

QArchSearch: A Scalable Quantum Architecture Search Package

ANKIT KULSHRESTHA, University of Delaware, USA

DANYLO LYKOV, Argonne National Laboratory, USA

ILYA SAFRO, University of Delaware, USA

YURI ALEXEEV, Argonne National Laboratory, USA

The current era of quantum computing has yielded several algorithms that promise high computational efficiency. While the algorithms are sound in theory and can provide potentially exponential speedup, there is little guidance on how to design proper quantum circuits to realize the appropriate unitary transformation to be applied to the input quantum state. In this paper, we present QArchSearch, an AI based quantum architecture search package with the QTensor library as a backend that provides a principled and automated approach to finding the best model given a task and input quantum state. We show that the search package is able to efficiently scale the search to large quantum circuits and enables the exploration of more complex models for different quantum applications. QArchSearch runs at scale and high efficiency on high-performance computing systems using a two-level parallelization scheme on both CPUs and GPUs, which has been demonstrated on the Polaris supercomputer.

ACM Reference Format:

Ankit Kulshrestha, Danylo Lykov, Ilya Safro, and Yuri Alexeev. 2023. QArchSearch: A Scalable Quantum Architecture Search Package. 1, 1 (October 2023), 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Quantum computing is a nascent and rapidly growing field that holds the promise of accomplishing tasks that were hitherto thought too be computationally intractable by classical computers. In the current era, we have access to noisy intermediate scale quantum computers (NISQ) that make it possible to run hybrid quantum algorithms to tackle problems in computational chemistry, finance [Herman et al. 2023], optimization [Ushijima-Mwesigwa et al. 2021] and related fields [Shaydulin et al. 2019]. These algorithms leverage a “variational” method in which quantum circuit parameters have to be trained through a classical optimization procedure (generally run on a classical co-processor).

In their most abstract form, variational quantum circuits are a form of linear parameterized unitary transformation $U(\theta)$ that map an input quantum state $|\psi\rangle_{in}$ to an output quantum state $|\psi\rangle_{out} = U(\theta) |\psi\rangle_{in}$. Currently, the goal in variational quantum algorithms (VQAs) is to find $\theta^* = \arg \min_{\theta} C(\theta)$ where $C(\theta)$ is some cost function that quantifies the quality of output. However, in this paper we consider an alternative problem: We aim to find the best possible circuit representing $U(\theta)$ given input $|\psi\rangle_{in}$ and cost function $C(\theta)$.

Finding an appropriate quantum circuit architecture for a given application is a computationally intensive search procedure that requires evaluation of several candidate quantum operations across different qubits and selecting the best performing circuit from amongst them. Our focus in this paper is to demonstrate how we scale such a search procedure across a state of the art HPC infrastructure to aid the search by training deep neural networks to suggest good circuit structures. We title our software as QArchSearch and include it as a part of widely available QTensor package.

Authors’ addresses: Ankit Kulshrestha, University of Delaware, Newark, DE, USA; Danylo Lykov, Argonne National Laboratory, Lemont, IL, USA; Ilya Safro, University of Delaware, Newark, DE, USA; Yuri Alexeev, Argonne National Laboratory, Lemont, IL, USA.

