

# Layered architecture for quantum computing

N. Cody Jones,<sup>1,\*</sup> Rodney Van Meter,<sup>2</sup> Austin G. Fowler,<sup>3</sup> Peter L. McMahon,<sup>1</sup> Jungsang Kim,<sup>4</sup> Thaddeus D. Ladd,<sup>1,5</sup> and Yoshihisa Yamamoto<sup>1,5</sup>

<sup>1</sup>*Edward L. Ginzton Laboratory, Stanford University, Stanford, California 94305-4088, USA*

<sup>2</sup>*Faculty of Environment and Information Studies, Keio University, Japan*

<sup>3</sup>*Centre for Quantum Computation and Communication Technology, University of Melbourne, Victoria, Australia*

<sup>4</sup>*Fitzpatrick Institute for Photonics, Duke University, Durham, NC, USA*

<sup>5</sup>*National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo 101-8430, Japan*

We develop a layered quantum computer architecture, which is a systematic framework for tackling the individual challenges of developing a quantum computer while constructing a cohesive device design. We discuss many of the prominent techniques for implementing circuit-model quantum computing and introduce several new methods, with an emphasis on employing surface code quantum error correction. In doing so, we propose a new quantum computer architecture based on optical control of quantum dots. The timescales of physical hardware operations and logical, error-corrected quantum gates differ by several orders of magnitude. By dividing functionality into layers, we can design and analyze subsystems independently, demonstrating the value of our layered architectural approach. Using this concrete hardware platform, we provide resource analysis for executing fault-tolerant quantum algorithms for integer factoring and quantum simulation, finding that the quantum dot architecture we study could solve such problems on the timescale of days.

## I. INTRODUCTION

Quantum computing as an engineering discipline is still in its infancy. Although the physics is well understood, developing devices which compute with quantum mechanics is technologically daunting. While experiments to date manipulate only a handful of quantum bits [1], we consider what effort is required to build a large-scale quantum computer. This objective demands more than a cursory estimate of the number of qubits and gates required for a given algorithm. One must consider the faulty quantum hardware, with errors caused by both the environment and deliberate control operations; when error correction is invoked, classical processing is required; constructing arbitrary gate sequences from a limited fault-tolerant set requires special treatment, and so on. This paper provides a framework to address the complete challenge of designing a quantum computer.

Many researchers have presented and examined components of large-scale quantum computing. We study here how these components may be combined in an efficient design, and we introduce new methods which improve the quantum computer we propose. This engineering pursuit is quantum computer architecture, which we develop here in layers. An *architecture* decomposes complex system behaviors into a manageable set of operations. A *layered architecture* does this through layers of abstraction where each embodies a critical set of related functions. For our purposes, each ascending layer brings the system closer to an ideal quantum computing environment.

The paper is organized as follows. The remainder of Section I provides a global view of a layered quantum

computer architecture, indicating how each of the topics we examine are connected. Section II enumerates the essential components of a quantum computer by examining a new hardware platform based on the optical control of quantum dots. Section III discusses control techniques for suppressing hardware errors prior to using active error correction. Section IV demonstrates how to implement and account for the resources of quantum error correction, with particular emphasis on the surface code [2]. Section V analyzes the necessary techniques for constructing universal quantum gates from the limited set of operations provided by error correction. Section VI calculates the computer resources necessary to implement two prominent quantum algorithms: integer factoring and quantum simulation. Section VII discusses timing issues which affect how the layers in the architecture interact with each other. Finally, Section VIII discusses how our findings are applicable to future work in quantum computing.

### A. Prior work on quantum computer architecture

Many different quantum computing technologies are under experimental investigation [1], but for each a scalable system architecture remains an open research problem. Since DiVincenzo introduced his fundamental criteria for a viable quantum computing technology [3] and Steane emphasized the difficulty of designing systems capable of running quantum error correction (QEC) adequately [4, 5], several groups of researchers have outlined various additional taxonomies addressing the architectural needs of large-scale systems [6, 7]. As an example, small-scale interconnects have been proposed for many technologies, but the problems of organizing subsystems using these techniques into a complete architecture for a large-scale system have been addressed by only a few

---

\* ncodyjones@gmail.com

researchers. In particular, the issue of heterogeneity in system architecture has received relatively little attention.

The most important subroutine in fault-tolerant quantum computers considered thus far is the preparation of ancilla states for fault-tolerant circuits, because these circuits often require very many ancillas. Taylor *et al.* proposed a design with alternating “ancilla blocks” and “data blocks” in the device layout [8]. Steane introduced the idea of “factories” for creating ancillas [9], which we examine for the case of the surface code in this work. Isailovic *et al.* [10] studied this problem for ion trap architectures and found that, for typical quantum circuits, approximately 90% of the quantum computer must be devoted to such factories in order to calculate “at the speed of data,” or where ancilla-production is not the rate-limiting process. The findings we present here are in close agreement with this estimate. Metodi *et al.* also considered production of ancillas in ion trap designs, focusing instead on a 3-qubit ancilla state used for the Toffoli gate [11], which is an alternative pathway to a universal fault-tolerant set of gates.

Some researchers have studied the difficulty of moving data in a quantum processor. Kiepinski *et al.* proposed a scalable ion trap technology utilizing separate memory and computing areas [12]. Because quantum error correction requires rapid cycling across all physical qubits in the system, this approach is best used as a unit cell replicated across a larger system. Other researchers have proposed homogeneous systems built around this basic concept. One common structure is a recursive H tree, which works well with a small number of layers of a Calderbank-Shor-Steane (CSS) code, targeted explicitly at ion trap systems [13, 14]. Oskin *et al.* [15], building on the Kane solid-state NMR technology [16], proposed a loose lattice of sites, explicitly considering the issues of classical control and movement of quantum data in scalable systems, but without a specific plan for QEC. In the case of quantum computing with superconducting circuits, the quantum von Neumann architecture specifically considers dedicated hardware for quantum memories, zeroing registers, and a quantum bus [17].

Long-range coupling and communication is a significant challenge for quantum computers. Cirac *et al.* proposed the use of photonic qubits to distribute entanglement between distant atoms [18], and other researchers have investigated the prospects for optically-mediated nonlocal gates [19–23]. Such photonic channels could be utilized to realize a modular, scalable distributed quantum computer [24]. Conversely, Metodi *et al.* consider how to use local gates and quantum teleportation to move logical qubits throughout their ion-trap QLA architecture [11]. Fowler *et al.* [25] investigated a Josephson junction flux qubit architecture considering the extreme difficulties of routing both the quantum couplers and large numbers of classical control lines, producing a structure with support for CSS codes and logical qubits organized in a line. Whitney *et al.* [26, 27] have investigated auto-

mated layout and optimization of circuit designs specifically for ion trap architectures, and Isailovic *et al.* [10, 28] have studied interconnection and data throughput issues in similar ion trap systems, with an emphasis on preparing ancillas for teleportation gates [29].

Other work has studied quantum computer architectures with only nearest-neighbor coupling between qubits in an array [30–34], which is appealing from a hardware design perspective. With the recent advances in the operation of the topological codes and their desirable characteristics such as having a high practical threshold and requiring only nearest-neighbor interactions, research effort has shifted toward architectures capable of building and maintaining large two- and three-dimensional cluster states [35–38]. These systems rely on topological error correction models [39], whose higher tolerance to error often comes at the cost of a larger physical system, relative to, for example, implementations based on the Steane code [40]. The surface code [2], which we examine in this work for its impact on architecture, belongs to the topological family of codes.

Recent attention has been directed at distributed models of quantum computing. Devitt *et al.* studied how to distribute a photonic cluster-state quantum computing network over different geographic regions [41]. The abstract framework of a quantum multicomputer recognizes that large-scale systems demand heterogeneous interconnects [42]; in most quantum computing technologies, it may not be possible to build monolithic systems that contain, couple, and control billions of physical qubits. Van Meter *et al.* [43] extended this architectural framework with a design based on nanophotonic coupling of electron spin quantum dots that explicitly uses multiple levels of interconnect with varying coupling fidelities (resulting in varying purification requirements), as well as the ability to operate with a very low yield of functional devices. Although that proposed system has many attractive features, concerns about the difficulty of fabricating adequately high quality optical components and the desire to reduce the surface code lattice cycle time led to the system design proposed in this paper.

