# Calibration-Aware Transpilation for Variational Quantum Optimization

Yanjun Ji
*Institute of Computer Architecture and Computer Engineering*
*University of Stuttgart*
Stuttgart, Germany
yanjun.ji@informatik.uni-stuttgart.de

Sebastian Brandhofer
*Institute of Computer Architecture and Computer Engineering*
*University of Stuttgart*
Stuttgart, Germany
sebastian.brandhofer@informatik.uni-stuttgart.de

Ilia Polian
*Institute of Computer Architecture and Computer Engineering*
*University of Stuttgart*
Stuttgart, Germany
ilia.polian@informatik.uni-stuttgart.de

*Abstract*—Today's Noisy Intermediate-Scale Quantum (NISQ) computers support only limited sets of available quantum gates and restricted connectivity. Therefore, quantum algorithms must be transpiled in order to become executable on a given NISQ computer; transpilation is a complex and computationally heavy process. Moreover, NISQ computers are affected by noise that changes over time, and periodic calibration provides relevant error rates that should be considered during transpilation. Variational algorithms, which form one main class of computations on NISQ platforms, produce a number of similar yet not identical quantum "ansatz" circuits. In this work, we present a transpilation methodology optimized for variational algorithms under potentially changing error rates. We divide transpilation into three steps: (1) noise-unaware and computationally heavy pre-transpilation; (2) fast noise-aware matching; and (3) fast decomposition followed by heuristic optimization. For a complete run of a variational algorithm under constant error rates, only step (3) needs to be executed for each new ansatz circuit. Step (2) is required only if the error rates reported by calibration have changed significantly since the beginning of the computation. The most expensive Step (1) is executed only once for the whole run. This distribution is helpful for incremental, calibration-aware transpilation when the variational algorithm adapts its own execution to changing error rates. Experimental results on IBM's quantum computer show the low latency and robust results obtained by calibration-aware transpilation.

*Index Terms*—Calibration-Aware, Transpilation, NISQ, QAOA, Benchmarking, Quantum Computing

## I. INTRODUCTION

Quantum computing promises fundamentally more efficient solutions for a number of hard real-world problems. In the current Noisy Intermediate-Scale Quantum (NISQ) era, variational algorithms [1] such as the Quantum Approximation Optimization Algorithm (QAOA) [2]–[8] or Variational Quantum Eigensolver (VQE) [9]–[11] are receiving significant attention, since they can cope with non-trivial error rates of NISQ computers. Variational algorithms interchange classical and quantum computations. One complete run of a variational algorithm executes a number of quantum "ansatz" circuits that are parameterized, i.e., have identical basic structure but differ in some specific parameters.

State-of-the-art NISQ computers come with limitations with respect to connectivity of their qubits and quantum operations supported. Moreover, they are affected by comparatively high noise levels that can strongly vary over time. For example, computers that are part of IBM Quantum Experience (IBM QX) undergo an hourly calibration, which includes error characterization, and the determined error rates are provided to their users. To illustrate the role of calibration, Fig. 1 shows the topology graph of the 27-qubit IBM QX machine *ibmq_ehningen* along with a snapshot of error rates for its components. It includes error rates for each single-qubit gate, all two-qubit (CNOT or cx) gates between qubits connected according to topology graph, and for readout operations.

We can see that error rates differ widely both across classes of errors (e.g., they are an order of magnitude higher for readout and two-qubit operations than for single-qubit gates) and also within one class. Hourly calibration data, which we collected on *ibmq_ehningen* over a period of 39 days, are reported in Fig. 2 and expose large-scale fluctuation in the temporal domain as well. The highest variations were observed for the cx gate between qubits 8 and 9 (Fig. 2a) and for the readout errors of qubit 15 (Fig. 2c).

In general, quantum circuits, including the ansatz circuits of variational algorithms, use operations and qubit interactions
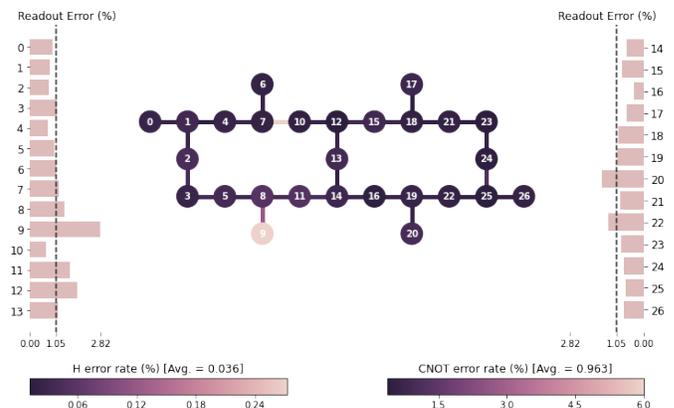


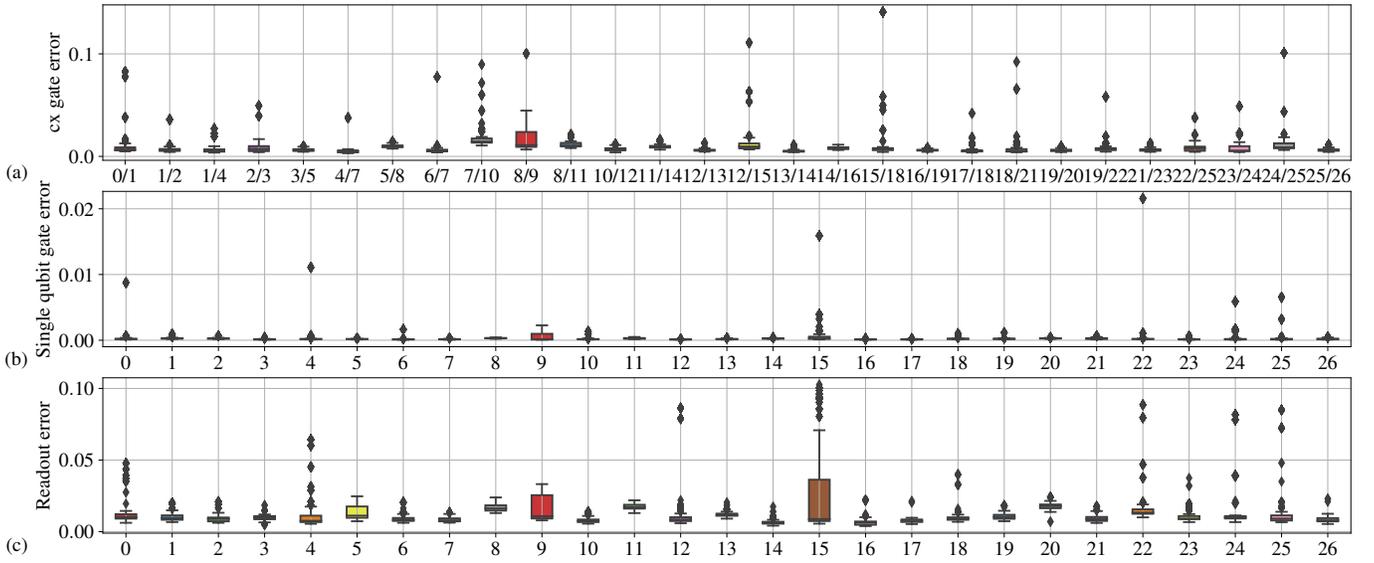Fig. 1: Topology graph of *ibmq_ehningen* with error information of single-qubit and two-qubit gates and readout.

Fig. 2: Error information of *ibmq_ehningen* over 39 days.



Fig. 3: Flowchart of Calibration-Aware (CA) Transpilation.

In this paper, we introduce calibration-aware transpilation, which is optimized for variational algorithms with parameterized ansatz circuits and avoids the need for a costly complete transpilation of each new circuit. The procedure is outlined in Fig. 3. The processed ansatz circuit $A$ is parameterized with values $\theta_i$, which, in case of QAOA, are angles determined by the classical optimization step. Transpiled circuits $A(\theta_1), A(\theta_2), \ldots$ differ only minimally, and their basic structure must be computed only once for $A$ and can be reused by all circuits. Moreover, transpilation takes error rates $\xi^*$ into account, which are determined by calibration.

After each re-calibration, the procedure checks whether new error rates $\xi^{new}$ are substantially different from $\xi^*$. If this is the case, transpilation is not repeated from scratch, but the initial, optimized solution is mapped to a different subset of qubits with the same sub-graph topology yet better fidelities. Overall, calibration-aware transpilation consists of three steps: Topology-Aware Pre-Transpilation (TAPT), executed once for the entire run of a variational algorithm and calculating parts of the solution applicable to all ansatz circuits; Noise-Aware Matching (NAM), invoked only when error rates have changed significantly; and Decomposition and Optimization (DO), which includes improvements for a specific ansatz circuit $A(\theta_i)$.