2023. Manuscript submitted to ACM

Manuscript submitted to ACM

1

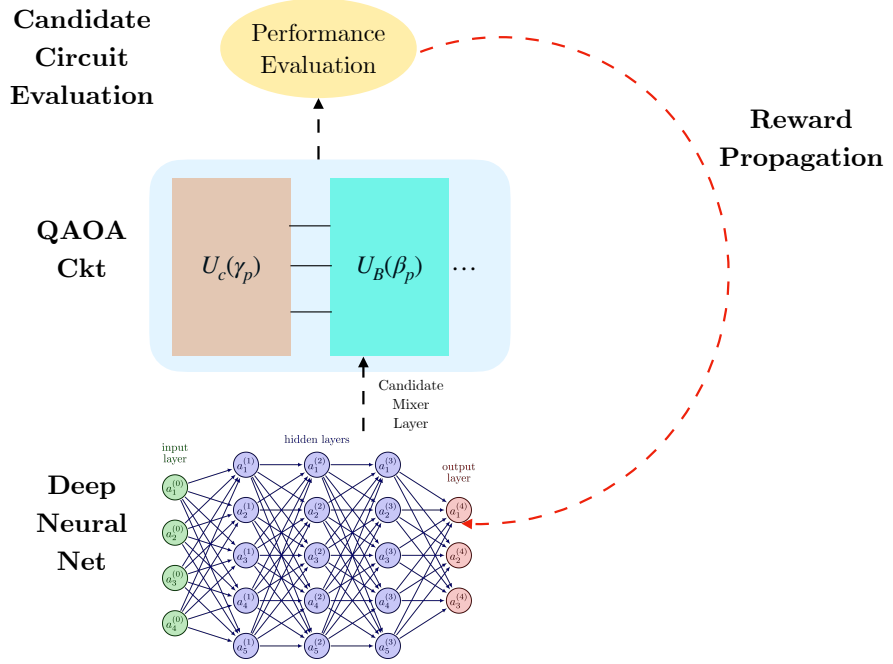


Fig. 1. An overview of the search process in QArchSearch software

In this work, we will use the Quantum Approximate Optimization Algorithm (QAOA) [Farhi et al. 2014] for the graph maxcut problem as the driver application of the QArchSearch package. Briefly, for a given simple undirected graph $G = (\mathcal{V}, \mathcal{E})$, the graph maxcut problem aims to find a maximum “cut set”, i.e., a partition of nodes into two disjoint parts to maximize the number (or the total weight in case the graph is weighted) of edges that span both parts. The cost function for max cut is given as [Farhi et al. 2014]:

$$C_{MC}(z) = \frac{1}{2} \sum_{(u,v) \in \mathcal{E}} (1 - z_u z_v), \quad (1)$$

where $z_i \in \{-1, +1\}$ is an indicator variable for node i that corresponds to the set membership the given node. In the QAOA setup, we start with an initial state $|s\rangle = |+\rangle^{\otimes n}$ where $|+\rangle = \frac{|0\rangle + |1\rangle}{\sqrt{2}}$. A p layer alternating ansatz is then applied to the initial input state:

$$|\boldsymbol{\gamma}, \boldsymbol{\beta}\rangle = e^{-i\beta_p B} e^{-i\gamma_p C} \dots e^{-i\beta_1 B} e^{-i\gamma_1 C} |s\rangle. \quad (2)$$

Here, $\boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^p$ are parameters of the cost operator C and the mixer operator B , respectively. The cost function is measured by computing $\langle \boldsymbol{\gamma}, \boldsymbol{\beta} | C(z) | \boldsymbol{\gamma}, \boldsymbol{\beta} \rangle$. In QAOA problems, the structure of the cost operator C is generally guided

by the problem we are interested in optimizing, but the structure of mixer operator is an open design problem. In our application, QArchSearch is responsible for searching low-depth mixers using the process depicted in Figure 1.

2 METHODOLOGY

2.1 QArchSearch

The QArchSearch software has three key components:

- Predictor module: This module accepts a tensor that represents the rotation gates and entanglement operators and generates a new circuit representation that is passed to the quantum builder module.
- Quantum Builder (a.k.a QBuilder): This module accepts the encoded tensor representation from the predictor module and generates the appropriate quantum circuit in an available quantum computing software. In our work, the circuits are generated using `Qiskit` software. The generated circuit is then passed to the evaluator.
- Evaluator Module: This module is responsible for training the generated quantum circuit on the QAOA cost function in Equation 1. The trained circuit is then evaluated and the reward is propagated back to the predictor module.

