# Assessing requirements to scale to practical quantum advantage

M. E. Beverland,[1] P. Murali,[1] M. Troyer,[1] K. M. Svore,[1] T. Hoefler,[2]
V. Kliuchnikov,[1] G. H. Low,[1] M. Soeken,[3] A. Sundaram,[1] and A. Vaschillo[1]

[1] *Microsoft Quantum, Redmond, WA 98052, USA*
[2] *ETH Zurich, Department of Computer Science, Zürich, 8006, Switzerland*
[3] *Microsoft Quantum, Zurich, Switzerland*
(Dated: November 19, 2022)

While quantum computers promise to solve some scientifically and commercially valuable problems thought intractable for classical machines, delivering on this promise will require a *large-scale* quantum machine. Understanding the impact of architecture design choices for a scaled quantum stack for specific applications, prior to full realization of the quantum system, is an important open challenge. To this end, we develop a framework for quantum resource estimation, abstracting the layers of the stack, to estimate resources required across these layers for large-scale quantum applications. Using a tool that implements this framework, we assess three scaled quantum applications and find that hundreds of thousands to millions of physical qubits are needed to achieve practical quantum advantage. We identify three qubit parameters, namely size, speed, and controllability, that are critical at scale to rendering these applications practical. A goal of our work is to accelerate progress towards practical quantum advantage by enabling the broader community to explore design choices across the stack, from algorithms to qubits.

## I. TOWARD QUANTUM APPLICATIONS WITH PRACTICAL IMPACT

With rapid progress in quantum computing, there is a shift in focus from scientific demonstrations on a small number of noisy qubits [3, 110] towards large quantum systems solving valuable business problems and resolving long-standing scientific questions that are classically intractable [35]. But just how large does a quantum computer need to be to achieve these forms of *practical quantum advantage*, and how long will such a computation take? Are some qubit technologies better suited than others for solving such problems? What are the best architecture choices across the hardware and software stacks to enable scaled quantum computation?

While many quantum algorithms have asymptotic speedups over their non-quantum counterparts [67], scaling up to problem sizes with a practical quantum advantage [4] will demand high-precision quantum computers. Quantum operations at the physical level are noisy, and so the long computations required for practical quantum advantage necessarily require error correction to achieve fault tolerance. Quantum computers are complex, and the overheads required to ensure fault-tolerant quantum computing will significantly outpace the resources required for fault-tolerant classical computing which is based on small, cheap, and reliable transistors. A fault-tolerant operation on a quantum computer requires orders of magnitude more space – including many transistors for qubit control and readout – and runs with much slower clock speeds than a classical computer. With these overheads, practical quantum advantage will be achieved, albeit only for algorithms with small I/O requirements and superquadratic (ideally exponential) speedups over their classical counterparts [4, 92, 138].

One algorithm with superquadratic speedup for which the cost of error correction is well studied is Shor's period finding algorithm [120]. Estimating the resources required for Shor's algorithm is important for assessing the vulnerability of some of today's public key cryptosystems to future quantum threats. With the fastest quantum hardware operations proposed to date, factoring a 2048-bit integer using Shor's algorithm could be done in minutes with an array of twenty five thousand *perfect, noiseless* qubits. Yet in reality, qubits are noisy and must have error correction

to enable long computation, and so as we will reconfirm later, the implementation cost increases to about a day with tens of millions of qubits [42].

Some of the most compelling quantum algorithms with scientific and commercial interest, however, are those which leverage the ability of a quantum computer to efficiently simulate quantum systems, with applications across chemistry, materials science, condensed matter, and nuclear physics. The exact simulation time of the dynamics of such quantum systems scales exponentially with classical algorithms, but has a favorable polynomial scaling for quantum algorithms [91]. The earliest application of scientific interest may be simulating the dynamics of around one hundred quantum spins in a quantum magnet [31]. The earliest commercially relevant applications will likely be quantum simulations of chemistry and materials science problems, where a quantum-accelerated elucidation of catalytic reaction mechanisms has applications to fertilizer production [113], carbon fixation [140], among many other problems [6, 34, 72, 140].

We propose a framework to understand the requirements of applications promising quantum advantage, determine their practicality, and assess changes in the underlying architecture to accelerate their implementation. The framework is composed of layers of abstractions capable of estimating the performance and required resources of quantum applications and algorithms on current and future quantum architecture designs. Starting with an algorithm implemented in a high-level programming language, the framework compiles the algorithm onto the architecture of a specified quantum device and models the resource requirements at the architecture level, considering qubit layout, allowed operations, and other design parameters. We specify models for each layer and implement a tool, the Azure Quantum Resource Estimator, that calculates based on those models an estimate of the resources required to implement a given quantum algorithm on a given architecture with an underlying qubit technology. We focus on *digital* quantum computers using gates and measurements for computation, rather than analog quantum simulators and annealers [31]; however, similar layers of abstraction may also be valuable to analyze these other systems.

We analyze the resources that would be required to implement three specific high-impact applications on various quantum architectures. The first considers the *quantum dynamics* of a simple quantum magnet, the so-called two-dimensional (2D) transverse field Ising model, where we consider a parameter regime that is on the boundary of what can be done by classical computation [107]. The second is for *quantum chemistry* where we analyze the activation energy of a catalyst for carbon fixation [140]. Lastly, we analyze *factoring* a large integer with Shor's algorithm. While we do not foresee running this application in practice, it is arguably the best-studied quantum algorithm and can help provide security requirements of our classical cryptosystems. Our analysis reveals that practical quantum advantage for scientific and commercial problems requires hundreds of thousands to millions of qubits which satisfy the following requirements.

**Requirements for scale.**— To achieve practical quantum advantage, quantum computers will require an underlying qubit technology that at scale is:

- *Controllable:* Practical quantum error correction requires reliable control of more than a million well-connected qubits, with parallel operations that fail in under one part in a thousand.

- *Fast:* To achieve a practical runtime of one month or less, while targeting a physical qubit count of around one million, operations will need to be performed in under a microsecond.

- *Small:* Scaling to a million and more qubits constrains the size of the qubit to tens of microns in diameter; this size is determined to avoid the complexity of coherent high-bandwidth quantum interconnects between qubits on different modules.

We see these requirements as necessary to scale to practical quantum advantage, and as valuable additions, for scale, to the criteria for building a quantum computer proposed by DiVincenzo in

the year 2000 [36]. They stem from incorporating research developments spanning the last two decades, insights across algorithms to qubits, and an empirical end-to-end resource analysis of several applications and qubit parameter settings. We hope they reveal new considerations and valuable insights in considering how to engineer and architect for scale.

## II. RESOURCE ESTIMATION FOR CLASSICAL AND QUANTUM COMPUTING

Estimating the resources required for computation, also known as performance modeling, plays a central role in classical high-performance computing. Elaborate techniques have been developed to understand application scaling characteristics, their portability across architectures, and resource consumption. These models consider some resources as a measure of cost, such as time, energy, number of instructions, and more recently, data movement. In this section, we first reflect on the mature field of classical computational resource estimation, and then leverage this perspective to present a framework for quantum resource estimation.

In the classical computing stack, the program is initially expressed in a high-level language at the top and is sequentially re-expressed in more and more explicit representations down the stack; see the left column of Fig. 1. We can consider each layer in the stack as having a language with a specific *instruction set* of allowed instructions, which can be broken down in terms of simpler instructions in the layer below, or combined together to form more abstract instructions for the layer above. For example the program could be expressed in C++, which is compiled into an intermediate representation (IR) for optimizations, e.g., LLVM's static single assignment form (SSA) intermediate representation [84]. The IR is then lowered to assembly code in a specific instruction-set architecture (ISA) that is the interface to the classical processor. Such assembly code could consist of x86 [60], ARM [45], or RISCV [114] instruction sets. Some complex instruction
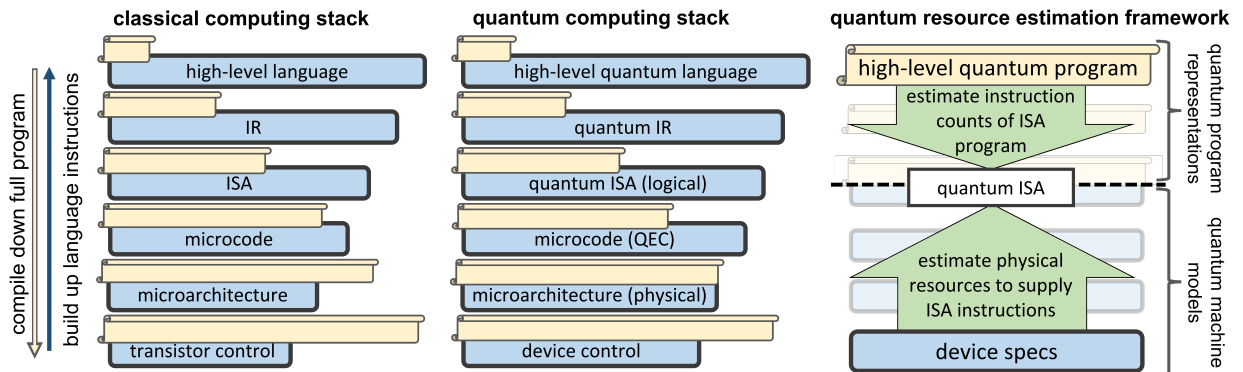


FIG. 1. The stacks for classical (left) and quantum (center) computing. At the top of the stack the program is expressed in a natural, high-level language such as C/C++ (classical) or Q# (quantum). This program is re-expressed as different program representations (manilla scrolls) in a sequence of languages (blue boxes) which are represented here as layers, with each layer being a lower-level language than that above, down to the physical signals controlling the device. A complementary viewpoint is that the instruction set available to the language at one layer can be combined to form a more elaborate instruction set for the layer above. The scrolls are shortest at the top, symbolizing that program representations are simplest when expressed in higher-level languages, while the blue boxes are smallest at the base of the stack, symbolizing that instruction sets are simplest at the lower levels. (right) Our quantum resource modeling framework is a modular representation of the layers of the quantum computing stack, which collects together the upper layers as *quantum program representations*, and the lower layers as *quantum machine models*, with the quantum ISA connecting both.

set computers such as Intel's Haswell translate the x86 code into a microcode representation for execution by the arithmetic logical units (ALUs). The microarchitecture then executes those micro-instructions and eventually translates them to electrical signals to be implemented by transistors. For some processors, microcode and microarchitecture are combined, but we illustrate them in separate layers here to draw parallels with the quantum stack.

The goal of resource estimation is to model the stack in order to estimate some physical resources such as the run time required to implement a high-level program. One approach could be to consider a chain of re-expressions of the program from the very top of the stack all the way to the bottom and then read off the physical resources needed. Near the top of the stack this is an appealing approach as it is very natural to think of a program being compiled down into a sequence of lower-level representations. Towards the bottom of the stack however this approach becomes very difficult because of the intricacies of the microarchitecture, runtime dependencies of the program execution, and the overall length of the representation itself. However, the instruction sets within which the program is expressed typically become simpler for lower layers, which can be exploited to model resources more conveniently. The strategy is then to choose a layer towards the middle of the stack, and model the layers above by re-expressing the high-level program down to the instruction set at this layer, and estimate the physical resources required to implement those instructions by building them up from the more basic instructions in the layers below.

There is of course a choice about where to draw the line between the top and the bottom in this view of the stack, and the best choice can depend on the precise goals of the model. For most purposes, the ISA layer has emerged as the natural interface because it provides an implementation-independent interface that can be used to ground modeling efforts. ISAs such as x86, ARM and RISCV represent program executions across a wide variety of classical processors. Performance modeling frameworks can leverage this powerful abstraction to model the resources of the full stack by first establishing the requirements of the program when expressed as instructions the ISA level, and then identifying the resources required at the device level to implement the ISA-level program.

The quantum computing stack is less established and abstractions are still in their definition phase. Some aspects are similar, for example while a quantum program includes instructions which are executed as quantum operations, the quantum program itself is just a sequence of instructions similar to a classical program. However, in stark contrast with the classical stack, which is built upon extremely reliable transistors, the quantum stack must account for the unavoidable fact that all forms of physical qubit implementations are inherently noisy [110]. To execute a quantum application successfully, quantum error correction (QEC) [47] must be used to build logical qubits that can be used to store and manipulate quantum information better than raw physical qubits. This QEC capability is central to scalable quantum computers, but the costs are formidable, often multiplying the number of qubits needed by a factor of thousands, and runtimes by a factor of hundreds. Therefore, the specific QEC approach used impacts design decisions across the quantum computing stack, both at the hardware layers which must be capable of implementing it, and at the software layers which must compile to logical operations compatible with it.

Prior works that address the design of a flexible quantum stack and resource estimation span vastly different viewpoints. For example, prior quantum architecture and compilation proposals [30, 40, 102, 119] develop stacks for near-term small quantum computers, but do not incorporate QEC. These works either ignore noise [40], or reduce its impact with error mitigation techniques [30] which unfortunately do not scale to larger systems [112]. On the other hand, works comparing resources for multiple QEC codes and logical operations, such as [10, 13, 23, 24, 90], typically address lower layers of the stack without incorporating the software layers. There is a growing number of informative end-to-end resource analyses, but typically these single out very specific algorithms and hardware and make very different assumptions across the stack, making direct comparisons of different approaches challenging [1, 7, 26, 38, 42, 49, 73, 115, 131].

## A.   A framework for quantum resource estimation

We present a quantum resource estimation framework that aims to capture the essential aspects of a full-stack quantum computer with the flexibility to allow direct and easy comparisons across a range of algorithms, software techniques, QEC codes, hardware specifications and other aspects. To form the skeleton of our framework, we propose the stack shown in Fig. 1 (center) for large-scale quantum computers. This framework builds upon the previous work mentioned above, but seeks to clearly delineate the different layers, clarify the functions of each layer and define how they inter-operate to attain the computation capabilities required for large-scale applications. We find that it is possible to draw broad parallels between the classical and quantum computing stacks, allowing us to leverage established resource modeling strategies from classical performance modeling.

In what follows we describe the layers in this stack, briefly stating the main characteristics of each layer, and listing some known examples from the quantum computing literature which naturally fit into it. As we will see, the layers are modular, allowing them to be designed either collectively or independently, provided that the interfaces between layers remain consistent. Consistency of the interfaces is ensured by providing explicit maps: compilation algorithms for program representations down to the ISA from above, and explicit constructions of instruction sets from those below up to the ISA from below. These maps are themselves a crucial part of the stack and the modeling framework. We start at the top of the stack, where the layers are conceptually similar to the classical case.

**High-level quantum language.**— This layer hosts the explicit representation of an algorithm as a program expressed in a high-level programming language. The language should enable the expression of any quantum algorithm. At this level of abstraction, users need not know the requirements of quantum error correction, and may assume their program and its operations are fault tolerant. Languages ideally include support for loops, functions, rich type systems, debugging and other functions that are common in classical computing in order to allow quantum applications to be developed and tested efficiently. Similarly, supporting hybrid quantum-classical computation which allow applications to use phases of classical and quantum operations seamlessly is valuable for practical applications. Examples of languages include Q# [132], Quipper [50], Scaffold [2], QWire [106], Quil [123] and others.

**Quantum IR.**— To execute a high-level quantum program, that program's instructions need to be *compiled* and expressed in terms of quantum ISA instructions. It is appealing to support executions of programs expressed in a diverse set of high-level languages across a range of hardware back ends that may have different quantum ISAs. In analogy with classical computing, we propose the use of a quantum intermediate representation (quantum IR) that expresses high-level operations in a language-agnostic and ISA-agnostic manner [61, 95, 111]. A compiler then includes three components — the front end, the optimization layer and the back end. The front end component translates high-level language instructions to the quantum IR. The back end component translates the quantum IR into the quantum ISA operations. The middle optimization layer performs transformations on the quantum IR to reduce resource requirements [61]. To support a new quantum language or quantum ISA, we simply need to add the corresponding front end or back end components, rather than reinventing the whole compiler. An example of a quantum IR is the aptly named QIR [111], which includes support for expressing logical operations, measurements and other quantum constructs that are required for applications.

The parallels of these two upper layers with the classical stack can be leveraged. For instance, later we consider high-level quantum programs expressed in Rust, and the classical IR LLVM has been applied to quantum compilation tasks Ref. [64, 96]. Next we address the layers at the bottom of the stack and describe how their instruction sets build up to form the instruction set of the logical qubits at the quantum ISA layer (which the quantum IR compiles down to).

**Device control.**— There is a wide range of hardware approaches to build the quantum device at the base of the stack, including those based on photonic [14, 21, 81], trapped atom [12, 62], trapped ion [29, 108], Majorana [71, 74, 117], spin [54, 65, 70], superconducting transmon [3, 57, 76, 82, 118, 127], and electro-acoustic [25, 51] platforms. Quantum devices themselves are not sufficient for computation, we require classical computing resources to bring up, calibrate and control the devices as well as the ability to transmit and process information from quantum measurements. The nature of how qubits are stored within the device, the range of control available, and the nature of that control are all very device-dependent and form the instruction set at the device level – the constraints on ion traps will be different than the constraints on transmon qubits, for example, simply because their physical realizations are fundamentally different. It is at this level that resources are typically most forthright - one may consider run time, number of qubits, control bandwidth, device size and other cost metrics.

**Microarchitecture (physical).**— In this layer, low-level details of the qubit device design are abstracted away, leading to a viewpoint consisting of a set of abstract physical qubits, along with a discrete set of physical instructions such as single-qubit measurements and CNOT gates. These instructions can be characterized by simple properties such as their action on the physical qubits, how long they take to apply, and their probability of failure (their so-called *error rate*). This abstraction allows microarchitectural components to be reused across a variety of qubit device implementations, while flexibly supporting the needs of different QEC schemes. Each of the operations in the physical instruction set is built from operations available at the device control level (for example, in ion traps a CNOT may be implemented by a sequence of Mølmer-Sørensen [99] gates, while in spin qubits it may be implemented by a sequence of exchange interactions [54]). In planar device platforms like superconducting qubits or Majorana qubits, the ability to support such device-level operations can lead to connectivity restrictions such as 2D nearest-neighbor connectivity at this level.

**Microcode (QEC).**— The microcode instruction set consists of the primitives required to implement a QEC strategy, which uses a QEC code to form reliable logical qubits using noisy physical qubits. These logical qubits support logical operations (which form the instruction set for the quantum ISA level above) that have better error rates than raw physical operations. For example, if the QEC strategy uses the surface code [16, 75], the primitives which form the instruction set will include *syndrome measurement circuits* expressed in the physical instruction set of the microcode layer below. For the surface code, syndrome measurements provide the information required to diagnose and correct faults at the physical level by measuring stabilizers, but can also induce logical operations by deforming the code [58, 89].

**Quantum ISA (logical).**— This layer forms the interface between the software and hardware layers. It abstracts the details of how QEC is implemented in the layer below, retaining only a set of fault-tolerant logical operations as its instruction set. It is crucial that the instruction set is *universal*, meaning it is complete for quantum computing. In this formulation of the stack, it is the map from the microcode level to the quantum ISA level that actually implements the error correction, using a classical algorithm known as a *decoder* to identify and correct faults while implementing specific logical operations. Clearly, the instructions in this layer must be chosen such that they can implemented by the underlying hardware; that is, when choosing an ISA, it is important to consider the constraints of the system architecture such as connectivity of logical qubit representations and ability to construct the logical operations with known sequence of micro-operations and physical qubit control operations. At the same time, the ISA operations should be chosen such that the software layers above can efficiently express common application operations using it. The planar quantum ISA is an example we will discuss later which arises from the logical operations of two different QEC approaches, one based on surface codes and the other on Hastings-Haah codes. For instance, the planar quantum ISA instruction set includes single-qubit

initializations, single-qubit measurements and multi-qubit measurements, which are all examples of Clifford operations. Clifford operations by themselves are not universal [46], but when combined with a fault-tolerant T state initialization, the resulting set becomes universal.

