# Entanglement-assisted circuit knitting

Shao-Hua Hu,[1, *] Po-Sung Liu,[2, †] and Jun-Yi Wu[3, 4, 5, ‡]

[1]*Department of Physics, National Tsing Hua University, Hsinchu 30013, Taiwan, ROC*

[2]*Department of Physics, National Cheng Kung University, Tainan 701, Taiwan, ROC*

[3]*Department of Physics, Tamkang University, New Taipei 25137, Taiwan, ROC*

[4]*Hon Hai Research Institute, Taipei, Taiwan, ROC*

[5]*Physics Division, National Center for Theoretical Sciences, Taipei, Taiwan, ROC*

## Abstract

Distributed quantum computing (DQC) provides a promising route toward scalable quantum computation, where entanglement-assisted LOCC and circuit knitting represent two complementary approaches. The former deterministically realizes nonlocal operations but demands extensive entanglement resources, whereas the latter requires no entanglement yet suffers from exponential sampling overhead. Here, we propose a hybrid framework that integrates these two paradigms by performing circuit knitting assisted with a limited amount of entanglement. We establish a general theoretical formulation that yields lower bounds on the optimal sampling overhead and present a constructive protocol demonstrating that a single shared Bell pair can reduce the overhead to the asymptotic limit of standard circuit knitting without requiring classical communication. This hybrid approach enhances both sampling and entanglement efficiency, enabling more resource-practical implementations of distributed quantum computation.

---

[*] shhphy@gmail.com

[†] sung920405@gmail.com

[‡] junyiwuphysics@gmail.com

**CONTENTS**

## I. INTRODUCTION

Quantum computing offers immense potential to outperform classical computation, making the realization of scalable, large-scale quantum devices a central challenge. Distributed quantum computing (DQC) [1–3] addresses this challenge by interconnecting multiple local quantum processing units (QPUs) to collectively realize a global computation. Equivalently, a global unitary operation can be divided into smaller segments executed on different QPUs. DQC implementations are generally classified into two approaches: entanglement-assisted local operation and classical communication (LOCC) [4–8] and circuit knitting [9–13].

Although entanglement-assisted LOCC approaches for DQC can deterministically implement a global unitary, they are highly resource-intensive in terms of the number of entangled pairs required. Such fully entanglement-assisted DQC schemes can be categorized into two types, namely quantum

state teleportation [14] and quantum telegate [5], which serve as fundamental building blocks for constructing distributed quantum processes [1–3, 6–8].

Similar to entanglement-assisted DQC, circuit knitting was proposed as a method to simulate a large quantum circuit using smaller subcircuits [9]. In contrast to entanglement-assisted approaches, circuit knitting requires no entanglement resources. Instead, it relies on a classical postprocessing technique known as quasi-probability simulation [10, 12, 15], which estimates the statistical outcomes of a quantum circuit rather than physically implementing the global operation. In this framework, the global operation is reconstructed by probabilistically sampling local operations, and the original statistics are recovered by assigning positive or negative weights to measurement outcomes in classical postprocessing. As a result, maintaining the same estimation accuracy incurs a sampling overhead, which increases the number of required circuit executions (measurement shots) and, consequently, the total runtime. Therefore, the central objective in the study of circuit knitting is to identify an optimal set of local operations that minimizes this sampling overhead and thus reduces the overall time cost.

These two approaches reveal a fundamental trade-off between entanglement consumption and execution time, representing the two extreme regimes of DQC. The entanglement-assisted approach requires a large amount of high-fidelity entanglement distributed across the quantum network of QPUs, but enables fast runtime. In contrast, circuit knitting requires no entanglement resources, yet incurs a significantly longer runtime due to its sampling overhead. In addition to entanglement resources, the role of classical communication in these approaches also deserves consideration. It has been shown that, for certain classes of unitary operations, the optimal sampling overhead can be achieved even without classical communication [13, 16–19]. Moreover, the virtual simulation of global gates offers an additional advantage of mitigating noise in quantum circuit implementations—another benefit of circuit knitting. However, in general, the sampling overhead increases exponentially with the number of nonlocal operations [20].

In this work, we aim to incorporate these two approaches into a hybrid solution for distributed quantum computing (DQC). In this scenario, local QPUs have access to a limited number of entangled pairs—insufficient to realize a fully entanglement-assisted DQC scheme. Consequently, a sampling overhead in time must be paid to simulate the target quantum operation via quasi-probability decomposition (QPD) sampling assisted by a partial entanglement resource. This hybrid approach leverages the advantages of both entanglement-assisted LOCC and circuit knitting. As illustrated in Fig. 1, introducing a finite amount of entanglement can significantly reduce the sampling overhead to a practical level, thereby preserving quantum advantage. Our goal is to

3

investigate this trade-off and identify the balance between entanglement consumption and sampling overhead in DQC.
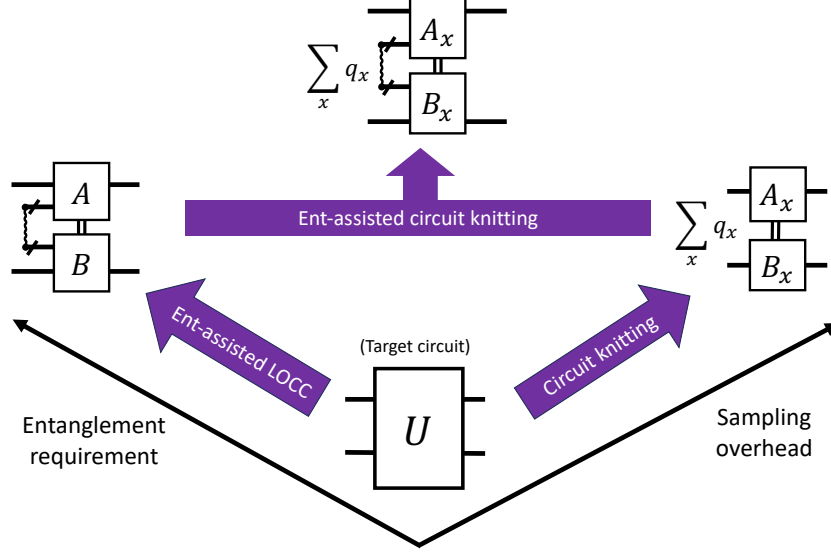


FIG. 1. This figure illustrates the main idea of the entangled-assisted circuit knitting, in which the trade-off of pre-shared entangled state and the sampling overhead is expected.

## II. RESOURCE-ASSISTED QUANTUM OPERATION

It is convenient to employ the vectorization formulation [21] in the Liouville space to describe density operators and their evolution under quantum channels. Here, we adopt the double ket notation to represent the vectorization of an operator. The vectorization of the identity operator is given by $|I_d\rangle\rangle := \sum_{i=0}^{d-1} |i,i\rangle$. Accordingly, we denote the vectorization of a general operator $\hat{O}$ by $|O\rangle\rangle := (\hat{O} \otimes \hat{I}_d) |I_d\rangle\rangle$. The density operator after the vectorization is then given by

$$|\rho\rangle\rangle = (\hat{\rho} \otimes \hat{I}_d) |I_d\rangle\rangle . \tag{1}$$

The unitary transformation of a vectorized state is denoted by the tilde mark and written as

$$\widetilde{U} := \hat{U} \otimes \hat{U}^*. \tag{2}$$

A general CPTP map $\widetilde{Q}$, that is decomposed as a sum of Kraus operators $\{\hat{K}_i\}_i$, can be then described by the sum of the vectorization of Kraus operators

$$\widetilde{Q} = \sum_i \widetilde{K}_i. \tag{3}$$

4

It is called the operator sum representation of $\widetilde{Q}$.

With an ancillary subspace, one can modulate the quantum operation with an ancillary input $|\rho_{anc}\rangle\rangle$ and a POVM measurement $\{\langle\langle M_m^{(anc)}|\}_m$ to construct an instrument $\{\widetilde{K}'_m\}_m$

$$\widetilde{K}'_m = \langle\langle M_m^{(\text{anc.})}| \, \widetilde{Q} \, |\rho_{\text{anc.}}\rangle\rangle \tag{4}$$

Such a construction of state-assisted quantum instruments has been employed in entanglement-assisted distributed quantum computing [5, 7, 8], in which maximally entangled states are employed as the ancillary.

We can describe the state preparation of $|\rho\rangle\rangle$ as a quantum operation mapping classical information to the ancillary Hilbert space,

$$\widetilde{r}_\rho^{(\text{anc.})} := |\rho_{\text{anc.}}\rangle\rangle . \tag{5}$$

In general, one can implement a pre-operation $\widetilde{P}$ before the state preparation. As a whole, one can construct a quantum instrument through

$$\widetilde{K}'_m = \langle\langle M_m^{(\text{anc.})}| \, \widetilde{Q} \circ \widetilde{r}_\rho^{(\text{anc.})} \circ \widetilde{P}. \tag{6}$$

In a more general framework, the operation $\widetilde{r}$ is not restricted to the initial preparation of ancillary states. Instead, it may represent any quantum operation that can be implemented with the assistance of a suitable quantum-state preparation. For instance, a quantum operation $\widetilde{R} = \sum_m \widetilde{K}'_m$, constructed from the set of operators $\widetilde{K}'_m$, can be recursively employed to generate a higher-level quantum instrument by replacing the ancillary-state preparation $\widetilde{r}_\rho^{(\text{anc.})}$ with $\widetilde{R}$.