## B. Layered framework

A good architecture must have a simple structure while also efficiently managing the complex array of resources in a quantum computer. Layered architectures are a conventional approach to solving such engineering problems in many fields of information technology, and Ref. [14] presents a layered architecture for quantum computer design software. Our architecture, which describes the physical design of the quantum computer, consists of five layers, where each layer has a prescribed set of duties to accomplish. The interface between two layers is defined by the services a lower layer provides to the one above it. To execute an operation, a layer must issue commands to the layer below and process the results. Designing a sys-

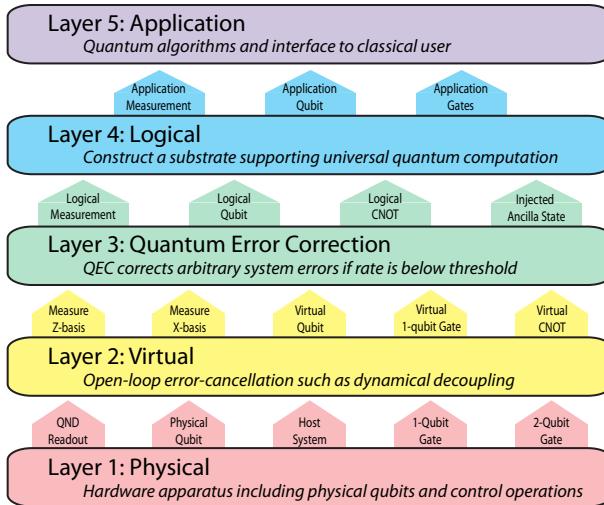


FIG. 1. Color. Layered control stack which forms the framework of a quantum computer architecture. Vertical arrows indicate services provided to a higher layer.

tem this way ensures that related operations are grouped together and that the system organization is hierarchical. Such an approach allows quantum engineers to focus on individual challenges, while also seeing how a process fits into the overall design. By organizing the architecture in layers, we deliberately create a *modular* design for the quantum computer.

The layered framework can be understood by a control stack composed of the five layers in the architecture. Fig. 1 shows an example of the control stack for the quantum dot architecture we propose here, but the particular interfaces between layers will vary according to the physical hardware, quantum error correction scheme, *etc.* that one chooses to implement. At the top of the control stack is the Application layer, where a quantum algorithm is implemented and results are provided to the user. The bottom Physical layer hosts the raw physical processes supporting the quantum computer. The layers between (Virtual, Quantum Error Correction, and Logical) are essential for shaping the faulty quantum processes in the Physical layer into a system of high-accuracy *fault-tolerant* [44] qubits and quantum gates at the Application layer.

### C. Interaction between layers

Two layers meet at an interface, which defines how they exchange instructions or the results of those instructions. Many different commands are being executed and processed simultaneously, so we must also consider how the layers interact dynamically. For the quantum computer to function efficiently, each layer must issue instructions to layers below in a tightly defined sequence. However, a robust system must also be able to handle

errors caused by faulty devices. To satisfy both criteria, a control loop must handle operations at all layers simultaneously while also processing syndrome measurement to correct errors that occur. A prototype for this control loop is shown in Fig. 2.

The primary control cycle defines the dynamic behavior of the quantum computer in this architecture since all operations must interact with this loop. The principal purpose of the control cycle is to successfully implement quantum error correction. The quantum computer must operate fast enough to correct errors; still, some control operations necessarily incur delays, so this cycle does not simply issue a single command and wait for the result before proceeding — pipelining is essential [10, 45]. A related issue is that operations in different layers occur on drastically different timescales, as discussed later in Section VII. Fig. 2 also describes the control structure needed for the quantum computer. Processors at each layer track the current operation and issue commands to lower layers. Layers 1 to 4 interact in the loop, whereas the Application layer interfaces only with the Logical layer since it is agnostic about the underlying design of the quantum computer, which is explained in Section VI.

## D. The QuDOS hardware platform

The layered framework for quantum computing was developed in tandem with a specific hardware platform, known as QuDOS (**q**uantum **D**ots with **o**ptically-**C**ontrolled **S**pins). The QuDOS platform uses electron spins within quantum dots for qubits. The quantum dots are arranged in a two-dimensional array; Fig. 3 shows a cut-away rendering of the quantum dot array inside an optical microcavity, which facilitates control of the electron spins with laser pulses. We demonstrate that the QuDOS design is a promising candidate for large-scale quantum computing, beginning with an analysis of the hardware in the Physical layer.

## II. LAYER 1: PHYSICAL

The essential requirements for the Physical layer are embodied by the DiVincenzo criteria [3], but we are also interested in performance of the quantum hardware. The timescale of operations and the degrees of errors, both systematic and random, are critical parameters which determine the size and speed of the computer. This section discusses the essential hardware components of a quantum computer, accompanied by the QuDOS platform we introduce as an example. We conclude by analyzing the performance of the QuDOS hardware. We caution that many of the required hardware elements are still under experimental development, but we choose those discussed below as examples to establish timescales which will impact higher layers of the architecture.

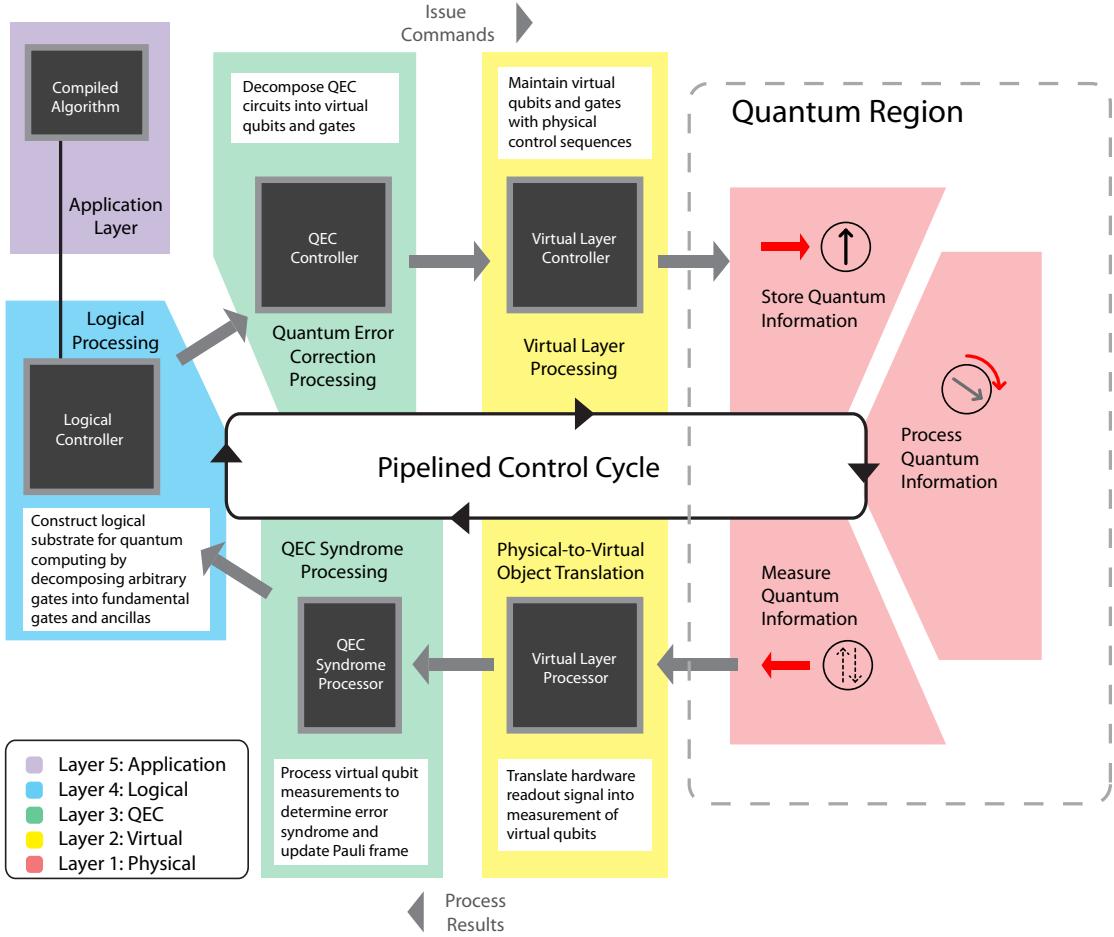


FIG. 2. Color. Primary control cycle of the layered architecture quantum computer. Whereas the control stack in Fig. 1 dictates the interfaces between layers, the control cycle determines the timing and sequencing of operations. The dashed box encircling the Physical layer indicates that all quantum processes happen exclusively here, and the layers above process and organize the operations of the Physical layer. The Application layer is external to the loop since it functions without any dependence on the specific quantum computer design.

