# AutoComm: A Framework for Enabling Efficient Communication in Distributed Quantum Programs

Anbang Wu
Department of Computer Science
University of California, Santa Barbara
anbang@ucsb.edu

Hezi Zhang
Department of Computer Science
University of California, Santa Barbara
hezi@ucsb.edu

Gushu Li
Department of Electrical & Computer Engineering
University of California, Santa Barbara
gushuli@ece.ucsb.edu

Alireza Shabani
Cisco Research
Los Angeles, California
ashabani@cisco.com

Yuan Xie
Department of Electrical & Computer Engineering
University of California, Santa Barbara
yuanxie@ucsb.edu

Yufei Ding
Department of Computer Science
University of California, Santa Barbara
yufeiding@cs.ucsb.edu

## ABSTRACT

Distributed quantum computing (DQC) is a promising approach to extending the computational power of near-term quantum devices. However, the non-local quantum communication between quantum devices is much more expensive and error-prone than the local quantum communication within each quantum device. Previous work on the DQC communication optimization focus on optimizing the communication protocol for each individual non-local gate and then adopt quantum compilation designs which are designed for local multi-qubit gates (such as controlled-x or CX gates) in a single quantum computer. The communication patterns in distributed quantum programs are not yet well studied, leading to a far-from-optimal communication cost. In this paper, we identify *burst communication*, a specific qubit-node communication pattern that widely exists in many distributed programs and can be leveraged to guide communication overhead optimization. We then propose AutoComm, an automatic compiler framework to first extract the burst communication patterns from the input programs, and then optimize the communication steps of burst communication discovered. Experimental results show that our proposed AutoComm can reduce the communication resource consumption and the program latency by 75.6% and 71.4% on average, respectively.

## 1 Introduction

Quantum computing is promising with its great potential of providing significant speedup to many problems, such as large-number factorization with an exponential speedup [1] and unordered database search with a quadratic speedup [2]. A large number of qubits is required in order to solve practical problems with quantum advantage and the qubit count requirement is even higher after taking quantum error correction [3] into consideration. However, it has turned out that extending the number of qubits on a single quantum processor is exceedingly difficult due to various hardware-level challenges such as crosstalk errors [4, 5], qubit addressability [6], fabrication difficulty [7], etc. The challenges usually increase with the size of quantum hardware and may limit the number of qubits accommodated by a single quantum processor.

Rather than relying on the advancement of a single quantum processor, an alternative way of increasing scalability is by distributed quantum computing (DQC), which integrates the computing resources of multiple modular quantum processors. For example, recent experiments have demonstrated an entanglement-based quantum network of three quantum processors [8]. Companies such as IBM also envision in their roadmap [9] a future of creating a large-scale quantum computer with quantum interconnects that link superconducting quantum processors. Similarly, the ion trap-based quantum computer also requires an optical network of multiple traps each with tens of qubits in-order to scale up, making DQC a path to realizing large-scale quantum computers [10].

In DQC, remote communication involving qubits in different computing nodes is essential yet far more expensive than the local communication on qubits within the same node (e.g., 5-100x time consumption and up to 40x accuracy degradation [11, 12]). There are two major schemes for remote quantum communication: one built upon the cat-entangler and cat-disentangler protocol [13], and the other based on the quantum teleportation [3]. In this paper, we denote the former scheme as Cat-Comm and the latter one as TP-Comm. Both schemes consume EPR pairs [14], which are pre-distributed entangled qubit pairs, as a resource to establish quantum communication. Cat-Comm can implement the remote CX gate [3] with only one EPR pair, but for general two-qubit gates like the SWAP gate [3], Cat-Comm requires up to three EPR pairs [15]. In contrast, TP-Comm conducts any remote two-qubit gate with two EPR pairs [14], making it more efficient for the SWAP gate. For a distributed program, more
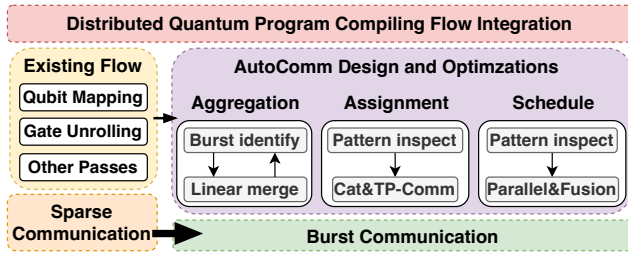
**Figure 1: AutoComm Overview.**

complex remote operations or more information getting transferred per EPR pair would lead to less communication cost.

The overall compiling flow for DQC is similar to that of single-node quantum programs, except with more emphasis on remote communication overhead. Ferrari et al. [15] propose a compiler design similar to single-node compilers [16, 17, 18, 19, 20] using Cat-Comm for each remote CX gate and TP-Comm for each remote SWAP gate. Unsatisfied with the low information of the remote CX gate, Baker et al. [11] eliminate all remote CX gates by using the remote SWAP gate, which only requires two EPR pairs for implementation but contains the information of three CX gates. Unfortunately, bounded by the information of a single two-qubit gate, these compilers cannot achieve higher throughput of information per EPR pair.

Eisert et al. [14] suggest higher throughput could be achieved by considering multi-qubit gates. Diadamo et al. [21] propose a specialized compiler for distributed VQE that uses Cat-Comm to implement controlled-unitary-unitary and controlled-controlled-unitary gates. However, their work can only optimize the gate written in the controlled-unitary form and thus cannot work with decomposed circuits. Moreover, their work cannot optimize programs lacking controlled-unitary blocks.

Besides increasing the 'height' (number of qubits) of remote operations, we observe that the throughput of information per EPR pair can also be significantly boosted up by expanding the 'width' (number of gates) of each remote communication. Specifically, we discover that a large amount of remote two-qubit gates in distributed quantum programs can be implemented collectively through one or two communication invocations. On top of the observation, we propose to optimize the communication overhead based on the *burst communication*, which denotes a group of continuous remote two-qubit gates between one qubit and one node. Burst communication is powerful as it is more information-intensive than single two-qubit gate and contains but not limited to controlled-unitary blocks. Burst communication is also flexible for optimization as it does not require specialized circuit representation and is available in decomposed circuits.

To this end, we develop the first burst-communication-centric optimization framework, *AutoComm* as shown in Figure 1. In contrast to existing compiling flows [11, 15, 16, 17, 18, 19, 20], where each remote CX gate is implemented independently (i.e., sparse communication), *AutoComm* greatly mitigates the communication bottleneck with burst communication and can be easily integrated into these existing compiling flows. Our framework consists of three key stages. Firstly, we perform a communication aggregation pass to group re-

mote gates and extract burst communication blocks. Due to the broad availability of burst communication in distributed quantum programs, this pass could generate a large amount of burst communication blocks for following optimizations. Secondly, we propose a hybrid communication scheme which examines the patterns of each burst communication block and assigns the optimal communication scheme for each block. The insight for this step is that, TP-Comm and Cat-Comm is more resource-efficient for different type of burst communications and considering only one communication scheme would incur extra resource consumption. Finally, we propose an adaptive communication schedule for burst communication blocks of different patterns to squeeze out the parallelism between them and thus reduce overall program latency. There are two critical observations for this optimization: it is possible to execute burst communication with shared qubits or nodes in parallel, and we can fuse some burst communication blocks to cut down the communication footprint.

Our contributions are summarized as follows:

- We identify the burst communication feature in distributed quantum computing and promote its importance in optimizing distributed quantum programs. We further propose the first communication optimization framework based on the burst communication.

- We propose a communication aggregation pass to expose burst communications of distributed quantum programs and then design a hybrid communication scheme, using both Cat-Comm and TP-Comm to accommodate different communication patterns.

- We propose an efficient communication scheduling method to optimize the latency adaptively squeezing out the parallelism of various patterns.

- Compared to the state-of-the-art baseline method [15], AutoComm significantly reduces the communication resource consumption and the program latency by 75.6% and 71.4% on average, respectively.

## 2 Background

In this section, we introduce necessary background to understand the distributed quantum computing and its communication. We do not cover the basic quantum computing concepts (e.g., qubit, gate, measurement) and recommend [3] for more details.

### 2.1 EPR Pair and Entanglement

*EPR entanglement* To establish quantum communication in a distributed quantum computer, we first need to generate a pair of qubits whose state is $\frac{1}{\sqrt{2}}(|00\rangle + |11\rangle)$, EPR entangled state. The two qubits such state is called EPR entanglement pair (Abbrev., EPR pair) [3]. The two qubits of an EPR pair can be distributed on different quantum devices, formulating a remote EPR pair [14]. The preparation of the remote EPR pair includes two stages: generation and purification. The generation stage generates and distributes EPR pairs but is very noisy, making the purification stage indispensable [22].

## 2.2 Distributed Quantum Computing

The development of quantum communication [8,23,24,25,26, 27,28,29,30,31,32,33] enables distributed computing over a series of quantum devices. As in classical distributed computing, remote communication between computing nodes is also the bottleneck of distributed quantum computing (DQC) and should be carefully optimized.

Different from the classical distributed computing system, quantum data cannot be easily shared across quantum nodes due to the no-cloning theorem [3]. The workaround is to exploit different communication schemes (e.g., *Cat-Comm* [14] and *TP-Comm* [3]) based on remote EPR entanglement, one of the key information resources in quantum processing. Figure 2 illustrates how to use these two schemes to implement one *remote CX gate*, with the control qubit $q_1$ residing in quantum nodes A and the target qubit $q'_1$ in node B. Qubits in Figure 2 fall into two categories. The first category of qubits is used to store quantum information and is called *data qubits*, e.g., $q_1$ and $q'_1$. The second category of qubits, called *communication qubits*, is used to hold the remote EPR entanglement required for quantum communication, e.g. $q_0$ and $q'_0$ in Figure 2.

