

# Optimizing Quantum Circuits, Fast and Slow

Amanda Xu

University of Wisconsin-Madison  
Madison, WI, USA  
axu44@wisc.edu

Swamit Tannu

University of Wisconsin-Madison  
Madison, WI, USA  
swamit@cs.wisc.edu

Abtin Molavi

University of Wisconsin-Madison  
Madison, WI, USA  
amolavi@wisc.edu

Aws Albarghouthi

University of Wisconsin-Madison  
Madison, WI, USA  
aws@cs.wisc.edu

## Abstract

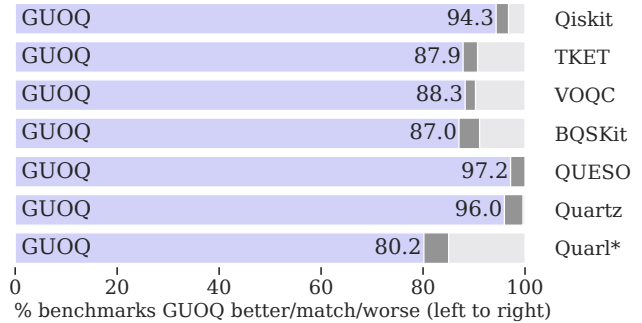
Optimizing quantum circuits is critical: the number of quantum operations needs to be minimized for a successful evaluation of a circuit on a quantum processor. In this paper we unify two disparate ideas for optimizing quantum circuits, *rewrite rules*, which are fast standard optimizer passes, and *unitary synthesis*, which is slow, requiring a search through the space of circuits. We present a clean, unifying framework for thinking of rewriting and resynthesis as abstract circuit transformations. We then present a radically simple algorithm, *GUOQ*, for optimizing quantum circuits that exploits the synergies of rewriting and resynthesis. Our extensive evaluation demonstrates the ability of *GUOQ* to strongly outperform existing optimizers on a wide range of benchmarks.

## 1 Introduction

Quantum computing enables efficient simulation of quantum mechanical phenomena, promising to catalyze advances in quantum physics, chemistry, materials science, and beyond. Near-term quantum computers with more than a thousand qubits operating in a noisy environment without error correction have already been deployed, marking the current era of *Noisy Intermediate Scale Quantum* (NISQ) computing [48]. Recent groundbreaking experiments have implemented error-corrected *logical* qubits and demonstrated potential for reducing *logical error* [7, 12]. Although many challenges remain, *fault-tolerant quantum computing* (FTQC) is on the horizon.

In both NISQ and FTQC, reducing errors is a critical obstacle to overcome. Every quantum operation has a probability of failure causing a quantum execution to quickly devolve into random noise. The NISQ paradigm aims to mitigate these errors in the absence of error correction primarily by reducing the number of operations. However, error correction in FTQC is not a panacea and introduces its own unique bottlenecks [9, 58], which can render the error correction scheme useless if left untamed. Especially in the near term, FTQC architectures may face challenges in handling large circuit depths due to physical imperfections such as two-level system (TLS) drift, qubit leakage, high-energy particle strikes,

GUOQ vs. State-of-the-Art Quantum Optimizers



**Figure 1.** Summary of *GUOQ* compared to state-of-the-art on 2-qubit-gate reduction for the IBMQ20 gate set. *GUOQ* and BQSKit are allowed to approximate the circuit up to  $\epsilon = 10^{-8}$ . \*Quarl requires an NVIDIA A100 (40GB) GPU to run.

and crosstalk [1, 7, 38]. Therefore, it is of utmost importance to reduce the number of operations for FTQC as well.

Current approaches tackling quantum circuit optimization primarily apply *peephole optimization* using a fixed set of *rewrite rules*. Some tools use a small set of hand-crafted rules [20, 29, 40], while others automatically synthesize rules [66, 67]. The idea is to apply rewrite rules in a schedule, transforming subcircuits to semantically equivalent ones with fewer operations. *Rewrite rules are fast* to apply—match a pattern and rewrite it—but inherently only perform local optimizations.

An orthogonal line of work has been studying the problem of *unitary synthesis*. A unitary matrix precisely represents the semantics of a quantum program. Some quantum algorithms are simple to state in the form of a unitary but nontrivial to decompose into elementary operations that can be executed on hardware [15, 18]. Thus, a large body of work has focused on synthesizing quantum circuits that implement a given unitary matrix [4, 13, 26, 43, 50, 51, 59, 62, 68]. Recent works [44, 65] have applied these algorithms to optimize quantum circuits by partitioning large circuits into manageably-sized *subcircuits* consisting of a few qubits at most and then *resynthesizing* each subcircuit to produce a new subcircuit whose unitary is equivalent, or close enough,

	Rewrite rules	Resynthesis
Fast	✓	✗
Limited by # gates	✓	✗
Limited by # qubits	✗	✓
Approximate	✗	✓

**Table 1.** Characteristics of rewrite rules and resynthesis.

to the original subcircuit’s unitary. *Unitary synthesis is slow*: usually a combinatorial search problem through the space of circuits, but can apply to deep subcircuits.

Rewrite rules and resynthesis have their own strengths and weaknesses—see Table 1. Individual rewrite rules are fast to apply but are limited to small patterns with few gates. On the other hand, resynthesis is slow but can support circuits with an arbitrary number of gates because the primary limiting factor is the number of qubits. Critically, resynthesis has the power to optimize deep subcircuits and it can *approximate* the solution to some degree. Rewrite rules are too small and rigid to find meaningful approximations. Approximations can unlock new circuit optimizations but must be applied carefully to avoid introducing too much error.

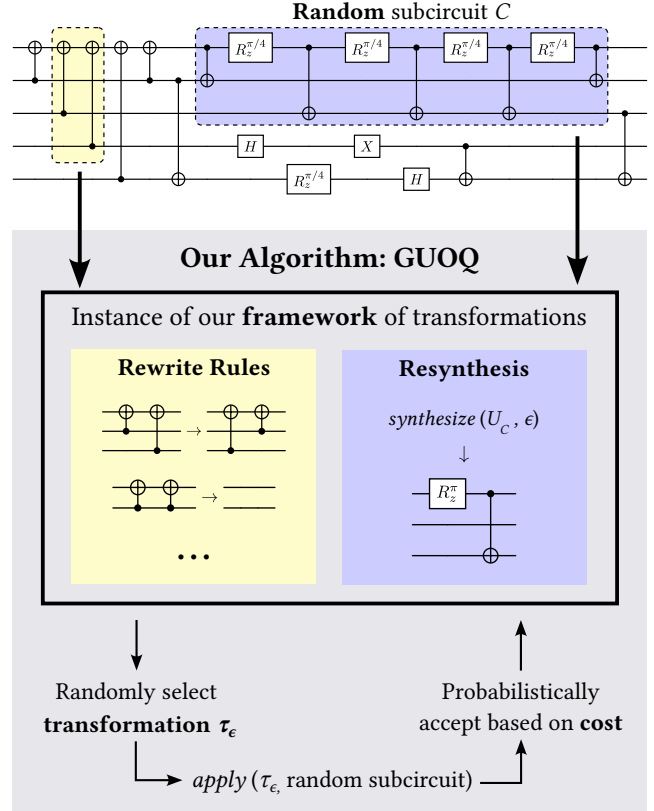
Inspired by how humans combine fast and slow modes of thinking [25], *Systems 1 and 2*, we ask the following question:

*Can we design an optimization approach that can synergistically combine the powers of optimizing quantum circuits fast, via rewrite rules, and slow, via resynthesis?*

We answer this question affirmatively: we demonstrate that we can unify rewriting and resynthesis and propose a simple optimization algorithm that significantly outperforms existing approaches in the literature.

**Our approach.** We propose a framework to unify rewrite rules and resynthesis for optimizing quantum circuits (see Fig. 2). The key insight is that we can abstract both optimizations into a closed-box circuit *transformation* with a degree of approximation  $\epsilon$ . Our flexible and generic framework allows arbitrary transformations, which can be applied freely. We prove a simple additive upper bound on the final approximation after applying a sequence of transformations.

We can instantiate our framework using a set of circuit transformations. The key challenge is deciding in what order to apply these transformations—this is the *phase ordering problem* that has plagued optimizing compilers for decades! We discover that, perhaps surprisingly, a simple and lightweight *simulated annealing*-like algorithm is the most effective solution, outperforming more clever heuristics or algorithms. The lack of structure in our problem lends itself to an approach that randomly and quickly alternates between fast and slow optimization, as opposed to sophisticated approaches guided by hand-crafted heuristics or reinforcement learning. Let this serve as another bitter lesson [61] that simple methods often prevail.

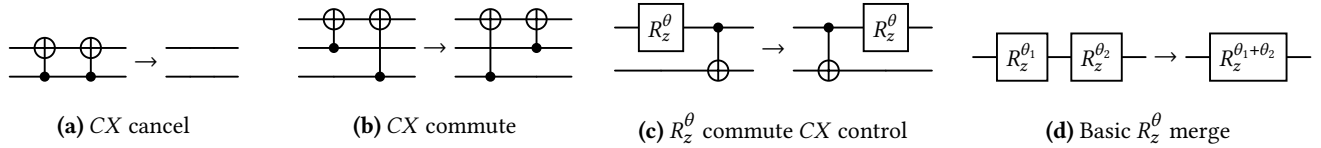


**Figure 2.** Overview of our approach

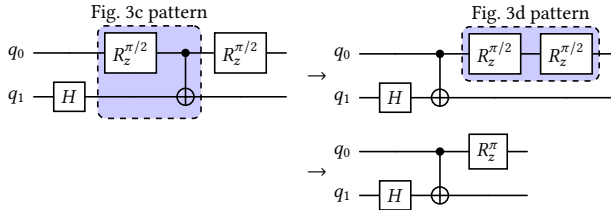
We implement our algorithm, GUOQ (Good Unified Optimizations for Quantum), and provide an extensive evaluation against state-of-the-art optimizers and *superoptimizers* using a benchmark suite with 247 diverse quantum circuits implementing near- and long-term algorithms. Our evaluation demonstrates the following: (1) GUOQ significantly outperforms state-of-the-art tools (see Fig. 1 for a summary), (2) GUOQ’s randomized search approach is critical for efficiently combining rewriting and resynthesis, and (3) GUOQ can flexibly perform well in the NISQ and FTQC regimes. For instance, GUOQ outperforms the recent superoptimizer, Quarl [32], in terms of 2-qubit-gate reduction on 80% of the benchmarks. This is despite the fact that Quarl uses a specialized deep reinforcement learning technique for quantum-circuit optimization, requires more computational resources (GPUs), and has been trained on portions of the benchmark suite.