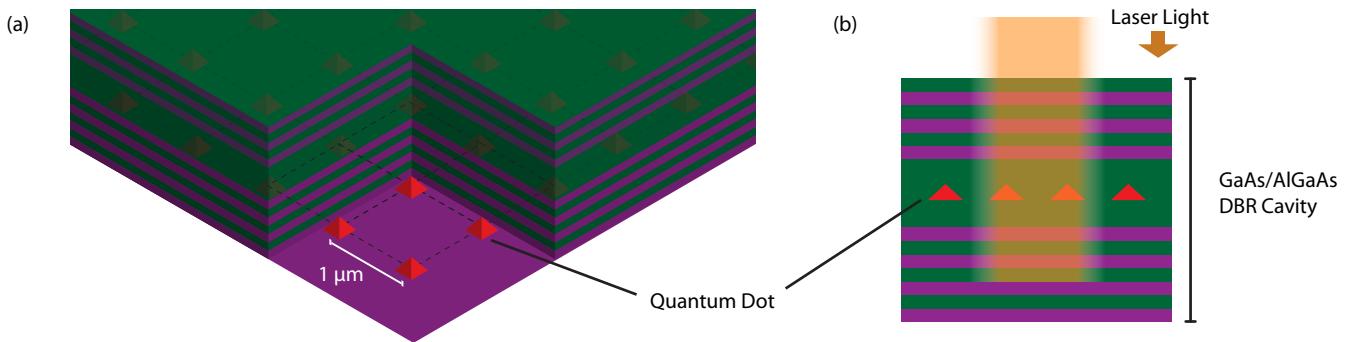


FIG. 3. Color. Quantum dots in a planar optical microcavity form the basis of the QuDOS hardware platform. (a) The quantum dots are arranged  $1\ \mu\text{m}$  apart in a two-dimensional square array. The quantum dots trap single electrons, whose spins will be used for quantum information processing. (b) Side view. The electron spins are manipulated with laser pulses sent into the optical cavity from above, and two neighboring quantum dots can be coupled by a laser optical field which overlaps them. The purple and green layers are AlGaAs and GaAs, grown by molecular beam epitaxy. The alternating layers form a distributed Bragg reflector (DBR) optical cavity which is planar, confining light in the vertical direction and extending across the entire system in horizontal directions.

### A. Physical qubit

A quantum computer must have the ability to store information between processing steps; the object fulfilling this role is conventionally known as the physical qubit. A physical qubit may be more complex than a two-level system, and this issue is addressed by Layer 2 in the architecture, where control operations are used to form a true quantum bit as an information unit (see Section III A). Examples of physical qubits include trapped ions, photon polarization modes, electron spins, and quantum states in superconducting circuits [1]. The remainder of the Physical layer is devoted to controlling and measuring the physical qubit.

The layered architecture design is flexible in the sense that the Physical layer can be tailored to a specific hardware, such as superconducting circuit qubits, with minimal change to higher layers such as error correction. The physical qubit we consider in QuDOS is the spin of an electron bound within an InGaAs self-assembled quantum dot (QD) surrounded by GaAs substrate [46–51]. These QDs can be optically excited to trion states (a bound electron and exciton), which emit light of wavelength  $\sim 900$  nm when they decay. A transverse magnetic field splits the spin levels into two metastable ground states [52], which will later form a two-level system for a virtual qubit in Layer 2. The energy separation of the spin states is important for two reasons related to controlling the electron spin. First, the energy splitting facilitates control with optical pulses as explained in Section II C. Second, there is continuous phase rotation between spin states  $|\uparrow\rangle$  and  $|\downarrow\rangle$  around the  $\hat{Z}$ -axis on the qubit Bloch sphere, which in conjunction with timed optical pulses provides complete unitary control of the electron spin vector.

### B. Host system

For our purposes, the *host system* is the engineered environment of the physical qubit which supports computing. Examples include the trapping fields in ion-trap designs, the waveguides in optical quantum computing, and the diamond crystal surrounding nitrogen-vacancy centers [1]. The host system will define the immediate environment of the physical qubit, which will be important for characterizing noise affecting quantum operations.

We noted above that, in QuDOS, the electron spin is bound within a quantum dot. These quantum dots are embedded in an optical microcavity, which will facilitate quantum gate operations via laser pulses. To accommodate the two-dimensional array of the surface code detailed in Layer 3, this microcavity must be planar in design, so the cavity is constructed from two distributed Bragg reflector (DBR) mirrors stacked vertically with a  $\lambda/2$  cavity layer in between, as shown in Fig. 3. This cavity is grown by molecular beam epitaxy (MBE). The QDs are embedded at the center of this cavity to maxi-

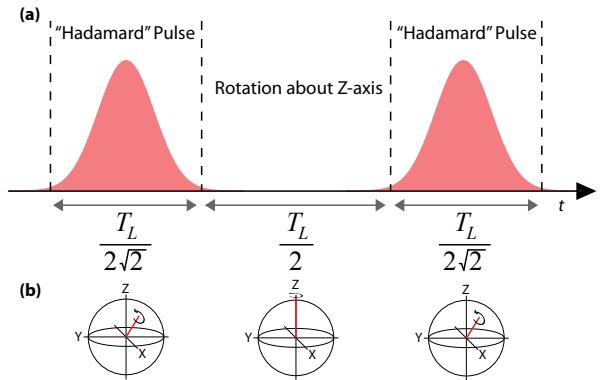


FIG. 4. Color. Hadamard pulses in QuDOS. (a) A short pulse sequence generates an  $\hat{X}$ -axis rotation on the spin Bloch sphere with two Hadamard pulses and  $\hat{Z}$ -axis precession from the magnetic field. The duration of the pulses and the delay between them is proportional to the Larmor period,  $T_L$ . (b) Bloch sphere diagrams showing the axis of rotation at each time during the sequence.

mize interaction with antinodes of the cavity field modes. Using MBE, high-quality ( $Q > 10^5$ ) microcavities can be grown with alternating layers of GaAs/AlAs [53]. The nuclei in the quantum dot and surrounding substrate have nonzero spin, which is an important source of noise (see Section II F).

### C. 1-qubit gate mechanism

The 1-qubit gate manipulates the state of a single physical qubit. This 1-qubit gate mechanism is still a physical process, and only later will these physical control operations be combined into a “virtual gate” (see Section III B). Still, it is important that the Physical layer delivers sufficient control of the physical qubit. Full unitary control of a qubit requires at least two adjustable degrees of freedom, such as rotation around two axes on the Bloch sphere, and three freely adjustable parameters [54].

The 1-qubit operations in QuDOS are developed using a transverse magnetic field and ultrafast laser pulses [51, 55]. The magnetic field provides a constant-angular-frequency rotation around the  $\hat{Z}$  axis, while laser pulses enact a power-dependent rotation around an orthogonal axis, which we label  $\hat{X}$ . The first non-ideal behavior we consider is that the laser pulse has some finite duration, so that  $\hat{X}$  and  $\hat{Z}$  precession happen concurrently, which impairs manipulation of the spin Bloch vector. To remedy this, we introduce “Hadamard pulses” [56]—one tunes the laser pulse power and duration to make the pulse-driven  $\hat{X}$ -axis rotation equal in angular frequency to the  $\hat{Z}$ -axis precession from the magnetic field, so that the axis of rotation becomes  $H = \frac{1}{\sqrt{2}}(\hat{X} + \hat{Z})$ . A “ $\pi$ -pulse” around this axis is a Hadamard gate, and by us-

ing two Hadamard pulses and rotation  $R_{\hat{Z}}(\theta)$  via free precession in the magnetic field, we can construct any  $\hat{X}$ -axis rotation by  $R_{\hat{X}}(\theta) = H \cdot R_{\hat{Z}}(\theta) \cdot H$ , as shown in Fig. 4. By implementing Hadamard pulses, we can obtain high-fidelity operations with pulses of finite duration. A challenging problem for QuDOS is how the system executes millions of control operations in parallel. We envision an optical imaging system consisting of an array of MEMS mirrors to individually steer laser control beams toward or away from quantum dots, along with electro-optic modulators to precisely control laser pulse timing; we discuss this approach in Appendix A, but rigorously engineering such a system is beyond our scope.

#### D. 2-qubit gate mechanism

The 2-qubit operation couples two physical qubits, which can generate entanglement. This mechanism is crucial for quantum computing, yet it is often difficult to implement experimentally. For example, entangling gates like CNOT are used frequently in quantum error correction, so developing fast, high-fidelity 2-qubit gate mechanisms is imperative for large-scale quantum information processing. In many cases, the 2-qubit gate is the process which defines the speed and accuracy of a quantum computer.

