High-Performance Graph Analytics Accelerator

Long Zheng, Pengcheng Yao Huazhong University of Science & Technology





Outline

Introduction of Graph Analytics

- Background of Graph Analytics
- The importance of Graph Analytics
- Why we need Graph accelerator
- □ Introduction of our work
 - Performance analysis
 - Computation Architecture Designs
 - Memory System Designs
- Conclusion



What is Graph?



> Origin: Seven Bridges of Königsberg



What is Graph?



- > Origin: Seven Bridges of Königsberg
- Graph: describing the relationships (edges) between different entities (vertices) with high flexibility and accuracy
 - ✓ Rest Area -> Vertex, Bridge -> Edge
 - ✓ An undirected graph with four vertices and seven edges



What is Graph?



Graph is ubiquitous!



Graph analytics: solving real-world problems by analyzing graphs



Each vertex in the input graph iteratively updates itself based on the properties of its neighbors

CGCL

✓ Graph processing

- Computing vertex/edge properties to analyze graph property, e.g., BFS, SSSP, PageRank
- ✓ Gather-Apply-Scatter: iterative computation based on neighbors



Graph analytics: solving real-world problems by analyzing graphs



Model how people interact and influence each other to predict COVID-19 outbreaks

CGCL

✓ Graph processing

- Computing vertex/edge properties to analyze graph property, e.g., BFS, SSSP, PageRank
- ✓ Gather-Apply-Scatter: iterative computation based on neighbors
- ✓ Widely adopted to analyze relationships



Graph analytics: solving real-world problems by analyzing graphs



Each **subgraph** in the input graph **iteratively extends itself** based on the **edges** of its **neighbors**

CGCL

✓ Graph mining

- Searching interesting graph patterns to analyze graph structure, e.g., subgraph matching, motif counting
- ✓ Extend-Filter-Process: iterative extension based on neighbors



Graph analytics: solving real-world problems by analyzing graphs



✓ Graph mining

- Searching interesting graph patterns to analyze graph structure, e.g., subgraph matching, motif counting
- ✓ Extend-Filter-Process: iterative extension based on neighbors
- ✓ Widely used in social network, finance, and biology







NLP

Image Processing







IOT

We are in a golden age of AI: Machine Learning (ML) technologies are ubiquitous in our life



ASIC Designs



EDA Tools



Neuroscience



Why we need to research graph analytics (GA), which seems to be out of date?







- Reason: the need of processing sparse data in real-world applications
- Both of ML and GA can be abstracted as matrix computation
 - ✓ <u>ML</u>: Dense matrix computation, where most of the elements are nonzero (e.g., 100% in DNN)
 - ✓ <u>GA</u>: Sparse matrix computation, where most of the elements are zero (e.g., 99.9999997% in Twitter Graph)





- > ML and GA are adopted in different areas
 - ✓ <u>ML</u>: a representative of regular applications which process dense data
 - ✓ **GA**: a representative of irregular applications which process sparse data
 - ✓ Both of them are necessary for solving real-world problems



Sparsity is ubiquitous



We only interact with a limited number of items

Sparsity also exists in dense applications



ML also has sparsity (Sparse NN, GNN)

Efficient processing technologies are missing



Existing SWs & HWs are designed for dense data

CGCL

Researches on sparsity are increasingly important

- Sparsity is ubiquitous in our life
- ✓ Even typical dense applications might have some kind of sparsity
- Most of existing SWs and HWs are designed for dense data
- Efficient processing technologies for sparse data are essential



Why We Need Graph Accelerators?

Sparsity leads to irregular computation





Machine Learning

Graph Analytics

- > ML: dense data process is regular and sequential
- GA: sparse data process is irregular and random

Irregularity is inefficient on generalpurpose architecture



1 Data contention on the same irregular data

- (2) Mismatched access granularity
- > Sparsity in graph data causes significant irregularity
 - ✓ Sparse data leads to irregular computations and memory accesses
 - General-purpose processors (e.g., CPU, GPU) are designed to efficiently process dense data
 - Sparse data cases significant data contentions and random accesses



Outline

Introduction of Graph Analytics

- Background of Graph Analytics
- The importance of Graph Analytics
- Why we need Graph accelerator

□ Introduction of our work

- Performance analysis
- Computation Architecture Designs
- Memory System Designs