The main advantage of calibration-aware transpilation is the significantly reduced effort, as all computationally heavy parts are accumulated in the TAPT step that is run only once, and the two remaining steps are lightweight. It is feasible to use NAM and DO in an incremental mode: whenever a new ansatz circuit is ready for execution on the quantum hardware, reserve the quantum computer, acquire calibration data, execute NAM (if needed) and DO, and then immediately execute the transpiled circuit on the reserved computer. This makes sure that the most recent calibration data are used for each ansatz circuit,

that are not supported by a given NISQ architecture. For this reason, they need to be transpiled: all their operations must be mapped to that architecture's quantum gates, and two-qubit gates must either be mapped to connected qubits or proximity must be established by adding `swap` gates. In addition, transpilation should be noise-aware, that is, try to use qubits with currently lowest error rates. Various transpilation [12]–[18] and heuristic [19]–[27] and exact [20], [22], [28]–[31] mapping methods have been proposed. Both transpilation and mapping are considered to be computationally hard problems.
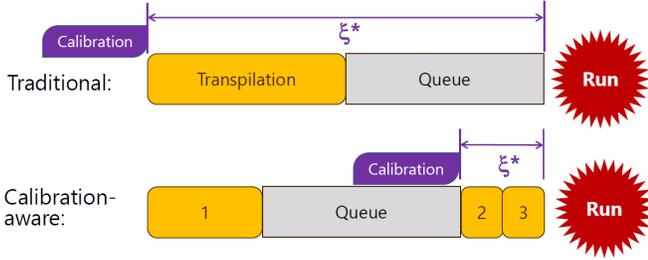
Fig. 4: Advantage for transpiling one circuit with calibration-aware transpilation compared to traditional: fresher error rate information. 1, 2 and 3 indicate TAPT, NAM and DO processes in CA transpilation, respectively.

while the amount of the quantum computer's time "wasted" during reservation is minimal. In the traditional approach with full transpilation of each ansatz circuit, reserving the quantum computer for its complete duration would be unrealistic. The computer would start processing other tasks, and the transpiled ansatz circuit, once it is ready, would be inserted into the regular queue and executed possibly at a time instant when the calibration data is outdated.

In addition, a potential advantage of transpiling one circuit with calibration-aware transpilation is shown in Fig. 4. The high-level idea is to submit the circuit directly to the queue after step 1, which is the most computationally heavy process. Steps 2 and 3, which take calibration data into account, are performed just before the circuit is due for execution, followed by immediate execution of the transpiled circuit. This ensures that CA has a fresher error rate information than the traditional approach where calibration data is acquired at the beginning of transpilation. This is crucial to the performance of algorithms on the NISQ computers as their errors change over time. Furthermore, based on CA's structures, we can significantly reduce effort and save time for a variational run with multiple circuits, see Fig. 5. While the traditional approach requires each circuit to be transpiled and passed into the queue before execution, CA only needs to pass the first circuit into the queue, and since steps 2 and 3 are fast, the remaining ansatz circuits can be run in one piece, resulting in a significant reduction in overheads.

The remainder of the paper is organized as follows. The next section reviews variational quantum algorithms with a focus on QAOA and includes some investigations of its behavior under noise using simulations. Section III provides details on the individual steps of calibration-aware transpilation. Section IV reports results of calibration-aware transpilation in comparison with other methods on several physical quantum computers, outlining its advantages in both: solution quality and runtime. Section V concludes the paper.

## II. VARIATIONAL QUANTUM OPTIMIZATION

### A. Quantum Approximation Optimization Algorithm (QAOA)

Using the quantum approximate optimization algorithm (QAOA), approximate solutions to computationally hard problems such as portfolio optimization can be computed. QAOA repeatedly performs two alternating steps. First, a set of parameters is chosen that are used to construct a parameterizable quantum circuit called ansatz circuit. The second step consists of the execution of such an ansatz circuit on a quantum computer to yield a set of measurement results that are evaluated subject to a problem-specific objective function. Then, again a set of parameters is chosen by a classical optimizer that uses the values of the previous objective functions and/or the gradient of that objective function to determine the next set of parameters. These two steps are repeated until the value of the objective function converges or the runtime budget is depleted.

We use QAOA for portfolio optimization to evaluate the quality of transpilation. Assume $n$ and $B$ are the number of available assets and the number of assets to be chosen, respectively. For each $i \in \{1, ..., n\}$, we introduce variables $z_i \in \{0, 1\}$ indicating whether this stock is picked or not. Approximation ratio (AR) is defined as:

$$\text{AR}(z_1, \ldots, z_n) = \begin{cases} \frac{F(z_1, \ldots, z_n) - F_{\max}}{F_{\min} - F_{\max}} & \text{if } \sum_i z_i = B \\ 0 & \text{if } \sum_i z_i \neq B \end{cases} \quad (1)$$

with $F$ the cost function [32]. Success probability is defined as the the probability of obtaining the optimal portfolio. We used $n = 5$, $B = 2$ with QAOA depth $p = 1$ in our experiments, i.e. the QAOA circuit has 5 qubits. Its depth is 19 and the total number of gates is 50, including 20 `cx` gates and 5 measurement gates.

The performance of QAOA depends on the initial values. Optimal parameters result in a better performance, i.e., a higher value of approximation ratio and/or success probability. The initial values of QAOA are usually obtained by classical optimizer that finds a local minimum in the area of attraction around the initial point of the probe. The optimal initial values of QAOA with $p = 1$ can be determined by grid search. Fig. 6 (a) and (b) show the optimization landscape of QAOA using qasm simulator in absence of noise and physical quantum computer *ibmq_ehningen*, respectively. The optimal initial values $\theta$, expectation values $E$, as well as approximation ratio and success probability are shown in Table I. We can see that the optimal initial values of QAOA and optimization landscape are hardly changed, i.e. this QAOA circuit is noise-tolerant and its parameters optimization is barely affected by noise.

### B. Simulation of QAOA with Noise Model

We simulate QAOA with original and optimal initial values using bit-flip, bit-phase flip and depolarizing error channels indicated by $\mathcal{E}_X$, $\mathcal{E}_Y$ and $\mathcal{E}_D$ [33] with error rate $\lambda$. The
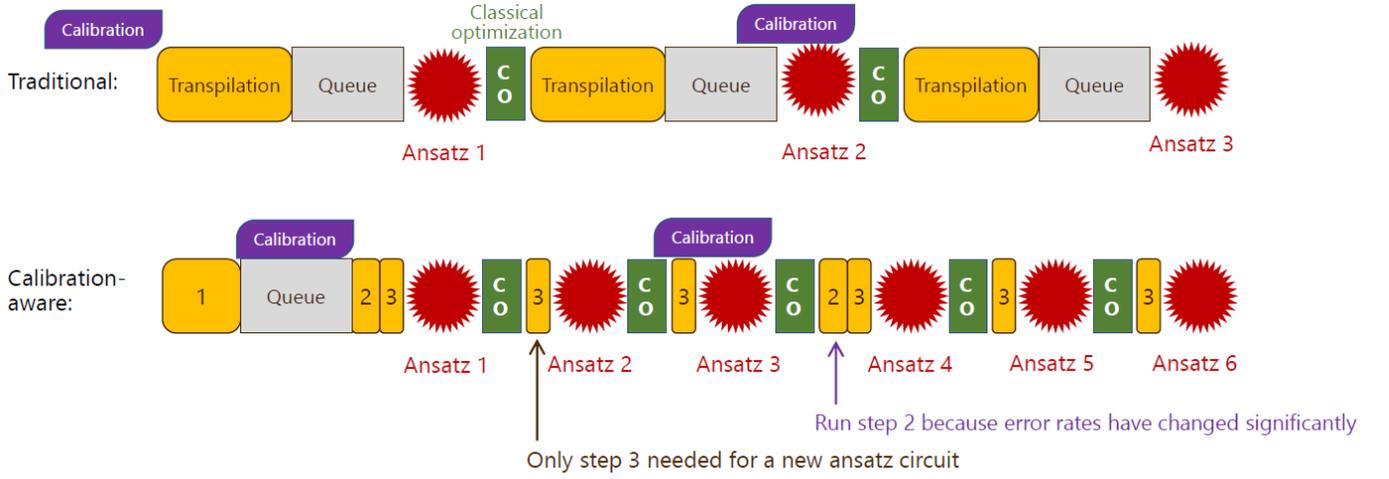
Fig. 5: Advantage of calibration-aware transpilation compared to traditional for a variational run with multiple ansatz circuits: effort reducing and time saving. Classical optimization is used to calculate new parameters of QAOA.
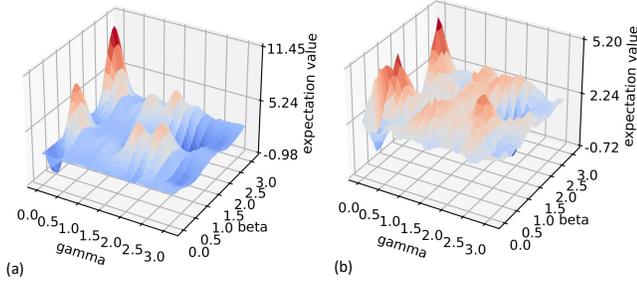


Fig. 6: Expectation values of QAOA with (a) Qasm simulator and (b) *ibmq_ehningen*.