The construction of a practical, scalable 2-qubit gate in QuDOS remains the most challenging element of the hardware, and various methods are currently under development. A fast, all-optically controlled 2-qubit gate would certainly be attractive, and early proposals [47] identified the importance of employing the nonlinearities of cavity QED. Ref. [47] suggests the application of two lasers for both single-qubit and 2-qubit control; more recent developments have indicated that both single-qubit gates [55, 57, 58] and 2-qubit gates [59] can be accomplished using only a single optical pulse.

We consider a 2-qubit gate via the dispersive interaction proposed in Ref. [59]. The critical figure of merit for the cavity QED system is the cooperativity factor  $C$ , which is proportional to the cavity quality factor  $Q$  divided by the cavity volume  $V$ . For QuDOS, we envision transverse cavity confinement entirely due to the extended microplanar microcavity arrangement, in which cooperativity factors are enhanced by the angle-dependence of the cavity response, an effect which is enlarged by high index of refraction contrast in the alternating mirrors of the DBR stack [46]. While existing cooperativity factors achieved this way are not estimated to be high enough to produce quantum gates with error rates sufficiently low for fault-tolerant quantum computing, advanced control techniques and multi-spin encodings (such as for “virtual qubits”; see Section III) may enable this technology to function with acceptable error rates. Ref. [59] estimates that this gate will require 10–100 ns to execute, and for the present analysis we as-

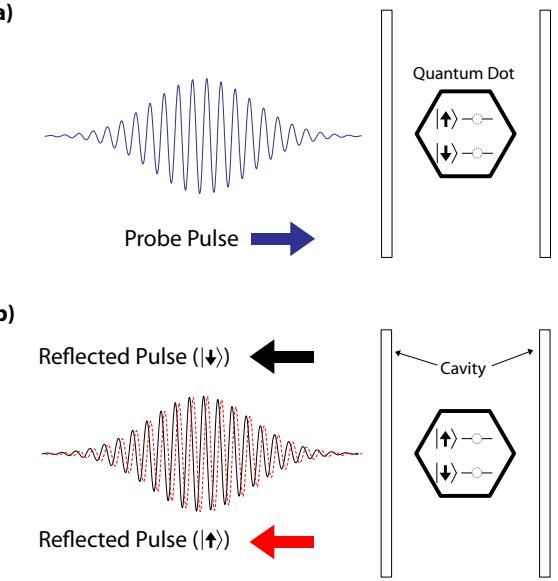


FIG. 5. Color. A dispersive quantum non-demolition (QND) readout scheme for QuDOS. (a) A probe pulse is sent into a microcavity containing a charged quantum dot. (b) The cavity-enhanced dispersive interaction between the pulse and the electron spin creates a state-dependent phase shift in the light which leaves the cavity. Measurement of the phase shift can perform projective measurement on the electron spin.

sume the value 32 ns, which coincides with a virtual gate in Section III B. Further enhancements to speed and/or gate fidelity may be available by introducing exchange interactions using microcavity polaritons [60]; studying this possibility is the subject of future work.

#### E. Measurement readout

Measurement is another essential component of quantum computing. At a bare minimum, one must be able to read the final result of a calculation, but typically measurement is used extensively in fault-tolerant quantum error correction. For this reason, quantum computers may require measurement which is comparable in speed and accuracy to the control operations. Moreover, many situations call for *quantum non-demolition* (QND) measurement, where the physical qubit is projected into an eigenstate of the measurement operator. To illustrate a counter-example, consider a qubit defined by the ground and first optically excited states of a quantum dot. A possible measurement scheme is to detect a photon emission, which would indicate the qubit was in the excited state. However, the final state of the qubit is the ground state for either measurement outcome, which is destructive measurement, and this procedure cannot be repeated. Conversely, QND measurement is highly desirable because it can be repeated, so that classical readout noise can be reduced by time-averaging.

QuDOS will require a QND measurement scheme which is still under experimental development. The proposed mechanism (shown in Fig. 5) is based on Faraday/Kerr rotation. The underlying physical principle is as follows: an off-resonant probe pulse impinges on a quantum dot, and it receives a different phase shift depending on whether the quantum dot electron is in the spin-up or spin-down state (these are separated in energy by the external magnetic field). Sensitive photodetectors combined with homodyne detection measure the phase shift to enact a projective QND measurement on the electron spin. Several results in recent years have demonstrated the promise of this mechanism for measurement: multi-shot experiments by Berezovsky *et al.* [61] and Atatüre *et al.* [62] have measured spin-dependent phase shifts in charged quantum dots, and Fushman *et al.* [63] observed a large phase shift induced by a neutral quantum dot in a photonic crystal cavity. Most recently, Young *et al.* observed a significantly enhanced phase shift from a quantum dot embedded in a micropillar cavity [64].

#### F. Noise sources and errors

Noise and decoherence are the biggest obstacles to scalable quantum computing. In general, the noise sources which corrupt the physical qubit or degrade the fidelity of control operations should be characterized as well as possible. For the present analysis, we consider the noisy environment for an electron spin in QuDOS. The primary noise in this system is dephasing, likely caused by the inhomogeneous distribution of nuclear spins in the quantum dot. The ensemble dephasing is characterized by  $T_2^* \approx 2$  ns, while the intrinsic dephasing is characterized by  $T_2 \approx 3 \mu\text{s}$  [65]. When the noise experienced by a qubit is dominated by dephasing, one can counteract decoherence with control sequences tailored to this noise source [66]. Section III A introduces a decoupling scheme designed specifically for QuDOS.

#### G. Hardware performance summary

We summarize the execution times for the essential Layer 1 operations in QuDOS in Table I. These are the quantum processes which are the building blocks of quantum information operations in Layers 2 and above. For a complete quantum processor, however, one would also have to consider the classical control hardware and the engineering concerns, such as delays, which may occur in a large system. For example, Ref. [34] considers the implications of classical control wires, such as routing concerns, signal timing, and the generation of heat in low-temperature devices. Although engineering of classical control hardware is an important problem, it lies outside the scope of our present analysis, and we reserve it for future work.

### III. LAYER 2: VIRTUAL

The Virtual layer is where quantum effects in the Physical layer are first cast into information primitives — virtual qubits and quantum gates. We use “virtual” as it is defined in the field of computer science, where a virtual object obeys a pre-determined set of behaviors, without specifying the structure of this object. As an example, a virtual qubit may be defined by a decoherence-free subspace [67–69] constructed from three electron spins; when considered as a whole, three spins have many more degrees of freedom than a single qubit. Similar behavior is seen in the quantum gates in QuDOS, which actually consist of a sequence of laser pulses. This transcription process of converting many physical elements into a virtual information unit is the task of Layer 2, and we clarify the functions of this layer below. Fig. 6 gives an overview of the Virtual layer processes in QuDOS.

In a general sense, the Virtual layer makes the Physical layer robust to systematic errors. This effect will be seen in both virtual qubits and gates, where we enforce symmetries in the system (by careful design of control operations) which cause correlated errors to cancel by interference. The simplest example of this behavior is the Hahn spin-echo sequence [70], and in fact decoupling techniques will play a prominent role in how we construct a virtual qubit.

#### A. Virtual qubit

The virtual qubit shapes the underlying physical qubit into a two-level system which approximates an ideal qubit. However, the virtual qubit is modeled as having some finite amount of decoherence, such as the depolarizing channel [54]. Where applicable, dynamical decoupling [71–73] and/or decoherence-free subspaces [67–69] are used to create long-lived virtual qubits, and the residual decoherence characterizes the lifetime of the virtual qubit. In what follows, we consider how to construct a virtual qubit with a charged quantum dot, including the mitigation of several non-ideal effects in this system.

In QuDOS, the virtual qubit is created from the two metastable spin states of an electron confined to a QD. As discussed in Section II F, the raw physical system has dephasing time  $T_2^* \approx 2$  ns [65] caused by an inhomogeneous distribution of nuclear spins in the environment of the electron. This dephasing time is insufficient for quantum error correction in Layer 3, so this system must be augmented with dynamical decoupling techniques [74, 75], which extend the dephasing time of the virtual qubit into the microsecond regime [65].

Constructing the virtual qubit in QuDOS requires Layer 2 to conceal the complexity of controlling the QD spin state. Because the physical qubit Bloch vector continuously rotates around the  $\hat{Z}$ -axis, control pulses must be accurately timed so that they perform the desired operation. Furthermore, control of the QD spin is com-

Operation	Mechanism	Duration	Notes
Spin phase precession ( $\hat{Z}$ -axis)	Magnetic field splitting of spin energy levels	40 ps	Inhomogeneous nuclear environment causes spectral broadening in Larmor frequency, which is the source of $T_2^*$ processes.
Spin state rotation pulse	Stimulated Raman transition with broadband optical pulse	14 ps	Red-detuned from spin ground state-trion transitions.
Entangling operation	Nonlinear phase shift of spin states via coupling to a common cavity mode	32 ns	CW laser signal modulated by an electro-optic modulator (EOM).
QND measurement	Dispersive phase-shift of light reflected from planar cavity	1 ns	CW laser signal modulated by an EOM.