Algorithm 1 shows the overall search procedure for searching best performing QAOA mixer circuit from a given gate alphabet \mathcal{A}_R . The current version of the search algorithm is an instance of random search which has shown to be a strong baseline in neural architecture search [Li and Talwalkar 2020]. We perform a search by varying the depth p from 1 to desired maximum depth. For each p we explore the best possible gate combination (Line 5) and construct a mixer circuit based on the nodes in the current graph (Line 6). We then instantiate the QAOA ansatz and run the variational algorithm for 200 steps with the COBYLA optimizer. The obtained energy is added to a global collection (Line 9). At the end of exploring all possible gate combinations, we select the best performing mixer circuit and compare it to the previously existing best performing mixer circuit if it exists (Line 10). The final best performing mixer circuit and corresponding cut energy are then returned to the main calling procedure.

Algorithm 1 QArchSearch for QAOA Mixer

Require: θ : parameters of quantum circuit, \mathcal{A}_R : Gate alphabet, p_{max} : depth of QAOA ansatz, K_{max} : maximum number of possible gate combinations G : input graph

```

1:  $U_B^{best} \leftarrow \phi$ 
2: for  $p : 1 \dots p_{max}$  do
3:   energies  $\leftarrow \{\}$ 
4:   for  $k : 1 \dots K_{max}$  do
5:     gate_comb  $\leftarrow$  GET_COMBINATIONS( $\mathcal{A}_R, k$ )
6:      $\hat{U}_B \leftarrow$  BUILD_MIXER_CKT( $G, \text{gate\_comb}$ )
7:      $U_{QAOA}(\theta) \leftarrow$  BUILD_QAOA_CKT( $\hat{U}_B, p$ )
8:      $\langle C \rangle \leftarrow$  SIMULATE_QAOA( $G, U_{QAOA}(\theta)$ )
9:     energies  $\leftarrow$  APPEND(energies,  $\langle C \rangle$ )
10:   $U_B^{best} \leftarrow$  SELECT_BEST(energies,  $U_B^{best}$ )
return
11:  $U_B^{best}, \langle C_{best} \rangle$ 

```

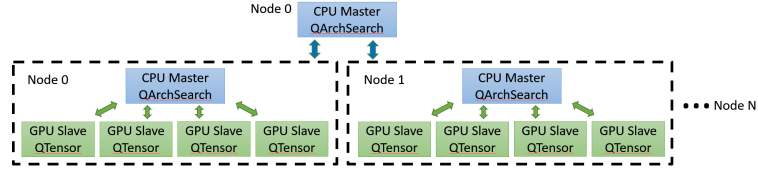


Fig. 2. The parallelization architecture within each node and between nodes using the QArchSearch and QTensor packages on the Polaris supercomputer.

2.2 QTensor

In this work, we used the Argonne-developed tensor network simulator [Lykov and Alexeev 2021; Lykov et al. 2022, 2023]. It is developed for running large-scale quantum circuit simulations using modern GPU-based supercomputers. It has been used to perform the largest QAOA simulations in the world. QTensor utilizes state-of-the-art heuristic tensor contraction order optimizers (third-party and own custom optimizers), which substantially reduce the simulation cost by minimizing the contraction width of the contraction sequence. We used a number of techniques to speed up simulations.

QTensor has support for a few tensor contraction libraries (backends) for contracting tensors efficiently. In this work, we used NumPy for tensor contraction on CPUs. The code is freely available on GitHub [QTe [n. d.]].

3 EXPERIMENTS AND RESULTS

In this section we present our results on the single and multi-core performance of QArchSearch. We then show that the circuit resulting from our search procedure generalizes to unseen graph instances and can achieve better max-cut energies even at low depths.

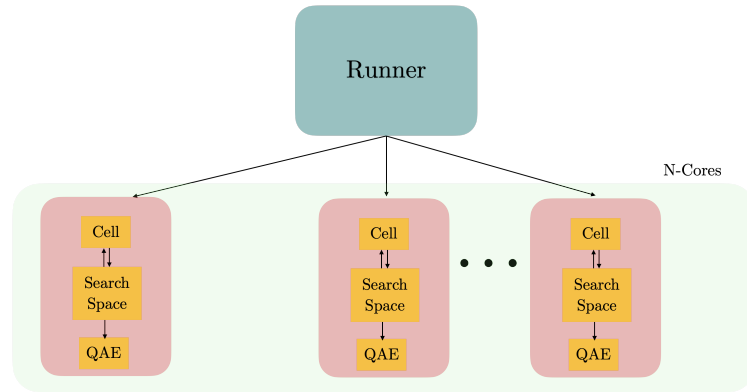


Fig. 3. The CPU level parallelism exposed by current version of QArchSearch.