TABLE I: Optimal solution founded by grid search with qasm simulator in absence of errors (error-free) and *ibmq_ehningen*. $\theta$: initial values. $E$: expectation values.

|  | *error-free* | *ibmq_ehningen* |
|---|---|---|
| $\theta$ | $(2.3, 2.1)$ | $(2.3, 2.3)$ |
| $E$ | $-0.957$ | $-0.726$ |
| AR | $0.417$ | $0.409$ |
| SP | $0.235$ | $0.295$ |

error channels act on qubits described by density matrix $\rho$ are defined as:

$$\mathcal{E}_X(\rho) = \lambda X \rho X + (1-\lambda)\rho$$
$$\mathcal{E}_Y(\rho) = \lambda Y \rho Y + (1-\lambda)\rho$$
$$\mathcal{E}_D(\rho) = \frac{\lambda}{4}(X\rho X + Y\rho Y + Z\rho Z) + (1 - \frac{3\lambda}{4})\rho.$$

We study state fidelity of QAOA final state with error rate up to $1\%$. Moreover, the behavior of approximation ratio and success probability with increased error rates is investigated.

The state fidelity of two quantum states is defined as

$$F(\rho_1\rho_2) = \text{Tr}\left[\sqrt{\sqrt{\rho_1}\rho_2\sqrt{\rho_1}}\right]^2 \qquad (2)$$

where $\rho_1$ and $\rho_2$ are density matrices of two quantum states. In our case, $\rho_1$ is the final QAOA state in absence of errors and $\rho_2$ is the state for QAOA with error rate $\lambda$. We consider discrete values $\lambda \in \{0, 0.1\%, 0.2\%, ..., 1\%\}$. The maximum fidelity 1 occurs at $\lambda = 0$.

In Fig. 6(a), "orig" stands for original initial values computed by the classical optimizer COBYLA and "opt" denotes optimal initial values determined by grid search. As shown in Fig. 7, optimal initial values of QAOA produce a better fidelity than original. The effect of bit-phase flips on fidelity is significant, while the fidelity under bit-flip errors varies only slightly. In addition, the type of initial values produces only a small difference under depolarizing error. At $\lambda = 1\%$, the fidelity of QAOA with bit-flip, depolarizing and bit-phase flip errors drops to about $0.95$, $0.75$ and $0.6$, respectively.

The approximation ratio and success probability of 10 QAOA runs with qasm simulator are shown in Fig. 8 (a) and (b), respectively. Without error, we achieved approximation ratios of around $0.42$ with optimal and $0.39$ with original initial values. The approximation ratio is strongly affected by bit-phase flip errors, like the fidelity in Fig. 7, and has values of around $0.33$ and $0.30$ with optimal and original initial values at $\lambda = 1\%$. QAOA under noise results in a better approximation ratio and a significantly better success probability when optimal (rather than original) initial values are used. The influence of the type of initial values on approximation ratio is smaller than success probability. We conclude from the simulation results that as the error rate increases, the fidelity decreases, leading to a lower approximation rate and a slight decrease in success probability.

## III. CALIBRATION-AWARE TRANSPILATION PROCEDURE

As has been discussed above (Fig. 3), calibration-aware (CA) transpilation is organized in three steps:

- Topology-Aware Pre-Transpilation (TAPT), which can be computationally complex and produces a high-quality (or even optimal) solution that is independent of error rates.
- Noise-Aware Matching (NAM), which takes the connectivity determined during the first step and maps it to a subgraph of the IBM QX's topology graph with the lowest error rates. This step is simple and needs to be repeated only if the error rates according to the calibration data have changed significantly.
- Decomposition and Optimization (DO), which is a collection of inexpensive procedures that take the ansatz circuit's parameters (for QAOA: angle $\theta_i$) into account. For example, certain quantum gates can be removed altogether for $\theta_i = 0$.

In the following, we provide details on the three steps.

### A. Topology-Aware Pre-Transpilation (TAPT)

Topology-aware pre-transpilation (TAPT) aims at satisfying the connectivity requirements of a quantum algorithm (here: ansatz circuit) at a given architecture. All two-qubit gates must either be mapped to connected qubits of the quantum hardware (e.g., qubits 10 and 12 in Fig. 1), or additional `swap` gates must be inserted such as to bring them onto neighboring qubits. On IBM's architecture used in this work, `swap` gates are rather expensive primitives, implemented by 3 `cx` gates. Therefore, TAPT aims at minimizing the number of required extra `swap` gates, and ultimately the total number of `cx` gates in the circuit. Note that in general, even an optimized transpiled circuit has more `cx` gates than before transpilation.

Qiskit's transpilation procedure includes randomization, and running it multiple times produces different solutions. To improve stability, we implemented an additional check that
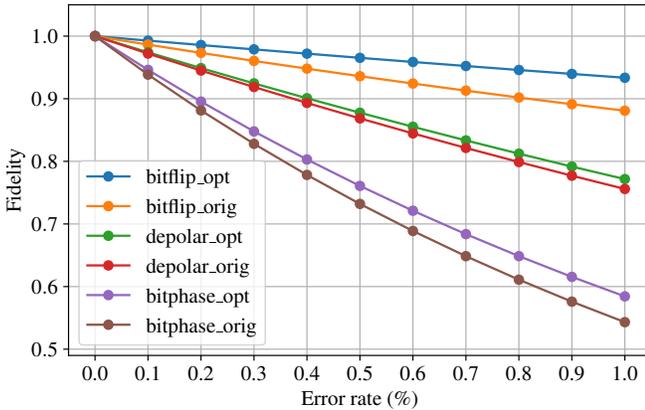


Fig. 7: Simulation of fidelity of QAOA with original (orig) and optimal (opt) initial values using qasm simulator under bitflip, depolarizing and bit-phase flip errors as a function of error rate from 0 to 1%.

bounds the maximum increase of `cx` gates $\text{CX}_{\text{inc}}^{\max}$, i.e., discards transpiled circuits with an increase of `cx` gates exceeding this threshold. Fig. 9 and Fig. 10 show the approximation ratio and, respectively, the success probability of 678 repetitions of QAOA on the physical quantum computer *ibmq_ehningen* with $\text{CX}_{\text{inc}}^{\max}$ set to $215\%, 185\%, 155\%, 140\%, 75\%$. For each such restriction, Fig. 9a and Fig. 10a show the $95\%$ confidence interval of approximation ratio and success probability, whereas Fig. 9b and Fig. 10b include the complete histograms. It can be seen that restricting the increase in `cx` gates tends to improve the performance of QAOA. Therefore, one main goal of our calibration-aware transpilation is to minimize the number of used `cx` gates (while also improving fidelity).