**Contributions.** We make the following contributions:

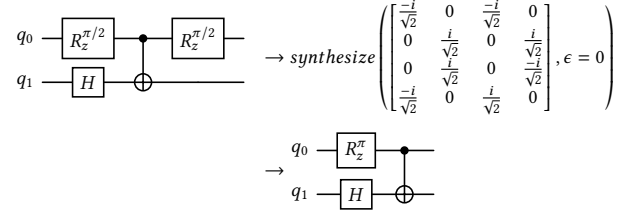
- A framework that abstracts the inner workings of rewriting and resynthesis into closed-box circuit *transformations* with *approximate* semantic guarantees.
- A lightweight algorithm, GUOQ, inspired by simulated annealing, that searches the space of transformations.
- An implementation of GUOQ and thorough evaluation considering both NISQ and FTQC that demonstrates its effectiveness compared to state-of-the-art optimizers.



**Figure 3.** Examples of rewrite rules. Observe how the rules with  $R_z^\theta$  use *symbolic*  $\theta$  angles.



**Figure 4.** Example of applying the rule from Fig. 3c followed by the rule from Fig. 3d.



**Figure 5.** Example of resynthesizing the initial Fig. 4 circuit.

## 2 Background and Overview

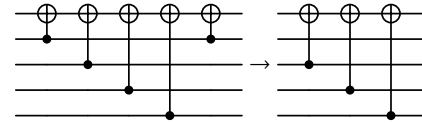
In this section, we provide the high-level background required for understanding our approach and use examples to highlight the differences between rewrite rules and resynthesis. We also describe an overview of our approach along with some concrete examples showing how rewrite rules and resynthesis can work together.

### 2.1 Quantum Circuits Background

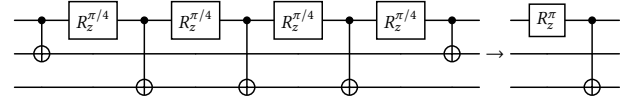
**Quantum circuits.** Quantum circuits are composed of combinations of quantum operations (or *gates*) applied to *qubits*. Some operations, like the *Hadamard* gate ( $H$ ), are only applied to a single qubit, whereas operations like the *controlled-NOT* gate ( $CX$ ) apply to an ordered pair of qubits. Another common class of quantum gates includes rotational gates parameterized on input angles, such as the  $R_z^\theta$  gate.

Consider the example circuit on the left sides of Figs. 4 and 5, where each horizontal wire corresponds to a qubit—e.g., the  $H$  gate is the first gate applied to qubit  $q_1$ .

**Rewrite rules.** A rewrite rule is a pair of semantically equivalent circuits. Although rewrite rules are in principle bidirectional, we will refer to the left-hand side as the *pattern* and the right-hand side as the *replacement* for simplicity. Fig. 3 shows examples of some commonly used rewrite rules. Applying rewrite rules to a circuit is simple. Begin by searching for a *match* for the pattern and if one exists, substitute it with the replacement. For example, in Fig. 4, there is a match for the pattern of the rule in Fig. 3c, shown by the highlighted subcircuit. Rewriting the match to the replacement allows the rule in Fig. 3d to apply and reduces the gate count by one. This general idea of composing rewrite rules is the heart of how we optimize quantum circuits using rewrite rules.



(a) Rewrite rules better



(b) Resynthesis better

**Figure 6.** Comparing rewrite rules and resynthesis.

**Resynthesis.** Circuit *resynthesis* takes advantage of the vast line of work done in *unitary synthesis* to resynthesize circuits according to some optimization objective. A circuit can be represented as a unitary matrix. Given a circuit's unitary matrix, unitary synthesis constructs a new circuit, with fewer gates, whose unitary is within  $\epsilon$  of the original unitary according to some distance metric. Note that the original circuit's structure is lost by converting it into a unitary and the synthesis algorithm needs to search for a new circuit structure from scratch. This is inherently a slow search through the space of circuits, e.g., the BQSKit compiler [69] performs a bottom-up search using two-qubit subcircuits.

Fig. 5 illustrates circuit resynthesis using the same initial circuit as Fig. 4 and no approximation by setting  $\epsilon = 0$ . Observe how the final circuit is equivalent to the result from applying rewrite rules because we can apply the rule in Fig. 3c to push the  $R_z^\pi$  gate across the control of the  $CX$ .

### 2.2 Comparing Rewriting and Resynthesis

Recall that resynthesis is limited by the number of qubits in the circuit, because the size of the unitary is exponential in

the number of qubits. Fig. 6a is an example of a circuit where it is better to apply rewrite rules. The structure closely resembles the circuit for the *quantum Fourier transform* (QFT) [11], a critical subroutine in many quantum algorithms. This circuit involves too many qubits for unitary synthesis to succeed in a reasonable amount of time. However, it only takes a few applications of two rewrite rules (Figs. 3a and 3b) to quickly get to the right-hand side. Even if we partition this circuit into more tractable 3-qubit subcircuits to resynthesize, we would not be guaranteed to reach the right-hand side. In fact, it would require a series of lucky coincidences over multiple rounds of partitioning the circuit.

Resynthesis involving fewer qubits though can compensate for the limited-sized patterns in rewrite rules. Fig. 6b is an example of a deep circuit where resynthesis can accelerate the search. Although we can achieve the optimized circuit using rewrite rules, it requires a long sequence of several of the rules from Fig. 3 in a very specific order. Resynthesis can discover the circuit on the right-hand side all at once because the circuit only involves 3 qubits.

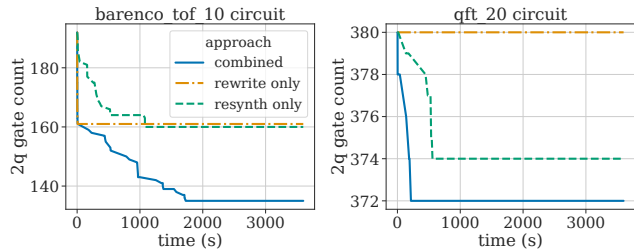
As we saw in the above examples, there are complementary qualities between rewrite rules and resynthesis.

### 2.3 Our Approach

**Unifying rewrite rules and resynthesis.** Our framework to unify rewrite rules and resynthesis introduces an abstraction for *transforming* circuits. We specify a function called a circuit *transformation*, denoted  $\tau_\epsilon$ , that returns a circuit that is semantically equivalent to the original up to the approximation  $\epsilon$ . Beyond this guarantee, transformations are closed-box. Given a set of rewrite rules and resynthesis algorithms, we can instantiate a set of transformations. Crucially, our framework allows us to apply circuit transformations in any order, and we prove an upper bound on the final approximation degree when applying an arbitrary sequence of transformations.

**Optimization objectives.** Optimizing quantum circuits requires diverse optimization objectives depending on the application. In NISQ, two-qubit gates are the dominant source of noise whereas in FTQC,  $T$  gates are the most costly to perform in an error-corrected fashion, followed by two-qubit gates. We view these gate-minimizing optimization objectives as *soft* constraints in our search. Our *hard* constraint when resynthesizing subcircuits is staying within the specified global error threshold  $\epsilon_f$ . Allowing more error can allow the synthesis algorithm to find a solution with fewer gates, so it is critical to find a balance between these two competing optimization objectives. For example, an objective for NISQ might be the following:  $\text{argmin}_{C'} 2Q\text{-COUNT}(C') \text{ s.t. } \epsilon_{C'} \leq \epsilon_f$ .

**The *GUOQ* algorithm.** The vast and discrete landscape for optimizing quantum circuits provides sparse navigation for traversing it. *GUOQ* is our simple algorithm inspired by simulated annealing [28] that rapidly and randomly searches the space of transformations. We find that an approach like



**Figure 7.** An example showing the two-qubit gate count of the current best solution for `barenco_tof_10` and `qft_20` over an hour of search using 1) only rewrite rules, 2) only resynthesis, and 3) rewrite rules and resynthesis combined.

*GUOQ* is well-suited for solving our problem because it has minimal memory requirements, is easy to implement, and explores the solution space much faster than other approaches. At a high level, the algorithm maintains a single candidate and applies randomly chosen transformations to randomly chosen subcircuits. Transformations with  $\epsilon = 0$  can be applied an unlimited number of times while transformations with  $\epsilon > 0$  are limited based on the target error threshold. If a transformation preserves or reduces gates related to the optimization objective, it is always accepted. Otherwise, it is only accepted with a small probability.

**A concrete example.** As a primer, we provide an example demonstrating the benefits of combining rewrite rules and resynthesis on the `barenco_tof_10` benchmark, which is an implementation of a multi-control Toffoli gate [5], and the `qft_20` benchmark, which implements the *quantum Fourier transform* (QFT) [11]. These benchmarks are critical building blocks for the famous Shor’s algorithm [56]. Fig. 7 shows the two-qubit gate count of the best solution over an hour of search for *GUOQ* using 1) rewrite rules only, 2) resynthesis only, and 3) rewrite rules and resynthesis combined.

As we can see, rewrite rules on their own quickly get stuck in a local minimum for almost the entire search, whereas resynthesis on its own is able to gradually make progress but is too slow. Combining rewrite rules and resynthesis allows resynthesis to progress the search when rewrite rules get stuck by mutating the circuit and “teleporting” the candidate solution to a different area of the optimization landscape. Working in tandem, rewrite rules and resynthesis can go beyond the capabilities of either alone.

## 3 Quantum Circuits and Approximations

We now formalize the necessary background for understanding quantum circuits and their semantics. The notion of approximate semantics and subcircuits will be critical for our framework.

**Semantics.** A quantum gate is a linear transformation of the state vector. A gate  $g$  acting on  $m$  qubits can be represented by a  $2^m \times 2^m$  unitary matrix  $U_g$ . Composing the

gates in a circuit with  $n$  qubits using matrix multiplication results in a  $2^n \times 2^n$  unitary matrix exactly representing the semantics of the circuit.

**Example 3.1.** The semantics of the  $T$  and  $CX$  gates are the following unitary matrices:

$$U_T := \begin{bmatrix} 1 & 0 \\ 0 & e^{i\pi/4} \end{bmatrix}, U_{CX} := \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

Consider the circuit  $C := T q_1; CX q_0 q_1$ . The unitary  $U_C$  representing the semantics of  $C$  is precisely  $U_{CX}(I \otimes U_T)$ , where the tensor product  $(I \otimes U_T)$  with the identity matrix denotes that the  $T$  gate is applied to the second qubit.

**Circuit equivalence.** Circuits are semantically equivalent if their unitaries are exactly equal. Additionally, some circuits may be equivalent up to *global phase* because two quantum states  $|\psi_1\rangle$  and  $|\psi_2\rangle = e^{i\phi} |\psi_1\rangle$  are observationally indistinguishable [42] for any angle  $\phi \in \mathbb{R}$ . Two circuits  $C$  and  $C'$  are semantically equivalent modulo global phase, denoted  $C \equiv C'$ , if and only if  $U_C = e^{i\phi} U_{C'}$  for some angle  $\phi$ .

