

Overhead for simulating a non-local channel with local channels by quasiprobability sampling

Kosuke Mitarai^{1,2,3} and Keisuke Fujii^{1,2,4}

¹Graduate School of Engineering Science, Osaka University, 1-3 Machikaneyama, Toyonaka, Osaka 560-8531, Japan.

²Center for Quantum Information and Quantum Biology, Institute for Open and Transdisciplinary Research Initiatives, Osaka University, Japan.

³JST, PRESTO, 4-1-8 Honcho, Kawaguchi, Saitama 332-0012, Japan.

⁴Center for Emergent Matter Science, RIKEN, Wako Saitama 351-0198, Japan

January 26, 2021

As the hardware technology for quantum computing advances, its possible applications are actively searched and developed. However, such applications still suffer from the noise on quantum devices, in particular when using two-qubit gates whose fidelity is relatively low. One way to overcome this difficulty is to substitute such non-local operations by local ones. Such substitution can be performed by decomposing a non-local channel into a linear combination of local channels and simulating the original channel with a quasiprobability-based method. In this work, we first define a quantity that we call channel robustness of non-locality, which quantifies the cost for the decomposition. While this quantity is challenging to calculate for a general non-local channel, we give an upper bound for a general two-qubit unitary channel by providing an explicit decomposition. The decomposition is obtained by generalizing our previous work whose application has been restricted to a certain form of two-qubit unitary. This work develops a framework for a resource reduction suitable for first-generation quantum devices.

1 Introduction

We now have a programmable quantum device whose dynamics cannot be simulated by a classical computer within its runtime [1]. However, the capability of such devices is rather limited because of the absence of the quantum error correction. They are frequently referred to as noisy intermediate scale quantum (NISQ) devices [2]. There has been a substantial amount of research efforts to develop useful applications of NISQ devices in recent years [3–9]. The weakness of NISQ devices is that the number of qubits, the fidelities of gates, and

the connectivity are limited. The gate fidelities are especially restricted for two-qubit entangling gates. One approach to circumvent such limitation is to use so-called variational quantum algorithms. They employ parametrized quantum circuits and optimize the parameters to perform a given task. In such algorithms, we frequently construct the largest possible circuit allowed on a device to maximize the advantage of the use of quantum devices.

While this approach is promising as it can in principle employ such circuits, such algorithms can still be improved if one can perform further resource reduction. For example, if we can reduce the number of qubits or two-qubit gates required to obtain an output from a certain quantum circuit, it would widen the range of circuits that can be used for variational algorithms. To this end, a few approaches have been proposed. One is to decompose a large circuit into smaller ones by “cutting” circuits using a tomography-like method [10]. Also, in Ref. [11], we have presented a method to “cut” a certain non-local gate by decomposing it into a linear combination of local operations. These approaches share the same property that the overhead for the decomposition, which in this context is defined by the number of circuit runs that is required to achieve a desired accuracy of the output, scales exponentially to the number of cuts performed.

They can be also understood as techniques for performing a quasiprobability decomposition of quantum channels. Quasiprobability distribution, which is defined by a set of complex numbers $\{q_i\}$ satisfying $\sum_i q_i = 1$, have recently found a wide range of applications in the area of quantum computing such as error mitigation for NISQ devices [12, 13] and classical simulation of near-Clifford quantum circuits [14–18]. In particular, Refs. [12, 13, 17, 18] considered a quasiprobability-based simulation of quantum channels; if a quantum channel Φ can be decomposed as $\Phi = \sum_i c_i \Phi_i$ where Φ_i and c_i are respectively a chan-

Kosuke Mitarai: mitarai@qc.ee.es.osaka-u.ac.jp

nel and a complex coefficient, Φ can be simulated by sampling Φ_i with probability proportional to $|c_i|$ and processing the phase of c_i with classical post-processing. The overhead of simulating the channel Φ using this decomposition is quantified by $\sum_i |q_i|$. If we perform such a decomposition multiple times, the overhead is quantified by the product of $\sum_i |q_i|$, thus leading to an exponential overhead to the number of decomposition performed. Refs. [12, 13] have developed techniques to build inverse channels of noise channels using an experimentally available set of quantum gates. As a technique for a classical simulation, Refs. [17, 18] has considered a quasiprobability decomposition of a non-Clifford channel into Clifford ones. In this context, we can view the decomposition performed in Ref. [10] as a quasiprobability decomposition of the identity channel into a measurement and state-preparation channel, and one in Ref. [11] as a quasiprobability decomposition of a non-local unitary channel into local ones.

In this work, we first define a quantity that we call channel robustness of non-locality in analog to the robustness of magic introduced in Ref. [15], which quantifies the minimal possible overhead that can be achieved for quasiprobabilistic simulation of a non-local channel by local channels. While this quantity is difficult to calculate in general, we show an analytic upper bound for general two-qubit unitary channels by constructing an explicit decomposition, generalizing the technique developed in Ref. [11]. Our previous work [11] has only considered decomposition of non-local gates expressed in the form of $e^{i\theta A_1 \otimes A_2}$ for Hermitian operators satisfying $A_1^2 = I$ and $A_2^2 = I$. In contrast, the decomposition developed in this work performs the cut of a general two-qubit gate in a single-step, leading to a substantially reduced overhead. Besides the reduced cost, the derivation of the decomposition is delivered more constructively than before which we believe is informative for further optimizations of this approach. While lower bounds of the defined robustness is also of theoretical interest that can characterize quantumness of a non-local channel, in this work, we focus on upper bounds obtained by explicit decompositions which enable us to actually simulate a nonlocal channel by local channels. This work develops a theoretical framework for a resource reduction suitable for first-generation quantum devices.

