



EcoRNN: Efficient Computing of LSTM RNN on GPUs

EcoSystem



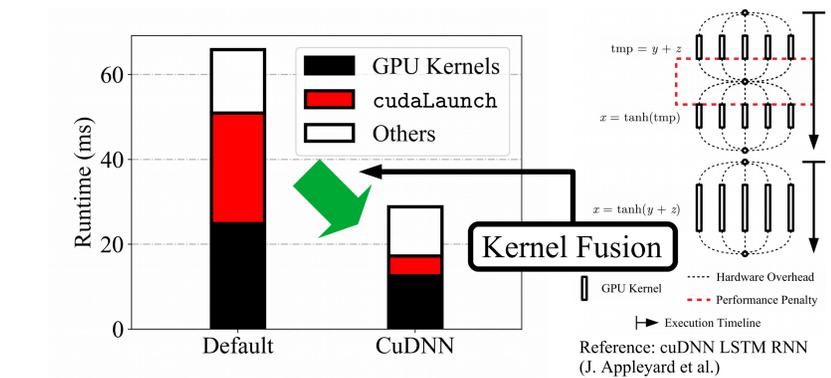
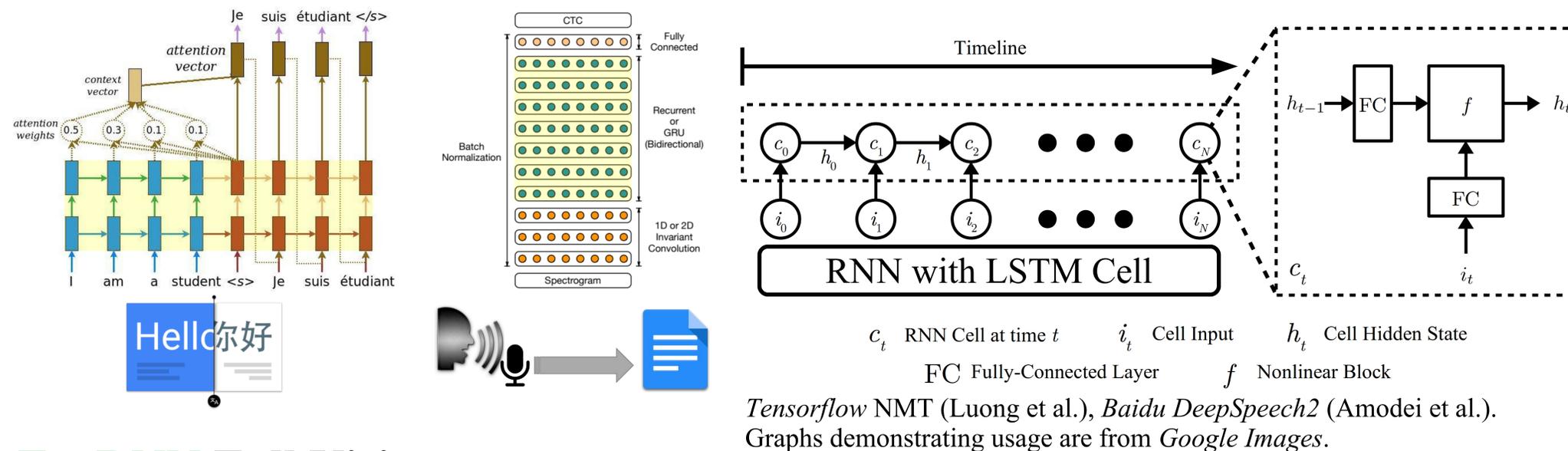
Bojian Zheng (M.Sc. Student), Gennady Pekhimenko (Advisor)

www.cs.toronto.edu/ecosystem

EcoSystem Research Group, Department of Computer Science, University of Toronto

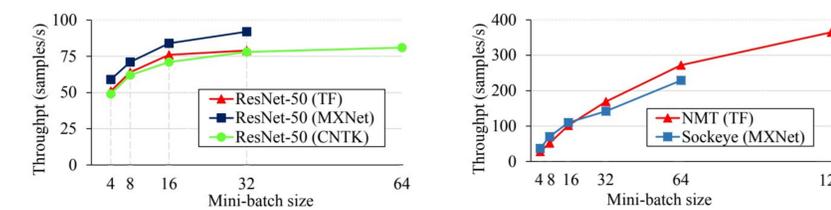
Background: Long-Short-Term-Memory Recurrent Neural Network

Problem Statement



✗ **Default** has **cudaLaunch overhead**.

