

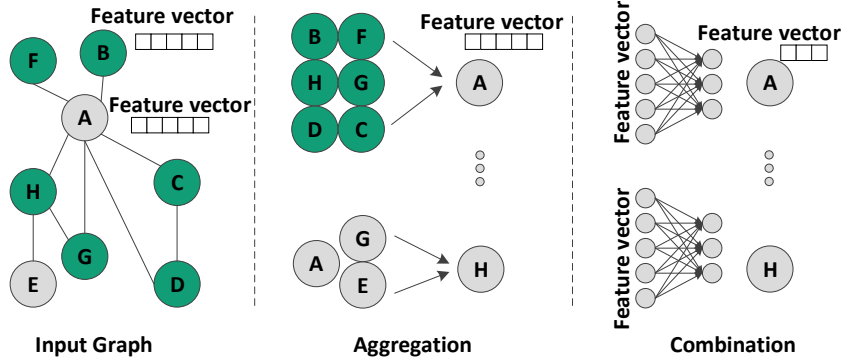
Li JJ, Wang K, Zheng H *et al.* GShuttle: Optimizing memory access efficiency for graph convolutional neural network accelerators. JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY 38(1): 115–127 Jan. 2023. DOI: 10.1007/s11390-023-2875-9

GShuttle: Optimizing Memory Access Efficiency for Graph Convolutional Neural Network Accelerators

Jia-Jun Li (李家军), Ke Wang (王可), Hao Zheng (郑皓)
and Ahmed Louri

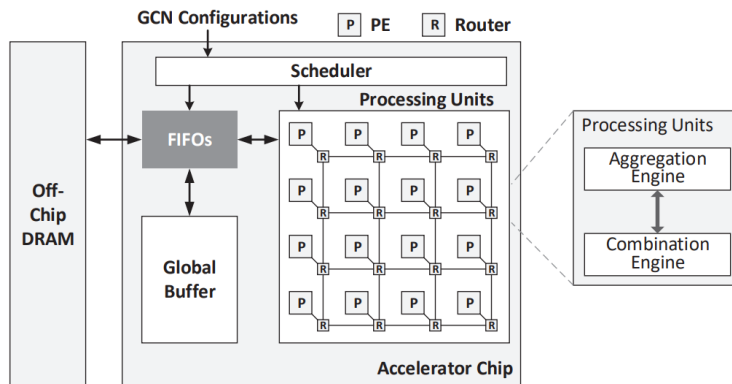
Research Objectives

GCNs are popular but poses a great challenge to the processing hardware



- **Graphs can be very large**
 - Billions of users in social networks
- **Stringent Latency requirements**
 - Power grid cascading failure prediction
- **High Throughput requirements**
 - E-commerce analysis on shopping season

GCN accelerators failed to exploit all the opportunity of data reuse



Optimizing memory access efficiency is the key to energy efficiency for GCN accelerators

A typical GCN accelerator

Research Method

- Define the optimization problem and build analytical models for memory accesses

$$\begin{aligned}
 & \underset{\mathbf{x}}{\text{Minimize}} && V = V_d(\mathcal{X}^t, \mathcal{X}^{oo}, \mathcal{X}^f) + \omega \cdot V_s(\mathcal{X}^u, \mathcal{X}^{oi}) \\
 & \text{s.t.} && 0 < T_m \leq M, \quad 0 < T_k \leq K \\
 & && 0 < T_{n0} \leq N, \quad 0 < T_{n1} \leq N \\
 & && 0 < T_{c0} \leq C, \quad 0 < T_{c1} \leq C \\
 & && S_X + S_W + S_{B1} \leq GLBsize \\
 & && S_A + S_O + S_{B2} \leq GLBsize \\
 & && P_{n0} \times P_{c0} \times P_k \leq \#PEs
 \end{aligned}
 \quad (4)$$

$$\begin{cases}
 S_X = \gamma_X \times T_{n0} \times T_k \\
 S_W = T_k \times T_{c0} \\
 S_{B1} = T_{n0} \times T_{c0} \\
 S_{B2} = T_{n1} \times T_{c1} \\
 S_A = \gamma_A \times T_m \times T_{n1} \\
 S_O = T_m \times T_{c1}
 \end{cases}
 \quad (4)$$