---

**Algorithm 1:** Topology-Aware Pre-Transpilation (TAPT)

**Input:** Original quantum circuit $qc$, Coupling map $G(V, E)$

**Output:** Topology-aware pre-transpiled circuit $pqc$

**begin**

  $M \leftarrow$ initial mapping;

  $U \leftarrow$ set of sub-circuits between two qubits in $qc$ with different structures;

  **for** $u \in U$ **do**

    transform all gates in $qc$ with the same structure as $u$ to $u$ with gate parameters;

  **end**

  Set initial mapping $M$;

  Route the re-constructed circuit by inserting `swap` gates using SMT based optimal algorithm to minimize depth;

  Decompose parameterized $u$ into basis gates of $qc$;

  **for** each `swap` gate in the circuit **do**

    **if** `swap` gate is before the measurement **then**

      remove `swap` gate and interchange the measurement of two qubits;

    **else**

      decompose `swap` into three `cx` gates and optimize with `cx` cancellation;

    **end**

  **end**

  Transform into a logic circuit by removing the idle wires;

  **return** $pqc$

**end**

---

Algorithm 1 shows the pseudocode of topology-aware pre-transpilation (TAPT). To obtain an efficient pre-transpiled circuit as a starting point, the algorithm starts with finding an initial mapping by *graph placement* [12]. This procedure identifies a sub-graph isomorphism between the graph of interacting logical qubits and the connectivity graph of the physical qubits. Before inserting `swap` gates, the circuit is partitioned into sub-circuits bounded by a `cx` gate on one or both sides. For QAOA ansatz circuits, such efficient sub-circuits have the form `cx rz cx` and implement the $Z \otimes Z$
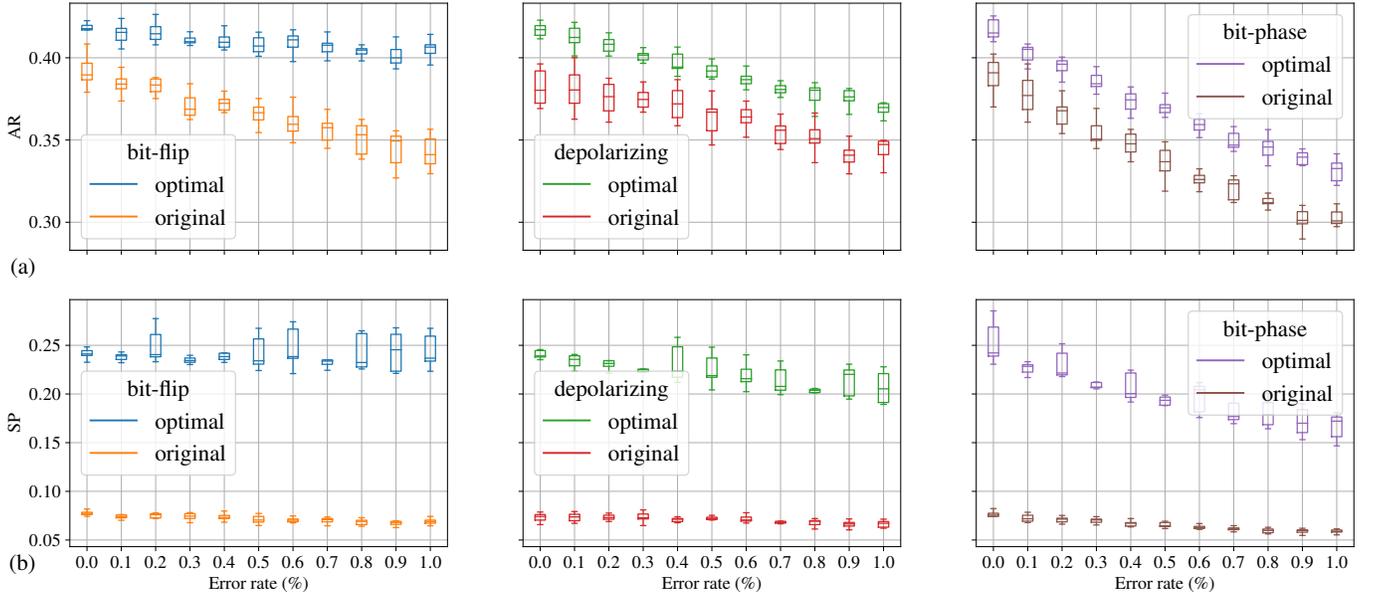
Fig. 8: Simulation of approximation ratio (a) and success probability (b) of QAOA with 10 repetitions using qasm simulator under bitflip, depolarizing and bit-phase flip errors as a function of error rate from 0 to 1%.
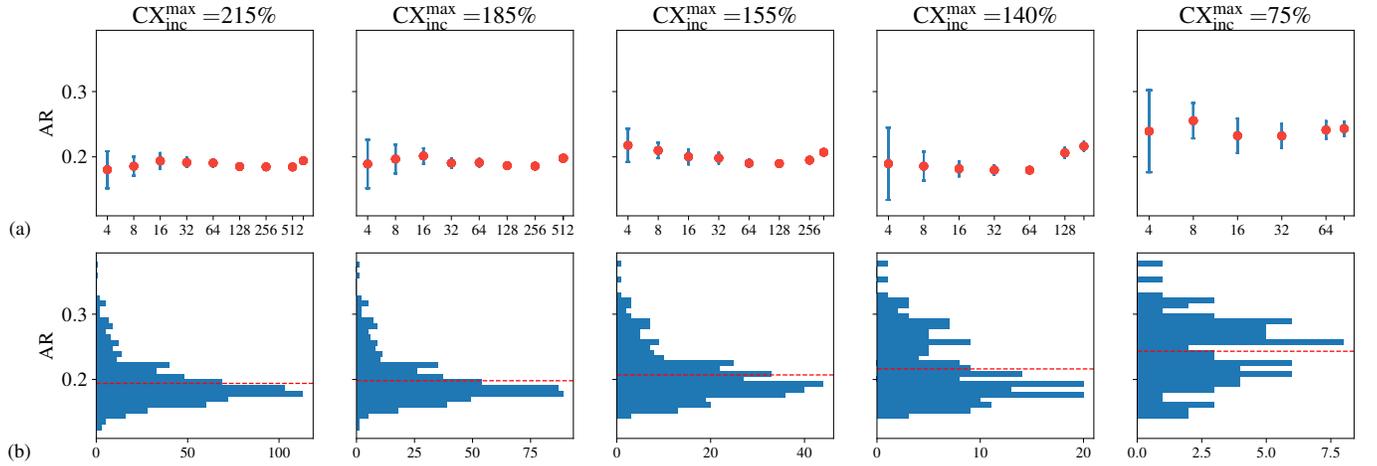


Fig. 9: 95% confidence interval (a) and histogram (b) for approximation ratio of 678 QAOA repetitions using Qiskit transpiler on *ibmq_ehningen*. With restriction of the amount of maximum increase in `cx` gates, the approximation ratio is improved.

interaction between two qubits.

To insert `swap` gates, we use an optimal method based on SMT (satisfiability modulo theory) [31] with circuit depth as the optimization objective to guarantee a high quality of pre-transpilation, as its runtime does not influence the performance of total process. The SMT method treats the sub-circuits identified as outlined above as primitive circuit elements. That is, `swap` gates are inserted only between sub-circuits. The rationale behind this procedure is the later application of `cx` cancellation, where `cx` gates on the boundaries of sub-circuits can be merged with `cx` gates that implement `swap` gates. In addition, considering sub-circuits reduces the problem complexity and the run-time of the SMT solver. Thereafter, the

sub-circuits and the inserted `swap` gates are decomposed into basis gates of the circuit and undergo optimization, including `cx` cancellation. The resulting circuit is depth-optimized and executable on (a sub-graph of) the target topology graph.

### B. Noise-Aware Matching (NAM)

The previous step maps the algorithm's qubits to a sub-graph of the topology graph, but it does not consider error rates of physical qubits in that sub-graph. At the same time, most topology graphs of today's larger-scale quantum computers have a large number of isomorphies and symmetries. For example, the topology graph from Fig. 1 can be understood as consisting of two "tiles" (physical qubits $0, \ldots, 14, 16$ and physical qubits $10, 12, \ldots, 26$). Any algorithm mapped to a
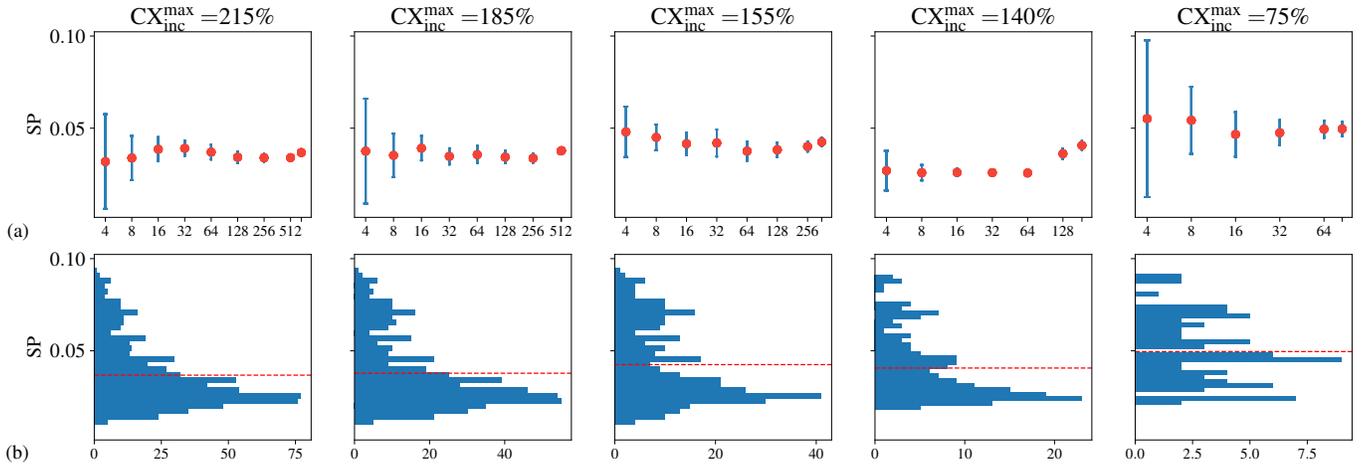
Fig. 10: 95% confidence interval (a) and histogram (b) for success probability of 678 QAOA repetitions using Qiskit transpiler on *ibmq_ehningen*. Success probability shows a similar behavior to approximation ratio in Fig. 9.

sub-graph from one "tile" can also run on the other. Moreover, each "tile" is symmetric. For instance, a 5-qubit algorithm mapped to physical qubits (0, 1, 4, 7, 6) is automatically valid for physical qubits (10, 12, 15, 18, 17); (15, 12, 10, 7, 6); (16, 14, 11, 8, 9); (26, 25, 22, 19, 20). Note that larger IBM computers have even more identical "tiles" and offer more valid alternatives for each result of step TAPT.