3.1 Performance Profiling Results

To perform the performance profiling, we first implemented a serially executed search procedure that examined every possible rotation gate combination and simulated the resulting circuit for depths $p = 1 \dots 4$. For each depth, we

performed a combination of $k = 1 \dots 4$ gates over the given rotation gate alphabet \mathcal{A}_R with $|\mathcal{A}_R| = 5$ leading to 2500 possible circuit combinations. All search profiling was performed on a dataset of 20, 10-node Erdos-Renyi graphs with varying degrees of connectivity.

Serial Search Process: We first profiled a serial search process that sequentially examined each possible gate combination for a given depth p . The expected run time of the algorithm was thus $O(pk)$ where k was the number of gates selected from \mathcal{A}_R .

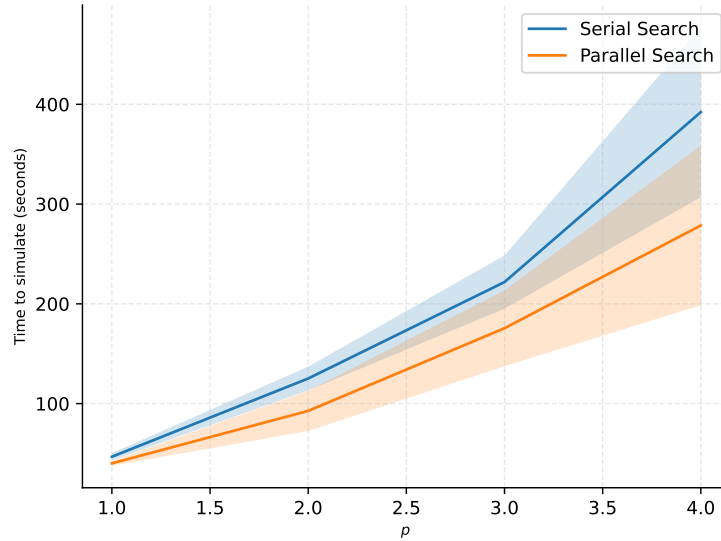


Fig. 4. Time to simulate circuits with serial and parallel quantum NAS procedure. The results are averaged over five separate runs of the NAS algorithm on different Erdos-Renyi Graphs

Parallelizing Architecture Search: To speedup the search process, it was necessary to parallelize the algorithm without causing a degradation in the quality of search results. We identified the sequential simulation of gate combinations for a given graph and depth as a major computational bottleneck. Hence, our focus was to improve run time by searching multiple possible gate combinations in parallel. This strategy is shown in Figure 3

To accomplish the aforementioned objective, we opted for process-level parallelism that can take advantage of multiple CPUs on a single node of a HPC cluster. We used Python’s multiprocessing library’s `starmap_async` method to create a pool of processes on different CPUs that executed the optimization objective of Equation 1 with different gate combinations in parallel. The run time was thus reduced from $O(pk)$ to $O(p)$ for a single graph.

Results: The results of our profiling experiments are shown in Figure 4 and Figure 5. Figure 4 shows the improvement in the run time of the algorithm with increasing depth of the QAOA ansatz. We note that in the serial case, the growth in the run time is quadratic as $p \approx k$. However, in the case of parallel the run time is improved by over 50% even when the depth approaches the maximum possible gate combinations.

