



Cyberscience

Center

# Quantum Compiler

(URL) https://www.sc.cc.tohoku.ac.jp (E-mail) yuta.sasaki.p4@dc.tohoku.ac.jp

### Automatic Vectorization Assisted by Quantum Annealer

Graduate School of Information Sciences & Cyberscience Center, Tohoku University

Yuta Sasaki, Michael Ryan Zielewski, Mulya Agung, Ryusuke Egawa, Hiroyuki Takizawa

c[i]=b[i]; //s3

Vectorized with

## Background

- Most modern processors support vector processing or similar techniques to improve performance.
- Compilers perform automatic loop vectorization after checking if the loop is safely vectorizable.
  - Programmers sometimes need to manually optimize a loop to help compiler.
  - When there is a cyclic data dependency in a loop, variable renaming is required to break it.
  - There are several kinds of variable renaming techniques for optimization, and the performance after vectorization depends on the techniques.

non-optimal renaming Vectorized with Reduced by 21% optimal renaming Original Loo 0.6Execution Time [sec/10,000 iterations] Fig.1 Speed up by vectorization

### **Problem – Combinatorial Optimization for Vectorization –**

- A loop containing many cyclic data dependencies ( $c_i$ ) needs several renaming techniques ( $R_i$ ).
- Finding an optimal combination of  $R_i$ s is a Set Cover Problem (SCP) [1]. To break all cyclic data dependencies (1) and to minimize runtime overhead (2).



- It is NP-complete. Computational time increases rapidly with the number of cyclic data dependencies.
- If an optimal combination instead of suboptimal ones, can be obtained within a practical time, compilers would perform better vectorization leading to higher performance.



### **Proposal** – Quantum Annealing for Compiler Support –

Find an optimal combination of R<sub>i</sub>s by using quantum annealing (QA), an emerging technology to solve combinatorial optimization problems.

- Model the SCP problem by PolySIMD [2], and solve the problem using QA.
  - A C source code as input, and an optimal combination of  $R_i$ s as output.
  - Need to convert SCP and its constraint to a QUBO (Quadratic Unconstrained Binary Optimization) problem, which D-Wave can process.



### Evaluation

Setup and Tools		QA solver	DW_2000Q_6	
Processor	Intel Core i7-9700K		QUBO generation tool	PyQUBO[3]
Memory	32GB		Embedding tool	minorminer
MIP solver (for comparison)	CBC[4]		Number of reads qubit	num_reads = 1

#### Input a loop containing eight dependency cycles (Fig. 2).

- The optimal solution is verified by matching the result with that of CBC, an exact solver.
  - Break all cycles (Fig. 3).
  - Execution time is was reduced by 21 % (Fig .1).

#### **Experimental results demonstrate the feasibility of QA** to assist automatic loop vectorization by compilers



### Processing time compared with CBC

- QA requires a longer time for this instance.
- Embedding is time-consuming.
  - CBC \* Only QA needs to generate QUBO and embedding. \* QA solution time is measured from D-Wave machine time, not including communication time.



#### Processing time related to problem size Solve SCP instances with $R_i$ and $c_i$ for i=1,...N.

- Solution time of CBC increases sharply in N > 50.
- QA preprocessing time (Generate QUBO and Embedding) also increases with N.
- Solution time of QA is not affected by N and constant for the number of reads.



If QA can process larger problems in the future, our proposal method may be able to solve SCP faster than the existing methods using only classical machines. > We need to reduce the overhead of preprocessing as problem size increases.

### **Conclusion and Future Work**

- Assuming offload to QA, we can implement compiler optimization without considering the number of cyclic dependencies.
- QA is promising not only to replace a combinatorial optimization part of HPC application execution but also to help automatic compiler optimization and hereby improve the performance even on conventional HPC applications.

.00

- Future Work : We will discuss how to reduce the embedding time, which accounts for a large proportion of the total execution. Precompute and store embeddings to reduce embedding time to hashing and lookup.
  - > Small QUBOs with 64 variables or less can be embedded as complete graphs, but this method is inefficient and not applicable to large problems.
  - $\succ$  Need to introduce parameters other than the number of variables.

#### The authors are grateful to Prasanth Chatarasi and Vivek Sarkar for willing to contribute the source code for PolySIMD, their tool Acknowledgement to extend vectorization technologies. Their study motivates our proposed method and is very helpful in its implementation.

This work is partially supported by MEXT Next Generation High-Performance Computing Infrastructures and Applications R&D Program "R&D of A Quantum-Annealing-Assisted" Next Generation HPC Infrastructure and its Applications," Grant-in-Aid for Scientific Research(B) #16H02822 and #17H01706.

[1] Pierre-Yves Calland, Alain Darte, Yves Robert, and Frédéric Vivien. On the removal of anti-and output-dependences. International Journal of Parallel Programming, 26(3):285–312, 1998. Cences [2] Prasanth Chatarasi, Jun Shirako, Albert Cohen, and Vivek Sarkar. A unified approach to variable renaming for enhanced vectorization. In International Workshop on Languages and Compilers for Parallel Computing, pages 1–20. Springer, 2018.

[3] Kotaro Tanahashi. Pyqubo. https://github.com/recruit-communications/pyqubo.

[4] Saltzman, Matthew J. "COIN-OR: an open-source library for optimization." Programming languages and systems in computational economics and finance. Springer, Boston, MA, 2002. 3-32.