**Estimating resources.**— Our goal here has been to define the layers of the quantum stack at a sufficiently high-level of abstraction, so that a broad range of known and yet-to-be discovered approaches can be readily incorporated. This forms the basis of a *modular* framework for quantum resource modeling, in which layers can be considered together or independently. We find it convenient to – similar to the classical case – split the quantum stack in two at the quantum ISA and analyze the upper layers from the top down in terms of programs, and the lower layers from bottom up in terms of instruction sets; see Fig. 1 (right). To use the framework to make explicit estimates, one models a consistent instantiation of each layer and each of the maps between layers with an explicit example. Later, we will specify examples of these layers and maps in more detail to answer questions regarding system requirements for scaling.

This framework enables the estimation of resources (such as the number of qubits, the run time and the power consumption) which would be required to implement a given quantum algorithm using a given qubit technology and with a fixed set of architectural choices. Within the framework, architectural choices can be viewed as explicit specifications for each layer and map in the hardware and software parts of the quantum stack. There is no single answer to questions on architectural design trade-offs, with different choices resulting in different estimation results which can be explored to make informed decisions as researchers and engineers work to scale up a complete system. One can, for example, trade off more qubits against shorter run times, or trade off faster qubit gate operations against lower fidelities. One can also use the framework much more broadly, to compare new error correction proposals, to evaluate the potential for quantum advantage for new or optimized quantum algorithms, and to motivate future research directions by identifying bottlenecks in the stack which generate large resource contributions.

The modular nature of the framework allows for future modifications and extensions. For some purposes it may be desirable to combine some layers together, and to split other layers into sub-layers. In this work our focus is on resource estimates for the quantum accelerated parts of the computation, however more layers can be added to model the hybrid aspects of the computation.

## B. Fault-tolerant design considerations

Before discussing specific implementations of the quantum resource estimation framework, we address several significant design considerations that arise from the need for QEC. Typically, a QEC approach is centered around a QEC code family which stores the logical qubits. A wide range of QEC codes have been discovered over the last three decades with wide ranging properties, including small codes such as the Shor [121] and Steane [125] codes, 2D local codes such as surface codes [16, 75], color codes [15], Hastings-Haah codes [55] and Bacon-Shor codes [5, 122], and more exotic families of positive-rate low-density parity check codes [52, 88, 104, 135]. The choice of QEC code has impact up and down the stack since different codes can place different implementation demands on the physical instructions provided by the microarchitecture level below, while also exposing different logical operations to the quantum ISA level above.

Here we focus on the requirements that QEC imposes on the physical operations at the microarchitecture level of the quantum stack, corresponding to the setting known as *circuit-level* analysis in the QEC literature. It is common in this setting to assume QEC is implemented with Clifford operations (state preparation, Clifford unitary gates, and Pauli measurement) for single qubits and across pairs of qubits coupled by a specified *connectivity*, and that those operations fail independently with a uniform error rate (see Appendix A). An important metric for a QEC approach in

this setting is its *threshold*, which specifies the maximum error rate that it can tolerate. In what follows, we leverage known threshold results for a range of QEC schemes to extract a number of key design parameters pertaining to error rates, parallelism and connectivity.

The theoretical upper limit for any threshold is not known, but the highest threshold discovered to date is about 3% using the scheme described in Ref. [80], and a range of schemes [11, 103, 137, 141] have been found with thresholds approaching 1%. To avoid prohibitive QEC overheads, error rates at least an order of magnitude smaller than the threshold are much preferred. For this reason, we require physical error rates on Clifford operations below 0.1%.

The aforementioned schemes obtain high threshold values only in the setting where operations can be applied in parallel, which may pose a significant hardware challenge for some platforms such as trapped ions [142]. However, the necessity of parallel operations has long been believed necessary for any QEC scheme to exhibit a non-zero threshold [48, 126]. We therefore require that physical operations can be applied in parallel.

Next we consider connectivity. Probably the most practically relevant classes of connectivity correspond to arrays of qubits with neighboring connections in one or two dimensions, which we refer to as *1D* and *2D connectivity* respectively. For connectivity with any fixed dimension, the thresholds of many of the more exotic code families are believed to drop to zero [20, 33]. With 2D connectivity the thresholds of some QEC schemes based on concatenating small codes are reduced [133]. Fortunately, many 2D local codes (including the surface code [141], the Hastings-Haah code [103] and the color code [11]) are naturally implementable with 2D connectivity and achieve high thresholds approaching 1% in this setting despite the connectivity restrictions. Further restricting to 1D connectivity leaves very few code families known to have finite thresholds, and those that do are around 0.01% and below [66, 129, 134]. Moreover, 1D connectivity is vulnerable to hardware defects that some qubits permanently unusable. To avoid such potentially fatal conditions, we require that qubits are well connected, for example with 2D connectivity but not 1D connectivity.

There are a number of important subtleties that we have ignored so far. First, the simple uniform discrete error model is not a true refection of most systems. Including some realistic aspects of noise, such as asymmetries in the probabilities of different error channels [25, 139] and richer information from noisy measurements [105] can improve error correction performance. On the other hand, including other realistic aspects, such as correlated errors [77] and permanent hardware defects which must be worked around to implement QEC [124, 130] can worsen error correction performance. To partially allow for these subtleties, we require that the *worst* error rates for Clifford operations that are employed by the QEC approach are below 0.1%. Next, we have required the qubits are *well connected* to enable practical QEC, and make it clear that 2D connectivity falls into this class, while 1D connectivity does not. However, we leave the formulation of a more complete definition which applies to more general connectives to future work.

In summary, we require physical qubits to be well connected such as having 2D connectivity, with the ability to perform parallel operations with error rates below one part in a thousand. This holds for all examples we consider in this paper, and forms the basis of one of our criteria for scale.

## III. ESTIMATING RESOURCES FOR THREE QUANTUM APPLICATIONS

In this section we introduce an implementation of the quantum resource modeling framework, hereafter referred to as 'the tool', which is publicly available in Microsoft Azure Quantum as the first version of the Azure Quantum Resource Estimator [97]. Our primary goal in this work is to identify architecture features which will be crucial to achieve practical quantum computing advantage. To do so, we estimate the resources for select applications with the potential for practical quantum

advantage (Section III A), using qubit parameters which are relevant for a number of prominent qubit technologies. We fix a set of consistent architectural options for the quantum stack based on established approaches from the literature.

Following our framework, estimates obtained using the tool can be viewed as consisting of two broad components as shown in Fig. 2. The first component (Section III B) involves an explicit front end compilation from a high-level language to a quantum IR, which is fed to the tool. The tool then models the back end compilation down to a quantum ISA, yielding ISA-level resource estimates including the number of logical qubits and the number of logical time steps. In the second component (Section III C), specific qubit parameters are fed to the tool which performs a bottom-up estimation of the hardware resources required to implement ISA instructions by modeling error correction with specific QEC codes. In Section III D we combine these components to obtain concrete resource estimates.

In the appendices we provide a more detailed and self-contained description of the material in this section, including our modeling assumptions. Our examples are available as samples at [98].

## A. Quantum applications

The first of our three example applications is a *quantum dynamics* simulation that is one of the smallest scientifically interesting problems that is out of reach for classical computation. Specifically, we consider a 2D transverse-field Ising model with 100 quantum spins, propagated for ten time steps using a fourth-order Trotter algorithm as described in Ref. [107]. This is among the simplest models to exhibit a quantum phase transition. Its dynamics are representative of strongly entangled quantum systems and as such is believed to be classically intractable [6].

Our second example is a *quantum chemistry* application to calculate the energy of a ruthenium-based catalyst for carbon fixation which could have implications for reversing the effects of global
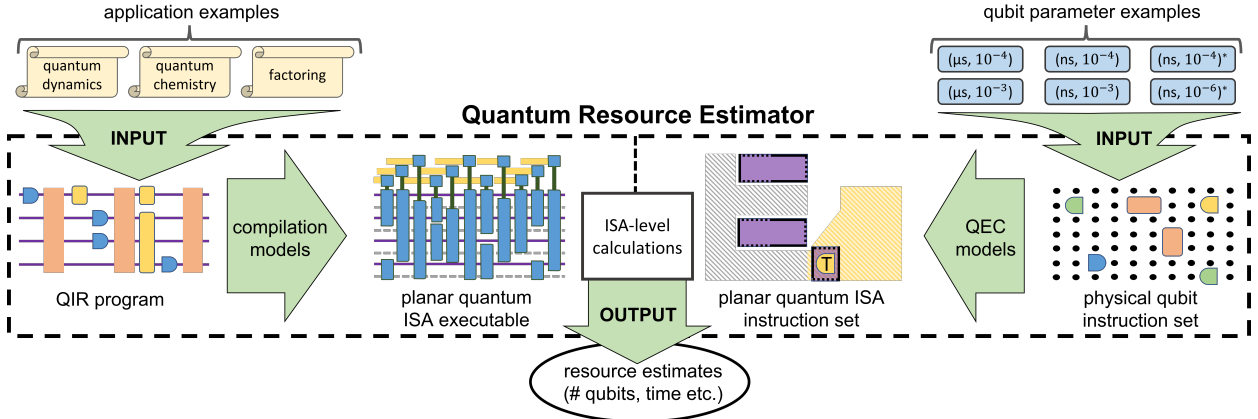


FIG. 2. A sketch of how the tool estimates resources for the three example applications and six qubit parameter examples we consider. Note the stack is represented left-to-right rather than top-to-bottom. Applications, addressed by a high-level quantum program, are first translated explicitly to a QIR program which is input into the tool, where its compilation down to an ISA-level executable is modeled. All our examples flow to the same planar quantum ISA, which has an instruction set consisting of logical surface code operations. On the hardware side, we input physical qubit parameters to the tool including the physical qubit instruction set and the times and error rates of those instructions. The tool then models how these noisy qubits are used to build up protected logical qubits and to fault-tolerantly implement the planar quantum ISA instructions using configurable QEC models. The tool outputs resource estimates such as the number of physical qubits and time required to run the application.

warming. More specifically, to estimate the energy of Complex XVIII in Ref. [140] to chemical accuracy of 1 mHa using the so-called 'double-factorized qubitization' algorithm described in Ref. [140].

Our third application is *factoring*, specifically to identify the pair of prime factors of a 2048-bit integer. Solving this problem would form the basis of an attack on widely used RSA-based encryption schemes. In our resource estimation model, we use a high-level program to implement Shor's factoring algorithm based on that described in Ref. [42].

To achieve the desired solution accuracy for each application, we also set a minimum probability that the overall algorithm succeeds, which we call the *algorithm execution accuracy*. For the quantum dynamics and quantum chemistry examples we assume that algorithm execution accuracies of 0.999 and 0.99 will be required. However factoring admits a lower algorithm execution accuracy since it is easy to check if a proposed solution is correct (by multiplying the purported prime factors to identify if they recover the original integer), and so we take it to be $2/3$ in this case.

## B. Requirements at the quantum ISA level

Here we describe how we obtain ISA-level resource estimates needed to target our three example applications. To identify the physical-level costs of implementing an ISA-level executable, we will not only need to know parameters capturing the size of the program, but also the required quality of the ISA operations needed to achieve the desired algorithm execution accuracy. This is because different quality requirements will result in different error correction solutions being preferred.

As a starting point, we assume explicit high-level implementations of the quantum-accelerated components of each of the three target applications, using Rust for quantum chemistry and factoring and Q# for quantum dynamics[1]. First we use an explicit front end compilation to a specific quantum IR, namely QIR [111] using the Q# compiler or the Rust compiler based on the application. The tool then models the back end stage of compilation which translates QIR into operations supported by the quantum ISA.

In this work we use the *planar quantum ISA*, which is based on the logical operations of the most well-established QEC code, namely the surface code [16, 75], and also applies to the recently developed Hastings-Haah codes [55, 103]. Logical qubits are stored in patches laid out in a 2D plane, along with ancilla regions. We sketch some key features of the planar quantum ISA instructions in Fig. 2, and provide a full specification in Fig. 6 in Appendix B. Only Clifford operations involving patches (purple) that are connected by ancilla regions (grey) are permitted, and two such operations can only be performed simultaneously only if the ancilla regions they require are disjoint. Non-Clifford T states are produced in dedicated T factories (yellow).

There are two primary conceptual challenges that need to be overcome to implement the back end compilation phase. Firstly, the QIR instruction set contains fine-angle rotation unitaries, which must be approximated with sufficient accuracy by sequences of operations from the planar quantum ISA, a process known as *synthesis*. Secondly, one requires methods to layout qubits and perform communication in order to map from the QIR, which has no connectivity restrictions, to the planar quantum ISA, which has geometric constraints. Standard layout and communication methods can be applied to address this second aspect of compilation, including SWAP-based schemes where qubits are moved by swapping neighbors, as well as teleportation-based schemes where qubits are moved using quantum teleportation. In Appendix D, we provide an alternative back end compilation scheme which we refer to as Parallel Synthesis Sequential Pauli Computation (PSSPC), which is based on a combination of favorable features from existing approaches [8, 22, 78, 89]. The

---

[1] These implementations are available as samples at [98]

tool first takes a trace of the QIR program and then uses accurate formulas to calculate the planar quantum ISA resource requirements that would result from back end compiling with PSSPC.

Table I presents the ISA-level resource requirements of our three example applications, which are seen to vary significantly across the applications. Our quantum dynamics example, selected as being among the smallest applications of scientific interest, requires only 200 logical qubits after mapping to the planar quantum ISA. However, even this application requires logical qubits with error rates that are below $10^{-11}$, which puts it beyond the capabilities of current noisy quantum hardware. Applications like chemistry and factoring require several thousand logical qubits, with each logical qubit having error rates in the range of $10^{-15}$ to $10^{-18}$ and T states with error rates in the range of $10^{-12}$ to $10^{-15}$. These applications require mature QEC implementations and large-scale, reliable quantum hardware which we discuss in the next subsection.

| application | algorithm execution accuracy $1 - \epsilon$ | quantum executable parameters | | | quality requirements | |
|---|---|---|---|---|---|---|
| | | $Q$ | $C_{\min}$ | $M$ | max $P$ | max $P_T$ |
| quantum dynamics | 0.999 | 230 | $1.5 \cdot 10^5$ | $2.4 \cdot 10^6$ | $9.7 \cdot 10^{-12}$ | $1.4 \cdot 10^{-10}$ |
| quantum chemistry | 0.99 | 2740 | $4.1 \cdot 10^{11}$ | $5.4 \cdot 10^{11}$ | $3.0 \cdot 10^{-17}$ | $6.1 \cdot 10^{-15}$ |
| factoring | 0.667 | 25481 | $1.2 \cdot 10^{10}$ | $1.5 \cdot 10^{10}$ | $3.5 \cdot 10^{-16}$ | $7.4 \cdot 10^{-12}$ |

TABLE I. Estimates of the ISA-level requirements to implement our applications. The application is compiled down from a high-level application program to a program expressed in the quantum ISA, with parameters as shown. These requirements include the number of logical qubits $Q$, the minimum number of logical time steps (also called logical time steps) $C_{\min}$ and the number of T states $M$ consumed by the program. To ensure that the algorithm fails with probability at most $\epsilon$, the tool also estimates the maximum allowed error rate for each logical qubit $P$ and the maximum allowed error rate for the distilled T states $P_T$.

## C. Requirements at the device level

To understand the high-level impact of physical qubit parameters on application resource estimates, we configure the tool to model the underlying qubit technology relatively abstractly in terms of operation times and error rates, and consider parameter configurations which are relevant for a range of prominent hardware approaches. These features are captured at the microarchitecture level of the quantum stack. Along with the operation time and error rate, the tool takes as input one of two physical operation sets: *gate-based instructions* (including CNOTs and single-qubit measurements) relevant for technologies such as superconducting qubits, trapped ion qubits etc., and *Majorana instructions* (including non-destructive two-qubit Pauli measurements) which are relevant for Majorana qubits [71, 109]. Both of these instruction sets have 2D connectivity and allow the parallel application of disjoint operations.

As listed in Table II, we consider six qubit parameter examples, each labelled with the unit of its operation time and error rate. In Fig. 2, we use an asterisk to differentiate the two examples which use Majorana-based instructions rather than gate-based instructions. With operation times in the microsecond range, the ($\mu$s, $10^{-3}$) and the ($\mu$s, $10^{-4}$) qubits are relevant for trapped ions [29], while the (ns, $10^{-3}$) and (ns, $10^{-4}$) qubits are more relevant for superconducting transmon qubits [82] or spin qubits [70]. We also include the (ns, $10^{-4}$) and (ns, $10^{-6}$) Majorana qubit examples, to account for future topological qubits based on Majorana zero modes [71]. These latter qubits are also expected to have operations in the nanosecond regime, and owing to topological protection in the hardware, they also have the potential for higher fidelities at the physical level. In Appendix A we provide more explanation of these qubit parameter examples and provide the full specifications of the aforementioned instruction sets in Fig. 5.

| qubit parameter examples | | operation times | | error rates | |
|---|---|---|---|---|---|
| | | gate | measurement | Clifford | non-Clifford |
| $(\boldsymbol{\mu}\textbf{s},\ \textbf{10}^{-3})$ qubit | | 100 $\mu$s | 100 $\mu$s | $10^{-3}$ | $10^{-6}$ |
| $(\boldsymbol{\mu}\textbf{s},\ \textbf{10}^{-4})$ qubit | | 100 $\mu$s | 100 $\mu$s | $10^{-4}$ | $10^{-6}$ |
| $(\textbf{ns},\ \textbf{10}^{-3})$ qubit | | 50 ns | 100 ns | $10^{-3}$ | $10^{-3}$ |
| $(\textbf{ns},\ \textbf{10}^{-4})$ qubit | | 50 ns | 100 ns | $10^{-4}$ | $10^{-4}$ |
| $(\textbf{ns},\ \textbf{10}^{-4})$ Majorana qubit | | 100 ns | 100 ns | $10^{-4}$ | 0.05 |
| $(\textbf{ns},\ \textbf{10}^{-6})$ Majorana qubit | | 100 ns | 100 ns | $10^{-6}$ | 0.01 |

TABLE II. We consider a set of examples of qubit parameters which represent various regimes of interest. Each example is labelled with the units of its operation times (either $\mu$s or ns), and the limiting error rate of its Clifford operations. The parameter values and instruction set chosen for each example is informed by published proposals for various hardware approaches. We assume that the first four examples have gate-based instruction sets, while the last two examples have Majorana instruction sets.

In this work, we assume that logical qubits are encoded in patches (purple in the ISA instructions of Fig. 2) of either the surface code [16, 75] or the recently discovered Hastings-Haah code [55]. We make this restriction primarily because these codes have high fault-tolerance thresholds, have relatively well-understood logical operations, and can be implemented using physical qubits with 2D connectivity. We model QEC to relate properties at the ISA (logical) level to properties at the microarchitecture (physical) level. For both the surface code and Hastings-Haah code, error suppression can be tuned by scaling the patch's code distance $d$. For fixed qubit parameters, increasing the distance costs quadratically more physical qubits $n(d)$ and has a linearly longer logical time step $\tau(d)$, but reduces the probability $P(d)$ of a logical qubit failing during a logical time step exponentially. A number of fault-tolerant logical Clifford operations have essentially no time cost in these codes, including state initialization, readout and Pauli gates. On the other hand, logical qubit movement and multi-qubit Pauli measurements are implemented via lattice surgery with the use of ancilla regions (grey in Fig. 2) and require a logical time step. For counting resources, it is convenient to partition the physical qubits into tiles, where each tile contains $n(d)$ physical qubits, and can either be used to encode a logical qubit in a code patch, or can form part of an ancilla region. The total number of physical qubits used for code patches and ancilla regions is then $Q \cdot n(d)$, where $Q$ is the number of tiles (which we also refer to as the number of logical qubits). See Appendix B for details.

For the ISA-level instruction set to be universal, we also require non-Clifford operations, such as a logical T gate. This is achieved via the generation of T states (also known as magic states) using T factories (yellow in Fig. 2), which typically requires many more resources than other quantum ISA operations. These T factories typically involve a sequence of rounds of distillation, where each round takes in many noisy T states encoded in a smaller distance code, processes them using a distillation unit, and outputs fewer less noisy T states encoded in a larger distance code, with the number of rounds, distillation units, and distances all being parameters which can be varied. There are therefore many options when selecting a T factory. We use the symbol $\mathcal{D}$ to label a particular T factory choice which requires $n(\mathcal{D})$ physical qubits, and takes a time $\tau(\mathcal{D})$ to produce a T state with error rate $P_T(\mathcal{D})$. We assume T factories that are space and time efficient implementations of the well-known 15-to-1 distillation circuit; see Appendix C for details.