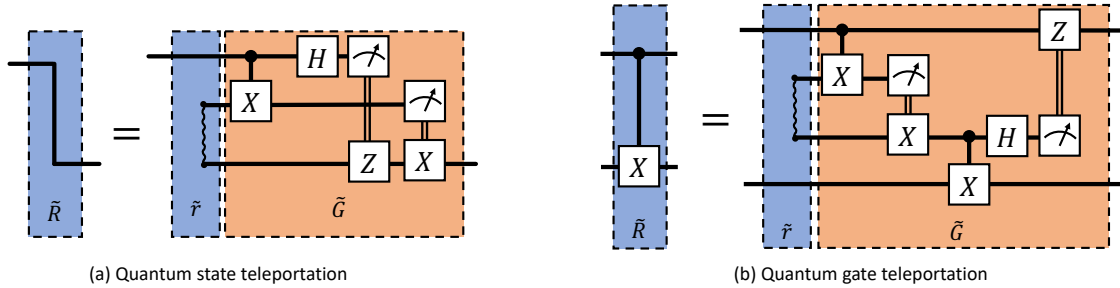(a) Quantum state teleportation      (b) Quantum gate teleportation

FIG. 2. Diagram of quantum state and gate teleportation. The blue one indicates the resource in the process, where the orange one indicates the free operation.

## A. Quasi-probability decomposition over free operations

Following the framework of quantum resource theories [22], let $\mathbb{F}$ be a set of free operations, which fulfills two conditions. First, it contains the identity map, $\widetilde{\mathbb{1}} \in \mathbb{F}$. Second, it is closed under

composition, i.e. $\widetilde{E}_a \circ \widetilde{E}_b \in \mathbb{F}$ if $\widetilde{E}_a$ and $\widetilde{E}_b$ are both free. Third, we consider the convex resource theories, that is, the set of free operations is convex. So taking a convex combination of the free operations remains free, i.e., $\forall \tilde{E}_a, \tilde{E}_b \in \mathbb{F}, p \in [0,1] \Rightarrow p\tilde{E}_a + (1-p)\tilde{E}_b \in \mathbb{F}$.

Let $\mathbb{X}$ be a random variable of bitstrings $x$ representing the classical information that is used to label the available free operations $\mathbb{F}_x$ One randomly implement the free maps $\mathbb{F}$ over $\mathbb{X}$ with a probability of $\{p_x\}_x$. Based on the classical information $x$, one randomly implements a free operation $\widetilde{F}_x$ with a probability of $p_x$. Following each free operation, a POVM measurements $\mathcal{M}_x = \{\langle\langle M_{m|x}^{(\text{anc.})}|\}_m$ conditional on $x$ is implemented on an ancillary subspace. In quasi-probability decomposition, one assigns a binary sign function $s_x(m) = \pm 1$ to each POVM operator $\langle\langle M_{m_x}|$ in the postprocessing of the measurements. As a whole, a quasi-probability decomposition over the free maps $\mathbb{F}$, can be formulated as follows.

**Definition 1 (Quasi-probability decomposition over free maps)** *Let $\mathbb{X}$ be a random variable of classical bitstrings $x$ associated with a probability distribution $\{p_x\}_x$. For each bitstring $x$, there are a free operation $\widetilde{F}_x \in \mathbb{F}$ and a free POVM $\mathcal{M}_x = \{\langle\langle M_{x,m}^{(\text{anc.})}|\}_m$ available. One can then construct a QPD of a quantum operation $\widetilde{R}$ by sampling $\mathbb{F}$ over $\mathbb{X}$ with the assignment $s_x$ to each POVM $\mathcal{M}_x$,*

$$\widetilde{R} = \gamma \sum_{x \in \mathbb{X}} p_x \sum_m s_x(m) \langle\langle M_{x,m}^{(\text{anc.})}| \, \widetilde{F}_x, \tag{7}$$

*where $\gamma$ is the normalization factor that normalizes the operation $\widetilde{R}$ to a CPTP map. The tuples $\mathcal{Q} = \{(p_x, \widetilde{F}_x; s_x, \mathcal{M}_x)\}_{x \in \mathbb{X}}$ is a QPD configuration of $\widetilde{R}$.*

The normalization factor for a configuration $\mathcal{Q}$ of $\widetilde{R}$ is denoted by $\gamma_\mathcal{Q}$. Since $s_m$ can be negative, without the normalization factor $\gamma_\mathcal{Q}$, the quantum operation is in general a CPTN (complete-positive trace-non-increasing). To ensure that $\widetilde{R}$ is a CPTP, it must fulfill $\frac{1}{d_{in}} \langle\langle \mathbb{1}_{out} | \widetilde{R} | \mathbb{1}_{in}\rangle\rangle = 1$, which determines the normalization factor as

$$\frac{1}{\gamma_\mathcal{Q}} = \frac{1}{d_{in}} \sum_{x \in \mathbb{X}} p_x \sum_m s_x(m) \langle\langle \mathbb{1}_{out} \otimes M_{x,m}^{(\text{anc.})} | \widetilde{F}_x | \mathbb{1}_{in}\rangle\rangle \leq 1 \tag{8}$$

In general, $\gamma_\mathcal{Q} \geq 1$, it has its minimum $\gamma_\mathcal{Q} = 1$, when $s_x(m) = 1 \ \forall x, m$.

Note that for arbitrary $\mathbb{F}$, a QPD configuration of $\tilde{R}$ may not exist. However, whenever such a decomposition is found, it can be used in the task of estimating expectation value with the form $\langle\langle O | \tilde{R} | \rho \rangle\rangle$, for any initial state $|\rho\rangle\rangle$ and observable $\langle\langle O|$. Given a QPD of $\tilde{R}$, we may rewrite the expectation value as

$$\langle\langle O | \tilde{R} | \rho \rangle\rangle = \gamma \sum_{x \in \mathbb{X}} p_x \sum_m s_x(m) \langle\langle M_{m|x}^{(\text{anc.})} | \, \widetilde{F}_x | \rho \rangle\rangle. \tag{9}$$

Such a QPD construction allows the evaluation of the $\widehat{O}$ through a Monte Carlo sampling simulation:

1. Implement the free map $\widetilde{F}_x$ on the input $|\rho\rangle\rangle$ with the probability $p_x$.

2. Implement the measurement $\mathcal{M}_x^{(\text{anc.})} = \{\langle\langle M_{m|x}^{(\text{anc.})}|\}_m$ on the ancillary qubits.

3. Assign the outcome with the sign $s_x(m)$.

However, one has to pay a price for using free operations to simulate a resource operation due to the additional normalization factor $\gamma_\mathcal{Q}$. Since it amplifies the statistical uncertainty, to compensate for the amplification of uncertainty, one needs to measure more samples in the measurements of $\widehat{O}$. More precisely, by Hoeffding's inequality, to estimate the outcome with the same accuracy, the total sampling number increases by a factor of $\gamma_\mathcal{Q}^2$ [9, 12, 15]. The normalization factor $\gamma_\mathcal{Q}$ is also called the *sampling overhead* of the QPD configuration $\mathcal{Q}$. Given a set of free maps $\mathbb{F}$, we can define the sampling overhead of the QPD of a resource operation $\widetilde{R}$ over $\mathbb{F}$ as the minimum overhead of QPD configurations constructed by free maps in $\mathbb{F}$

$$\gamma_\mathbb{F}(\tilde{R}) := \inf\left\{\gamma_\mathcal{Q} : \mathcal{Q} \text{ is a QPD configuration for } \widetilde{R} \text{ with } \{\widetilde{F}_x\}_x \subseteq \mathbb{F}.\right\} \tag{10}$$

### B. Resource-free circuit knitting

The conventional circuit knitting technique [9, 13, 17–19] can be viewed as an instance of quasi-probability decomposition (QPD) of global unitaries over LOCC. It provides a practical framework for realizing distributed modular quantum computing [8], where large-scale quantum circuits are executed using multiple small-scale quantum processing units (QPUs). The QPD of global unitaries over LOCC across local modular QPUs thus offers a natural strategy for resource-free distributed quantum computation (DQC).

As illustrated in Fig. 3, circuit knitting typically involves two types of circuit partitioning: wire cutting in the time domain and gate cutting in the spatial domain. These two forms of circuit cutting can be interpreted as resource-free QPD simulation of state teleportation and gate teleportation, respectively, both serving as fundamental building blocks of entanglement-assisted DQC.

Consider a bipartite system that consists of two QPUs. The target resource quantum operations of circuit knitting are global bipartite unitaries. In general, a global bipartite unitary can be written in the local unitary decomposition (LUD) with the following definition.
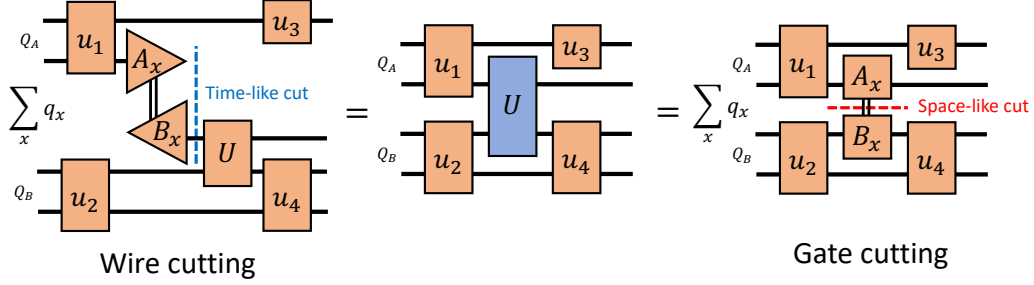
FIG. 3. This figure illustrates two types of circuit knitting, namely wire cutting and gate cutting.

