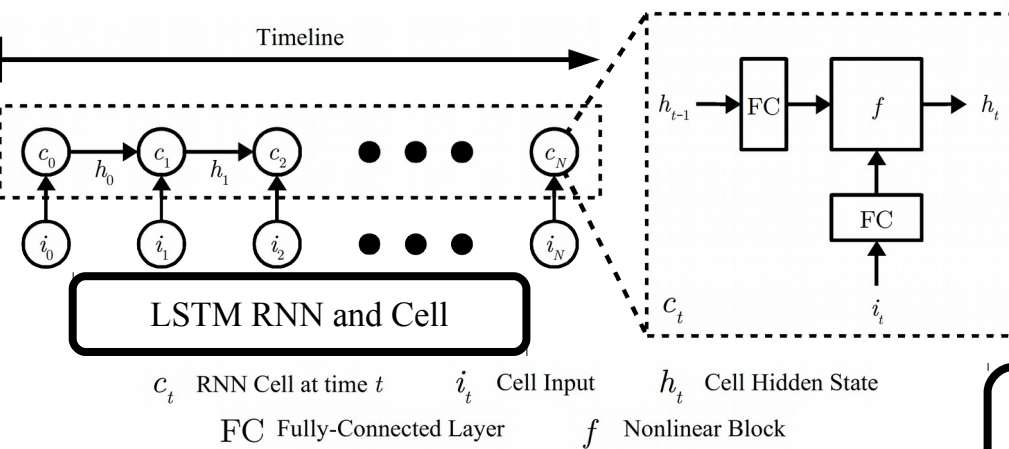


EcoRNN: Fused LSTM RNN Implementation with Data Layout Optimization

Bojian Zheng¹, Akshay Nair¹, Qionsi Wu¹, Nandita Vijaykumar², Gennady Pekhimenko¹

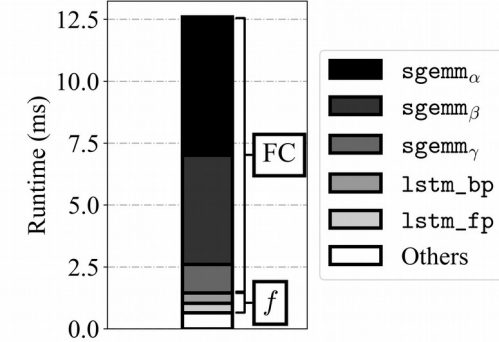
¹University of Toronto, ²Carnegie Mellon University

Background



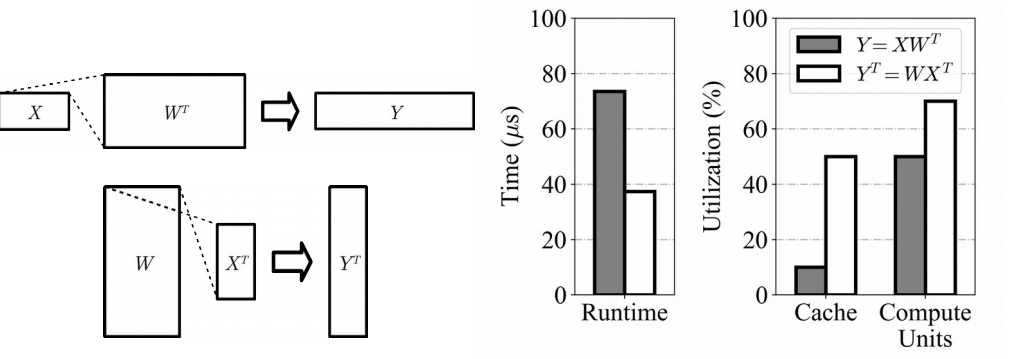
MXNet default implementation has **cudaLaunch** overhead. cuDNN implementation is **closed-source**.

Observation



The runtime bottleneck of LSTM RNN is FC layers.

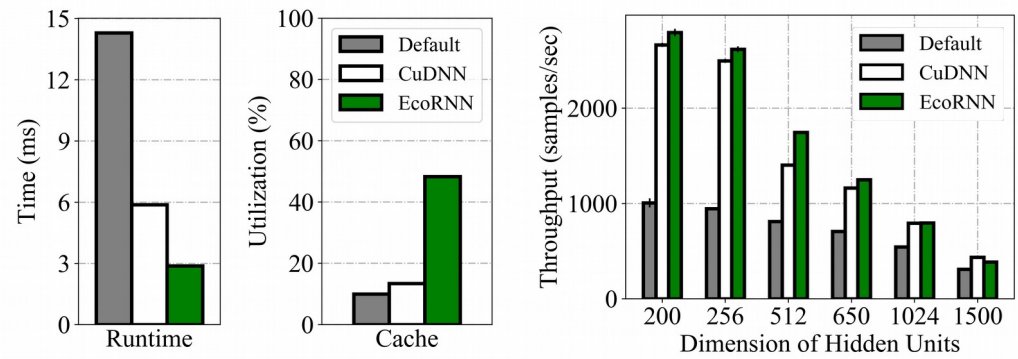
Data Layout Optimization



$Y = XW^T$ vs. $Y^T = WX^T$

Optimizing data layout improves **cache utilization** of FC layers.

Results



Performance comparison of LSTM RNN with different backends

Training throughput on MXNet language modeling benchmark