As shown in Figure 2(a), the first communication scheme Cat-Comm utilizes cat-entangler to transfer the state of the control qubit $q_1$ to node B, execute the target CX gate, and then use cat-disentangler to transfer the state back to node A. While TP-Comm, the second communication scheme in Figure 2(b), employs quantum teleportation [3] to transfer the state of $q_1$, and then execute the target CX gate. Though Cat-Comm and TP-Comm both require one EPR pair and two bits of classical communication, Cat-Comm is more widely-used than TP-Comm in DQC compilers [15,21]. This is mainly due to the dirty side-effect of TP-Comm. We would need another invocation of TP-Comm to release the occupation of the communication qubit (e.g., $q'_0$ in Figure 2(b)), which would be later used for other quantum communications. As a result, two EPR pairs are actually required to implement a single remote CX gate by TP-Comm, with one pair for handling the dirty side-effect.
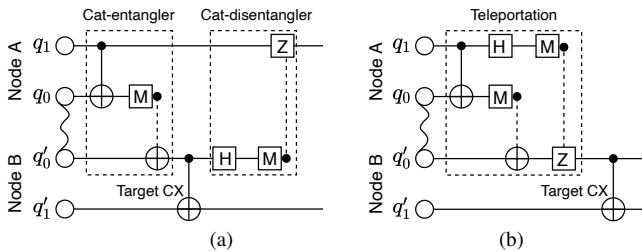


Figure 2: The implementation of one remote CX. (a) The Cat-Comm version. (b) The TP-Comm version. Each wavy line denotes an EPR pair between qubits, and each dashed line denotes one bit of classical communication. M denotes measurement.

In Figure 2, we only show how to implement one individual CX gate. To implement complex remote interactions between quantum nodes, one simple strategy is to first decompose the remote interaction into several remote CX gates and implement each remote CX gate as in Figure 2. However, this strat-
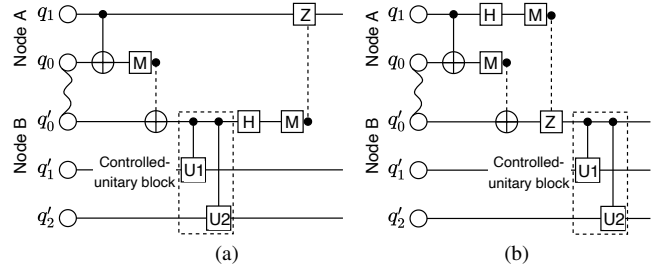


Figure 3: The optimized implementation of the controlled-unitary block $C - U1 - U2$. (a) The Cat-Comm version. (b) The TP-Comm version.
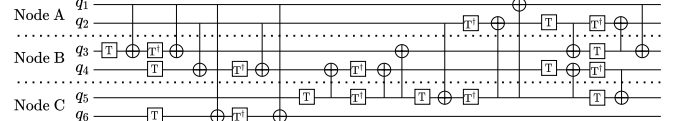


Figure 4: Program snippet extracted from quantum arithmetic circuits [36].

egy may incur heavy communication costs. Prior work [13] spots a more efficient way to implement a controlled-unitary block between two quantum nodes. Figure 3 provides the optimized implementation of the controlled-block $C - U1 - U2$, where $U1$ and $U2$ are some unitary quantum operations. The implementation in Figure 3 only requires one EPR pair, fewer than implementing each remote two-qubit gate independently.

Besides the controlled-unitary block, we discover that plenty of quantum communications in distributed quantum programs can be transformed into a group of remote interactions between one qubit and one quantum node. We name such a group of remote interactions *burst communication*. Different from the single CX case, Cat-Comm and TP-Comm each has its own advantage for burst communication of various patterns. Unfortunately, existing DQC compilers [11,34] either do not take advantage of the burst communication or only consider the basic controlled-unitary case [21].

In later sections, we would use *one remote EPR pair* and *one remote communication* interchangeably, because for either Cat-Comm or TP-Comm, one invocation just requires one remote EPR pair.

## 3 Problem and Motivation

In this section, we first introduce the communication problem in distributed quantum programs and then identify the optimization opportunities by considering burst communication.

For the rest of the discussions, we assume quantum communication can be established between any two quantum nodes, a typical assumption in data-center distributed computing [35]. We also assume that each quantum node has only two communication qubits, which is realistic for near-term DQC [15].

### 3.1 Communication Problem

The example distributed program in Figure 4 is modified from quantum arithmetic circuits [36]. This program contains many remote CX gates whose control qubit and target

qubit reside in different quantum nodes, e.g., $CX\ q_1, q_3$. Remote CX gates are inevitable in DQC especially when the program's qubit number is substantially larger than each quantum node's. To make the distributed program executable, we should transfer the states of qubits in remote CX gates to make them locally executable temporarily. The state transfer involves remote communication between quantum nodes, which can be accomplished by Cat-Comm or TP-Comm. Due to the noisy nature of quantum communication, remote operations are far more error-prone than local quantum gates. The long runtime of quantum communication would also lead to the decoherence of quantum states. As a result, to produce high fidelity outcome, we hope the number of remote communication to be as small as possible, so is the latency induced.

As indicated in Section 2, one remote CX gate requires at least one remote communication. While there is little room for optimizing the communication cost of one remote CX gate, there is a large optimization space when considering burst communication, which involves a group of remote CX gates. For example, we can execute the first two CX gates on $q_1, q_3$ in Figure 4 collectively, with only one communication by using the circuit in Figure 3(a). From the perspective of information theory, burst communication is more informative than the communication with only one remote CX. The overall communication cost and latency would be considerably lowered if handling all remote CX gates in this burst manner.

Fortunately, as we see in the next section, burst communication is prevalent in diverse distributed quantum programs.

## 3.2 Burst Communication in DQC

Aside from the arithmetic program shown in Figure 4, we also see burst communication in a variety of quantum programs. As examples, we examine the burst communication of the Quantum Fourier Transform (QFT) program [3] and the Quantum Approximate Optimization Algorithm (QAOA) [37] by hand. These two represent different categories of quantum programs: QAOA is one of the most important applications in near-term quantum computing whereas QFT is the building block circuit of quantum algorithms.

We first give a formal definition of the burst communication in DQC. In this paper, we refer to a group of continuous remote two-qubit gates between one qubit $q$ and one node as *burst communication*. For two remote two-qubit gates $g_1$ and $g_2$, the continuity of these two gates means there are no other remote gates between $g_1$ and $g_2$.

To characterize the burst communication of a distributed program $dprog$, for a remote gate $g$ in $dprog$, we define function $\varepsilon(g)$ to be the largest burst communication block that contains $g$. The gate order of $dprog$ may affect the burst communication block found. $\varepsilon(g)$ is defined to be the largest over all functional-equivalent gate order of $dprog$. We then define $len(\varepsilon(g))$ to be the number of remote CX gates in $\varepsilon(g)$ if compiled to the CX+U3 basis [17]. Finally, we are ready to define the inverse-burst distribution as follows:

$$P(x) = \frac{|\{g|len(\varepsilon(g)) < x\}|}{\#g}. \tag{1}$$

A lower $P(x)$ suggests more burst communication.

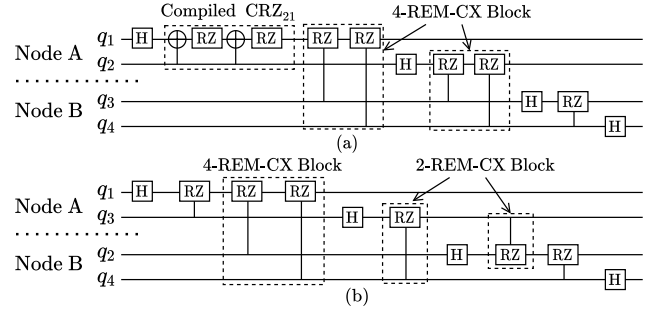We begin by examining the QFT program using the aforementioned definition. We assume the total qubit number is



Figure 5: **(a) QFT program with two nodes and two qubits per node. (b) The layout for the maximal $P_4$. Parameters omitted for simplicity. For demonstration, we do not combine $CRZ_{43}$ and $CRZ_{32}$ to form a 4-REM-CX block.**

$n$, the quantum node number is $k$, and qubits are evenly distributed across all nodes, with $t = \frac{n}{k}$ qubits per node. Figure 5 shows the QFT program with $k = 2$ and $t = 2$. For the QFT program, as shown in Figure 5, each $q_i$ is controlled by all qubits $q_j$ (through the CRZ gate) that satisfies $j > i$ [3]. First, we have $P(2) = 0$ because each CRZ gate in QFT is compiled into two CX gates, as illustrated in Figure 5(a). Now, we consider $P(4)$. For the $i$-th qubit satisfies $i \leq n - k$, the number of $j$ s.t. $\varepsilon(CRZ_{ji}) < 4$ is at most $\lfloor \frac{i-1}{t-1} \rfloor$ because for one node, if at least two of its qubits have subscripts $> i$, this node would have at least two qubits being interacted by qubit $i$. Since CRZ gates are commutable with each other, we could form a communication block with at least 4 CX gates. On the other hand, if $i > n - k$, then the $i$-th qubit is at most interacted with $n - i$ qubits, thus the number of $j$ s.t. $\varepsilon(CRZ_{ji}) < 4$ is at most $n - i$. Therefore, we have

$$P(4) \leq \frac{\sum_{i=1}^{n-k} \lfloor \frac{i-1}{t-1} \rfloor + \sum_{i=n-k+1}^{n}(n-k)}{\sum_{i=1}^{n}(n-i) - k\sum_{l=1}^{t}(t-l)} = \frac{1}{t}.$$

This indicates there are $1 - P(4) = 1 - \frac{1}{t}$ remote gates within a communication block that possesses more than 4 CX gates. Generally, we can prove that $P(2m) \leq \frac{m-1}{t}$. This upper bound is quite promising when $t$ is large and it is actually loose. For Figure 5(b) which corresponds to the upper bound of $P(4)$, there may be $\frac{1}{t}$ of remote CRZ gates, i.e., $CRZ_{43}$ and $CRZ_{32}$ not in a block with 4 remote CX gates at the first glance. But we can actually combine $CRZ_{43}$ and $CRZ_{32}$ to form a 4-REM-CX block since there are no other remote gates between them. This indicates that QFT has more abundant burst communication than the upper bound suggests.