2 Decomposition of non-local channels into local channels

2.1 Notation

We use the notation $|\rho\rangle\rangle$ to express a density matrix ρ to stress that ρ can also be seen as a vector. Bold-font symbols are to express a quantum channel corresponding to a gate-like operation represented by a normal font. For example, a unitary channel \mathbf{U} acts on a state $|\rho\rangle\rangle$ as $\mathbf{U}|\rho\rangle\rangle = |U\rho U^\dagger\rangle\rangle$ where U is a unitary matrix. Inner product between two operators $|A\rangle\rangle$ and $|B\rangle\rangle$ is defined as $\langle\langle A|B\rangle\rangle = \text{Tr}(A^\dagger B)$.

2.2 Channel robustness of non-locality

In standard experimental platforms including superconducting qubits and ion traps, it is often thought that the arbitrary single-qubit rotation characterized by an axis $n = (n_1, n_2, n_3)$ and an angle θ , $R(n, \theta) = \exp[-i\theta(\sum_\alpha n_\alpha \sigma_\alpha)]$, and the single-qubit projective measurements along any axis are somewhat easier operations than two-qubit entangling operations. Experimentally, the projective measurement is realized by rotating the axis by $R(n, \theta)$ and performing the projective measurement along z -axis. The quantum channel $\mathbf{M}(n)$ corresponding to the projective measurement is a probabilistic map; when applied to a state $|\rho\rangle\rangle$, it returns a state $\mathbf{\Pi}(\pm n)|\rho\rangle\rangle/p_\pm$ with some probability p_\pm , where $\mathbf{\Pi}(\pm n)$ is a projector to an eigenstate of $\pm \sum_\alpha n_\alpha \sigma_\alpha$ with eigenvalue ± 1 .

To implement $\mathbf{\Pi}(n)$ itself, we can define a probabilistic map $\tilde{\mathbf{\Pi}}(n)$ that takes a state $|\rho\rangle\rangle$ to $\mathbf{\Pi}(n)|\rho\rangle\rangle/p_+$ with probability p_+ and to $|0\rangle\rangle$ with probability p_- where $|0\rangle\rangle$ corresponds to the zero matrix. The map to $|0\rangle\rangle$ means simply to ignore the case when the measurement resulted in -1 . However, just discarding the -1 case is inefficient, especially when we also want to perform $\mathbf{\Pi}(-n)$. To resolve this issue, we define a probabilistic map $\tilde{\mathbf{\Pi}}(n, c_+, c_-)$ that takes a state ρ to $c_\pm \mathbf{\Pi}(\pm n)|\rho\rangle\rangle/p_\pm$ with probability p_\pm , where $c_\pm \in \{0\} \cup \{e^{i\phi} | \phi \in [0, 2\pi]\}$. Let us define the expected value of a random vector $|\sigma\rangle\rangle$ which becomes $|\sigma_i\rangle\rangle$ with a probability p_i as $\mathbb{E}[|\sigma\rangle\rangle] := \sum_i p_i |\sigma_i\rangle\rangle$. Observe that the following holds for any state ρ ,

$$\mathbb{E}[\tilde{\mathbf{\Pi}}(n, c_+, c_-)|\rho\rangle\rangle] = c_+ \mathbf{\Pi}(n)|\rho\rangle\rangle + c_- \mathbf{\Pi}(-n)|\rho\rangle\rangle. \quad (1)$$

We write $\mathbb{E}[\tilde{\mathbf{\Pi}}(n, c_+, c_-)] = c_+ \mathbf{\Pi}(n) + c_- \mathbf{\Pi}(-n)$ in this sense and henceforth use the notation like this. $\tilde{\mathbf{\Pi}}(n, c_+, c_-)$ includes the both of the cases which we mentioned earlier; if we want to apply only $\mathbf{\Pi}(n)$ we can set $c_- = 0$, and we can also apply both of $\mathbf{\Pi}(\pm n)$ simultaneously with different coefficients. The reason we restrict $|c_\pm| = 1$ is to assure $|\text{Tr}[\tilde{\mathbf{\Pi}}(n, c_+, c_-)\rho]| \leq 1$

for any state ρ and any realization of $\tilde{\Pi}(n, c_+, c_-)$, thus preventing the decomposition overhead to occur at this stage.

With the above consideration, available local operations in practice, which we denote as \mathbf{L}_i , are the ones that can be written as an arbitrary product of $\mathbf{R}(n, \theta)$ and $\tilde{\Pi}(n)$ and their tensor products. We denote a set of such possible \mathbf{L}_i by \mathcal{L} . The most general form of decomposition that we aim to build for a given non-local quantum channel Φ is,

$$\Phi = \sum_i c_i \mathbf{L}_i, \quad (2)$$

where $\mathbf{L}_i \in \mathcal{L}$.

Given a decomposition above, Φ can be ‘‘simulated’’ in a Monte-Carlo manner by sampling \mathbf{L}_i with probability proportional to $|c_i|$. More concretely, let us define a probabilistic map $\hat{\Phi}$ such that it becomes $\frac{c_i}{|c_i|} \mathbf{L}_i$ with probability $p_i = |c_i|/W(\Phi)$ where $W(\Phi) = \sum_i |c_i|$. Then,

$$\begin{aligned} \mathbb{E}[W(\Phi)\hat{\Phi}] &= W(\Phi) \times \sum_i \frac{|c_i|}{W(\Phi)} \frac{c_i}{|c_i|} \mathbf{L}_i \\ &= \Phi, \end{aligned} \quad (3)$$

which shows that $W(\Phi)\hat{\Phi}$ becomes equal to Φ when executed for many times. This algorithm involves only local operations with classical communication (LOCC). However, note that the above protocol is not a simple probabilistic mixture of LOCC as it multiplies the complex coefficient $c_i/|c_i|$ to each channel \mathbf{L}_i .

