Parallelization of All-Pairs-Shortest-Path Algorithms in Unweighted Graph

Masahiro Nakao, Hitoshi Murai, Mitsuhisa Sato

(RIKEN Center for Computational Science)

Background (1/3)

- Network topology of a large-scale parallel computer system affects the overall performance significantly
 - Supercomputer
 - Data center



- For a parallel application where the communication pattern is NOT known in advance (e.g. Big data), some researches consider that a network topology with randomness is effective
- The reason is that diameter and average hops of the random graph are smaller than those of a regular graph (e.g. k-ary ncube) due to the small-world effect

Background (2/3)

- By describing the calculation node as "vertex" and the network wiring as "edge", the network topology is represented as an unweighted graph with number of vertices (n) and degree (d)
- Diameter and average hops are important metrics of network topology, they are calculated using All-Pairs-Shortest-Path (APSP) algorithm



Background (3/3)



- To create a random graph with small diameter and average hops, Simulated Annealing (SA) are often used
- SA repeats APSP calculation many times
 - However, the calculation cost of APSP is very high !!
 - The complexity is proportional to the square of the number of vertices (n)

e.g. For a problem (n, d) = (1M, 32), the time required for one APSP is **about 37 hours** by the methods based on BFS on Intel Xeon Gold 6126

Objective and a part of results

- Our previous research in HPC Asia 2019 provides SA algorithm where APSP is calculated using Breadth-First Search (BFS-APSP)
 - BFS-APSP is parallelized by MPI+OpenMP
- This research introduces another APSP algorithm based on adjacency matrix (ADJ-APSP), and compares with BFS-APSP
 - ADJ-APSP is parallelized by MPI+OpenMP for multi-core cluster
 - ADJ-APSP is parallelized by MPI+CUDA for GPU cluster

BFS-APSP (about 37 hours) \rightarrow ADJ-APSP (3,880 sec.) \rightarrow ADJ-APSP by MPI+OpenMP on 64 CPUs x 12 Cores (6.77 sec.) \rightarrow ADJ-APSP by MPI+CUDA on 128 GPUs (0.28 sec.)

You can download programs from

https://github.com/mnakao/APSP/

Agenda

- Background
- BFS-APSP
- ADJ-APSP
- Performance
- Summary

Serial BFS-APSP

- BFS can find the hops from one vertex to others
 - APSP can be obtained by executing BFS from all vertices
 - Top-down approach is used
- The computational complexity of BFS is O(nd).
 - The computational complexity of BFS-APSP is O(n²d) because BFS is repeated n times



Parallel BFS-APSP

- Multiple BFSs are performed in parallel using MPI from different vertices, and one BFS is executed using OpenMP
- MPI parallelization,
 - Starting points are assigned to each MPI process evenly
 - Thus, the maximum number of MPI processes is *n*
 - Communication time is small because only scalar values of diameter and average hops in each MPI process are collected at the end of the program
- OpenMP parallelization,
 - Each OpenMP thread searches not-visited vertices
 - The implementation requires exclusive control to update own list of notvisited vertices

Agenda

- Background
- BFS-APSP
- ADJ-APSP
- Performance
- Summary

Serial ADJ-APSP(1/3)

- Let A be an adjacency matrix of a graph
- If the value of an element a_{i, j} in A^k is 1, it means that the vertex i can reach the vertex j within k hops



Serial ADJ-APSP(1/3)

- Let A be an adjacency matrix of a graph
- If the value of an element a_{i, j} in A^k is 1, it means that the vertex i can reach the vertex j within k hops



Serial ADJ-APSP(2/3)



- When all elements are 1, the value of k is Diameter
- Average hops can be calculated by summing all elements whose value for a_{i, j} is 0 divided by number of elements (170/90 ≒ 1.89)

Serial ADJ-APSP(3/3)



Parallel ADJ-APSP



Parallel ADJ-APSP



Parallel ADJ-APSP for GPU



Comparison between BFS-APSP and ADJ-APSP

	BFS-APSP	ADJ-APSP	
Computational complexity	O(n ² d)	O(n ² d <mark>D/E</mark>)	n: vertices d: degree D: diameter *
Maximum number of MPI processes	n	n/E**	E: bits in element (64 (we use uint64_t)
OpenMP exclusive control	critical directive	(none)	
For GPU	\bigtriangleup	0	
Communication	MPI_Allreduce() for scalar x 2		

* In general, the value of D of graphs with randomness is small due to the small-world effect.