Similarly, for the QAOA program, we assume $k$ nodes and $t$ qubits per node. We also suppose $r$ remote ZZ interactions between any two nodes. Figure 6 shows the QAOA program with $k = 2$ and $t = 3$. Likewise, $P(2) = 0$ since each ZZ interaction is compiled into two CX gates, as shown in Figure 6(a). For every two nodes, the qubit layout to minimize $len(\varepsilon(ZZ))$ for each ZZ interaction is to make every two ZZ interactions have no shared qubits, i.e., not adjacent. However, this layout at most accommodates $t$ ZZ interactions. If $r > t$, the number of ZZ interactions s.t. $len(\varepsilon(ZZ)) < 4$ is at most $t - 2(r \bmod t)$ by examining the gate adjacency.
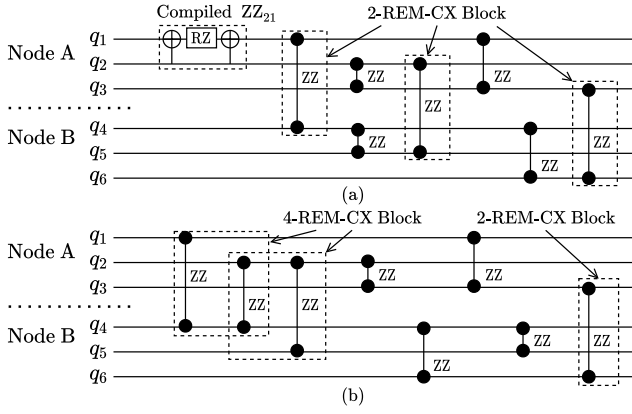
**Figure 6: QAOA program with two nodes and three qubits per node. Parameters omitted for simplicity. (a) inter-node communication number $r = 3$. (b) $r = 4$.**

Thus, $P(4) \leq \frac{t-2(r \bmod t)}{r}$. For example in Figure 6(b), only $\frac{t-2(r \bmod t)}{r} = \frac{1}{4}$ of remote ZZ interactions are not in a 4-REM-CX block. Generally, if $r > st$ for some integer $s$, $P(2(s+1)) \leq \frac{t-2(r \bmod t)}{r} < \frac{1}{s}$. This study reveals that burst communication is broadly available in the QAOA program.

We could derive a similar analysis for other programs. Further numerical evidence for the richness of burst communication in various programs is shown in Figure 15. The next step is to figure out how to utilize the abundant burst communication in distributed programs to optimize the communication overhead, as discussed in the next section.

## 3.3 Optimization Opportunities

To exploit burst communication in distributed quantum programs, we need to answer three key questions:

*How to unveil the burst communication?* The burst communication is high-level program information and cannot be deduced simply from the low-level circuit language, especially when the remote interactions between multiple nodes are all mixed together. For example in Figure 4, gate $CX\ q_2; q_4$ between node A and node B is followed by $CX\ q_1; q_6$, which is the interaction between node A and node C. To maximize the benefits of burst communication, we need to discover groups of remote gates in disordered quantum circuits.

*How to select the best communication scheme?* Burst communication comes in various forms. Cat-Comm may not always be better than TP-Comm for burst communication, unlike the single CX case. For example in Figure 4, if we use Cat-Comm to implement the last three remote CX gates between $q_3$ and node A, three EPR pairs are needed. However, with TP-Comm to teleport $q_3$ to node A, at most two EPR pairs are needed. Thus, to reduce the communication cost, we should examine the pattern of burst communication and choose the communication scheme wisely.

*How to schedule burst communication?* Finally, we need to schedule the execution of burst communication blocks. If we arrange all burst communication in a sequential way, the large time overhead would impose non-negligible decoherence errors on quantum states. As a result, we should
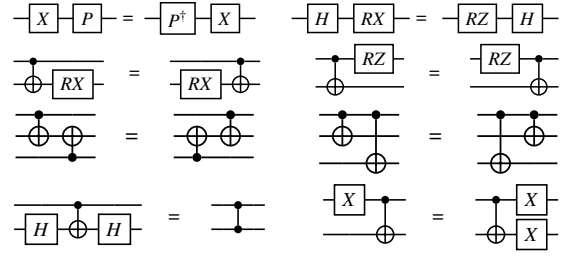


**Figure 7: Phase gates P includes $Z, RZ, S, T$, etc. X-rotation-centered rules for gate commutation.**

maximize the parallelism in burst communication to generate high-fidelity output. To achieve this goal, we must first identify the relationships between communication blocks and then reduce the gaps caused by them adaptively.

## 4 AutoComm Framework

In this section, we first give an overview of the AutoComm framework and then introduce each component in detail.

### 4.1 Design overview

We propose the *AutoComm* framework as shown in Figure 1. AutoComm focuses on the communication optimization of distributed quantum programs and serves as the back-end of front compiling flows like mapping qubits to quantum nodes. We would adopt existing technologies for these front compiling stages, as we would see in Section 5.

To optimize the communication overhead in distributed programs, AutoComm comes with three stages to utilize the burst communication. First, it aggregates remote two-qubit gates by gate commutation. Gate commutation is quite common in quantum programs [38]. Commutable gates, on the one hand, may be ordered arbitrarily and hide the burst communication. On the other hand, we could also utilize gate commutation to uncover burst communication blocks. In this stage, a pre-processing step is used to identify burst communication, and a linear merge step is employed to combine isolated burst communication blocks.

Second, it assigns an optimal communication scheme for each burst communication. We observe that the pattern of burst communication impacts the efficiency of communication schemes. Cat-Comm is less expensive for some patterns, while TP-Comm may be more cost-effective for others. It is thus important to examine the communication patterns and consider both Cat-Comm and TP-Comm for hybrid communication, rather than focusing on one scheme.

Third, it performs a block-level schedule of burst communication. It is possible to run communication blocks with shared nodes or qubits concurrently or shorten the quantum state transfer path across quantum nodes for specified communication patterns. Combined with these optimizations, a greedy schedule is effective for burst communication blocks.

### 4.2 Communication Aggregating

Burst communication is prevalent in distributed programs, but may not be immediately available due to two factors: CX gates may be scattered across the program, and whether CX gates are remote depends on the qubit mapping to quantum

5

**Algorithm 1:** Linear merge procedure

**Input:** An array of communication blocks *blk_list*
**Output:** Merged communication blocks *blk_list_merge*

1 *blk_list_merge* = [ ] ;
2 *blk* = *blk_list*[0] ;
3 **while** *there are blocks in blk_list not visited* **do**
4    *non_commute_gates* = [ ];
5    **for** *blk_next in unvisited blocks of blk_list* **do**
6       // Attempt merge *blk* to *blk_next*
      **for** *gate between blk and blk_next* **do**
7          **if** *gate is single-qubit and not commutes with blk* **then**
8             *non_commute_gates*.append(*gate*);
9          **if** *gate is two-qubit* **then**
10             check if *gate* is commutable with *non_commute_gates* and *blk*;
11             **if** *not commutable* **then**
12                **if** *gate is in-node two-qubit* **then**
13                   *non_commute_gates*.append(*gate*);
14                **else**
15                   break;
16       **end**
17       *blk* = merge *blk*, *non_commute_gates* and *blk_next*;
18    **end**
19    **if** *the above merge failed* **then**
20       Try to merge *blk_next* to *blk* similarly;
21       **if** *succeeds* **then**
22          *blk* = merge *blk*, *non_commute_gates* and *blk_next*;
23       **else**
24          *blk* = *blk_next*;
25 **end**
26 output the merged blocks and adjust the order of commutable gates;



**Figure 8: Communication aggregation for the example program in Figure 4. (a) Preprocessing. (b) Linear merge. (c) Iterative refinement.**

nodes. To uncover hidden burst communications, we need to rewrite the circuit and aggregate remote CX gates.

Figure 7 summarizes the X-rotation-centered rewriting rules used for gate commutation. Rules for other rotation axes can be obtained by similar transformation. Below are the main steps to aggregate remote gates based on these rules.

*Preprocessing:* The first step is to identify the qubit-node pair of burst communication. We start with the qubit-node pair associated with the most remote gates as it would likely lead to a large burst communication block. For example in Figure 4, the chosen qubit-node pair is ($q_3$, node A) as it is associated with 5 remote CX gates. We then search for consecutive remote CX gates related to this qubit-node pair. This step would result in many isolated communication blocks, for example in Figure 8(a), we obtain four small blocks.

*Linear merge:* The next step is to merge isolated small communication blocks obtained in the preprocessing. As illustrated in Algorithm 1, we merge related communication blocks in a linear and greedy manner. For communication blocks ①, ②, ③, ④ in Figure 8(a), we can easily merge block ① and ② since only single-qubit gates exist between those two blocks. However, we can not merge block ② and block ③ because gate $CX\ q_5, q_3$ is commutable with neither block ② nor block ③. Finally, as shown in Figure 8(b), we obtain two larger communication blocks.
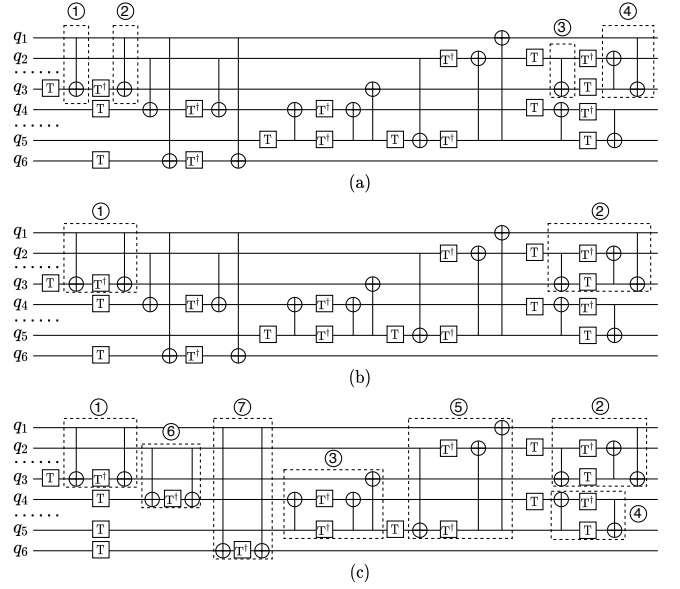
*Iterative refinement:* Then we merge communication blocks of other qubit-node pairs in descending order of their number of remote gates until no improvement is made. The final result of communication aggregation is shown in Figure 8(c).

## 4.3 Communication Assignment

With burst communication blocks, the next optimization is to find the best way to execute them. We address this problem by first examining the pros and cons of Cat-Comm and TP-Comm, and then assigning the optimal communication scheme based on the pattern analysis of burst communication blocks. Since we assume only two communication qubits in each quantum node, the communication patterns discussed here center on interactions between one qubit and one node. Extending burst communication to the node-to-node situation is promising when communication qubits are plentiful. We leave it for future work.

