

# Realization of Thread Level Parallelism on Quantum Devices

Keren Li,<sup>1,2,\*</sup> Zidong Lin,<sup>3</sup> Zheng An,<sup>4</sup> Guanru Feng,<sup>3</sup> Zipeng Wu,<sup>3,†</sup> Shiyao Hou,<sup>5,‡</sup> and Jingen Xiang<sup>3,§</sup>

<sup>1</sup>College of Physics and Optoelectronic Engineering, Shenzhen University, Shenzhen 518060, China

<sup>2</sup>Quantum Science Center of Guangdong-Hong Kong-Macao Greater Bay Area (Guangdong), Shenzhen 518045, China

<sup>3</sup>Shenzhen SpinQ Technology Co., Ltd., 518043, Shenzhen, China

<sup>4</sup>Department of Physics, The Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong, China

<sup>5</sup>College of Physics and Electronic Engineering, Center for Computational Sciences,

Sichuan Normal University, Chengdu 610068, China

(Dated: November 10, 2025)

Scaling up quantum devices is a central challenge for realizing practical quantum computation. Modular quantum architectures promise scalability, yet experiments to date have relied on either  $\lesssim 10^3$ -qubit monolithic chips or fragile interconnects with high loss. Here, we introduce a classical linkage scheme that merges multiple independent quantum processing units (QPUs) into a single logical device, enabling thread-level parallelism (TLP). Theoretically, we show that quantum routines with product-state inputs and low-rank entangling layers can be re-expressed in an efficient parallelizable form. Experimentally, we validate this architecture on clusters comprising up to sixteen benchtop nuclear magnetic resonance (NMR) quantum nodes. A four-qubit Greenberger-Horne-Zeilinger (GHZ) state is partitioned into parallel two-qubit subcircuits, achieving a fidelity of 93.8% with respect to the ideal state. A non-Hermitian evolution, implemented via a truncated Cauchy integral on Hermitian Hamiltonians, reproduces exact observables with high accuracy. Our results demonstrate that classical links suffice to scale up the logical size of quantum computations and realize general, non-unitary channels on today's hardware, opening an experimentally accessible route toward software-defined, clustered quantum accelerators.

*Introduction.* For half a century, advances in classical computing have depended critically on increasing parallelism. When single-core performance reached physical and architectural limits, research shifted toward extracting parallelism at multiple levels. Modern architectures use three typical forms of parallelism [see Fig. 1(a)]: instruction-level parallelism (ILP) for optimizing execution pipelines, data-level parallelism (DLP) for processing multiple data elements simultaneously, and thread-level parallelism (TLP) for concurrent execution of multiple threads, each handling different data or tasks. Central processing units (CPUs) exploit ILP through pipelining and out-of-order execution, whereas graphics processing units (GPUs), originally designed for rendering, enable DLP and TLP by integrating thousands of lightweight cores onto a single chip [1–3]. This evolution establishes software-hardware co-design, whereby compute-intensive workloads are automatically partitioned and distributed.

Quantum computing now approaches a similar inflection point. Although the accessible state space of an  $n$ -qubit device scales as  $2^n$ , the qubit count of a fabricated chip is fixed; expanding it requires entirely new fabrication procedures. Increasing qubits sharply elevates risk from materials non-uniformity, wiring and cryogenic I/O scaling, packaging challenges, and yield degradation, rendering large monolithic QPUs economically and techno-

logically daunting. An alternative is to interconnect existing QPUs, forming clusters that collectively address larger Hilbert spaces and support broader algorithms without redesigning the quantum hardware. Yet, unlike classical workloads, quantum computation must preserve coherence and entanglement across partitions. Naive partitioning disrupts quantum correlations and is therefore nontrivial. Recent proposals in distributed quantum computing and quantum networking use long-range entanglement to interconnect distant QPUs [4–6]. Despite rapid experimental progress [7–10], maintaining reliable quantum links remains challenging. An alternative is thus to coordinate multiple QPUs using only classical communication, for example, with circuit cutting techniques [11–13]. Notably, a 142-qubit superconducting demonstration linked two 127-qubit chips via classical interconnects [14], and software stacks such as CUDA-Q enable programming across multiple QPUs [15]. However, existing approaches do not fully exploit parallelism between distinct QPUs, and practical, hardware-level architectures with demonstrations remain scarce [see Appendix A].

Here, we propose a framework that brings TLP to QPUs. Conceptually, we develop a modular, parallelizable implementation of completely positive trace-preserving (CPTP) channels together with a decomposition that distributes a channel across multiple QPUs with bounded classical coordination overhead. Architecturally, we realize a scalable cluster that links QPUs using only classical communication and instantiate it with nuclear magnetic resonance (NMR) quantum nodes. Experimentally, we demonstrate two hallmarks of TLP. First, we prepare a four-qubit GHZ state using three-

\* likr@szu.edu.cn

† z2wu@spinq.cn

‡ hshiyao@sicnu.edu.cn

§ jxiang@spinq.cn

qubit NMR nodes, showing that TLP enables emulation of logical systems larger than any individual device. Second, we realize non-Hermitian Hamiltonian simulations, demonstrating the ability to reproduce classes of operations otherwise inaccessible to unitary devices of comparable size. Together, these results indicate that TLP can substantially extend the algorithmic reach of current quantum hardware and offers a concrete route toward scalable, practical quantum computation.

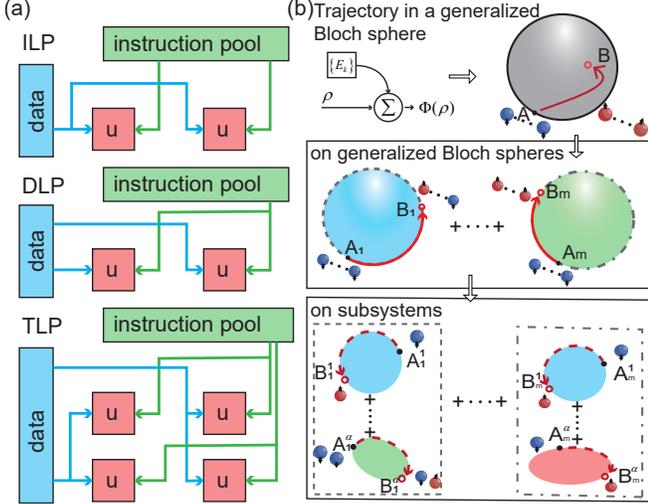


FIG. 1: (a) Instruction-, data-, and thread-level parallelism in a classical processor. (b) Factorized evaluation on modular quantum hardware: Trajectory in a generalized Bloch sphere is decomposed into blocks that factorize across subsystems. Local trajectories (on the generalized Bloch sphere of subsystem) are measured independently, and then classically aggregated to recover the global result.

*Architecture.* In quantum computing, a common task is to compute the expectation value of an observable after applying a quantum channel to an input state. This is specified as

$$\text{Tr}(O \cdot \Phi(\rho)), \quad (1)$$

where  $\Phi$  is a quantum channel,  $\rho$  is the input state, and  $O$  is the observable of interest. Note that  $\rho$  and  $O$  play symmetric roles in Eq. (1), enabling parallelism in analogy to classical computing: ILP, by parallelizing operations within  $\Phi$ ; DLP, by running the same operations on multiple copies of  $\rho$  or measuring different  $O$  in parallel; TLP, by executing multiple quantum operations concurrently, each with different inputs and measurements. Here we focus on TLP and map a quantum process onto a cluster of modular QPUs.

Frequently,  $\rho$  and  $O$  factorize across subsystems, i.e.,  $\rho = \bigotimes_{a=1}^A \rho^a$  and  $O = \bigotimes_{a=1}^A O^a$ . Consider a general form,  $\Phi(\cdot) = \sum_{p=1}^q (\sum_{i=1}^m c_{p,i} U_{p,i})(\cdot)(\sum_{j=1}^m c_{p,j} U_{p,j})^\dagger$ , which subsumes standard Kraus and linear combination of unitaries decompositions, defines a CPTP map under

appropriate coefficient constraints, and thereby naturally exposes ILP. Moreover, assume each  $U_{p,k}$  admits the factorized expansion,

$$U_{p,k} = \sum_{\alpha=1}^{\ell} \bigotimes_{a=1}^A U_{p,k\alpha}^a, \quad (2)$$

the evaluation thus decouples across subsystems. In particular, Eq. (1) can be rewritten as

$$\sum_{p=1}^q \sum_{i,j=1}^m c_{p,i} c_{p,j}^* \sum_{\alpha,\alpha'=1}^{\ell} \prod_{a=1}^A \text{Tr}(U_{p,i\alpha}^a \rho^a U_{p,j\alpha'}^{a\dagger} O^a), \quad (3)$$

reducing the computation to a sum of overlaps, each measurable on a small subsystem [see Fig. 1(b)]. The final result is obtained by classically aggregating all outputs from the  $q \times m^2 \times \ell^2 \times \mathcal{A}$  evaluations [see Appendix B for details]. This decomposition provides the theoretical basis for our architecture and experiments. With varying  $O^a$ , it yields a systematic, scalable route to TLP on modular QPU clusters, where heterogeneous operations and measurements are executed in parallel and orchestrated by a classical controller.

As for requirements, beyond tensor-separable  $\rho$  and  $O$ , the derivation of Eq. (3) assumes a factorization of the global unitary into sums of local operators, as in Eq. (2). A brute-force operator-Schmidt decomposition is possible in principle, but the cost scales as  $O(d^{2n})$  for  $n$  qubits and quickly becomes impractical. Fortunately, layered circuit structures are common in practice. If  $U_{p,k}$  admits a layered-circuit form, the factorization can be obtained efficiently and even optimized by recasting the task as a min-cut problem on the associated circuit graph [16, 17]. Consider the bisection min-cut problem, i.e.,  $\mathcal{A} = 2$ . A layered circuit  $U_{p,k}$  can be written as  $\prod_{t=1}^{m_d} (\sum_{\alpha_t=1}^{\ell_t} \bigotimes_{a=1}^2 g_a^{\alpha_t} U_a^{(t)})$ , where the  $t$ -th entangling layer factorizes as  $\sum_{\alpha_t=1}^{\ell_t} \bigotimes_{a=1}^2 g_{t\alpha_t}^a$ . In fact, some  $g_{t\alpha_t}^a$  may be the identity. In this case,  $m_d$ , originally set as the number of circuit layers, can be reduced. The goal is to bipartition the circuit while minimizing the number of inter-partition entangling gates. Applying a circuit-graph min-cut algorithm [complexity  $O(n^2)$  in the number  $n$  of graph vertices, specified in Appendix C.], one finds a cut of size  $m' \leq m_d$  such that the global expectation in Eq. (3) factorizes into at most  $q \times m^2 \times (\prod_{t=1}^{m'} \ell_t)^2 \times 2$  independent single-subsystem traces, each executable on an individual QPU node. Furthermore, each trace of overlap in Eq. (3) can be measured efficiently using a single-ancilla protocol, requiring only one additional qubit per subsystem; see Appendix C, especially Lemmas C.1 and C.2, for details.