Noise-Aware Matching (NAM) considers up to $N$ alternative sub-graphs and selects the one with the highest effective average fidelity. $N$ is a user-defined constant, which trades the likelihood of finding a good matching against the number of necessary computations; the latter can be important when calibration-aware transpilation is used in the incremental mode and the quantum computer waits until the new sub-graph is identified. The effective average fidelity of quantum circuit $qc$ with the calibration data $\xi = (f_u, f_{cx}, f_d)$ as:

$$get\_fidelity(qc, \xi) = \frac{1}{3}(\prod_{u \in qc} f_u + \prod_{cx \in qc} f_{cx} + \prod_{d \in qc} f_d) \quad (3)$$

where $f_u$, $f_{cx}$ and $f_d$ are fidelities of single qubit, cx gate and readout gate, respectively.

The noise-aware matching (NAM) is described in Algorithm 2. As input we have pre-transpiled circuit, which has satisfied the connectivity of sub-graph of topology graph, coupling map, the latest calibration data, the function *get_fidelity* to calculate the fidelity of circuit, and the number of trials $N$ (we used $N = 15$ in our experiments). In order to select high-fidelity qubits, we perform $N$ trial matchings with Qiskit's transpile. With the randomization of Qiskit's transpilation procedure, the pre-transpiled circuit is matched to different physical qubits of quantum computer. Effective average fidelities of the matched circuits on physical qubits are calculated and the circuit with the highest fidelity is picked. This circuit contains the information of selected physical qubits in $N$ trials.

---

**Algorithm 2:** Noise-Aware Matching (NAM)

**Input:** Topology-aware pre-transpiled circuit $pqc$,
   Coupling map $G(V, E)$, Calibration data $\xi$,
   Fidelity computation function *get_fidelity*,
   Number of trials $N$
**Output:** Selected physical qubits
**begin**
   $mqc \leftarrow$ matched $pqc$ with Qiskit;
   $c_m \leftarrow$ number of cx gates in $mqc$;
   $f_m \leftarrow get\_fidelity(mqc, \xi)$;
   $j \leftarrow 0$;
   **while** $j \neq N$ **do**
      $rqc \leftarrow$ re-matched $pqc$ with Qiskit;
      $c_r \leftarrow$ number of cx gates in $rqc$;
      $f_r \leftarrow get\_fidelity(rqc, \xi)$;
      **if** $c_r \leq c_m$ *and* $f_r > f_m$ **then**
         | $mqc \leftarrow rqc$
      **end**
      $j \leftarrow j + 1$;
   **end**
   **return** Physical qubits used in $mqc$
**end**

---

*C. Decomposition and Optimization (DO)*

After the NAM process, the target qubits used to run the algorithm are fixed. Then the quantum algorithm is decomposed into the native gate set supported by IBM QX. In this work, we are using IBM's computers that support single-qubit gates and the cx gate as an entangling gate. After decomposition, we apply optimization techniques *Optimize1qGates*, *CommutationAnalysis*, *CommutativeCancellation*, *CXCancellation*, *RemoveDiagonalGatesBeforeMeasure* and *RemoveResetInZeroState* provided by Qiskit. We repeat this process 15 times aiming to obtain the final transpiled circuit with the least number of cx gates. Note that this step is architecture-specific
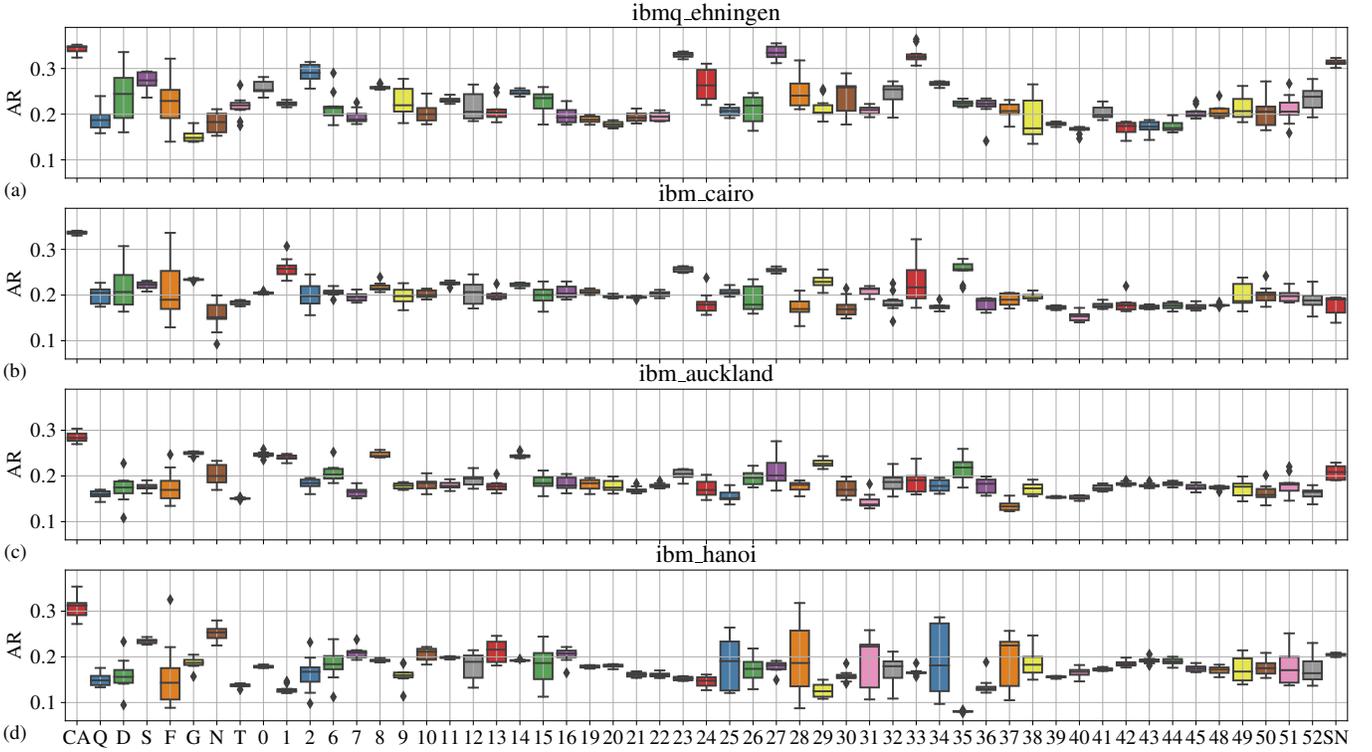
Fig. 11: Approximation ratio of 10 QAOA runs with different transpilation methods on four IBM QX computers. CA: Calibration-aware (this paper); Q: Qiskit's built-in transpilation procedure with *optimization_level* 3; D/S/F: SMT-based methods [31] to minimize depth (D), minimize the number of `swap` gates (S), to maximize fidelity (F); G: t|ket⟩ with initial mapping based on *graph placement* [12]; N: t|ket⟩ with *noise aware placement* [12]; T: staq [13]; 00...52: methods composed of different initial mapping and routing procedures, some including ZX-calculus optimization [34]; SN: swap network [14]. 10 runs of QAOA per data point.

and would need to be adapted for a different platform; for example, Google's computers use `cz` gates as entangled gates.