*Cat-Comm vs. TP-Comm:* Suppose we have a burst communication block between a qubit $q_1$ in node A and several qubits in node B, with a total of *n* remote CX gates in the block. If the block can be executed by a single call to Cat-Comm, the savings on EPR pairs would be up to *n* times, compared to executing each remote CX gate individually. However, as discussed below, not all communication blocks can be cheaply executed via Cat-Comm. Compared to Cat-Comm, the savings on ERP pairs with TP-Comm is at most $\frac{n}{2}$ times as TP-Comm requires two EPR pairs to execute any burst communication block: one to teleport $q_1$ to node B, the other to release the occupancy of $q_1$ on the communication qubit in node B. For simplicity, we use the other EPR pair to teleport $q_1$ back to node A. We postpone to Section 4.4 to handle the case that teleporting $q_1$ to some other node is better than moving back. Overall, Cat-Comm provides higher ERP pair savings for specific burst communication blocks, while
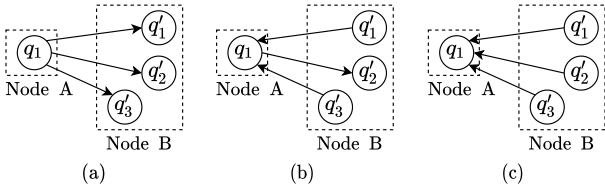
**Figure 9: Two primitive communication patterns (a)(b) and the variant (c).**
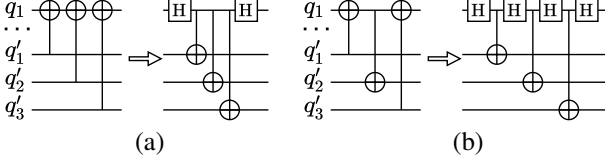


**Figure 10: The transformation between communication patterns by using Hadamard gates.**

TP-Comm can handle an arbitrary communication block with up to two EPR pairs.

*Pattern analysis:* Figure 9(a)(b) shows two primitive patterns for qubit-to-node burst communication. For the unidirectional communication pattern in Figure 9(a) where one qubit (i.e., $q_1$) always serves as the control qubit, the communication block can be implemented by Cat-Comm with only one EPR pair if no single-qubit gate on the control qubit separates two-qubit gates [13]. For example, one call of Cat-Comm can handle the gate sequence $CX\,q_1, q_1'; CX\,q_1, q_2'$, but cannot address $CX\,q_1, q_1'; H\,q_1; CX\,q_1, q_2'$ due to the middle H gate. To optimize this communication pattern with Cat-Comm, we should remove single-qubit gates on the control qubit. When they are not removable, we resort to TP-Comm.

A varied unidirectional pattern in which $q_1$ always serves as the target qubit, as shown in Figure 9(c), also occurs frequently in distributed quantum programs. This pattern can be transformed into the pattern in Figure 9(a) by applying a series of Hadamard gates, as shown in Figure 10(a).

In contrast to unidirectional patterns, Figure 9(b) shows a bidirectional pattern in which $q_1$ serves as both control qubit and target qubit. A block in this pattern cannot be executed by a single call of Cat-Comm as Cat-Comm cannot transfer the state of target qubits. Even if we transform it to the unidirectional pattern in Figure 10(b) with Hadamard gates, single-qubit gates on the control qubit still prevent a cheap implementation by Cat-Comm. In fact, for the block in Figure 10(b), TP-Comm is more efficient as it only requires two EPR pairs, while Cat-Comm requires three EPR pairs.

To summarize, for unidirectional patterns in Figure 9(a)(c), we will try Cat-Comm first, while for the bidirectional pattern in Figure 9(b), TP-Comm is preferred.

*Communication assignment:* Now, we are ready to assign an optimal communication scheme, either Cat-Comm or TP-Comm, to each burst communication block. Considering Figure 8(c) as an example, we assign Cat-Comm to unidirectional blocks ①, ⑥ and ⑦, and assign TP-Comm to bidirectional blocks ②, ④ and ⑤. For ③, although being unidirectional, it cannot be executed by one call of Cat-Comm

as there is a $T^\dagger$ gate on the control qubit between two CX gates. Since executing it with either Cat-Comm or TP-Comm requires two EPR pairs, we set the TP-Comm assignment as default. The finalized assignment is shown in Figure 11(a).

## 4.4 Communication Scheduling

After optimizing the count of remote communications, we then schedule the execution of burst communication blocks to reduce the total execution time of the distributed program and reduce the impact of decoherence. Based on the quantitative data shown in Table 1, the preparation of remote EPR pairs is the most time-consuming one among various operations and hence should be carefully optimized to hide its latency. While the quantitative data may vary across quantum devices, the schedule design in this section should be also effective.

| Operation | Variable Name | Latency |
|---|---|---|
| Single-qubit gates | $t_{1q}$ | $\sim 0.1$ CX |
| CX and CZ gates | $t_{2q}$ | 1 CX |
| Measure | $t_{ms}$ | 5 CX |
| EPR preparation | $t_{ep}$ | $\sim 12$ CX |
| One-bit classical comm | $t_{cb}$ | $\sim 1$ CX |

**Table 1: The quantitative latency data of operations in distributed quantum programs, extracted from [22, 39]. All latencies are normalized to CX counts.**

The designs here aim to maximize block-level parallelism and shorten the latency of sequential execution by fusion.

*More block-level parallelism:* The essence of scheduling is to maximize the parallelism in a circuit. For burst communication blocks without nodes or qubits in common, they can be concurrently executed in nature. For blocks with shared nodes or qubits, their parallelism is limited by their commutability, as well as the communication resources each node holds. With the constraint that each node can establish only two communications in parallel, there is little room for lazy operations, and we adopt a greedy strategy to execute commutable blocks, i.e., execute as many blocks as possible simultaneously, as soon as EPR pairs are prepared.

For Cat-Comm blocks, we can execute two commutable blocks in parallel at most if they share nodes, as shown in Figure 12. For TP-Comm blocks, the situation is complex as each TP-Comm blocks require two EPR pairs. For two commutable TP-blocks, rather than prioritizing the completion of one TP-comm as in Figure 13(a), we observe that parallelism can be enabled by communication alignment, as shown in Figure 13(b). Compared to Figure 13(a), Figure 13(b) aligns the qubit teleportation of the two blocks, leading to a latency saving of $t_{block} + 2t_{tele}$. This TP-Comm alignment technique can be generalized to the case of $n$ commutable TP-Comm blocks (any two blocks may share common nodes). With TP-Comm alignment, the total latency saving can be up to $(n-1)(t_{block} + 2t_{tele})$ (e.g., if those TP-Comm blocks are on nodes $\{A_1, A_2\}, \{A_2, A_3\}, \cdots, \{A_n, A_{n+1}\}$ respectively).

*Fusion of sequential blocks:* Sometimes communication blocks have to be executed in sequence. However, if the teleported qubits in TP-Comm blocks are the same, we can optimize their executions by fusing the teleportations, as shown in Figure 14. Figure 14(a) shows a simple schedule where each TP-Comm is executed independently. As each
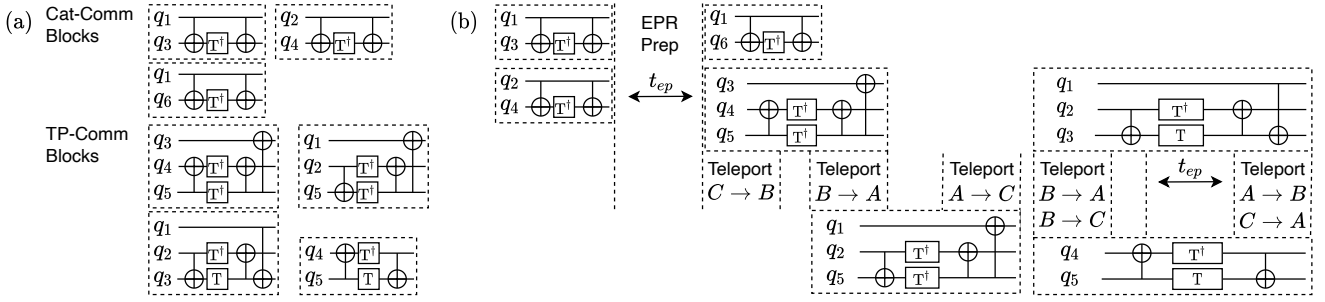
**Figure 11: (a) The result of the communication assignment pass. (b) The result of the communication schedule pass.**
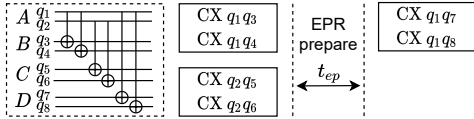


**Figure 12: The schedule optimization for commutable Cat-Comm blocks, with shared qubit or node.**



**Figure 13: The schedule optimization for TP-Comm blocks. Aligned qubit teleportation in (b) is better than the independent qubit teleportation in (a).**
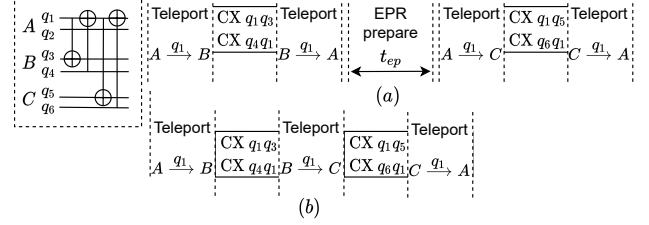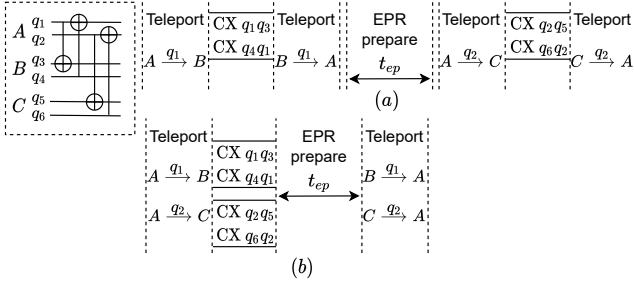


**Figure 14: The schedule optimization for TP-Comm blocks. Cyclic qubit teleportation in (b) is better than the SWAP-style qubit teleportation in (a).**

achieved compared to executing each remote CX gate independently.

## 5 Evaluation

In this section, we first compare the performance of Auto-Comm to the baseline method and then evaluate the effect of each optimization pass in AutoComm. We finally perform a sensitivity analysis on AutoComm to study how its performance evolves as the program configuration changes.