** The number of elements of columns in an adjacency matrix

Agenda

- Background
- BFS-APSP
- ADJ-APSP
- Performance
- Summary

Experiment environment

The K computer in RIKEN R-CCS



Cygnus in CCS, Univ. of Tsukuba



CPU	SPARC64 VIIIfx (8Cores, 2.0GHz)
Memory	DDR3 (64GB/s, 16GB)
Network	Torus fusion six-dimensional mesh/torus network, 5GB/s \times 10
Software	Fujitsu Compiler K-1.2.0-25

For MPI+OpenMP versions

CPU	Intel Xeon Gold 6126 (12Cores, 2.6GHz) \times 2
Memory	DDR4 (128GB/s \times 2, 192GB)
GPU	NVIDIA Tesla V100 (900GB/s, 32GB) \times 4
Network	InfiniBand HDR100 (12.5GB/s) \times 4
Software	intel/19.0.3, mvapich/2.3.1, cuda/10.1

For MPI+OpenMP versions For MPI+CUDA version

Serial algorithm



- (n, d, D) = (50, 4, 5), (1726, 30, 3), and (64Ki, 6, 9) *64Ki is 65,536.
- Test programs are available at http://research.nii.ac.jp/graphgolf/
- ADJ-APSP is always faster than BFS-APSP
 - 8.08 to 29.49 times better

Parallel algorithm by OpenMP



- (n, d, D) = (64Ki, 6, 9) and (1M, 32, 5)
- ADJ-APSP is always faster than BFS-APSP too
 - Improvement ratio of ADJ-APSP is better than that of BFS-APSP
 - In case of (64Ki, 6, 9) using 8 threads on the K computer, improvement ratio of ADJ-APSP is 7.18, and that of BFS-APSP is 3.54
 - In case of (1M, 32, 5) using 12 threads on Cygnus system, ADJ-APSP(1threads) : 3,880sec. → ADJ-APSP(12 threads) : 475sec.

Parallel algorithm by MPI+OpenMP



- The number of threads is set to the maximum value
- The maximum number of processes in (64Ki, 6, 9) is 65,536 for BFS-APSP and 1,024 for ADJ-APSP, respectively
- The maximum number of processes in (1M, 32, 5) is 1,000,000 for BFS-APSP and 15,625 for ADJ-APSP, respectively
- ADJ-APSP is faster than BFS-APSP in the same number of processes
- ADJ-APSP(1CPU) : 475sec. \rightarrow ADJ-APSP(64CPUs) : 6.77sec.

Parallel algorithm by MPI+OpenMP



- The number of threads is set to the maximum value
- The maximum number of processes in (64Ki, 6, 9) is 65,536 for BFS-APSP and 1,024 for ADJ-APSP, respectively
- The maximum number of processes in (1M, 32, 5) is 1,000,000 for BFS-APSP and 15,625 for ADJ-APSP, respectively
- ADJ-APSP is faster than BFS-APSP in the same number of processes
- ADJ-APSP(1CPU) : 475sec. \rightarrow ADJ-APSP(64CPUs) : 6.77 sec.

Parallel algorithm by MPI+CUDA



- 1GPU v.s. 1CPU
 - (64Ki, 6, 9): 0.751 sec. (1CPU) \rightarrow 0.061 sec. (1GPU) x 12.6
 - (1M, 32, 5): 475 sec. (1CPU) \rightarrow 28.7 sec. (1GPU) x 16.5
- Multiple GPUs
 - (64Ki, 6, 9): 0.0608 sec. (1GPU) \rightarrow 0.0015 sec. (64GPUs) x 38.4
 - (1M, 32, 5): 28.7 sec. (1GPU) → 0.28 sec. (128GPUs) x 101.1

Parallel algorithm by MPI+CUDA



- 1GPU v.s. 1CPU
 - (64Ki, 6, 9): 0.751 sec. (1CPU) → 0.061 sec. (1GPU) x 12.6
 - (1M, 32, 5): 475 sec. (1CPU) → 28.7 sec. (1GPU)
 x 16.5
- Multiple GPUs
 - (64Ki, 6, 9): 0.0608 sec. (1GPU) → 0.0015 sec. (64GPUs) x 38.4
 - (1M, 32, 5): 28.7 sec. (1GPU) → 0.28 sec. (128GPUs) x 101.1

Agenda

- Background
- BFS-APSP
- ADJ-APSP
- Performance
- Summary

Summary

- We parallelize BFS-APSP and ADJ-APSP using MPI+OpenMP and MPI+CUDA
- ADJ-APSP has a better performance in the serial algorithm and threaded algorithm than BFS-APSP on a single CPU
 - approx. 37hours (BFS, 1core) \rightarrow 3,880sec. (ADJ, 1core) \rightarrow 475sec. (ADJ, 12cores)
- However, because the maximum number of processes of BFS-APSP is larger than that of ADJ-APSP, the performance of BFS-APSP on a large computational resource may be better than that of ADJ-APSP
- We achieved further speedup by parallelizing ADJ-APSP using GPUs
 - 28.7sec. (ADJ, 1GPU) \rightarrow 0.28 sec.(ADJ, 128GPUs)