✗ **CuDNN** is **closed-source, limits innovation**.

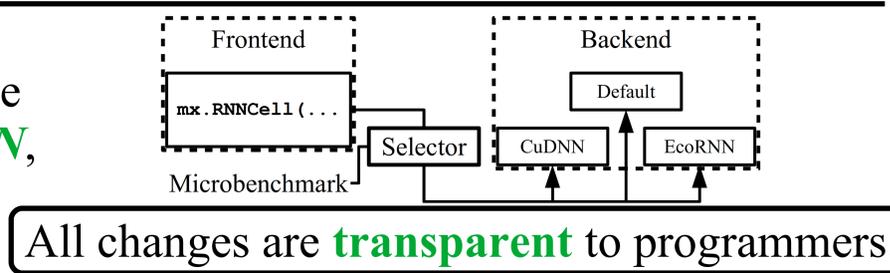


✗ RNN Training is **Memory Capacity**-bounded.

TBD DNN Training Benchmark Suite (Zhu et al., IISWC'18).
tbd-suite.ai

EcoRNN Full Vision

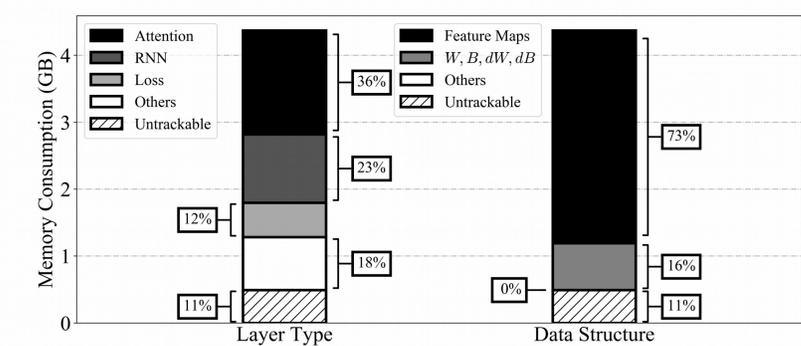
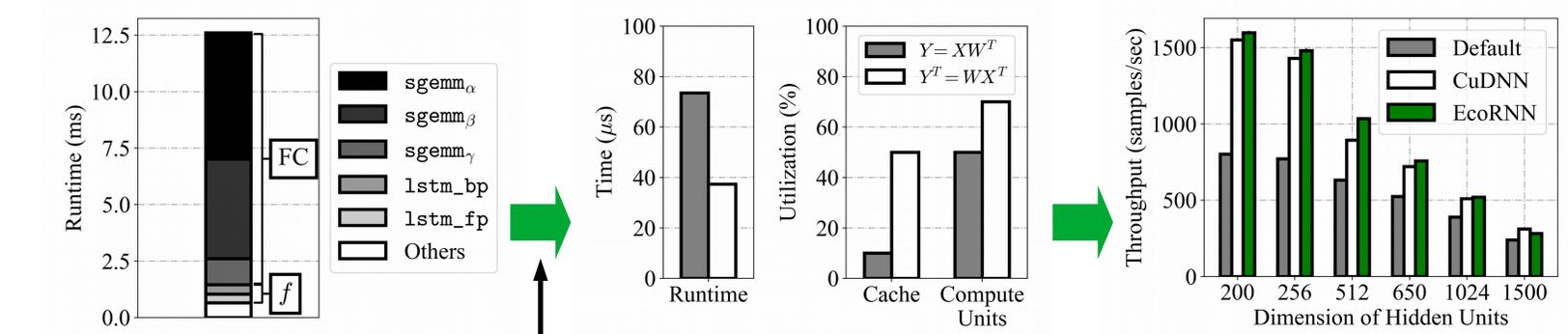
EcoRNN is an **open-source** implementation that has runtime performance **comparable with or even better than CuDNN**, but **consumes less memory** and **supports auto-tuning**.



Preliminary Results: Performance

Memory Consumption

Future Work



Baidu persistent RNN

✓ Weight Parameter Reuse

✓ High Performance when Batch Size is Small

✗ **Inflexible** (hard to be ported to new GPUs and cell types)

The runtime bottleneck is **FC layers**.

Data layout optimization **improves cache hit rate**.

Training Throughput on MXNet Language Modeling benchmark

The memory bottleneck is **Feature Maps of Attention and RNN**.

ML Compilers (e.g., TVM, XLA)

Gist (Jain et al., ISCA'18)