For each application, the code distance $d$ and T factory $\mathcal{D}$ must be chosen to ensure that the ISA requirements from Table I are satisfied, namely that $P(d) \leq \max P$ and $P_T(\mathcal{D}) \leq \max P_T$, but with minimal cost. As an example, Table III shows the resource costs to satisfy the ISA-level requirements for the factoring application across different physical qubit parameters. We can see that the number of physical qubits and the duration of a logical time step reduce as physical error rates improve. For this example application, improving physical error rates from

$10^{-3}$ to $10^{-4}$ can reduce the number of qubits by a factor of four, and can halve the logical time step. Comparing qubits with microsecond and nanosecond operating times, we see the impact of physical operation time at the ISA level — logical time steps are several orders of magnitude higher than physical operating times, underscoring the importance of fast physical operations. T factories incur significant physical overheads, requiring several thousand physical qubits and only producing new T states once every 10 to 15 logical time steps, motivating the need for multiple parallel T factories to meet application T state demand.

| qubit parameter examples | QEC code selected | logical qubit parameters for $P(d) \leq 3.5 \cdot 10^{-16}$ | | | T factory parameters for $P_T(\mathcal{D}) \leq 7.4 \cdot 10^{-12}$ | |
|---|---|---|---|---|---|---|
| | | distance $d$ | # qubits $n(d)$ | logical time step $\tau(d)$ | # qubits $n(\mathcal{D})$ | duration $\tau(\mathcal{D})$ |
| ($\mu$s, $10^{-3}$) qubit | surface | 27 | 1458 | 16 ms | 17640 | 163 ms |
| ($\mu$s, $10^{-4}$) qubit | surface | 13 | 338 | 7 ms | 4840 | 85 ms |
| (ns, $10^{-3}$) qubit | surface | 27 | 1458 | 10 $\mu$s | 33320 | 128 $\mu$s |
| (ns, $10^{-4}$) qubit | surface | 13 | 338 | 5 $\mu$s | 5760 | 72 $\mu$s |
| (ns, $10^{-4}$) Majorana qubit | Hastings-Haah | 15 | 1012 | 4 $\mu$s | 21840 | 52 $\mu$s |
| (ns, $10^{-6}$) Majorana qubit | Hastings-Haah | 7 | 244 | 2 $\mu$s | 16416 | 23 $\mu$s |

TABLE III. Physical resources required to implement a single logical qubit and to produce a single T state with a T factory for factoring across the qubit parameter examples specified in Table II. For the logical qubit we also show which QEC code was selected and the corresponding code distance.

### D. Combining the requirements into application resource estimates

We now outline how the tool combines the aspects described in Section III B and Section III C to estimate thew physical resources required our three application examples with our six qubit parameter examples. According to Section III B, the tool provides detailed estimates of the ISA-level executable, including the number of logical qubits, the minimum number of logical time steps, and also number of logical T states required. The tool also estimates the maximum logical error rates for logical qubits and logical T states that achieve the required algorithm execution accuracy.

Then, according to Section III C, the tool selects the QEC code and code distance for logical qubits and the T factory configuration which have the smallest physical space-time footprints while achieving the required logical error rate and logical T state error rate respectively. This choice depends on the physical qubit parameters, and sets the number of physical qubits per logical qubit and per T factory, and also the duration of a logical time step and the T factory. Next, the tool selects the number of T factories to balance the rate of production of T states against the rate of T state consumption demanded by the algorithm[2]. Finally, the tool computes the total number of physical qubits as the sum of physical qubits required to implement all logical qubits and all T factories. It computes the overall computation time using the number of logical time steps, and the duration of a single logical time step.

As can be seen in Fig. 3, qubit counts and runtimes vary by orders of magnitude based on the physical qubit parameters for all three applications. Even the smallest practical applications such as quantum dynamics involve devices with more than 100K physical qubits and applications in chemistry and factoring require upwards of 1M qubits. Unsurprisingly, the three orders of magnitude longer operation time for the ($\mu$s, $10^{-3}$) and ($\mu$s, $10^{-4}$) qubit models compared with

---

[2] In some cases, the number of T factories required to meet demand may be very large, prompting the exploration of a trade-off by slowing down the algorithm allowing for a reduced number of T factories and lower qubit count.
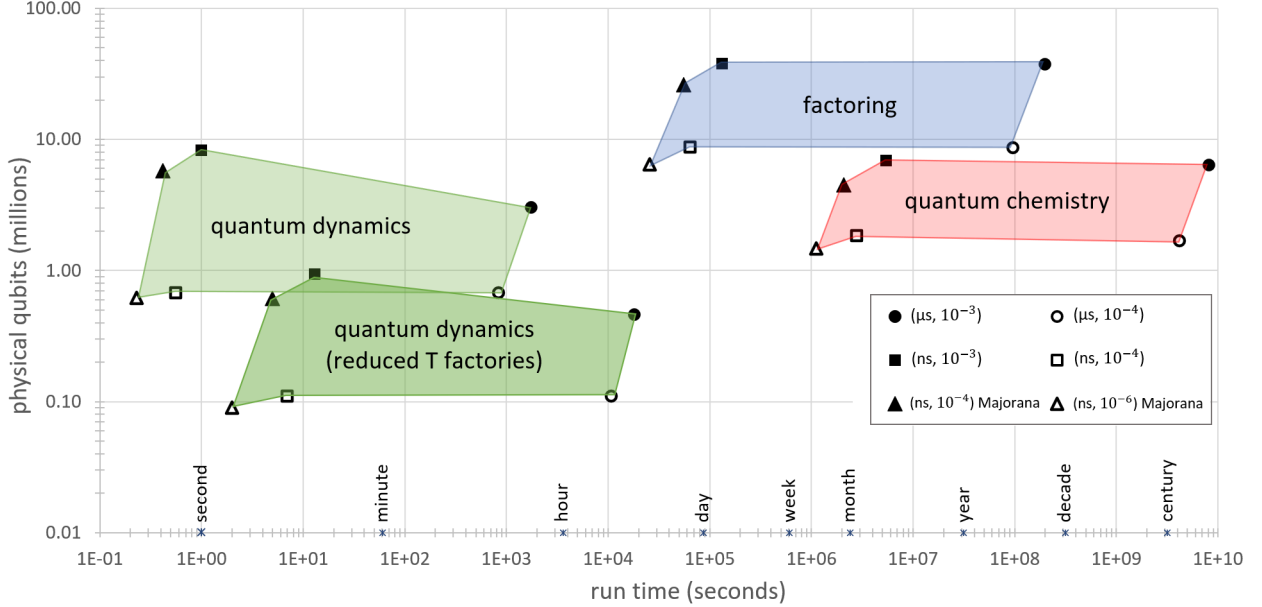
FIG. 3. Estimates of the resources required to implement three applications, assuming the qubit parameter examples specified in Table II. We explore a trade-off in the quantum dynamics application by considering two implementations: one which uses sufficient T factories to supply the needs of the shortest-depth algorithm and another which slows the algorithm down, allowing for a reduced number of T factories.

the other models results in about three orders of magnitude longer application run times. For quantum chemistry, this leads to impractical runtimes of more than a century.

To build the tool and obtain these estimates we have made many choices including what algorithmic, compilation and QEC options to include and also what approximations and assumptions to make. Other resource estimation works have made different choices, which leads to small differences in estimates [42]. In the appendices, we point out more explicitly these choices, assumptions and approximations, and in particular, collect together our primary assumptions in Appendix G. We anticipate that future work will extend and improve the tool and framework in two ways. Firstly, the estimates for a given stack will become more accurate as assumptions and approximations are honed and made more realistic, for example by including more detailed and nuanced noise models. Secondly, the estimates will become more favorable as improved solutions and optimizations are included in the stack and resource models, such as algorithmic improvements and hand-optimized compilation of important subroutines. We expect the broad conclusions that we draw from these results in the next section to hold true despite these choices, approximations and assumptions. This is because our conclusions are relatively insensitive to order of magnitude changes in resource estimates.

## IV. TECHNOLOGICAL IMPLICATIONS AND CONCLUSIONS

More than two decades ago, DiVincenzo [36] specified a set of fundamental requirements that any usable quantum computer should satisfy. For example, DiVincenzo identified the necessity of low error rates, by requiring *long relevant decoherence times, much longer than the gate operation time.* Since then, a variety of qubit technologies that satisfy DiVincenzo's criteria have been developed, including technologies such as superconducting and trapped ion qubits. However, it is an open question as to what additional conditions beyond DiVincenzo's criteria a qubit technology

must satisfy to scale to practical quantum advantage, and whether or not current technologies are on the path to do so.

Our work sheds light on this question, by estimating and analyzing the resources needed for quantum applications with practical value implemented on fault-tolerant architecture designs assuming a variety of underlying qubit technologies. Our analysis rests on several assumptions including qubit technology trends, QEC schemes, instruction sets, and choices across other parts of the stack. We have attempted to reflect the community's current best understanding of the quantum stack and its future evolution in these assumptions. Our resource estimates indicate that while scaling requirements to achieve the first scientifically interesting demonstrations of quantum advantage are not too stringent, scaling to commercially valuable applications, such as those using computational catalysis in chemistry, presents significant challenges. We deduce based on the results that to overcome these challenges and achieve scale, a qubit technology will need to be:

**Controllable.**— While most quantum hardware is controllable at small qubit counts, scaling qubit control to millions of fast and reliable qubit operations without introducing new noise channels is important. Controllability is at the core of scaling up, and has implications on several aspects of the physical qubit device, including its connectivity, error rates, operation speed and size.

Regarding connectivity, it has been shown that 2D connectivity is sufficient to implement fault-tolerant quantum error correction and build logical qubits [38, 103, 133], however such an array may present challenges for control and readout wiring. In contrast, while 1D connectivity is simpler for control wiring, it requires prohibitively low error rates [66] and seems unlikely to be viable with fabrication defects which permanently eliminates some qubits and operations. All-to-all connectivity is challenging to realize in many architectures, and may not be practically beneficial for large-scale implementations unless QEC schemes with good fault tolerant logical operations are developed to take advantage of such properties.

At the heart of fault-tolerant QEC is the requirement that the error rate of each physical operation (including qubit preparations, measurements, and gates) is below a threshold value. With no connectivity restrictions, the highest known threshold [80] is 3%, and 1% can be achieved with 2D connectivity [38, 103]. The qubit overhead of QEC goes down significantly as error rates are reduced further below this threshold value; however the overheads may still be prohibitive if error rates are not an order of magnitude better, say below around 0.1%. Even so, error rates have to be well below threshold while allowing for *parallel* operations at a scale of *millions of qubits*.

Controllability means being able both to control those millions of parallel operations with the desired error rates, and to readout out those millions of qubits in parallel to enable decoding of the errors at speed; all while ensuring the overarching logical clock time is fast enough to complete the computation within a month runtime or less. To execute syndrome measurements on these qubits and communicate the quantum measurements to the decoder, we require large quantum-classical bandwidth and processing power for decoding. The exact estimates of bandwidth requirements depend on the choice of QEC code, system size and physical operation times, but roughly, for system with a few million qubits, we estimate that several terabytes per second of bandwidth will be required between the quantum and classical plane. Furthermore, processing these measurements at a rate that is sufficient to effectively correct errors demands petascale classical computing resources that are tightly integrated with the quantum machine.

**Fast.**— The run times identified in our study help determine the desired speed of physical operations for a given quantum architecture to be practical. We find that logical gate times under 10 $\mu$s, in turn requiring physical gate times around 100 ns, would be needed to complete the quantum chemistry algorithm within a month, using a few million physical qubits. While these gate times are realistic for solid state or photonic qubits [14], this is several orders of magnitude faster than current proposals for large ion trap or neutral atom-based devices [85, 116]. For some algorithms, it may be possible to compensate for slower qubit speeds by using parallelization

techniques, most likely resulting in increasing the required number of physical qubits.

**Small.**— Bringing together a million and more physical qubits along with their control and readout systems presents another challenge: the size of the quantum computer. We foresee a monolithic, single-wafer approach as desirable and estimate that a linear extent of a micron or more is needed for wiring, defining a 'sweet spot' for the size of a qubit of around 10 microns [39].

To understand this sweet spot, consider estimating the approximate footprint of the quantum plane by multiplying the number of physical qubits by the size of a physical qubit, including control electronics. While precise estimates on the size of leading qubit platforms are not available in the literature, we estimate that a system with 10M superconducting transmon qubits would require an area of $\sim 10$ square meters based on current designs [3]. At such scale, a dilution refrigerator that provides the operating environment for the qubit plane will no longer be able to accommodate a monolithic quantum module. While modular quantum architectures that span multiple dilution refrigerators have been proposed [17], the reliability of the components, such as high bandwidth interconnects and low latency coherent quantum networks with associated distributed fault-tolerance schemes, and the ability to scale up while maintaining low error rates is not yet known. In contrast, the pursuit of a monolithic architecture in which ten million physical qubits must fit onto a single wafer suggests that shrinking the qubit to a linear extent of less than 100 microns would be needed. This presents another set of challenges, including not making the qubit too small. If qubits are too tightly packed, then scalable control and readout of arrays of qubits will become hard.

While no qubit technology currently implemented satisfies all of these requirements, two architecture proposals, while not yet implemented, appear to satisfy the requirements from a design perspective. These include recent proposals of electro-acoustic qubits [25, 51], and the topological qubit approach based on Majorana Zero Modes [71]. Understanding these requirements on the path to scaling up, for these and other existing technologies, illuminates the opportunities and advancements needed to achieve practical quantum advantage. It will be through studying the tradeoffs across the stack, from the algorithm to the qubit device, and then verifying them in implemnetation, that we will accelerate the development of scalable quantum computers. Both the framework and Azure Resource Estimator tool pave the way for the community to better understand the path to scaling up, the dependencies across the stack on the ultimate qubit numbers and runtimes, and the opportunities to bring down the costs at every layer. Practical quantum advantage is on the horizon; it may be accelerated through breakthrough techniques, including improved quantum error correction protocols, faster and more accurate decoders enabling improvements in the thresholds, richer logical operations and logical connectivity, improved algorithm compilation and optimization, and the discovery of low-cost, useful quantum algorithms with the promise of advantage. These research directions, and more, can be studied in the context of resource estimation, and just may be key to unlocking quantum at scale.

**APPENDICES**

In these appendices, we provide a self-contained, more detailed description of how the resource estimates presented in Section III of the main text are obtained. Since these estimates are calculated using the first release version of the Azure Quantum Resource Estimator, hereafter called 'the tool', these appendices also serve to specify the choices, assumptions and approximations of the tool.

Fig. 4 summarizes the set of examples and major options that form the stacks that we model to estimate resources, and Table IV collects together the input, output, and intermediate data used to calculate the estimates for each application example using each qubit parameter example, which are used to produce Fig. 3 in the main text. The appendices methodically address each aspect of this estimation calculation.
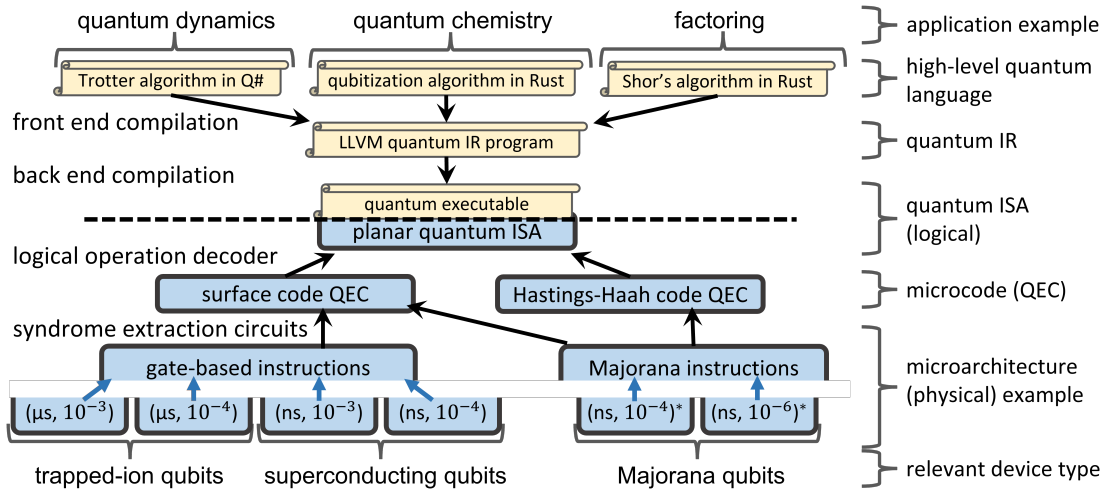


FIG. 4. The examples and options included in the quantum stack for our resource estimates. We consider three application examples at the top, and a range of hardware parameter examples relevant for a variety of hardware approaches at the base of the stack. We label layers of the quantum stack on the right, and the maps between these layers on the left. To highlight our qubit parameter examples, we separate the microarchecture layer into two sub-layers. All our examples flow through the same planar quantum ISA. Applications are translated down the software stack and the resulting ISA-level executable is input to the tool. The tool also takes as input configurable architecture models that include fault-tolerance details and physical qubit models. The tool outputs resource estimates such as the number of physical qubits or time required to run the application.

While it is impossible to fully model every aspect of a future large-scale quantum computing system, here we strive to model system components that are expected significantly impact performance; we abstract away or neglect aspects which are expected to have small contributions to the cost. In this tool we have made a number of assumptions to enable simple and efficient, yet effective, resource estimation. The formulas used in our estimates are expressed to the lowest order in parameters which are expected to be small (such as error rates). These formulas are therefore only valid for sufficiently small values of those parameters. (Various checks are included in the tool to flag when this is not the case given user inputs.) We anticipate that the resource estimates we present here will be further improved by incorporating additional details and relaxing a number of our assumptions. The primary assumptions that are made throughout the various aspects of this analysis are collected together and listed in Appendix G.

