An Evaluation of a Hierarchical Clustering Method Using Quantum Annealing

Introduction – Quantum Annealing and Clustering –

- ✓ Digital computers (e.g. CPU, Vector Processer Unit (VPU), GPU)
- Multiplicity of various problems
- × Performance limited by the Moore's law
- × Huge amounts of power consumption
- Quantum annealing (QA) (on Quantum Processer Unit (QPU))
 Power efficient
- Specialized to Combinatorial Optimization Problems
- × Limited multiplicity of some problems
- \rightarrow Hybrid computing with digital computers and QA
- Hierarchical clustering using QA
- o Possibility of accelerating clustering
- $\times\,$ Small problem size due to the limitation of qubits in QA

- Non-hierarchical clustering (like K-means)
- Low computational complexity
- × Need to know the number of clusters in advance
- Hierarchical clustering
- × High computational complexity
- No need to know the number of clusters in advance
- → Need to use the appropriate clustering method depending on the situation

<u>Objective</u>

Need to clarify features of each processor and each clustering method by comparing the execution time and quality of the clustering method





✓ How to choose representative points by MWIS

• A chunk having $\{x_1, ..., x_n\}$ and each weight of w_i (1) Similarity matrix $N_{ij}^{(\epsilon)} = \begin{cases} 1 & if \ distance(x_i, x_j) < \epsilon \ @ \ Quadratically \ Constrained \ Quadratic \ Program \ (QCQP) \ given \ by \ @ \ 1. \ QCQP \ is being \ solved \ by \ Greedy \ algorithm \ maximize \ \sum_{s=(0,1)^n}^n s_i w_i \ subject \ to \ \sum_{i=1}^n S_i N_{ij}^{(\epsilon)} s_j = 0 \end{cases}$ Binary variables S obtained as solutions of MWIS Binary variables S obtained as solutions of MWIS 2. QCQP is being solved by Greedy algorithm 2. QCQP is transformed into QUB0[2] to solve by SA or QA

Performance Evaluation

Data set : MoCap Hand Postures Quality (Calinski Harabasz score) 5 types of hand postures from 14 users in a motion QA achieves the best quality in HAC capture environment The guality of solutions of QCQP is important for HAC Number of data is 78095 and number of features used score for experiments in 9 of 36 features QA can search the whole search space by quantum Environments Harabasz fluctuations Method Hardware Software SA cannot improve the quality by increasing the annealing time to explore the search space CPU Intel Xeon Gold 6126 Scikit-learn v0.22. 102 Greedy falls a local optimum and cannot find the best VPU NEC Vector Engine Type 10B Frovedis v0.9. Calinski solution GPU NVIDIA Tesla V100 Rapidsai v0.12.0a KMeans(CPU) Greedy K-means is always superior to HAC - SA SA OpenJij v0.0.9 101 KMeans(VPU) Intel Xeon Gold 6126 - OA KMeans(GPU) K-means does not agglomerate data and can avoid Greedy loss of information of the original data QA D-Wave 2000Q Ocean SDK v1.4.0 Number of clusters k Execution time Breakdown of HAC using QA (k=5) Quite low performance of QA QA still needs to be accelerated KMeans(CPU) Greedy Steady execution times of HAC (Greedy, SA, and QA) 10 KMeans(VPU) SA Computation time[s] KMeans(GPU) It is because the number of processed data decreases by 10 agglomerating data for each hierarchy Solution organizin 0.626% The execution times of K-means increase as the Othe 101 1.687 number of clusters increases The execution time of K-means of CPU, VPU, and GPU is 10 long in the cases that the number of cluster is large HAC is effective when the number of clusters is large, Dat 0.013% even if the number of clusters is known Number of clusters k The time of QA is short Discussion : towards the high performance clustering QUBO embedding executed on a CPU HAC ; The number of clusters is unknown dominates the whole execution time K-means; After the number of clusters is decided by HAC QUBO embedding duplicates data to make up for Combination of the hierarchical clustering on QA and non-hierarchical clustering on digital missing connections because of not all qubits computers can be the promising connected in QPU **Conclusions and Future Work** References

- By using QA in HAC, the quality of the clustering results becomes higher than those of SA and Greedy.
- HAC is effective not only when the number of clusters is unknown but also when the number of clusters is known and large.
- As future work, we will combine hierarchical clustering and non-hierarchical clustering for the large-scale clustering.

[1] Tim, J. et al. :A Quantum Annealing-Based Approach to Extreme Clustering, FICC 2020, 2020, DOI: 10.1007/978-3-030-39442-4_15.

[2] Boros, E. et al. :Local search heuristics for quadratic unconstrained binary optimization (QUBO), *Journal of Heuristics*, 2007, 13.2: 99-132.