**Definition 2 (Local unitary decomposition)** *Let $\widehat{U}$ be a unitary act on two subsystems $A|B$. A local unitary decomposition (LUD) of $\widehat{U}$ is given by*

$$\widehat{U} = \sum_i \lambda_i \widehat{A}_i \otimes \widehat{B}_i, \tag{11}$$

*where $\hat{A}_i$ and $\hat{B}_i$ are all unitary and the coefficient $s_i$ is a positive real number. We call a unitary operator is KAK-like, if there exists a LUD such that two sets of local unitaries $\{\hat{A}_i\}_i$ and $\{\hat{B}_i\}_i$ are both orthogonal [23].*

It is worth noting that all two-qubit unitary operators are KAK-like. The tensor product of two KAK-like unitaries is also a KAK-like unitary. It has been shown that the QPD overhead of a KAK-like unitary over LOCC is equivalent to the one over LO [18, 19].

**Proposition 3** *[18, 19] Let $\hat{U}$ be a bipartite unitary with the LUD $\hat{U} = \sum_i \lambda_i \hat{A}_i \otimes \hat{B}_i$ with a coefficient $\boldsymbol{\lambda} = (\lambda_1, ..., \lambda_K)$. It has a QPD overhead over LO $\gamma_{\mathrm{LO}}(\tilde{U}) \leq 2||\boldsymbol{\lambda}||_1^2 - ||\boldsymbol{\lambda}||_2^2$. Moreover, when $\hat{U}$ is KAK-like unitary, we have*

$$\gamma_{\mathrm{LOCC}}(\tilde{U}) = \gamma_{\mathrm{LO}}(\tilde{U}) = 2||\boldsymbol{\lambda}||_1^2 - 1 \tag{12}$$

The overhead of a resource-free QPD is proportional to the norm of the Schmidt coefficient $\boldsymbol{s}$, which increases exponentially with the entangling power of a unitary.

### C.  Resources-assisted quasi-probability decomposition

To enhance the practical feasibility of circuit knitting, we incorporate entangled-state preparation into the QPD framework, thereby formulating an entanglement-assisted QPD that reduces sampling overhead through the utility of entanglement resources. Conversely, from an equivalent

perspective, one may incorporate QPD into entanglement-assisted DQC, thereby reducing the required entanglement resources at the cost of additional sampling overhead. We refer to this unified framework as resource-assisted QPD, which is implemented by the following sampling process:

1. Initialize a set of random classical bitstrings $\mathbb{X}$ with the probability distribution $\{p_x\}_{x \in \mathbb{X}}$.

2. With a probability of $p_x$, one implements a free operation $\widetilde{F}_x$ before the utility of the assisting resource operation $\widetilde{r}$, followed by a free operation $\widetilde{G}_x$. The free operations $\widetilde{F}_x$ and $\widetilde{G}_x$ are the pre- and post-operation of the assisting resource $\widetilde{r}$, respectively, which build up a $\widetilde{r}$-assisted quantum operation $\widetilde{F}_x \circ \widetilde{r} \circ \widetilde{G}_x$ labeled by $x$.

3. Implement the measurement $\mathcal{M}_x^{(\text{anc.})} = \{\langle\langle M_{m|x}^{(\text{anc.})}|\}_m$ on the ancillary qubits to construct a quantum instrument.

4. Assign each outcome $m$ with the sign $s_x(m)$.

The formal definition of a resource-assisted QPD is formulated as follows.

**Definition 4 (Resource-assisted quasi-probability decomposition over free maps)** *Let $\mathbb{X}$ be a random variable of classical bitstrings $x$ associated with a probability distribution $\{p_x\}_x$. For each bitstring $x$, there are free operations $\widetilde{F}_x$ and $\widetilde{G}_x$, and a free POVM $\mathcal{M}_x = \{\langle\langle M_{x,m}^{(\text{anc.})}|\}_m$ available. In addition, one has a resource operation $\widetilde{r}$ available. One can then construct a $\widetilde{r}$-assisted QPD of a quantum operation $\widetilde{R}$ by sampling $\widetilde{G}_x \circ \widetilde{r} \circ \widetilde{F}_x$ over $\mathbb{X}$ with the assignment $s_x$ to each POVM $\mathcal{M}_x$,*

$$\widetilde{R} = \gamma_\mathcal{Q} \sum_{x \in \mathbb{X}} p_x \sum_m s_x(m) \langle\langle M_{m|x}^{(\text{anc.})}| \widetilde{G}_x \circ \widetilde{r} \circ \widetilde{F}_x. \tag{13}$$

*where $\gamma_\mathcal{Q}$ is the normalization factor that normalizes the operation $\widetilde{R}$ to a CPTP map. The tuples $\mathcal{Q}(\widetilde{r} \to \widetilde{R}) = \{(p_x, \widetilde{F}_x, \widetilde{G}_x; s_x, \mathcal{M}_x)\}_{x \in \mathbb{X}}$ is a $\widetilde{r}$-assisted QPD configuration for the quantum operation $\widetilde{R}$.*

The resource-assisted QPD can be used to describe both resource-free circuit knitting and fully entanglement-assisted DQC, which are two extremum cases with resource-free $\widetilde{r} = |\mathbb{1}_{\text{anc.}}\rangle\rangle$ and QPD-free $s_x = 1$, respectively. For example, the quantum state teleportation and quantum telegate protocols shown in Fig. 4 are two fully entanglement-assisted DQC protocols implemented with entanglement-assisted LOCC. It is a special case of Bell-state-assisted QPD with the configuration $\mathcal{Q} = \left(p_x = 1, \widetilde{F}_x = \widetilde{\mathbb{1}}, \widetilde{G}_x; s_x = 1, \mathcal{M}_x = \langle\langle\mathbb{1}|\right)$, where the resource operations are highlighted in blue and $\widetilde{G}_x$ is the LOCC highlighted in orange.

(a) Quantum state teleportation   (b) Quantum gate teleportation
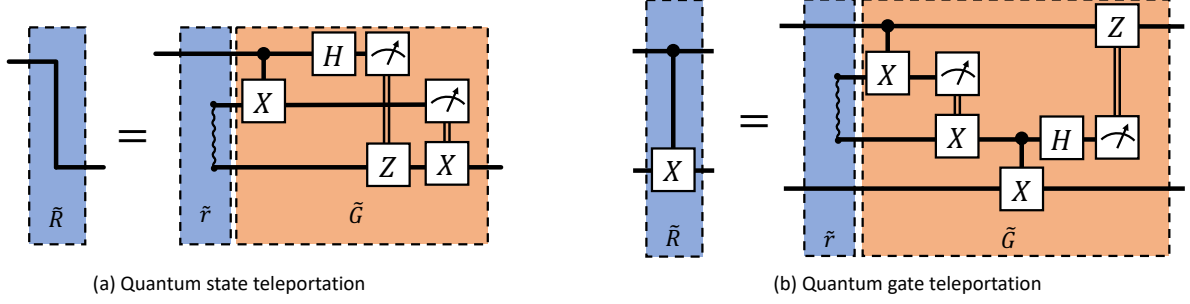
FIG. 4. Fully entanglement-assisted DQC. (a) Quantum state teleportation. (b) Quantum telegate implemented by entanglement-assisted LOCC [5].

**Definition 5** *The normalization factor $\gamma$ in Eq. (13) is called the $\widetilde{r}$-assisted $\mathcal{Q}$-sampling overhead for $\widetilde{R}$, which is denoted and determined by*

$$\gamma_{\mathcal{Q}}(\widetilde{r} \to \widetilde{R}) = \left( \frac{1}{d_{in}} \sum_{x \in \mathbb{X}} p_x \sum_m s_x(m) \langle\langle \mathbb{1}_{out} \otimes M_{m|x}^{(anc.)} | \widetilde{G}_x \circ \widetilde{r} \circ \widetilde{F}_x | \mathbb{1}_{in} \rangle\rangle \right)^{-1}. \tag{14}$$

*Given a set of free maps $\mathbb{F}$, the $\widetilde{r}$-assisted $\mathbb{F}$-sampling overhead for a quantum operation $\widetilde{R}$ is defined as the minimum sampling overhead of the $\widetilde{r}$-assisted QPD configurations constructed from the free maps in $\mathbb{F}$,*

$$\gamma_{\mathbb{F}}(\widetilde{r} \to \widetilde{R}) := \inf \left\{ \gamma_{\mathcal{Q}} : \mathcal{Q}(\widetilde{r} \to \widetilde{R}) \text{ is a } \widetilde{r}\text{-assisted QPD for } \widetilde{R} \text{ with } \widetilde{F}_x \in \mathbb{F} \text{ and } \widetilde{G}_x \in \mathbb{F} \right\}. \tag{15}$$

The arrow "$\to$" indicates the conversion of the initial resource operation $\widetilde{r}$ to the target resource operation $\widetilde{R}$. For consistence, we omit the arrow in the QPD overhead in Eq. (10),

$$\gamma_{\mathbb{F}}(\tilde{R}) := \gamma_{\mathbb{F}}(\tilde{r} \to \tilde{R}), \tag{16}$$

if the initial resource is free, $\tilde{r} \in \mathbb{F}$. The overheads of resource-assisted $\mathbb{F}$-QPDs have the following properties.

**Lemma 6** *Let $\widetilde{A}, \widetilde{B}, \widetilde{C}$ be three quantum operations. The resource-assisted QPDs over a set of free maps $\mathbb{F}$ fulfills the following properties:*

**Sub-multiplicity:** $\gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{C}) \le \gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{B}) \cdot \gamma_{\mathbb{F}}(\widetilde{B} \to \widetilde{C})$.