The remaining appendices are divided into three parts. The first part, Appendix A-Appendix D, describes models for different components of a large-scale quantum computer. In Appendix E, we

use these models to estimate physical resources (qubit counts, runtime) required for executing quantum applications. The third part, Appendix F, details the applications and algorithmic optimizations used to reduce resource requirements.

| application | planar quantum ISA requirements | | | | qubit parameters | QEC optimization | | | resource requirements | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $Q$ | $C_{\min}$ | $C$ | $M$ | | $d$ | $F$ | factory ratio | physical qubits | physical run time |
| quantum dynamics | 230 | $1.5 \cdot 10^5$ | $1.5 \cdot 10^5$ $1.5 \cdot 10^5$ $1.5 \cdot 10^5$ $1.5 \cdot 10^5$ $1.5 \cdot 10^5$ $1.5 \cdot 10^5$ | $2.4 \cdot 10^6$ | ($\mu$s, $10^{-3}$) ($\mu$s, $10^{-4}$) (ns, $10^{-3}$) (ns, $10^{-4}$) (ns, $10^{-4}$)* (ns, $10^{-6}$)* | 19 9 19 9 9 5 | 199 199 242 199 260 224 | 95% 95% 98% 95% 99% 95% | 3.0M 0.68M 8.2M 0.68M 5.8M 0.62M | 29 mins 14 mins 1.1 secs 0.56 secs 0.42 secs 0.23 secs |
| quantum dynamics (reduced T factories) | 230 | $1.5 \cdot 10^5$ | $1.5 \cdot 10^6$ $1.5 \cdot 10^6$ $1.5 \cdot 10^6$ $1.5 \cdot 10^6$ $1.5 \cdot 10^6$ $1.5 \cdot 10^6$ | $2.4 \cdot 10^6$ | ($\mu$s, $10^{-3}$) ($\mu$s, $10^{-4}$) (ns, $10^{-3}$) (ns, $10^{-4}$) (ns, $10^{-4}$)* (ns, $10^{-6}$)* | 21 11 21 11 11 5 | 18 17 22 17 22 23 | 56% 50% 78% 50% 79% 66% | 0.46M 0.11M 0.94M 0.11M 0.61M 0.09M | 5.3 hours 2.8 hours 13 secs 6.7 secs 5.0 secs 2.3 secs |
| quantum chemistry | 2740 | $4.1 \cdot 10^{11}$ | $4.1 \cdot 10^{11}$ $4.1 \cdot 10^{11}$ $4.1 \cdot 10^{11}$ $4.1 \cdot 10^{11}$ $4.1 \cdot 10^{11}$ $4.1 \cdot 10^{11}$ | $5.4 \cdot 10^{11}$ | ($\mu$s, $10^{-3}$) ($\mu$s, $10^{-4}$) (ns, $10^{-3}$) (ns, $10^{-4}$) (ns, $10^{-4}$)* (ns, $10^{-6}$)* | 33 17 33 17 17 9 | 15 14 17 17 19 19 | 6.9% 5.9% 14% 15% 22% 22% | 6.4M 1.6M 6.9M 1.9M 4.5M 1.3M | 260 years 130 years 2.0 months 1.0 month 24 mins 12 days |
| factoring | 25481 | $1.2 \cdot 10^{10}$ | $1.2 \cdot 10^{10}$ $1.2 \cdot 10^{10}$ $1.2 \cdot 10^{10}$ $1.2 \cdot 10^{10}$ $1.2 \cdot 10^{10}$ $1.2 \cdot 10^{10}$ | $1.5 \cdot 10^{10}$ | ($\mu$s, $10^{-3}$) ($\mu$s, $10^{-4}$) (ns, $10^{-3}$) (ns, $10^{-4}$) (ns, $10^{-4}$)* (ns, $10^{-6}$)* | 27 13 27 13 15 7 | 13 14 15 18 15 13 | 0.6% 0.8% 1.3% 1.1% 0.9% 1.2% | 37M 8.6M 37M 8.7M 26M 6.2M | 6.2 years 3.0 years 1.5 days 18 hours 15 hours 7.1 hours |

TABLE IV. Estimated resource costs to implement some example applications using hardware devices with various parameters. The application is compiled down from a high-level quantum program to a program expressed in the planar quantum ISA, with parameters as shown. To supply the planar quantum ISA with the capability to run this program, QEC is used. The code distance $d$ of the error correcting code and the number of T state distillation factories $F$ are carefully optimized to minimize the resources needed while providing enough fault protection that the overall computation runs reliably. We also show the percentage of the physical qubits which are used for T state distillation factories. We use an asterisk to indicate Majorana qubits. In the second implementation of the quantum dynamics example we slow the algorithm down by setting $C = 10C_{\min}$, which allows for a reduction in the number of T factories, thereby reducing the qubit overhead.

## Appendix A: Physical qubit models

We begin by specifying the instruction sets and noise models for the physical qubits representation in our analysis, at the level of the microarchitecture in our quantum stack. We also provide a set of example qubit parameters.

We assume that physical qubits are capable of implementing one of two instruction sets; see Fig. 5. The *gate-based* instruction set implements unitary operations such as CNOT or CZ as its

native entangling operations. On the other hand, the *Majorana* instruction set implements parity checks such as ZZ and XX measurements as its native entangling operations.[3] With Majorana instruction sets, unitary operations are realized using a sequence of measurements. For both gate-based and Majorana instruction sets in Fig. 5, we specify the allowed operations by visualizing the qubits laid out in a array with nearest neighbor connectivity. This visualization can either represent the true layout of the qubits in the hardware, such as for solid-state based approaches or can be considered as no more than a visualization to ensure that all operations in the instruction set are available, such as for trapped ions. Note that a particular hardware approach may have other allowed operations not in one of these instruction sets — these additional operations are not relevant for our analysis as we assume only operations from the instruction sets are used.
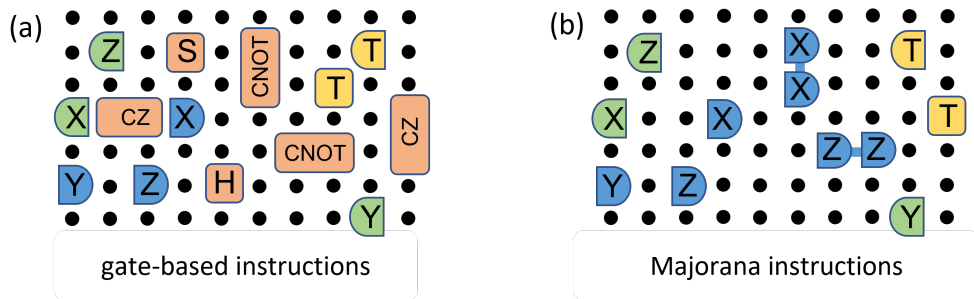


FIG. 5. We consider gate-based and Majorana instruction sets for physical qubits. Both sets include state preparation in the Pauli basis (green), T state preparation (yellow), T gate application (yellow), and measurement in all Pauli basis (blue). (a) The gate-based instruction set also includes Hadamard H and S gates, and CNOT and CZ entangling gates between adjacent qubits (all orange). (b) The Majorana instruction set also allows non-destructive joint Pauli measurements of adjacent qubits (blue).

We assume the standard noise model known as *circuit noise*, wherein each operation on the physical qubits fails independently. More specifically, we assume each Clifford operation fails with probability $p$, where Clifford operations include the Hadamard and phase gates, the Controlled NOT (CNOT) and Controlled Z (CZ) gates, single-qubit basis state preparation, and one-qubit and two-qubit Pauli measurements. We consider an idle qubit for a single time step to be an identity operation which can also fail with probability $p$. The preparation of a single-qubit T state, which is a non-Clifford operation, is assumed to fail with probability $p_T$. Failures are modeled by introducing random Pauli operators on qubits targeted by the operation. When the operation which fails is a measurement, the outcome is randomly flipped in addition to applying a Pauli operator to the support of the measurement. We also use the terminology of error rate in place of infidelity when describing noisy states, by making the simplifying assumption that the noise on the state arose from circuit noise (or logical Pauli noise where appropriate) with a given error rate.

In Table II, we consider a range of example qubit parameters which abstractly represent different regimes of operation fidelity and time. These values of parameters are selected to be relevant for different hardware approaches. For readability, we label each example by the approximate operation time, Clifford operation error rate and instruction set. For example, the first qubit has a gate time $t_{\text{gate}}$ and measurement time $t_{\text{meas}}$ of 100 $\mu$s, a Clifford error rate $p$ of $10^{-3}$ with a gate-based instruction set; therefore, we label it ($\mu$s, $10^{-3}$). The first four examples are gate-based qubits and the remaining two are Majorana qubits. To fit the simple circuit noise model which we assume, we choose a uniform error rate for all Clifford operations, which can be assumed to equal the highest of the error rates for any Clifford operation in the instruction set (including single-qubit gates and

---

[3] Note that in some literature, the Majorana instruction set is known as 'measurement-based', but this can cause confusion with quantum computing using cluster states, which is also commonly referred to as measurement-based.

measurements and two-qubit gates and measurements). We select the parameters for each example as follows:

- **($\mu$s, $10^{-3}$) qubit** and **($\mu$s, $10^{-4}$) qubit** may be relevant for future versions of trapped ion qubits [29, 108], which typically have operations times in the microsecond regime. We choose 100 $\mu$s as the gate and measurement time based on typical assumptions for ion qubits [87, 101]. We assume high-quality single qubit gates with an error rate of $10^{-6}$ and two cases for two-qubit gates with error rates of $10^{-3}$ and $10^{-4}$ based on the architectural assumptions in Ref. [85].

- **(ns, $10^{-3}$) qubit** and **(ns, $10^{-4}$) qubit** may be relevant for future versions of superconducting transmon qubits [3, 57, 76, 82, 118, 127] or spin qubits [54, 65, 70], which typically have operation times in the nanosecond regime. We project 50 ns gate times and 100 ns measurement times for future systems based on recent system implementations [37, 68], which can be viewed as relatively optimistic (particularly for measurements). We evaluate two cases, with $10^{-3}$ and $10^{-4}$ two-qubit gate error rates, respectively, as realistic and optimistic targets for a scaled up system [42, 68].

- **(ns, $10^{-4}$) Majorana qubit** and **(ns, $10^{-6}$) Majorana qubit** may be relevant for future Majorana qubits [71, 74, 117]. For these qubits, we assume that measurements and the physical T gate each take 100 ns. Owing to topological protection in the hardware, we assume single and two-qubit measurement error rates (Clifford error rates) of $10^{-4}$ as a realistic target and $10^{-6}$ as an optimistic target. Non-Clifford operations in this architecture do not have topological protection, so we assume a 5% and 1% error rate for non-Clifford physical T gates for the realistic and optimistic models respectively.

Together, these examples cover a range of operation times and error rates, enabling sufficient exploration of the resource costs anticipated to enable practical quantum applications. Note that in Fig. 4 and Table IV, we use an asterisk to indicate the Majorana qubit for notational convenience.

## Appendix B: Quantum error correction and the planar quantum ISA

Here we review two quantum error correction (QEC) schemes, focusing on estimates of the resources required for their implementation and for applying fault-tolerant logical operations on the encoded information. For qubits with a gate-based instruction set, we assume the surface code. It is the best-understood QEC scheme for this class of qubits and offers a high threshold for practical implementation. For qubits with a Majorana instruction set, we consider both the surface code and also the Hastings-Haah code, which is a recently developed QEC scheme that offers better space-time costs than surface codes on Majorana qubits in many regimes [55, 103]. It is also in principle possible to implement the Hastings-Haah code with qubits that use a gate-based instruction set, but the overhead is higher than the surface code in all regimes of interest in this case so we do not include it [44].

In Table V we provide formulas to estimate the physical qubit and time overheads required to implement logical qubits with the surface code and the Hastings-Haah code as a function of qubit design parameters. We estimate the logical error rate $P(d)$ of a patch of surface code or Hastings-Haah code of distance $d$ with physical error rate $p$ using the formula

$$P(d) = a \left( \frac{p}{p^*} \right)^{\frac{d+1}{2}}, \tag{B1}$$

| QEC scheme | logical error rate | qubits per logical qubit | logical time step |
|---|---|---|---|
| surface code (gate-based qubits) | $P_{\mathrm{sur}}(d) = 0.03 \left(\frac{p}{0.01}\right)^{\frac{d+1}{2}}$ | $n_{\mathrm{sur}}(d) = 2d^2$ | $\tau_{\mathrm{sur}}(d) = (4t_{\mathrm{gate}} + 2t_{\mathrm{meas}})d$ |
| surface code (meas-based qubits) | $P_{\mathrm{sur,meas}}(d) = 0.08 \left(\frac{p}{0.0015}\right)^{\frac{d+1}{2}}$ | $n_{\mathrm{sur,meas}}(d) = 2d^2$ | $\tau_{\mathrm{sur,meas}}(d) = 20t_{\mathrm{meas}}d$ |
| Hastings-Haah code (meas-based qubits) | $P_{\mathrm{HH}}(d) = 0.07 \left(\frac{p}{0.01}\right)^{\frac{d+1}{2}}$ | $n_{\mathrm{HH}}(d) = 4d^2 + 8(d-1)$ | $\tau_{\mathrm{HH}}(d) = 3t_{\mathrm{meas}}d$ |

TABLE V. Options of error correction schemes. We model each logical qubit as having a probability $P$ of failing during a single logical time step, which can be tuned by changing the code distance $d$, which is an odd integer. Increasing $d$ increases the number of qubits $n(d)$ required to store each logical qubit and the time $\tau(d)$ required for each logical time step. The formulas are based on Refs. [38] and [141] for gate-based surface codes, on Refs. [136] and [27] for surface codes implemented with Majorana instruction sets (replacing 8 steps to measure a single stabilizer in Ref. [136] by 20 steps to measure all stabilizers), and on Ref. [103] for Hastings-Haah codes.

where the pre-factor $a$ and threshold value $p^*$ can be extracted numerically from simulations.

Protected logical operations can be applied to logical qubits stored in the surface code or the Hastings-Haah code. We assume precisely the same types of logical patches, set of logical operations and costs (in units of number of tiles and logical time steps) for surface codes and Hastings-Haah codes, which forms the logical instruction set which we call the *planar quantum ISA* shown in Fig. 6. The differences that we account for between these two QEC codes are all captured by the different formulas for the logical error rate $P(d)$, number $n(d)$ of physical qubits in a tile and the time $\tau(d)$ for a logical time step for a given distance $d$, as specified in Table V. The logical operations can be inexpensive, almost perfect and instantaneous (Pauli unitaries, single-qubit preparations and measurements), medium cost, requiring ancilla tiles and one logical time step (multi-qubit Pauli measurements), or expensive, such as T state preparation for the execution of so-called non-Clifford gates, which requires a significant number of time steps and auxiliary qubits, which is discussed in detail in Appendix C. The patch choices and logical Clifford operations and costs are largely based on Refs. [13, 89].

As a simplification, we do not account for the overhead contributions associated with a number of known challenges. The circuits expressed in the physical-level instruction set which are used to implement QEC codes and logical operations fault-tolerantly are referred to as *syndrome measurement circuits*. For surface codes, the two-tile patches may require additional technical challenges to implement the necessary syndrome measurement circuits, such as measuring higher-weight stabilizer generators [22, 23]. Two-tile patches for Hastings-Haah codes have not been studied in detail in the literature but may pose similar implementation challenges. Moreover, the size and shape of a tile may need to be altered to accommodate the non-standard patches (such as the two-tile patch) and to account for ancillas between patches that can aid lattice surgery. We also assume that a Clifford operation which takes $A$ logical time steps to implement and involves $B$ tiles fails with probability $A \cdot B \cdot P(d)$.

## Appendix C: T state distillation factories

The non-Clifford T state preparation in the planar quantum ISA in Fig. 6 is crucial because the other operations in that quantum ISA, which are all Clifford operations, and are not sufficient for universal quantum computation [46]. To implement non-Clifford operations for practical-scale algorithms, we require low error rate T gates or T states, however these can be difficult to directly implement on logical qubits [9, 19], and can also be difficult for some physical qubits [71, 117].
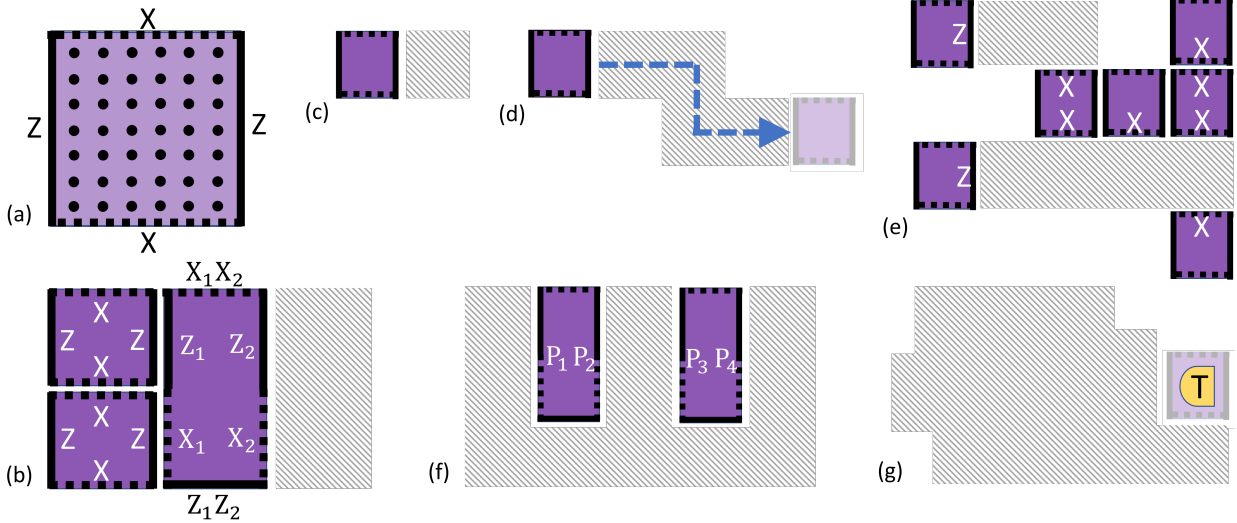
FIG. 6. **Fault-tolerant logical operations: instruction set of the planar quantum ISA**. (a) Patches of physical qubits encode logical qubits in either the surface code or the Hastings-Haah code, with error correction performed using operations from the instruction set. We group the array of physical qubits into tiles, each containing $n(d)$ physical qubits, and count logical operations in units of logical time steps $\tau(d)$ as defined in Table V. (b) We focus here on *one-tile patches*, and *two-tile patches*, which encode one and two qubits into $n(d)$ and $2n(d)$ physical qubits respectively. We label sections of the boundary of these patches which are relevant for the logical operations of the planar quantum ISA. **Preparation** of qubits in one-tile and two-tile patches in computational basis states can be applied in $0\tau(d)$. Arbitrary **Pauli unitaries** can be applied to qubits in one-tile and two-tile patches in $0\tau(d)$. Destructive **measurement** of a qubit in a one-tile patch in the X or Z basis and of both of the qubits in a two-tile patch (either both in the X basis or both in the Z basis) can be done in $0\tau(d)$. (c) **Hadamard and phase unitaries** can be applied to a qubit in a one-tile patch with an adjacent ancilla tile in $3\tau(d)$ and $2\tau(d)$ respectively. The ancilla can be on any boundary for the Hadamard, and next to a Z boundary for the phase. (d) A **move** of a qubit in a one-tile patch to a different tile along a path connecting any boundaries of the starting and ending tiles in $1\tau(d)$. (e) A non-destructive **multi-qubit X, Z type measurement** of the qubits in a set of one-tile patches can be implemented in $1\tau(d)$. This can be done by forming a connected region of tiles, where each tile hosting a qubit to be measured in the X (Z) basis connects through one or both of its X (Z) boundaries, and connections can be between adjacent qubit tiles, or mediated by ancilla tiles. (h) A non-destructive **multi-qubit arbitrary Pauli measurement** of a set of qubits stored in two-tile patches can be performed in $1\tau(d)$. This can be achieved with a connected region of ancilla tiles with access to all of the X and Z boundaries of each of the two-tile patches. (i) **T state preparation** in either a one-tile patch or a two-tile patch can be achieved using a distillation factory as described in Appendix C, using a non-integer number of tiles and logical time steps which depends on the quality requirements of the state.

Instead, the required low error rate T states are produced using a *T state distillation factory* [18, 79], which we sometimes refer to as a 'T state factory', a 'distillation factory' or even just a 'factory' for short. Here we consider the resources required to implement factories which produce T states in distance-$d$ surface or Hastings-Haah codes with error rate $P_T$.

The operation of a T state factory typically begins by first producing imperfect T states using some means, for example, a low-quality physical T gate acting on a physical qubit. These T states are then typically transferred to a logical qubit and is refined using a distillation unit, which uses only Clifford operations. This procedure can be iterated, where the output T states of one round are fed into the next round as inputs. In the error rate of physical T gates $p_T$ is much larger than that of physical Clifford operations $p$, it can make sense to do some distillation in physical qubits before transferring the state to a logical qubit.
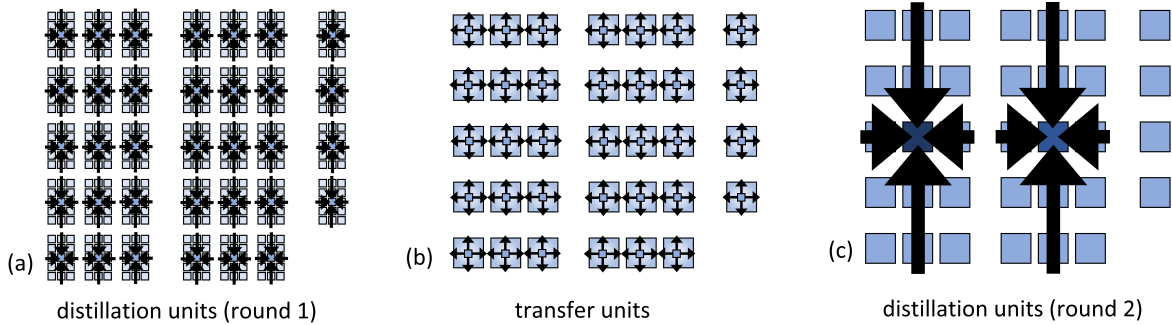
FIG. 7. An example of a two-round T state distillation factory. First, a transfer unit encodes $c_0 = 34 \cdot 15 = 510$ noisy T states into distance-$d_1$ code patches (not shown). (a) Next, in the first distillation round, each of these T states are input into one of $c_1 = 34$ copies of a 15-to-1 distillation unit are applied on the distance-$d_1$ code patches. It is possible that not all of the 34 copies of the distillation unit are accepted. (b) For those copies which are accepted, the code patch hosting the T state produced is expanded from distance $d_1$ to distance $d_2$, and moved into place for the next distillation round using a transfer unit. (c) In the second distillation round, 30 T states (provided at least that many remain) are fed into $c_2 = 2$ copies of a 15-to-1 distillation unit. If fewer than 30 T states remain (but more than 14), then just one of the $c_2 = 2$ copies of the distillation unit can be implemented. The T states encoded in distance-$d_2$ codes produced by round two form the output of the distillation factory.