### 5.1 Experiment Setup

*Metric* The first metric we considered is the number of issued remote communications. Each remote communication would consume one remote EPR pair for both Cat-Comm and TP-Comm. To avoid the ambiguity on the cost of TP-Comm, we say TP-Comm needs two communications (i.e., two EPR pairs) to execute one burst communication block, with one of the communications handling its dirty side-effect. The number of remote communications models the resource overhead of executing distributed programs and a lower value is favored.

The second metric is the maximum number of remote two-qubit gates executed through one communication. For TP-Comm blocks, this number is averaged on two communications. We denote this metric by 'Peak # REM CX'. This metric models the communication throughput of information and a higher value is preferred.

Finally, we consider two metrics that model the relative performance, in communication cost and program latency respectively, of AutoComm to baselines. The first one is 'improv. factor', which is defined to be 'total communication # by baseline/total communication # by AutoComm'. The second one is 'LAT-DEC factor' and is defined to be 'program

node has only two communication qubits, we need to wait for $t_{ep}$ before executing the next TP-Comm block. In contrast, Figure 14(b) fuses the teleportations between quantum nodes, forming a cycle: $A \to B \to C \to A$. With fusion, the number of teleportations is reduced by one and the overall execution time is reduced by $t_{ep} + t_{tele}$, where $t_{tele}$ is the time to teleport one qubit (about 8 CX time as shown in Figure 2(b)). Generally, if we have *n* TP-Comm blocks with the same teleported qubit, the total number of teleportation would be reduced by $n - 1$, and the saving of overall latency would be $(n - 1)(t_{ep} + t_{tele})$.

From another view, the fusion also optimizes the token passing problem in classical distributed computing [35], which also appears in Section 4.3, about whether to move the teleported qubit in TP-Comm back or to another node.

With the designs above, the communication schedule pass should apply block-level commutation analysis to unveil the patterns discussed above and then apply corresponding optimizations. We omit the details since this procedure is very similar to the communication aggregation except working at the block level. With all those optimizations applied, Figure 11(b) shows the optimized communication schedule for the program in Figure 4. In total, 2.4x latency saving is

| Type | Name | # qubit | # node | # gate | # CX | # REM CX |
|---|---|---|---|---|---|---|
| Build-ing Blocks | Multi-Controlled Gate (MCTR) | 100 | 10 | 10640 | 4560 | 1680 |
| | | 200 | 20 | 21840 | 9360 | 3568 |
| | | 300 | 30 | 33040 | 14160 | 5632 |
| | Ripple-Carry Adder (RCA) | 100 | 10 | 1569 | 785 | 220 |
| | | 200 | 20 | 3169 | 1585 | 662 |
| | | 300 | 30 | 4769 | 2385 | 820 |
| | Quantum Fourier Transform (QFT) | 100 | 10 | 40100 | 20000 | 9000 |
| | | 200 | 20 | 160200 | 80000 | 38000 |
| | | 300 | 30 | 360300 | 180000 | 87000 |
| Real World Appli-cations | Bernstein Vazirani (BV) | 100 | 10 | 265 | 65 | 56 |
| | | 200 | 20 | 535 | 135 | 126 |
| | | 300 | 30 | 803 | 203 | 194 |
| | QAOA | 100 | 10 | 6000 | 4000 | 3144 |
| | | 200 | 20 | 24000 | 16000 | 14076 |
| | | 300 | 30 | 54000 | 36000 | 32896 |
| | UCCSD | 8 | 4 | 3129 | 1420 | 900 |
| | | 12 | 6 | 40659 | 19142 | 15136 |
| | | 16 | 8 | 129829 | 64956 | 53426 |

**Table 2: Benchmark programs. The qubits are evenly distributed across quantum nodes. The number of remote CX gates (# REM CX) is computed on the qubit mapping by 'Static Overall Extreme Exchange' [11].**

latency by baseline/program latency by AutoComm'. We hope these two metrics to be as large as possible.

*Baseline* For the baseline method, we implement the compiler [15] which only exploits the Cat-Comm scheme for remote CX gates and does not consider burst communication. To reduce the program latency, the baseline adopts a greedy scheduling method, i.e., executing operations as soon as possible. For both the baseline and AutoComm, we map qubits to distributed quantum nodes by the 'Static Overall Extreme Exchange' strategy studied in [11].

*Platforms* We perform all experiments on a Ubuntu 18.04 server with a 6-core Intel E5-2603v4 CPU and 32GB RAM. Other software includes Python 3.8.3 and Qiskit 0.18.3 [17].

*Benchmark programs* We consider two categories of benchmark programs, as shown in Table 2. The first category of benchmarks are function-specific, i.e., they focus on implementing specific elementary functions, e.g., arithmetic operations and Fourier transformation. These quantum programs are often used as the building blocks of large quantum applications. The second category of benchmarks are application-driven. These programs usually target at solving real-world problems, e.g., Bernstein-Vazirani (BV) algorithm, Quantum Approximate Optimization algorithm (QAOA), and Unitary Coupled Cluster ansatzes (UCCSD). Specifically, we choose the graph maxcut problem for the QAOA test programs, and for the UCCSD programs, we select molecules $LiH, BeH_2$, and $CH_4$ which correspond to 8, 12, and 16 particles (thus qubits), respectively. All benchmark programs used in the evaluation are collected from IBM Qiskit [17] and RevLib [36].

## 5.2 Compared to Baseline

We evaluate both AutoComm and the baseline method on benchmark programs in Table 2. The results of AutoComm and its relative performance to the baseline are shown in Table 3.
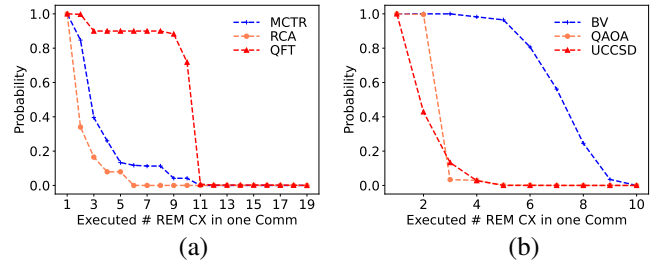


**Figure 15: Burst communications by AutoComm: Pr[X] = Pr[one communication carries >= X REM-CXs].**

*Burst communication statistics:* Figure 15 shows the distribution of burst communications assembled by AutoComm. This distribution is closely related to the inverse-burst distribution discussed in Section 3.2 but is easier to compute. We can see that burst communications exist widely in various distributed quantum programs, no matter in building-block circuits (Figure 15(a)) or in real-world applications (Figure 15(b)). Moreover, Figure 15 demonstrates the effectiveness of AutoComm in unveiling burst communications. In Figure 15, the communications that each carries $\geq 2$ remote CX gates account for 76.8% of the total remote communications, on average.

*Communication cost:* AutoComm achieves significant communication cost reduction on the benchmark programs. Compared to the baseline method, AutoComm reduces the number of remote communications by a factor of 4.1x on average, up to 9.2x. The peak communication throughput (i.e., 'Peak # REM CX') by AutoComm is 8.8x on average and up to 18x of that by the baseline. These improvements indicate that AutoComm can efficiently discover and utilize burst communications in distributed quantum programs, transferring more information in each communication than the baseline method.

The good communication performance of AutoComm comes from two factors: the aggregation of remote CX gates and the hybrid implementation of burst communications by using both Cat-Comm and TP-Comm. We will further elaborate on this point in Section 5.4.

*Latency:* AutoComm also achieves significant latency reduction on benchmark programs. Compared to the baseline method, AutoComm reduces the program execution time by a factor of 3.5x on average, up to 7.1x, as shown in Table 3. The trend of latency reduction is closely related to the trend of communication cost reduction. This is as expected because AutoComm keeps the local parallelism in the program when aggregating remote interactions.

## 5.3 Compared to GP-based Compiler

We further compare AutoComm to the graph-partition-based (GP-based) compiler [11]. A GP-based compiler converts remote interactions to local interactions by swapping qubits with a strategy derived from graph partition algorithms. To reduce the communication cost and program latency of the GP-based compiler, we utilize TP-Comm for swapping qubits since TP-Comm requires only two communications for one remote SWAP gate, one communication less than using Cat-

| Name | Tot Comm | TP-Comm | Peak # REM CX | Improv. factor | LAT-DEC factor |
|---|---|---|---|---|---|
| MCTR-100-10 | 533 | 220 | 10 | 3.15 | 3.27 |
| MCTR-200-20 | 972 | 418 | 10 | 3.67 | 3.83 |
| MCTR-300-30 | 2044 | 1112 | 10 | 2.76 | 2.88 |
| RCA-100-10 | 79 | 54 | 5.5 | 2.78 | 3.34 |
| RCA-200-20 | 469 | 224 | 5.5 | 1.41 | 2.10 |
| RCA-300-30 | 410 | 204 | 5.5 | 2.00 | 3.30 |
| QFT-100-10 | 2068 | 1784 | 18 | 8.70 | 6.53 |
| QFT-200-20 | 8351 | 7566 | 18 | 9.10 | 6.98 |
| QFT-300-30 | 18835 | 17348 | 18 | 9.24 | 7.13 |
| BV-100-10 | 9 | 0 | 8 | 6.22 | 4.33 |
| BV-200-20 | 19 | 0 | 8 | 6.63 | 4.63 |
| BV-300-30 | 29 | 0 | 8 | 6.69 | 4.69 |
| QAOA-100-10 | 1448 | 266 | 6 | 2.17 | 1.83 |
| QAOA-200-20 | 6787 | 728 | 8 | 2.07 | 1.79 |
| QAOA-300-30 | 16053 | 1138 | 6 | 2.05 | 1.69 |
| UCCSD-8-4 | 464 | 0 | 4 | 1.94 | 1.74 |
| UCCSD-12-6 | 8973 | 0 | 4 | 1.69 | 1.55 |
| UCCSD-16-8 | 33303 | 0 | 5 | 1.60 | 1.50 |

**Table 3: The results of AutoComm and its relative performance to the baseline. The name column are acronyms of test programs in Table 2.**
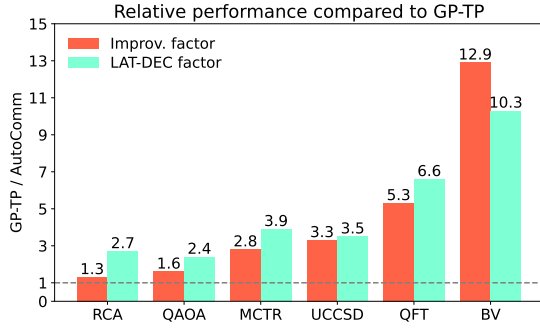


**Figure 16: Compared to GP-TP. Results are averaged over different configurations of # qubit and # node in Table 2.**

Comm. We denote this version of the GP-based compiler by *GP-TP*. Once again, for GP-TP, we adopt the as-soon-as-possible schedule strategy in [15].