Let us now consider the overhead associated with the decomposition. In many cases, the output from a quantum system that is evolved with a channel Φ is an expectation value of an observable O , which can be written as $\langle\langle O|\Phi|\rho\rangle\rangle$. $\langle\langle O|\Phi|\rho\rangle\rangle$ is usually estimated by sampling eigenvalues of O from the final state $\Phi|\rho\rangle$. Let the sampled S eigenvalues be $\{o_s\}_{s=1}^S$. Normally, we construct an estimator $\widehat{\langle O \rangle}$ as $\widehat{\langle O \rangle} = \frac{1}{S} \sum_s o_s$. Let us assume that absolute value of eigenvalues of O is bounded by o_{\max} and thus $|o_s| \leq o_{\max}$. Then, by Hoeffding’s bound, we can assure that $|\widehat{\langle O \rangle} - \langle\langle O|\Phi|\rho\rangle\rangle| \leq \epsilon$ with probability at least $1 - \delta$ if we take $S = 2(o_{\max}/\epsilon)^2 \ln[1/(2\delta)]$.

The number of samples required to achieve the same accuracy increases if one tries to simulate Φ with $\hat{\Phi}$. Since $\mathbb{E}[W(\Phi)\hat{\Phi}] = \Phi$, $\mathbb{E}[W(\Phi)\langle\langle O|\hat{\Phi}|\rho\rangle\rangle] = \langle\langle O|\Phi|\rho\rangle\rangle$. We can construct an estimator $\widehat{\langle O \rangle}$ by $\widehat{\langle O \rangle} = \frac{1}{S} \sum_s W(\Phi)o'_s$ where o'_s is a sample drawn from $\hat{\Phi}|\rho\rangle$ with a single realization of $\hat{\Phi}$. The application of $\hat{\Phi}$ introduced in the last section involves many stochastic processes; it means to stochastically apply \mathbf{L}_i with probability p_i , and \mathbf{L}_i itself is a stochastic map involving $\tilde{\Pi}(n, c_+, c_-)$. However, in the end, any realization

of $\hat{\Phi}$ becomes a single-qubit operation that preserves the magnitude of the trace of ρ or maps the state to $|0\rangle$. Therefore, it is guaranteed that the absolute value of a sample o'_s obtained by measuring O of $\hat{\Phi}|\rho\rangle$ is also bounded by o_{\max} . Again by Hoeffding’s bound, $|\widehat{\langle O \rangle} - \langle\langle O|\Phi|\rho\rangle\rangle| \leq \epsilon$ with probability at least $1 - \delta$ if we take $S = 2(W(\Phi)o_{\max}/\epsilon)^2 \ln[1/(2\delta)]$. We can see that $W(\Phi)^2$ amounts to the overhead of the decomposition.

The above discussion leads us to define the following quantity which we call the channel robustness of non-locality,

$$\widetilde{W}(\Phi) = \min_{\{c_i|\Phi = \sum_i c_i \mathbf{L}_i, \mathbf{L}_i \in \mathcal{L}\}} \sum_i |c_i|. \quad (4)$$

$\widetilde{W}(\Phi)$ quantifies the minimum amount of cost when we perform the simulation of a non-local channel Φ by probabilistic application of the local, experimentally feasible operations. $\widetilde{W}(\Phi)$ is submultiplicative, i.e., $\widetilde{W}(\Phi_2\Phi_1) \leq \widetilde{W}(\Phi_2)\widetilde{W}(\Phi_1)$, which is proved in Appendix. This allows us to upper-bound the overhead caused by the decomposition of a chain of quantum channels, $\Phi_N \cdots \Phi_2\Phi_1$ by $\prod_{n=1}^N \widetilde{W}(\Phi_n)$.

Note that if we change the available set of operations to some other ones from \mathcal{L} , Eq. (4) quantifies the overhead of the decomposition in that case. For example, the overhead of the decomposition of the identity gate presented in Ref. [10] can be quantified by setting the available decomposition to be measure-and-prepare channels. Another example is the decomposition of non-Clifford circuits into stabilizer-preserving channels considered in Refs. [17, 18]. The cost for a family of the error mitigation technique called probabilistic error cancellation [12, 13] is also in relation to this quantity; it is quantified by substituting the target channel Φ with an inverse of a noise channel.

As \mathcal{L} consists of operations with continuous parameters, we can also define $\widetilde{W}(\Phi)$ using an integral instead of a discrete sum. Formally, we can write,

$$\widetilde{W}(\Phi) = \min_{\{c|\Phi = \int c(\lambda)\mathbf{L}(\lambda)d\lambda, \mathbf{L}(\lambda) \in \mathcal{L}\}} \int |c(\lambda)|d\lambda, \quad (5)$$

where λ denotes some continuous parameters that specifies an element in \mathcal{L} .

The calculation of $\widetilde{W}(\Phi)$ for a general channel Φ is challenging as it involves a complex minimization procedure. Nevertheless, in the next section, we give an upper bound of $\widetilde{W}(\Phi)$ for a general two-qubit unitary channel Φ by explicitly constructing a decomposition using a complete but not overcomplete basis in \mathcal{L} .