We consider T state distillation factories which are implemented in a sequence of rounds, where each round consists of a set of identical distillation units run in parallel; see Fig. 7. We focus on factories composed of two types of units. *Distillation units* (see Table VI) take in a set of T states with a given quality, and output a smaller set of T states, but which typically have a higher quality. We assume that distillation units have the same kind of qubits as input and output, for example distance-$d$ surface codes. Distillation units typically involves a post selection test - if it is passed, then the output is accepted, otherwise the output is discarded. *Transfer units* take T states stored in qubits of one type in one location and transfer them to T states stored in qubits of a different type at another location. For example a transfer unit could take a set of T states stored in physical Majorana qubits, and transfer them into a set of distance-$d$ Hastings-Haah codes positioned where they are needed as inputs to a distillation unit in the next round. Another example would take qubits in distance-$d_1$ surface codes at the locations output from a previous round of distillation, and transfer them into distance-$d_2$ surface codes at the locations where they are needed as inputs to a distillation unit in the next round. In our analysis we neglect additional overheads required for implementing the transfer units and any noise they introduce. We specify a distillation factory $\mathcal{D}$ by stating: the number of rounds it uses $R$, and for each round $r \in \{1, 2, \ldots R\}$, both the distillation unit $M_r$ (including the qubit type, such as a specific physical qubit example, or a logical surface code or logical Hastings-Haah code, and code distance) and the number of copies $c_r$ of that unit. We do not specify the transfer units between rounds because we do not include the associated costs in our resource estimates. See Table VII for some example distillation factories.

The outputs of each distillation round are fed into the next, which leads to the following sequence of calculations which can be used to identify the overall properties of the factory. First we set the input error rate $Q_0$ for the first round to be the physical error rate of the T, i.e., $Q_0 = p_T$. Then let $Q_{r-1}$ be the error rate of the encoded T state that is input into round $r$, and let $P_r$ be the Clifford error rate of the unit at that round (where $P_r$ is calculated from the physical qubit type and the distance and type of error correcting code used by the unit). The acceptance probability $P_r^{\mathrm{acc}}$ and output T state error rate $Q_r$ for each unit in round $r$ are then calculated by setting $P_T$ to $Q_{r-1}$ and setting $P$ to $P_r$. In Table VI we show the relevant formulas for the four explicit distillation

| distillation unit | # input Ts | # output Ts | acceptance probability | # qubits | time | output error-rate |
|---|---|---|---|---|---|---|
| 15-to-1 space-eff. physical | 15 | 1 | $1 - 15p_T - 356p$ | 12 | $46t_{\text{meas}}$ | $35p_T^3 + 7.1p$ |
| 15-to-1 space-eff. logical | 15 | 1 | $1 - 15P_T - 356P$ | $20n(d)$ | $13\tau(d)$ | $35P_T^3 + 7.1P$ |
| 15-to-1 RM prep. physical | 15 | 1 | $1 - 15p_T - 356p$ | 31 | $23t_{\text{meas}}$ | $35p_T^3 + 7.1p$ |
| 15-to-1 RM prep. logical | 15 | 1 | $1 - 15P_T - 356P$ | $31n(d)$ | $13\tau(d)$ | $35P_T^3 + 7.1P$ |

TABLE VI. Distillation units. Physical distillation units assumes Majorana qubits with error rates of $p$ and $p_T$ for Clifford and input T states respectively. Logical distillation units assume either the surface code or the Hastings-Haah code, with each logical qubit requiring $n$ physical qubits, each logical time step requiring a time $\tau$, and logical error rates of $P$ and $P_T$ for Clifford and input T states respectively. The units are implemented by circuits of allowed operations as shown in Fig. 8 and Fig. 9. The coefficients for $p$ and $P$ are estimates based on numerical results for similar circuits, but may be inaccurate.

units specified in Fig. 8 and Fig. 9, which are all different explicit implementations of the standard 15-to-1 protocol [18].

The overall output T states of a distillation factory $\mathcal{D}$ are those output from its last round, and therefore have error rate

$$P_T(\mathcal{D}) = Q_R. \tag{C1}$$

We estimate the run time and qubit requirements as:

$$\tau(\mathcal{D}) = \sum_{r=1}^{R} \tau(M_r), \tag{C2}$$

$$n(\mathcal{D}) = \max_{r \in \{1,2,\ldots R\}} (c_r n(M_r)), \tag{C3}$$

where $n(M_r)$ is the number of qubits needed for a single copy of distillation unit $M_r$ and $c_r$ is the number of copies of the unit in round $r$. Each unit in a factory can fail, therefore the number of T states output by a factory is a random variable. We choose to address this by defining the number of output T states from the factory $M(\mathcal{D})$ as follows:

$$M(\mathcal{D}) = \# \text{ T states output in at least 99\% of runs of } \mathcal{D}. \tag{C4}$$

In principle, many factories have $M(\mathcal{D}) = 0$ according to this definition, which would occur for example if the final round consists of a single unit with an acceptance probability below 99%. Therefore many of the distillation factories that we consider over-provision units so that after accounting for some units failing the required number of T states is passed on to the next round. We assume that such decisions are specified to our framework by setting the appropriate distillation factory parameters such as $R$, $c_r$ and choice of unit types for each round.
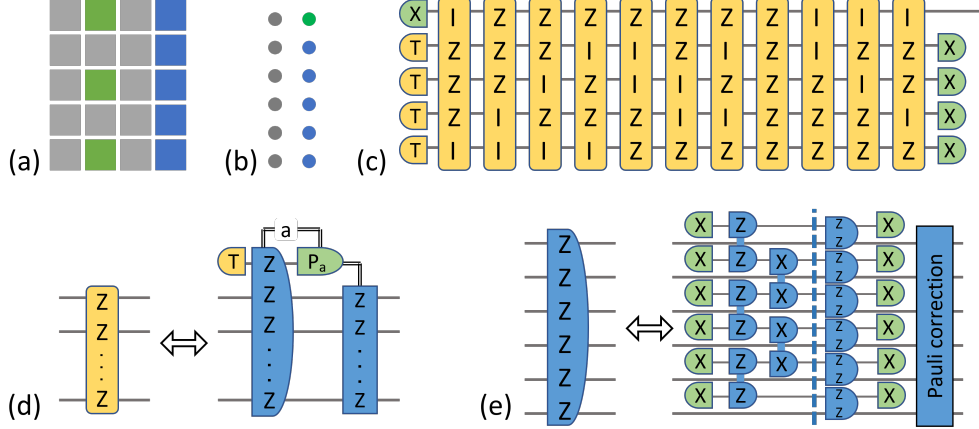
FIG. 8. Space-efficient 15-to-1 distillation units. (a) and (b) show the logical and physical layouts using 20 patches and 12 physical qubits respectively. (c) A unitary version of the distillation circuit based on Ref. [90]. Four T states are fed in at the beginning of the circuit, and then a sequence of diagonal non-Clifford rotations are applied (yellow boxes, with the rotation $\exp(-iP\frac{\pi}{8})$ applied for a box containing Pauli $P$). (d) A diagonal non-Clifford rotation $\exp(-iP\frac{\pi}{8})$ can be implemented using a small Clifford circuit given an ancilla prepared in the T state. First, a multi-qubit Z measurement is applied with the support of $P$ extended to include the ancilla, the outcome is $a$. Then the ancilla is measured in the basis $P_a$, where $P_0 = X$ and $P_1 = Y$. A Pauli correction is applied conditioned on the outcome of this ancilla measurement. (e) For physical-level distillation, the multi-qubit Pauli measurement must be broken down further into smaller pieces. Up until the dashed line involves the construction of a cat state on neighboring ancilla qubits. Note that if the support of the multi-qubit Z measurement is not over all six qubits, we replace the upcoming ZZ measurement for any qubits not in the support, and replace it by a single Z measurement. The Pauli correction at the end depends on the outcomes of the measurements in the circuit. **Logical distillation.**— These operations allow the distillation circuit to be re-expressed as surface code operations using 20 tiles to form a logical distillation unit. The five blue tiles store the qubits shown in the original circuit, and as such, four of them contain T states at the start of the distillation unit. The three green tiles also host T states, and these will be refreshed as the circuit proceeds. Each of the diagonal non-Clifford rotations is implemented using the gadget from (d), where the T state is held in one of the green tiles, and the five grey tiles represent ancillas that are used to facilitate the multi-qubit Z measurement, which takes a single logical time step. Following the multi-qubit Z measurement, that green ancilla is measured in either the X (taking zero logical time steps) or the Y basis (taking one logical time step with the aid of three adjacent ancilla tiles). Following this, another T state is transferred to the green tile, which we assume can be done in a single logical time step. As we have three green tiles, we can overlap this procedure for consecutive rotations such that the 11 non-Clifford rotations are each implemented in a total of 13 logical time steps. This represents the full time of the distillation unit since the initial X basis preparation and final single-qubit measurements take zero logical time steps. **Physical distillation.**— The physical distillation unit can be run on 12 physical qubits laid out as shown. Each of the 11 non-Clifford rotations requires five time steps, but the ancilla preparation in the X basis for all but the first round since the X basis at the end of the previous round deems it unnecessary. As the operations at the very start and end of the circuit are not supported on the ancillas, the can be done in parallel to the first and last steps of the non-Clifford rotations, such that this physical distillation unit requires $5 + 10 \cdot 4 + 1 = 46$ physical time steps.

## Appendix D: Compilation from QIR to the planar quantum ISA

The planar quantum ISA, consisting of the Clifford operations defined in Fig. 6 and non-Clifford T states distilled as described in Appendix C, forms a set of basic instructions that is sufficient for universal quantum computation. Here we describe how the program expressed in QIR, is mapped to the instructions of the planar quantum ISA, which we refer to as back end compilation.
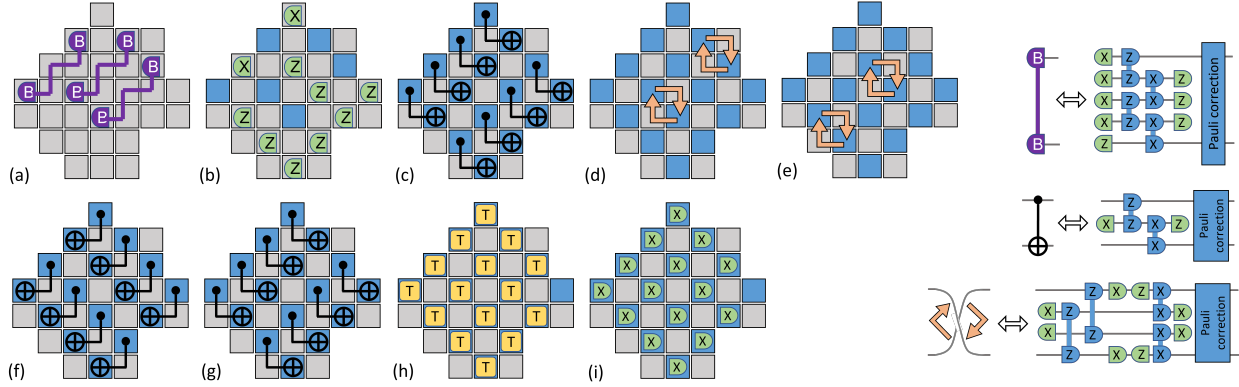
FIG. 9. Reed-Muller state preparation 15-to-1 distillation units. Each of the 31 squares represents a physical Majorana qubit or a logical qubit in a one-tile patch. In steps (a-g), a Reed-Muller state is prepared, using the gadgets shown on the right, which are available for both physical Majorana qubits and logical operations. We consider the grey squares ancillas, and the blue squares are those qubits which have been set and acted upon and host the Reed-Muller state. A naive counting of these first seven steps requires $(4+1+4+6+6+4+4) = 29$ physical time steps, however overlapping consecutive steps with disjoint support reduces this to 21 physical time steps. At the logical level, these seven steps require $(2+0+2+2+2+2+2) = 12$ logical time steps. (h) Next a T gate is applied to all but one of the qubits in the Reed-Muller state. For physical qubits, this can be done directly in a single physical time step. For logical qubits, this is instead achieved by transferring a T state to one of the grey qubits, and applying a $ZZ$ measurement as in Fig. 8(b). As we do not account for the transfer time, this requires a single logical time step. (i) Finally, each of the qubits to which T gates had been applied are measured in the X basis, which takes a single physical time step, or zero logical time steps. The overall time for the physical distillation unit is therefore 21+1+1= 23 physical time steps, and for the logical distillation unit is 12+1+0=13 logical time steps.

| $P_T$ | round 1 | | | | round 2 | | | | resources | |
|---|---|---|---|---|---|---|---|---|---|---|
| | unit | qubits | time | copies | unit | qubits | time | copies | qubits | time |
| $5.6 \cdot 10^{-11}$ | 15-to-1 space-eff. $(d=9)$ | 3240 | 46.8 $\mu$s | 1 | | | | | 3240 | 46.8 $\mu$s |
| $2.1 \cdot 10^{-15}$ | 15-to-1 space-eff. $(d=5)$ | $16000 = 16 \cdot 1000$ | 26 $\mu$s | 16 | 15-to-1 RM prep. $(d=13)$ | 10478 | 57.2 $\mu$s | 1 | 16000 | 83.2 $\mu$s |
| $5.51 \cdot 10^{-13}$ | 15-to-1 space-eff. $(d=3)$ | $5760 = 16 \cdot 360$ | 15.6 $\mu$s | 16 | 15-to-1 space-eff. $(d=11)$ | 4840 | 57.2 $\mu$s | 1 | 5760 | 72.8 $\mu$s |

TABLE VII. T state distillation factories to target the required error rates for our three application examples with the (ns, $10^{-4}$) qubit parameter example.

Here we describe the back end compilation scheme that we assume in this paper at a high level. Our compilation scheme is formed by combining sequential Pauli-based computation (SPC), as described in Ref. [89] and elaborated upon in Ref. [22], with an approach to synthesize sets of diagonal non-Clifford unitaries in parallel, as was done in Ref. [8]. We call this compilation scheme Parallel Synthesis Sequential Pauli Computation (PSSPC).

We assume that the input is a QIR program which has first been re-expressed as a sequence of Clifford layers and non-Clifford layers, which we refer to as the input circuit. A Clifford layer of the input circuit consists of a set of basic instructions, which are either Clifford unitaries, basis state preparations, or Pauli measurements. Non-Clifford layers of the input circuit require T states to

implement, and include operations such as T gates, Toffoli gates, and single-qubit arbitrary angle rotations (where the angle is not an integer multiple of $\pi/4$) around either the X, Y or Z axis. The input circuit can be dynamic in the sense that the measurement outcomes in one Clifford layer can affect the instructions applied in subsequent layers of the input circuit. In the following description, we describe a compilation strategy which can be applied to compile a dynamic input circuit layer by layer[4] to produce an equivalent output circuit which is expressed in terms of the logical instruction set in Fig. 6.

To orchestrate the execution of this algorithm, PSSPC performs the following circuit transformations (see Fig. 10).

**Step 0.**— We first re-express the input algorithm such that it consists of a sequence of Clifford layers and diagonal non-Clifford layers. This is straight-forward: a rotation around the X or Y axis can be converted by conjugation of a Clifford into a rotation around the Z axis $R_z(\theta) = |0\rangle\langle0| + e^{i\theta}|1\rangle\langle1|$. The Cliffords that are used to implement this transformation are then absorbed into the Clifford layers before and after the non-Clifford layer. Similarly, a Toffoli gate can be re-expressed as a CCZ gate with a Hadamard applied before and after on the target qubit.

**Step 1.**— Delegating expensive operations to synthesis qubits to increase parallelism: Non-Clifford rotations $R_z(\theta)$ are frequently used in quantum algorithms, but are among the most resource intensive operations when we consider implementation costs. To improve algorithm performance, PSSPC seeks to delegate these operations to ancilla qubits, referred to as synthesis qubits. By delegating these operations, multiple rotations can be synthesized in parallel to increase the T consumption rate and improve algorithm runtime.

As shown in Fig. 10(b), PSSPC first re-expresses rotation operations on algorithm qubits using ancillary synthesis qubits. To implement a rotation such as $R_z(\theta)$ on an algorithm qubit, it is first entangled with an ancilla by measuring a Pauli supported on both that qubit and the ancilla. Then a sequence of operations are performed on the ancilla to apply the rotation's phase using Clifford gates and also T gates. The phase is kicked back to the algorithm qubit by measuring out the ancilla and applying a Pauli correction to the algorithm qubit. All other diagonal non-Clifford unitaries (including the T gate and the CCZ) are applied similarly (see Figure 20 in [8]).

**Step 2.**— Eliminating Clifford unitaries: After delegating the synthesis of rotations and Toffoli gates to ancillas, operations on the algorithm qubits are a series of Clifford unitaries and Pauli measurements. While these operations can be implemented using the instruction set in Fig. 6, we can further optimize the circuit by eliminating Clifford unitaries.

Clifford unitaries map Pauli operations to other Pauli operations. Using this property, PSSPC commutes each Clifford unitary through the subsequent Pauli operations. This transforms the local Pauli operations to multi-qubit Pauli measurements and Pauli corrections. This optimization allows us eliminate the Cliffords at the cost of executing multi-qubit Pauli measurements.

Like SPC, a major benefit of PSSPC is that all Clifford unitaries in the input algorithm are eliminated and do not need to be applied. SPC differs from PSSPC by synthesizing the non-Clifford gates in place rather than in external ancillas, and the Ts resulting algorithm which has the drawback of serializing all T gates in the circuit. This drawback is partially overcome by

---

[4] When estimating resources to implement a dynamic algorithm, we randomly pick outcomes of measurements in the algorithm such that each layer is specified for explicit analysis.

PSSPC by parallelizing the synthesis of rotations which appear in the same layer of the input circuit, but the phase kickback of the synthesized rotations within a layer is still serial. Other layout and compilation approaches have been proposed [8, 28, 32, 53, 63, 83, 86, 100, 128, 143], but are more involved to implement and analyze in detail.

We estimate the resources required to implement an algorithm using PSSPC in terms of the following parameters of the input algorithm. Let $Q_{\text{alg}}$ be the number of logical qubits used by the algorithm. Let $M_R$ be the number of single-qubit rotations, $M_T$ be the number of T gates, and $M_{\text{Tof}}$ be the number of Toffoli gates, and $M_{\text{meas}}$ be the number of Pauli measurements of the input algorithm. Let $D_R$ be the number of non-Clifford layers in which there is at least one arbitrary angle rotation (note this can be smaller than the total number of non-Clifford layers since it excludes those layers consisting entirely of T gates and Toffoli gates). Let the target error budget for synthesis in the overall program be $\epsilon_{\text{syn}}$ (measured in the diamond norm).