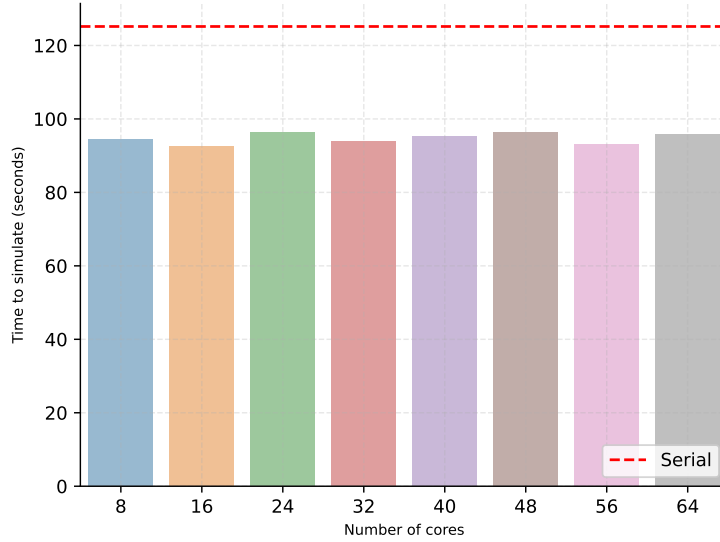


Fig. 5. Time to simulate a graph with $p = 2$ with different number of cores available on a HPC cluster. The dashed red line indicates the time to simulate the same graph with serial search.

Figure 5 shows the time to simulate for a graph with $p = 2$ and the number of available CPUs varied from 8 to 64 in increments of 8. We can see that our parallel version can efficiently utilize the available CPUs and are 0.76 times faster than a serial algorithm for the same graph and p .

3.2 Performance of Discovered Circuit

Once the search procedure was run, we evaluated the possible discovered combinations of the mixer layer on a separate dataset of 20, 10 node random 4-regular graphs. The best performing mixer circuit is shown in Figure 6. For each discovered circuit, we calculated the approximation ratio r defined as:

$$r = \frac{\langle C_{max} \rangle}{C_{classical}} \quad (3)$$

Where $\langle C_{max} \rangle$ is the expected energy of the largest cut discovered by the given quantum circuit. The approximation ratio measures the quality of solutions discovered by the quantum procedure as compared to a classical one. The results are shown in Figure 7. We can clearly see that the best performing mixer layer combination achieves the highest approximation ratio for low p value.

We further compare the performance of the searched mixer circuit on the ER and random regular graphs with the default mixer choice for maxcut QAOA. Figure 8 shows the results for average r obtained by the searched mixer and baseline mixer. These results were averaged by computing energies with $p = 1, 2, 3$. We can see that the searched mixer yields a higher average approximation ratio on ER random graphs.

In the case of random regular graph both mixer circuits perform comparably at all values of p . These results are shown in Figure 9. We show individual r since the aggregated values over p are equal (1.0).