2.3 Upper bound for two-qubit unitary channel

It is well-known [19, 20] that the non-local part of two-qubit gates can always be written as,

$$U = \exp \left[i \left(\sum_{\alpha=1}^3 \theta_{\alpha} \sigma_{\alpha} \otimes \sigma_{\alpha} \right) \right] \\ = \sum_{\alpha=0}^3 u_{\alpha} \sigma_{\alpha} \otimes \sigma_{\alpha}, \quad (6)$$

where σ_0 is the 2×2 identity operator, and σ_1, σ_2 and σ_3 are Pauli x, y and z operators, respectively. θ_{α} is a real parameter, and u_{α} is a coefficient that is determined from $\{\theta_{\alpha}\}$. It leads to the following expression of U ,

$$U|\rho\rangle\rangle = \sum_{\alpha, \alpha'} u_{\alpha} u_{\alpha'}^* |(\sigma_{\alpha} \otimes \sigma_{\alpha})\rho(\sigma_{\alpha'} \otimes \sigma_{\alpha'})\rangle\rangle. \quad (7)$$

Note that $\sum_{\alpha} |u_{\alpha}|^2 = 1$ follows from the unitarity.

First, we expand the general two-qubit unitary defined in Eq. (7) using $|\sigma_{\beta}\rangle\rangle$ as a single-qubit basis vector as follows:

$$\langle\langle \sigma_{\beta'} \otimes \sigma_{\gamma'} | U | \sigma_{\beta} \otimes \sigma_{\gamma} \rangle\rangle \\ = \sum_{\alpha, \alpha'} u_{\alpha} u_{\alpha'}^* \text{Tr} [\sigma_{\beta'} \sigma_{\alpha} \sigma_{\beta} \sigma_{\alpha'}] \text{Tr} [\sigma_{\gamma'} \sigma_{\alpha} \sigma_{\gamma} \sigma_{\alpha'}]. \quad (8)$$

From this expression, it is clear that if we can construct a single-qubit channel $U_{\alpha\alpha'}$ such that $U_{\alpha\alpha'}\rho = \sigma_{\alpha}\rho\sigma_{\alpha'}$ for any ρ , we can write the above as,

$$\langle\langle \sigma_{\beta'} \otimes \sigma_{\gamma'} | U | \sigma_{\beta} \otimes \sigma_{\gamma} \rangle\rangle \\ = \sum_{\alpha, \alpha'} u_{\alpha} u_{\alpha'}^* \langle\langle \sigma_{\beta'} | U_{\alpha\alpha'} | \sigma_{\beta} \rangle\rangle \langle\langle \sigma_{\gamma'} | U_{\alpha\alpha'} | \sigma_{\gamma} \rangle\rangle \\ = \sum_{\alpha, \alpha'} u_{\alpha} u_{\alpha'}^* \langle\langle \sigma_{\beta'} \otimes \sigma_{\gamma'} | U_{\alpha\alpha'}^{\otimes 2} | \sigma_{\beta} \otimes \sigma_{\gamma} \rangle\rangle. \quad (9)$$

Therefore, we conclude $U = \sum_{\alpha, \alpha'} u_{\alpha} u_{\alpha'}^* U_{\alpha\alpha'}^{\otimes 2}$.

Now, we construct $U_{\alpha\alpha'}$ with available single-qubit operations. Observe that,

$$\sigma_{\alpha}\rho\sigma_{\alpha'} = \frac{1}{2}(\sigma_{\alpha}\rho\sigma_{\alpha'} + \sigma_{\alpha'}\rho\sigma_{\alpha}) + \frac{1}{2}(\sigma_{\alpha}\rho\sigma_{\alpha'} - \sigma_{\alpha'}\rho\sigma_{\alpha}). \quad (10)$$

Let us define the following operators $A_{\alpha\alpha', \pm}$ and $B_{\alpha\alpha', \pm}$ which can be implemented through single-qubit operations:

$$A_{\alpha\alpha', \pm} = \frac{1}{2}(\sigma_{\alpha} \pm \sigma_{\alpha'}), \quad (11)$$

$$B_{\alpha\alpha', \pm} = \frac{1}{2}(\sigma_{\alpha} \pm i\sigma_{\alpha'}). \quad (12)$$

The corresponding channels $A_{\alpha\alpha', \pm}$ and $B_{\alpha\alpha', \pm}$ act on a single-qubit density matrix ρ like $A_{\alpha\alpha', \pm}\rho A_{\alpha\alpha', \pm}^{\dagger}$.

Building on $A_{\alpha\alpha', \pm}$ and $B_{\alpha\alpha', \pm}$, we further define the following channels:

$$A_{\alpha\alpha'} = A_{\alpha\alpha', +} - A_{\alpha\alpha', -}, \quad (13)$$

$$B_{\alpha\alpha'} = B_{\alpha\alpha', +} - B_{\alpha\alpha', -}. \quad (14)$$

With simple algebra, we can see that,

$$A_{\alpha\alpha'}\rho = \frac{1}{2}(\sigma_{\alpha}\rho\sigma_{\alpha'} + \sigma_{\alpha'}\rho\sigma_{\alpha}), \quad (15)$$

$$B_{\alpha\alpha'}\rho = \frac{1}{2i}(\sigma_{\alpha}\rho\sigma_{\alpha'} - \sigma_{\alpha'}\rho\sigma_{\alpha}). \quad (16)$$