TABLE I. Parameters for Layer 1 quantum operations. Spin phase precession is determined by the spin-state energy splitting due to an external magnetic field. To implement a Hadamard gate, the broadband pulse time is  $1/\sqrt{8}$  of the Larmor period ( $T_{\text{Larmor}}$ ). Times for entangling operation and QND measurement are estimated from simulation.

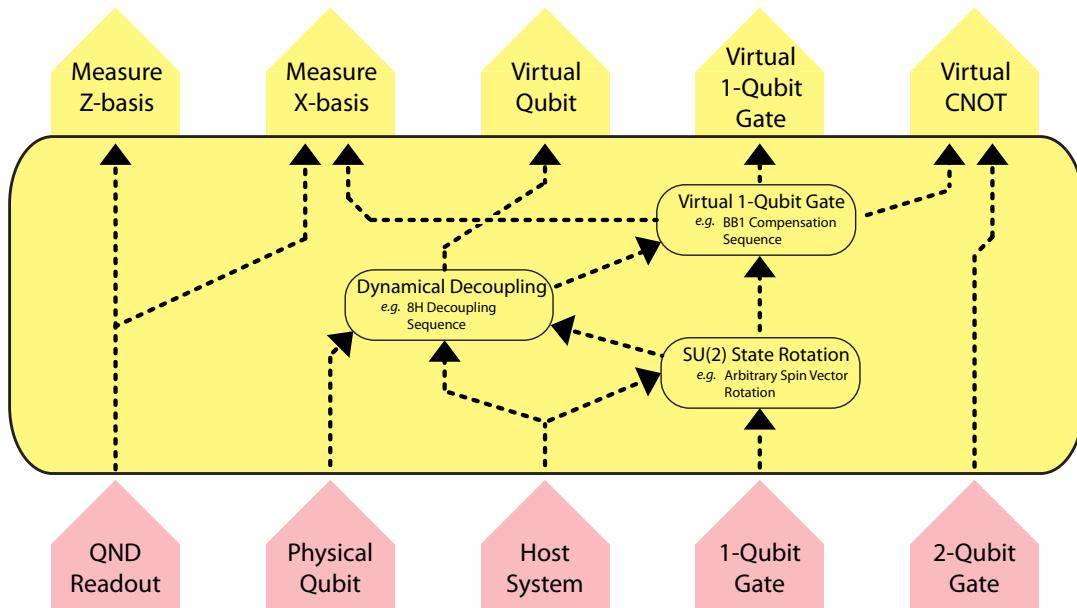


FIG. 6. Color. The mechanics of the Virtual layer. The outputs of Layer 1 are combined in controlled sequences to produce virtual qubits and gates. Arrows indicate how the output of one process is used by another process.

plicated by the inhomogeneous nuclear-spin environment which causes the  $\hat{Z}$ -axis rotation to proceed at a somewhat uncertain angular frequency. This problem is mitigated by a dynamical decoupling (DD) sequence, so that the system is decoupled from environmental noise and brought into a precisely controlled reference frame at a predictable time. Fig. 7a illustrates the “8H” decoupling sequence (so named because it uses eight Hadamard pulses), which is appropriate for use in QuDOS. This control sequence is designed both to decouple a qubit from dephasing noise and to compensate for systematic pulse errors in the presence of a strong but slowly-fluctuating drift term in the qubit Hamiltonian, which is the case for optically-controlled quantum dots in a strong magnetic field. Although longer sequences consisting of more pulses may in theory decouple to higher fidelity, we have chosen a sequence of just eight Hadamard pulses to

minimize execution time. Instead of using a more common sequence like Carr-Purcell (CP) [76, 77] or Uhrig dynamical decoupling (UDD) [78], the sequence in Fig. 7 is custom-designed to eliminate to first-order the errors which occur in both the free evolution and control of the virtual qubit (CP and UDD cannot accomplish the latter). We note, however, that the 8H sequence does have a structure similar to the CP sequence.

Fig. 8 shows the simulated effectiveness of 8H as compared to CP and UDD. We have selected  $\tau = 1$  ns (from Fig. 7a), so one iteration of the sequence requires 8 ns. Because this sequence is specifically designed to account for the errors particular to QuDOS, the performance exceeds that of the more common dynamical decoupling schemes. Nonetheless, 8H may be very effective in other quantum information systems where the physical qubit states are separated in energy and the control pulses have

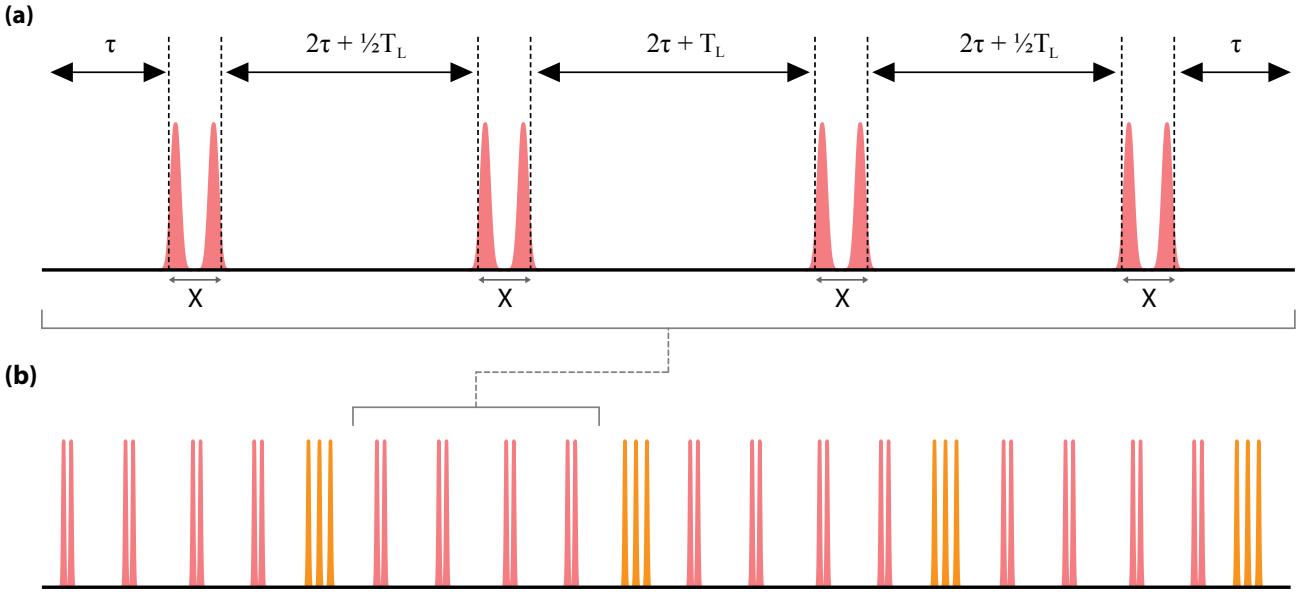


FIG. 7. Color. A special dynamical decoupling sequence for QuDOS, known as 8H since it requires eight Hadamard pulses.  $T_L$  is the Larmor period determined by the external magnetic field (see Table I). (a) Timing specification for the 8H sequence, where  $\tau$  is an arbitrary time. Each of the pulse pairs enacts a  $\pi$ -rotation around the  $X$ -axis of the virtual qubit Bloch sphere, as shown in Fig. 4. For 8H to work efficiently,  $\tau \ll T_2$ . (b) Four 8H sequences in a row interleaved with arbitrary gates formed from three Hadamard pulses (orange). The overall sequence forms a virtual gate by way of a BB1 compensation sequence.