**Ordering:** $\gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{B}) = 1 \Rightarrow \gamma_{\mathbb{F}}(\widetilde{A}) \ge \gamma_{\mathbb{F}}(\widetilde{C})$.

**Right-convexity:** $\gamma_{\mathbb{F}}(\widetilde{A} \to p\widetilde{B} + (1-p)\widetilde{C}) \le p\gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{B}) + (1-p)\gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{C}) \ \forall p \in [0,1]$.

**Left-convexity:** $p\gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{C}) + (1-p)\gamma_{\mathbb{F}}(\widetilde{B} \to \widetilde{C}) \le \gamma_{\mathbb{F}}(p\widetilde{A} + (1-p)\widetilde{B} \to \widetilde{C}) \ \forall p \in [0,1]$.

**Right-monotonicity:** $\widetilde{f} \in \mathbb{F} \Rightarrow \max \left\{ \gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{f} \circ \widetilde{B}), \gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{B} \circ \widetilde{f}) \right\} \leq \gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{B}).$

**Left-monotonicity:** $\widetilde{f} \in \mathbb{F} \Rightarrow \gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{B}) \leq \min \left\{ \gamma_{\mathbb{F}}(\widetilde{f} \circ \widetilde{A} \to \widetilde{B}), \gamma_{\mathbb{F}}(\widetilde{A} \circ \widetilde{f} \to \widetilde{B}) \right\}.$

*Proof: see Appendix A* ∎

## III.   ENTANGLEMENT-ASSISTED CIRCUIT KNITTING

In distributed quantum computing (DQC), one considers the implementation of a global unitary operation across two QPUs assisted by shared entanglement resources. The free operations in this setting are separable state preparation and local operations with classical communication (LOCC). In this section, we incorporate entanglement-assisted LOCC into the QPD framework to investigate the relationship between the available entanglement resources and the resulting QPD overhead.

### A.   Entanglement-assisted wire cutting

With LOCC as the free maps, we consider the wire cutting with a pre-established entangled state $|\rho\rangle\rangle$. The $\rho$-assisted QPD overhead for the identity channel from $A$ to $B$ is equal to 1, if the assisting state is a perfect Bell state.

**Lemma 7** *For any dimension d, it holds that*

$$\gamma_{\mathrm{LOCC}} \left( |\Phi_d\rangle\rangle \to \widetilde{\mathbb{1}}_d^{(A \to B)} \right) = \gamma_{\mathrm{LOCC}} \left( \widetilde{\mathbb{1}}_d^{(A \to B)} \to |\Phi_d\rangle\rangle \right) = 1 \tag{17}$$

*Proof: From the general qudit teleportation protocol, we have $\gamma_{\mathrm{LOCC}} \left( |\Phi_d\rangle\rangle \to \widetilde{\mathbb{1}}_d^{(A \to B)} \right) = 1$. Conversely, if the channel $\widetilde{\mathbb{1}}_d^{(A \to B)}$ is available, one can establish $|\Phi_d\rangle\rangle$ by sending half of a locally prepared maximally entangled state from A to B. This shows that both conversions have sampling overhead equal to 1.* ∎

The optimal sampling overhead for the wire cutting derived in [24, 25] can be directly obtained as a corollary of the sub-multiplicity of Lemma 6 as follows.

**Corollary 8 (Resource-assisted wire cutting [24, 25])** *The $\tilde{r}$-assisted QPD overhead over LOCC for an identity channel is given by*

$$\gamma_{\mathrm{LOCC}} \left( \tilde{r} \to \tilde{\mathbb{1}}_d^{(A \to B)} \right) = \gamma_{\mathrm{LOCC}} \left( \tilde{r} \to |\Phi_d\rangle\rangle \right) \tag{18}$$

*Moreover, if $\tilde{r} = |\rho\rangle\rangle$ is a state preparation, then the optimal overhead determined by the fully entangled fraction $F_d(\rho)$,*

$$\gamma_{\text{LOCC}}\left(|\rho\rangle\rangle \to \tilde{\mathbb{1}}_d^{(A \to B)}\right) = \frac{2}{F_d(\rho)} - 1, \tag{19}$$

*where $F_d(\rho)$ is the fully entangled fraction of the state $\rho$,*

$$F_d(\rho) := \max_{\tilde{\epsilon} \in \text{LOCC}} \left\{ \langle\langle \Phi_d | \, \tilde{\epsilon} \, | \rho \rangle\rangle \right\}, \tag{20}$$

*and $|\Phi_d\rangle$ is a $d$-dimensional maximally entangled state.*

Proof: Employing the sub-multiplicity in Lemma 6, one can first prove that the $\tilde{r}$-QPD overheads of the identity channel and the preparation of a maximally entangled state are identical, $\gamma_{LOCC}(|\rho\rangle\rangle \to |\Phi_d\rangle\rangle) = \gamma_{LOCC}(|\rho\rangle\rangle \to \tilde{\mathbb{1}}_d^{(A \to B)})$, as follows

$$\gamma_{LOCC}(\tilde{r} \to |\Phi_d\rangle\rangle) \le \gamma_{LOCC}(\tilde{r} \to \tilde{\mathbb{1}}_d^{(A \to B)}) \, \gamma_{LOCC}(\tilde{\mathbb{1}}_d^{(A \to B)} \to |\Phi_d\rangle\rangle)$$

$$= \gamma_{LOCC}(\tilde{r} \to \tilde{\mathbb{1}}_d^{(A \to B)})$$

$$\le \gamma_{LOCC}(\tilde{r} \to |\Phi_d\rangle\rangle) \, \gamma_{LOCC}(|\Phi_d\rangle\rangle \to \tilde{\mathbb{1}}_d^{(A \to B)})$$

$$= \gamma_{LOCC}(\tilde{r} \to |\Phi_d\rangle\rangle). \tag{21}$$

According to [26, 27], for the case $\tilde{r} = |\rho\rangle\rangle$, the overhead of virtual distillation of $\rho$ for a $d$-dimensional maximally entangled state $|\Phi_d\rangle$ is determined by the fully entangled fraction of $\rho$, where

$$\gamma_{LOCC}(|\rho\rangle\rangle \to |\Phi_d\rangle\rangle) = \frac{2}{F_d(\rho)} - 1. \tag{22}$$

∎

Note that if $|\rho\rangle\rangle$ is a separable state, which has the fully entangled fraction $F_d(|\rho\rangle\rangle) = \frac{1}{d}$, the formula in Eq. (19) reproduces the resource-free QPD overhead $\gamma_{\text{LOCC}} = 2d - 1$ for wire cutting derived in [28–30].

## B. Entanglement-assisted gate cutting

Intuitively, the sampling overhead of a resource-free QPD implementation of a quantum operation is lower-bounded by the entanglement of its Choi state [20]. In the presence of an entangled resource, the bound is relaxed to the difference between the entanglement of the Choi state and that of the resource state.

**Corollary 9** *Let $\hat{U}$ be a unitary across a bipartite system. For any pure state $|\rho\rangle\rangle$ with $\mathcal{E}(|\Phi_U\rangle\rangle) \geq \mathcal{E}(|\rho\rangle\rangle)$, where $\mathcal{E}$ is the entanglement entropy, the $\rho$-assisted QPD overhead over LOCC is lower bounded by*

$$\gamma_{\text{LOCC}}(|\rho\rangle\rangle \to \tilde{U}) \geq 2^{\mathcal{E}(|\Phi_U\rangle\rangle) - \mathcal{E}(|\rho\rangle\rangle)}. \tag{23}$$

Proof: The lower bound is a straightforward result of the right-monotonicity of overhead in Lemma 6. Since the Choi state can be generated from the $\widetilde{U}$ operation acting on the a state $|\Phi_{A,R_A} \otimes \Phi_{B,R_B}\rangle\rangle$, which is maximally entangled with respect to the bipartition $(A,B)|(RA,RB)$, while is separable with respect to the partition $(A, R_A)|(B, R_B)$. The state preparation of $|\Phi_{A,R_A} \otimes \Phi_{B,R_B}\rangle\rangle$ is therefore free for $(A, R_A)|(B, R_B)$-LOCC. According to the right-monotonicity of overhead, it holds then $\gamma_{\text{LOCC}}(|\rho\rangle\rangle \to |\Phi_U\rangle\rangle) \geq \gamma_{\text{LOCC}}(|\rho\rangle\rangle \to \mathbb{1}_{R_A,R_B} \otimes \widetilde{U})$. ∎

In general, it is difficult to determine the minimum overhead for a general unitary. However, for the unitary that can be extracted from its Choi state using LOCC, i.e., $\gamma_{\text{LOCC}}(|\Phi_U\rangle\rangle \to \tilde{U}) = 1$, its overhead over LOCC can be well characterized by the overhead of the Choi state preparation.