## IV. BENCHMARKING WITH QAOA

In this section, we benchmark the calibration-aware (CA) transpilation with QAOA on four IBM QX computers, *ibmq_ehningen*, *ibm_cairo*, *ibm_auckland* and *ibm_hanoi*, all of which have 27 qubits and the same topology graph, as shown in Fig. 1. The initial values of QAOA used here are determined by COBYLA, i.e. the original values labeled in Fig. 8. We use three sets of metrics: the quality of the transpilation process (quantified by percental increase of the circuit's depth $\Delta d\%$, its total number of gates $\Delta g\%$ and number of its `cx` gates $\Delta g_{cx}\%$ as a result of transpilation); the runtime of the transpilation procedure; and the quality of QAOA in terms of approximation ratio and success probability achieved on a physical quantum computer.

We start with a comparison of CA with a large number of different transpilation methods with respect to approximation ratio and success probability. Then, we pick one of the best methods observed and compare it with CA in-depth. Finally, we outline the potential benefits of CA for runtime of a complete QAOA algorithm.

### A. Approximation Ratio, Success Probability, cx Gate Count

We implemented a number of transpilation methods and compared the achieved approximation ratio in Fig. 11. In addition to CA, this figure includes Qiskit's built-in transpilation procedure, SMT-based methods [31], two variants of t|ket⟩ [12], staq [13], a total of 53 composite methods depending on different initial mapping and routing procedures, and swap network (SN) [14]. CA either outperforms other methods or is on a par with the best of them for all four quantum computers.

We believe that this advantage is due to CA's NAM step considering a number of possible sub-graphs, selecting the one with the best fidelity according to a more up-to-data calibration data than other methods. At the same time, the TAPT step produces a robust depth-optimized solution for the connectivity constraints, thus limiting errors that stem from excessive gate-count. We observed that purely fidelity-oriented transpilation can incur strong variability in gate count for different ansatz circuits $A(\theta_1), A(\theta_2), \ldots$; CA's DO step is applying only minimal modifications to the basic solution from the TAPT step, thus leading to well-aligned transpilation results for different circuits.

The percental increase in `cx` gates after transpilation is reported in Fig. 12. The numbers differ only minimally among
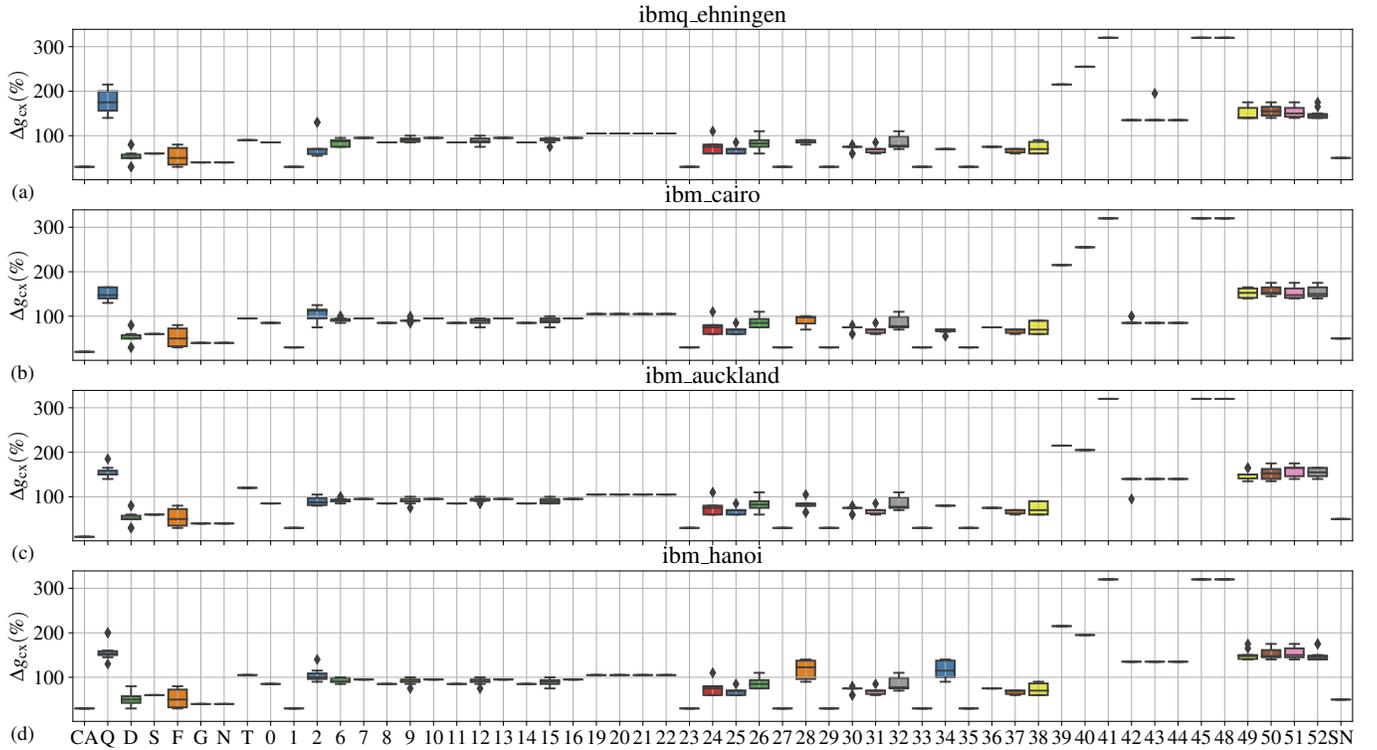
Fig. 12: Percental increase in the number of `cx` gates after transpilation (10 QAOA runs, same methods as in Fig. 11).

TABLE II: Comparison of percentage increase in depth $\Delta d\%$, the total number of gates $\Delta g\%$ and the number of `cx` gates $\Delta g_{\text{cx}}\%$ after transpilation with CA and SF. $\mu$: average value. $\sigma$: standard deviation.

| | $\Delta d\%$ | | | | $\Delta g\%$ | | | | $\Delta g_{\mathbf{cx}}\%$ | | | |
| | CA | | SF | | CA | | SF | | CA | | SF | |
| **IBM QX** | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *ibmq_ehningen* | 126.32 | 0.0 | 116.84 | 31.47 | 170.00 | 0.0 | 104.00 | 66.78 | 30.00 | 0.0 | 42.00 | 17.20 |
| *ibm_cairo* | 126.32 | 0.0 | 110.00 | 27.49 | 170.00 | 0.0 | 93.20 | 63.03 | 30.00 | 0.0 | 48.00 | 18.87 |
| *ibm_auckland* | 126.32 | 0.0 | 116.84 | 31.47 | 170.00 | 0.0 | 104.00 | 66.78 | 30.00 | 0.0 | 42.00 | 17.20 |
| *ibm_hanoi* | 126.32 | 0.0 | 116.84 | 31.47 | 170.00 | 0.0 | 104.00 | 66.78 | 30.00 | 0.0 | 42.00 | 17.20 |

the four quantum computers. Again, CA is consistently best or among the best methods with respect to this metric, while other methods produce an increase of up to 330% (more than four times) in `cx` gates. CA's outcome is also much more stable, since all its ansatz circuits are based on the same high-quality basic solution provided by the TAPT step, whereas, e.g., Qiskit (Q) produces quite different transpilation results for each QAOA run and quantum computer.

Table II shows a detailed comparison of CA with SF in terms of increase in depth $\Delta d\%$, total number of gates $\Delta g\%$ and number of `cx` gates $\Delta g_{\text{cx}}\%$. The table shows that SF is a good transpilation method with a slightly better $\Delta d\%$, significantly better $\Delta g\%$, but significantly worse $\Delta g_{\text{cx}}\%$. SF exposes large-scale variability whereas CA's results are repeatable.

To assess the significance of calibration data and the NAM step, we report in Fig. 13 the approximation ratio and the success probability of three methods: topology-aware transpilation

TA, which is CA that stopped after the TAPT step and did not incorporate any calibration data; SMT-based transpilation SF [31] that maximizes the circuit's fidelity and is used for reference; and the full CA procedure with all its three steps. CA by far outperforms TA and is also consistently better than SF, while both TA and CA are less affected by the variability of the obtained results. We conclude that all three steps of CA are needed to obtain a high-quality solution.