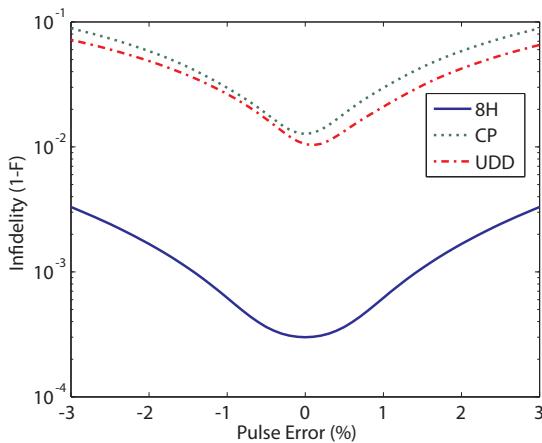


FIG. 8. Color. Simulation of the decoupling effectiveness of the 8H sequence compared to CP and UDD (each using 4  $X$  gates) in the presence of dephasing noise and control errors. Here, “pulse error” is a systematic, relative deviation in the energy of every pulse. In all cases, two Hadamard pulses are combined to produce an approximate  $X$  gate, as in Fig. 4. The vertical axis is infidelity after evolution of the sequence in Fig. 7a with  $\tau = 1$  ns; here infidelity is  $1 - F = 1 - \chi_{II}$ , where  $\chi_{II}$  is the identity-to-identity matrix element in process tomography for the decoupling gate sequence with random noise. Since we aim to execute virtual gates with  $1 - F < 10^{-3}$ , laser pulse errors must be less than 1% in order for the virtual qubit memory error rate to be adequately low.

a duration which is comparable to the free precession (e.g. Larmor) period of the qubit Bloch vector.

## B. Virtual gate

Virtual gates manipulate the state of the virtual qubit by combining physical control operations in Layer 1 in a manner which creates destructive interference of control errors. Quantum operations must be implemented by physical hardware which is ultimately faulty to some extent. Many errors are *systematic*, so that they are correlated in time, even if they are unknown to the quantum computer designer. Virtual gates suppress systematic errors as much as possible in order to satisfy the demands of the error correction system in Layer 3.

Efficient schemes exist for eliminating systematic errors. Compensation sequences can correct correlated errors in the gate operations in Layer 1 [79, 80]. This situation arises often for errors due to imperfections in the control operations, such as laser intensity fluctuations or the coupling strength of a quantum dot electron to an optical field (caused by fabrication imperfections). If these errors are correlated on timescales longer than operations in this architecture, a compensation sequence is effective for generating a virtual gate with lower net error than each of the constituent gates in the sequence. Many compensation sequences are quite general, so that error reduction works without knowledge of the type or

magnitude of error. Dynamically corrected gates are an alternative scheme where one tunes the time-dependent Hamiltonian of the control operations [81]. Beyond such open-loop control techniques, it is also desirable to characterize the accuracy of operations in the Virtual layer, especially multi-qubit gates and entangled states. Systematically evaluating quantum operations is an important component of a research program to develop quantum computers and merits further investigation; however, it is beyond our present scope.

In QuDOS, the ultrafast pulses in Layer 1 would ideally induce a state rotation in the spin basis (two-level system), but inevitably the physical system will suffer from some loss of fidelity by both systematic and random processes. We attempt to cause destructive interference of any systematic errors—both from the environment and control pulses—by embedding a BB1 compensation sequence within a train of 8H dynamical decoupling sequences, as shown in Fig. 7b. This approach is motivated by the properties of the physical qubit. The electron spin has a strong but slowly fluctuating drift term in its Hamiltonian because of the magnetic field and the nuclear spin environment. The 8H sequence brings the qubit “into focus” (analogous to a “spin echo”) only at prescribed instants, which are when the BB1 pulses are applied. This approach is more accurate than a BB1 sequence without refocusing because of the time required to implement rotations on the physical qubit Bloch sphere using Hadamard pulses, for the same reasons that 8H is more effective at decoupling than CP or UDD sequences, as in Figure 8. The BB1 compensation sequence requires four arbitrary gates [79]; hence the virtual gate with error cancellation requires 32 ns.

### C. Measurement of virtual qubits

Measurement is a crucial operation which must also be applied to the virtual qubit in a manner consistent with other control processes. For example, dynamical decoupling prevents measurement by isolating a qubit from environment interactions, so DD may have to be suspended during readout. Since measurement plays a crucial role in error correction, this mechanism should be made as fast and efficient as possible, and a slow measurement process may suffer loss of fidelity if the physical qubit decoheres quickly without DD.

Even if the measurement process is much faster than qubit decoherence, classical noise in the measurement readout signal could be a concern. If the Physical layer provides QND measurement, then the Virtual layer can repeat the measurement of a virtual qubit multiple times and overcome noise in readout circuitry by a majority poll of discrete measurement outcomes. This is a simple yet robust way to suppress measurement errors. For example, if the optical measurement pulse in QuDOS requires 1 ns, then measurement could be repeated about 30 times in the same window of time as a virtual gate.

Another possibility is to couple a virtual qubit to one or more ancilla qubits which facilitate measurement [82]. In such a scheme, the measurement process at the Physical layer could be destructive, but since only the ancilla is destroyed, the back-action on the original qubit is QND measurement, which can be repeated.

Measurement of the virtual qubit in QuDOS requires that the DD sequence be halted, because the 8H sequence interferes with readout. Since the measurement pulse is in the  $\hat{Z}$ -basis, rotations around the  $\hat{Z}$ -axis from the magnetic environment do not affect the outcome. Neglecting DD during measurement is acceptable because the longitudinal ( $T_1$ ) relaxation time is very long compared with the measurement pulse duration [83, 84].

## IV. LAYER 3: QUANTUM ERROR CORRECTION

Fault-tolerant quantum error correction (QEC) is essential for large-scale quantum computing. In Section VI B we analyze an implementation of Shor’s factoring algorithm requiring an error-per-gate of order  $10^{-15}$ , which is simply infeasible on faulty hardware, even using Layer 2 techniques like dynamical decoupling. The action of error correction on a quantum information system is to pump entropy out in the form of an error syndrome; in the process, new resources—*logical* qubits and gates—are created. Whereas Layer 2 causes correlated errors to cancel, Layer 3 isolates and removes arbitrary errors, so long as the error rate is below a threshold [85]. If this condition is met, QEC can in principle produce arbitrarily low-error logical qubits and gates. Such complete error suppression is necessary because quantum algorithms in the Application layer assume logical qubits and gates are error-free.

The field of quantum error correction has become too broad to cover in its entirety [54, 86, 87]. Instead, we analyze the case of stabilizer codes [88], and we specifically consider the surface code [2, 89, 90] for QuDOS. We select the surface code for its high threshold and two-dimensional nearest-neighbor interaction geometry. This section focuses on the aspects of quantum error correction that are relevant for a quantum computer architecture, such as determining the size of a code sufficient for a certain application, as well as how the errors are tracked by Pauli frames in classical hardware. Fig. 9 shows the surface code operations in Layer 3 for QuDOS.

Other error-correction schemes besides the surface code could also be implemented in a layered architecture. As examples, the  $C4/C6$  code [91, 92] and Bacon-Shor codes [32–34, 93] have also received significant attention as viable schemes for fault-tolerant quantum computation. Because the layered architecture is modular, replacing the surface code with another QEC scheme is possible as long as the Virtual layer supports the necessary operations of the new code. In general, the QEC code chosen is likely to impact many aspects of a quantum computing

system, such as device geometry, connectivity, and sensitivity to defective components, so that the structure and behavior of the computer is defined in large part by the selected code. For this reason, much attention should be devoted to optimizing Layer 3 in any quantum computer architecture.

One key message is worth emphasizing. The *threshold error rate* of an error-correcting code is defined as the error rate at which error correction begins to show a net gain in protecting information. A functioning quantum error correction system must operate below threshold, and a *practical* system must operate well below threshold. We show in this section that the resources required for error correction become manageable when the hardware error rate is about an order of magnitude below the threshold of the chosen code.

### A. Estimating the strength of error correction needed

We consider how to estimate the degree of error correction required for a given application because this determines the necessary amount of resources in the computer. Quantum error correction schemes generate protected codespaces within a larger Hilbert space formed from many qubits. The tradeoff for reducing logical errors is that instead of requiring a single qubit, the quantum computer now requires many virtual qubits to produce a logical qubit. The number of virtual qubits required for a single logical qubit is an important resource-usage quantity, and it depends on the performance aspects of the quantum computer:

- error per virtual gate ( $\varepsilon_V$ ), which is an input to Layer 3 from Layer 2,
- threshold error per virtual gate of the error-correcting code ( $\varepsilon_{\text{thresh}}$ ),
- distance ( $d$ ) of the code,
- maximum error per logical gate ( $\varepsilon_L$ ), which is upper-bounded by the performance requirements of the quantum algorithm in Layer 5.