**Approximating circuits.** We can approximate the semantics of a circuit  $C$  arbitrarily by defining a notion of *distance*, which is a function of the original and approximated circuits' unitaries. The *Hilbert–Schmidt distance* ( $\Delta$ ) is a convenient distance function used in prior work [43, 44] due to its ability to handle equivalence modulo global phase and ease of computation. The Hilbert–Schmidt distance between two unitary matrices is formally defined in Def. 3.2. Using this definition, we can define approximate circuit equivalence as shown in Def. 3.3.

**Definition 3.2** (Hilbert–Schmidt distance). Let  $U$  and  $U'$  be two  $N \times N$  unitary matrices.  $\Delta(U, U') := \sqrt{1 - \frac{\|Tr(U^\dagger U')\|^2}{N^2}}$ .

**Definition 3.3** (Approximate circuit equivalence). Two circuits  $C$  and  $C'$  are  $\epsilon$ -equivalent, denoted  $C \equiv_\epsilon C'$ , if and only if  $\Delta(U_C, U_{C'}) \leq \epsilon$ .

**Subcircuits.** To define a subcircuit, it is best to consider a DAG representation of a circuit, where the nodes are gates and the wires between gates are directed from left to right. Then a subcircuit is precisely a *convex* subgraph of this DAG. As defined in prior work [67], a *convex* subgraph is a subgraph that contains every path that exists between its nodes in the original graph. Intuitively, this requirement enforces continuous qubit wires in the subcircuit. For example, Fig. 2 (top) shows two highlighted subcircuits.

## 4 A Unified Optimization Framework

We are ready to formally describe our framework unifying rewrite rules and resynthesis. Our framework introduces abstract *transformations* with an associated error tolerance  $\epsilon$ . We show how to represent rewrite rules and resynthesis in

this unifying framework and prove an upper bound on the error when composing these transformations in arbitrary sequences. This underlies the design of `GUOQ`, which applies transformations in arbitrary orders.

### 4.1 Circuit Transformations

The following definition presents an abstraction of a circuit transformation as a function that takes a circuit  $C$  and produces an  $\epsilon$ -equivalent  $C'$ .

**Definition 4.1** (Transformation). Let  $C$  represent a circuit type. A transformation  $\tau_\epsilon : C \rightarrow C$  accepts a circuit  $C$  and returns a circuit  $C'$  such that  $C \equiv_\epsilon C'$ .

Rewrite rules and resynthesis can be represented in this framework as transformations over subcircuits. Consider a rewrite rule  $C_1 \rightarrow C_2$ . The transformation

$$\tau_0(C) := \begin{cases} C_2 & \text{if } C = C_1 \\ C & \text{otherwise} \end{cases}$$

captures the rewrite rule by transforming circuits that are syntactically equivalent to  $C_1$  (up to qubit renaming) and acting as the identity function otherwise. Observe how this transformation carries an  $\epsilon = 0$  because rewrite rules preserve exact semantic equivalence.

Similarly, we can define a transformation representing resynthesis. Assume there exists a circuit resynthesis function  $\text{RESYNTH} : (C \times \mathbb{R}) \rightarrow C$  that given a circuit  $C$  and error tolerance  $\epsilon$ , returns a circuit  $C'$  such that  $C \equiv_\epsilon C'$ . The transformation is simply  $\tau_\epsilon(C) := \text{RESYNTH}(C, \epsilon)$ . An example of such a circuit resynthesis function is a thin wrapper around a unitary synthesis function that computes the input circuit's unitary before invoking unitary synthesis.

### 4.2 Composing Transformations

Composing transformations is not as straightforward as composing rewrite rules because transformations are allowed to approximate the circuit. Prior work [44] shows how to upper bound the error when approximating *disjoint* partitions of a circuit. We present a flexible and generic framework that allows us to apply transformations in an arbitrary fashion. In particular, we can apply a transformation to a subcircuit that only *partially* contains a previously transformed subcircuit.

Formalized in Thm. 4.2, we prove that the upper bound on the error when composing a finite sequence of transformations is the sum of all the errors from each transformation. Without loss of generality, we can assume all transformations have the same error.

**Theorem 4.2.** *Suppose we are given a set of transformations  $\tau_\epsilon^1, \dots, \tau_\epsilon^n$ . Let  $C_0, \dots, C_n$  be a sequence of circuits such that  $C_i$  is the result of applying transformation  $\tau_\epsilon^i$  to a subcircuit of  $C_{i-1}$  for all  $1 \leq i \leq n$ . Then,  $C_0 \equiv_{n\epsilon} C_n$ .*

## 5 GUOQ: A Stochastic Algorithm

In this section, we begin by formally stating the quantum-circuit optimization problem. Optimizing quantum circuits is hard because of the large search space and the difficulty of simulating quantum circuits. Next we describe our algorithm GUOQ for solving this problem given an instantiation of our framework. GUOQ is fast, flexible, and easy to implement. It applies a given set of transformations in a randomized fashion inspired by *simulated annealing*, a classic algorithm for solving discrete optimization problems. Finally, we discuss implementation details for optimizing our algorithm.

### 5.1 Optimization objectives for quantum circuits

Different quantum computing hardware and paradigms will require unique optimization objectives. For example, on NISQ hardware it is critical to reduce two-qubit gate count. Other optimization objectives include  $T$  count and circuit *depth*, rather than gate count. Our approach is flexible and we can define any cost function,  $\text{COST} : C \rightarrow \mathbb{R}$ , to minimize, where  $C$  is the set of all circuits.

**Example 5.1** (Multiple optimization objectives). Consider the FTQC setting where we want to reduce primarily  $T$  gates, followed by  $CX$  gates. We can model this optimization objective by defining  $\text{COST}$  as  $2 \cdot \#_T(C) + \#_{CX}(C)$ , where  $\#_T(C)$  and  $\#_{CX}(C)$  are the  $T$  and  $CX$  gate counts, respectively.

Transformations in our framework can be approximate, so it is natural to accept as input an error tolerance that the result should not exceed. Using this error tolerance as a hard constraint and  $\text{COST}$  as a soft constraint we can formulate the problem as a succinct constrained optimization problem.

**Definition 5.2** (Quantum-Circuit Optimization Problem). Given a circuit  $C$  and an *error tolerance*  $\epsilon_f \geq 0$ , the quantum-circuit optimization problem is the following constrained optimization problem:

$$\underset{C'}{\text{argmin}} \text{COST}(C') \quad \text{s.t.} \quad \Delta(C, C') \leq \epsilon_f$$

### 5.2 The GUOQ Algorithm

We propose an algorithm inspired by simulated annealing. Simulated annealing is a general algorithm for solving optimization problems with large search spaces. It has many nice properties such as being fast, memory-efficient, easy to implement, and interruptible at any time to obtain a partial solution. These properties, inherited in our algorithm, unlock an effective approach for solving a problem that is incredibly difficult to craft or learn predictive heuristics for.

Alg. 1 shows the pseudocode for our algorithm. The inputs to the algorithm are the inputs to the quantum-circuit optimization problem and a set of transformations  $\mathcal{T}$ . Given a set of rewrite rules and resynthesis methods, we can instantiate  $\mathcal{T}$  as discussed in § 4. The “moves” we are allowed to make to modify the current solution are precisely the transformations in  $\mathcal{T}$ . The heart of the algorithm simply randomly

---

### Algorithm 1 The GUOQ Algorithm

---

```

1: procedure GUOQ(circuit  $C$ , error  $\epsilon_f$ , transformations  $\mathcal{T}$ )
2:   initialize  $C_{best}$  and  $C_{curr}$  to  $C$ 
3:   initialize error $_{best}$  and error $_{curr}$  to 0
4:   while within time limit do
5:     Randomly select transformation  $\tau_\epsilon$  in  $\mathcal{T}$ 
6:     if error $_{curr} + \epsilon > \epsilon_f$  then
7:       continue
8:     Randomly select subcircuit  $C_s$  in  $C_{curr}$ 
9:      $C_{curr}^\tau \leftarrow$  result of replacing  $C_s$  with  $\tau_\epsilon(C_s)$  in  $C_{curr}$ 
10:    if  $\text{COST}(C_{curr}^\tau) \leq \text{COST}(C_{curr})$  then
11:       $C_{curr} \leftarrow C_{curr}^\tau$ 
12:      error $_{curr} \leftarrow$  error $_{curr} + \epsilon$ 
13:    else with probability  $\exp\left(-t \frac{\text{COST}(C_{curr}^\tau)}{\text{COST}(C_{curr})}\right)$ 
14:       $C_{curr} \leftarrow C_{curr}^\tau$ 
15:      error $_{curr} \leftarrow$  error $_{curr} + \epsilon$ 
16:    if  $\text{COST}(C_{curr}) < \text{COST}(C_{best})$  then
17:       $C_{best} \leftarrow C_{curr}$ 
18:      error $_{best} \leftarrow$  error $_{curr}$ 
19:  return  $C_{best}$ 

```

---

samples a transformation and randomly samples a subcircuit of the current solution circuit to apply the transformation to. This move is always accepted if it improves or preserves the quality of the solution with respect to  $\text{COST}$ . Otherwise, it is accepted with some small probability that can be tuned using the temperature hyperparameter  $t$ . We adapt the standard acceptance probability for simulated annealing [28], which approaches 0 as the candidate solution cost increases.

The remainder of the algorithm ensures the upper bound on the error in the final solution does not exceed the specified tolerance  $\epsilon_f$ . Using Thm. 4.2, we can simply keep track of the sum of all the errors across all transformations applied. If applying a transformation would exceed the error bound, then we abstain and skip to the next iteration where we have the opportunity to sample a rewrite rule transformation with  $\epsilon = 0$ . Thm. 5.3 states the correctness of GUOQ.

**Theorem 5.3** (Correctness of GUOQ). *Let  $C'$  be the result of  $\text{GUOQ}(C, \epsilon_f, \mathcal{T})$  for any circuit  $C$ , error tolerance  $\epsilon_f \geq 0$ , and set of transformations  $\mathcal{T}$ . Then,  $C \equiv_{\epsilon_f} C'$ .*

### 5.3 How to implement GUOQ efficiently

We now discuss key implementation ideas that improve the performance of GUOQ in practice.

**Weighing fast & slow.** Applying a transformation to a circuit is decomposed into selecting a transformation to apply and a location in the circuit to apply it to. In practice, we limit the probability of choosing resynthesis to 1.5% of the time, since it is expensive. In the remaining 98.5% of the time, we uniformly sample one of the rewrite rules.