Therefore, $U_{\alpha\alpha'}$ can be written as,

$$U_{\alpha\alpha'} = A_{\alpha\alpha'} + iB_{\alpha\alpha'}. \quad (17)$$

The above decomposition of $U_{\alpha\alpha'}$ leads us to the following decomposition of U :

$$U = \sum_{\alpha\alpha'} u_{\alpha} u_{\alpha'}^* (A_{\alpha\alpha'} + iB_{\alpha\alpha'})^{\otimes 2}. \quad (18)$$

Note that there are symmetries $A_{\alpha\alpha'} = A_{\alpha'\alpha}$ and $B_{\alpha\alpha'} = -B_{\alpha'\alpha}$. Using them, we rewrite the expression for later convenience,

$$U = \sum_{\alpha} |u_{\alpha}|^2 \sigma_{\alpha}^{\otimes 2} \\ + \sum_{\alpha < \alpha'} (u_{\alpha} u_{\alpha'}^* + u_{\alpha'} u_{\alpha}^*) (A_{\alpha\alpha'}^{\otimes 2} - B_{\alpha\alpha'}^{\otimes 2}) \\ + \sum_{\alpha < \alpha'} i(u_{\alpha} u_{\alpha'}^* - u_{\alpha'} u_{\alpha}^*) (A_{\alpha\alpha'} \otimes B_{\alpha\alpha'} + B_{\alpha\alpha'} \otimes A_{\alpha\alpha'}). \quad (19)$$

To calculate upper bound for $\widetilde{W}(U)$, we need to formulate Eq. (19) to fit in the form of Eq. (2). σ_{α} , which constitutes the first term of the decomposition, is trivially in \mathcal{L} . Let us now consider $A_{\alpha\alpha'}$. We note that from the symmetry it suffices to consider the case where $\alpha < \alpha'$. When $\alpha = 0$, $A_{\alpha\alpha', \pm}$ becomes a projector $\Pi(\pm n)$ where $n_{\alpha''} = \delta_{\alpha'\alpha''}$. Therefore, $A_{\alpha\alpha'}$ takes the form of $\widetilde{\Pi}(n, 1, -1)$, which means $A_{0\alpha'} \in \mathcal{L}$. For $\alpha \neq 0$, $A_{\alpha\alpha', \pm}$ is proportional to a single-qubit rotation that swaps the α -axis and α' -axis. More concretely, $2A_{\alpha\alpha', \pm} \in \mathcal{L}$ for $\alpha \neq 0$ and $\alpha < \alpha'$. As for $B_{\alpha\alpha'}$, when $\alpha = 0$, $B_{\alpha\alpha', \pm}$ becomes proportional to a single-qubit rotation around α' -axis. Likewise to the previous case, $2B_{\alpha\alpha', \pm} \in \mathcal{L}$. For $\alpha \neq 0$, $B_{\alpha\alpha', \pm}$ can be implemented by a projector followed by a flip; for example, $\frac{1}{2}(\sigma_1 + i\sigma_2) = \frac{1}{2}\sigma_1(\sigma_0 - \sigma_3)$. With this observation, we can see that the channel $B_{\alpha\alpha'}$ in this case can be written as a product of $\widetilde{\Pi}$ and σ_{α} which makes $B_{\alpha\alpha'} \in \mathcal{L}$ for $\alpha \neq 0$ and $\alpha < \alpha'$.

Combining the above properties, we can calculate $W(U) = \sum_i |c_i|$ for the decomposition given in Eq. (19)

as,

$$W(\mathbf{U}) = 1 + \sum_{\alpha \neq \alpha'} (|u_\alpha u_{\alpha'}^* + u_{\alpha'} u_\alpha^*| + |u_\alpha u_{\alpha'}^* - u_{\alpha'} u_\alpha^*|), \quad (20)$$

which gives an upper bound of $\widetilde{W}(\mathbf{U})$. We note that the operations used in the proposed decomposition, namely σ_α ($\alpha \in \{0, 1, 2, 3\}$), $\mathbf{A}_{\alpha\alpha'}$ and $\mathbf{B}_{\alpha\alpha'}$ with $\alpha < \alpha'$ are 16 linearly independent single-qubit channel and thus form a complete basis in the space of single-qubit superoperators. This means $W(\mathbf{U})$ is uniquely determined as long as the same basis set is used.

2.4 Behaviour of $W(\mathbf{U})$

Here, we numerically investigate the behavior of $W(\mathbf{U})$ defined in Eq. (20), restricting the domain of $\{\theta_\alpha\}$ in which each point is not locally equivalent, meaning that a two-qubit unitary represented by a point $(\theta_1, \theta_2, \theta_3)$ cannot be translated to another point in the domain by transforming it with single-qubit unitaries, according to Ref. [20]. In Fig. 1, we depict such a domain of $\{\theta_\alpha\}$ ¹. Note that there are exceptional local-equivalence in the domain; every point $A_1 A_2 A_3$ and $O A_2 A_3$ is locally equivalent to $A'_1 A'_2 A'_3$ and $O A'_2 A'_3$, respectively.

Since $W(\mathbf{U})$ is symmetric to the reflection of θ_x , we only investigate the tetrahedron $O A_1 A_2 A_3$. In Fig. 2, we show the behavior of $W(\mathbf{U})$ on the surfaces and edges of the domain. We numerically found that $W(\mathbf{U})$ is maximized at $(\theta_1, \theta_2, \theta_3) \approx (\pi/4, 0.202\pi, 0.136\pi)$ which lies on the surface $A_1 A_2 A_3$ with its value being approximately 8.87. The behavior of $W(\mathbf{U})$ seems to be unrelated to other measures such as entangling power of \mathbf{U} [19, 21]; for example, while the point A_1 corresponds to controlled- σ_α gates which can produce the maximal amount of entanglement and has $W(\mathbf{U}) = 3$, A_3 which corresponds to the swap gate has $W(\mathbf{U}) = 7$. Although we believe the decomposition given in this work is close to optimal, this counter-intuitive result might be caused by the non-optimality.