To determine  $\varepsilon_L$ , the simplest approach,  $KQ$  product, assumes the worst case. If the quantum algorithm has a circuit with logical depth  $K$  acting on  $Q$  logical qubits, then the maximum failure probability is given by

$$P_{\text{fail}} = 1 - (1 - \varepsilon_L)^{KQ} \approx KQ\varepsilon_L \quad (1)$$

for small  $\varepsilon_L$ . Therefore, we demand that  $\varepsilon_L \ll 1/KQ$ . Given these quantities, the average error per logical gate in a code operating well below threshold may be closely approximated [54, 85, 86, 94–96] by

$$\varepsilon_L \approx C_1 \left( C_2 \frac{\varepsilon_V}{\varepsilon_{\text{thresh}}} \right)^{\lfloor \frac{d+1}{2} \rfloor}, \quad (2)$$

Parameter	Symbol	Value
Threshold error per virtual gate [96]	$\varepsilon_{\text{thresh}}$	$9 \times 10^{-3}$
Error per virtual gate	$\varepsilon_V$	$1 \times 10^{-3}$
Circuit depth (lattice refresh cycles)	$K$	$1.6 \times 10^{11}$
Logical qubits (“Shor”, Section VI B)	$Q$	72708
Error per lattice refresh cycle	$\varepsilon_L$	$2.6 \times 10^{-20}$
Surface code distance	$d$	31
Virtual qubits per logical qubit	$n_V/n_L$	6240

TABLE II. Parameters determining the size of the surface code in QuDOS for an implementation of Shor’s factoring algorithm.

where  $C_1$  is a constant determined by the specific implementation of the code,  $C_2 \sim 1$ , and, by assumption,  $\varepsilon_V \ll \varepsilon_{\text{thresh}}$ . The data in Ref. [96] suggests  $C_1 \approx 0.13$  and  $C_2 \approx 0.61$  for the surface code, which we now use as an example. Given a known  $\varepsilon_V$  and code-specific quantities  $\{\varepsilon_{\text{thresh}}, C_1, C_2\}$ , one can determine the necessary distance  $d$  such that the probability of failure of an entire quantum algorithm is sufficiently small. For comparison, Aliferis presents similar analysis for concatenated codes such as the Bacon-Shor code [97].

Equation (2) illustrates that the error per virtual gate should be  $\varepsilon_V < 0.2\varepsilon_{\text{thresh}}$ ; otherwise, the code distance, and hence size of the quantum computer, will be impractically large. Table II provides an example of these calculations for the QuDOS quantum computer. Error per virtual gate ( $\varepsilon_V$ ) is also assumed, and the  $K$  and  $Q$  values are for Shor’s algorithm factoring a 1024-bit integer (see Section VI B). We require that  $\varepsilon_L \leq 10^{-2}/KQ$ , so that the logical error probability of the quantum algorithm is less than 1%.

Determining the necessary strength of error correction also indicates how large the quantum computer is in terms of qubits. We can estimate the number of virtual qubits per logical qubit, or  $n_V/n_L$ , by considering the minimum area needed for the two lattice defects, which make up a logical qubit, separated by distance  $d$  in the surface code [2]. For a typical set of parameters, as might be required in a large-scale computing application such as Shor’s factoring algorithm [98], 6240 virtual qubits are needed to construct a logical qubit. This is a nontrivial overhead, because quantum algorithms require a substantial number of logical qubits, as we discuss in greater detail in subsequent sections. For example, quantum simulation algorithms may require  $\sim 1000$ – $10,000$  logical qubits [99, 100] and integer factoring may require 100,000 logical qubits or more, depending on the methods of calculating arithmetic [101]; more detail on why so many logical qubits are necessary is given in Section V B. Combining the size of quantum computations with the requirements of error correction means that large-scale quantum computing architectures will require millions or billions of virtual qubits (and hence physical qubits).

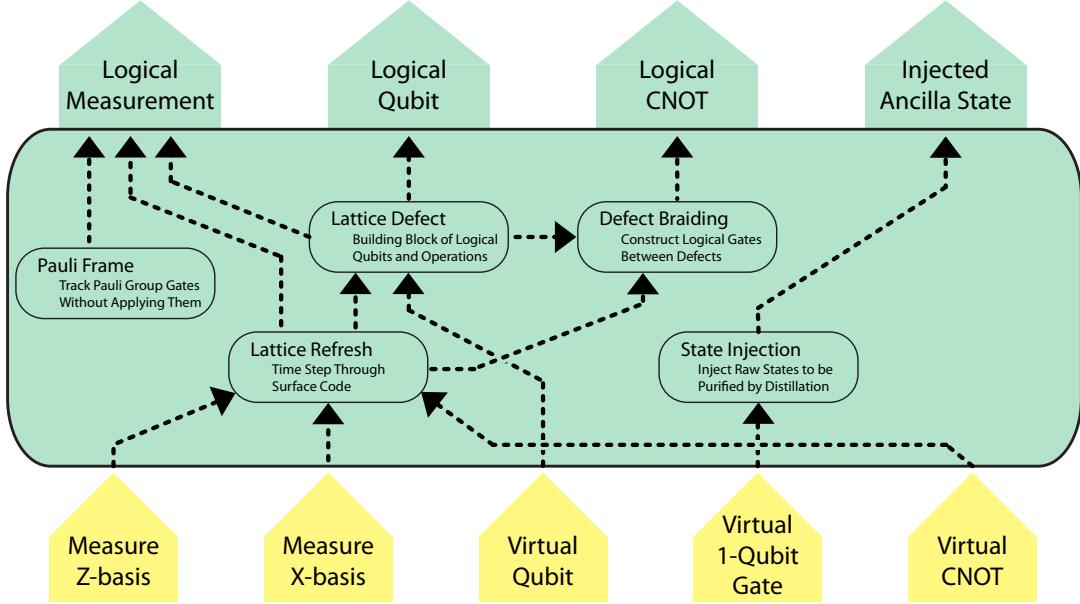


FIG. 9. Color. Process translation in Layer 3 in QuDOS. A surface code is constructed with virtual qubits and gates, ultimately yielding logical qubits and operations. The arrows in yellow along the bottom are outputs of Layer 2, whereas the green arrows at the top are the outputs of Layer 3. Small dashed arrows indicate that the output of one process is used by another process.

### B. Pauli frames

A Pauli frame [91, 102] is a simple and efficient classical computing technique to track the result of applying a series of Pauli gates ( $\mathbb{X}$ ,  $\mathbb{Y}$ , or  $\mathbb{Z}$ ) to single qubits. The Gottesman-Knill Theorem implies that tracking Pauli gates can be done efficiently on a classical computer [103]. Many quantum error correction codes, such as the surface code, project the encoded state into a perturbed codeword with erroneous single-qubit Pauli gates applied (relative to states within the codespace). The syndrome reveals what these Pauli errors are, up to undetectable stabilizers and logical operators, and error correction is achieved by applying those same Pauli gates to the appropriate qubits (since Pauli gates are Hermitian and unitary). However, quantum gates are faulty, and applying additional gates may introduce more errors into our system.

Rather than applying every correction operation, one can keep track of what Pauli correction operation *would be applied*, and continue with the computation. This is possible because the operations needed for error correction are in the Clifford group. When a measurement in a Pauli  $\mathbb{X}$ ,  $\mathbb{Y}$ , or  $\mathbb{Z}$  basis is finally made on a qubit, the result is modified based on the corresponding Pauli gate which should have been applied earlier, as in Fig. 10. This stored Pauli gate is called the Pauli frame [91, 102], since instead of applying a Pauli gate, the quantum computer *changes the reference frame for the qubit*, which can be understood by remapping the axes on the Bloch sphere, rather than moving the Bloch vector.

The Pauli frame is maintained as follows. Denote the Pauli frame at time  $t$  as  $F_t$ :

$$F_t = \bigotimes_j P_t(j), \quad (3)$$

where  $P_t(j) = \{\mathbb{I}, \mathbb{X}, \mathbb{Y}, \mathbb{Z}\}$  is an element from the Pauli group corresponding to qubit  $j$  at time  $t$ . Any Pauli gate in the quantum circuit is multiplied into the Pauli frame and *is not implemented* in hardware, so  $F_{t+1} = \left( \bigotimes_j U_{\{\mathbb{I}, \mathbb{X}, \mathbb{Y}, \mathbb{Z}\}} \right) F_t$  for all Pauli gates  $U_{\{\mathbb{I}, \mathbb{X}, \mathbb{Y}, \mathbb{Z}\}}$  in the circuit at time  $t$ . Other gates  $U_C$  in the Clifford group *are implemented*, but they will transform the Pauli frame by

$$F_{t+1} = U_C F_t U_C^\dagger. \quad (4)$$

The quantum computer operations proceed normally, with the only change being how the final measurement of that qubit is interpreted. The set of Clifford gates is sufficient for Layer 3, though the next section describes another Pauli frame for non-Clifford logical operations.