As shown in Figure 16, AutoComm achieves significant reduction in both communication cost and program latency, compared to GP-TP. Specifically, AutoComm reduces the communication cost by a factor of 3.3x on average, up to 12.9x. It also reduces the program execution time by a factor of 4.3x on average, up to 10.3x. On the side of information theory, AutoComm improves the performance by enabling a higher throughput of information. Each remote communication in GP-TP carries less than two remote CX gates which is much smaller than AutoComm. On the algorithmic side, AutoComm reduces unnecessary qubit movement by taking advantage of burst communication. For example, for a potential burst communication between $q_1$ and node B, if there are some commutable remote CX gates between $q_1$ and node C lying in between and interrupting the communication block between $q_1$ and node B, the GP-TP method needs to move $q_1$ to node B first, then to node C and back to node B again. However, with burst communication, we only need to first move $q_1$ to node B, and then to node C.

## 5.4 Optimization Analysis

In this section, we further explore and analyze the effect of each optimization in AutoComm. For each analysis, we change only one component of AutoComm at a time, with other components fixed, to isolate the effect of each component/optimization.

*The effect of communication aggregation:* Table 3 demonstrates the benefit of communication aggregation compared to the baseline. Here we further demonstrate the necessity of considering gate commutation in the aggregation pass. Figure 17(a) shows the communication cost comparison between the aggregation without gate commutation and the aggregation used in AutoComm. For the programs in Figure 17(a), AutoComm reduces the communication cost by a factor of 5.5x on average, up to 6.7x, compared to the aggregation without gate commutation. Gate commutation is indispensable for discovering burst communications, not only because multi-qubit gates are often scattered in quantum circuits, but also due to the uncertainty of qubit mapping to quantum nodes (the uncertainty of whether a CX is remote or local).

*The effect of hybrid communication assignment:* We further demonstrate the importance of considering both Cat-Comm and TP-Comm for burst communication. Figure 17(b) shows the communication cost comparison between the communication assignment with Cat-Comm only and the hybrid assignment scheme in AutoComm. The Cat-Comm only method is extended from the specialized compiler [21] for distributed VQE. For the programs in Figure 17(b), AutoComm reduces the communication cost by a factor of 2.8x on average, up to 4.6x, compared to the Cat-Comm only method. The key enabler for the hybrid scheme in AutoComm is that Cat-Comm only applies to few communication patterns and for the cases that Cat-Comm cannot apply, TP-Comm would be more efficient.

*The effect of communication scheduling:* We then study the effect of the communication scheduling optimization in AutoComm. Figure 17(c) shows the latency comparison between AutoComm's scheduling, denoted by burst-greedy, and the greedy (as-soon-as-possible) scheduling for communication blocks. For the programs in Figure 17(c), the burst-greedy method reduces the program latency by a factor of 1.4x on average, up to 1.6x, compared to the general greedy schedule. The effectiveness of AutoComm for scheduling burst communication stems from its smart utilization of communication qubits, especially for TP-Comm blocks, as discussed in Section 4.4.

## 5.5 Sensitivity Analysis

The performance of AutoComm may be affected by some external factors, e.g., the number of input qubits and the number of computing nodes. In this section, we study how the performance of AutoComm varies with those factors. We focus on 'improv. factor' here, and the variation of 'LAT-DEC factor' would follow a similar trend.

*The effect of # qubit:* Figure 17(d) shows how the improv. factor of AutoComm changes with the number of qubits. As shown in the figure, the improv. factor converges when # qubit/# node is large. This may be due to the fact that the
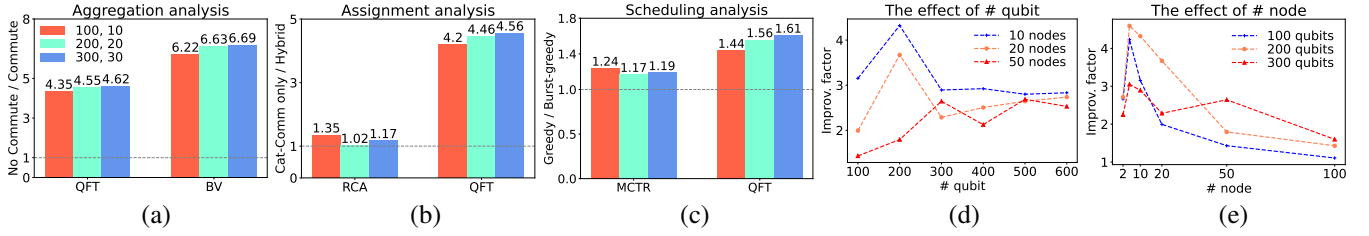
**Figure 17: (a)-(c) The effects of the proposed optimizations. Bars with different colors denote different configurations of (#qubit, #node). For (a)(b), the y-axis is the ratio of # remote communications. For (c), the y-axis is the ratio of program latency. (d)(e) The effects of # qubit and # node on the improv. factor of AutoComm. The test program in (d)(e) is MCTR.**

number of burst communication blocks also increases when the total number of remote multi-qubit gates grows with the number of qubits. Such behavior is preferable because it illustrates that AutoComm can provide a consistent reduction for the communication cost as the number of qubits grows.

*The effect of # node:* Figure 17(e) shows how the improv. factor of AutoComm changes with the number of nodes. In this figure, the performance of AutoComm deteriorates when # qubit/# node is small. This is because it is harder to find large communication blocks when the number of qubits in each node is limited to be small. Therefore we should avoid using too many nodes for distributing quantum programs because in such a case the remote multi-qubit gates would proliferate and there is little chance to execute those remote interactions collectively, given the fact that the number of communication qubits in each node is only two.

## 6  Discussion and Future Work

To the best of our knowledge, this paper is the first attempt that formalizes burst communication in distributed quantum programs. We discover a large number of burst communications hidden in various distributed quantum programs and propose the first modular framework to uncover these burst communications and use them to optimize the communication overhead. We argue that the formalization of burst communication and the modular solution proposed in this paper unveil new opportunities for communication optimization in DQC and would potentially inspire a series of works for overcoming DQC's communication problem.

Although we show that the proposed framework significantly surpasses existing works in optimizing the communication overhead of distributed quantum programs, there is still much space left for potential improvements.

*Extending to general collective communication* This paper only considers the near-term DQC where communication qubits are supposed to be limited. In such a case, we are restricted to studying the qubit-to-node burst communication, which is a special case of the general collective communication, involving a group of nodes. Assuming the availability of more communication qubits in the future, we could consider node-to-node collective communication which offers a potential optimization opportunity as we can now aggregate small qubit-to-node burst communication blocks into a larger one. Besides, for the fusion operation in the communication schedule optimization, we can also extend it to node-to-node

communication blocks.

*Co-designing with front compiling stages* The proposed framework is designed to be easily pluggable into existing compiling flows. But we could also couple it with front compiling stages to achieve further optimization. For example, existing compilers include a pass to add SWAP gates to change the qubit layout to optimize circuit metrics. We could co-design with this pass to maximize the number and size of burst communications. Besides, in the case where burst communication is deeply hidden, we could also consider using unitary synthesis to create burst communication in the gate decomposition pass. Finally, we could co-design with the qubit mapping pass to achieve a balance of communication overhead and device utilization rate, as shown in Figure 17(d)(e).

*Combining with quantum error correction* Since DQC involves quantum communication which is far more noisy than local quantum gates, reinforcing the whole distributed quantum system with quantum error correction (QEC) becomes vital for future DQC. One promising way to implement QEC in DQC is to encode one logical qubit in each node, and use quantum communication to implement logical operations between logical qubits. In this case, the CX gate between logical qubits would involve a large number of physical qubits simultaneously and provide great opportunities for burst communication optimization. Besides, communications coming from magic state distillation are also worth considering.

## 7  Related Work

Most existing quantum compilers [16, 17, 18, 19, 20] focus on the compilation of programs within a single quantum computer. Extending these works to DQC cannot achieve high information throughput per quantum communication, as in the compiler proposed by Ferrari et al. [15]. Baker et al. [11] propose using the more informative remote SWAP gates to replace all remote CX gates in distributed quantum programs and obtain a higher throughput. Diadamo et al. [21] further increase the communication throughput by considering multiple-qubit control-unitary blocks. However, their work requires specialized circuit representation and cannot optimize general quantum programs. Moreover, all these works do not consider the burst communication and related optimizations proposed in this paper.

Another line of work executes distributed quantum programs in a hybrid way. Tang et al. [40] propose a way to

execute quantum programs in distributed computing nodes but without inter-node communication. To overcome the expressibility loss due to no inter-node communication, their work relies heavily on classical post-processing techniques and cannot be extended to large-scale quantum programs.

Other quantum communication-related works focus on building robust quantum communication networks [41, 42, 43, 44] or reducing the resource consumption of existing quantum communication techniques [45, 46, 47, 48, 49]. These works are orthogonal to this paper.

## 8 Conclusion

As in classical distributed computing, the inter-node communication overhead bottlenecks distributed quantum computing. Existing compilers [11, 15, 21] for distributed programs either treat the inter-node communication like the in-node communication or only provide optimization for gates in the control-unitary form. These works fail to utilize the hidden communication patterns in distributed quantum programs. To overcome the shortcomings of existing compilers, this paper explores various distributed quantum programs and identifies burst communication for the first time. Burst communication is a qubit-node communication pattern that widely exists in many distributed programs. Based on burst communication, we propose the framework, AutoComm, which is proved to be efficient in cutting down inter-node communication overhead, by comprehensive evaluations on diverse distributed benchmarks. The proposed framework can be easily integrated into existing compiling flows of quantum programs and would benefit near-term distributed quantum computing.

## REFERENCES

[1] Peter W Shor. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. *SIAM review*, 41(2):303–332, 1999.

[2] Lov K Grover. A fast quantum mechanical algorithm for database search. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, pages 212–219, 1996.

[3] Michael A Nielsen and Isaac Chuang. Quantum computation and quantum information, 2002.