### B. Runtime

One central objective of our calibration-aware (CA) transpilation approach is to reduce the runtime of iterative variational algorithms. In this section, we compare CA with SMT-based method in [31] that maximizes the circuit's fidelity (SF). Table III reports the average runtimes of CA's three steps, their sum ($\mu$) and standard deviation ($\sigma$), along with the average runtime and standard deviation for SF, which is not partitioned into steps. It can be seen that the overall runtimes of CA and SF are comparable, suggesting a similar amount of computational
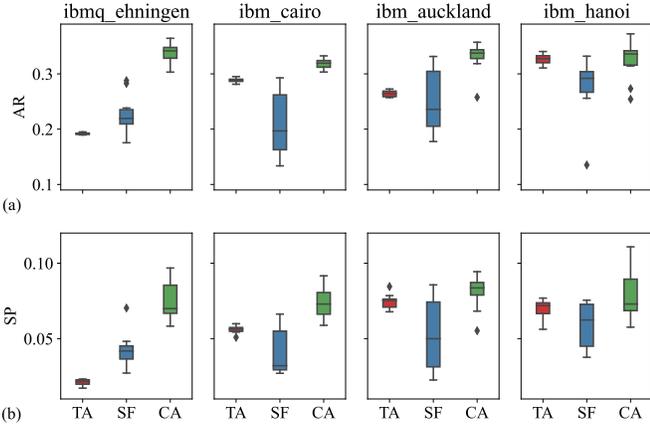
Fig. 13: Approximation ratio and success probability with IBM QX computers. TA: topology-aware (noise-unaware) transpilation (TAPT step of CA); SF: SMT-based fidelity maximization; CA: full calibration-aware method (TAPT, NAM, DO).

TABLE III: Comparison of runtime of transpilation with CA and SF for 10 QAOA runs (in seconds). $\mu$: average runtime. $\sigma$: standard deviation.

| | CA | | | | | SF | |
|---|---|---|---|---|---|---|---|
| IBM QX | TAPT | NAM | DO | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ |
| ibmq_ehningen | 18.13 | 4.41 | 2.26 | 24.80 | 0.22 | 30.00 | 3.65 |
| ibm_cairo | 23.07 | 3.59 | 3.13 | 29.79 | 0.26 | 25.90 | 1.74 |
| ibm_auckland | 22.92 | 3.44 | 2.31 | 28.67 | 0.14 | 27.94 | 5.46 |
| ibm_hanoi | 17.44 | 4.37 | 2.50 | 24.31 | 0.31 | 29.89 | 3.20 |

effort being invested. However, CA manages to serialize its computation without a major deterioration of overall runtime. Furthermore, the standard deviation of total runtime with CA is an order of magnitude smaller than with SF, which means less latency due to transpilation.

To illustrate the runtime advantage enabled by CA, assume we have $N_A$ QAOA ansatz circuits to execute. The expectation value of the total runtime for $N_A$ circuits with SF is

$$\mu_{\text{SF}}(N_A) = \mu_{\text{SF}} \times N_A \tag{4}$$

with $\mu_{\text{SF}}$ being the average runtime for one circuit. If the calibration data changes every $m$ iterations, the expected total runtime with CA is

$$\mu_{\text{CA}}(N_A) = \mu_{\text{TAPT}} + \mu_{\text{NAM}} \times \left\lceil \frac{N_A}{m} \right\rceil + \mu_{\text{DO}} \times N_A \tag{5}$$

where $\mu_{\text{TAPT}}$, $\mu_{\text{NAM}}$ and $\mu_{\text{DO}}$ are average runtimes for TAPT, NAM and DO process, respectively. For the special case that calibration data is fixed, we have $m = N_A$ and only the last step DO needs to be executed each time.

The projected runtimes of a complete QAOA run with $N_A \in \{5, 100\}$ ansatz circuits are shown in Table IV. The data assumes two scenarios: changing error rate (CER), where the calibration data changes after 5 iterations ($m = 5$), and fixed error rate (FER), where the calibration data remains unchanged ($m = N_A$). We see an improvement of up to one order of

TABLE IV: Comparison of runtime (in seconds) of CA and SF for $N_A = 5, 100$. CER: Changing error rate after $m = 5$ iterations. FER: Fixed error rate ($m = N_A$). $\Delta\mu$: average time savings compared to SF.

| $N_A$ | IBM QX | CA | | SF | $\Delta\mu(\%)$ | |
|---|---|---|---|---|---|---|
| | | CER | FER | | CER | FER |
| 5 | ibmq_ehningen | 33.83 | 33.83 | 149.98 | -77.44 | -77.44 |
| | ibm_auckland | 37.91 | 37.91 | 139.71 | -72.87 | -72.87 |
| | ibm_cairo | 42.33 | 42.33 | 129.48 | -67.31 | -67.31 |
| | ibm_hanoi | 34.32 | 34.32 | 149.43 | -77.03 | -77.03 |
| 100 | ibmq_ehningen | 332.23 | 248.46 | 2999.66 | -88.92 | -91.72 |
| | ibm_auckland | 322.67 | 257.31 | 2794.29 | -88.45 | -90.79 |
| | ibm_cairo | 408.21 | 339.99 | 2589.61 | -84.24 | -86.87 |
| | ibm_hanoi | 355.05 | 272.03 | 2988.57 | -88.12 | -90.90 |

magnitude due to CA's three-step structure where the most expensive part of the calculation is executed only once.

### C. Conclusion and Comparison

For evaluation, we compared CA approach with several other methods. The experiments show that applying our approach yields better and more stable results. The main advantage of CA is that TAPT needs to be performed only once and a number of trials in the fast NAM processing qualify that the implementation is performed on qubits with high fidelity. Another highlight is that the run time for $N_A$ ansatz circuits is significantly reduced: up to 88.92% with CER and 91.72% with FER for $N_A = 100$. Moreover, with CA we have stable properties for transpiled circuit, which are shown by the increase of depth, number of gates and number of cx gates. All this guarantees that CA produces consistently high quality on four IBM quantum computers.

## V. CONCLUSIONS AND FUTURE WORK

The decisive role of variational algorithms during the NISQ era justifies a specialized transpilation approach for such algorithms. Calibration-aware transpilation leverages the knowledge that subsequent ansatz circuits have the same basic structure and differ only in their parameters. It naturally adapts itself to abrupt changes in the error rates of the quantum computer executing the algorithm, which is a reality today. Our results show that calibration-aware transpilation strikes a good balance between quality and stability of transpilation results and saves time thanks to offloading the heaviest computation to a procedure that is run one time for all ansatz circuits.

Our findings are confirmed by results for QAOA obtained on four physical quantum computers with a similar architecture. They are compared with an extensive set of previous transpilation procedures executed on the same computers and are put into perspective with simulations assuming standard error models. We believe that calibration-aware transpilation is particularly attractive for today's quantum cloud computers with their potentially long queuing times: if a circuit is executed long after it has been transpiled, its calibration data can become outdated and the actual error rate can get much worse than expected. Calibration-aware transpilation enables

incremental operation, where the execution starts almost instantly after calibration, with only lightweight parts of the transpilation procedure being completed in between.

For the future, we are interested in further improving the algorithm's performance, especially for emerging NISQ computers with 100s or 1000s qubits. Here, quick variants of the NAM step that leverage the architecture's symmetries are essential. Moreover, the TAPT step must be further evaluated for stability for other variational algorithms and quantum architectures. Another interesting question is whether we can make calibration-aware transpilation provably optimal with respect to one of the targets, despite being divided into three independent steps. This would make optimal approaches available for time-critical quantum circuit execution, as the expensive TAPT step can be performed on a classical computer before any access to a quantum computer.

## REFERENCES

[1] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio, and P. J. Coles, "Variational quantum algorithms," *CoRR*, vol. abs/2012.09265, 2020. [Online]. Available: https://arxiv.org/abs/2012.09265

[2] E. Farhi, J. Goldstone, and S. Gutmann, "A quantum approximate optimization algorithm," *arXiv preprint arXiv:1411.4028*, 2014.

[3] E. Farhi and A. W. Harrow, "Quantum supremacy through the quantum approximate optimization algorithm," *arXiv preprint arXiv:1602.07674*, 2016.

[4] J. Choi and J. Kim, "A tutorial on quantum approximate optimization algorithm (qaoa): Fundamentals and applications," in *2019 International Conference on Information and Communication Technology Convergence (ICTC)*, 2019, pp. 138–142.

[5] L. Zhu, H. L. Tang, G. S. Barron, F. Calderon-Vargas, N. J. Mayhall, E. Barnes, and S. E. Economou, "An adaptive quantum approximate optimization algorithm for solving combinatorial problems on a quantum computer," *arXiv preprint arXiv:2005.10258*, 2020.

[6] L. Zhou, S.-T. Wang, S. Choi, H. Pichler, and M. D. Lukin, "Quantum approximate optimization algorithm: Performance, mechanism, and implementation on near-term devices," *Physical Review X*, vol. 10, no. 2, p. 021067, 2020.