*Hardware.* Fig. 2(a) shows a cluster with loose coupling, where multiple QPUs are integrated. At the core of the system is a classical controller that assigns computational tasks, synchronizes all QPU nodes, generates pulse operations, and coordinates instruction dis-

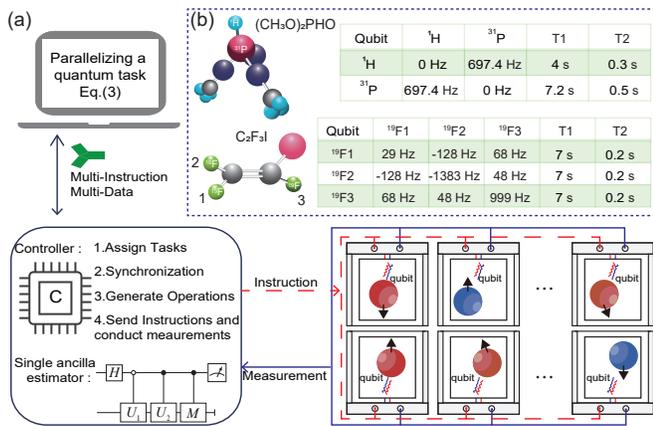


FIG. 2: Hardware overview. (a) Clustered-QPU architecture enabling TLP in quantum computing. (b) Physical sample assembly employed in this work.

patch and measurement collection. Operating in a multi-instruction, multi-data mode, the controller ensures both flexibility and high throughput for realizing TLP on quantum devices. Each QPU node can, in principle, be instantiated on any platform, such as nuclear-spin qubits, photonic qubits, superconducting circuits, or others, provided it supports universal control and high-fidelity state preparation and measurement. This modular design enables seamless integration of heterogeneous quantum processors within a unified computational platform, with all inter-processor signaling carried classically.

During operation, the central scheduler decomposes tasks according to the theoretical framework, dispatches the resulting subtasks to each node, and orchestrates control-pulse timing and data acquisition to maintain precise synchronization. The cluster is fully modular and plug-and-play: nodes can be added, removed, or replaced on the fly, with the controller automatically updating task assignments and preserving system-wide logical consistency. Because all inter-node communication is classical, the architecture is inherently robust and highly scalable.

As a demonstration, we construct clusters comprising 8 or 16 benchtop NMR QPUs [18, 19], where each node hosts 3 or 2 nuclear-spin qubits [Fig. 2(b)]. The physical realization employs different nuclear species,  $^1\text{H}$ ,  $^{19}\text{F}$ , and  $^{31}\text{P}$ , each with distinct Larmor frequencies, chemical shifts, and coherence properties. Key parameters, including  $T_1$  and  $T_2$ , are summarized in Fig. 2(b) and Appendix D, highlighting multi-second coherence that supports reliable thread-level parallel computation. This platform validates our architecture along two axes. First, clusters of small QPUs collectively emulate circuits that exceed the size of any single node. Second, general quantum channels are realized efficiently by parallelizing simple unitary evolutions. We next present two representative experimental tasks.

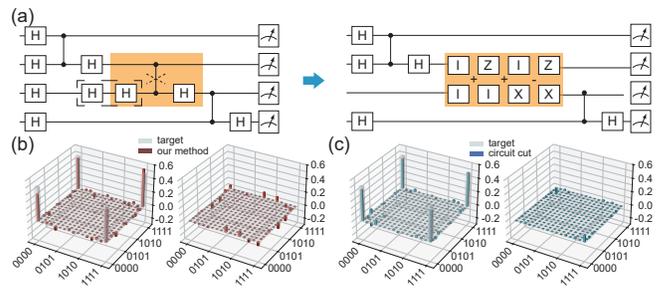


FIG. 3: (a) Experimental quantum circuit for preparing a four-qubit GHZ state. The CZ gate between qubits 2 and 3 has been “cut” and implemented as a linear combination of local operations. (b) and (c) depict real and imaginary parts of the reconstructed density matrices for our method and the circuit-cut state, with the target value is shown with transparent blocks.

*GHZ state preparation.* We first demonstrate the preparation of a four-qubit GHZ state, a canonical benchmark in quantum information [20, 21]. As shown in Fig. 3(a), the CNOT between qubits 2 and 3 is “cut” and replaced by a linear combination of local operations,  $U_{cnot} = \frac{1}{2}(\mathbb{1} \otimes \mathbb{1} + \sigma_z \otimes \mathbb{1} + \mathbb{1} \otimes \sigma_x - \sigma_z \otimes \sigma_x)$ , where  $\sigma_{x,z}$  denote Pauli operators and  $\mathbb{1}$  is the identity. This allows the output state to be expressed as a sum of product substates, and the original 4-qubit GHZ circuit is thus decomposed into a superposition of two-qubit subcircuits, making the protocol compatible with small-scale hardware. Using the single-ancilla estimator, we measure the overlaps and implement parallelization on two kinds of circuits, shown in Appendix E. Additionally, we show compatibility of the hardware with another technique invoking parallelization, quantum circuit cutting [22], which decomposes the four-qubit circuit into a sum of ten two-qubit subcircuits.

Both approaches are implemented on clusters with either eight 3-qubit nodes or sixteen 2-qubit nodes. To obtain full information on the prepared GHZ state, we perform 128 observable evaluations, extracting  $\langle \sigma_x \rangle$ ,  $\langle \sigma_y \rangle$ , and the projectors  $|0\rangle\langle 0|$ ,  $|1\rangle\langle 1|$  on the ancilla qubit across 32 circuits derived from two template circuits. In the circuit-cutting scheme, full state tomography is employed, yielding 160 observable evaluations from 10 circuits. The reconstructed density matrices are shown in Fig. 3(b,c), where techniques [23] are employed to mitigate incoherent errors as quite different effects caused by homo ( $(\text{CH}_3\text{O})_2\text{PHO}$ ) and hetero ( $\text{C}_2\text{F}_3\text{I}$ ) nuclear systems, and fidelities are calculated as 93.8% and 92.3% relative to the ideal GHZ state, respectively. By executing heterogeneous operations and measurements in parallel, these results validate TLP on our clusters and enable the preparation of target states beyond the qubit capacity of any single node. Full experimental details are provided in Appendix E.

For comparison, circuit cutting aggregates subresults classically and therefore requires no ancilla to probe over-

laps. However, it entails more subtasks (160 vs. 128 in our experiment) and heavier preprocessing to decompose a general entangling gate. Under the condition that the entangling gate is fixed as CNOT, as the number  $m$  of layers or inter-partition entangling gates increases, the subtask counts scale as  $16 \times 10^m$  vs.  $8 \times 2^m$ , underscoring the efficiency advantage of our approach to TLP. Moreover, our architecture, with single-ancilla estimator, benefits the implementation of general quantum channels, as demonstrated below.

*Non-Hermitian system simulation.* Efficient Hamiltonian simulation on quantum devices promises deep insights into quantum dynamics with broad applications [24, 25]. The GHZ preparation above illustrates a lightweight use of our multi-QPU cluster (of the order of  $10^2$  subtasks). Here we present a more demanding application involving the order of  $10^3$  subtasks executed across the  $16 \times 2$ -qubit cluster. Focusing on non-Hermitian Hamiltonians and imaginary-time evolution, we employ TLP across QPUs, implementing the linear-combination Hamiltonian simulation (LCHS) method [26].

Consider  $A(t) = H(t) - iL(t)$  with the Hermitian components  $H(t)$  and  $L(t)$ . LCHS is built on the identity, shown in Fig. 4(a),  $\mathcal{T} \exp\left(-i \int_0^t A(s) ds\right) = \int_{\mathbb{R}} \frac{dk}{\pi(1+k^2)} \mathcal{T} \exp\left(-i \int_0^t [H(s) + kL(s)] ds\right)$ , where  $\mathcal{T}$  denotes time ordering and the Cauchy-Lorentz kernel  $1/(\pi(1+k^2))$  arises from the Fourier representation of  $e^{-|x|}$ . Discretizing the  $k$ -integral on  $[-K, K]$  with uniform step  $\Delta k$ , quadrature nodes  $\{k_j\}$ , and trapezoidal weights  $\{w_j\}$  yields

$$u(t) \approx \sum_j c_j U_j(t), \quad (4)$$

with  $U_j(t) = \mathcal{T} \exp(-i \int_0^t [H(s) + k_j L(s)] ds)$  and  $c_j = w_j / (\pi(1+k_j^2))$ . Hence the quantum operation governed by  $A(t)$  fits as a special case of Eq. (3). If only expectation values are needed,  $\langle O \rangle_t = \langle u(t) | O | u(t) \rangle \approx \sum_{k, k'} c_k^* c_{k'} \langle u_0 | U_k^\dagger(t) O U_{k'} | u_0 \rangle$ , which is a subset of Eq. (3), with  $|u_0\rangle$  the initial state and  $|u(t)\rangle$  the time-evolved state. In our TLP setting, these terms are dispatched across QPU threads and aggregated classically.

For the experiment, we first consider dynamics generated by a single-qubit and reproduce time-independent non-Hermitian Hamiltonian simulation with  $H = \sigma_x$  and  $L = \mathbb{1} + \sigma_z$  [27, 28]. Here  $\sigma_{x,z}$  are Pauli matrices and  $\mathbb{1}$  is the identity. Initializing the system in  $|0\rangle$ , we record  $\langle \sigma_{x,z} \rangle$  as functions of the evolution time.

Based on LCHS, the non-Hermitian Hamiltonian dynamics decomposes into a collection of parallel subtasks (approximately  $400T^2$ , with  $T$  the evolution time) implemented using a single-ancilla estimator circuit. Each subtask corresponds to one experiment and is executed on a node of the  $16 \times 2$ -qubit NMR-QPU cluster. In total, we perform over 1700 experimental runs (about 106 per node) to resolve the time evolution ( $T =$

$0.1k$  ( $k = 1, \dots, 10$ )), thereby demonstrating TLP. For benchmarking, we also carry out numerical simulations using both the LCHS method and direct integration of the Schrödinger's equation. The two approaches agree with fidelity 99.5%, corroborating the accuracy of LCHS. Fig. 4(b) shows the time evolution of  $\langle \sigma_{y,z} \rangle$ , obtained from LCHS experiments, LCHS simulations, and direct Schrödinger integration. The absolute deviation is  $0.129 \pm 0.070$ . The close agreement across all traces demonstrates the feasibility and reliability of parallel non-Hermitian Hamiltonian simulation on our QPU cluster. Further evidence is provided in Fig. 4(d), which reports the state fidelity of each experimental snapshot relative to theory: the median fidelity is 93.1%, and a 90% threshold is indicated. See Appendix F for details.