[4] Kenneth R Brown, Jungsang Kim, and Christopher Monroe. Co-designing a scalable quantum computer with trapped atomic ions. *npj Quantum Information*, 2(1):1–10, 2016.

[5] Colin D Bruzewicz, John Chiaverini, Robert McConnell, and Jeremy M Sage. Trapped-ion quantum computing: Progress and challenges. *Applied Physics Reviews*, 6(2):021314, 2019.

[6] Colin D Bruzewicz, John Chiaverini, Robert McConnell, and Jeremy M Sage. Trapped-ion quantum computing: Progress and challenges. *Applied Physics Reviews*, 6(2):021314, 2019.

[7] Markus Brink, Jerry M Chow, Jared Hertzberg, Easwar Magesan, and Sami Rosenblatt. Device challenges for near term superconducting quantum processors: frequency collisions. In *2018 IEEE International Electron Devices Meeting (IEDM)*, pages 6–1. IEEE, 2018.

[8] Matteo Pompili, Sophie LN Hermans, Simon Baier, Hans KC Beukers, Peter C Humphreys, Raymond N Schouten, Raymond FL Vermeulen, Marijn J Tiggelman, Laura dos Santos Martins, Bas Dirkse, et al. Realization of a multinode quantum network of remote solid-state qubits. *Science*, 372(6539):259–264, 2021.

[9] J. Gambetta. Ibm's roadmap for scaling quantum technology. https://www.ibm.com/blogs/research/2020/09/ibm-quantum-roadmap/.

[10] C. Monroe, R. Raussendorf, A. Ruthven, K. R. Brown, P. Maunz, L.-M. Duan, and J. Kim. Large scale modular quantum computer architecture with atomic memory and photonic interconnects. *Physical Review A*, 62(5):052317, 2000.

[11] Jonathan M. Baker, Casey Duckering, Alexander Hoover, and Frederic T. Chong. Time-sliced quantum circuit partitioning for modular architectures. *Proceedings of the 17th ACM International Conference on Computing Frontiers*, 2020.

[12] Christopher Young, Akbar Safari, Preston Huft, J. Zhang, Eun Oh, Ravikumar Chinnarasu, and Mark Saffman. An architecture for quantum networking of neutral atom processors. 2022.

[13] Anocha Yimsiriwattana and Samuel J Lomonaco Jr. Generalized ghz states and distributed quantum computing. *arXiv preprint quant-ph/0402148*, 2004.

[14] Jens Eisert, Kurt Jacobs, Polykarpos Papadopoulos, and Martin B Plenio. Optimal local implementation of nonlocal quantum gates. *Physical Review A*, 62(5):052317, 2000.

[15] Davide Ferrari, Angela Sara Cacciapuoti, Michele Amoretti, and Marcello Caleffi. Compiler design for distributed quantum computing. *IEEE Transactions on Quantum Engineering*, 2:1–20, 2021.

[16] Gushu Li, Yufei Ding, and Yuan Xie. Tackling the qubit mapping problem for nisq-era quantum devices. *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2019.

[17] MD SAJID ANIS, Héctor Abraham, AduOffei, Rochisha Agarwal, Gabriele Agliardi, Merav Aharoni, Ismail Yunus Akhalwaya, Gadi Aleksandrowicz, Thomas Alexander, Matthew Amy, Sashwat Anagolum, Eli Arbel, Abraham Asfaw, Anish Athalye, Artur Avkhadiev, Carlos Azaustre, PRATHAMESH BHOLE, Abhik Banerjee, Santanu Banerjee, Will Bang, Aman Bansal, Panagiotis Barkoutsos, Ashish Barnawal, George Barron, George S. Barron, Luciano Bello, Yael Ben-Haim, M. Chandler Bennett, Daniel Bevenius, Dhruv Bhatnagar, Arjun Bhobe, Paolo Bianchini, Lev S. Bishop, Carsten Blank, Sorin Bolos, Soham Bopardikar, Samuel Bosch, Sebastian Brandhofer, Brandon, Sergey Bravyi, Nick Bronn, Bryce-Fuller, David Bucher, Artemiy Burov, Fran Cabrera, Padraic Calpin, Lauren Capelluto, Jorge Carballo, Ginés Carrascal, Adam Carriker, Ivan Carvalho, Adrian Chen, Chun-Fu Chen, Edward Chen, Jielun (Chris) Chen, Richard Chen, Franck Chevallier, Kartik Chinda, Rathish Cholarajan, Jerry M. Chow, Spencer Churchill, CisterMoke, Christian Claus, Christian Clauss, Caleb Clothier, Romilly Cocking, Ryan Cocuzzo, Jordan Connor, Filipe Correa, Abigail J. Cross, Andrew W. Cross, Simon Cross, Juan Cruz-Benito, Chris Culver, Antonio D. Córcoles-Gonzales, Navaneeth D, Sean Dague, Tareq El Dandachi, Animesh N Dangwal, Jonathan Daniel, Marcus Daniels, Matthieu Dartiailh, Abdón Rodríguez Davila, Faisal Debouni, Anton Dekusar, Amol Deshmukh, Mohit Deshpande, Delton Ding, Jun Doi, Eli M. Dow, Eric Drechsler, Eugene Dumitrescu, Karel Dumon, Ivan Duran, Kareem EL-Safty, Eric Eastman, Grant Eberle, Amir Ebrahimi, Pieter Eendebak, Daniel Egger, ElePT, Emilio, Alberto Espiricueta, Mark Everitt, Davide Facoetti, Farida, Paco Martín Fernández, Samuele Ferracin, Davide Ferrari, Axel Hernández Ferrera, Romain Fouilland, Albert Frisch, Andreas Fuhrer, Bryce Fuller, MELVIN GEORGE, Julien Gacon, Borja Godoy Gago, Claudio Gambella, Jay M. Gambetta, Adhisha Gammanpila, Luis Garcia, Tanya Garg, Shelly Garion, James R. Garrison, Tim Gates, Leron Gil, Austin Gilliam, Aditya Giridharan, Juan Gomez-Mosquera, Gonzalo, Salvador de la Puente González, Jesse Gorzinski, Ian Gould, Donny Greenberg, Dmitry Grinko, Wen Guan, Dani Guijo, John A. Gunnels, Harshit Gupta, Naman Gupta, Jakob M. Günther, Mikael Haglund, Isabel Haide, Ikko Hamamura, Omar Costa Hamido, Frank Harkins, Kevin Hartman, Areeq Hasan, Vojtech Havlicek, Joe Hellmers, Łukasz Herok, Stefan Hillmich, Hiroshi Horii, Connor Howington, Shaohan Hu, Wei Hu, Junye Huang, Rolf Huisman, Haruki Imai, Takashi Imamichi, Kazuaki Ishizaki, Ishwor, Raban Iten, Toshinari Itoko, Alexander Ivrii, Ali Javadi, Ali Javadi-Abhari, Wahaj Javed, Qian Jianhua, Madhav Jivrajani, Kiran Johns, Scott Johnstun, Jonathan-Shoemaker, JosDenmark, JoshDumo, John Judge, Tal Kachmann, Akshay Kale, Naoki Kanazawa, Kang-Bae, Annanay Kapila, Anton Karazeev, Paul Kassebaum, Josh Kelso, Scott Kelso, Vismai Khanderao, Spencer King, Yuri Kobayashi, Kovi11Day, Arseny Kovyrshin, Rajiv Krishnakumar, Vivek Krishnan, Kevin Krsulich, Prasad Kumkar, Gawel Kus, Ryan LaRose, Enrique Lacal, Raphaël Lambert, Haggai Landa, John Lapeyre, Joe Latone, Scott Lawrence, Christina Lee, Gushu Li, Jake Lishman, Dennis Liu, Peng Liu, Abhishek K M, Liam Madden, Yunho Maeng, Saurav Maheshkar, Kahan Majmudar, Aleksei Malyshev, Mohamed El Mandouh, Joshua Manela, Manjula, Jakub Marecek, Manoel Marques, Kunal Marwaha, Dmitri Maslov, Paweł Maszota, Dolph Mathews, Atsushi Matsuo, Farai Mazhandu, Doug McClure, Maureen McElaney, Cameron McGarry, David McKay, Dan McPherson, Srujan Meesala, Dekel Meirom, Corey Mendell, Thomas Metcalfe, Martin Mevissen, Andrew Meyer, Antonio Mezzacapo, Rohit Midha, Daniel Miller, Zlatko