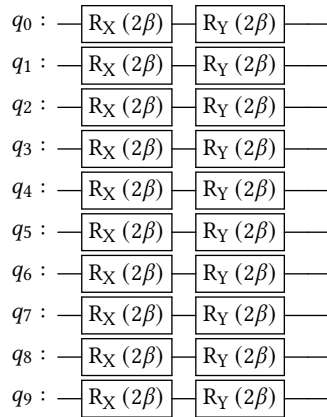


Fig. 6. Best performing searched mixer circuit for Max-cut QAOA

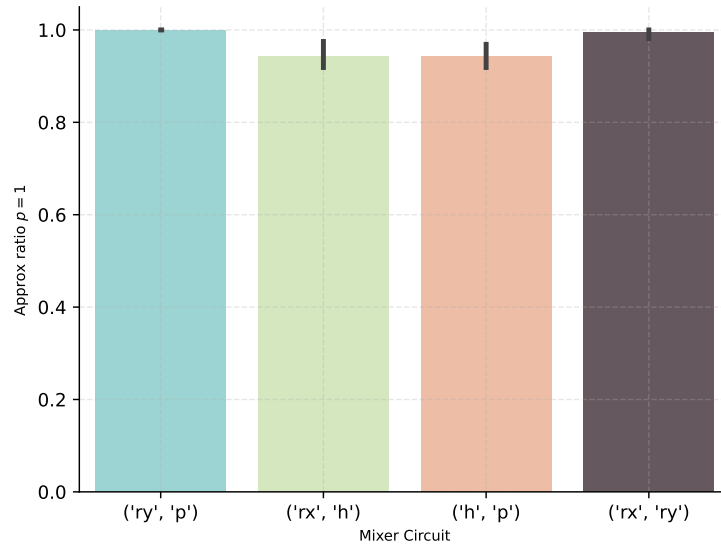


Fig. 7. Approximation ratios obtained for $p = 1$ on 4-regular random graphs. All parameterized gates in the mixer circuit share the same parameter and hence do not incur additional computational cost.

Overall, the mixer found by our search algorithm generally performs better with lower resource usage on different types of random graphs. Moreover, we show that our algorithm is able to extract the best *general* performing circuit given data and some evaluation metric.

4 KEY CHALLENGES AND FUTURE WORK

In this paper we have demonstrated that parallelizing the search algorithm is extremely important and requires careful design to obtain a meaningful speedup for large problem instances. We now discuss some important directions that are

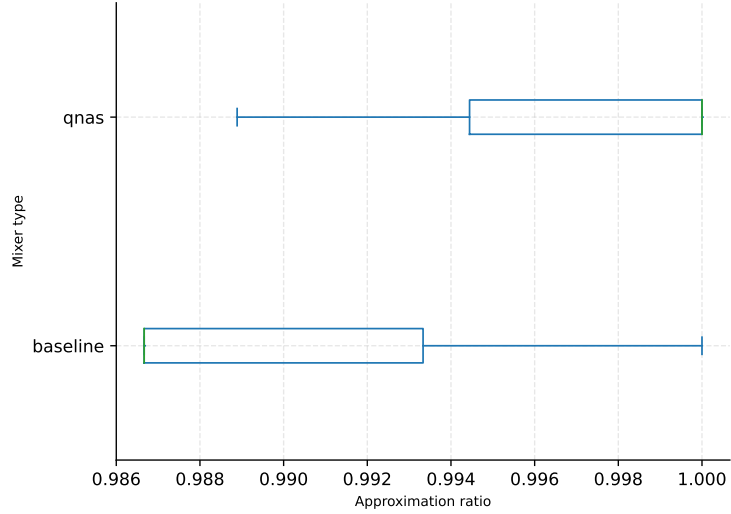


Fig. 8. Comparison of r obtained by the baseline and searched (qnas) mixer circuits

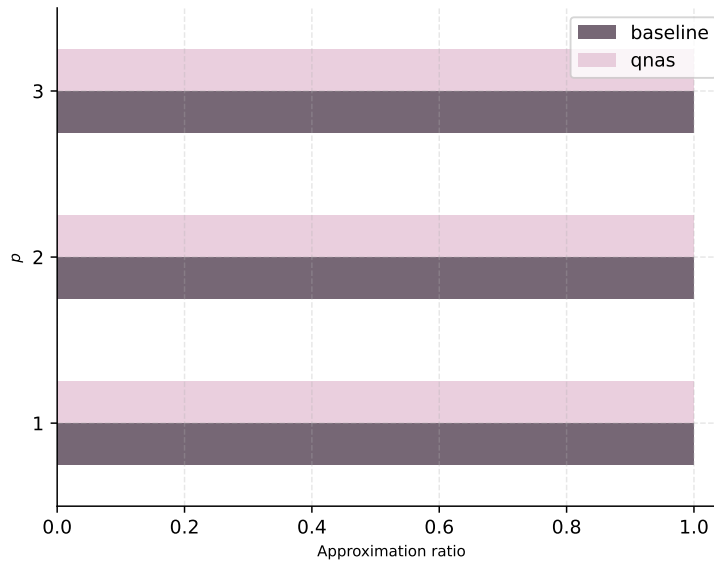


Fig. 9. Approximation ratios obtained by baseline and qnas mixer circuits on 10 node random regular graph of degree 4.

currently in development for QArchSearch.