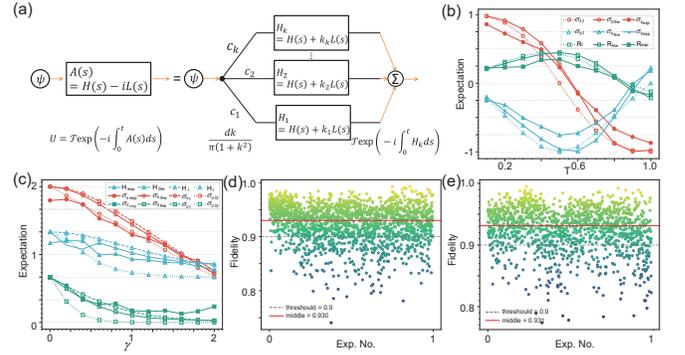


FIG. 4: (a) sketches the basic idea of linear combination of Hamiltonian simulations. (b) Time evolution of  $\langle \sigma_{y,z} \rangle$  and a randomly generated Hermitian observable under a non-Hermitian Hamiltonian. (c) Measured  $\langle H(\gamma) \rangle$ , and  $\langle \sigma_{x,z} \rangle$  after imaginary-time evolution of  $H(\gamma)$ . Solid lines with filled points indicate experimental data; hollow points with solid represents simulation result via experimental method; while dashed lines represent results via first principles calculations. In (c) dashed-dotted lines represent results from exact diagonalization and longer-time imaginary-time evolution. (d) and (e) show the fidelity of the experimentally prepared states relative to the numerical predictions, while the median fidelity and the 90% threshold are highlighted.

Then we demonstrate a ground-energy estimation via imaginary-time evolution for  $H(\gamma) = 2\mathbb{1} + \gamma\sigma_x$ ,  $\gamma \in [0, 2]$ , whose analytic ground energy is  $E_0(\gamma) = 2 - |\gamma|$ . Starting from  $|0\rangle$ , we apply non-unitary evolution  $e^{-H(\gamma)T}$  followed by renormalization, which approximates the ground state for sufficiently large  $T$ . The ground energy is then estimated by the expectation  $\langle H(\gamma) \rangle$ . This setting provides a clean testbed for imaginary-time simulation, given the closed-form  $E_0(\gamma)$  for benchmarking.

Similar to LCHS and more precisely via a truncated Cauchy integral representation [29], we approximate the imaginary-time propagator as a collection of parallel subtasks ( $\sim 121$  per  $\gamma$ ). We demonstrate TLP by executing over 1200 subtasks with a single-ancilla estimator circuit

on the  $16 \times 2$ -qubit NMR cluster. In experiments, we let  $\gamma = 0.2k$  ( $k = 0, \dots, 10$ ), and measure  $\langle H(\gamma) \rangle$ ,  $\langle \sigma_{x,y} \rangle$  with  $T = 0.5$ . For benchmarking, experimental results are compared against three references: a LCHS simulation, direct imaginary-time evolution with  $T = 0.5$  and exact diagonalization of  $H(\gamma)$  (and, as a tighter baseline, imaginary-time evolution with  $T = 1.5$ ). Fig. 4(c) shows the measured and simulated expectation values of  $\langle \sigma_{x,z} \rangle$ , and the ground energy as a function of  $\gamma$ . Experimental results, simulations of the experimental method, numerical imaginary-time evolution, and exact diagonalization or long-time (e.g.,  $T = 1.5$ ) imaginary-time evolution are shown with the absolute deviation  $0.136 \pm 0.088$ , which confirms that our parallel QPU cluster can efficiently and accurately realize imaginary-time evolution and ground-state energy estimation. Additional support is given by Fig. 4(e), which summarizes the fidelity between the experimentally prepared states and theoretical ground states for each  $\gamma$ . With a moderate total evolution time  $T = 0.5$ , all data points achieve fidelities exceeding 99%. Increasing  $T$  (e.g.,  $T = 1.5$ ) further improves the ground-state approximation but requires more experimental runs, reflecting the tradeoff between accuracy and experimental cost [See Appendix F for details].

*Conclusion.* While building large, general-purpose QPUs remains technologically and economically challenging, clustering existing devices enables significant performance enhancement without waiting for next-generation hardware. Even as large, fault-tolerant chips become available, their re-fabrication will likely remain costly, making QPU clustering a viable “performance boost” strategy, analogous to overclocking in classical systems, but without imposing extra physical strain on individual chips.

Our approach harnesses the strengths of current quan-

tum processors while respecting their limitations, enabling flexible parallel execution of diverse quantum operations and measurements. We show that quantum routines with product-state inputs and low-rank entangling layers can be re-expressed in an efficient parallelizable form. By orchestrating many small QPUs in parallel, we achieve TLP and scalability unattainable by a single device. This architecture was validated experimentally using NMR clusters of eight and sixteen nodes. Lightweight benchmarks distributed a four-qubit GHZ circuit across two-qubit nodes, reproducing the target state with 93.8% fidelity. To further stress-test the system, we implemented a resource-intensive protocol involving thousands of parallel executions for LCHS, successfully emulating both non-Hermitian and imaginary-time dynamics—constituting the first such demonstration on a modular quantum platform.

In summary, TLP on quantum devices immediately enlarges the logical register and enables complex channel-level operations inaccessible to standalone chips. This framework thus benefits other near-term quantum applications, such as variational quantum eigensolvers [30, 31]. Since our scheme relies only on classical interconnects, it is compatible with current fabrication capabilities and can be scaled, offering a pragmatic and flexible bridge from today’s few-qubit devices to future large-scale, fault-tolerant quantum accelerators.

## ACKNOWLEDGMENTS

We thank texra.ai for helpful suggestions on writing and language. K.L. and S.H. acknowledge the Scientific Foundation for Youth Scholars of Shenzhen University, Guangdong Provincial Quantum Science Strategic Initiative (GDZX2403001, GDZX2303001).

- 
- [1] J. L. Hennessy and D. A. Patterson, *Computer Architecture, Fifth Edition: A Quantitative Approach*, 5th ed. (Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2011).
  - [2] D. B. Kirk and W.-m. W. Hwu, *Programming Massively Parallel Processors: A Hands-on Approach*, 1st ed. (Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2010).
  - [3] J. Nickolls and W. J. Dally, *IEEE Micro* **30**, 56 (2010).
  - [4] M. Caleffi, M. Amoretti, D. Ferrari, J. Illiano, A. Manzalini, and A. S. Cacciapuoti, *Computer Networks* **254**, 110672 (2024).
  - [5] C. Monroe, R. Raussendorf, A. Ruthven, K. R. Brown, P. Maunz, L.-M. Duan, and J. Kim, *Phys. Rev. A* **89**, 022317 (2014).
  - [6] S. Wehner, D. Elkouss, and R. Hanson, *Science* **362**, eaam9288 (2018).
  - [7] L. Li, L. D. Santis, I. B. Harris, K. C. Chen, Y. Gao, I. Christen, H. Choi, M. Trusheim, Y. Song, C. Errando-Herranz, *et al.*, *Nature* **630**, 70 (2024).
  - [8] D. Main, P. Drmota, D. P. Nadlinger, E. M. Ainley, A. Agrawal, B. C. Nichol, R. Srinivas, G. Araneda, and D. M. Lucas, *Nature* **638**, 383–388 (2025).
  - [9] M. Mollenhauer, A. Irfan, X. Cao, S. Mandal, and W. Pfaff, *Nature Electronics* **8**, 610 (2025).
  - [10] X. Liu, X.-M. Hu, T.-X. Zhu, C. Zhang, Y.-X. Xiao, J.-L. Miao, Z.-W. Ou, P.-Y. Li, B.-H. Liu, Z.-Q. Zhou, *et al.*, *Nature Communications* **15**, 8529 (2024).
  - [11] C. Piveteau and D. Sutter, *IEEE Transactions on Information Theory* **70**, 2734–2745 (2024).
  - [12] A. W. Harrow and A. Lowe, *PRX Quantum* **6**, 010316 (2025).
  - [13] E. Bäumer, V. Tripathi, D. S. Wang, P. Rall, E. H. Chen, S. Majumder, A. Seif, and Z. K. Mineev, *PRX Quantum* **5**, 030339 (2024).
  - [14] A. Carrera Vazquez, C. Tornow, D. Ristè, S. Woerner, M. Takita, and D. J. Egger, *Nature* **636**, 75 (2024).
  - [15] T. C.-Q. development team, *Cuda-q* (2024).
  - [16] S. Arora, S. Rao, and U. Vazirani, *Journal of the ACM (JACM)* **56**, 1 (2009).

- [17] M. Henzinger, J. Li, S. Rao, and D. Wang, Deterministic near-linear time minimum cut in weighted graphs, in *Proceedings of the 2024 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)* (Society for Industrial and Applied Mathematics, 2024) pp. 3089–3139.
- [18] S.-Y. Hou, G. Feng, Z. Wu, H. Zou, W. Shi, J. Zeng, C. Cao, S. Yu, Z. Sheng, X. Rao, *et al.*, *EPJ Quantum Technology* **8**, 1 (2021).
- [19] G. Feng, S.-Y. Hou, H. Zou, W. Shi, S. Yu, Z. Sheng, X. Rao, K. Ma, C. Chen, B. Ren, G. Miao, J. Xiang, and B. Zeng, *IEEE Nanotechnology Magazine* **16**, 20 (2022).
- [20] M. Pont, G. Corrielli, A. Fyrrillas, I. Agresti, G. Carvacho, N. Maring, P.-E. Emeriau, F. Ceccarelli, R. Albiero, P. H. Dias Ferreira, *et al.*, *npj Quantum Information* **10**, 50 (2024).
- [21] Z. Bao, S. Xu, Z. Song, K. Wang, L. Xiang, Z. Zhu, J. Chen, F. Jin, X. Zhu, Y. Gao, *et al.*, *Nature Communications* **15**, 8823 (2024).
- [22] K. Mitarai and K. Fujii, *New Journal of Physics* **23**, 023021 (2021).
- [23] W. J. Huggins, S. McArdle, T. E. O’Brien, J. Lee, N. C. Rubin, S. Boixo, K. B. Whaley, R. Babbush, and J. R. McClean, *Phys. Rev. X* **11**, 041036 (2021).
- [24] K. Wang, L. Xiao, J. C. Budich, W. Yi, and P. Xue, *Phys. Rev. Lett.* **127**, 026404 (2021).
- [25] J. L. Bosse, A. M. Childs, C. Derby, F. M. Gambetta, A. Montanaro, and R. A. Santos, *Nature Communications* **16**, 2673 (2025).
- [26] D. An, J.-P. Liu, and L. Lin, *Phys. Rev. Lett.* **131**, 150603 (2023).
- [27] C. M. Bender, D. C. Brody, and H. F. Jones, *American Journal of Physics* **71**, 1095 (2003).
- [28] Y. Wu, W. Liu, J. Geng, X. Song, X. Ye, C.-K. Duan, X. Rong, and J. Du, *Science* **364**, 878 (2019).
- [29] M. Huo and Y. Li, *Quantum* **7**, 916 (2023).
- [30] L. Mineh and A. Montanaro, *Quantum Science and Technology* **8**, 035012 (2023).
- [31] M. Cattelan, S. Yarkoni, and W. Lechner, *npj Quantum Information* **11**, 27 (2025).