Gate set	Gates	Architecture
IBMQ20 [2]	$U1^\theta, U2^{\theta_1, \theta_2}, U3^{\theta_1, \theta_2, \theta_3}, CX$	Supercond.
IBM-EAGLE [2]	$R_z^\theta, SX, X, CX$	Supercond.
IONQ [60]	$R_x^\theta, R_y^\theta, R_z^\theta, R_{xx}^\theta$	Ion Trap
Nam [40]	$R_z^\theta, H, X, CX$	None
Clifford + $T$ [17]	$T, T^\dagger, S, S^\dagger, H, X, CX$	Fault Tolerant

**Table 2.** Summary of gate sets.

**Randomly selecting subcircuits.** We choose a random subcircuit to apply the transformation to by picking a node uniformly at random in the circuit DAG to begin constructing a subcircuit from. For rewrite rules transformations, completely random subcircuits will likely not be transformed nontrivially. We optimize this step by starting at a random node and performing a full pass through the circuit, replacing every disjoint match of the left-hand side with the right-hand side of the rule. For resynthesis transformations, we start at a random node and grow a subcircuit greedily until we cannot add more nodes without exceeding the qubit limit. We only apply resynthesis to a single subcircuit per iteration.

**Applying resynthesis asynchronously.** Invoking a unitary synthesis subroutine, even for a circuit with only 3 qubits, is slow and takes on the order of seconds or minutes to return a solution. To make better use of our time, we choose to make these calls asynchronously so we can apply rewrite rules concurrently. If the result of resynthesis is accepted, we effectively discard all modifications from rewrite rules made in the interim.

## 6 Implementation and Evaluation

We implemented GUOQ in Java. GUOQ interfaces with existing resynthesis tools [43, 69] and can be instantiated with arbitrary rewrite rules and gate sets. We designed our evaluation of GUOQ to answer the following research questions:

- Q1 How does GUOQ compare to state-of-the-art tools?
- Q2 What’s the effect of unifying rewriting & resynthesis?
- Q3 What’s the best way to apply rewriting & resynthesis?
- Q4 Does GUOQ extend to fault-tolerant computing (FTQC)?

We focus on NISQ in the first three questions using the diverse gate sets and tools available and then explore FTQC in Q4.

**Gate sets.** Our approach is flexible and can handle arbitrary gate sets. We evaluate on a variety of gate sets for promising quantum architectures, summarized in Table 2. The Nam gate set is not realized directly on hardware but is studied extensively in prior work due to its resemblance to the Clifford +  $T$  gate set.

**Benchmarks.** We consider all the benchmarks used in prior optimization work [32, 66] as well as benchmarks used in approximate optimization [44]. Prior optimization work primarily focuses on circuits with fewer than 2,000 gates. We expand the suite to larger application circuits considered

Tool	Superoptimize	Approach
Qiskit [2]	✗	fixed sequence of passes
TKET [57]	✗	fixed sequence of passes
VOQC [20]	✗	fixed sequence of passes
BQSKit [69]	✓	partition + resynthesis
QUESO [66]	✓	beam search + rewrite rules
Quartz [67]	✓	beam search + rewrite rules
Quarl [32]	✓	reinf. learning + rewrite rules

**Table 3.** State-of-the-art optimizers.

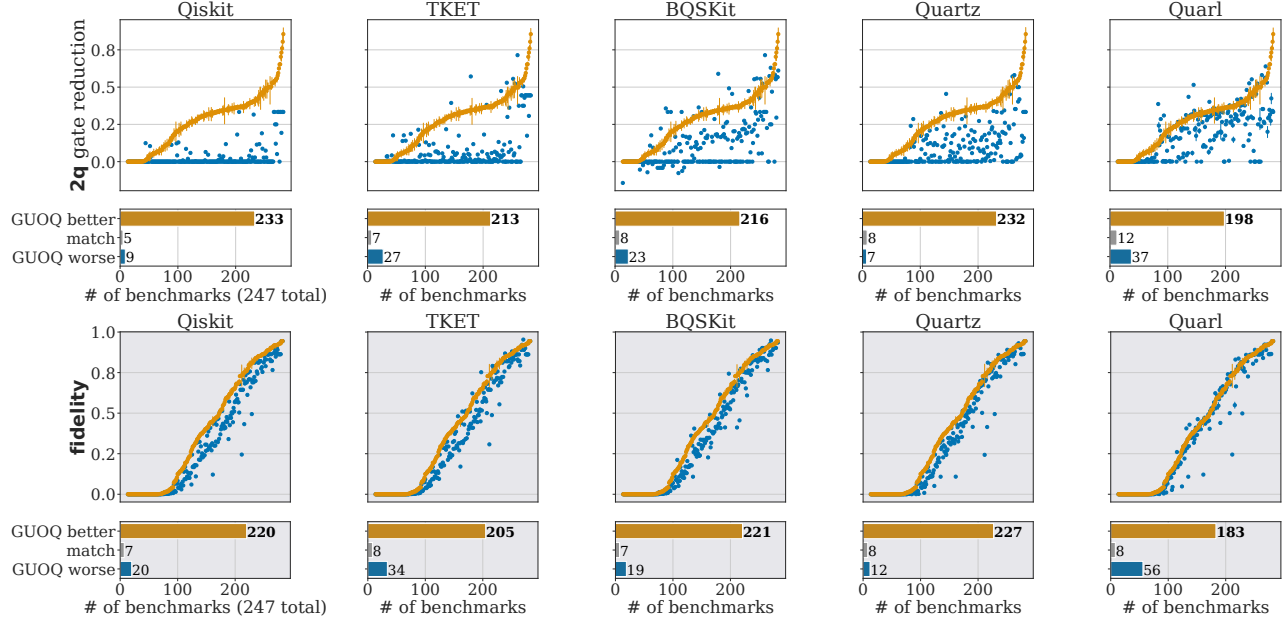
in prior mapping-and-routing work [72], and circuits implementing standard algorithms. Experimenting with larger circuits is key because total gate count is the primary metric that affects the scalability of optimizers for circuits with more than a few qubits. Our benchmark suite of 247 circuits includes important quantum algorithms in the near and long term such as QAOA [14], VQE [46], QPE [30], QFT [11], Grover’s [18], and Shor’s [56]. To ensure a fair comparison between each tool’s optimization phase, the input circuit throughout this evaluation is always already decomposed into the target gate set. Fig. 15 in Appendix B summarizes the total gate counts of all the input circuits. The benchmark circuits act on 4 to 36 qubits.

**Metrics.** For NISQ, we focus on two-qubit gate reduction because two-qubit gates have orders of magnitude higher error rates compared to single-qubit gates. Gate reduction is computed as  $1 - \frac{\text{optimized count}}{\text{original count}}$ . We also compute the circuit *fidelity*, or success probability, to emphasize that two-qubit gates are the dominant source of error. The fidelity of a gate is  $1 - \text{its error rate}$  and the fidelity of a circuit is the product of its gate fidelities. For IBMQ20 and IBM-EAGLE, we use the calibration data for the IBM Washington device available in Qiskit; for the Ion gate set, we use data for the IonQ Forte device [22]. In Q4, we consider different metrics for FTQC.

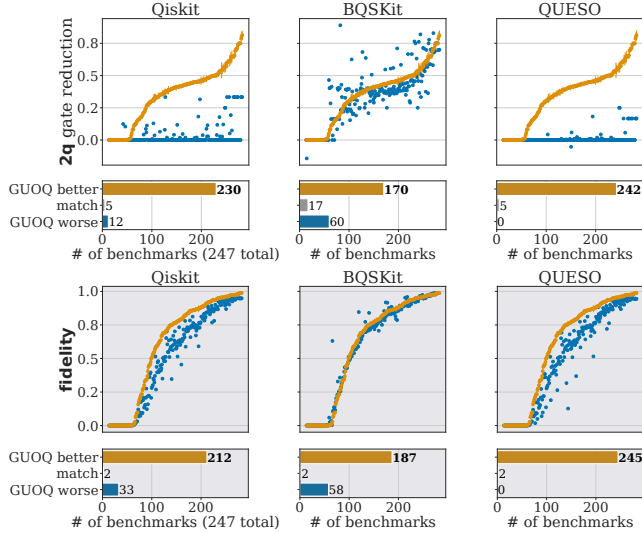
*In our plots, each point corresponds to a benchmark. For each benchmark circuit and tool, we compute the mean metric for a number of runs of the tool (10 trials for GUOQ) and a 95% confidence interval.* For readability, we present the benchmarks in increasing order sorted based on GUOQ. A point where GUOQ lies above the respective point for the other tool implies that GUOQ outperforms the other tool. The bar plot below each scatter plot summarizes the number of benchmarks GUOQ on average outperforms, matches, and underperforms, respectively, the tool in the title.

### Q1: Comparison with state-of-the-art optimizers

**State of the art.** We compared GUOQ to 7 state-of-the-art optimizers listed in Table 3. Our goal is to compare the optimization phase of various tools so we do not invoke mapping and routing and tools are allowed to change the input circuit’s connectivity if they support this feature. We exclude PyZX [29] from this research question because its primary



**Figure 8.** Comparison against state-of-the-art optimizers on the IBM-EAGLE gate set. Each graph has a bar chart summarizing the number of benchmarks GUOQ outperforms, matches, or underperforms the tool in the title. In each graph, orange is GUOQ and blue is the tool in the title. Each point is the mean metric and a 95% confidence interval for a number of runs on a single benchmark. The points are sorted based on GUOQ and a point where the orange lies above the blue indicates GUOQ is better.



**Figure 9.** GUOQ vs state-of-the-art on the IONQ gate set.

optimization objective is to reduce  $T$  gate count and often either increases or preserves the two-qubit gate count. We compare against PyZX in Q4 where we explore  $T$  reduction.

**Instantiation of GUOQ.** GUOQ uses rewrite rules generated by QUESO [66] and does not consider any size-increasing rewrite rules. To reduce two-qubit gate count, GUOQ uses BQSKit [69] for resynthesis and the optimization objective is to maximize fidelity. We limit random subcircuits to have

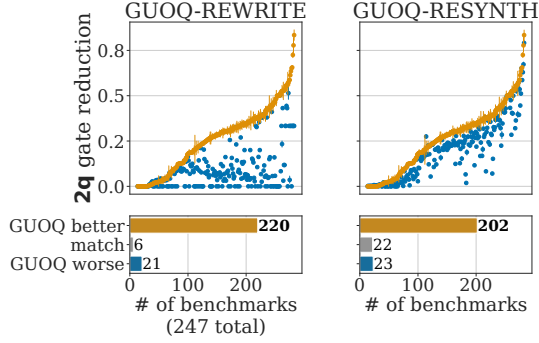
at most 3 qubits. The temperature hyperparameter  $t$  is set to 10, corresponding to a very small probability of accepting a worse solution. We chose this value empirically by performing a sweep of values from 0 (always accept) to 10.

**Experimental setup.** Unless otherwise indicated, we allocated each tool 1 hour, 32GB of RAM, and 1 CPU core on an AMD EPYC 7763 64-Core Processor. Quarl was run on a cluster of machines and was allocated 64GB of RAM, 1 NVIDIA A100-SXM4-40GB or 80GB GPU, and 1 CPU core. We only evaluate a tool on its supported gate sets. For approximate tools, we enforce an error upper bound of  $10^{-8}$ , which is (1) many orders of magnitude smaller than the error rate of a single two-qubit gate in NISQ ( $10^{-3}$  [27]) and (2) on-par with the logical error rate of a single error correction cycle for FTQC ( $10^{-6}$  to  $10^{-9}$  [41]) or an arbitrary-angle approximation.