**Theorem 10** *For a unitary $\widetilde{U}$ that can be constructed from its Choi state using LOCC, i.e. $\gamma_{\text{LOCC}}(|\Phi_U\rangle\rangle \to \tilde{U}) = 1$, the overhead of a $\tilde{r}$-assisted QPD for the unitary $\widetilde{U}$ is equivalent to the one for Choi state preparation $|\Phi_U\rangle\rangle$,*

$$\gamma_{\text{LOCC}}(\tilde{r} \to \tilde{U}) = \gamma_{\text{LOCC}}(\tilde{r} \to |\Phi_U\rangle\rangle). \tag{24}$$

*If $\tilde{r} = |\rho\rangle\rangle$ is a state preparation, the overhead is given by its $d_U$-dimensional fully entangled fraction $F_{d_U}(\rho)$, where $d_U$ is the operator Schmidt rank of $\widetilde{U}$,*

$$\gamma_{\text{LOCC}}(|\rho\rangle\rangle \to \tilde{U}) = \frac{2}{F_{d_U}(|\rho\rangle\rangle)} - 1. \tag{25}$$

Proof: Since one can also create the Choi state $|\Phi_U\rangle\rangle$ of a unitary $\widetilde{U}$ via separable state preparation, which is a free map, the $\widetilde{U}$-assisted QPD $|\Phi_U\rangle\rangle$ is $\gamma_{\text{LOCC}}(\tilde{U} \to |\Phi_U\rangle\rangle) = 1$. If $\widetilde{U}$ is reversely implementable with the Choi state $|\Phi_U\rangle\rangle$, i.e. $\gamma_{\text{LOCC}}(|\Phi_U\rangle\rangle \to \tilde{U}) = 1$, the state $|\Phi_U\rangle\rangle$ and operation $\tilde{U}$ are now equally resourceful. As a result

$$\gamma_{\text{LOCC}}(\tilde{r} \to |\Phi_U\rangle\rangle) = \gamma_{\text{LOCC}}(\tilde{r} \to \tilde{U})$$

Furthermore, according to [31], $\gamma_{\text{LOCC}}(|\Phi_U\rangle\rangle \to \tilde{U}) = 1$ implies that $|\Phi_U\rangle\rangle$ has uniform Schmidt coefficients, which means $|\Phi_U\rangle\rangle$ is a maximally entangled state with a Schmidt rank of $d_U$. According to [26, 27], the overhead of virtual distillation of $\rho$ for a $d_U$-dimensional maximally entangled

state $|\Phi_d\rangle$ is determined by the fully entangled fraction of $\rho$,

$$\gamma_{LOCC}(|\rho\rangle\rangle \to |\Phi_{d_U}\rangle\rangle) = \frac{2}{F_{d_U}(\rho)} - 1. \tag{26}$$

This completed the proof. ■

Although the overhead of a general unitary is difficult to determine, one can consider a smaller set of free maps consisting of local operations (LO) without classical communication, which admits a well-characterized upper bound on the overhead of a general unitary.

**Theorem 11** *Let $\hat{U}$ be a bipartite unitary with the LUD $\hat{U} = \sum_i \lambda_i \hat{A}_i \otimes \hat{B}_i$, and $|\Phi_2\rangle\rangle$ be a two-qubit maximally entangled state. The overhead of the $|\Phi_2\rangle\rangle$-assisted QPD for $\widetilde{U}$ over LO is upper bounded by*

$$\gamma_{\mathrm{LO}}(|\Phi_2\rangle\rangle \to \tilde{U}) \le ||\boldsymbol{\lambda}||_1^2. \tag{27}$$

*This upper bound can be achieved via the following explicit QPD configuration:*

$$\tilde{U} = \gamma_{\mathcal{Q}} \sum_{i,j} p_{i,j} \sum_{\boldsymbol{m}\in\{0,1\}^{\otimes 2}} (-1)^{|\boldsymbol{m}|} \langle\langle \boldsymbol{m}^{(anc.)}| \, \tilde{F}_{i,j} \, |\Phi_2\rangle\rangle. \tag{28}$$

*In which $p_{i,j} = \frac{\lambda_i \lambda_j}{||\boldsymbol{\lambda}||_1^2}$.*

Proof: See Appendix B. ■

The explicit implementation of the QPD in Eq. (28) is illustrated in Fig. 5. One can rewrite the QPD in Eq (28) to the QPD given in [18, 19] (see Appendix C for details). It is worth noting that the overhead of the QPD in Eq. (28) exhibits a multiplicative structure in its sampling overhead, since the one-norm is multiplicative under the tensor product of the vector, i.e., $||s \otimes s'||_1 = ||s||_1 \cdot ||s'||_1$. Moreover, for a KAK-like unitary, this sampling overhead coincides with the regularized optimal sampling overhead [13, 18], which is given by

$$\gamma_{\mathrm{LOCC}}^{\infty}(\tilde{U}) := \lim_{n\to\infty} \sqrt[n]{\gamma_{\mathrm{LOCC}}(\tilde{U}^{\otimes n})} = ||\boldsymbol{\lambda}||_1^2.$$

We can therefore conclude that, for a target KAK-like unitary, the overhead of a $\Phi_2$-state-assisted-QPD over local operations (LO) is upper-bounded by the regularized optimal sampling overhead of resource-free QPD over the LOCC. we can conclude that

$$\gamma_{\mathrm{LO}}(|\Phi_2\rangle\rangle \to \tilde{U}) \le \gamma_{\mathrm{LOCC}}^{\infty}(\tilde{U}).$$

This result for a Bell-state-assisted QPD can be extended to the case where the pre-shared entangled resource is not maximally entangled.
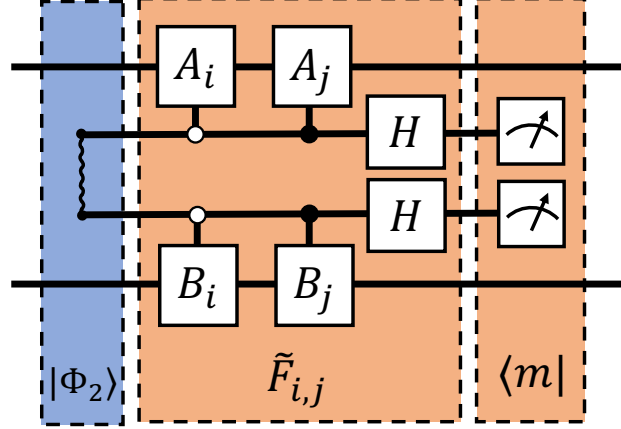
FIG. 5. Here we used the wiggle line to denote the Bell state $|\Phi_2\rangle = \frac{1}{\sqrt{2}}(|0,0\rangle + |1,1\rangle)$. The solid dot denotes the control unitary operates at $|1\rangle\langle 1|$, and the white dot denotes the control unitary operates at $|0\rangle\langle 0|$.

**Corollary 12** *Let $|\psi(r)\rangle = \sqrt{\frac{1+r}{2}}|0,0\rangle + \sqrt{\frac{1-r}{2}}|1,1\rangle$ with $r \in [0,1]$. For an arbitrary target unitary $\hat{U} = \sum_i \lambda_i \hat{A}_i \otimes \hat{B}_i$, it has a QPD overhead upper bounded as follows*

$$\gamma_{\text{LO}}(|\psi(r)\rangle\rangle \to \tilde{U}) \leq ||\boldsymbol{\lambda}||_2^2 + \min\{\frac{1}{\sqrt{1-r^2}}, 2\}(||\boldsymbol{\lambda}||_1^2 - ||\boldsymbol{\lambda}||_2^2). \tag{29}$$

Proof: With the same QPD in Eq. (28) shown in Fig 5, one can now replace the Bell state by $|\psi(r)\rangle$,

$$\sqrt{1-r^2} \sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}|\, \tilde{F}_{i,j}|\Phi_2\rangle\rangle = \sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}|\, \tilde{F}_{i,j}|\psi(r)\rangle\rangle \forall i,j \tag{30}$$

For the $i = j$ terms, we can just leave the entangled state, and then implement the local unitary $\tilde{A}_i \otimes \tilde{B}_i$ directly. As a result, the QPD becomes

$$\tilde{U} = \sum_i \lambda_i^2 \tilde{A}_i \otimes \tilde{B}_i + \frac{2}{\sqrt{1-r^2}} \sum_{i>j} \lambda_i \lambda_j \sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}|\, \tilde{F}_{i,j}|\psi(r)\rangle\rangle. \tag{31}$$

So we can conclude that

$$\gamma_{\text{LO}}(|\psi(r)\rangle\rangle \to \tilde{U}) \leq ||\boldsymbol{\lambda}||_2^2 + \frac{(||\boldsymbol{\lambda}||_1^2 - ||\boldsymbol{\lambda}||_2^2)}{\sqrt{1-r^2}}. \tag{32}$$

But one can see that, when $r$ is too small, it may be larger than the overhead given by proposition 3. So we add a min function, such that when $2 \leq \frac{1}{\sqrt{1-r^2}}$, we just stay with the usual gate cut. ∎

In addition to the upper bound, we can establish a lower bound on the sampling overhead of gate cutting assisted by a single Bell pair.