3 Discussion

3.1 Comparison with gate-based decomposition approach

If we can measure $\langle \psi_1 | \psi_2 \rangle$ for some fixed state $|\psi_1\rangle$ and $|\psi_2\rangle$, we can directly utilize the fact that a two-qubit gate is decomposed as $\sum_{\alpha \in \{I, x, y, z\}} u_\alpha \sigma_\alpha \otimes \sigma_\alpha$. As we discuss later, this measurement can be demanding

¹It slightly differs from Ref. [20]. We shift half of the tetrahedron presented in Fig. 2 of Ref. [20] corresponding to the region $\theta_x \geq \pi/4$ to $\theta_x \leq \pi/0$ using the periodicity of θ_x .

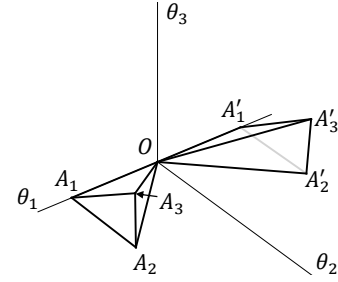


Figure 1: Domain of $(\theta_1, \theta_2, \theta_3)$ in which a two-qubit unitary represented by each point is not locally equivalent to each other. In the figure, $O = (0, 0, 0)$, $A_1 = (\pi/4, 0, 0)$, $A_2 = (\pi/4, \pi/4, 0)$, $A_3 = (\pi/4, \pi/4, \pi/4)$, $A'_1 = (-\pi/4, 0, 0)$, $A'_2 = (-\pi/4, \pi/4, 0)$ and $A'_3 = (-\pi/4, \pi/4, \pi/4)$.

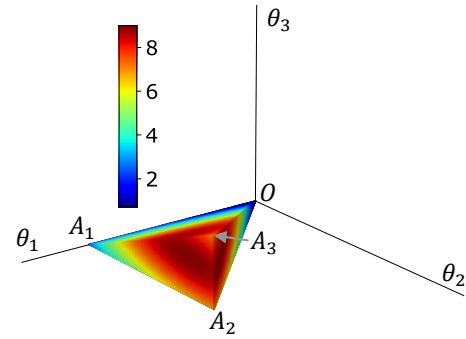


Figure 2: Behaviour of $W(\mathbf{U})$ on the surface of the tetrahedron $O A_1 A_2 A_3$.

for early days quantum computers. Let V be a sequence of gates consisting of alternating layers of single-qubit and two-qubit gates. Note that any quantum circuit can be written in this form. V can be written as $V = D_L S_L \cdots D_2 S_2 D_1 S_1$ where D_i 's and S_i 's are two-qubit and single-qubit gates, respectively. We assume $D_i = \sum_{\alpha} d_{\alpha_i} \sigma_{\alpha_i}^{(a_i)} \otimes \sigma_{\alpha_i}^{(b_i)}$ where $\sigma_{\alpha}^{(a)}$ is a Pauli matrix acting on the a -th qubit. Now, focusing on the i -th two-qubit gate, we can express an expectation value of an observable O at the end of the circuit as,

$$\langle 0|V^{\dagger}OV|0\rangle = \sum_{\alpha_i} d_{\alpha_i}^* d_{\alpha_i} \langle 0|V_{i,\alpha_i}^{\dagger}OV_{i,\alpha_i}|0\rangle, \quad (21)$$

where,

$$V_{i,\alpha_i} = D_L S_L \cdots \sigma_{\alpha_i}^{(a_i)} \otimes \sigma_{\alpha_i}^{(b_i)} \cdots D_2 S_2 D_1 S_1. \quad (22)$$

This decomposition also allows us to perform a ‘‘virtual’’ two-qubit gate on a quantum circuit in the sense that, in V_{i,α_i} , the i -th two-qubit gate in V is replaced by $\sigma_{\alpha_i}^{(a_i)} \otimes \sigma_{\alpha_i}^{(b_i)}$ which is a tensor product of local operations. We can do this by the following algorithm. Let us assume that O is written as $O = \sum_i c_i P_i$, where P_i is a tensor product of Pauli operators. With this assumption, we can evaluate $\langle 0|V_{i,\alpha_i}^{\dagger}OV_{i,\alpha_i}|0\rangle$ by $\sum_k c_k \langle 0|V_{i,\alpha_i}^{\dagger}P_k V_{i,\alpha_i}|0\rangle$. More concretely, we define $|\psi_{i,\alpha_i}\rangle = V_{i,\alpha_i}|0\rangle$ and $|\psi_{k,i,\alpha_i}\rangle = P_k V_{i,\alpha_i}|\psi_{i,\alpha_i}\rangle$ and then measure $\langle \psi_{i,\alpha_i}|\psi_{k,i,\alpha_i}\rangle$ which is possible by the assumption. If we are to perform the sum of Eq. (21) in a Monte-Carlo manner, we can sample α_i' and α_i with a probability proportional to $|d_{\alpha_i'}^* d_{\alpha_i}|$. This leads us to define $G(D_i) := \sum_{\alpha_i',\alpha_i} |d_{\alpha_i'}^* d_{\alpha_i}|$ which quantifies the overhead of the decomposition, that is, we need $G(D_i)^2$ times more samples to reach a desired error compared to the decomposition-free case.