## Appendix A: Supplementary information for multi-QPU strategies

When several quantum-processing units (QPUs) are asked to execute a single workload, two fundamentally different resources can be used to tie them together:

1. *Quantum correlated links* — Entanglement is distributed between modules and exploited through gate teleportation, remote measurement, or lattice surgery. Theoretical blueprints date back to modular ion-trap architectures with photonic routers [5] and underpin the long-term vision of distributed quantum computing [7–9] or a quantum internet [6, 10].
2. *Classically correlated links* — Each QPU is measured locally; the global observable is then reconstructed from the measurement records with “divide-and-conquer” methods. This family includes wire- and gate-based quantum circuit cutting [12], circuit knitting [11] and dynamic circuits techniques [13, 14].

Both routes trade resources differently. Quantum links preserve coherence across modules and avoid exponential sampling overhead, but they impose stringent latency and loss budgets on the interconnect. Classical links need no fragile quantum channel and can already incorporate non-unitary dynamics, yet their classical post-processing overhead grows with every wire or gate that is cut. Specifically, Table S1 briefly summarizes these representative multi-QPU strategies and their key features.

TABLE S1: Representative multi-QPU strategies and their key features

Work	Year	Correlation	Interconnect	Non-unitary	Key feature	Hardware scale
Distributed Quantum Computing						
Monroe <i>et al.</i> [5]	2014	Quantum	Probabilistic photonic interface	n/a	Modular blueprint	Proposal
Spin-photon CMOS [7]	2024	Quantum	photonics	No	Wafer-scale spin-photon node	1 mm <sup>2</sup> chip
Trapped-ion link [8]	2025	Quantum	Heralded photons	No	Deterministic gate teleport	2 ions module (~ 2m)
Microwave cable [9]	2025	Quantum	Plug-and-play coaxial cable	No	99 % swap in 100ns	2 transmons(≤ 1m)
Quantum Internet						
Optical 7 km gate [10]	2024	Quantum	Telecom fiber	No	Metro-scale entangling gate	2 memories
Dynamic circuits						
Dynamic circuits [13]	2024	Classical	Conditional feed-forward operations	No	101-qubit CNOT teleport	127-qubit chip
Dynamic-circuit cutting [14]	2024	Classical	Real-time classical link	No	142-qubit graph via cutting	2×127 qubits
Quantum circuit cutting						
Optimal circuit cuts [12]	2025	Classical	Theory	No	Tight sample-cost bounds	Theoretical
Quantum circuit knitting						
Circuit knitting [11]	2024	Classical	Software only	No	Reducing gates overhead	Simulator
Others						
CUDA-Q multi-QPU [15]	2024	Classical	Software only	n/a	Async kernels; More tests required	Simulator
<b>This work</b>	2025	Classical	Ethernet	<b>Yes</b>	Efficient in a large class of CPTP map.	16×2/8 × 3 spins

### a. Quantum-correlated prototypes.

- Spin-photon CMOS integration [7] - Li *et al.* heterogeneously bond diamond colour-centre qubits onto 45 nm CMOS photonics, enabling wafer-scale quantum light sources and detectors.
- Trapped-ion optical link [8] - Two ion chains 2 m apart execute a distributed Grover search via deterministic photonic gate teleportation, reaching 86 % gate fidelity.
- Plug-and-play microwave cable [9] - A detachable coaxial bus swaps microwave photons between two superconducting chips in < 100 ns with > 99% efficiency.
- 7 km non-local photonic gate [10] - Quantum memories plus telecom photons realise entangling gates over a 7 km fibre loop.

### b. Classically correlated prototypes.

- Dynamic-circuit cutting - Shallow measure-and-feed-forward protocols teleport CNOTs across 100+ qubits [13], while two 127-qubit processors linked only by real-time classical feed-forward prepare 142-qubit graph states [14].
- Circuit knitting with LOCC [11] - Information-theoretic analysis shows classical messaging reduces sampling overhead for knitted circuits.
- Optimal quantum circuit cuts [12] - Tight bounds on sample cost and an application to clustered Hamiltonian simulation.

*c. This work - thread-level parallelism for quantum channels.* We introduce a method that factorizes an arbitrary completely-positive trace-preserving map into independent Kraus blocks, dispatches those blocks to a cluster of QPUs over standard Ethernet, and recombines the measurement records in  $\mathcal{O}(1)$  classical time. Unlike prior “divide-and-conquer” methods that target only unitary circuits [13, 14], our scheme natively supports non-unitary dynamics such as non-hermitian Hamiltonian evolution and imaginary-time propagation.

*Hardware demonstration*— A  $16 \times 2$ -qubit (or  $8 \times 3$ -qubit) NMR cluster prepares a 4-qubit GHZ state with a tomographic fidelity of 93.8%; emulates non-Hermitian Hamiltonian and imaginary-time dynamics. These results establish that purely classical networking can extend both the size and the algorithmic richness of NISQ workloads beyond what any single device can achieve today, delivering the first channel-level TLP benchmark on real hardware.

## Appendix B: Supplementary information for realizing thread level parallelism

We start from the generic quantum-computing objective

$$\mathrm{Tr}[O\Phi(\rho)],$$

where  $\Phi$  is a quantum channel,  $\rho$  an input state and  $O$  the observable of interest.

$\Phi$  can be a general structure that admits a linear-combination-of-unitaries form

$$\Phi(\cdot) = \sum_{p=1}^q \left( \sum_{i=1}^m c_{p,i} U_{p,i} \right) (\cdot) \left( \sum_{j=1}^m c_{p,j} U_{p,j} \right)^\dagger, \quad (\text{B1})$$

with complex amplitudes  $c_{p,i}$  and unitaries  $\{U_{p,i}\}_{i=1}^m$ . Substituting Eq. (B1) into the trace and expanding gives

$$\mathrm{Tr}[O\Phi(\rho)] = \sum_{p=1}^q \sum_{i,j=1}^m c_{p,i} c_{p,j}^* \mathrm{Tr}[OU_{p,i}\rho U_{p,j}^\dagger]. \quad (\text{B2})$$

Here, we employ a without losing generosity assumption that both the state and the observable are product states over  $\mathcal{A}$  identical subsystems,

$$\rho = \bigotimes_{a=1}^{\mathcal{A}} \rho^a, \quad O = \bigotimes_{a=1}^{\mathcal{A}} O^a.$$

And a non-trivial requirement that every unitary  $U_{p,i}$  have an explicit expression of the summation of tensor-product decomposition,

$$U_{p,i} = \sum_{\alpha=1}^{\ell} \bigotimes_{a=1}^{\mathcal{A}} U_{p,i\alpha}^a, \quad (\text{B3})$$

with the same  $\ell$  for all indices. Insert Eq. (B3) (and its Hermitian conjugate) into Eq. (B2); linearity yields

$$\mathrm{Tr}[O\Phi(\rho)] = \sum_{p=1}^q \sum_{i,j=1}^m c_{p,i} c_{p,j}^* \sum_{\alpha,\alpha'=1}^{\ell} \mathrm{Tr}\left[O\left(\bigotimes_a U_{p,i\alpha}^a\right)\rho\left(\bigotimes_a U_{p,j\alpha'}^a\right)^\dagger\right]. \quad (\text{B4})$$

Because both the state, observable and each tensor factor act on disjoint Hilbert spaces, the global trace factorizes into a product of local traces:

$$\mathrm{Tr}\left[O\left(\bigotimes_a U_{p,i\alpha}^a\right)\rho\left(\bigotimes_a U_{p,j\alpha'}^a\right)^\dagger\right] = \prod_{a=1}^{\mathcal{A}} \mathrm{Tr}\left[U_{p,i\alpha}^a \rho^a U_{p,j\alpha'}^{a\dagger} O^a\right]. \quad (\text{B5})$$

Substituting this identity into Eq. (B4) yields the fully factorised expression

$$\sum_{p=1}^q \sum_{i,j=1}^m c_{p,i} c_{p,j}^* \sum_{\alpha,\alpha'=1}^{\ell} \prod_{a=1}^{\mathcal{A}} \text{Tr}[U_{p,i\alpha}^a \rho^a U_{p,j\alpha'}^{a\dagger} O^a], \quad (\text{B6})$$

which is exactly Eq. (3) in the main text. Each factor inside the product is measurable on a single subsystem with programmable unitary transformations and measurements, while the outer sums can be accumulated classically. This indicates thread-level parallelism, TLP, can be realized across quantum processors.

In summary, every quantum-information task of the form  $\text{Tr}[O \Phi(\rho)]$  can be split into  $q \times m^2 \times \ell^2 \times \mathcal{A}$  independent sub-tasks, provided that the state  $\rho$  and observable  $O$  factorize over  $\mathcal{A}$  subsystems, and each operator  $U_{p,k}$  admits the tensor-product expansion of Eq. (B3). The integers  $q, m, \ell$  and  $\mathcal{A}$  are precisely those introduced in Eqs. (B1) and (B3). Furthermore, if the channel reduces to an incoherent sum of unitaries,  $\Phi(\rho) = \sum_{k=1}^q U_k \rho U_k^\dagger$ , and every  $U_k$  further decomposes as  $U_k = \sum_{\alpha=1}^{\ell} \bigotimes_{a=1}^{\mathcal{A}} U_{k\alpha}^a$ , then the workload factorizes into only  $q \times \ell^2 \times \mathcal{A}$  sub-tasks, because the double index  $(i, j)$  collapses to a single index  $k$  ( $m = 1$ ).

### Appendix C: Supplementary information for pre-conditions in architecture

The first is an efficient factorization that rewrites the global unitary into sums of local operators. A brute-force operator-Schmidt decomposition would accomplish factorization step, its cost grows as  $\Theta(d^{2n})$  for  $n$  qubits and is therefore intractable.

Instead, we exploit the layered structure of quantum circuits to recursively apply precomputed local factorizations. In general, given a depth- $m_d$  circuit

$$U = \prod_{i=1}^{m_d} \left( G_e^{(i)} \otimes_{a=1}^{\mathcal{A}} U_a^{(i)} \right),$$

whose  $i$ -th layer contains a single entangling block  $G_e^{(i)}$  acting across the  $\mathcal{A}$  subsystems, and a factorized expansion of each entangling block,

$$G_e^{(i)} = \sum_{\alpha_i=1}^{\ell_i} \bigotimes_{a=1}^{\mathcal{A}} g_{i\alpha_i}^a, \quad g_{i\alpha_i}^a \in \text{U}(\dim \mathcal{H}_a),$$

we can rewrite the entire circuit as

$$U = \prod_{i=1}^{m_d} \sum_{\alpha_i=1}^{\ell_i} \bigotimes_{a=1}^{\mathcal{A}} g_{i\alpha_i}^a U_a^{(i)},$$

i.e. as a nested linear combination of  $\prod_{i=1}^{m_d} \ell_i$  purely local tensor-product operators.

Bisecting the circuit is a good start of the problem. It is trivial that one have a way to bisect a unitary into sub circuits. For current popular circuit, entangling block can be chosen as CNOT or CZ, which have a fixed factorized expansion. Therefore,  $m_d$  appears,  $\prod_{i=1}^{m_d} \ell_i = \ell$ , which influences the number of sub tasks take this form. To optimize it, we introduce an algorithm of layer-wise decomposition (Lemma C.1) to reduce it.