GPU Integration: In this work, we noted that simulating quantum circuits was another computational bottleneck. However, improving the run time of quantum circuit simulations requires a different strategy since state vectors cannot be arbitrarily chunked and passed to multiple CPUs in a cluster. One possible way to improve the runtime is then to

Manuscript submitted to ACM

consider running the simulations on a GPU device. The future versions of QArchSearch will tightly integrate with QTensor to allow a user to seamlessly select a GPU backend whenever possible.

Deep Neural Network based Search: In this work we employed a version of random search to search for possible combinations of mixer circuits for the maxcut QAOA problem. Since this was a less complex problem than searching for full quantum circuits, random search returned strong generalized mixer circuits. However, our aim is to discover best quantum circuits for *any* given dataset and performance measure.

In the upcoming version of QArchSearch we will integrate several deep neural network based search algorithms like [Zhou et al. 2018; Zoph and Le 2016].

5 RELATED WORK

Quantum architecture search is a very important and active area of research and different works have considered the problem from different angles.

Fosel *et al* [Fösel et al. 2021] consider the problem of optimizing the design of quantum circuits by first proposing inefficient circuits and then training a DNN to optimize the circuit given a desired circuit metric (e.g. number of gates). Ostaszewski *et al* [Ostaszewski et al. 2021] also propose to use deep reinforcement learning (RL) to obtain a good circuit for solving VQE [Peruzzo et al. 2014] problem. Another hybrid method is considered by [Duong et al. 2022] where they propose to use Bayesian Optimization (BO) to discover optimal circuit architectures for a QNN given a particular dataset and loss function. Finally, a pure quantum architecture search is proposed by Du *et al* [Du et al. 2022] where they utilize a quantum “supercircuit” to search for child quantum circuits that satisfy a given metric. To reduce computational cost, the parameters are shared amongst all child circuits.

In the hybrid approach (i.e. using DNN to discover circuits) a major bottleneck is the sample-inefficient nature of RL algorithms. Typically, it takes days if not weeks to find good candidate architecture on a given dataset. We note that our proposed software also falls in the hybrid category of algorithms. Our objective with QArchSearch is to reduce this search time to a couple of hours. One of the reasons we do not opt for a pure quantum architecture search procedure is due to the inherent issues of scalability. For instance [Du et al. 2022] note that they are unable to search for circuits beyond 2 or 3 qubits. In order for a general architecture search package to be useful, we desire that it scale to arbitrarily as many qubits as desired.

6 CONCLUSIONS

In this work, we demonstrated the implementation of QArchSearch package that finds short-depth compact quantum circuits and architectures for a given objective function using quantum simulator QTensor. QArchSearch runs at scale and high efficiency on high-performance computing systems using two-level parallelization scheme on both CPUs and GPUs, which has been demonstrated on the 44-Petaflop supercomputer Polaris located in Argonne Leadership Computing Facility [Pol [n. d.]].

Our software satisfies a critical need in the quantum computing community - a scalable software that automates searching of candidate quantum architectures for a given problem. Our software can also incorporate arbitrary constraints in the search procedure and thus deliver custom architectures that exceed performance of manually designed ones. In order to satisfy the needs of scalability and speed, we leverage reinforcement learning techniques running

on GPUs and parallelize the search process on a large scale HPC system. Our belief is that our software will enable quantum computing researchers to find shorter-depth circuits for various advanced applications in the field.

7 ACKNOWLEDGEMENTS

This work used in part the resources of the Argonne Leadership Computing Facility, which is a Department of Energy Office of Science User Facility supported under Contract DE-AC02-06CH11357. The views, opinions and/or findings expressed are those of the authors and should not be interpreted as representing the official views or policies of the Department of Energy or the U.S. Government. This work was supported in part with funding from the Defense Advanced Research Projects Agency (DARPA).