Conclusion



Building A General-Purpose Graph Analytics Computer





ĊĠCL

Our Work in Accelerator Designs





Performance Analysis — ICDCS'17

		Prior Work	Our Work
Algorithm	Not Stalled	Memory Sys.	Inside Core
Breadth-First Search	25.486	38.022	36.492
PageRank	23.704	46.683	29.613
Connected Components	20.691	43.751	35.558
Triangle Counting	23.674	37.750	38.576

Understanding the Performance bottlenecks

- Existing work believes that graph analytics is bottlenecked by memory accesses
- ✓ We find that most (i.e., 41%) of the CPU cycles are stalled on memory accesses
- ✓ There are still 35% of the CPU cycles are stalled on core execution



Performance Analysis — ICDCS'17





Performance Analysis — ICDCS'17



Limitation of Out-of-Order Execution

- Out-of-order execution, but sequential retirement
- When facing a long latency instruction, OOO buffers can soon be clogged by succeeding instructions that are already executed
- \checkmark OOO buffers account for 50.3% of the total resource-related stalls
- Long latency instructions in graph analytics
 - ✓ Computation: atomic operations; Memory: random vertex/edge accesses

CGCL

More details in "Towards Dataflow-based Graph accelerator"



- Inefficiencies of atomic operations in graph analytics
 - ✓ Sparsity cause irregular computations, which lead to data contention





Inefficiencies of atomic operations in graph analytics

- ✓ Sparsity cause irregular computations, which lead to data contention
- To ensure the correctness of results, existing work uses atomic operations to serialize the process of conflict data





Inefficiencies of atomic operations in graph analytics

- ✓ Sparsity cause irregular computations, which lead to data contention
- To ensure the correctness of results, existing work uses atomic operations to sequentially access the conflict data
- ✓ Overheads of atomic operations: 45%



Whether we can parallelize the process of conflict data while ensuring the accuracy of results?



Insight: Parallel accumulation

Algorithm	Operation Type	
Breadth-First Search	CAS if less	
Weakly Connected Components	CAS if less	
Shortest Path	CAS if less	
PageRank	Atomic add	
Triangle Counting	Atomic add	
Degree Centrality	Atomic add	
Collaborative Filtering	Atomic add	

 Characteristics of graph atomic operations: following commutative law and associative law



 Insight: accumulating the atomic operations in parallel before writing back the results

ĊĠCL



The insight sounds to be promising, but not easy to be efficiently implemented



Challenge: Irregular vertex degrees



Traditional accumulator (adder tree)



Suboptimal performance of adder tree

- Different vertex has different number of degrees (i.e., the data received in each cycle belongs to multiple vertices)
- Traditional adder tree can only accumulate for one vertex, leading to suboptimal performance
- Using multiple adder trees is kind of a cop-out. It works, but significantly increases hardware fanouts



- Problem Formulization
 - ✓ Original (integer programming): $p_j = \sum_{1 \le i \le N} a_i \cdot b_{ij}, 1 \le j \le M_i$
 - ✓ Characteristic

partial order in pull model, i.e., all requests are sent in an ascending order of destination vertex

 $\square \quad b_{ij} \le b_{kj}, \text{ if } i \le k$

✓ **Simplified (Prefix Sum)**: $p_j = f(c_j^2)$, where

$$f(i) = \begin{cases} f(i-1) + a_i, & i \notin \{c_1^1, c_2^1, \dots, c_M^1\} \\ \\ a_i, & i \in \{c_1^1, c_2^1, \dots, c_M^1\} \end{cases}$$

- Differences between traditional prefix sum and our problem
 - ✓ Breakpoint
 - \square f(i) is initialized to a_i when a break point c_m^1 appears
 - **D** Accumulator needs to find all dynamically changing c_m^1
 - ✓ Filter:
 - **D** Only $f(c_j^2)$ is the valid results, where c_j^2 also dynamically changes
 - Accumulator needs to precisely select the accumulated results