[7] S. Hadfield, Z. Wang, B. O'Gorman, E. G. Rieffel, D. Venturelli, and R. Biswas, "From the quantum approximate optimization algorithm to a quantum alternating operator ansatz," *Algorithms*, vol. 12, no. 2, p. 34, 2019. [Online]. Available: https://doi.org/10.3390/a12020034

[8] M. Fernández-Pendás, E. F. Combarro, S. Vallecorsa, J. Ranilla, and I. F. Rúa, "A study of the performance of classical minimizers in the quantum approximate optimization algorithm," *J. Comput. Appl. Math.*, vol. 404, p. 113388, 2022. [Online]. Available: https://doi.org/10.1016/j.cam.2021.113388

[9] A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, X.-Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O'brien, "A variational eigenvalue solver on a photonic quantum processor," *Nature communications*, vol. 5, no. 1, pp. 1–7, 2014.

[10] J. R. McClean, J. Romero, R. Babbush, and A. Aspuru-Guzik, "The theory of variational hybrid quantum-classical algorithms," *New Journal of Physics*, vol. 18, no. 2, p. 023023, 2016.

[11] J. Romero, R. Babbush, J. R. McClean, C. Hempel, P. J. Love, and A. Aspuru-Guzik, "Strategies for quantum computing molecular energies using the unitary coupled cluster ansatz," *Quantum Science and Technology*, vol. 4, no. 1, p. 014008, 2018.

[12] S. Sivarajah, S. Dilkes, A. Cowtan, W. Simmons, A. Edgington, and R. Duncan, "t—ket⟩: a retargetable compiler for nisq devices," *Quantum Science and Technology*, 2020.

[13] M. Amy and V. Gheorghiu, "staq—a full-stack quantum processing toolkit," *Quantum Science and Technology*, vol. 5, no. 3, p. 034016, 2020.

[14] M. P. Harrigan, K. J. Sung, M. Neeley, K. J. Satzinger, F. Arute, K. Arya, J. Atalaya, J. C. Bardin, R. Barends, S. Boixo *et al.*, "Quantum approximate optimization of non-planar graph problems on a planar superconducting processor," *Nature Physics*, vol. 17, no. 3, pp. 332–336, 2021.

[15] D. S. Steiger, T. Häner, and M. Troyer, "Projectq: an open source software framework for quantum computing," *Quantum*, vol. 2, p. 49, 2018.

[16] P. Murali, N. M. Linke, M. Martonosi, A. J. Abhari, N. H. Nguyen, and C. H. Alderete, "Full-stack, real-system quantum computer studies: Architectural comparisons and design insights," in *2019 ACM/IEEE 46th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2019, pp. 527–540.

[17] N. Khammassi, I. Ashraf, J. van Someren, R. Nane, A. M. Krol, M. A. Rol, L. Lao, K. Bertels, and C. G. Almudéver, "Openql: A portable quantum programming framework for quantum accelerators," *ACM J. Emerg. Technol. Comput. Syst.*, vol. 18, no. 1, pp. 13:1–13:24, 2022. [Online]. Available: https://doi.org/10.1145/3474222

[18] V. Bergholm, J. A. Izaac, M. Schuld, C. Gogolin, and N. Killoran, "Pennylane: Automatic differentiation of hybrid quantum-classical computations," *CoRR*, vol. abs/1811.04968, 2018. [Online]. Available: http://arxiv.org/abs/1811.04968

[19] M. Y. Siraichi, V. F. dos Santos, C. Collange, and F. M. Q. Pereira, "Qubit allocation as a combination of subgraph isomorphism and token swapping," *Proc. ACM Program. Lang.*, vol. 3, no. OOPSLA, pp. 120:1–120:29, 2019. [Online]. Available: https://doi.org/10.1145/3360546

[20] M. Y. Siraichi, V. F. dos Santos, S. Collange, and F. M. Q. Pereira, "Qubit allocation," in *Proceedings of the 2018 International Symposium on Code Generation and Optimization, CGO 2018, Vösendorf / Vienna, Austria, February 24-28, 2018*, J. Knoop, M. Schordan, T. Johnson, and M. F. P. O'Boyle, Eds. ACM, 2018, pp. 113–125. [Online]. Available: https://doi.org/10.1145/3168822

[21] A. Zulehner, A. Paler, and R. Wille, "An efficient methodology for mapping quantum circuits to the IBM QX architectures," *IEEE Trans. Comput. Aided Des. Integr. Circuits Syst.*, vol. 38, no. 7, pp. 1226–1236, 2019. [Online]. Available: https://doi.org/10.1109/TCAD.2018.2846658

[22] P. Murali, J. M. Baker, A. Javadi-Abhari, F. T. Chong, and M. Martonosi, "Noise-adaptive compiler mappings for noisy intermediate-scale quantum computers," *CoRR*, vol. abs/1901.11054, 2019. [Online]. Available: http://arxiv.org/abs/1901.11054

[23] S. S. Tannu and M. K. Qureshi, "Not all qubits are created equal: A case for variability-aware policies for nisq-era quantum computers," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2019, Providence, RI, USA, April 13-17, 2019*, I. Bahar, M. Herlihy, E. Witchel, and A. R. Lebeck, Eds. ACM, 2019, pp. 987–999. [Online]. Available: https://doi.org/10.1145/3297858.3304007

[24] G. Li, Y. Ding, and Y. Xie, "Tackling the qubit mapping problem for nisq-era quantum devices," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2019, Providence, RI, USA, April 13-17, 2019*, I. Bahar, M. Herlihy, E. Witchel, and A. R. Lebeck, Eds. ACM, 2019, pp. 1001–1014. [Online]. Available: https://doi.org/10.1145/3297858.3304023

[25] A. M. Childs, E. Schoute, and C. M. Unsal, "Circuit transformations for quantum architectures," *CoRR*, vol. abs/1902.09102, 2019. [Online]. Available: http://arxiv.org/abs/1902.09102

[26] Z. Li, F. Meng, Z. Zhang, and X. Yu, "Qubits' mapping and routing for NISQ on variability of quantum gates," *Quantum Inf. Process.*, vol. 19, no. 10, p. 378, 2020. [Online]. Available: https://doi.org/10.1007/s11128-020-02873-5

[27] S. Niu, A. Suau, G. Staffelbach, and A. Todri-Sanial, "A hardware-aware heuristic for the qubit mapping problem in the NISQ era," *CoRR*, vol. abs/2010.03397, 2020. [Online]. Available: https://arxiv.org/abs/2010.03397

[28] D. Bhattacharjee and A. Chattopadhyay, "Depth-optimal quantum circuit placement for arbitrary topologies," *CoRR*, vol. abs/1703.08540, 2017. [Online]. Available: http://arxiv.org/abs/1703.08540

[29] A. Shafaei, M. Saeedi, and M. Pedram, "Qubit placement to minimize communication overhead in 2d quantum architectures," in *19th Asia and South Pacific Design Automation Conference, ASP-DAC 2014, Singapore, January 20-23, 2014*. IEEE, 2014, pp. 495–500. [Online]. Available: https://doi.org/10.1109/ASPDAC.2014.6742940

[30] M. Pedram and A. Shafaei, "Layout optimization for quantum circuits with linear nearest neighbor architectures," *IEEE Circuits and Systems Magazine*, vol. 16, no. 2, pp. 62–74, 2016.

[31] B. Tan and J. Cong, "Optimal layout synthesis for quantum computing," in *IEEE/ACM International Conference On Computer Aided Design, ICCAD 2020, San Diego, CA, USA, November 2-5, 2020*. IEEE, 2020, pp. 137:1–137:9. [Online]. Available: https://doi.org/10.1145/3400302.3415620

[32] M. Hodson, B. Ruck, H. Ong, D. Garvin, and S. Dulman, "Portfolio rebalancing experiments using the quantum alternating operator ansatz," *arXiv preprint arXiv:1911.05296*, 2019.

[33] M. A. Nielsen and I. L. Chuang, *Quantum Computation and Quantum Information (10th Anniversary edition)*. Cambridge University Press, 2016. [Online]. Available: https://www.cambridge.org/de/academic/subjects/physics/ quantum-physics-quantum-information-and-quantum-computation/ quantum-computation-and-quantum-information-10th-anniversary-edition? format=HB

[34] A. Kissinger and J. van de Wetering, "PyZX: Large Scale Automated Diagrammatic Reasoning," in Proceedings 16th International Conference on *Quantum Physics and Logic,* Chapman University, Orange, CA, USA., 10-14 June 2019, ser. Electronic Proceedings in Theoretical Computer Science, B. Coecke and M. Leifer, Eds., vol. 318. Open Publishing Association, 2020, pp. 229–241.