**Lemma C.1** (Layer-wise decomposition). *Let*

$$U = \prod_{i=1}^{m_d} \left( G_e^{(i)} \bigotimes_{a=1}^n U_a^{(i)} \right)$$

be an  $n$ -qubit circuit of depth  $m_d$ , where at each layer  $i$ ,  $G_e^{(i)}$  is a set of two-qubit gates that may couple any pair of qubits, and  $U_a^{(i)}$  acts locally on qubit  $a$ .

- (i) Viewing the circuit as a weighted graph  $G = (V, E)$  with  $|V| = n$  and edge weights  $w_{uv}$  equal to the number of two-qubit gates acting on  $\{u, v\}$ , one can find in  $\tilde{O}(|E|) \subseteq \tilde{O}(n^2)$  time (Henzinger-Li-Rao-Wang [17]) a minimum cut. The cut removes  $m' \leq m_d$  crossing gates  $\{g^{(i)}\}_{i=1}^{m'} \subseteq \{G_e^{(i)}\}$  and partitions the qubits into two subsystems, hereafter indexed by  $a \in \{1, 2\}$ .

(ii) Suppose each crossing gate admits a known factorization  $g^{(i)} = \sum_{\alpha_i=1}^{\ell_i} V_{i\alpha_i}^{(1)} \otimes V_{i\alpha_i}^{(2)}$ , where all  $V_{i\alpha_i}^{(a)}$  are unitary ( $\ell_i \leq 4$  for CNOT/CZ). Denoting  $\tilde{U}_a^{(i)}$  as the local operation within subsystem  $a$ , the entire circuit decomposes as

$$U = \prod_{i=1}^{m'} \sum_{\alpha_i=1}^{\ell_i} \bigotimes_{a=1}^2 V_{i\alpha_i}^{(a)} \tilde{U}_a^{(i)},$$

which contains at most  $\prod_{i=1}^{m'} \ell_i$  tensor terms, independent of the dimensions of the two subsystems.

Remarkably, if each side of the cut must contain roughly  $n/2$  qubits, the partition problem becomes NP-hard; an  $O(\sqrt{\log n})$ -approximation is available via the SDP + metric-embedding technique of Arora-Rao-Vazirani [16]. Because  $\ell_i$  is a small constant (at most 4 for most common two-qubit gates), the total number of tensor terms grows with the cut size  $m'$ , i.e.,  $\prod_{i=1}^{m'} \ell_i = \ell$ , not with the full Hilbert-space dimension  $2^n$ .

**Result C.1** (Thread-level parallelism (TLP) criterion). *Consider a workload whose objective value is*

$$\text{Tr}[O \Phi(\rho)], \quad \Phi(\cdot) = \sum_{p=1}^q \left( \sum_{i=1}^m c_{p,i} U_{p,i} \right) (\cdot) \left( \sum_{j=1}^m c_{p,j} U_{p,j} \right)^\dagger,$$

and assume

1. the data are tensor-separable,  $\rho = \bigotimes_{a=1}^A \rho^a$  and  $O = \bigotimes_{a=1}^A O^a$ ;
2. each  $U_{p,i}$  is realized by a depth- $m_d$  circuit  $U_{p,i} = \prod_{t=1}^{m_d} (G_e^{(t)} \otimes_a U_a^{(t)})$  in which every entangling layer factorizes as  $G_e^{(t)} = \sum_{\alpha_t=1}^{\ell_t} \bigotimes_{a=1}^A g_{t\alpha_t}^a$ .

Then there exists a cut size  $m' \leq m_d$  such that the global expectation value factorizes into at most

$$q \times m^2 \times \left( \prod_{t=1}^{m'} \ell_t \right)^2 \times \mathcal{A}$$

independent single-subsystem traces. Each trace can be executed on an individual QPU node, and the final result is obtained by classical post-processing.

**Implications.** Whenever the input state and observable are already tensor-separable, and all non-local gates admit a known low-rank tensor expansion (e.g. CNOT, CZ), the entire computation can be mapped onto a QPU cluster with one thread per tensor term—no quantum interconnect is required, and the required single-ancilla estimator is depicted as follows.

**Lemma C.2** (Single-ancilla estimator). *Let  $\rho^a = |\psi_0^a\rangle\langle\psi_0^a|$  be a pure state on sub-system  $a$ , let  $U_{i\alpha}^a$  and  $U_{j\alpha'}^a$  be arbitrary unitaries on  $\mathcal{H}_a$ , and let  $O^a$  be any unitary observable. The complex overlap  $\text{Tr}(U_{i\alpha}^a \rho^a U_{j\alpha'}^{a\dagger} O^a)$  is obtained with a single ancillary qubit by*

$$\langle \sigma_x \rangle_{\text{anc}} + i \langle \sigma_y \rangle_{\text{anc}},$$

where the expectation values are measured in the interferometric circuit, which is shown in Fig. S1. Setting  $j = i, \alpha' = \alpha$  recovers channel observables of the form  $U_{i\alpha}^a \rho^a U_{i\alpha}^{a\dagger}$ .

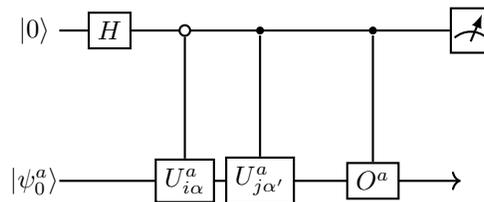


FIG. S1: Single-ancilla estimator: an ancillary qubit is used to estimate overlaps.

*Proof sketch*— Prepare the ancilla in  $|0\rangle$  and subsystem  $a$  in  $|\psi_0^a\rangle$ . Apply the controlled unitary  $|0\rangle\langle 0| \otimes U_{i\alpha}^a + |1\rangle\langle 1| \otimes U_{j\alpha'}^a$  to obtain

$$|\Psi_{i,j,\alpha,\alpha'}^{(a)}\rangle = \frac{1}{\sqrt{2}}(|0\rangle \otimes U_{i\alpha}^a |\psi_0^a\rangle + |1\rangle \otimes U_{j\alpha'}^a |\psi_0^a\rangle).$$

Next, insert  $O^a$  only on the  $|1\rangle$  branch:

$$|\Psi_{i,j,\alpha,\alpha'}^{O,(a)}\rangle = \frac{1}{\sqrt{2}}(|0\rangle \otimes U_{i\alpha}^a |\psi_0^a\rangle + |1\rangle \otimes O^a U_{j\alpha'}^a |\psi_0^a\rangle).$$

Tracing out subsystem  $a$  leaves the ancilla density matrix

$$\rho_{\text{anc}} = \frac{1}{2} \begin{pmatrix} 1 & \langle \psi_0^a | U_{i\alpha}^{a\dagger} O^a U_{j\alpha'}^a | \psi_0^a \rangle \\ \langle \psi_0^a | U_{j\alpha'}^a O^a U_{i\alpha}^{a\dagger} | \psi_0^a \rangle & 1 \end{pmatrix}.$$

The off-diagonal element is precisely the desired overlap. Projective measurements of the ancilla in the  $\sigma_x$  and  $\sigma_y$  bases read out its real and imaginary parts. For channels of the form  $U_{i,\alpha} \rho U_{i,\alpha}^\dagger$  one simply sets  $j, \alpha' = i, \alpha$ , recovering the same protocol.  $\blacksquare$

Together, Lemma C.1 and C.2 meet both prerequisites of Eq. (3): an efficient factorization whose cost depends only on the interaction cut, and a practical estimator that requires just one extra qubit per subsystem.

## Appendix D: Supplementary information for configuration of current Hardware

In this work, we implement clusters comprising either 8 or 16 benchtop NMR quantum processors, with each QPU node hosting two or three nuclear spin qubits.

First, we introduce a cluster of  $8 \times 3$ -qubit QPUs based on the Triangulum platform, as detailed in [19]. As shown in Fig. 2(b), each device consists of three qubits, realized by the  $^{19}\text{F}$  nuclei in  $\text{C}_2\text{F}_3\text{I}$  molecules.

The relevant relaxation times,  $T_1$  (longitudinal) and  $T_2$  (transverse), are measured for all three nuclei and determine the system's coherence properties. The free evolution of each three-qubit NMR system is governed by the internal Hamiltonian:

$$\mathcal{H}_{\text{int}} = \pi\nu_1\sigma_z^1 + \pi\nu_2\sigma_z^2 + \pi\nu_3\sigma_z^3 + \frac{\pi}{2}J_{12}\sigma_z^1\sigma_z^2 + \frac{\pi}{2}J_{23}\sigma_z^2\sigma_z^3 + \frac{\pi}{2}J_{13}\sigma_z^3\sigma_z^1, \quad (\text{D1})$$

where  $\nu_i$  ( $i = 1, 2, 3$ ) denote the Larmor frequencies, and  $J_{ij}$  are the  $J$ -coupling constants between the spins.

Universal quantum control is achieved by applying transverse radio-frequency (RF) pulses, described by

$$\mathcal{H}_{\text{rf}} = -\frac{1}{2}\omega_1 \sum_{i=1}^3 [\cos(\omega_{\text{rf}}t + \phi)\sigma_x^i + \sin(\omega_{\text{rf}}t + \phi)\sigma_y^i]. \quad (\text{D2})$$

By tuning the RF field parameters—amplitude  $\omega_1$ , phase  $\phi$ , frequency  $\omega_{\text{rf}}$ , and pulse duration—together with the system's internal evolution, arbitrary three-qubit quantum gates can be realized.

In our experiments, we employ 8 Triangulum units, with their main specifications summarized in Table S2. Device-to-device variations in magnetic field strength, chemical shifts, and other parameters are hard to directly compensated via software calibration and by operating in the rotating frame. Therefore, variant circuit or shaped pulses are employed for robust control or generating standard quantum gates.

Second, we introduce  $16 \times 2$  qubit QPUs, which uses Gemini, referenced in [18]. As shown in Fig. 2(b), the device comprises two qubits, represented by the nuclei  $^1\text{H}$  and  $^{31}\text{P}$  in Dimethylphosphite ( $(\text{CH}_3\text{O})_2\text{PH}$ ) molecules.  $T_1$  and  $T_2$  represent the longitudinal and transverse relaxation times, respectively. The free evolution of this 2-qubit system is primarily governed by the internal Hamiltonian,

$$\mathcal{H}_{\text{int}} = \pi\nu_1\sigma_z^1 + \pi\nu_2\sigma_z^2 + \frac{\pi}{2}J_0\sigma_z^1\sigma_z^2, \quad (\text{D3})$$

where  $\nu_1$  and  $\nu_2$  can be adjusted to 0 Hz in a rotating frame, and  $J_0 = 697.4$  Hz denotes the resonance frequency of the  $J$ -coupling strength between the spins. Further details are available in [18]. To control the system's evolution,

TABLE S2: Specifications of the SPINQ Triangulum Units

Model	Serial No.	Magnetic Field (T)	Homogeneity FWHM (ppm)	$^{19}\text{F}$ Freq. (MHz)
Triangulum	T2023089	0.886	0.8	35.470
	T2023043	0.905	0.8	36.249
	T20250617	0.895	0.7	35.858
	T20250618	0.895	0.7	35.858
	T20231213	0.944	0.85	37.809
	T20250208	0.888	0.6	35.581
	T20250504	0.896	0.7	35.875
	T20250420	0.866	0.6	34.684

transverse radio-frequency (r.f.) pulses serve as the control field, expressed as,

$$\mathcal{H}_{rf} = -\frac{1}{2} \sum_{i=1}^2 \omega_1^i (\cos(\omega_{rf}t + \phi^i) \sigma_x^i + \sin(\omega_{rf}t + \phi^i) \sigma_y^i). \quad (\text{D4})$$

By adjusting the parameters in the r.f. field [Eq. (D4)], such as intensity  $\omega_1$ , phase  $\phi$ , frequency  $\omega_{rf}$ , and duration, the theoretical achievement of two-qubit universal quantum gates is possible through the combination of the system's internal dynamics.