It is trivial that $G(D_i)$ is always smaller than $W(D_i)$. Therefore, if we can measure $\langle \psi_{i,\alpha_i'}|\psi_{k,i,\alpha_i}\rangle$, it is always better to use this approach. For example, in a classical simulation we can easily calculate $\langle \psi_{i,\alpha_i'}|\psi_{k,i,\alpha_i}\rangle$. However, it is not the case for a quantum computer, in particular for a NISQ device. Measurement of the overlap $\langle \psi_{i,\alpha_i'}|\psi_{k,i,\alpha_i}\rangle$, including its phase, is a demanding task. One way of performing this task is to use a controlled- V_{i,α_i} as mentioned in e.g. Refs. [16, 22], which is unlikely to be implemented on a NISQ device due to its complexity. The original motivation of this work and our previous works [11, 23] has been to avoid such complex operations. Note that the famous swap test [24, 25] cannot be applied to this task since it can only evaluate $|\langle \psi_{i,\alpha_i'}|\psi_{k,i,\alpha_i}\rangle|^2$. Investigations on the relation between $\widetilde{W}(D_i)$ and $G(D_i)$ are left for the future work.

3.2 Comparison with the previous work

In the previous work [11], we have proposed the decomposition for an gate in the form $e^{i\theta A_1 \otimes A_2}$ for Hermitian operators A_1 and A_2 satisfying $A_1^2 = I$ and $A_2^2 = I$. It is a special case of this work, which is recovered by setting $u_0 = \cos \theta$ and $u_{\alpha} = i \sin \theta$ for one chosen $\alpha \in \{1, 2, 3\}$. Therefore, the cost overhead of this special case is determined by $1 + 2|u_0 u_{\alpha}^* - u_{\alpha} u_0^*|$, which takes maximum at $\theta = \pi/4$. If we are to decompose a general two-qubit gate in the form of $\exp\left[i\left(\sum_{\alpha=1}^3 \theta_{\alpha} \sigma_{\alpha} \otimes \sigma_{\alpha}\right)\right]$ using this technique, we decompose each of $\exp[i\theta_{\alpha} \sigma_{\alpha} \otimes \sigma_{\alpha}]$. Then, the overhead is quantified by the product of $1 + 2|u_I u_{\alpha}^* - u_{\alpha} u_I|$, which reaches its maximum $3^3 = 27$ at $\theta_{\alpha} = \pi/4$ for all α . On the other hand, $W_{\mathcal{U}}$ defined in Eq. (20), which quantifies the overhead required by the present approach, becomes 7, showing substantial improvement.

While we believe that the decomposition given in this work is close to optimal, there can be better decompositions with smaller $W_{\mathcal{U}}$. The search for optimal decomposition will require some form of numerical search. In the context of classical simulation of near Clifford circuits, Ref. [16] has performed such a search. However, the optimization of the decomposition considered in this work will be more complicated than the aforementioned work, since the number of available operations is infinitely many as can be seen from Eq. (5). We believe the decomposition proposed in this work can be a good starting point of the optimization if it is not optimal and leave it as future work.

4 conclusion

We have introduced a quantity called channel robustness of non-locality which quantifies the minimal amount of overhead required for decomposing non-local channels into local ones with a quasiprobability-based method. While the calculation of the quantity for general non-local channels is difficult due to the need for a complicated optimization, we have successfully established an upper bound for a general two-qubit unitary channel. The upper bound is obtained by constructively deriving an explicit decomposition. Its overhead is substantially lowered compared to the previous work [11]. While we believe the present decomposition is close to optimal, there might be a better decomposition of a general two-qubit channel than the one presented in this work, which we leave as possible future work. This formalism of decomposing an experimentally challenging channel into a linear combination of experimentally-easy channels allows us to readily perform the decomposition using a quantum device.

Acknowledgments

KM is supported by JST PRESTO JPMJPR2019 and KAKENHI No. 20K22330. KF is supported by KAKENHI No.16H02211, JST ERATO JPMJER1601, and JST CREST JPMJCR1673. This work is supported by MEXT Quantum Leap Flagship Program (MEXT Q-LEAP) Grant Number JPMXS0118067394. Program code for generating Fig. 2 is available at <https://github.com/kosukemtr/nonlocal-local-decomposition>.

A Submultiplicability of $\widetilde{W}(\Phi)$

Lemma 1 *Let Φ_1 and Φ_2 be any quantum channels and $\Phi_{21} = \Phi_2\Phi_1$. Then, $\widetilde{W}(\Phi_{21}) \leq \widetilde{W}(\Phi_2)\widetilde{W}(\Phi_1)$.*

proof— Let

$$\Phi_\mu = \sum_i c_{\mu i} \mathbf{L}_{\mu i} \quad (23)$$

and $\sum_i |c_{\mu i}| = \widetilde{W}(\Phi_\mu)$. Then, Φ_{21} can be decomposed as,

$$\Phi_{21} = \sum_{ij} c_{2i} c_{1j} \mathbf{L}_{2i} \mathbf{L}_{1j}. \quad (24)$$

Because $\mathbf{L}_{2i} \mathbf{L}_{1j} \in \mathcal{L}$, the above gives a decomposition of Φ_{21} in the form of Eq. 2. Therefore,

$$\begin{aligned} \widetilde{W}(\Phi_{21}) &\leq \sum_{ij} |c_{2i} c_{1j}| \\ &= \sum_i |c_{2i}| \sum_j |c_{1j}| \\ &= \widetilde{W}(\Phi_2) \widetilde{W}(\Phi_1). \end{aligned} \quad (25)$$