We run Quarl for 3 trials with rotation merging, following their paper’s experimental setup. For Qiskit and BQSKit, we use the most powerful optimization levels, which are 3 and 4, respectively. We report the best solution found within the time and memory limits for all tools. That is, we use the partial solutions that GUOQ, QUESO, Quartz, and Quarl provide and use the original circuit for the other tools.

**Results.** The results for the IBMQ20, IBM-EAGLE, and Nam gate sets are all similar, so we only show the plots for the IBM-EAGLE gate set. Fig. 8 shows the comparison between GUOQ and state-of-the-art with respect to two-qubit gate reduction and fidelity on the IBM-EAGLE gate set. Recall that





**Figure 10.** GUOQ vs using only rewrite rules or resynthesis.

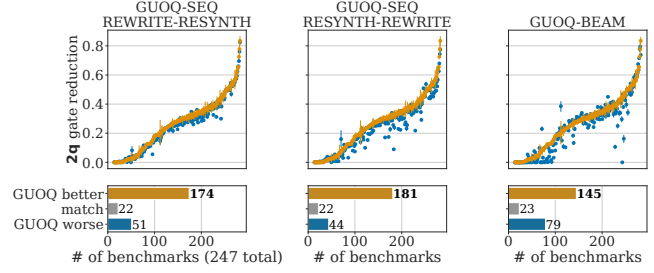
benchmarks where the orange point lies above the blue point are ones where GUOQ outperforms the tool in the title. For example, consider the left-most column where GUOQ outperforms Qiskit with respect to two-qubit gate reduction on 233/247 of the benchmarks, matches on 5, and underperforms on 9. We observe a very similar story in the fidelity graph on the bottom. This result holds in general against all other tools and GUOQ outperforms state-of-the-art on *at least* 80% and 74% of the benchmarks with respect to two-qubit gate reduction and fidelity, respectively. Recall that Quarl requires a GPU to run its specialized reinforcement learning and trains on a subset of the benchmark suite. GUOQ reduces two-qubit gate count by 28% on average while the next best tool, Quarl, has an average reduction of 18% and the best industrial toolkit, TKET, achieves 7% average reduction.

The results in Fig. 9 on the IONQ gate set depict a similar story. As we can see in the comparison against QUESO, a tool which synthesizes rewrite rules, the IONQ gate set is challenging for QUESO because rewrite rules are limited to patterns with a maximum of 3 gates to limit the combinatorial explosion of rules. GUOQ performs well because resynthesis can compensate for the limited rewrite rules. We will see in Q4 an example of the opposite effect. This demonstrates how the effectiveness of rewrite rules and resynthesis varies across different gate sets. Thus, unifying them provides a generic approach for optimizing diverse circuits.

**Q1 summary.** GUOQ significantly outperforms state-of-the-art across all gate sets. In the *worst* case, GUOQ outperforms other tools on 69% and 74% of the benchmarks with respect to two-qubit gate reduction and fidelity, respectively. In particular, GUOQ achieves an average of 28% two-qubit gate reduction on the IBM-EAGLE gate set while the state-of-the-art superoptimizer (requires GPU) and industrial toolkit achieve average reductions of 18% and 7%, respectively.

## Q2: Effect of combining rewriting and resynthesis

We explored this question by fixing the IBMQ20 gate set and running GUOQ with three different sets of transformations.



**Figure 11.** Comparing GUOQ against other search algorithms

GUOQ-REWRITE only uses rewrite rules synthesized by QUESO, GUOQ-RESYNTH only uses resynthesis, and GUOQ uses both.

**Results.** Fig. 10 shows that removing either rewrite rules or resynthesis is overall detrimental to the performance of GUOQ. Similar to Fig. 8, the title of each plot is the baseline and points below the curve are benchmarks where GUOQ outperforms the baseline. Observe how most of reduction comes from using resynthesis because it is a powerful optimization on its own. Interleaving with rewrite rules, which can take care of simple optimizations quickly, pushes the reduction even further.

**Q2 summary.** We can exploit the synergy between rewrite rules and resynthesis to achieve well beyond the capabilities of either alone.

## Q3: How to combine rewriting and resynthesis?

To answer this question, we fixed the IBMQ20 gate set and set of rewrite rule and resynthesis transformations while varying the search algorithm. We compare GUOQ against three alternate search algorithms for combining rewrite rules and resynthesis: (1) GUOQ-SEQ-REWRITE-RESYNTH, which spends the first half of the allotted time running GUOQ with rewrite rules only, then switches to running with resynthesis only, (2) GUOQ-SEQ-RESYNTH-REWRITE does the opposite, and (3) GUOQ-BEAM uses the MaxBeam algorithm of QUESO [66] to instantiate our framework.

**Results.** Fig. 11 shows the results. GUOQ outperforms the coarse interleaving GUOQ-SEQ-RESYNTH-REWRITE and GUOQ-SEQ-REWRITE-RESYNTH use on a majority of the benchmarks. This implies that tightly interleaving these different transformations is preferable to choosing a fixed ordering. We also see that relative to each other: GUOQ-SEQ-RESYNTH-REWRITE and GUOQ-SEQ-REWRITE-RESYNTH result in different solutions, which is further evidence that the ordering matters.

We now turn our attention to GUOQ-BEAM, which allows arbitrary orderings of rewrite rules and resynthesis but does not randomly sample transformations like GUOQ. Instead, GUOQ-BEAM maintains a large bounded priority queue of candidates and attempts to apply every transformation in each iteration (one for each transformation successfully applied),

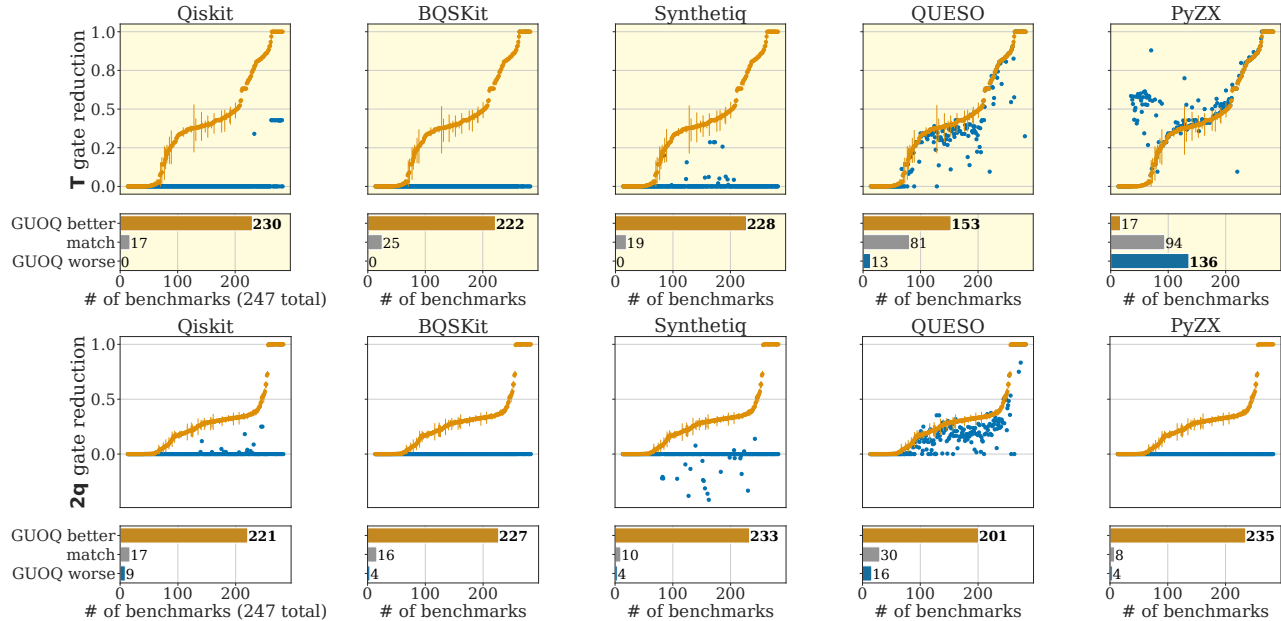


Figure 12. Comparison against state-of-the-art optimizers on the Clifford +  $T$  gate set.

which saturates the queue quickly with solutions of the same cost. In fact, the solutions are generally a few local transformations away from one another so the search makes much slower progress compared to GUOQ. The influx of candidates to the bounded queue also causes many solutions to be pruned, effectively wasting the time spent generating those candidates. Especially for larger circuits, the sizable queue is memory intensive, causing further slowdowns. In summary, the benefit of beam search considering many candidates to avoid local minima is lost in this problem setting.

**Q3 summary.** GUOQ achieves the best results by tightly interleaving rewrite rules and resynthesis using a simple, lightweight randomized algorithm.

#### Q4: Does GUOQ extend to FTQC?

In this research question, we shift our focus to the fault-tolerant Clifford +  $T$  gate set, where the desired optimization objective is an amalgamation of two NP-hard optimization problems [64]. We want to primarily reduce  $T$  gates [9] and reducing two-qubit gates is secondary, but still critical. Error correction is not perfect and the longer a quantum computation runs, the higher the risk of accruing uncorrectable logical error. Two-qubit gates, specifically  $CX$  in the Clifford +  $T$  gate set, increase the circuit runtime disproportionately because they inherently require more time compared to single-qubit Clifford gates [34] and architectural constraints can limit parallel execution [6, 21]. Furthermore, we anticipate that the problem of  $CX$  congestion will be exacerbated in compact FTQC architectures [33] with less routing space.

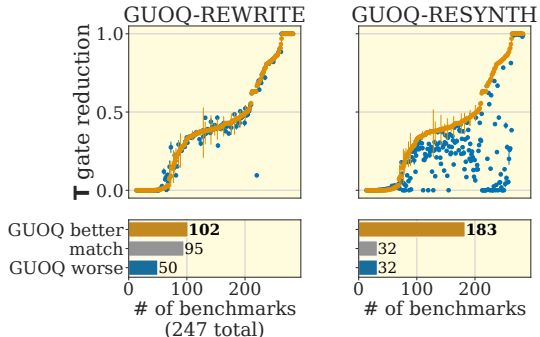


Figure 13. Revisiting Q2 for the Clifford +  $T$  gate set.

Lastly, recent work [16] has significantly reduced the space-time cost of preparing a  $T$  magic state, indicating that the focus on solely reducing  $T$  gates may soon need to be recalibrated.

We instantiate GUOQ with SynthetiQ [43], a state-of-the-art unitary synthesis algorithm for finite gate sets, and rewrite rules generated by QUESO [66]. We consider two additional tools in this comparison: (1) PyZX [29], a state-of-the-art optimizer for reducing  $T$  count using the rewrite-rule-based ZX-calculus and (2) our implementation of a BQSKit-style partitioning optimizer [44] that uses SynthetiQ.