**Proposition 13** *The overhead of a Bell-state-assisted gate cutting of a KAK-like unitary over LOCC is lower bounded by*

$$||\boldsymbol{\lambda}||_1^2 - 2\lambda_1 \sum_{i>1} \lambda_i \leq \gamma_{\text{LOCC}}(|\Phi_2\rangle\rangle \to \tilde{U}), \tag{33}$$

*where $\lambda_i$ is the Schmidt coefficient of $\tilde{U}$ with decreasing order.*

*Proof:   First, we have*

$$\gamma_{\text{LOCC}}(|\Phi_2\rangle\rangle \to |\Phi_U\rangle\rangle) \leq \gamma_{\text{LOCC}}(|\Phi_2\rangle\rangle \to \tilde{U}). \tag{34}$$

*Since the LOCC cannot increase the Schmidt rank, we can use the robustness of the Schmidt rank [32]. Which is given by*

$$\gamma_{\text{LOCC}}(|\Phi_2\rangle\rangle \to |\Phi_U\rangle\rangle) = \min_{\mathcal{R}_s(|\rho_\pm\rangle\rangle) \leq 2} \left\{ 1 + 2t \Big| |\rho_+\rangle\rangle = \frac{|\Phi_U\rangle\rangle + t|\rho_-\rangle\rangle}{1+t} \right\}$$

$$= (\sum_i \lambda_i)^2 - 2\lambda_1 \sum_{i>1} \lambda_i \leq \gamma_{\text{LOCC}}(|\Phi_2\rangle\rangle \to \tilde{U}). \tag{35}$$

*This completes the proof.* ∎

For illustration, we numerically generate two parallel Haar-random two-qubit gates and show the probability distribution corresponding to different bounds in Fig. 6.
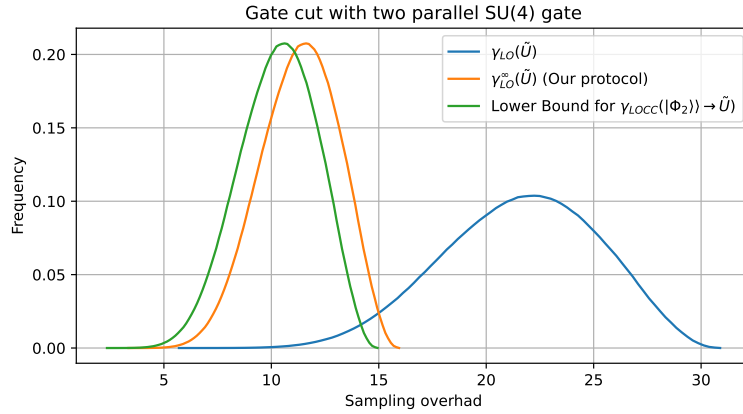


FIG. 6.   We randomly sample $10^8$ of two parallel SU(4) gates under the Haar measure, and plot the probability distribution of their overhead factors. It shows that with a pair of Bell states, the overhead can be reduced roughly by half.

From another perspective, beyond reducing sampling overhead, our QPD framework also enhances entanglement efficiency in distributed quantum computing (DQC). In particular, consider the controlled-phase rotation gate $\hat{U}(\theta) = \text{diag}(1, 1, 1, e^{i\theta})$. It is known that this gate requires one

Bell pair per application when implemented via gate teleportation using LOCC, achieving this with zero sampling overhead [5].

In contrast, under our QPD in Eq. (B4), the Bell state is needed only in the off-diagonal terms $(i > j)$, while the diagonal terms $(i = j)$ can be implemented with local unitaries alone. Thus, the probability of using a Bell state in each round is

$$P_B = 1 - \frac{1}{||\boldsymbol{\lambda}||_1^2}. \tag{36}$$

Taking the sampling overhead into account, the expected number of Bell states used per execution becomes

$$\langle \#\Phi_2 \rangle = P_B \cdot ||\boldsymbol{\lambda}||_1^4 = |\sin\theta| + |\sin\theta|^2, \tag{37}$$

for the controlled rotation gate. Here, we further multiply the sampling overhead to make a fair comparison with the entanglement-assisted DQC. This result reveals that for $\theta \leq \sin^{-1}\left(\frac{\sqrt{5}-1}{2}\right) \approx 0.42\pi$, our QPD yields both: (1) a lower sampling overhead than standard gate cutting, and (2) a lower entanglement cost than gate teleportation (since $\langle \#\Phi_2 \rangle \leq 1$). Moreover, no classical communication is required.

This example illustrates the hybrid nature of our approach, which combines the advantages of gate cutting and gate teleportation. It also reveals an intrinsic trade-off between sampling overhead and entanglement consumption, providing valuable insight even though the presented decomposition is not necessarily optimal.

## IV.  CONCLUSION AND DISCUSSION

This work introduces a framework of entanglement-assisted circuit knitting formulated through the concept of virtual distillation of quantum operations, unifying entanglement-assisted DQC and circuit knitting. Within this framework, LOCC and separable state preparation are treated as free maps. Building on Refs. [24, 25], we show that for certain classes of unitaries, the optimal sampling overhead can be achieved using arbitrary pre-shared entangled states. We further establish a general lower bound on this overhead determined by the entanglement entropy of the resource, revealing the intrinsic limitation of entanglement-assisted gate cutting.

We then further consider the case where only local operations (LO) are treated as free maps and provide a constructive example demonstrating how a single Bell state can be used to implement a gate cut. In particular, the resulting sampling overhead matches the regularized optimal overhead

of the standard gate-cutting protocol without entanglement, clearly showing that entanglement assistance can enhance sampling efficiency. Moreover, for certain classes of unitary gates, our protocol achieves higher entanglement efficiency than conventional gate-teleportation schemes.

Consistent with other recent studies in circuit knitting, our present analysis remains limited to certain classes of target unitaries. Future extensions to general bipartite unitary operations would be both conceptually and practically significant.

Another important question that remains open in this work concerns the advantage provided by classical communication. Specifically, what is the precise role of classical communication in entanglement-assisted circuit knitting? In entanglement-assisted DQC, classical communication is known to be essential. However, for gate cutting, several studies have shown that it does not reduce the sampling overhead [13, 18, 19].

A natural extension of this question is to include catalytic resources [33, 34]. For the case of entanglement with LOCC as the free map, an advantage from catalytic resources may indeed emerge, analogous to the catalytic effect observed in state teleportation [35].

### Appendix A: Optimal sampling overhead with pre-established resource

The proofs of the properties of resource-assisted QPD overhead in Lemma 6 are provided as follows.

**Sub-multiplicity:** $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{C}) \leq \gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B})\gamma_{\mathbb{F}}(\tilde{B} \to \tilde{C})$.

Proof: Based on Eq (13), suppose we have the QPD for $\tilde{A} \to \tilde{B}$ and $\tilde{B} \to \tilde{C}$ as

$$\tilde{B} = \mathcal{N} \sum_x p_x \sum_a \wp_{a,x} \langle\langle M_{a|x}| \circ \tilde{F}_x^{(1)} \circ \tilde{A} \circ \tilde{F}_x^{(0)} \tag{A1}$$

and

$$\tilde{C} = \mathcal{N}' \sum_y p'_y \sum_b \wp'_{b,y} \langle\langle W_{b|y}| \circ \tilde{L}_y^{(1)} \circ \tilde{B} \circ \tilde{L}_y^{(0)} \tag{A2}$$

Then we can immediately construct a QPD for $\tilde{A} \to \tilde{C}$ as

$$\tilde{C} = \mathcal{N}\mathcal{N}' \sum_{x,y} p_x p'_y \sum_{a,b} \wp_{a,x} \wp'_{b,y} \langle\langle M_{a|x} \otimes W_{b|y}| \circ (\tilde{L}_y^{(1)} \circ \tilde{F}_x^{(1)}) \circ \tilde{A} \circ (\tilde{F}_x^{(0)} \circ \tilde{L}_y^{(0)}) \qquad (A3)$$

So by taking QPDs that achieve the $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B})$ and $\gamma_{\mathbb{F}}(\tilde{B} \to \tilde{C})$, we obtain a QPD of $\tilde{A} \to \tilde{C}$ with sampling overhead $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B})\gamma_{\mathbb{F}}(\tilde{B} \to \tilde{C})$. This completed the proof. ■

**Ordering:** $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B}) = 1 \Rightarrow \gamma_{\mathbb{F}}(\tilde{A}) \geq \gamma_{\mathbb{F}}(\tilde{B})$.

Proof: With the sub-multiplicity, we can take $\gamma_{\mathbb{F}}(\tilde{f} \to \tilde{B}) \leq \gamma_{\mathbb{F}}(\tilde{f} \to \tilde{A})\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B})$, in which $\tilde{f} \in \mathbb{F}$. So if we have $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B}) = 1$, the inequality becomes $\gamma_{\mathbb{F}}(\tilde{B}) = \gamma_{\mathbb{F}}(\tilde{f} \to \tilde{B}) \leq \gamma_{\mathbb{F}}(\tilde{f} \to \tilde{A}) = \gamma_{\mathbb{F}}(\tilde{A})$. Notice that the "$\Leftarrow$" direction does not hold in general. ■

**Right-convexity:** $\gamma_{\mathbb{F}}(\tilde{A} \to p\tilde{B} + (1-p)\tilde{C}) \leq p\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B}) + (1-p)\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{C}) \; \forall p \in [0,1]$.

Proof: With the QPD for $\tilde{A} \to \tilde{B}$ as

$$\tilde{B} = \mathcal{N} \sum_x p_x \sum_a \wp_{a,x} \langle\langle M_{a|x}| \circ \tilde{F}_x^{(1)} \circ \tilde{A} \circ \tilde{F}_x^{(0)}. \qquad (A4)$$

and QPD for $\tilde{A} \to \tilde{C}$ as

$$\tilde{C} = \mathcal{N}' \sum_y p'_y \sum_b s'_y(b) \langle\langle W_{b|y}| \circ \tilde{G}_x^{(1)} \circ \tilde{A} \circ \tilde{G}_x^{(0)}. \qquad (A5)$$

We immaterially obtain

$$p\tilde{B} + (1-p)\tilde{C} = p\left(\mathcal{N} \sum_x p_x \sum_a \wp_{a,x} \langle\langle M_{a|x}| \circ \tilde{F}_x^{(1)} \circ \tilde{A} \circ \tilde{F}_x^{(0)}\right)$$
$$+ (1-p)\left(\mathcal{N}' \sum_y p'_y \sum_b s'_y(b) \langle\langle W_{b|y}| \circ \tilde{G}_x^{(1)} \circ \tilde{A} \circ \tilde{G}_x^{(0)}\right) \qquad (A6)$$

for all $p \in [0,1]$, hence we conclude that

$$\gamma_{\mathbb{F}}(\tilde{A} \to p\tilde{B} + (1-p)\tilde{C}) \leq p\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B}) + (1-p)\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{C}) \; \forall p \in [0,1]. \qquad (A7)$$

Similarly, for $\tilde{A} \to \tilde{f} \circ \tilde{B}$, it follows that

$$\tilde{f} \circ \tilde{B} = \mathcal{N} \sum_x p_x \sum_a \wp_{a,x} \langle\langle M_{a|x}| \circ (\tilde{f} \circ \tilde{F}_x^{(1)}) \circ \tilde{A} \circ \tilde{F}_x^{(0)} \qquad (A8)$$

Hence both $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{f} \circ \tilde{B})$ and $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B} \circ \tilde{f})$ should be smaller then $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B})$, which is the desired results. ■

**Left-convexity:** $p\gamma_{\mathbb{F}}(\widetilde{A} \to \widetilde{C}) + (1-p)\gamma_{\mathbb{F}}(\widetilde{B} \to \widetilde{C}) \leq \gamma_{\mathbb{F}}(p\widetilde{A} + (1-p)\widetilde{B} \to \widetilde{C}) \; \forall p \in [0,1]$.