□

References

- [1] F. Arute, K. Arya, R. Babbush, D. Bacon, J. C. Bardin, R. Barends, R. Biswas, S. Boixo, F. G. S. L. Brandao, D. A. Buell, B. Burkett, Y. Chen, Z. Chen, B. Chiaro, R. Collins, W. Courtney, A. Dunsworth, E. Farhi, B. Foxen, A. Fowler, C. Gidney, M. Giustina, R. Graff, K. Guerin, S. Habegger, M. P. Harrigan, M. J. Hartmann, A. Ho, M. Hoffmann, T. Huang, T. S. Humble, S. V. Isakov, E. Jeffrey, Z. Jiang, D. Kafri, K. Kechedzhi, J. Kelly, P. V. Klimov, S. Knysh, A. Korotkov, F. Kostritsa, D. Landhuis, M. Lindmark, E. Lucero, D. Lyakh, S. Mandrà, J. R. McClean, M. McEwen, A. Megrant, X. Mi, K. Michielsen, M. Mohseni, J. Mutus, O. Naaman, M. Neeley, C. Neill, M. Y. Niu, E. Ostby, A. Petukhov, J. C. Platt, C. Quintana, E. G. Rieffel, P. Roushan, N. C. Rubin, D. Sank, K. J. Satzinger, V. Smelyanskiy, K. J. Sung, M. D. Trevithick, A. Vainsencher, B. Villalonga, T. White, Z. J. Yao, P. Yeh, A. Zalcman, H. Neven, and J. M. Martinis, *Nature* **574**, 505 (2019).
- [2] J. Preskill, *Quantum* **2**, 79 (2018).
- [3] A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, X.-Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O’Brien, *Nature Communications* **5**, 4213 (2014).
- [4] S. McArdle, S. Endo, A. Aspuru-Guzik, S. C. Benjamin, and X. Yuan, *Rev. Mod. Phys.* **92**, 015003 (2020).
- [5] E. Farhi, J. Goldstone, and S. Gutmann, “A quantum approximate optimization algorithm,” (2014), [arXiv:1411.4028 \[quant-ph\]](https://arxiv.org/abs/1411.4028).
- [6] K. Mitarai, M. Negoro, M. Kitagawa, and K. Fujii, *Phys. Rev. A* **98**, 032309 (2018).
- [7] E. Farhi and H. Neven, “Classification with quantum neural networks on near term processors,” (2018), [arXiv:1802.06002 \[quant-ph\]](https://arxiv.org/abs/1802.06002).
- [8] C. Bravo-Prieto, R. LaRose, M. Cerezo, Y. Subasi, L. Cincio, and P. J. Coles, “Variational quantum linear solver,” (2019), [arXiv:1909.05820 \[quant-ph\]](https://arxiv.org/abs/1909.05820).
- [9] R. LaRose, A. Tikku, E. O’Neel-Judy, L. Cincio, and P. J. Coles, *npj Quantum Information* **5**, 8 (2019).
- [10] T. Peng, A. W. Harrow, M. Ozols, and X. Wu, *Phys. Rev. Lett.* **125**, 150504 (2020).
- [11] K. Mitarai and K. Fujii, *New Journal of Physics*, accepted (2020).
- [12] K. Temme, S. Bravyi, and J. M. Gambetta, *Phys. Rev. Lett.* **119**, 180509 (2017).
- [13] S. Endo, S. C. Benjamin, and Y. Li, *Phys. Rev. X* **8**, 031027 (2018).
- [14] H. Pashayan, J. J. Wallman, and S. D. Bartlett, *Phys. Rev. Lett.* **115**, 070501 (2015).
- [15] M. Howard and E. Campbell, *Phys. Rev. Lett.* **118**, 090501 (2017).
- [16] S. Bravyi, G. Smith, and J. A. Smolin, *Phys. Rev. X* **6**, 021043 (2016).
- [17] R. S. Bennink, E. M. Ferragut, T. S. Humble, J. A. Laska, J. J. Nutaro, M. G. Pleszkoch, and R. C. Pooser, *Phys. Rev. A* **95**, 062337 (2017).
- [18] J. R. Seddon and E. T. Campbell, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* **475**, 20190251 (2019).
- [19] B. Kraus and J. I. Cirac, *Phys. Rev. A* **63**, 062309 (2001).
- [20] J. Zhang, J. Vala, S. Sastry, and K. B. Whaley, *Phys. Rev. A* **67**, 042313 (2003).

- [21] F.-Z. Kong, J.-L. Zhao, M. Yang, and Z.-L. Cao, *Phys. Rev. A* **92**, 012127 (2015).
- [22] Y. Ibe, Y. O. Nakagawa, T. Yamamoto, K. Mitarai, Q. Gao, and T. Kobayashi, “Calculating transition amplitudes by variational quantum eigensolvers,” (2020), [arXiv:2002.11724 \[quant-ph\]](https://arxiv.org/abs/2002.11724) .
- [23] K. Mitarai and K. Fujii, *Phys. Rev. Research* **1**, 013006 (2019).
- [24] H. Buhrman, R. Cleve, J. Watrous, and R. de Wolf, *Phys. Rev. Lett.* **87**, 167902 (2001).
- [25] J. C. Garcia-Escartin and P. Chamorro-Posada, *Phys. Rev. A* **87**, 052330 (2013).