REFERENCES

- [n. d.]. Polaris supercomputer. <https://www.alcf.anl.gov/polaris>. Accessed: 2023-09-12.
- [n. d.]. QTensor simulator on GitHub. <https://github.com/danlkv/QTensor>. Accessed: 2023-07-20.
- Yuxuan Du, Tao Huang, Shan You, Min-Hsiu Hsieh, and Dacheng Tao. 2022. Quantum circuit architecture search for variational quantum algorithms. *npj Quantum Information* 8, 1 (2022), 62.
- Trong Duong, Sang T Truong, Minh Tam, Bao Bach, Ju-Young Ryu, and June-Koo Kevin Rhee. 2022. Quantum neural architecture search with quantum circuits metric and bayesian optimization. *arXiv preprint arXiv:2206.14115* (2022).
- Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. 2014. A Quantum Approximate Optimization Algorithm Applied to a Bounded Occurrence Constraint Problem. *arXiv preprint arXiv:1412.6062* (2014). <https://arxiv.org/abs/1412.6062>
- Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. 2014. A Quantum Approximate Optimization Algorithm. *arXiv:1411.4028* (2014). <https://arxiv.org/abs/1411.4028>
- Thomas Fösel, Murphy Yuezhen Niu, Florian Marquardt, and Li Li. 2021. Quantum circuit optimization with deep reinforcement learning. *arXiv preprint arXiv:2103.07585* (2021).
- Dylan Herman, Cody Googin, Xiaoyuan Liu, Yue Sun, Alexey Galda, Ilya Safro, Marco Pistoia, and Yuri Alexeev. 2023. Quantum computing for finance. *Nature Reviews Physics* (2023), 1–16.
- Liam Li and Ameet Talwalkar. 2020. Random search and reproducibility for neural architecture search. In *Uncertainty in artificial intelligence*. PMLR, 367–377.
- Danylo Lykov and Yuri Alexeev. 2021. Importance of Diagonal Gates in Tensor Network Simulations. In *2021 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*. 447–452. <https://doi.org/10.1109/ISVLSI51109.2021.00088>
- Danylo Lykov, Roman Schutski, Alexey Galda, Valeri Vinokur, and Yuri Alexeev. 2022. Tensor Network Quantum Simulator With Step-Dependent Parallelization. In *2022 IEEE International Conference on Quantum Computing and Engineering (QCE)*. IEEE. <https://doi.org/10.1109/qce53715.2022.00081>
- Danylo Lykov, Jonathan Wurtz, Cody Poole, Mark Saffman, Tom Noel, and Yuri Alexeev. 2023. Sampling frequency thresholds for the quantum advantage of the quantum approximate optimization algorithm. *npj Quantum Information* 9, 1 (July 2023). <https://doi.org/10.1038/s41534-023-00718-4>
- Mateusz Ostaszewski, Lea M Trenkwalder, Wojciech Masarczyk, Eleanor Scerri, and Vedran Dunjko. 2021. Reinforcement learning for optimization of variational quantum circuit architectures. *Advances in Neural Information Processing Systems* 34 (2021), 18182–18194.
- A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, X.-Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O’Brien. 2014. A variational eigenvalue solver on a photonic quantum processor. *Nature Communications* 5, Article 4213 (2014), 4213 pages. <https://www.nature.com/articles/ncomms5213>
- Ruslan Shayduln, Hayato Ushijima-Mwesigwa, Ilya Safro, Susan Mniszewski, and Yuri Alexeev. 2019. Network community detection on small quantum computers. *Advanced Quantum Technologies* 2, 9 (2019), 1900029.
- Hayato Ushijima-Mwesigwa, Ruslan Shayduln, Christian FA Negre, Susan M Mniszewski, Yuri Alexeev, and Ilya Safro. 2021. Multilevel combinatorial optimization across quantum architectures. *ACM Transactions on Quantum Computing* 2, 1 (2021), 1–29.
- Yanqi Zhou, Siavash Ebrahimi, Sercan Ö Arık, Haonan Yu, Hairong Liu, and Greg Diamos. 2018. Resource-efficient neural architect. *arXiv preprint arXiv:1806.07912* (2018).
- Barret Zoph and Quoc V Le. 2016. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578* (2016).