In our experiments, we employ 16 Gemini units, with their main specifications summarized in Table S3. While there are device-to-device variations in magnetic field strength, chemical shifts, and other parameters, these differences are compensated in software by appropriate calibration and by working in the rotating frame. Standard NMR techniques such as hard and shaped pulses are employed for robust control.

TABLE S3: Specifications of the SPINQ Gemini Lab Units

Model	Serial No.	Magnetic Field (T)	Homogeneity FWHM (ppm)	$^1\text{H}$ Freq. (MHz)	$^{31}\text{P}$ Freq. (MHz)
Gemini Lab	L20250104	0.654	0.6	27.848	11.273
	L20250607	0.644	0.5	27.436	11.106
	L20250415	0.661	0.3	28.138	11.390
	L20250414	0.657	0.5	27.984	11.328
	L20250419	0.656	0.4	27.949	11.314
	L20250413	0.660	0.3	28.097	11.375
	L20250412	0.661	0.4	28.124	11.385
	L20250418	0.645	0.4	27.474	11.121
	L20250417	0.653	1.0	27.806	11.256
	L20250411	0.638	0.4	27.172	10.999
	L20250410	0.644	0.5	27.405	11.094
	L20250602	0.651	0.7	27.721	11.221
	L20250421	0.636	0.6	27.061	10.954
	L20250416	0.654	0.5	27.865	11.280
	L20250620	0.648	0.5	27.578	11.163
	L20250606	0.621	0.7	26.449	10.706

## Appendix E: Supplementary information for GHZ state preparation

### 1. Our method

Fig. S2 provides circuit to generate a 4-qubit GHZ state using a 4 qubit unitary transformation. According to Lemma C.1, the 4-qubit GHZ circuit is decomposed into parallel two-qubit subcircuits, each suitable for small-scale

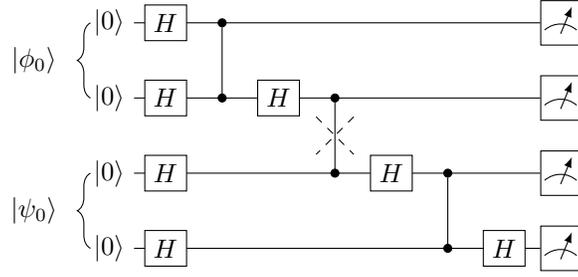


FIG. S2: Typical quantum circuit for 4 qubit GHZ state, where control-Z gate between the 2nd and 3rd qubits is cut to generate sub-circuits.

hardware. As the CNOT gate between qubits 2 and 3 can be “cut” and replaced by a sum of local operations:

$$U_{cnot} = \frac{1}{2} (\mathbb{1} \otimes \mathbb{1} + \sigma_z \otimes \mathbb{1} + \mathbb{1} \otimes \sigma_x - \sigma_z \otimes \sigma_x),$$

with Pauli operators  $\sigma_{x,y,z}$  and identity  $\mathbb{1}$ , we can generate a replaced circuits, shown in Fig. S3.

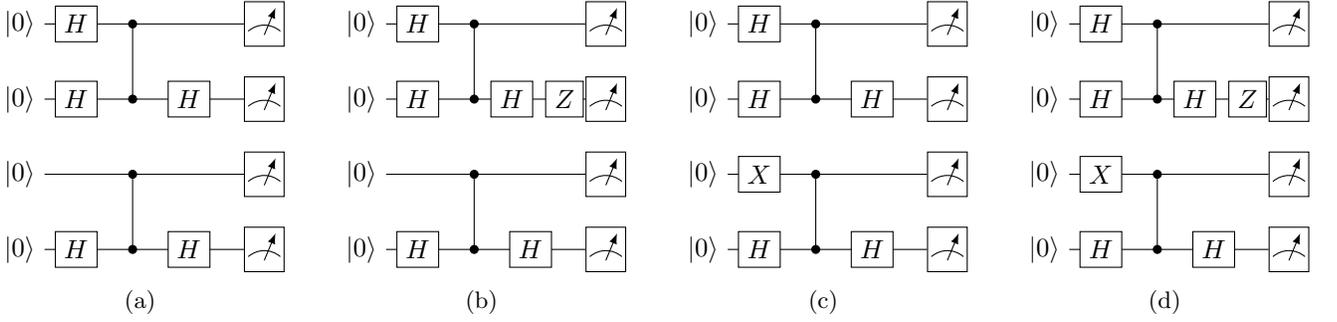


FIG. S3: Experimental circuit for 4-qubit GHZ state. circuits in (a-d) are generated by cutting the control-Z gate between the 2nd and 3rd qubits and replace it by a sum of local operations. Those circuit are summed in a superposition way.

Mathematically, the output state is a sum of product of sub-states

$$\sum_{\{j,k\}} (-1)^{j-1 \cdot k-1} |\phi_j\rangle \otimes |\psi_k\rangle,$$

where  $j, k = 1, 2$ .  $|\phi_1\rangle = U_{cnot}(H \otimes I)|00\rangle$ ,  $|\phi_2\rangle = I \otimes \sigma_z |\phi_1\rangle$ ,  $|\psi_1\rangle = U_{cnot}|00\rangle$ , and  $|\psi_2\rangle = U_{cnot}(\sigma_x \otimes I)|00\rangle$ . This indicates the entire preparation process can be parallelized. Via combining single-ancilla estimator, we can form the experimental circuit depicted in Fig. S4, and measure the ancilla qubit with  $\sigma_x$  and  $\sigma_y$ . Through varying  $M$  in the circuit for all two-qubit Pauli bases, we can get sufficient information for reconstructing a GHZ state.

In this experiment, two types of circuits are employed, with  $M$  varied to cover all Pauli bases required for full tomography. If every expectation value for a given Pauli basis is counted as an independent measurement, a total of  $4 \times 16$  measurements are required. This is notably reduced compared to conventional four-qubit experiments, as both real and imaginary parts are separately recorded using the ancilla-assisted circuit.

Each node of the Triangulum-pro desktop spectrometer hosts three  $^{19}\text{F}$  nuclear spins in a  $\text{C}_2\text{F}_3\text{I}$  molecule, with typical coherence times of  $T_2 \approx 500 \text{ ms}$  and a calibrated maximum Rabi rate of  $|\omega_1|/2\pi = 8.3 \text{ kHz}$ . Eight such nodes are synchronized by a common clock and exchange shot data via Gigabit Ethernet. We focus here on one representative experiment to illustrate the procedure:

*Initialization* — A spatial operation sequence is applied to transform the thermal equilibrium state into a pseudo-pure ground state  $|000\rangle$ . The resulting state fidelity reaches 98.5% (see spectrum in Fig. S5).

*Gate Implementation:* The unitary operator represented by the sub-circuits (shown in Figure S4 dashed boxes) is realized by a single shaped pulse of 28 ms duration, discretized into  $M = 800$  time slices ( $\Delta t = 35 \mu\text{s}$ ).

*Pulse Optimization* — An initial guess is generated using random spline envelopes  $\{B_x(t), B_y(t)\}$  constrained by  $|B| \leq 0.4|\omega_1|$ , and accepted only if the simulated unitary fidelity exceeds 20%. The GRAPE algorithm is then

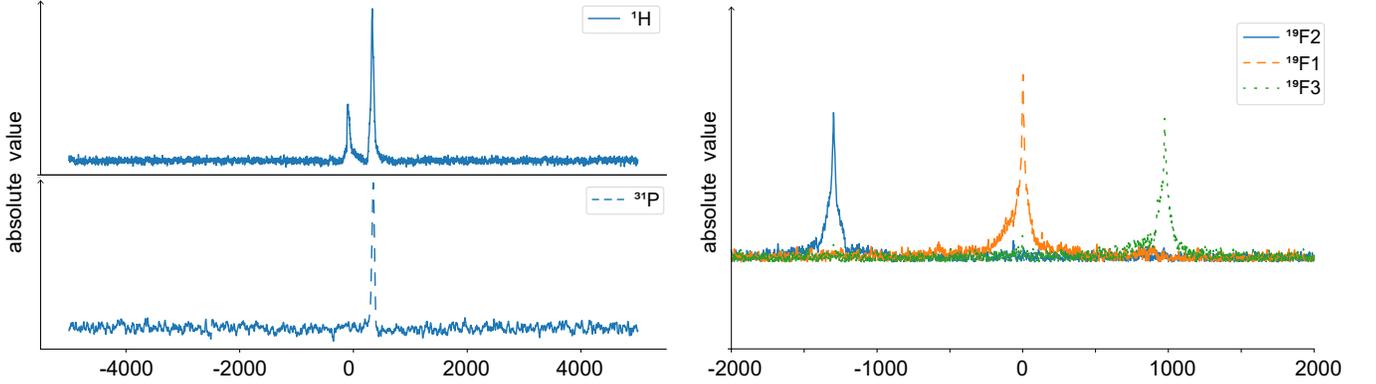
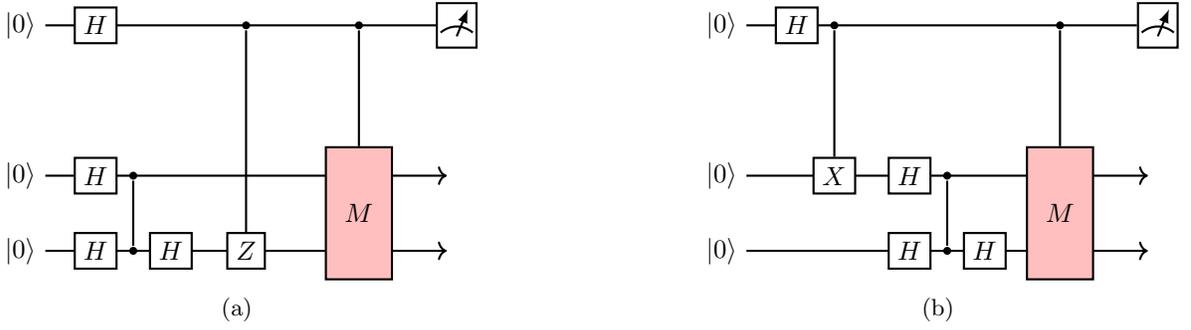


FIG. S5: Spectrum of all spins when system is at pseudo-pure state for  $(\text{CH}_3\text{O})_2\text{PHO}$  (left) and  $\text{C}_2\text{F}_3\text{I}$  (right). The horizontal axis is the chemical shift while vertical axis is absolute value of signal.

employed to maximize the worst-case channel fidelity  $\mathcal{F}_{\min}$  across nine error configurations, spanning three static  $B_0$  offsets ( $-20$ ,  $0$ , and  $+20$  Hz) and three RF amplitude miscalibrations ( $0.95$ ,  $1.00$ ,  $1.05$ ). Optimization typically converges within 15–20 L-BFGS iterations. The final control amplitudes and phases are discretized to 0.01% resolution and exported as a waveform file with total pulse duration of 28 ms. GRAPE simulations yield an average gate fidelity of  $\overline{\mathcal{F}}_{\text{gate}} = 99.3\%$  over all error conditions.