Proof: Follows similarly to the right-convexity. ∎

**Right-monotonicity:** $\tilde{f} \in \mathbb{F} \Rightarrow \max\{\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{f} \circ \tilde{B}), \gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B} \circ \tilde{f})\} \leq \gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B})$.

Proof: With the QPD for $\tilde{A} \to \tilde{B}$ as

$$\tilde{B} = \mathcal{N} \sum_x p_x \sum_a \wp_{a,x} \langle\langle M_{a|x}| \circ \tilde{F}_x^{(1)} \circ \tilde{A} \circ \tilde{F}_x^{(0)}. \tag{A9}$$

For $\tilde{A} \to \tilde{B} \circ \tilde{f}$, it follows that

$$\tilde{B} \circ \tilde{f} = \mathcal{N} \sum_x p_x \sum_a \wp_{a,x} \langle\langle M_{a|x}| \circ \tilde{F}_x^{(1)} \circ \tilde{A} \circ (\tilde{F}_x^{(0)} \circ \tilde{f}). \tag{A10}$$

Similarly, for $\tilde{A} \to \tilde{f} \circ \tilde{B}$, it follows that

$$\tilde{f} \circ \tilde{B} = \mathcal{N} \sum_x p_x \sum_a \wp_{a,x} \langle\langle M_{a|x}| \circ (\tilde{f} \circ \tilde{F}_x^{(1)}) \circ \tilde{A} \circ \tilde{F}_x^{(0)} \tag{A11}$$

Hence both $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{f} \circ \tilde{B})$ and $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B} \circ \tilde{f})$ should be smaller then $\gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B})$, which is the desired results. ∎

**Left-monotonicity:** $\tilde{f} \in \mathbb{F} \Rightarrow \gamma_{\mathbb{F}}(\tilde{A} \to \tilde{B}) \leq \min\{\gamma_{\mathbb{F}}(\tilde{f} \circ \tilde{A} \to \tilde{B}), \gamma_{\mathbb{F}}(\tilde{A} \circ \tilde{f} \to \tilde{B})\}$.

Proof: Follows similarly to the Right-monotonicity. ∎

### Appendix B: Deviation on the QPD with Bell state

We will prove Theorem 11 by the directed contraction of the circuit-cutting protocol. For the simplicity, we denote $\hat{\Lambda}_i = \hat{A}_i \otimes \hat{B}_i$, where the LUD of the target unitary as:

$$\tilde{U} = \sum_{i,j} \lambda_i \lambda_j (\hat{A}_i \otimes \hat{B}_i) \otimes (\hat{A}_j \otimes \hat{B}_j)^* = \sum_{i,j} \lambda_i \lambda_j \hat{\Lambda}_i \otimes \hat{\Lambda}_j^*$$

$$= \sum_i \lambda_i^2 \tilde{\Lambda}_i + \sum_{i \neq j} \lambda_i \lambda_j \hat{\Lambda}_i \otimes \hat{\Lambda}_j^* \tag{B1}$$

The $i = j$ part is already a combination of local channels, so our main problem now is to find a way to deal with the $i \neq j$ part, that is

$$\sum_{i \neq j} \lambda_i \lambda_j \hat{\Lambda}_i \otimes \hat{\Lambda}_j^* = \sum_{i > j} \lambda_i \lambda_j \left( \hat{\Lambda}_i \otimes \hat{\Lambda}_j^* + \hat{\Lambda}_j \otimes \hat{\Lambda}_i^* \right) \tag{B2}$$

By define a rank-2 operator as $\hat{\Pi}_{i,j}^{(\pm)} = \frac{1}{2}\left(\hat{\Lambda}_i \pm \hat{\Lambda}_j\right)$, one can rewire the "non-diagonal" parts as

$$\tilde{\Pi}_{i,j}^{(\pm)} = \frac{1}{4}\left(\tilde{\Lambda}_i + \tilde{\Lambda}_j\right) \pm \frac{1}{4}\left(\hat{\Lambda}_i \otimes \hat{\Lambda}_j^* + \hat{\Lambda}_j \otimes \hat{\Lambda}_i^*\right)$$

$$\Rightarrow \tilde{\Pi}_{i,j}^{(+)} - \tilde{\Pi}_{i,j}^{(-)} = \frac{1}{2}\left(\hat{\Lambda}_i \otimes \hat{\Lambda}_j^* + \hat{\Lambda}_j \otimes \hat{\Lambda}_i^*\right) \tag{B3}$$

Therefore, we obtain the following decomposition

$$\tilde{U} = \sum_i \lambda_i^2 \tilde{\Lambda}_i + 2\sum_{i>j} \lambda_i \lambda_j (\tilde{\Pi}_{i,j}^{(+)} - \tilde{\Pi}_{i,j}^{(-)}). \tag{B4}$$

Now we will show that it can be implemented through a local operation with one Bell state. First, we set our initial state as $|Q_a\rangle \otimes |\Phi_2\rangle \otimes |Q_b\rangle$. Then apply the local unitary $\hat{U}_{AB} = \hat{U}_A \otimes \hat{U}_B$, with

$$\hat{U}_{i,j}^{(A)} = \hat{A}_i \otimes |+\rangle\langle+| + \hat{A}_j \otimes |-\rangle\langle-| \tag{B5}$$

$$\hat{U}_{i,j}^{(B)} = |+\rangle\langle+| \otimes \hat{B}_i + |-\rangle\langle-| \otimes \hat{B}_j \tag{B6}$$

and our free map is just $\tilde{F}_{i,j} = \hat{U}_{i,j}^{(A)} \otimes \hat{U}_{i,j}^{(B)} \in \text{LO}$. Finally, we measure the two auxiliary qubits in the middle with the computational basis $\langle\langle m|$. The final state becomes

$$\langle\langle m| \tilde{F}_{i,j} |Q_a \otimes \Phi_2 \otimes Q_b\rangle\rangle = \begin{cases} \frac{1}{2}\tilde{\Pi}_{i,j}^{(+)}|Q_a, Q_b\rangle\rangle & \text{if } m \in \{00, 11\} \\ \frac{1}{2}\tilde{\Pi}_{i,j}^{(-)}|Q_a, Q_b\rangle\rangle & \text{if } m \in \{01, 10\} \end{cases}. \tag{B7}$$

So by setting the post-processing function

$$s_{i,j}(m) := (-1)^{m_a + m_b} = (-1)^{|m|} \tag{B8}$$

We obtain

$$(\tilde{\Pi}_{i,j}^{(+)} - \tilde{\Pi}_{i,j}^{(-)})|Q_a \otimes Q_b\rangle\rangle = \sum_{m \in \{0,1\}^2} (-1)^{|m|} \langle\langle m| \tilde{F}_{i,j} |Q_a \otimes \Phi_2 \otimes Q_b\rangle\rangle. \tag{B9}$$

Figure 5 gives a circuit diagram for the above equation. Note that the probability is equal for getting each outcome $m$.

Finally, we have

$$\tilde{U}|Q_a \otimes Q_b\rangle\rangle = \sum_{i,j} \lambda_i \lambda_j \sum_{m \in \{0,1\}^2} (-1)^{|m|} \langle\langle m| \tilde{F}_{i,j} |Q_a \otimes \Phi_2 \otimes Q_b\rangle\rangle$$

$$= \gamma_Q \sum_{i,j} p_{i,j} \sum_{m \in \{0,1\}^2} (-1)^{|m|} \langle\langle m| \tilde{F}_{i,j} |Q_a \otimes \Phi_2 \otimes Q_b\rangle\rangle \tag{B10}$$

$$\Rightarrow \tilde{U} = \gamma_Q \sum_{i,j} p_{i,j} \sum_{m \in \{0,1\}^2} (-1)^{|m|} \langle\langle m| \tilde{F}_{i,j} |\Phi_2\rangle\rangle \tag{B11}$$

In which

$$p_{i,j} = \gamma_{\mathcal{Q}}^{-1} \lambda_i \lambda_j \text{ and } \gamma_{\mathcal{Q}} = ||\boldsymbol{\lambda}||_1^2 \tag{B12}$$

Hence, we obtain the QPD with overhead $||\boldsymbol{\lambda}||_1^2$, which completes the proof.