We assume the fast block layout from Ref. [89] (also shown in Fig. 11), such that each pair of algorithm qubits is stored in a two-tile surface code patch, surrounded by ancilla tiles which can be used to measure an arbitrary Pauli operator on the algorithm qubits in one logical time step.
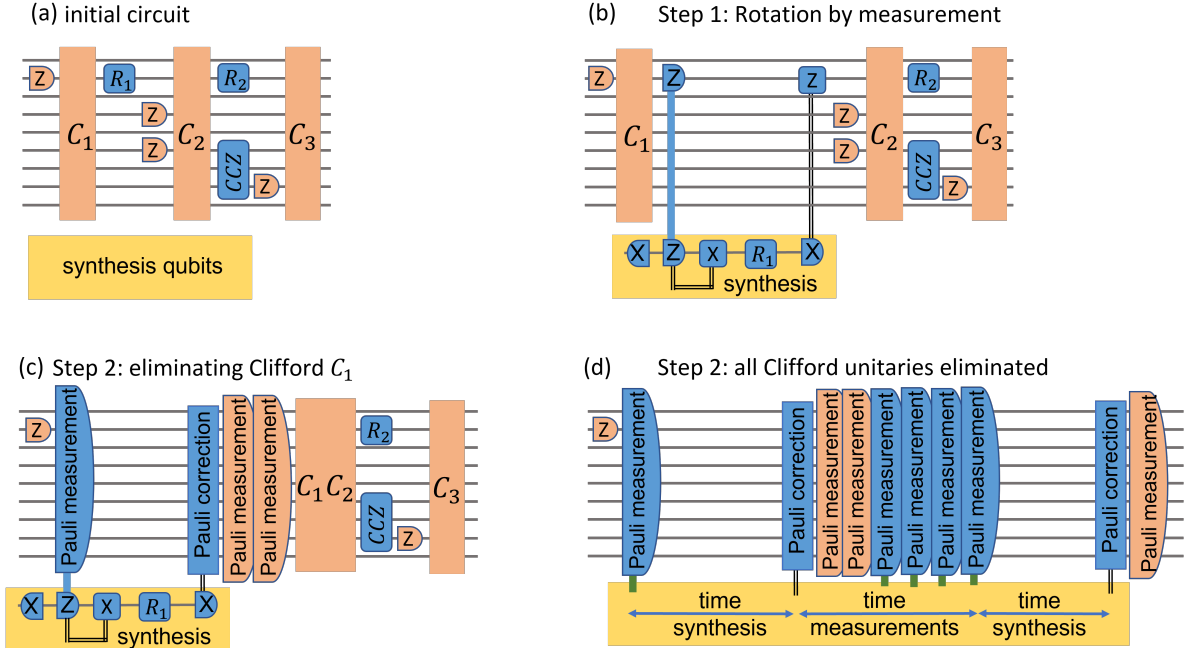


FIG. 10. (a) The initial circuit consists of alternating Clifford (orange) and diagonal non-Clifford layers (blue). For illustrating the principles behind PSSPC, we show the implementation details for the first non-Clifford (the single rotation $R_1$ between Clifford unitaries $C_1$ and $C_2$), but the concepts extend to parallel distinct rotations. (b) $R_1$ is delegated to an ancilla qubit, where the operation is synthesized using T states and the phases are kicked back to the algorithm qubit. (c) $C_1$ is eliminated by commuting it through the Pauli operations that follow it, transforming them into multi-qubit Pauli operations. (d) After eliminating all Clifford unitaries, the circuit consists of a sequence of Pauli measurements and Pauli corrections, along with synthesis on ancilla qubits. For the second set of non-Cliffords (between $C_2$ and $C_3$ in the input circuit), we have a series of Pauli measurements and synthesis, followed by a final Pauli correction after synthesis and then the Pauli measurements from the original circuit. We assume that the Pauli correction for the current non-Clifford layer is completed before moving on to the next layer.

This results in the number of logical qubits[5] $Q$ given by

$$Q = 2Q_{\text{alg}} + \left\lceil \sqrt{8Q_{\text{alg}}} \right\rceil + 1. \tag{D1}$$

We account for the resources required for distillation separately in Appendix E, and in this section we assume that T states are produced on command and simply count how many T states are needed to implement the algorithm.

Synthesizing a single-qubit rotation to diamond-norm accuracy $\epsilon'$ re-expresses the unitary as Cliffords and at most $R_T(\epsilon')$ non-Clifford T gates, where

$$R_T(\epsilon') = \left\lceil A \log_2(1/\epsilon') + B \right\rceil, \tag{D2}$$

for constants $A$ and $B$. We assume the Clifford+T synthesis in Table 1 of Ref. [78], which results in $A = 0.53$ and $B = 5.3$.[6] To meet the accuracy requirement for the overall algorithm, which involves $M_R$ single-qubit rotations, we set $\epsilon' = \epsilon_{\text{syn}}/M_R$.

As described in Fig. 10(d), the computation proceeds in layers. Each diagonal non-Clifford unitary is applied as part of a layer of non-Clifford rotations from the input circuit. The first step to implement such a unitary is to couple the algorithm qubits to the synthesis ancilla qubits. This requires one multi-qubit Pauli measurement for each algorithm qubit in the support of the unitary. Therefore one logical time step is used for a T gate or a single-qubit rotation $R_z(\theta)$, while each Toffoli requires three logical time steps. We assume that T gates and CCZ gates require no additional synthesis time beyond the time required for these measurements, justified by techniques in Ref. [43, 59]. To estimate the time required to synthesize an arbitrary angle rotation $R_z(\theta)$, we note that the synthesis consists of $R_T(\epsilon_{\text{syn}}/M_R)$ non-Clifford T gates interleaved with Clifford unitaries. We assume this can be implemented by a sequence of $R_T(\epsilon_{\text{syn}}/M_R)$ Pauli measurements supported on the synthesis ancilla and other ancillas to which the T states are delivered from distillation factories, taking $R_T(\epsilon_{\text{syn}}/M_R)$ logical time steps. The final step to implement a rotation, after the synthesis has finished, is to apply a measurement of the ancilla and then a Pauli correction to the algorithm qubits, both of which can be done instantaneously.
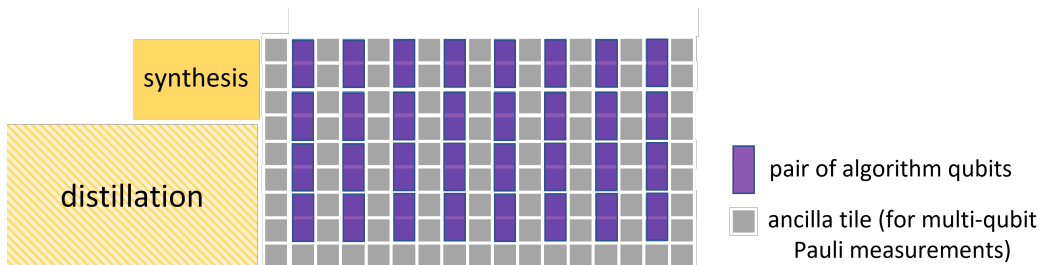


FIG. 11. Overall layout of qubits in the quantum computer, using the PSSPC layout. The main block of logical qubits is purple and grey. Each purple rectangle denotes a pair of algorithm qubits encoded in a two-tile patch of surface or Hastings-Haah code, while the grey boxes are tiles of ancillas to facilitate multi-qubit Pauli measurements. A block of qubits is dedicated for distilling T states which are used to apply T gates and Toffolis, or are passed on to the synthesis block to synthesize rotations. We neglect the qubits needed for synthesis. In the overall program, we use $Q$ to denote the number of tiles in the main block, $M$ to denote the number of T states that are produced, and $C$ to denote the number of logical time steps needed to run the program.

---

[5] More precisely, $Q$ is the number of tiles, some of which are allocated to ancillas, and some of which are allocated to store algorithm qubits in two-tile patches.

[6] For simplicity, we use this formula for all single-qubit arbitrary angle rotations, and do not distinguish between best, worst and average cases. We also do not incorporate further improvements from [78] that use Clifford+$\sqrt{\text{T}}$ synthesis and known improvements for special cases, for example when many rotations are applied in parallel with the same angle [41].

The time required to implement a non-Clifford layer depends on whether the layer consists entirely of Toffoli and T gates, or whether it also contains some arbitrary angle rotations. In the former case, the number of logical time steps required is simply the number of measurements needed to entangle the algorithm and synthesis ancillas. In the latter case, in addition to the time needed for these entangling measurements, one also needs to add the time for synthesis. It is possible to further compress this by overlapping synthesis for one non-Clifford layer with the measurements for the following Clifford and non-Clifford layers, but for simplicity we assume that consecutive layers are executed sequentially without overlap. Within each layer however, the synthesis of rotations that were parallel in the input quantum IR circuit is performed in parallel. Since all Clifford unitaries in the Clifford layer are eliminated, the time required for the Clifford layer is simply the time required to implement the measurements it contains. The total number $C_{\min}$ of logical time steps is therefore

$$C_{\min} = (M_{\text{meas}} + M_R + M_T) + \lceil A \log_2(M_R/\epsilon_{\text{syn}}) + B \rceil D_R + 3M_{\text{Tof}}. \tag{D3}$$

Note that this is our estimate of the number of logical time steps that would be needed if the T state factories produce T states at a sufficient rate so as to not limit the algorithm speed. In practice, the time $C \geq C_{\min}$, and one could increase $C$ to slow the algorithm down if for example it puts a strain on the qubit resources available to produce T states.

Finally, using Eq. (D2) and noting that four T states are sufficient to implement Toffoli, the number of T states $M$ is given by

$$M = \lceil A \log_2(M_R/\epsilon_{\text{syn}}) + B \rceil M_R + 4M_{\text{Tof}} + M_T. \tag{D4}$$

We neglect the ancilla qubits used for unitary synthesis and those that would be required to transport the T states from distillation factories to where they are consumed for synthesis.

## Appendix E: Putting the pieces together to estimate resources of an algorithm

This appendix puts together components described over previous appendices in order to estimate the resources required to implement a high-level language quantum algorithm on a specified fault-tolerant architecture. After compilation of the algorithm to the instruction set described in Fig. 6, the problem reduces to estimating the resources required to implement a circuit with $Q_{\text{alg}}$ algorithm qubits, with $M_R$ single-qubit rotations spread among $D_R$ layers, $M_T$ T gates, and $M_{\text{Tof}}$ Toffoli gates. Furthermore, let $1 - \epsilon$ be the desired probability that the overall algorithm succeeds.[7]

The value of $\epsilon$, which we refer to as the *error budget*, depends on how the samples from algorithm instances are to be post-processed, and its value is set by taking the desired algorithm execution accuracy to be $1 - \epsilon$. For example, if one is running Shor's algorithm for factoring integers, a large value of $\epsilon$ may be tolerated as one can check that the output are indeed the prime factors of the input. On the other hand, a much smaller $\epsilon$ may be needed for an algorithm solving a problem with a solution which cannot be efficiently verified. For the target architecture, we assume that a specific physical qubit model and quantum error correction scheme has been chosen from Appendix A and Appendix B. The target architecture implements the logical instruction set described in Fig. 6.

To achieve the required algorithm success probability, we require that

$$\epsilon_{\text{log}} + \epsilon_{\text{dis}} + \epsilon_{\text{syn}} \leq \epsilon, \tag{E1}$$

---

[7] More formally, $\epsilon$ is the total variational distance between the probability distribution of the final output bit strings of the compiled circuit and the algorithm circuit. Commonly, the circuits correspond to channels and the required total variational distance can be bounded using the diamond norm distance between the compiled channel and algorithm's channel as explained in Appendix B in [78].

where $\epsilon_{\mathrm{log}}$, $\epsilon_{\mathrm{dis}}$ and $\epsilon_{\mathrm{syn}}$ are the probabilities of at least one logical failure, at least one faulty T distillation, and at least one failed rotation synthesis respectively. In our estimates, we retain only the first order terms in these small probabilities. We ensure Eq. (E1) is satisfied by requiring that each of $\epsilon_{\mathrm{log}}$, $\epsilon_{\mathrm{dis}}$ and $\epsilon_{\mathrm{syn}}$ are at most $\epsilon/3$.

We assume the PSSPC compilation strategy from Appendix D, which sets the number of logical qubits $Q$ and the number of T states consumed $M$. This also specifies a minimum possible logical runtime of the algorithm $C_{\mathrm{min}}$. In what follows we establish the code distance $d$ required for error correction, the choice of distillation factory $\mathcal{D}$, the number of distillation factories $F$, and the number of logical time steps $C$. We then use these to identify the number of physical qubits $q$ and runtime $t$.

Before computing these quantities, let's consider the need to balance the speed of the algorithm vs. the rate of T state production. From Appendix D, we know that the algorithm requires at least $C_{\mathrm{min}}$ time steps. That is, we are free to run the algorithm at any desired speed by choosing $C \geq C_{\mathrm{min}}$. Similarly, we are allowed to choose the number of factories $F$ and control how many T states are available for the algorithm per logical time step. However, we require that for any chosen value of $C$, $F$ must be chosen large enough to produce all the required T states for the algorithm within $C$ time steps. Therefore, with large $C$, as we slow down the algorithm, we can afford to reduce the number of T factories and save qubits required for their implementation[8]. Similarly, as we speed up the algorithm to finish in the least number of logical time steps possible, we require a large number of T factories to quickly supply the required T states, and incur higher qubit overheads. To navigate this space-time tradeoff, we fix a scenario such as space-optimal or time-optimal or an intermediate case and choose suitable values for $C$ and $F$. In what follows, we present the analysis for the general case where some $C \geq C_{\mathrm{min}}$ has been selected.

**Estimating resources for logical operations:** We determine the number of physical qubits and time required for the logical operations by selecting an appropriate code distance used for the logical qubits. From Eq. (B1), we see that $d = 2\log(a/P)/\log(p^*/p) - 1$ for constants $a$ and $p^*$ that depend on the type of physical qubit chosen in the architecture. Seeking the smallest error correction overhead that achieves an acceptably low probability that logical failure occurs in the implementation of the algorithm, i.e., that $QCP = \epsilon_{\mathrm{log}} \leq \epsilon/3$, we select the following distance,

$$d = \left\lceil \frac{2\log(a\epsilon/(3QC))}{\log(p^*/p)} - 1 \right\rceil_{\mathrm{odd}}. \tag{E2}$$

We then obtain the physical run time by multiplying the number of logical time steps $C$ by the time per logical time step $\tau(d)$ as defined in Appendix B,

$$t = \tau(d)C. \tag{E3}$$

Similarly, the number of physical qubits required for the logical operations is $Qn(d)$. Note that the code distance is a function of $C$. If $C$ is chosen to be much higher than $C_{min}$ to reduce the number of factories, we may require a higher code distance to protect the algorithm qubits for the longer execution duration.

**Estimating resources for T distillation:** Next, we determine $\mathcal{D}$, the number of factories and the physical resources required for distillation. Towards this, we first require the number of T states required by the algorithm at the given error budget. That is, we need to satisfy Eq. (D4) with $\epsilon_{\mathrm{syn}} \leq \epsilon/3$, which allows us to calculate $M$, the minimum number of T states required. Next, we determine the quality of these T states. In Table III, we show the number of physical qubits $n(\mathcal{D})$

---

[8] Note that slowing down the algorithm may not always result in a reduction in overall qubit counts. As the algorithm is slowed down, we can use lesser T factories and incur less qubit overhead for them. But for the algorithm's logical qubits, the code distance depends on the number of logical time steps (see Eq. (E2)). We may require a higher code distance to protect the algorithm qubits during a slower execution

and the time $\tau(\mathcal{D})$ to distill a single T state with error rate $P_T$ with various choices of distillation factory $\mathcal{D}$. As we aim for the $M$ injected T states to fail with low probability, we require that $MP_T(\mathcal{D}) = \epsilon_{\text{dis}} \leq \frac{\epsilon}{3}$. Given a list of all the considered factories with the specified hardware type $\{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \ldots, \mathcal{D}_m\}$, we therefore select $\mathcal{D}$ as the factory with the smallest space-time footprint which satisfies this inequality, i.e.,

$$\mathcal{D} = \underset{\mathcal{D}_i \in \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \ldots, \mathcal{D}_m\}}{\arg\min} \{n(\mathcal{D}_i)\tau(\mathcal{D}_i)|P_T(\mathcal{D}_i) \leq \epsilon/3M\}. \tag{E4}$$

Finally, the smallest number $F$ of distillation factories capable of producing the demanded $M$ T states during the algorithm's run time is selected

$$F = \left\lceil \frac{M \cdot \tau(\mathcal{D})}{M(\mathcal{D}) \cdot t} \right\rceil. \tag{E5}$$

At this point a check should be made to ensure that the factories execute during the runtime of the algorithm, i.e., that $\tau(d)C \geq \tau(\mathcal{D})$. This catches the case when the runtime of an algorithm is less than the runtime of a single T factory. If this is not satisfied, $C$ should be increased and the analysis repeated.

**Total resources:** The total number of qubits $q$ for the quantum computation is then found from adding those qubits for distillation to those required to build the logical qubits,

$$q = Fn(\mathcal{D}) + Qn(d). \tag{E6}$$

The total runtime of the algorithm is given by Eq. (E3).

## Appendix F: Applications

Here we provide further details of the three example applications, and the calculations that result in the resource estimates in Table IV.

Before getting into specifics about each of three application examples, we make a remark about quantum algorithm sampling. Throughout this work, we are focused on the resources required to implement a full quantum algorithm once, however it is often the case that quantum algorithms must be repeated (even when run perfectly). This is because the output of some quantum algorithms is not a deterministic bit string, but instead is a drawn from a probability distribution over bit strings, and the information from the algorithm may require learning details of the distribution, requiring a set of samples. We do not account for the resource costs that algorithms may need to be sampled many times in this work. This could in principle be achieved by either re-running the algorithm consecutively on the available quantum computer, or by having many independent quantum computers running in parallel, or some combination.

**Quantum dynamics**.— Here we outline the calculation that produced the resource requirements of the quantum dynamics example in Table IV, following the approach in Ref. [107]. We are interested in the time evolution $e^{-iHt}$ by the 2D $\sqrt{N} \times \sqrt{N}$ transverse field Ising Hamiltonian

$$H = -J \underbrace{\sum_{\langle j,k \rangle} Z_j Z_k}_{A} + g \underbrace{\sum_j X_j}_{B} \tag{F1}$$

with nearest-neighbour interactions is accomplished using fourth-order product formulas. Trotter & Suzuki [56] define a recursive construction using

$$U_2(\Delta) = e^{-iA\Delta/2}e^{-iB\Delta}e^{-iA\Delta/2} = e^{-iH\Delta} + \mathcal{O}(\Delta^3), \tag{F2}$$

$$U_4(\Delta) = U_2(\gamma\Delta)U_2(\gamma\Delta)U_2((1-4\gamma)\Delta)U_2(\gamma\Delta)U_2(\gamma\Delta) = e^{-iH\Delta} + \mathcal{O}(\Delta^5),$$

$$\gamma = (4 - 4^{1/3})^{-1}.$$

When applying $T > 1$ time steps steps, the first and last terms may be merged, leading to $10T + 1$ exponentials in Eq. (F2). Note that $A$ and $B$ may be interchanged depending on which is more expensive. These product formulas reduce the resource costs of quantum dynamics to evaluating the resource cost of applying $e^{-iA\Delta}$ and $e^{-iB\Delta}$. All terms in $A$ commute and similarly for $B$, and $e^{-iB\Delta}$ is the product of $N$ single-qubit rotations $\prod_j e^{-iX_j g\Delta}$. The simulation can be carried out by assigning a qubit to each of the $N$ sites in the lattice, such that $Q_{\text{alg}} = N$. By conjugating with a CNOT gate, each nearest-neighbour $Z_j Z_k$ term is transformed in to a single-qubit Pauli $Z_j$. Thus $e^{-iA\Delta}$ can be implemented by eight depth-1 layers of $N/2$ CNOT gates and four depth-1 layers of $N/2$ single qubit rotations $\prod_j e^{iZ_j J\Delta}$. The number of rotations is therefore $M_R = (5T + 1)N + 10TN = (15T + 1)N$. Finally, the depth of the rotation gates over the $T$ time steps is $D_R = 5T + 1 + 4(5T) = 25T + 1$.

In the example we consider, we take $\epsilon = 0.001$, $N = 100$ and $T = 20$ such that when the algorithm, which is initially written as a Q# program, is explicitly compiled down to a QIR-level program with $Q_{\text{alg}} = 100$, $M_R = 30100$, $M_T = 0$, $M_{\text{Tof}} = 0$ and $M_{\text{meas}} = 1400000$ and $D_R = 501$. From Eq. (D1), Eq. (D3) and Eq. (D4), we calculate that with the PSSPC compilation and layout, the ISA-level executable that this QIR program would compile down to has $Q = 230$ logical qubits, which are needed for at least $C_{\text{min}} = 1.5 \cdot 10^5$ logical time steps, consuming $M = 2.4 \cdot 10^6$ T states.