**Results.** The top and bottom rows of Fig. 12 show the comparison against other tools with respect to  $T$  and  $CX$  gate reduction respectively. GUOQ outperforms all tools except PyZX with respect to  $T$  reduction and outperforms all tools with respect to  $CX$  gate reduction. Observe how reducing  $T$  gates is hard and a general-purpose tool, like Qiskit, only

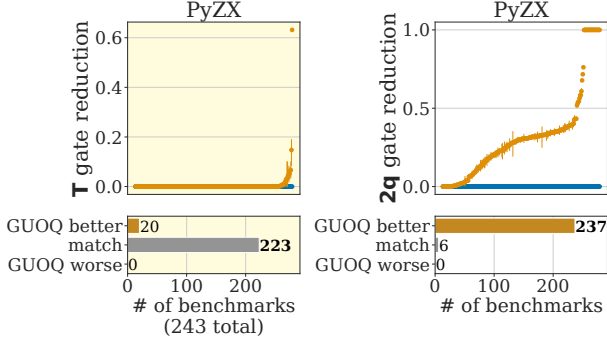


Figure 14. Running GUOQ on PyZX output.

reduces the number of  $T$  gates on 5% of the benchmarks. PyZX is a domain-specific optimizer that uses a powerful graph-based theory, the ZX-calculus, which is specialized for reducing phase gates, but does not reduce  $CX$  gates at all.

Digging deeper, we find that GUOQ is unable to surpass PyZX in  $T$  gate reduction because unitary synthesis for finite gate sets is much harder than for continuous gate sets. Fig. 13 shows the results from the same ablation as in Q2 and we find that rewrite rules contribute more than resynthesis.

To get a sense for how “good” PyZX’s solution is, we ran GUOQ on the output of PyZX for the 243 benchmarks it provided a solution for within the time and memory limits. We discovered that GUOQ can drastically reduce the  $CX$  gate count without increasing  $T$  gate count! Fig. 14 shows this result, which is exciting because PyZX on its own does not reduce  $CX$  gate count. Extending PyZX with GUOQ pushes the boundaries for the multifaceted FTQC optimization objective.

**Q4 summary.** GUOQ outperforms all tools with respect to  $T$  gate reduction except PyZX, which it only outperforms or matches on 45% of the benchmarks. However, GUOQ vastly outperforms PyZX (and other tools) with respect to  $CX$  gate reduction on at least 81% of the benchmarks, which is also critical in FTQC. Additionally, when run on the output of PyZX, GUOQ can reduce  $CX$  gate count on average by 32% without increasing  $T$  gate count.

## 7 Related Work

**Quantum optimizers.** Traditional quantum-circuit optimizers primarily use a fixed set of hand-crafted optimizations [2, 3, 8, 20, 40, 45, 57] applied in a fixed sequence. PyZX [29] takes a variant of this approach: representing a circuit as a ZX-diagram and applying the graphical rewrite rules of the ZX-calculus. GUOQ instead performs a fine-grained search over arbitrary optimizations.

**(Quantum) superoptimizers.** The idea of *superoptimizing* classical programs has been around for decades [37, 63].

This vast line of work [10, 23, 36, 39, 47, 53, 54] stems from the idea of finding the *optimal* solution for small programs that can later be applied in peephole optimizers. Approaches like STOKE [54] use MCMC [19] to superoptimize x86 assembly by randomly mutating the program. The ergodic theory behind MCMC lends itself to applications where testing correctness is easy. However, quantum circuits cannot be efficiently simulated on classical hardware [42].

Rewrite rules and unitary synthesis have both been used as the basis of quantum-circuit *superoptimizers*. Quartz [67] and QUESO [66] synthesize rewrite rules for a given gate set and use beam search to explore the space of rule application schedules. Quarl [32] applies reinforcement learning to schedule the application of rules generated by Quartz. None of these apply approximate circuit transformations. On the other hand, resynthesis-based superoptimizers [44, 69] optimize circuits by performing a single pass of partitioning into subcircuits, followed by applying unitary synthesis to each subcircuit. This approach circumvents the qubit count limitations of unitary synthesis, but is rigid and misses potential optimization opportunities that straddle the boundary between two adjacent partitions. In contrast, GUOQ is not limited to resynthesizing disjoint subcircuits of the original circuit. GUOQ can freely choose subcircuits to resynthesize by using Thm. 4.2 to bound the error when composing applications of resynthesis.

In a similar spirit to Quarl, other recent approaches have applied reinforcement learning to superoptimize quantum circuits. MQTPredictor [49] predicts the optimal passes and device with respect to an optimization objective but currently only considers a subset of Qiskit and tket passes. AlphaTensor-Quantum [52] is a closed-source approach that uses reinforcement learning for tensor decomposition to optimize  $T$  count.

**Domain-specific optimizers.** Other work targets specific applications like Hamiltonian simulation [31, 35] or variational algorithms [24]. Some tools operate at a lower level of abstraction than we do by considering gate *pulses* [55] or a higher level by starting from a program written in a high-level quantum programming language [70, 71]. In contrast, GUOQ is designed to be flexible for diverse quantum assembly circuits and architectures.

**Unitary synthesis.** Extensive prior work has considered the unitary synthesis problem for both finite and parameterized gate sets. For finite gate sets, some approaches [4, 62] provide theoretical guarantees of optimality in terms of circuit size, whereas others [26, 43] sacrifice optimality for improved runtime. For parameterized gate sets, several techniques [13, 50, 51, 59, 68] use numerical optimization to instantiate template circuits. However, all of these algorithms can only be applied to circuits with a handful of qubits.

## 8 Conclusions

We have described a generic and flexible framework for unifying rewriting and resynthesis for quantum-circuit optimization along with a simple and effective algorithm parameterized on an instantiation of this framework. Our approach, GUOQ, outperforms state-of-the-art optimizers in both near (NISQ) and long (FTQC) term quantum computing paradigms. For future work, we are interested in developing *symbolic* unitary synthesis so we can learn general transformations on the fly—as opposed to ones with highly specific angles—that will be more likely to apply later in the search.

## A Proofs

**Thm. 4.2.** By induction on  $n$ . *Base case:* Trivially,  $C_0 \equiv_0 C_0$ . *Induction case:* Assume  $C_0 \equiv_{k\epsilon} C_k$  for  $k \geq 0$  as our inductive hypothesis. We will show  $C_0 \equiv_{(k+1)\epsilon} C_{k+1}$ . We have  $C_k \equiv_\epsilon C_{k+1}$  by the proof of [44, §3.8] for disjoint partitions. Let  $U := U_{C_0}$ ,  $U' := U_{C_k}$ ,  $U'' := U_{C_{k+1}}$ ,  $\epsilon_1 := k\epsilon$ , and  $\epsilon_2 := \epsilon$ . Now it suffices to show  $\Delta(U, U'') \leq \epsilon_1 + \epsilon_2 = (k+1)\epsilon$ .

$$\begin{aligned}
 & \Delta(U, U'') \\
 &= \sqrt{1 - \frac{\|Tr(U^\dagger U'')\|^2}{N^2}} && \text{Def. 3.2} \\
 &= \sqrt{1 - \frac{\|Tr(U^\dagger U' U' U'')\|^2}{N^2}} && U' U'^\dagger = I \\
 &= \sqrt{1 - \frac{\|Tr[(U^\dagger U')(U' U'')]\|^2}{N^2}} && \text{Grouping terms} \\
 &\leq \sqrt{1 - \frac{\|Tr(U^\dagger U')\|^2}{N^2}} + \sqrt{1 - \frac{\|Tr(U' U'')\|^2}{N^2}} && [44, §3.8] \\
 &\leq \epsilon_1 + \epsilon_2 && \text{Defs. 3.2 and 3.3}
 \end{aligned}$$

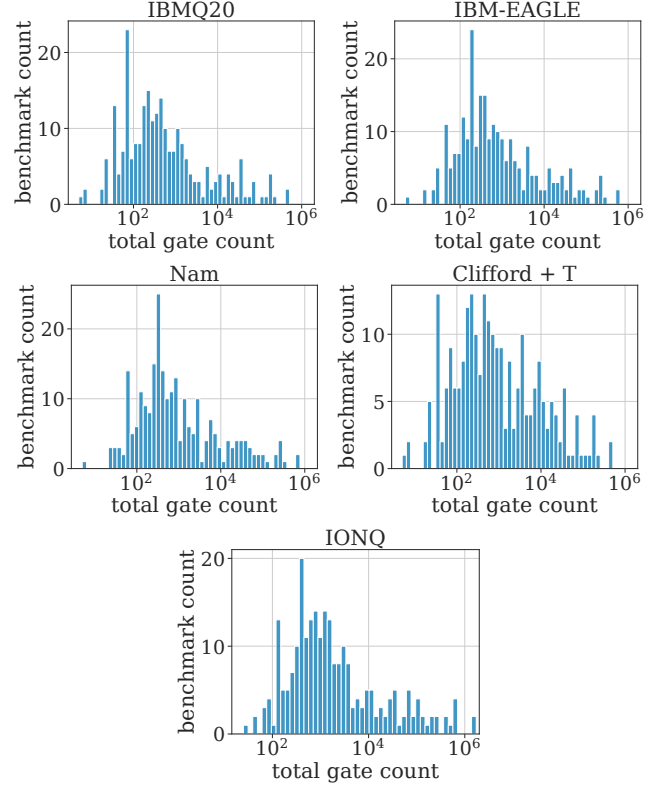
By definition of  $\equiv_\epsilon$ ,  $C_0 \equiv_{(k+1)\epsilon} C_{k+1}$ , as desired.  $\square$

**Thm. 5.3.** Follows directly from Thm. 4.2 and Alg. 1, line 6.

## B Benchmark Data

### References

- [1] Rajeev Acharya, Igor Aleiner, Richard Allen, Trond I. Andersen, Markus Ansmann, Frank Arute, Kunal Arya, Abraham Asfaw, Juan Atalaya, Ryan Babbush, Dave Bacon, Joseph C. Bardin, Joao Basso, Andreas Bengtsson, Sergio Boixo, Gina Bortoli, Alexandre Bourassa, Jenna Bovaird, Leon Brill, Michael Broughton, Bob B. Buckley, David A. Buell, Tim Burger, Brian Burkett, Nicholas Bushnell, Yu Chen, Zijun Chen, Ben Chiaro, Josh Cogan, Roberto Collins, Paul Conner, William Courtney, Alexander L. Crook, Ben Curtin, Dripto M. Debroy, Alexander Del Toro Barba, Sean Demura, Andrew Dunsworth, Daniel Eppens, Catherine Erickson, Lara Faoro, Edward Farhi, Reza Fatemi, Leslie Flores Burgos, Ebrahim Forati, Austin G. Fowler, Brooks Foxen, William Giang, Craig Gidney, Dar Gilboa, Marissa Giustina, Alejandro Grajales Dau, Jonathan A. Gross, Steve Habegger, Michael C. Hamilton, Matthew P. Harrigan, Sean D. Harrington, Oscar Higgott, Jeremy Hilton, Markus Hoffmann, Sabrina Hong, Trent Huang, Ashley Huff, William J. Huggins, Lev B. Ioffe, Sergei V. Isakov, Justin Iveland, Evan Jeffrey, Zhang Jiang, Cody Jones, Pavol Juhas, Dvir Kafri, Kostyantyn Kechedzhi, Julian Kelly, Tanuj Khattar, Mostafa Khezri, Mária Kieferová, Seon Kim, Alexei Kitaev, Paul V. Klimov, Andrew R. Klots, Alexander N. Korotkov, Fedor Kostritsa, John Mark Kreikebaum, David Landhuis, Pavel Laptev, Kim-Ming Lau, Lily Laws, Joonho Lee, Kenny