*GHZ State Preparation* — Finally, the GRAPE-compiled pulse is executed concurrently on all nodes in the cluster. Preparing GHZ State across two neighboring nodes yields a state fidelity of 93.8%.

## 2. Quantum Circuit Cutting

Similarly, we implement quantum circuit cutting to decompose a 4-qubit GHZ preparation circuit into a sum of eleven 2-qubit subcircuits, following the procedure of [22]. Unlike our architecture, which allows coherent superpositions and therefore requires ancillary systems for coherent control, quantum circuit cutting operates entirely classically. That is, the summation in circuit cutting is incoherent—no entanglement or interference is preserved—thereby avoiding the need for ancilla but increasing the classical processing overhead.

Quantum circuit cutting is a powerful strategy to simulate large quantum circuits on smaller hardware by dividing a global circuit into independently executable fragments. This technique is particularly valuable for near-term quantum devices constrained by limited qubit counts, coherence times, and gate fidelities. It enables logical circuit execution beyond physical limitations, relying on classical coordination and tomography-based post-processing.

In the spatial (or state-space) variant of circuit cutting, we approximate a two-subsystem unitary channel as

$$U_{ab}\rho U_{ab}^\dagger \approx \sum_i c_i U_{ai}\rho_a U_{ai}^\dagger \otimes U_{bi}\rho_b U_{bi}^\dagger, \quad (\text{E1})$$

assuming an initially factorized input state  $\rho = \rho_a \otimes \rho_b$ . This expansion allows the joint evolution under  $U_{ab}$  to

be mimicked by local operations  $U_{ai}$  and  $U_{bi}$  on separate devices, followed by classical post-processing. Figure S6 illustrates this decomposition, where the bipartite gate is expressed as a convex combination of local unitaries applied independently to subsystems  $a$  and  $b$ .

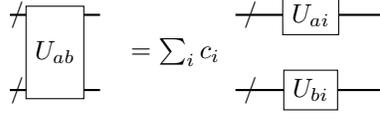


FIG. S6: Schematic of quantum circuit cutting from an entire system to a bi-partition one.

In our demonstration, we apply this scheme to the 4-qubit GHZ circuit shown in Fig. S2. The key component—namely, the controlled-Z gate connecting the two partitions—is decomposed into a sum over tensor products of single-qubit operators. This decomposition reduces the circuit into 2-qubit subcircuits, which can then be distributed across multiple 2-qubit quantum processors.

To reconstruct the final GHZ state, each subcircuit is executed separately and its output measured in all required Pauli bases. In total, 10 distinct subcircuits are used, and each requires measurements in 16 different Pauli settings for full state tomography. Therefore, a total of  $10 \times 16 = 160$  measurement configurations are needed for complete reconstruction.

$$\frac{1}{2} \begin{array}{c} \boxed{e^{i\pi Z/4}} \\ \boxed{e^{i\pi Z/4}} \end{array} + \frac{1}{2} \begin{array}{c} \boxed{e^{-i\pi Z/4}} \\ \boxed{e^{-i\pi Z/4}} \end{array} - \frac{1}{2} \sum_{\alpha \in \{\pm 1\}^2} \alpha_1 \alpha_2 \left\{ \begin{array}{c} \boxed{(I + \alpha_1 Z)/2} \\ \boxed{e^{i(\alpha_2 + 1)\pi Z/4}} \\ \boxed{e^{i(\alpha_2 + 1)\pi Z/4}} \\ \boxed{(I + \alpha_2 Z)/2} \end{array} \right\}$$

FIG. S7: Decompositions of controlled-Z gate into a sequence of single-qubit operations.

Each node of the Gemini-Lab desktop spectrometer hosts a  $^1\text{H}$  nuclear spin and a  $^{31}\text{P}$  nuclear spins in a Dimethyl Phosphite molecule  $(\text{CH}_3\text{O})_2\text{PH}$ , with typical coherence times of  $T_2 \approx 500$  ms and a calibrated maximum Rabi rate of  $|\omega_1|/2\pi = 6.25$  kHz. Eight such nodes are synchronized by a common clock and exchange shot data via Gigabit Ethernet. We focus here on one representative experiment to illustrate the procedure:

*Initialization* — An operation sequence named relaxation method is applied to transform the thermal equilibrium state into a pseudo-pure ground state  $|00\rangle$ . The resulting state fidelity reaches 99.5% (See Fig. S5).

*Gate Implementation:* A parameterized two-qubit circuit is used to realize the sub-circuits shown in Figure S4 (dashed boxes). The parameterized quantum circuit contains single-qubit arbitrary angle rotation gate, rotation parameters and two-qubit gate (controlled-not gate, CNOT gate), which is shown in Fig. S8. We use this one as this is a fully compiled gate set. Each parameter is optimized through the gradient descent algorithm. The duration of each single-qubit rotation gates is about 0.3 ms, and the duration of CNOT gate is about 2 ms. Their theoretical fidelity are both greater than 99%.

*GHZ State Preparation* — Finally, the circuits are executed concurrently on all nodes in the cluster. Preparing GHZ State across two neighboring nodes yields a state fidelity of 92.3%.

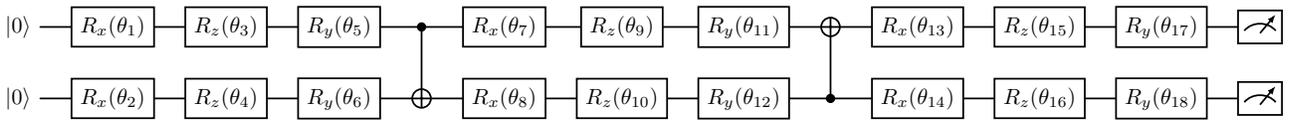


FIG. S8: Parametrized circuit to implement the entire unitary operations in system of  $(\text{CH}_3\text{O})_2\text{PHO}$ .

## Appendix F: Supplementary information for non-hermitian Hamiltonian Simulation

Simulating non-Hermitian Hamiltonian dynamics—central to open quantum systems,  $\mathcal{PT}$ -symmetric optics, and gain-loss processes—remains a significant challenge in quantum science. Unlike conventional unitary evolution, non-

Hermitian generators produce norm-non-preserving dynamics and complex spectra, rendering classical simulation especially demanding. Existing quantum algorithms for non-Hermitian channels often require complex ancillary structures. Here, we establish a thread-level parallel protocol, in which modular subtasks can be distributed across multiple quantum processors. This approach systematically enables scalable, parallel simulation of non-Hermitian dynamics, addressing a key bottleneck in both theory and experiment.

Many scientific and engineering tasks require simulating non-unitary evolution. Recent advances, such as linear combination of unitaries (LCU) and dilation methods, use ancillary qubits to enable non-Hermitian Hamiltonian simulation. As a representative protocol, the LCU framework approximates the time-evolution operator as a sum of unitaries:

$$e^{-iHt} \approx \sum_{k=0}^K \gamma_k U_k,$$

with coefficients  $\gamma_k$  and unitary circuits  $U_k$  constructed by, e.g., Taylor or Chebyshev expansion. For example, using the Taylor series expansion up to order  $K$ , we have

$$e^{-iHt} = \sum_{n=0}^N \frac{(-iHt)^n}{n!} + \mathcal{O}\left(\frac{(Ht)^{N+1}}{(N+1)!}\right).$$

By substituting  $H = \sum_{m=1}^M \beta_m V_m$ , where  $\beta_m$  are coefficients and  $V_m$  are unitary operators,

$$e^{-iHt} \approx \sum_{n=0}^N \frac{(-it)^n}{n!} \left( \sum_{m=1}^M \beta_m V_m \right)^n.$$

This results in a form of  $\sum_{k=0}^K \gamma_k U_k$ , with

$$\begin{aligned} U_k &= V_{m_1} V_{m_2} \dots V_{m_n}, \\ \gamma_k &= \frac{(-it)^n}{n!} \beta_{m_1} \beta_{m_2} \dots \beta_{m_n}, \end{aligned}$$

where  $m_1, m_2, \dots, m_n \in [1, M]$ , and  $K = (M^{N+1} - 1)/(M - 1)$ .

For time-dependent or non-Hermitian Hamiltonians, a linear combination of Hamiltonian simulation (LCHS) protocol (see Algorithm 1) is employed [26, 29], which naturally decomposes the evolution into independent unitary tasks suitable for parallel execution across QPU nodes. The method also needs ancilla systems, in our setting, we use cluster QPU, which leads a plug-and-play paradigm.

---

### Algorithm 1 Hamiltonian Simulation for Non-Hermitian Dynamics

---

**Input:**  $A(t) = H(t) + iL(t)$ , initial state  $u_0$ , simulation time  $T$ , precision epsilon.

**Output:** Final state  $u(T)$

- 1: Decompose  $A(t)$  into Hermitian part  $H(t)$  and anti-Hermitian part  $L(t)$ .
  - 2: Define function  $V(t, k) = T \exp(-i \int_0^t (H(s) + kL(s)) ds)$ .
  - 3: Truncate the integral at  $K = \frac{c}{\epsilon}$ , where  $c$  is a constant.
  - 4: Discretize the integral using a trapezoidal rule:
  - 5:  $K = \frac{c}{\epsilon}$  and  $M = \lfloor \frac{2KT}{\epsilon} \rfloor$
  - 6: **for**  $j = 0$  to  $M$  **do**
  - 7:    $k_j = -K + \frac{2jK}{M}$
  - 8:    $w_j = (2 - (\text{if } j == 0 \text{ or } j == M)) \frac{K}{M}$
  - 9:    $c_j = \frac{w_j}{\pi(1+k_j^2)}$
  - 10:    $U_j(t) = T \exp(-i \int_0^t (H(s) + k_j L(s)) ds)$
  - 11:    $v_j = c_j U_j(T) u_0$
  - 12: **end for**
  - 13: Sum up the contributions:
  - 14:  $u(T) = \sum_{j=0}^M v_j$
  - 15: **return**  $u(T)$
-

### 1. Non-Hermitian Hamiltonian

We consider the open-system evolution of a single qubit governed by a time-independent non-Hermitian operator [27, 28]:

$$\mathcal{A}(t) = H(t) - iL(t), \quad H(t) = \sigma_x, \quad L(t) = \mathbb{1} + \sigma_z, \quad (\text{F1})$$

where  $\sigma_{x,y,z}$  denote the Pauli matrices and  $\mathbb{1}$  is the identity operator. The system is initialized in the ground state  $|\psi_0\rangle = (1, 0)$ , and we measure the expectation values of  $\sigma_y$  and  $\sigma_z$ .