At the physical level, here we illustrate the analysis with the (ns, $10^{-4}$) qubit example (with the other qubit parameter examples being analyzed analogously). In this example we will ultimately increase $C$ so that it is significantly larger than $C_{\text{min}}$ to reduce the qubit count as a trade-off paid for by an increase in run time. To see why we do so, let us first consider the case where $C = C_{\text{min}}$. We use Eq. (E2) to find the distance $d = 9$, such that $n(d) = 2d^2 = 162$ and $\tau(d) = (4t_{\text{gate}} + 2t_{\text{meas}})d = 400$ ns. This sets the algorithm run time via Eq. (E3) to be $t = 0.55$ s. From Eq. (E4), a distillation factory is selected. There are $F = 199$ factories. Each factory has a single round of logical distillation (15-to-1 space efficient), with one unit of code distance 9. Each factory produces $M(\mathcal{D}) = 1$ encoded T state with error rate $5.6 \cdot 10^{-11}$, and contains $n(\mathcal{D}) = 3240$ physical qubits, and requires $\tau(\mathcal{D}) = 46.8$ $\mu$s to run. The total number of physical qubits using is then Eq. (E6) is $q = 0.68$M, with distillation accounting for 94.5% of these qubits.

Because of the very high fraction of qubits allocated to distillation when $C = C_{\text{min}}$, we consider a slowed down version of the algorithm where we set $C = 10C_{\text{min}}$, therefore easing the requirement on the number of T factories required. With the (ns, $10^{-4}$) qubit example, we use Eq. (E2) to find the distance increases to $d = 11$ since the logical qubits need to be protected for longer. The number of physical qubits per logical qubit is therefore $n(d) = 2d^2 = 242$ and $\tau(d) = (4t_{\text{gate}} + 2t_{\text{meas}})d = 4.4$ $\mu$s. This increases the algorithm run time via Eq. (E3) to be $t = 6.78$ s. Precisely the same factory is selected, namely a single round of logical distillation (15-to-1 space efficient), with one unit of code distance 9. Each factory produces $M(\mathcal{D}) = 1$ encoded T state with error rate $5.6 \cdot 10^{-11}$, and contains $n(\mathcal{D}) = 3240$ physical qubits, and requires $\tau(\mathcal{D}) = 46.8$ $\mu$s to run. However, now only $F = 17$ factories are required. The total number of physical qubits using is then Eq. (E6) drops to $q = 0.11$M, with distillation accounting for 49.7% of these qubits.

**Quantum chemistry**.— We assume the so-called 'double-factorized qubitization' algorithm described in Ref. [93, 140] is used. The *qubitization* approach is based on quantum phase estimation, but instead of constructing the standard $U = \exp(-iH/\alpha)$ from the Hamiltonian matrix $H$, one instead takes $U \approx \exp(-i\sin^{-1}(H/\alpha))$, which can typically be implemented with fewer resources by qubitization [94]. Using *double-factorization*, the naive representation of $H$ is compressed into fewer terms and also with a smaller $\alpha$.

We take $\epsilon = 0.01$. The algorithm, initially expressed as a Rust program, is explicitly compiled down to a QIR-level program with $Q_{\text{alg}} = 1318$ logical qubits with $M_T = 5.53 \cdot 10^7$, $D_R = 2.05 \cdot 10^8$, $M_R = 2.06 \cdot 10^8$, $M_{\text{Tof}} = 1.35 \cdot 10^{11}$ and $M_{\text{meas}} = 1.37 \cdot 10^9$. From Eq. (D1), Eq. (D3) and Eq. (D4), we calculate that with the PSSPC compilation and layout, the ISA-level executable that this

QIR program would compile down to has $Q = 2740$ logical qubits, which are needed for at least $C_{\min} = 4.10 \cdot 10^{11}$ logical time steps, consuming $M = 5.44 \cdot 10^{11}$ T states. We set $C = C_{\min}$.

At the physical level, here we illustrate the analysis with the (ns, $10^{-4}$) qubit example (with the other qubit parameter examples being analyzed analogously). We use Eq. (E2) to find the distance $d = 17$, such that $n(d) = 578$ and $\tau(d) = 6.8$ $\mu$s. This sets the algorithm run time via Eq. (E3) to be $t = 1$ month and one day. From Eq. (E4), the following distillation factory is selected: There are $F = 17$ factories. Each factory has two rounds of logical distillation, with 16 units of code distance 5 (15-to-1 space-eff.) in the first round and 1 unit of code distance 13 (15-to-1 RM prep.) in the second round. Each factory requires $n(\mathcal{D}) = 16000$ physical qubits, $\tau(\mathcal{D}) = 83.2$ $\mu$s run time, and produces $M(\mathcal{D}) = 1$ encoded T states with error rate $2.13 \cdot 10^{-15}$. The total number of physical qubits using Eq. (E6) is $q = 1.86$M.

**Factoring**.— We use a custom optimized implementation of Shor's algorithm [120]. We follow the logical algorithm implementation in Ref. [42], applying all of their optimizations (including coset representation, windowing over exponents and multiplicands, iterative phase estimation), except that we do not implement carry runways. As an input number to factorize, we use the RSA-2048 number from the RSA factoring challenge [69], and set the trial generator to 7. Further, we set both the windowing parameters $c_{\mathrm{mul}}$ and $c_{\exp}$ to 5.

First note that by the verifiable nature of the factoring problem, namely that once a solution has been found it can be easily verified, we can tolerate a large probability of the algorithm failing due to a logical fault since if the output is not valid the algorithm can be re-run. For this application example, we take $\epsilon = 1/3$.

In the example we consider of factoring a 2048-bit integer, the algorithm is initially expressed as a Rust program, and is explicitly compiled down to a QIR-level program with $Q_{\mathrm{alg}} = 12581$ algorithm qubits. Moreover, we find that $M_T = 12$, $D_R = 12$, $M_R = 12$, $M_{\mathrm{Tof}} = 3.73 \cdot 10^{10}$ and $M_{\mathrm{meas}} = 1.08 \cdot 10^9$. From Eq. (D1), Eq. (D3) and Eq. (D4), we calculate that with the PSSPC compilation and layout, the ISA-level executable that this QIR program would compile down to has $Q = 25481$ logical qubits, which are needed for at least $C_{\min} = 1.23 \cdot 10^{10}$ logical time steps, consuming $M = 1.49 \cdot 10^{10}$ T states.

At the physical level, here we illustrate the analysis with the (ns, $10^{-4}$) qubit example (with the other qubit parameter examples being analyzed analogously). We use Eq. (E2) to find the distance $d = 13$, such that $n(d) = 2d^2 = 338$ and $\tau(d) = (4t_{\mathrm{gate}} + 2t_{\mathrm{meas}})d = 5.2$ $\mu$s. This sets the algorithm run time via Eq. (E3) to be $t = 17$ hours 43 mins. From Eq. (E4), a distillation factory is selected. There are $F = 18$ factories. Each factory has two rounds of logical distillation (both 15-to-1 space efficient), with 16 units of code distance 3 in the first round and 1 unit of code distance 11 in the second round, which produces $M(\mathcal{D}) = 1$ encoded T state with error rate $5.51 \cdot 10^{-13}$, and which contains $n(\mathcal{D}) = 5760$ physical qubits, and requires $\tau(\mathcal{D}) = 72.8$ $\mu$s to run. The total number of physical qubits using Eq. (E6) is then $q = 8.72$M.

In the introduction, we state that factoring a 2048-bit integer using Shor's algorithm could be done in minutes with an array of twenty five thousand *perfect, noiseless* qubits. This estimate is obtained by assuming that the perfect qubits replace the ISA-level logical qubits in Table I. With perfect operations, there would be no need for distillation, and so the only cost would be the qubits 25481 algorithm and ancilla qubits in the array. Taking the $1.2 \cdot 10^{10}$ logical time steps to be physical time steps, each lasting 100 ns, results in a run time of 20 mins.

## Appendix G: Primary assumptions

Here we summarize our primary assumptions. These are the subset of assumptions which are not justified to our satisfaction in the current literature, and which we believe will significantly

increase our resource estimates if relaxed.

> **Uniform independent physical noise**.— We assume that the noise on physical qubits and physical qubit operations is the standard circuit noise model. In particular we assume error events at different space-time locations are independent and that error rates are uniform across the system in time and space.

> **Efficient classical computation**.— We assume that classical overhead (compilation, control, feedback, readout, decoding, etc.) does not dominate the overall cost of implementing the full quantum algorithm.

> **Syndrome measurement circuits for planar quantum ISA**.— We assume that syndrome measurement circuits with similar depth and error correction performance to those for standard surface and Hastings-Haah code patches can be constructed to implement all operations of the planar quantum ISA (as specified in Fig. 6).

> **Uniform independent logical noise**.— We assume that the error rate of a logical operation is approximately equal to its space-time volume (the number of tiles multiplied by the number of logical time steps) multiplied by the error rate of a logical qubit in a standard one-tile patch in one logical time step.

> **Negligible Clifford costs for synthesis**.— We assume that the space overhead for synthesis and space and time overhead for transport of T states within T state factories and to synthesis qubits are all negligible.

> **Smooth T state consumption rate**.— We assume that the rate of T state consumption throughout the compiled algorithm is almost constant, or can be made almost constant without significantly increasing the number of logical time steps for the algorithm.

[1] Distributed quantum computation architecture using semiconductor nanophotonics. *International Journal of Quantum Information*, 8(1-2):295–323, 2010.

[2] Ali J Abhari, Arvin Faruque, Mohammad J Dousti, Lukas Svec, Oana Catu, Amlan Chakrabati, Chen-Fu Chiang, Seth Vanderwilt, John Black, and Fred Chong. Scaffold: Quantum programming language. Technical report, Princeton Univ NJ Dept of Computer Science, 2012.

[3] Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C Bardin, Rami Barends, Rupak Biswas, Sergio Boixo, Fernando GSL Brandao, David A Buell, et al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779):505–510, 2019.

[4] Ryan Babbush, Jarrod McClean, Craig Gidney, Sergio Boixo, and Hartmut Neven. Focus beyond quadratic speedups for error-corrected quantum advantage, 2020.

[5] Dave Bacon. Operator quantum error-correcting subsystems for self-correcting quantum memories. *Physical Review A*, 73(1):012340, 2006.

[6] Bela Bauer, Sergey Bravyi, Mario Motta, and Garnet Kin-Lic Chan. Quantum algorithms for quantum chemistry and quantum materials science. *Chemical Reviews*, 120(22):12685–12717, 2020.

[7] Stephane Beauregard. Circuit for shor's algorithm using 2n+3 qubits. *Quantum Info. Comput.*, 3(2):175–185, mar 2003.

[8] Michael Beverland, Vadym Kliuchnikov, and Eddie Schoute. Surface code compilation via edge-disjoint paths. *PRX Quantum*, 3:020342, May 2022.

[9] Michael E Beverland, Oliver Buerschaper, Robert Koenig, Fernando Pastawski, John Preskill, and Sumit Sijher. Protected gates for topological quantum field theories. *Journal of Mathematical Physics*, 57(2):022201, 2016.

[10] Michael E. Beverland, Aleksander Kubica, and Krysta M. Svore. The cost of universality: A comparative study of the overhead of state distillation and code switching with color codes, 2021.

[11] Michael Edward Beverland. *Toward realizable quantum computers*. PhD thesis, California Institute of Technology, 2016.

[12] Dolev Bluvstein, Harry Levine, Giulia Semeghini, Tout T Wang, Sepehr Ebadi, Marcin Kalinowski, Alexander Keesling, Nishad Maskara, Hannes Pichler, Markus Greiner, et al. A quantum processor based on coherent transport of entangled atom arrays. *Nature*, 604(7906):451–456, 2022.

[13] Hector Bombin, Chris Dawson, Ryan V. Mishmash, Naomi Nickerson, Fernando Pastawski, and Sam Roberts. Logical blocks for fault-tolerant topological quantum computation, 2021.

[14] Hector Bombin, Isaac H Kim, Daniel Litinski, Naomi Nickerson, Mihir Pant, Fernando Pastawski, Sam Roberts, and Terry Rudolph. Interleaving: Modular architectures for fault-tolerant photonic quantum computing. *arXiv preprint arXiv:2103.08612*, 2021.

[15] Hector Bombin and Miguel Angel Martin-Delgado. Topological quantum distillation. *Physical review letters*, 97(18):180501, 2006.

[16] S. B. Bravyi and A. Yu. Kitaev. Quantum codes on a lattice with boundary, 1998.

[17] Sergey Bravyi, Oliver Dial, Jay M. Gambetta, Dario Gil, and Zaira Nazario. The future of quantum computing with superconducting qubits. 2022.

[18] Sergey Bravyi and Alexei Kitaev. Universal quantum computation with ideal clifford gates and noisy ancillas. *Physical Review A*, 71(2):022316, 2005.

[19] Sergey Bravyi and Robert König. Classification of topologically protected gates for local stabilizer codes. *Phys. Rev. Lett.*, 110:170503, Apr 2013.

[20] Sergey Bravyi, David Poulin, and Barbara Terhal. Tradeoffs for reliable quantum information storage in 2d systems. *Phys. Rev. Lett.*, 104:050503, Feb 2010.

[21] Jacques Carolan, Christopher Harrold, Chris Sparrow, Enrique Martín-López, Nicholas J Russell, Joshua W Silverstone, Peter J Shadbolt, Nobuyuki Matsuda, Manabu Oguma, Mikitaka Itoh, et al. Universal linear optics. *Science*, 349(6249):711–716, 2015.

[22] Christopher Chamberland and Earl T. Campbell. Universal quantum computing with twist-free and temporally encoded lattice surgery. 2021.

[23] Christopher Chamberland and Earl T Campbell. Circuit-level protocol and analysis for twist-based lattice surgery. *Physical Review Research*, 4:023090, May 2022.

[24] Christopher Chamberland and Andrew W. Cross. Fault-tolerant magic state preparation with flag qubits. *Quantum*, 3:143, May 2019.

[25] Christopher Chamberland, Kyungjoo Noh, Patricio Arrangoiz-Arriola, Earl T. Campbell, Connor T. Hann, Joseph Iverson, Harald Putterman, Thomas C. Bohdanowicz, Steven T. Flammia, Andrew Keller, Gil Refael, John Preskill, Liang Jiang, Amir H. Safavi-Naeini, Oskar Painter, and Fernando G.S.L. Brandão. Building a fault-tolerant quantum computer using concatenated cat codes. *PRX Quantum*, 3:010329, Feb 2022.

[26] Christopher Chamberland, Kyungjoo Noh, Patricio Arrangoiz-Arriola, Earl T. Campbell, Connor T. Hann, Joseph Iverson, Harald Putterman, Thomas C. Bohdanowicz, Steven T. Flammia, Andrew Keller, Gil Refael, John Preskill, Liang Jiang, Amir H. Safavi-Naeini, Oskar Painter, and Fernando G. S. L. Brandão. Building a fault-tolerant quantum computer using concatenated cat codes, 2020.

[27] Rui Chao, Michael E. Beverland, Nicolas Delfosse, and Jeongwan Haah. Optimization of the surface code design for Majorana-based qubits. *Quantum*, 4:352, October 2020.

[28] Andrew M. Childs, Eddie Schoute, and Cem M. Unsal. Circuit transformations for quantum architectures. In Wim van Dam and Laura Mancinska, editors, *14th Conference on the Theory of Quantum Computation, Communication and Cryptography (TQC 2019)*, volume 135 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 3:1–3:24, Dagstuhl, Germany, 2019. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.

[29] J. I. Cirac and P. Zoller. Quantum computations with cold trapped ions. *Phys. Rev. Lett.*, 74:4091–4094, May 1995.

[30] Antonio D. Córcoles, Abhinav Kandala, Ali Javadi-Abhari, Douglas T. McClure, Andrew W. Cross, Kristan Temme, Paul D. Nation, Matthias Steffen, and Jay M. Gambetta. Challenges and opportunities of near-term quantum computing systems. *Proceedings of the IEEE*, 108(8):1338–1352, 2020.

[31] Andrew J Daley, Immanuel Bloch, Christian Kokail, Stuart Flannigan, Natalie Pearson, Matthias Troyer, and Peter Zoller. Practical quantum advantage in quantum simulation. *Nature*, 607(7920):667–

676, 2022.

[32] Niel de Beaudrap and Steven Herbert. Quantum linear network coding for entanglement distribution in restricted architectures. 4:356.

[33] Nicolas Delfosse, Michael E Beverland, and Maxime A Tremblay. Bounds on stabilizer measurement circuits and obstructions to local implementations of quantum ldpc codes. *arXiv preprint arXiv:2109.14599*, 2021.

[34] Alain Delgado, Pablo A. M. Casares, Roberto dos Reis, Modjtaba Shokrian Zini, Roberto Campos, Norge Cruz-Hernández, Arne-Christian Voigt, Angus Lowe, Soran Jahangiri, M. A. Martin-Delgado, Jonathan E. Mueller, and Juan Miguel Arrazola. Simulating key properties of lithium-ion batteries with a fault-tolerant quantum computer. *Phys. Rev. A*, 106:032428, Sep 2022.

[35] M. H. Devoret and R. J. Schoelkopf. Superconducting circuits for quantum information: An outlook. *Science*, 339(6124):1169–1174, 2013.

[36] David P. Divincenzo. The Physical Implementation of Quantum Computation. *Fortschritte der Physik*, 48(9-11):771–783, January 2000.

[37] Rajeev Acharya et al. Suppressing quantum errors by scaling a surface code logical qubit, 2022.

[38] Austin G. Fowler, Matteo Mariantoni, John M. Martinis, and Andrew N. Cleland. Surface codes: Towards practical large-scale quantum computation. *Phys. Rev. A*, 86:032324, Sep 2012.

[39] D.P. Franke, J.S. Clarke, L.M.K. Vandersypen, and M. Veldhorst. Rent's rule and extensibility in quantum computing. *Microprocessors and Microsystems*, 67:1–7, 2019.

[40] X. Fu, L. Riesebos, M. A. Rol, J. van Straten, J. van Someren, N. Khammassi, I. Ashraf, R. F. L. Vermeulen, V. Newsum, K. K. L. Loh, J. C. de Sterke, W. J. Vlothuizen, R. N. Schouten, C. G. Almudever, L. DiCarlo, and K. Bertels. eqasm: An executable quantum instruction set architecture. 2018.

[41] Craig Gidney. Halving the cost of quantum addition. *Quantum*, 2:74, June 2018.

[42] Craig Gidney and Martin Ekerå. How to factor 2048 bit RSA integers in 8 hours using 20 million noisy qubits. *Quantum*, 5:433, April 2021.

[43] Craig Gidney and Austin G. Fowler. Flexible layout of surface code computations using autoccz states, 2019.

[44] Craig Gidney, Michael Newman, and Matt McEwen. Benchmarking the Planar Honeycomb Code. *Quantum*, 6:813, September 2022.

[45] John Goodacre. Technology Preview: The ARMv8 Architecture, 2011.

[46] D Gottesman. The heisenberg representation of quantum computers. 6 1998.

[47] Daniel Gottesman. Theory of fault-tolerant quantum computation. *Physical Review A*, 57(1):127, 1998.

[48] Daniel Gottesman. An introduction to quantum error correction and fault-tolerant quantum computation. In *Quantum information science and its contributions to mathematics, Proceedings of Symposia in Applied Mathematics*, volume 68, pages 13–58, 2010.

[49] Elie Gouzien and Nicolas Sangouard. Factoring 2048-bit rsa integers in 177 days with 13 436 qubits and a multimode memory. *Physical Review Letters*, 127(14):140503, 2021.

[50] Alexander S Green, Peter LeFanu Lumsdaine, Neil J Ross, Peter Selinger, and Benoît Valiron. Quipper: a scalable quantum programming language. In *Proceedings of the 34th ACM SIGPLAN conference on Programming language design and implementation*, pages 333–342, 2013.