Methodology

- Optimize traditional prefix sum adders to mitigate the gaps
- We choose Ladner-Fischer adders for minimized latency and resources
- ✓ **Partial order** in $c_m^1(c_i^1 \le c_j^1)$ if $i \le j$)
- Adding log(|Pipeline|) bits to each update as vertex ID
- ✓ c_m^1 can be generated in runtime by comparing the vertex ID of consecutive updates



$$p_{j} = f(c_{j}^{2}), \text{ where }$$

$$f(i) = \begin{cases} f(i-1) + a_{i}, & i \notin \{c_{1}^{1}, c_{2}^{1}, \dots, c_{M}^{1}\} \\ a_{i}, & i \in \{c_{1}^{1}, c_{2}^{1}, \dots, c_{M}^{1}\} \end{cases} \qquad f(i) = \begin{cases} f(i-1) + a_{i}, & \upsilon_{i} = \upsilon_{i-1} \\ a_{i}, & \upsilon_{i} \neq \upsilon_{i-1} \end{cases}$$





Hardware Designs

- Breakpoint recognition: replacing original adder logic with the <u>conditional</u> adding logic
- Data filtering: select the accumulated results in the location of each vertex's <u>last</u> <u>update</u>





> Example

- Fully pipelined accumulation
- ✓ Accumulator receives 8 updates in the 1st cycle
- Adder directly sends the second data when finding a breakpoint in the 2nd and 3rd cycle
- Because the location of the last updates of vertex 1 is 3, accumulator select the data in the 3rd port as the results for vertex 1



Evaluation

- ✓ **Platform:** Xilinx U250 Acceleration card, use 2 channel memory only
- Performance (170MHz ~ 250MHz):
 1.23 1.86 GTEPs (Overall) or 3.69 5.58 GTEPs (Single Iteration) for BFS
 3.22 3.98 GTEPs for PageRank
- ✓ **Normalized Speedup: 3.6-6.2x** speedup comparing to state-of-the-art
- 3.96x speedups for finance anti-fraud applications from Ping An Technology



More details in "An Efficient Graph Accelerator with Parallel Data Conflict Management"





Memory Access Characteristics

- ✓ Graph processing: random vertex access
- Graph Mining: random vertex and edge access
- Existing accelerators statically hold all vertex data on-chip
- While this mechanism is efficient for graph processing, it is not practical to hold all edge data for graph mining with limited on-chip memory capacity



Insight: Exploration Locality



An intuition is whether we can achieve considerable performance by storing only a small portion of data?

Exploration locality: the power-law distribution in natural graphs is amplified during the process. 5% data can generate at most 95% memory accesses





- Methodology: Divide-and-Conquer
 - Statically retaining frequently accessed data, and dynamically replacing the left
- > Challenge
 - ✓ Frequency calculation: hard to precisely locate the frequently accessed data
 - ✓ Intermediate results: hard to be stored in off-chip memory





victim = arg max{*Rank*(ON₁(*v*)) + $\lambda Rec(v)$ }, $v \in V_{Low}$

Hardware Design

- <u>Memory hierarchy</u>: static memory (scratchpad memory) + dynamic memory (cache)
- Observation: access frequency of a vertex is mainly determined by the degrees of itself and its 1-hop neighbors
- ✓ <u>Data addressing</u>: predict access frequency with high accuracy and low overheads

CGCL

✓ <u>Replacement policy</u>: Locality-aware replacement



Hardware Design

- ✓ <u>DFS-based architecture</u>: extending the subgraphs in DFS manner to avoid writing back intermediate results
- ✓ <u>Slot-based pipeline</u>: efficiently process different subgraphs in parallel in a dataflow manner
- ✓ <u>Data compression</u>: compacting the parameters of DFS traverse
- ✓ <u>Load-balance</u>: aggressive and precise work-stealing mechanism







- Evaluation
 - ✓ **Platform:** Xilinx U250 FPGA acceleration card
 - ✓ **Performance:** ≥ 2.4x speedups comparing to state-of-the-art Fractal
 - ✓ **Energy efficiency:** \ge 33.6x comparing to state-of-the-art Fractal



More details in "A Locality-Aware Energy-Efficient Accelerator for Graph Mining Applications"



Outline

Introduction of Graph Analytics

- Background of Graph Analytics
- The importance of Graph Analytics
- Why we need Graph accelerator
- □ Introduction of our work
 - Performance analysis
 - Computation Architecture Designs
 - Memory System Designs

□ Conclusion



Conclusion

- Graph is ubiquitous in our life
- Researching graph analytics is important
- Graph processing is bottlenecked by both computation and memory access
- Atomic operations in graph analytics can be accumulated in parallel to avoid pipeline stalls
- Exploration locality can be leveraged to achieve consider memory performance on large volume of graph data





Thanks!