**Figure 15.** Summary of the benchmarks’ original total gate counts across all gate sets (log-scale x-axis).

Lee, Brian J. Lester, Alexander Lill, Wayne Liu, Aditya Locharla, Erik Lucero, Fionn D. Malone, Jeffrey Marshall, Orion Martin, Jarrod R. McClean, Trevor McCourt, Matt McEwen, Anthony Megrant, Bernardo Meurer Costa, Xiao Mi, Kevin C. Miao, Masoud Mohseni, Shirin Montazeri, Alexis Morvan, Emily Mount, Wojciech Mruzckiewicz, Ofer Naaman, Matthew Neeley, Charles Neill, Ani Nersisyan, Hartmut Neven, Michael Newman, Jiun How Ng, Anthony Nguyen, Murray Nguyen, Murphy Yuezhen Niu, Thomas E. O’Brien, Alex Opremcak, John Platt, Andre Petukhov, Rebecca Potter, Leonid P. Pryadko, Chris Quintana, Pedram Roushan, Nicholas C. Rubin, Negar Saei, Daniel Sank, Kannan Sankaragomathi, Kevin J. Satzinger, Henry F. Schurkus, Christopher Schuster, Michael J. Shearn, Aaron Shorter, Vladimir Shvarts, Jindra Skrzyny, Vadim Smelyanskiy, W. Clarke Smith, George Sterling, Doug Strain, Marco Szalay, Alfredo Torres, Guifre Vidal, Benjamin Villalonga, Catherine Vollgraf Heidweiller, Theodore White, Cheng Xing, Z. Jamie Yao, Ping Yeh, Juhwan Yoo, Grayson Young, Adam Zalcman, Yaxing Zhang, Ningfeng Zhu, and Google Quantum AI. Suppressing quantum errors by scaling a surface code logical qubit. *Nature*, 614(7949):676–681, 2023.

- [2] Gadi Aleksandrowicz, Thomas Alexander, Panagiotis Barkoutsos, Luciano Bello, Yael Ben-Haim, David Bucher, Francisco Jose Cabrera-Hernández, Jorge Carballo-Franquis, Adrian Chen, Chun-Fu Chen, Jerry M. Chow, Antonio D. Córcoles-Gonzales, Abigail J. Cross, Andrew Cross, Juan Cruz-Benito, Chris Culver, Salvador De La Puente González, Enrique De La Torre, Delton Ding, Eugene Dumitrescu, Ivan Duran, Pieter Eendebak, Mark Everitt, Ismael Faro Sertage, Albert Frisch, Andreas Fuhrer, Jay Gambetta, Borja Godoy Gago, Juan Gomez-Mosquera, Donny Greenberg, Ikko Hamamura, Vojtech Havlicek, Joe Hellmers, Łukasz Herok, Hiroshi Horii, Shaohan Hu, Takashi

- Imamichi, Toshinari Itoko, Ali Javadi-Abhari, Naoki Kanazawa, Anton Karazeev, Kevin Krsulich, Peng Liu, Yang Luh, Yunho Maeng, Manoel Marques, Francisco Jose Martín-Fernández, Douglas T. McClure, David McKay, Srujan Meesala, Antonio Mezzacapo, Nikolaj Moll, Diego Moreda Rodríguez, Giacomo Nannicini, Paul Nation, Pauline Ollitrault, Lee James O’Riordan, Hanhee Paik, Jesús Pérez, Anna Phan, Marco Pistoia, Viktor Prutyantov, Max Reuter, Julia Rice, Abdón Rodríguez Davila, Raymond Harry Putra Rudy, Mingi Ryu, Ninad Sathaye, Chris Schnabel, Eddie Schoute, Kanav Setia, Yunong Shi, Adenilton Silva, Yukio Siraichi, Seyon Sivarajah, John A. Smolin, Mathias Soeken, Hitomi Takahashi, Ivano Tavernelli, Charles Taylor, Pete Taylour, Kenso Trabing, Matthew Treinish, Wes Turner, Desiree Vogt-Lee, Christophe Vuillot, Jonathan A. Wildstrom, Jessica Wilson, Erick Winston, Christopher Wood, Stephen Wood, Stefan Wörner, Ismail Yunus Akhalwaya, and Christa Zoufal. Qiskit: An Open-source Framework for Quantum Computing, January 2019.
- [3] Matthew Amy and Vlad Gheorghiu. staq—a full-stack quantum processing toolkit. *Quantum Science and Technology*, 5(3):034016, jun 2020.
  - [4] Matthew Amy, Dmitri Maslov, Michele Mosca, and Martin Roetteler. A meet-in-the-middle algorithm for fast synthesis of depth-optimal quantum circuits. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 32(6):818–830, 2013.
  - [5] Adriano Barenco, Charles H. Bennett, Richard Cleve, David P. DiVincenzo, Norman Margolus, Peter Shor, Tycho Sleator, John A. Smolin, and Harald Weinfurter. Elementary gates for quantum computation. *Phys. Rev. A*, 52:3457–3467, Nov 1995.
  - [6] Michael Beverland, Vadym Kliuchnikov, and Eddie Schoute. Surface code compilation via edge-disjoint paths. *PRX Quantum*, 3:020342, May 2022.
  - [7] Dolev Bluvstein, Simon J. Evered, Alexandra A. Geim, Sophie H. Li, Hengyun Zhou, Tom Manovitz, Sepehr Ebadi, Madelyn Cain, Marcin Kalinowski, Dominik Hangleiter, J. Pablo Bonilla Ataides, Nishad Maskara, Iris Cong, Xun Gao, Pedro Sales Rodriguez, Thomas Karolyshyn, Giulia Semeghini, Michael J. Gullans, Markus Greiner, Vladan Vuletić, and Mikhail D. Lukin. Logical quantum processor based on reconfigurable atom arrays. *Nature*, 626(7997):58–65, 2024.
  - [8] C. Campbell, F. T. Chong, D. Dahl, P. Frederick, P. Goiporia, P. Gokhale, B. Hall, S. Issa, E. Jones, S. Lee, A. Litteken, V. Omole, D. Owusu-Antwi, M. A. Perlin, R. Rines, K. N. Smith, N. Goss, A. Hashim, R. Naik, E. Younis, D. Lobser, C. G. Yale, B. Huang, and J. Liu. Superstaq: Deep optimization of quantum programs. In *2023 IEEE International Conference on Quantum Computing and Engineering (QCE)*, pages 1020–1032, Los Alamitos, CA, USA, sep 2023. IEEE Computer Society.
  - [9] Earl T. Campbell, Barbara M. Terhal, and Christophe Vuillot. Roads towards fault-tolerant universal quantum computation. *Nature*, 549(7671):172–179, 2017.
  - [10] Berkeley Churchill, Rahul Sharma, JF Bastien, and Alex Aiken. Sound loop superoptimization for google native client. In *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS ’17, page 313–326, New York, NY, USA, 2017. Association for Computing Machinery.
  - [11] Don Coppersmith. An approximate fourier transform useful in quantum factoring. *arXiv preprint quant-ph/0201067*, 2002.
  - [12] MP da Silva, C Ryan-Anderson, JM Bello-Rivas, A Chernoguzov, JM Dreiling, C Foltz, JP Gaebler, TM Gatterman, D Hayes, N Hewitt, et al. Demonstration of logical qubits and repeated error correction with better-than-physical error rates. *arXiv e-prints*, pages arXiv–2404, 2024.
  - [13] Marc G. Davis, Ethan Smith, Ana Tudor, Koushik Sen, Irfan Siddiqi, and Costin Iancu. Towards optimal topology aware quantum circuit synthesis. In *2020 IEEE International Conference on Quantum Computing and Engineering (QCE)*, pages 223–234, 2020.
  - [14] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. A quantum approximate optimization algorithm, 2014.
  - [15] Richard P. Feynman. Simulating physics with computers. *International Journal of Theoretical Physics*, 21(6):467–488, 1982.
  - [16] Craig Gidney, Noah ShuTTY, and Cody Jones. Magic state cultivation: growing t states as cheap as cnot gates, 2024.
  - [17] Daniel Gottesman. Theory of fault-tolerant quantum computation. *Phys. Rev. A*, 57:127–137, Jan 1998.
  - [18] Lov K. Grover. A fast quantum mechanical algorithm for database search. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing*, STOC ’96, page 212–219, New York, NY, USA, 1996. Association for Computing Machinery.
  - [19] W. K. Hastings. Monte carlo sampling methods using markov chains and their applications. *Biometrika*, 57(1):97–109, 1970.
  - [20] Kesha Hietala, Robert Rand, Shih-Han Hung, Xiaodi Wu, and Michael Hicks. A verified optimizer for quantum circuits. *Proc. ACM Program. Lang.*, 5(POPL), jan 2021.
  - [21] Fei Hua, Yanhao Chen, Yuwei Jin, Chi Zhang, Ari Hayes, Youtao Zhang, and Eddy Z. Zhang. Autobraid: A framework for enabling efficient surface code communication in quantum computing. In *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO ’21, page 925–936, New York, NY, USA, 2021. Association for Computing Machinery.
  - [22] IonQ. Ionq forte. <https://ionq.com/quantum-systems/forte>, 2024.
  - [23] Zhihao Jia, Oded Padon, James Thomas, Todd Warszawski, Matei Zaharia, and Alex Aiken. Taso: optimizing deep learning computation with automatic generation of graph substitutions. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, SOSP ’19, page 47–62, New York, NY, USA, 2019. Association for Computing Machinery.
  - [24] Yuwei Jin, Zirui Li, Fei Hua, Tianyi Hao, Huiyang Zhou, Yipeng Huang, and Eddy Z. Zhang. Tetris: A Compilation Framework for VQA Applications in Quantum Computing. 9 2023.
  - [25] Daniel Kahneman. *Thinking, fast and slow*. macmillan, 2011.
  - [26] Chan Gu Kang and Hakjoo Oh. Modular component-based quantum circuit synthesis. *Proc. ACM Program. Lang.*, 7(OOPSLA1), apr 2023.
  - [27] Ilyas Khan and Jenni Strabley. Quantinuum extends its significant lead in quantum computing, achieving historic milestones for hardware fidelity and quantum volume, 4 2024.
  - [28] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
  - [29] Aleks Kissinger and John van de Wetering. PyZX: Large Scale Automated Diagrammatic Reasoning. In Bob Coecke and Matthew Leifer, editors, *Proceedings 16th International Conference on Quantum Physics and Logic*, Chapman University, Orange, CA, USA., 10-14 June 2019, volume 318 of *Electronic Proceedings in Theoretical Computer Science*, pages 229–241. Open Publishing Association, 2020.
  - [30] A Yu Kitaev. Quantum measurements and the abelian stabilizer problem. *arXiv preprint quant-ph/9511026*, 1995.
  - [31] Gushu Li, Anbang Wu, Yunong Shi, Ali Javadi-Abhari, Yufei Ding, and Yuan Xie. Paulihedral: a generalized block-wise compiler optimization framework for quantum simulation kernels. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS ’22, page 554–569, New York, NY, USA, 2022. Association for Computing Machinery.
  - [32] Zikun Li, Jinjun Peng, Yixuan Mei, Sina Lin, Yi Wu, Oded Padon, and Zhihao Jia. Quarl: A learning-based quantum circuit optimizer. *Proc. ACM Program. Lang.*, 8(OOPSLA1), apr 2024.
  - [33] Daniel Litinski. A game of surface codes: Large-scale quantum computing with lattice surgery. *Quantum*, 2018.
  - [34] Daniel Litinski and Felix von Oppen. Lattice Surgery with a Twist: Simplifying Clifford Gates of Surface Codes. *Quantum*, 2:62, May 2018.
  - [35] Yuhao Liu, Shize Che, Junyu Zhou, Yunong Shi, and Gushu Li. Fermihedral: On the optimal compilation for fermion-to-qubit encoding. In

- Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3, ASPLOS '24*, page 382–397, New York, NY, USA, 2024. Association for Computing Machinery.
- [36] Zhengyang Liu, Stefan Mada, and John Regehr. Minotaur: A simd-oriented synthesizing superoptimizer, 2023.
- [37] Henry Massalin. Superoptimizer: a look at the smallest program. *SIGARCH Comput. Archit. News*, 15(5):122–126, oct 1987.
- [38] Matt McEwen, Lara Faoro, Kunal Arya, Andrew Dunsworth, Trent Huang, Seon Kim, Brian Burkett, Austin Fowler, Frank Arute, Joseph C. Bardin, Andreas Bengtsson, Alexander Bilmes, Bob B. Buckley, Nicholas Bushnell, Zijun Chen, Roberto Collins, Sean Demura, Alan R. Derk, Catherine Erickson, Marissa Giustina, Sean D. Harrington, Sabrina Hong, Evan Jeffrey, Julian Kelly, Paul V. Klimov, Fedor Kostritsa, Pavel Laptev, Aditya Locharla, Xiao Mi, Kevin C. Miao, Shirin Montazeri, Josh Mutus, Ofer Naaman, Matthew Neeley, Charles Neill, Alex Opremcak, Chris Quintana, Nicholas Redd, Pedram Roushan, Daniel Sank, Kevin J. Satzinger, Vladimir Shvarts, Theodore White, Z. Jamie Yao, Ping Yeh, Juhwan Yoo, Yu Chen, Vadim Smelyanskiy, John M. Martinis, Hartmut Neven, Anthony Megrant, Lev Ioffe, and Rami Barends. Resolving catastrophic error bursts from cosmic rays in large arrays of superconducting qubits. *Nature Physics*, 18(1):107–111, 2022.
- [39] Manasij Mukherjee, Pranav Kant, Zhengyang Liu, and John Regehr. Dataflow-based pruning for speeding up superoptimization. *Proc. ACM Program. Lang.*, 4(OOPSLA), nov 2020.
- [40] Yunseong Nam, Neil J Ross, Yuan Su, Andrew M Childs, and Dmitri Maslov. Automated optimization of large quantum circuits with continuous parameters. *npj Quantum Information*, 4(1):1–12, 2018.
- [41] Hartmut Neven and Julian Kelly. Suppressing quantum errors by scaling a surface code logical qubit, 2 2023.
- [42] Michael A Nielsen and Isaac Chuang. Quantum computation and quantum information, 2002.
- [43] Anouk Paradis, Jasper Dekoninck, Benjamin Bichsel, and Martin Vechev. Synthetiq: Fast and versatile quantum circuit synthesis. *Proc. ACM Program. Lang.*, 8(OOPSLA1), apr 2024.
- [44] Tirthak Patel, Ed Younis, Costin Iancu, Wibe de Jong, and Devesh Tiwari. Quest: systematically approximating quantum circuits for higher output fidelity. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '22*, page 514–528, New York, NY, USA, 2022. Association for Computing Machinery.
- [45] J. Paykin, A. T. Schmitz, M. Ibrahim, X. Wu, and A. Y. Matsuura. Pcoast: A pauli-based quantum circuit optimization framework. In *2023 IEEE International Conference on Quantum Computing and Engineering (QCE)*, pages 715–726, Los Alamitos, CA, USA, sep 2023. IEEE Computer Society.
- [46] Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J. Love, Alán Aspuru-Guzik, and Jeremy L. O'Brien. A variational eigenvalue solver on a photonic quantum processor. *Nature Communications*, 5(1):4213, 2014.
- [47] Phitchaya Mangpo Phothilimthana, Aditya Thakur, Rastislav Bodik, and Dinakar Dhurjati. Scaling up superoptimization. *SIGARCH Comput. Archit. News*, 44(2):297–310, mar 2016.
- [48] John Preskill. Quantum Computing in the NISQ era and beyond. *Quantum*, 2:79, August 2018.
- [49] Nils Quetschlich, Lukas Burgholzer, and Robert Wille. Mqt predictor: Automatic device selection with device-specific circuit compilation for quantum computing. *ACM Transactions on Quantum Computing*, jun 2024. Just Accepted.
- [50] Péter Rakyta and Zoltán Zimborás. Approaching the theoretical limit in quantum gate decomposition. *Quantum*, 6:710, May 2022.
- [51] Péter Rakyta and Zoltán Zimborás. Efficient quantum gate decomposition via adaptive circuit compression. 3 2022.
- [52] Francisco J. R. Ruiz, Tuomas Laakkonen, Johannes Bausch, Matej Balog, Mohammadamin Barekatain, Francisco J. H. Heras, Alexander Novikov, Nathan Fitzpatrick, Bernardino Romera-Paredes, John van de Wetering, Alhusein Fawzi, Konstantinos Meichanetzidis, and Pushmeet Kohli. Quantum circuit optimization with alphasor, 2024.
- [53] Raimondas Sasnauskas, Yang Chen, Peter Collingbourne, Jeroen Ketema, Jubi Taneja, and John Regehr. Souper: A synthesizing superoptimizer. 2017.
- [54] Eric Schkufza, Rahul Sharma, and Alex Aiken. Stochastic superoptimization. *SIGPLAN Not.*, 48(4):305–316, mar 2013.
- [55] Yunong Shi, Nelson Leung, Pranav Gokhale, Zane Rossi, David I. Schuster, Henry Hoffmann, and Frederic T. Chong. Optimized compilation of aggregated instructions for realistic quantum computers. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '19*, page 1031–1044, New York, NY, USA, 2019. Association for Computing Machinery.
- [56] P.W. Shor. Algorithms for quantum computation: discrete logarithms and factoring. In *Proceedings 35th Annual Symposium on Foundations of Computer Science*, pages 124–134, 1994.
- [57] Seyon Sivarajah, Silas Dilkes, Alexander Cowtan, Will Simmons, Alec Edgington, and Ross Duncan. t|ket>: a retargetable compiler for nisq devices. *Quantum Science and Technology*, 6(1):014003, 2020.
- [58] Luka Skoric, Dan E. Browne, Kenton M. Barnes, Neil I. Gillespie, and Earl T. Campbell. Parallel window decoding enables scalable fault tolerant quantum computation. *Nature Communications*, 14(1):7040, 2023.
- [59] Ethan Smith, Marc Grau Davis, Jeffrey Larson, Ed Younis, Lindsay Bassman Oftelie, Wim Lavrijsen, and Costin Iancu. Leap: Scaling numerical optimization based synthesis using an incremental approach. *ACM Transactions on Quantum Computing*, 4(1), feb 2023.
- [60] IonQ staff. Getting started with native gates. <https://ionq.com/docs/getting-started-with-native-gates>, 2024.
- [61] Rich Sutton. The bitter lesson, 2019. Available at <http://www.incompleteideas.net/InIdeas/BitterLesson.html>.
- [62] Robert R. Tucci. An introduction to cartan's kak decomposition for qc programmers, 2005.
- [63] Valentin F. Turchin. The concept of a supercompiler. *ACM Trans. Program. Lang. Syst.*, 8(3):292–325, June 1986.
- [64] John van de Wetering and Matt Amy. Optimising quantum circuits is generally hard, 2024.
- [65] Xin-Chuan Wu, Marc Grau Davis, Frederic T. Chong, and Costin Iancu. Qgo: Scalable quantum circuit optimization using automated synthesis, 2022.
- [66] Amanda Xu, Abtin Molavi, Lauren Pick, Swamit Tannu, and Aws Albarghouthi. Synthesizing quantum-circuit optimizers. *Proc. ACM Program. Lang.*, 7(PLDI), jun 2023.
- [67] Mingkuan Xu, Zikun Li, Oded Padon, Sina Lin, Jessica Pointing, Auguste Hirth, Henry Ma, Jens Palsberg, Alex Aiken, Umut A. Acar, and Zhihao Jia. Quartz: superoptimization of quantum circuits. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2022*, page 625–640, New York, NY, USA, 2022. Association for Computing Machinery.
- [68] E. Younis, K. Sen, K. Yelick, and C. Iancu. Qfast: Conflating search and numerical optimization for scalable quantum circuit synthesis. In *2021 IEEE International Conference on Quantum Computing and Engineering (QCE)*, pages 232–243, Los Alamitos, CA, USA, oct 2021. IEEE Computer Society.
- [69] Ed Younis, Costin C Iancu, Wim Lavrijsen, Marc Davis, Ethan Smith, and USDOE. Berkeley quantum synthesis toolkit (bqskit) v1, 4 2021.
- [70] Charles Yuan and Michael Carbin. Tower: data structures in quantum superposition. *Proc. ACM Program. Lang.*, 6(OOPSLA2), oct 2022.

[71] Charles Yuan and Michael Carbin. The t-complexity costs of error correction for control flow in quantum computation. *Proc. ACM Program. Lang.*, 8(PLDI), jun 2024.

[72] Alwin Zulehner, Alexandru Paler, and Robert Wille. An efficient methodology for mapping quantum circuits to the ibm qx architectures. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 38(7):1226–1236, 2019.