[51] Jérémie Guillaud and Mazyar Mirrahimi. Repetition cat qubits for fault-tolerant quantum computation. *Phys. Rev. X*, 9:041053, Dec 2019.

[52] Larry Guth and Alexander Lubotzky. Quantum error correcting codes and 4-dimensional arithmetic hyperbolic manifolds. *Journal of Mathematical Physics*, 55(8):082202, 2014.

[53] Frederik Hahn, A. Pappa, and Jens Eisert. Quantum network routing and local complementation. *npj Quantum Information*, 5:1–7, 09 2019.

[54] R. Hanson, L. P. Kouwenhoven, J. R. Petta, S. Tarucha, and L. M. K. Vandersypen. Spins in few-electron quantum dots. *Rev. Mod. Phys.*, 79:1217–1265, Oct 2007.

[55] Matthew B. Hastings and Jeongwan Haah. Dynamically Generated Logical Qubits. *Quantum*, 5:564, October 2021.

[56] Naomichi Hatano and Masuo Suzuki. *Finding Exponential Product Formulas of Higher Orders*, pages 37–68. Springer Berlin Heidelberg, Berlin, Heidelberg, 2005.

[57] Sabrina S Hong, Alexander T Papageorge, Prasahnt Sivarajah, Genya Crossman, Nicolas Didier,

Anthony M Polloreno, Eyob A Sete, Stefan W Turkowski, Marcus P da Silva, and Blake R Johnson. Demonstration of a parametrically activated entangling gate protected from flux noise. *Physical Review A*, 101(1):012302, 2020.

[58] Clare Horsman, Austin G Fowler, Simon Devitt, and Rodney Van Meter. Surface code quantum computing by lattice surgery. *New Journal of Physics*, 14(12):123011, dec 2012.

[59] Thomas Häner, Vadym Kliuchnikov, Martin Roetteler, and Mathias Soeken. Space-time optimized table lookup, 2022.

[60] Intel. Intel 64 and IA-32 Architectures Software Developer's Manual. https://www.intel.co.uk/content/www/uk/en/architecture-and-technology/64-ia-32-architectures-software-developer-vol-1-manual.html, 2016.

[61] David Ittah, Thomas Häner, Vadym Kliuchnikov, and Torsten Hoefler. Qiro: A static single assignment-based quantum program representation for optimization. *ACM Transactions on Quantum Computing*, 3(3), jun 2022.

[62] D. Jaksch, H.-J. Briegel, J. I. Cirac, C. W. Gardiner, and P. Zoller. Entanglement of atoms via cold controlled collisions. *Phys. Rev. Lett.*, 82:1975–1978, Mar 1999.

[63] Ali Javadi-Abhari, Pranav Gokhale, Adam Holmes, Diana Franklin, Kenneth R. Brown, Margaret Martonosi, and Frederic T. Chong. Optimized surface code communication in superconducting quantum computers. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*. ACM.

[64] Ali JavadiAbhari, Shruti Patil, Daniel Kudrow, Jeff Heckey, Alexey Lvov, Frederic T Chong, and Margaret Martonosi. Scaffcc: A framework for compilation and analysis of quantum computing programs. In *Proceedings of the 11th ACM Conference on Computing Frontiers*, pages 1–10, 2014.

[65] Hamza Jnane, Brennan Undseth, Zhenyu Cai, Simon C. Benjamin, and Bálint Koczor. Multicore quantum computing. *Phys. Rev. Applied*, 18:044064, Oct 2022.

[66] Cody Jones, Michael A Fogarty, Andrea Morello, Mark F Gyure, Andrew S Dzurak, and Thaddeus D Ladd. Logical qubit in a linear array of semiconductor quantum dots. *Physical Review X*, 8(2):021058, 2018.

[67] Stephen Jordan et al. Quantum algorithm zoo. *Retrieved June*, 27:2013, 2011.

[68] Petar Jurcevic, Ali Javadi-Abhari, Lev S Bishop, Isaac Lauer, Daniela F Bogorin, Markus Brink, Lauren Capelluto, Oktay Günlük, Toshinari Itoko, Naoki Kanazawa, Abhinav Kandala, George A Keefe, Kevin Krsulich, William Landers, Eric P Lewandowski, Douglas T McClure, Giacomo Nannicini, Adinath Narasgond, Hasan M Nayfeh, Emily Pritchett, Mary Beth Rothwell, Srikanth Srinivasan, Neereja Sundaresan, Cindy Wang, Ken X Wei, Christopher J Wood, Jeng-Bang Yau, Eric J Zhang, Oliver E Dial, Jerry M Chow, and Jay M Gambetta. Demonstration of quantum volume 64 on a superconducting quantum computing system. *Quantum Science and Technology*, 6(2):025020, mar 2021.

[69] Burt Kaliski. Announcement of rsa factoring challenge. *RSA Laboratories*, 1991.

[70] Bruce E Kane. A silicon-based nuclear spin quantum computer. *nature*, 393(6681):133–137, 1998.

[71] Torsten Karzig, Christina Knapp, Roman M. Lutchyn, Parsa Bonderson, Matthew B. Hastings, Chetan Nayak, Jason Alicea, Karsten Flensberg, Stephan Plugge, Yuval Oreg, Charles M. Marcus, and Michael H. Freedman. Scalable designs for quasiparticle-poisoning-protected topological quantum computation with majorana zero modes. *Phys. Rev. B*, 95:235305, Jun 2017.

[72] Isaac H Kim, Ye-Hua Liu, Sam Pallister, William Pol, Sam Roberts, and Eunseok Lee. Fault-tolerant resource estimate for quantum chemical simulations: Case study on li-ion battery electrolyte molecules. *Physical Review Research*, 4(2):023019, 2022.

[73] Isaac H. Kim, Ye-Hua Liu, Sam Pallister, William Pol, Sam Roberts, and Eunseok Lee. Fault-tolerant resource estimate for quantum chemical simulations: Case study on li-ion battery electrolyte molecules. *Phys. Rev. Research*, 4:023019, Apr 2022.

[74] A Yu Kitaev. Unpaired majorana fermions in quantum wires. *Physics-uspekhi*, 44(10S):131, 2001.

[75] A.Yu. Kitaev. Fault-tolerant quantum computation by anyons. *Annals of Physics*, 303(1):2–30, 2003.

[76] Morten Kjaergaard, Mollie E Schwartz, Jochen Braumüller, Philip Krantz, Joel I-J Wang, Simon Gustavsson, and William D Oliver. Superconducting qubits: Current state of play. *Annual Review of Condensed Matter Physics*, 11:369–395, 2020.

[77] Rochus Klesse and Sandra Frank. Quantum error correction in spatially correlated quantum noise. *Physical review letters*, 95(23):230503, 2005.

[78] Vadym Kliuchnikov, Kristin Lauter, Romy Minko, Adam Paetznick, and Christophe Petit. Shorter quantum circuits, 2022.

[79] Emanuel Knill. Fault-tolerant postselected quantum computation: Schemes. *arXiv preprint quant-ph/0402171*, 2004.

[80] Emanuel Knill. Quantum computing with realistically noisy devices. *Nature*, 434(7029):39–44, 2005.

[81] Emanuel Knill, Raymond Laflamme, and Gerald J Milburn. A scheme for efficient quantum computation with linear optics. *nature*, 409(6816):46–52, 2001.

[82] Jens Koch, Terri M. Yu, Jay Gambetta, A. A. Houck, D. I. Schuster, J. Majer, Alexandre Blais, M. H. Devoret, S. M. Girvin, and R. J. Schoelkopf. Charge-insensitive qubit design derived from the cooper pair box. *Phys. Rev. A*, 76:042319, Oct 2007.

[83] L. Lao, B. van Wee, I. Ashraf, J. van Someren, N. Khammassi, K. Bertels, and C. G. Almudever. Mapping of lattice surgery-based quantum circuits on surface code architectures. 4(1):015005.

[84] C. Lattner and V. Adve. llvm: a compilation framework for lifelong program analysis and transformation. In *International Symposium on Code Generation and Optimization, 2004. CGO 2004.*

[85] Bjoern Lekitsch, Sebastian Weidt, Austin G. Fowler, Klaus Mølmer, Simon J. Devitt, Christof Wunderlich, and Winfried K. Hensinger. Blueprint for a microwave trapped ion quantum computer. *Science Advances*, 3(2):e1601540, 2017.

[86] Debbie Leung, Jonathan Oppenheim, and Andreas Winter. Quantum network communication—the butterfly and beyond. *IEEE Transactions on Information Theory*, 56(7):3478–3490, 7 2010.

[87] Pak Hong Leung, Kevin A. Landsman, Caroline Figgatt, Norbert M. Linke, Christopher Monroe, and Kenneth R. Brown. Robust 2-qubit gates in a linear ion crystal using a frequency-modulated driving force. *Phys. Rev. Lett.*, 120:020501, Jan 2018.

[88] Anthony Leverrier, Jean-Pierre Tillich, and Gilles Zémor. Quantum expander codes. In *2015 IEEE 56th Annual Symposium on Foundations of Computer Science*, pages 810–824. IEEE, 2015.

[89] Daniel Litinski. A Game of Surface Codes: Large-Scale Quantum Computing with Lattice Surgery. *Quantum*, 3:128, March 2019.

[90] Daniel Litinski. Magic State Distillation: Not as Costly as You Think. *Quantum*, 3:205, December 2019.

[91] Seth Lloyd. Universal quantum simulators. *Science*, 273(5278):1073–1078, 1996.

[92] Julie Love. Myth vs. reality: a practical perspective on quantum computing, 2020.

[93] Guang Hao Low. Halving the cost of quantum multiplexed rotations, 2021.

[94] Guang Hao Low and Isaac L. Chuang. Hamiltonian Simulation by Qubitization. *Quantum*, 3:163, July 2019.

[95] Alexander McCaskey and Thien Nguyen. A mlir dialect for quantum assembly languages. In *2021 IEEE International Conference on Quantum Computing and Engineering (QCE)*, pages 255–264. IEEE, 2021.

[96] Giulia Meuli, Mathias Soeken, Martin Roetteler, and Thomas Häner. Enabling accuracy-aware quantum compilers using symbolic resource estimation. *Proc. ACM Program. Lang.*, 4(OOPSLA), nov 2020.

[97] Microsoft. Azure Quantum Resource Estimator. `https://aka.ms/AQ/RE`, 2022.

[98] Microsoft. Resource estimation sample notebooks. `https://github.com/microsoft/Quantum/tree/main/samples/azure-quantum/resource-estimation`, 2022.

[99] Klaus Mølmer and Anders Sørensen. Multiparticle entanglement of hot trapped ions. *Phys. Rev. Lett.*, 82:1835–1838, Mar 1999.

[100] Prakash Murali, Jonathan M. Baker, Ali Javadi Abhari, Frederic T. Chong, and Margaret Martonosi. Noise-adaptive compiler mappings for noisy intermediate-scale quantum computers. pages 1015–1029, 2019.

[101] Prakash Murali, Dripto M. Debroy, Kenneth R. Brown, and Margaret Martonosi. Architecting noisy intermediate-scale trapped ion quantum computers. In *Proceedings of the ACM/IEEE 47th Annual International Symposium on Computer Architecture*, ISCA '20, page 529–542. IEEE Press, 2020.

[102] Prakash Murali, Norbert Matthias Linke, Margaret Martonosi, Ali Javadi Abhari, Nhung Hong Nguyen, and Cinthia Huerta Alderete. Full-stack, real-system quantum computer studies: Architectural comparisons and design insights. In *Proceedings of the 46th International Symposium on Computer Architecture*, ISCA '19, page 527–540, New York, NY, USA, 2019. Association for Computing Machinery.

[103] Adam Paetznick, Christina Knapp, Nicolas Delfosse, Bela Bauer, Jeongwan Haah, Matthew B. Hastings, and Marcus P. da Silva. Performance of planar floquet codes with majorana-based qubits, 2022.

[104] Pavel Panteleev and Gleb Kalachev. Asymptotically good quantum and locally testable classical ldpc codes. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 375–388, 2022.

[105] Christopher A Pattison, Michael E Beverland, Marcus P da Silva, and Nicolas Delfosse. Improved quantum error correction using soft information. *arXiv preprint arXiv:2107.13589*, 2021.

[106] Jennifer Paykin, Robert Rand, and Steve Zdancewic. Qwire: a core language for quantum circuits. *ACM SIGPLAN Notices*, 52(1):846–858, 2017.

[107] Natalie Pearson. Simulating many-body quantum systems: Quantum algorithms and experimental realisation. 2020.

[108] Juan M Pino, Jennifer M Dreiling, Caroline Figgatt, John P Gaebler, Steven A Moses, MS Allman, CH Baldwin, M Foss-Feig, D Hayes, K Mayer, et al. Demonstration of the trapped-ion quantum ccd computer architecture. *Nature*, 592(7853):209–213, 2021.

[109] S. Plugge, L. A. Landau, E. Sela, A. Altland, K. Flensberg, and R. Egger. Roadmap to majorana surface codes. *Phys. Rev. B*, 94:174514, Nov 2016.

[110] John Preskill. Quantum Computing in the NISQ era and beyond. *Quantum*, 2:79, August 2018.

[111] QIR Alliance. QIR Alliance. https://www.qir-alliance.org/, 2022.

[112] Yihui Quek, Daniel Stilck França, Sumeet Khatri, Johannes Jakob Meyer, and Jens Eisert. Exponentially tighter bounds on limitations of quantum error mitigation. *arXiv preprint arXiv:2210.11505*, 2022.

[113] Markus Reiher, Nathan Wiebe, Krysta M. Svore, Dave Wecker, and Matthias Troyer. Elucidating reaction mechanisms on quantum computers. *Proceedings of the National Academy of Sciences*, 114(29):7555–7560, 2017.

[114] RISCV. RISC-V instruction set architecture (ISA) and related specifications, 2021.

[115] Martin Roetteler, Michael Naehrig, Krysta M. Svore, and Kristin Lauter. Quantum resource estimates for computing elliptic curve discrete logarithms. In Tsuyoshi Takagi and Thomas Peyrin, editors, *Advances in Cryptology – ASIACRYPT 2017*, pages 241–270, Cham, 2017. Springer International Publishing.

[116] M Saffman. Quantum computing with atomic qubits and rydberg interactions: progress and challenges. *Journal of Physics B: Atomic, Molecular and Optical Physics*, 49(20):202001, oct 2016.

[117] Sankar Das Sarma, Michael Freedman, and Chetan Nayak. Majorana zero modes and topological quantum computation. *npj Quantum Information*, 1(1):1–13, 2015.

[118] Joseph A Schreier, Andrew A Houck, Jens Koch, David I Schuster, Bradley R Johnson, Jerry M Chow, Jay M Gambetta, J Majer, Luigi Frunzio, Michel H Devoret, et al. Suppressing charge noise decoherence in superconducting charge qubits. *Physical Review B*, 77(18):180502, 2008.

[119] Yunong Shi, Nelson Leung, Pranav Gokhale, Zane Rossi, David I. Schuster, Henry Hoffmann, and Frederic T. Chong. Optimized Compilation of Aggregated Instructions for Realistic Quantum Computers. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '19, page 1031–1044, New York, NY, USA, 2019. Association for Computing Machinery.

[120] P. W. Shor. Algorithms for quantum computation: discrete logarithms and factoring. In *Proceedings 35th Annual Symposium on Foundations of Computer Science*, pages 124–134, 1994.

[121] Peter W Shor. Scheme for reducing decoherence in quantum computer memory. *Physical review A*, 52(4):R2493, 1995.

[122] Peter W Shor. Fault-tolerant quantum computation. In *Proceedings of 37th conference on foundations of computer science*, pages 56–65. IEEE, 1996.

[123] Robert S. Smith, Michael J. Curtis, and William J. Zeng. A Practical Quantum Instruction Set Architecture, 2016.

[124] Thomas M. Stace, Sean D. Barrett, and Andrew C. Doherty. Thresholds for topological codes in the presence of loss. *Phys. Rev. Lett.*, 102:200501, May 2009.

[125] Andrew Steane. Multiple-particle interference and quantum error correction. *Proceedings of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 452(1954):2551–2577, 1996.

[126] Andrew M Steane. Space, time, parallelism and noise requirements for reliable quantum computing.

*Fortschritte der Physik: Progress of Physics*, 46(4-5):443–457, 1998.

[127] M. Steffen, D. P. DiVincenzo, J. M. Chow, T. N. Theis, and M. B. Ketchen. Quantum computing: An ibm perspective. *IBM Journal of Research and Development*, 55(5):13:1–13:11, 2011.

[128] Damian S. Steiger, Thomas Häner, and Matthias Troyer. Advantages of a modular high-level quantum programming framework. 66:81–89.

[129] Ashley M Stephens and Zachary WE Evans. Accuracy threshold for concatenated error detection in one dimension. *Physical Review A*, 80(2):022313, 2009.

[130] Armands Strikis, Simon C Benjamin, and Benjamin J Brown. Quantum computing is scalable on a planar array of qubits with fabrication defects. *arXiv preprint arXiv:2111.06432*, 2021.

[131] Martin Suchara, John Kubiatowicz, Arvin Faruque, Frederic T. Chong, Ching-Yi Lai, and Gerardo Paz. Qure: The quantum resource estimator toolbox. In *2013 IEEE 31st International Conference on Computer Design (ICCD)*, pages 419–426, 2013.

[132] Krysta Svore, Alan Geller, Matthias Troyer, John Azariah, Christopher Granade, Bettina Heim, Vadym Kliuchnikov, Mariia Mykhailova, Andres Paz, and Martin Roetteler. Q# enabling scalable quantum computing and development with a high-level dsl. In *Proceedings of the real world domain specific languages workshop 2018*, pages 1–10, 2018.

[133] Krysta M Svore, Barbara M Terhal, and David P DiVincenzo. Local fault-tolerant quantum computation. *Physical Review A*, 72(2):022317, 2005.

[134] Thomas Szkopek, P Oscar Boykin, Heng Fan, Vwani P Roychowdhury, Eli Yablonovitch, Geoffrey Simms, Mark Gyure, and Bryan Fong. Threshold error penalty for fault-tolerant quantum computation with nearest neighbor communication. *IEEE transactions on nanotechnology*, 5(1):42–49, 2006.

[135] Jean-Pierre Tillich and Gilles Zémor. Quantum ldpc codes with positive rate and minimum distance proportional to the square root of the blocklength. *IEEE Transactions on Information Theory*, 60(2):1193–1202, 2014.

[136] Alan Tran, Alex Bocharov, Bela Bauer, and Parsa Bonderson. Optimizing Clifford gate generation for measurement-only topological quantum computation with Majorana zero modes. *SciPost Phys.*, 8:91, 2020.

[137] Maxime A Tremblay, Nicolas Delfosse, and Michael E Beverland. Constant-overhead quantum error correction with thin planar connectivity. *Physical Review Letters*, 129(5):050504, 2022.

[138] Matthias Troyer. Disentangling hype from reality: Achieving practical quantum advantage. Q2B Practical Quantum Computing Conference, 2020.

[139] David K. Tuckett, Stephen D. Bartlett, Steven T. Flammia, and Benjamin J. Brown. Fault-tolerant thresholds for the surface code in excess of 5% under biased noise. *Phys. Rev. Lett.*, 124:130501, Mar 2020.

[140] Vera von Burg, Guang Hao Low, Thomas Häner, Damian S. Steiger, Markus Reiher, Martin Roetteler, and Matthias Troyer. Quantum computing enhanced computational catalysis. *Phys. Rev. Research*, 3:033055, Jul 2021.

[141] David S. Wang, Austin G. Fowler, and Lloyd C. L. Hollenberg. Surface code quantum computing with error rates over 1 *Phys. Rev. A*, 83:020302, Feb 2011.

[142] Ye Wang, Mark Um, Junhua Zhang, Shuoming An, Ming Lyu, Jing-Ning Zhang, L.-M. Duan, Dahyun Yum, and Kihwan Kim. Single-qubit quantum memory exceeding ten-minute coherence time. *Nature Photonics*, 11(10):646–650, Oct 2017.

[143] Alwin Zulehner and Robert Wille. Compiling SU(4) quantum circuits to IBM QX architectures. pages 185–190.