Minev, Abby Mitchell, Nikolaj Moll, Alejandro Montanez, Gabriel Monteiro, Michael Duane Mooring, Renier Morales, Niall Moran, David Morcuende, Seif Mostafa, Mario Motta, Romain Moyard, Prakash Murali, Jan Müggenburg, Tristan NEMOZ, David Nadlinger, Ken Nakanishi, Giacomo Nannicini, Paul Nation, Edwin Navarro, Yehuda Naveh, Scott Wyman Neagle, Patrick Neuweiler, Aziz Ngoueya, Johan Nicander, Nick-Singstock, Pradeep Niroula, Hassi Norlen, NuoWenLei, Lee James O'Riordan, Oluwatobi Ogunbayo, Pauline Ollitrault, Tamiya Onodera, Raul Otaolea, Steven Oud, Dan Padilha, Hanhee Paik, Soham Pal, Yuchen Pang, Ashish Panigrahi, Vincent R. Pascuzzi, Simone Perriello, Eric Peterson, Anna Phan, Kuba Pilch, Francesco Piro, Marco Pistoia, Christophe Piveteau, Julia Plewa, Pierre Pocreau, Alejandro Pozas-Kerstjens, Rafał Pracht, Milos Prokop, Viktor Prutyanov, Sumit Puri, Daniel Puzzuoli, Jesús Pérez, Quant02, Quintiii, Rafey Iqbal Rahman, Arun Raja, Roshan Rajeev, Isha Rajput, Nipun Ramagiri, Anirudh Rao, Rudy Raymond, Oliver Reardon-Smith, Rafael Martín-Cuevas Redondo, Max Reuter, Julia Rice, Matt Riedemann, Rietesh, Drew Risinger, Marcello La Rocca, Diego M. Rodríguez, RohithKarur, Ben Rosand, Max Rossmannek, Mingi Ryu, Tharrmashastha SAPV, Nahum Rosa Cruz Sa, Arijit Saha, Abdullah Ash-Saki, Sankalp Sanand, Martin Sandberg, Hirmay Sandesara, Ritvik Sapra, Hayk Sargsyan, Aniruddha Sarkar, Ninad Sathaye, Bruno Schmitt, Chris Schnabel, Zachary Schoenfeld, Travis L. Scholten, Eddie Schoute, Mark Schulterbrandt, Joachim Schwarm, James Seaward, Sergi, Ismael Faro Sertage, Kanav Setia, Freya Shah, Nathan Shammah, Rohan Sharma, Yunong Shi, Jonathan Shoemaker, Adenilton Silva, Andrea Simonetto, Deeksha Singh, Divyanshu Singh, Parmeet Singh, Phattharaporn Singkanipa, Yukio Siraichi, Siri, Jesús Sistos, Iskandar Sitdikov, Seyon Sivarajah, Magnus Berg Sletfjerding, John A. Smolin, Mathias Soeken, Igor Olegovich Sokolov, Igor Sokolov, Vicente P. Soloviev, SooluThomas, Starfish, Dominik Steenken, Matt Stypulkoski, Adrien Suau, Shaojun Sun, Kevin J. Sung, Makoto Suwama, Oskar Słowik, Hitomi Takahashi, Tanvesh Takawale, Ivano Tavernelli, Charles Taylor, Pete Taylour, Soolu Thomas, Kevin Tian, Mathieu Tillet, Maddy Tod, Miroslav Tomasik, Caroline Tornow, Enrique de la Torre, Juan Luis Sánchez Toural, Kenso Trabing, Matthew Treinish, Dimitar Trenev, TrishaPe, Felix Truger, Georgios Tsilimigkounakis, Davindra Tulsi, Wes Turner, Yotam Vaknin, Carmen Recio Valcarce, Francois Varchon, Adish Vartak, Almudena Carrera Vazquez, Prajjwal Vijaywargiya, Victor Villar, Bhargav Vishnu, Desiree Vogt-Lee, Christophe Vuillot, James Weaver, Johannes Weidenfeller, Rafal Wieczorek, Jonathan A. Wildstrom, Jessica Wilson, Erick Winston, WinterSoldier, Jack J. Woehr, Stefan Woerner, Ryan Woo, Christopher J. Wood, Ryan Wood, Steve Wood, James Wootton, Matt Wright, Lucy Xing, Jintao YU, Bo Yang, Unchun Yang, Daniyar Yeralin, Ryota Yonekura, David Yonge-Mallo, Ryuhei Yoshida, Richard Young, Jessie Yu, Lebin Yu, Christopher Zachow, Laura Zdanski, Helena Zhang, Iulia Zidaru, and Christa Zoufal. Qiskit: An open-source framework for quantum computing, 2021.

[18] Matthew Amy and Vlad Gheorghiu. staq—a full-stack quantum processing toolkit. *arXiv: Quantum Physics*, 2019.

[19] Nader Khammassi, Imran Ashraf, J. van Someren, Răzvan Nane, A. M. Krol, M. A. Rol, Lingling Lao, Koen Bertels, and Carmen Garcia Almudever. Openql : A portable quantum programming framework for quantum accelerators. *ACM J. Emerg. Technol. Comput. Syst.*, 18:13:1–13:24, 2022.

[20] 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*, 2020.

[21] Stephen Diadamo, Marco Ghibaudi, and James R. Cruise. Distributed quantum computing and network control for accelerated vqe. *IEEE Transactions on Quantum Engineering*, 2:1–21, 2021.

[22] Nemanja Isailovic, Yatish Patel, Mark Whitney, and John Kubiatowicz. Interconnection networks for scalable quantum computers. In *33rd International Symposium on Computer Architecture (ISCA'06)*, pages 366–377. IEEE, 2006.

[23] David L Moehring, Peter Maunz, Steve Olmschenk, Kelly C Younge, Dzmitry N Matsukevich, L-M Duan, and Christopher Monroe. Entanglement of single-atom quantum bits at a distance. *Nature*, 449(7158):68–71, 2007.

[24] Stephan Ritter, Christian Nölleke, Carolin Hahn, Andreas Reiserer, Andreas Neuzner, Manuel Uphoff, Martin Mücke, Eden Figueroa, Joerg Bochmann, and Gerhard Rempe. An elementary quantum network of single atoms in optical cavities. *Nature*, 484(7393):195–200, 2012.

[25] Julian Hofmann, Michael Krug, Norbert Ortegel, Lea Gérard, Markus Weber, Wenjamin Rosenfeld, and Harald Weinfurter. Heralded entanglement between widely separated atoms. *Science*, 337(6090):72–75, 2012.

[26] LJ Stephenson, DP Nadlinger, BC Nichol, S An, P Drmota, TG Ballance, K Thirumalai, JF Goodwin, DM Lucas, and CJ Ballance. High-rate, high-fidelity entanglement of qubits across an elementary quantum network. *Physical review letters*, 124(11):110501, 2020.

[27] Hannes Bernien, Bas Hensen, Wolfgang Pfaff, Gerwin Koolstra, Machiel S Blok, Lucio Robledo, Tim H Taminiau, Matthew Markham, Daniel J Twitchen, Lilian Childress, et al. Heralded entanglement between solid-state qubits separated by three metres. *Nature*, 497(7447):86–90, 2013.

[28] Peter C Humphreys, Norbert Kalb, Jaco PJ Morits, Raymond N Schouten, Raymond FL Vermeulen, Daniel J Twitchen, Matthew Markham, and Ronald Hanson. Deterministic delivery of remote entanglement on a quantum network. *Nature*, 558(7709):268–273, 2018.

[29] Aymeric Delteil, Zhe Sun, Wei-bo Gao, Emre Togan, Stefan Faelt, and Ataç Imamoğlu. Generation of heralded entanglement between distant hole spins. *Nature Physics*, 12(3):218–223, 2016.

[30] Robert Stockill, MJ Stanley, Lukas Huthmacher, E Clarke, M Hugues, AJ Miller, C Matthiesen, Claire Le Gall, and Mete Atatüre. Phase-tuned entangled state generation between distant spin qubits. *Physical review letters*, 119(1):010503, 2017.

[31] P Maunz, S Olmschenk, D Hayes, DN Matsukevich, L-M Duan, and C Monroe. Heralded quantum gate between remote quantum memories. *Physical review letters*, 102(25):250502, 2009.

[32] Severin Daiss, Stefan Langenfeld, Stephan Welte, Emanuele Distante, Philip Thomas, Lukas Hartung, Olivier Morin, and Gerhard Rempe. A quantum-logic gate between distant quantum-network modules. *Science*, 371(6529):614–617, 2021.

[33] Norbert Kalb, Andreas A Reiserer, Peter C Humphreys, Jacob JW Bakermans, Sten J Kamerling, Naomi H Nickerson, Simon C Benjamin, Daniel J Twitchen, Matthew Markham, and Ronald Hanson. Entanglement distillation between solid-state quantum network nodes. *Science*, 356(6341):928–932, 2017.

[34] Daniele Cuomo, Marcello Caleffi, Kevin Krsulich, Filippo Tramonto, Gabriele Agliardi, Enrico Prati, and Angela Sara Cacciapuoti. Optimized compiler for distributed quantum computing. *ArXiv*, abs/2112.14139, 2021.

[35] Maarten Van Steen and A Tanenbaum. Distributed systems principles and paradigms. *Network*, 2:28, 2002.

[36] R. Wille, D. Große, L. Teuber, G. W. Dueck, and R. Drechsler. RevLib: An online resource for reversible functions and reversible circuits. In *Int'l Symp. on Multi-Valued Logic*, pages 220–225, 2008. RevLib is available at http://www.revlib.org.

[37] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. A quantum approximate optimization algorithm. *arXiv: Quantum Physics*, 2014.

[38] Yun Seong Nam, Neil J. Ross, Yuan Su, Andrew M. Childs, and Dmitrii L. Maslov. Automated optimization of large quantum circuits with continuous parameters. *npj Quantum Information*, 4:1–12, 2017.

[39] Roberto Sanchez Correa and Jean Pierre David. Ultra-low latency communication channels for fpga-based hpc cluster. *Integration*, 63:41–55, 2018.

[40] Wei Tang, Teague Tomesh, Martin Suchara, Jeffrey Larson, and Margaret Martonosi. Cutqc: using small quantum computers for large quantum circuit evaluations. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 473–486, 2021.

[41] Mihir Pant, Hari Krovi, Don Towsley, Leandros Tassiulas, Liang Jiang, Prithwish Basu, Dirk Englund, and Saikat Guha. Routing entanglement in the quantum internet. *npj Quantum Information*, 5(1):1–9, 2019.

[42] Kaushik Chakraborty, Filip Rozpedek, Axel Dahlberg, and Stephanie Wehner. Distributed routing in a quantum internet. *arXiv preprint arXiv:1907.11630*, 2019.

[43] Frederik Hahn, A Pappa, and Jens Eisert. Quantum network routing and local complementation. *npj Quantum Information*, 5(1):1–7, 2019.

[44] Changhao Li, Tianyi Li, Yi-Xiang Liu, and Paola Cappellaro. Effective routing design for remote entanglement generation on quantum networks. *npj Quantum Information*, 7(1):1–12, 2021.

[45] H-J Briegel, Wolfgang Dür, Juan I Cirac, and Peter Zoller. Quantum repeaters: the role of imperfect local operations in quantum communication. *Physical Review Letters*, 81(26):5932, 1998.

[46] L-M Duan, Mikhail D Lukin, J Ignacio Cirac, and Peter Zoller. Long-distance quantum communication with atomic ensembles and linear optics. *Nature*, 414(6862):413–418, 2001.

[47] Lilian Childress, JM Taylor, Anders Søndberg Sørensen, and Mikhail D Lukin. Fault-tolerant quantum repeaters with minimal physical resources and implementations based on single-photon emitters. *Physical Review A*, 72(5):052330, 2005.

[48] Nicolas Sangouard, Christoph Simon, Hugues De Riedmatten, and Nicolas Gisin. Quantum repeaters based on atomic ensembles and linear optics. *Reviews of Modern Physics*, 83(1):33, 2011.

[49] Sreraman Muralidharan, Linshu Li, Jungsang Kim, Norbert Lütkenhaus, Mikhail D Lukin, and Liang Jiang. Optimal architectures for long distance quantum communication. *Scientific reports*, 6(1):1–10, 2016.