## Appendix C: Deviation on the QPD without Bell state

Here, we will going to demonstrate how to use the QPD in Eq. 28 to construct the QPD over LO given by [18, 19]. First, the Bell state has the QPD given by

$$|\Phi_2\rangle\rangle = \frac{1}{2}(|00\rangle\rangle + |11\rangle\rangle\rangle) + \frac{1}{2}(|\sigma_x \otimes \sigma_x\rangle\rangle - |\sigma_y \otimes \sigma_y\rangle\rangle). \tag{C1}$$

So for the $i, j$ terms, it has

$$\sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}| \tilde{F}_{i,j}|\Phi_2\rangle\rangle = \frac{1}{2} \sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}| \tilde{F}_{i,j}(|00\rangle\rangle + |11\rangle\rangle\rangle)$$
$$+ \frac{1}{2} \sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}| \tilde{F}_{i,j}(|\sigma_x \otimes \sigma_x\rangle\rangle - |\sigma_y \otimes \sigma_y\rangle\rangle) \tag{C2}$$

However, from the Figure 5, one can see that

$$\sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}| \tilde{F}_{i,j}(|00\rangle\rangle + |11\rangle\rangle\rangle) = 0. \tag{C3}$$

Which makes

$$\tilde{U} = \sum_i \lambda_i^2 \tilde{\Lambda}_i + \sum_{i\geq j} \lambda_i \lambda_j \sum_{\boldsymbol{m}\in\{0,1\}^2} (-1)^{|\boldsymbol{m}|} \langle\langle\boldsymbol{m}| \tilde{F}_{i,j}(|\sigma_x \otimes \sigma_x\rangle\rangle - |\sigma_y \otimes \sigma_y\rangle\rangle) \tag{C4}$$

Notice that, to make $(|\sigma_x \otimes \sigma_x\rangle\rangle - |\sigma_y \otimes \sigma_y\rangle\rangle)$ a QPD of density matrix, the corresponding normalization factor (negativity) is 4. So the total sampling overhead is

$$\gamma_{\mathcal{Q}} = ||\boldsymbol{\lambda}||_2^2 + 4\sum_{i>j} \lambda_i \lambda_j = 2||\lambda||_1^2 - ||\lambda||_2^2. \tag{C5}$$

which can be shown to be optimal [18, 19], if $\hat{U}$ is KAK like.

---

[1] Marcello Caleffi, Michele Amoretti, Davide Ferrari, Jessica Illiano, Antonio Manzalini, and Angela Sara Cacciapuoti. Distributed quantum computing: A survey. *Computer Networks*, 254:110672, December 2024.

[2] David Barral, F. Javier Cardama, Guillermo Díaz-Camacho, Daniel Faílde, Iago F. Llovo, Mariamo Mussa-Juane, Jorge Vázquez-Pérez, Juan Villasuso, César Piñeiro, Natalia Costas, Juan C. Pichel, Tomás F. Pena, and Andrés Gómez. Review of distributed quantum computing: From single qpu to high performance quantum computing. *Computer Science Review*, 57:100747, August 2025.

[3] Johannes Knörzer, Xiaoyu Liu, Benjamin F. Schiffer, and Jordi Tura. Distributed quantum information processing: A review of recent progress, 2025.

[4] Daniel Gottesman and Isaac L. Chuang. Demonstrating the viability of universal quantum computation using teleportation and single-qubit operations. *Nature*, 402:390–393, 1999.

[5] J. Eisert, K. Jacobs, P. Papadopoulos, and M. B. Plenio. Optimal local implementation of nonlocal quantum gates. *Physical Review A*, 62(5):052317, October 2000.

[6] Pablo Andrés-Martínez and Chris Heunen. Automated distribution of quantum circuits via hypergraph partitioning. *Physical Review A*, 100:032308, Sep 2019.

[7] Jun-Yi Wu, Kosuke Matsui, Tim Forrer, Akihito Soeda, Pablo Andrés-Martínez, Daniel Mills, Luciana Henaut, and Mio Murao. Entanglement-efficient bipartite-distributed quantum computing. *Quantum*, 7:1196, December 2023.

[8] Pablo Andres-Martinez, Tim Forrer, Daniel Mills, Jun-Yi Wu, Luciana Henaut, Kentaro Yamamoto, Mio Murao, and Ross Duncan. Distributing circuits over heterogeneous, modular quantum computing network architectures. *Quantum Science and Technology*, 9(4):045021, aug 2024.

[9] Tianyi Peng, Aram W. Harrow, Maris Ozols, and Xiaodi Wu. Simulating large quantum circuits on a small quantum computer. *Physical Review Letters*, 125(15):150504, October 2020.

[10] Kosuke Mitarai and Keisuke Fujii. Overhead for simulating a non-local channel with local channels by quasiprobability sampling. *Quantum 5, 388 (2021)*, 5:388, January 2020.

[11] Christophe Piveteau, David Sutter, and Stefan Woerner. Quasiprobability decompositions with reduced sampling overhead. *npj Quantum Information*, 8(1), February 2022.

[12] Kosuke Mitarai and Keisuke Fujii. Constructing a virtual two-qubit gate by sampling single-qubit operations. *New Journal of Physics*, 23(2):023021, February 2021.

[13] Christophe Piveteau and David Sutter. Circuit knitting with classical communication. *IEEE Transactions on Information Theory*, 70(4):2734–2745, April 2024.

[14] Charles H. Bennett, Gilles Brassard, Claude Crépeau, Richard Jozsa, Asher Peres, and William K. Wootters. Teleporting an unknown quantum state via dual classical and einstein-podolsky-rosen channels. *Physical Review Letters*, 70:1895–1899, Mar 1993.

[15] Hakop Pashayan, Joel J. Wallman, and Stephen D. Bartlett. Estimating outcome probabilities of quantum circuits using quasiprobabilities. *Physical Review Letters*, 115(7):070501, August 2015.

[16] Christian Ufrecht, Maniraman Periyasamy, Sebastian Rietsch, Daniel D. Scherer, Axel Plinge, and Christopher Mutschler. Cutting multi-control quantum gates with zx calculus. *Quantum*, 7:1147, October 2023.

[17] Christian Ufrecht, Laura S. Herzog, Daniel D. Scherer, Maniraman Periyasamy, Sebastian Rietsch, Axel Plinge, and Christopher Mutschler. Optimal joint cutting of two-qubit rotation gates. *Physical Review A*, 109(5):052440, May 2024.

[18] Lukas Schmitt, Christophe Piveteau, and David Sutter. Cutting circuits with multiple two-qubit unitaries. *Quantum*, 9:1634, February 2025.

[19] Aram W. Harrow and Angus Lowe. Optimal quantum circuit cuts with application to clustered hamiltonian simulation. *PRX Quantum*, 6(1):010316, January 2025.

[20] Mingrui Jing, Chengkai Zhu, and Xin Wang. Circuit knitting facing exponential sampling-overhead scaling bounded by entanglement cost. *Physical Review A*, 111(1):012433, January 2025.

[21] Giacomo Mauro D'Ariano, Giulio Chiribella, and Paolo Perinotti. *Quantum Theory from First Principles: An Informational Approach*. Cambridge University Press, November 2016.

[22] Eric Chitambar and Gilad Gour. Quantum resource theories. *Reviews of Modern Physics*, 91(2):025001, April 2019.

[23] We defined the orthogonality between operators, by the Hilbert-Schmidt inner product.

[24] Marvin Bechtold, Johanna Barzen, Frank Leymann, and Alexander Mandl. Cutting a wire with non-maximally entangled states. In *2024 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages 1136–1145. IEEE, May 2024.

[25] Marvin Bechtold, Johanna Barzen, Frank Leymann, Alexander Mandl, and Felix Truger. Joint wire cutting with non-maximally entangled states. *Advanced Quantum Technologies*, 8(5), January 2025.

[26] Xiao Yuan, Bartosz Regula, Ryuji Takagi, and Mile Gu. Virtual quantum resource distillation. *Physical Review Letters*, 132(5):050203, February 2024.

[27] Ryuji Takagi, Xiao Yuan, Bartosz Regula, and Mile Gu. Virtual quantum resource distillation: General framework and applications. *Physical Review A*, 109(2):022403, February 2024.

[28] Lukas Brenner, Christophe Piveteau, and David Sutter. Optimal wire cutting with classical communication. *IEEE Transactions on Information Theory*, 71(10):7742–7752, October 2025.

[29] Edwin Pednault. An alternative approach to optimal wire cutting without ancilla qubits. arXiv, March 2023.

[30] Hiroyuki Harada, Kaito Wada, and Naoki Yamamoto. Doubly optimal parallel wire cutting without ancilla qubits. *PRX Quantum*, 5(4):040308, October 2024.

[31] Dan Stahlke and Robert B. Griffiths. Entanglement requirements for implementing bipartite unitary operations. *Physical Review A*, 84(3):032316, September 2011.

[32] Nathaniel Johnston, Chi-Kwong Li, Sarah Plosker, Yiu-Tung Poon, and Bartosz Regula. Evaluating the robustness of k-coherence and k-entanglement. *Physical Review A*, 98(2):022328, August 2018.

[33] Chandan Datta, Tulja Varun Kondra, Marek Miller, and Alexander Streltsov. Catalysis of entanglement and other quantum resources. *Reports on Progress in Physics*, 86(11):116002, October 2023.

[34] Patryk Lipka-Bartosik, Henrik Wilming, and Nelly H.Y. Ng. Catalysis in quantum information theory. *Reviews of Modern Physics*, 96(2):025005, June 2024.

[35] Patryk Lipka-Bartosik and Paul Skrzypczyk. Catalytic quantum teleportation. *Physical Review Letters*, 127(8):080502, August 2021.