This non-Hermitian evolution is experimentally realized on a cluster of  $16 \times 2$ -qubit nodes, with a total of over 1700 experimental runs (approximately 106 per node). The simulation protocol is based on the linear combination of Hamiltonian simulation (LCHS), which approximates the non-unitary evolution using a truncated Cauchy-integral formula, as detailed in Algorithm 1. In addition, numerical simulations are performed using both this methodology and the exact integration of Schrödinger’s equation via a first-order Trotter step.

Symbol	Meaning	Value
$T$	Total evolution time	0.1-1.0 (step 0.1)
$\Delta t$	Trotter step size	0.01
$\varepsilon$	LCHS quadrature step size	0.2
$c$	Constant in cutoff	0.5
$K$	Integration cutoff	$\lfloor c/\varepsilon \rfloor = 2$
$M$	Number of nodes	$\lfloor 2KT/\varepsilon \rfloor$

TABLE S4: Numerical parameters used in the simulation of non-Hermitian Hamiltonian evolution.

Table S4 summarizes the numerical parameters. Here,  $M$  is the number of nodes and thus sets the count of decomposed unitary terms. Employing the single-ancilla estimator, illustrated in Fig. S9, each value of  $T$  requires  $1 + \lfloor 2KT/\varepsilon \rfloor$  individual experiments per observable. We sweep

$$T \in \{0.1, 0.2, \dots, 1.0\},$$

and for each  $T$  record the expectation of a random Hermitian operator  $R$  as well as of  $\sigma_x$  and  $\sigma_z$ . Altogether, measuring any one observable entails

$$\sum_{j=1}^{10} (1 + \lfloor 2KT_j/\varepsilon \rfloor)^2$$

distinct experiments. Although the ideal count is 1770, stepwise truncation in the computation of  $M$  leads to a slight discrepancy, 1747. The reason is one calculating the  $T = 0.6$  case. The truncation on machine made the  $M = 11$ , which makes little influence.

The parameters above are selected via considering the efficiency of experiments and accuracy. Throughout this protocol, the fidelity between the Trotter approximation and the LCHS implementation (with  $\Delta t = 0.01$ ) mostly exceeds 99.2%, confirming the accuracy of our linear-combination approach.

TABLE S5: Summary of  $T$ , nodes  $M$ , and resulting fidelity for each sampling point.

$k$	1	2	3	4	5	6	7	8	9	10
$T$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
$M$	2	4	6	8	10	11	14	16	18	20
Fidelity	0.9999	0.9987	0.9945	0.9870	0.9808	0.9828	0.9929	0.9999	0.9964	0.9871

In Table S6, we report the deviations between the experimentally measured expectation values of  $\sigma_y$  and  $\sigma_z$  and their numerically simulated counterparts, as well as those for a randomly generated Hermitian observable  $R$ , over the course of a non-Hermitian Hamiltonian evolution. The mean absolute deviation of 0.116, quoted in the main text, is obtained by averaging all of the absolute values listed in this table.

TABLE S6: Deviations between experimental measurements and numerical simulations for  $\langle\sigma_y\rangle$ ,  $\langle\sigma_z\rangle$ , and a random Hermitian observable  $R$  during non-Hermitian Hamiltonian evolution.

	1	2	3	4	5	6	7	8	9	10
$\sigma_y$	-0.1208	-0.1901	-0.1898	-0.0930	-0.0456	-0.0188	0.1648	0.2057	0.2123	0.1044
$\sigma_z$	-0.0497	0.0891	0.1104	0.1780	0.2054	0.2645	0.2498	0.1037	0.0895	-0.0411
$R$	0.0167	-0.0563	-0.1193	-0.0824	-0.1006	-0.1625	-0.1040	-0.0245	-0.0385	0.0471

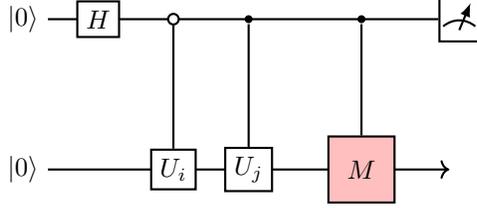


FIG. S9: Ancilla-assisted quantum circuit for measuring overlaps  $\langle U_j^\dagger M U_i \rangle$  required in the LCHS protocol. By varying  $M$ , we obtain the desired real and imaginary parts of the relevant observables.

## 2. Imaginary time evolution

We further investigate the ground-state energy problem for a single-qubit Hamiltonian parameterized by a real coupling  $\gamma \in [0, 2]$ ,

$$H(\gamma) = 2 \mathbb{1} + \gamma \sigma_x, \quad (\text{F2})$$

where  $\mathbb{1}$  is the identity and  $\sigma_x$  is the Pauli- $x$  operator. The analytic ground-state energy is  $E_0(\gamma) = 2 - |\gamma|$ , which serves as an independent benchmark. The ground state can also be found via imaginary-time evolution: the qubit is initialized in  $|\psi_0\rangle = (1, 0)^T$ , propagated in a non-unitary fashion, and the ground-state energy is estimated as  $E(\tau) = \langle \psi(\tau) | H(\gamma) | \psi(\tau) \rangle$ .

Experimentally, this imaginary-time evolution is implemented as a linear combination of unitary evolutions via the truncated Cauchy-integral formula. On the cluster of  $16 \times 2$ -qubit nodes, imaginary-time evolution of  $H(\gamma)$  for a total time  $T = 0.5$  is performed with more than 1331 experimental runs for obtaining single observable (about 83 per node). In addition to the ground-state energy, the expectation values of  $\sigma_x$  and  $\sigma_z$  are also measured for the resulting ground states.

To validate our experimental results, we benchmark three numerical simulations. First is the exact imaginary-time Trotter integration, in which the state is evolved by repeated application of  $\exp[-H(\gamma)\Delta t]$  with a fixed step size  $\Delta t = 0.01$  and  $t = 0.5$ . The second is the LCHS-imaginary protocol, where the non-unitary propagator  $\exp[-H(\gamma)\tau]$  is approximated via a truncated Cauchy integral with quadrature step  $\varepsilon = 0.3$ , shown in Algorithm 1. The third is the exact diagonalization on  $H(\gamma)$ , which yields the ground-state energy, to get a more approximated wavefunction, we use numerical imaginary-time evolution with  $T = 1.5$ .

Symbol	Meaning	Value
$\gamma$	parameter in Hamiltonian	0.0-2.0 (step 0.2)
$T$	total evolution time	0.5
$\Delta t$	Trotter step	0.01
$\varepsilon$	LCHS quadrature step	0.3
$c$	a constant	1
$K$	integration cutoff	$\lceil c/\varepsilon \rceil = 3$
$M$	number of nodes	$\lceil 2KT/\varepsilon \rceil = 10$

TABLE S7: Numerical parameters used in all simulation of imaginary-time evolution.

The simulation parameters are summarized in Table S7. In particular,  $M$  denotes the number of nodes, which in turn determines the total number of decomposed unitary terms,  $U_i, U_j$ . To extract the required expectation values, we design a single-ancilla estimator, illustrated in Fig. S9. Thus, this decomposition requires  $11 \times 11$  ( $(M+1) \times (M+1)$ )

individual experiments. We sweep the parameter  $\gamma$  across eleven uniformly spaced values,

$$\gamma_k = 0.2k, \quad k = 0, 1, \dots, 10,$$

and for each  $\gamma_k$  we record the expectation values of the Hamiltonian  $H(\gamma)$  as well as of  $\sigma_x$  and  $\sigma_z$  on the corresponding ground state. Consequently, measuring any single observable entails

$$11 \times 11 \times 11 = 1331$$

distinct experiments.

At the conclusion of the imaginary-time evolution ( $T = 0.5$ ), the fidelity between the Trotter approximation and the LCHS protocol exceeds 99.5% for most sampling points, thereby confirming the accuracy of the linear-combination approach. Detailed fidelity values are listed in Table S8.

TABLE S8: Fidelity between the Trotter method and the LCHS protocol after imaginary-time evolution ( $T = 0.5$ ) for  $\gamma_k = 0.2k$  (indices  $k = 1, \dots, 10$ ).

	1	2	3	4	5	6	7	8	9	10
fidelity	0.9990	0.9969	0.9950	0.9944	0.9952	0.9967	0.9983	0.9995	0.99998	0.99977

In Table S9, we report the deviations between the experimentally measured expectation values of  $\langle\sigma_x\rangle$  and  $\langle\sigma_z\rangle$  and their numerically simulated counterparts, as well as those for the Hamiltonian observable  $H(\gamma)$ , at the conclusion of an imaginary-time evolution ( $T = 0.5$ ). The overall mean absolute deviation of 0.136, quoted in the main text, is obtained by averaging all of the absolute values listed in this table.

TABLE S9: Deviations between experimental measurements and numerical simulations for  $\langle\sigma_x\rangle$ ,  $\langle\sigma_z\rangle$ , and the Hamiltonian  $H(\gamma)$  after imaginary-time evolution ( $T = 0.5$ ).  $\gamma$  is set as  $0.2k$ , where  $k$  label this horizontal axis.

	1	2	3	4	5	6	7	8	9	10	11
$\sigma_1$	-0.3055	-0.2313	-0.2008	-0.3122	-0.2248	-0.2227	-0.0779	-0.0072	-0.0351	0.0761	0.0498
$\sigma_2$	-0.2459	-0.1625	-0.0745	-0.2681	-0.1279	-0.1384	-0.0746	-0.0652	-0.0296	-0.0692	-0.1333
$\sigma_3$	-0.0108	0.0705	-0.0115	0.0529	0.1826	0.1067	0.1317	0.1807	0.1328	0.1955	0.3062

Finally, we present the experimental implementation on hardware. This is carried out in a two-qubit cluster of the Gemini-Lab desktop spectrometer, as described in the Appendix D.

The overall procedure follows a structure similar to the quantum circuit cutting experiments detailed in Appendix E, and here we focus on a single representative thread.

As in previous sections, a relaxation-based initialization sequence is applied to convert the thermal equilibrium state into a pseudo-pure state that approximates  $|00\rangle$ . The resulting state fidelity reaches 99.5%.

Next, a parameterized two-qubit quantum circuit is used to implement the sub-circuits illustrated in Fig. S9. The circuit consists of arbitrary angle single-qubit rotation gates and a two-qubit controlled-NOT (CNOT) gate. All parameters are optimized using a gradient descent algorithm. Each single-qubit rotation has a duration of approximately 0.3 ms, while the CNOT gate duration is about 2 ms. Theoretical fidelities for both types of gates exceed 99%. Finally, the compiled circuit is executed and the measurement is performed to extract the desired observables.