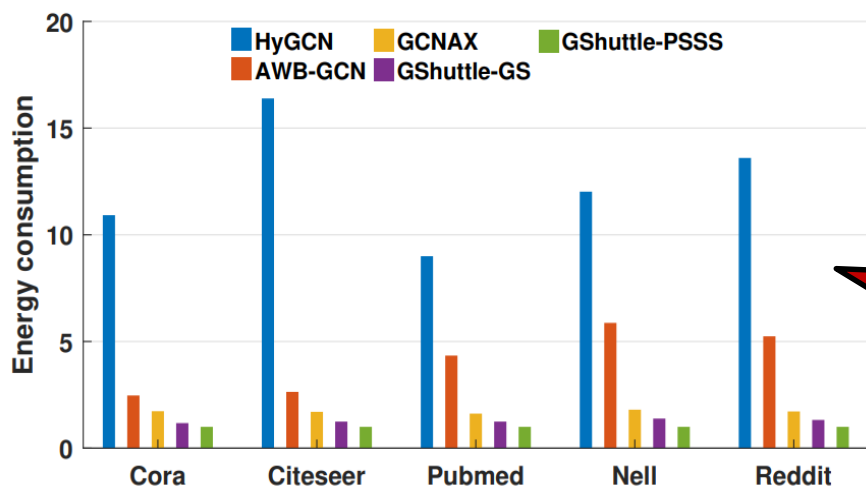
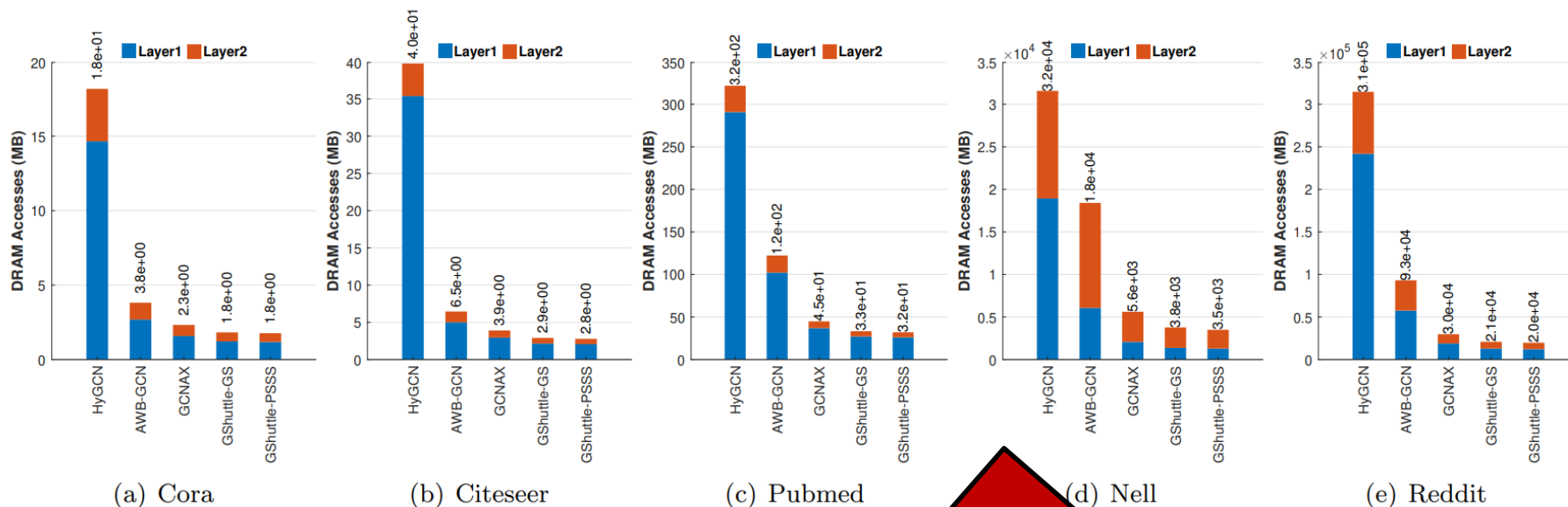
$$\begin{cases}
 \alpha_X = \alpha_W = \frac{N}{T_{n0}} \times \frac{C}{T_{c0}} \times \frac{K}{T_k} \\
 \alpha_{B1} = \frac{N}{T_{n0}} \times \frac{C}{T_{c0}} \\
 \alpha_{B2} = \alpha_A = \frac{M}{T_m} \times \frac{C}{T_{c1}} \times \frac{N}{T_{n1}} \\
 \alpha_O = \frac{M}{T_m} \times \frac{C}{T_{c1}}
 \end{cases}
 \quad (5)$$

- Develop two algorithms to solve the problem

Table 2. The greedy search algorithm to determine the design variables.

Conditions	Loop Fusion	Inter-tiling Loop Order	Tile Size Setting Priority
$N \cdot C \geq GLB_{size}$	No	$n_0 \rightarrow c_0 \rightarrow k, m \rightarrow c_1 \rightarrow n_1$	① T_{n0}, T_m ② T_{c0}, T_{c1} ③ T_{n1}, T_k
$N \cdot C < GLB_{size}$	Yes	$n_0 \rightarrow c_0 \rightarrow k \rightarrow m$	① T_{n0}, T_{n1} ② T_{c0}, T_{c1} ③ T_m, T_k

Research Results



GShuttle reduces DRAM accesses up to a factor of 1.7X compared to the SOTA approaches

GShuttle improves energy efficiency by 1.4X compared to the SOTA approaches

Research Conclusions

- GShuttle can find the optimal design variables of GCN dataflow under certain design constraints.
- The results show that GShuttle could significantly reduce the number of DRAM and SRAM accesses for GCN Accelerators.
- we expect that GShuttle can be applied to many existing GCN accelerators such as HyGCN and AWB-GCN.