We emphasize that the Pauli frame is a *classical object* stored in the digital circuitry that handles error correction. Pauli frames are nonetheless very important to the functioning of a surface code quantum computer. Layer 3 uses a Pauli frame with an entry for each virtual qubit in the error-correcting code. As errors occur, the syndrome processing step identifies a most-likely pattern of Pauli errors. Instead of applying the recovery step directly, the Pauli frame is updated in classical memory. The Pauli gates form a closed group under multiplication (and global phase of the quantum state is unimportant),

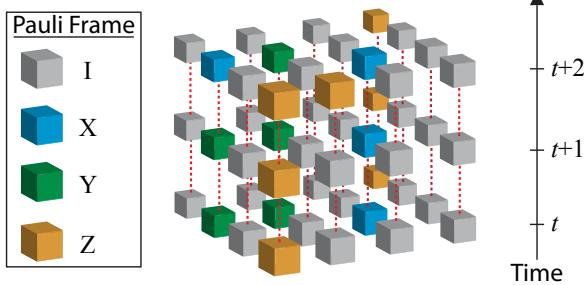


FIG. 10. Color. Example of a Pauli frame evolving over time with entries corresponding to virtual qubits forming a surface code. Each horizontal slice is the Pauli frame at that time. For example, if the qubit in the top-front-corner position is measured in the  $X$  basis, the interpreted result is the negation of the observed outcome, because the Pauli frame  $Z$  anticommutes with this measurement basis.

so the Pauli frame only tracks one of four values ( $X$ ,  $Y$ ,  $Z$ , or  $I$ ) for each virtual qubit in the lattice.

## V. LAYER 4: LOGICAL

The Logical layer takes the fault-tolerant resources from Layer 3 and creates a logical substrate for universal quantum computing. This task requires additional processing of error-corrected gates and qubits to produce any arbitrary gate required in the Application layer, as shown in Fig. 11. Quantum error correction provides only a limited set of gates — to see why, consider that no finite number of syndrome bits can distinguish arbitrarily small rotation gate errors. A common set of gates provided by QEC is the Clifford group; although circuits from this set can be simulated efficiently on a classical computer by the Gottesman-Knill Theorem [54], the Clifford group forms the backbone of quantum circuits. Still, some QEC schemes, such as the surface code, do not provide the full Clifford group without some sort of ancilla. We identify the set of fault-tolerant gates generated by Layer 3 without the use of ancillas as the *fundamental gates*. The Logical layer then constructs arbitrary gates from circuits of fundamental gates and ancillas injected into the error-correcting code. For example, surface code architectures inject and purify the ancillas  $|Y\rangle = \frac{1}{\sqrt{2}}(|0\rangle + i|1\rangle)$  and  $|A\rangle = \frac{1}{\sqrt{2}}(|0\rangle + e^{i\pi/4}|1\rangle)$ ; then the surface code consumes these ancillas in quantum circuits to produce  $S = e^{i(\pi/4)\sigma_Z}$  and  $T = e^{i(\pi/8)\sigma_Z}$  gates, respectively [2, 54]. This section discusses the important functions of the Logical layer: implementing logical Pauli frames; distilling ancilla states like  $|Y\rangle$  and  $|A\rangle$ ; implementing the full Clifford group in the surface code without measurement; and approximating arbitrary quantum gates for the Application layer.

Gate	Implementation	Execution Time (Lattice Steps)
$X, Y, Z$	Pauli frame	Instantaneous
CNOT	Defect braiding	$13\lceil d/4 \rceil$
$H$ (Hadamard)	Shift lattice	$13\lceil d/8 \rceil$
$MX, MZ$ (Measurement)	Measure stabilizers	1

TABLE III. Fundamental gates in QuDOS using surface code QEC (Ref. [2]). The execution time here is one possible implementation, but in many cases the surface code computation can be deformed into other topologically equivalent circuits which yield faster execution at the expense of more spatial resources, or vice versa.

### A. Fundamental gates and logical Pauli frame

Fundamental gates are provided natively by the error-correcting code in Layer 3. For example, Table III shows the fundamental gates used in QuDOS. In practice, Pauli gates are implemented with a *logical* Pauli frame, which is qualitatively the same as the Pauli frame in Layer 3 for virtual qubits (Section IV B). However, in Layer 4 we may also need to apply gates  $U_{NC}$  outside the Clifford group. The gate we actually implement,  $U'_{NC}$ , results from a Pauli frame transform:

$$U'_{NC} = F_t U_{NC} F_t^\dagger. \quad (5)$$

Note the distinction between this expression and Eq.(4), where the Pauli frame is changed by a Clifford gate.

The fundamental gate set in Table III is particular to the surface code, and it is not the full Clifford group because it is missing the phase gate  $S$ ; Section V C illustrates a method to construct  $S$  using an ancilla, without measurement, which is efficient since it uses a small number of fundamental gates and the ancilla can be re-used. The remaining logical gates to produce a universal set in the surface code will require ancilla states which are injected and distilled [2].

### B. Magic state distillation

The conventional method for making a universal set of quantum gates in a *fault-tolerant* manner is to produce a certain ancilla state and use it in a quantum circuit equivalent to the desired logical gate [44, 54, 97]. In some cases these circuits require measurement that consumes the ancilla, so that the number of ancilla states required is proportional to the number of gates in the quantum algorithm. For example, the algorithms discussed in Section VI require around  $10^{12}$  or more ancilla states, which are typically manufactured on an as-needed basis.

To complicate matters, ancillas must be produced by methods that are not fault-tolerant, such as initializing a virtual qubit and applying the appropriate virtual gate. This ancilla state can then be *injected* into a QEC code in

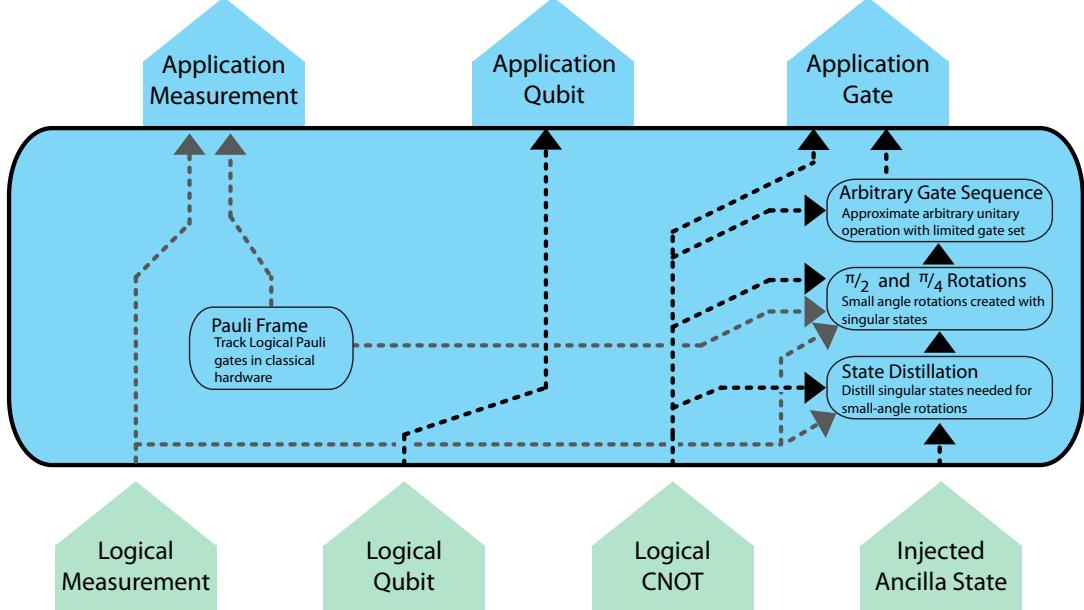


FIG. 11. Color. Organization of processes in the Logical layer. Logical qubits from Layer 3 are unaltered, but faulty singular states are distilled into high-fidelity states  $|Y\rangle = \frac{1}{\sqrt{2}} (|0\rangle + i|1\rangle)$  and  $|A\rangle = \frac{1}{\sqrt{2}} (|0\rangle + e^{i\frac{\pi}{4}}|1\rangle)$ . The distilled states are used to create arbitrary gates with specialized quantum circuits [39, 104–106].

Layer 3 [2], but it carries with it the errors in its production. Fortunately, a few “magic states” can be distilled by using several low-fidelity ancillas and fundamental gates to produce one high-fidelity ancilla. Once the ancilla fidelity is higher than the necessary logical gate fidelity, we may construct arbitrary fault-tolerant logical gates. We examine here the resource costs for this process; each distillation is expensive, and very many ancillas must be distilled. We characterize the performance of the magic state distillation because it will probably dominate the resource costs of any quantum computer that uses it. Accordingly, this is an important area for future optimizations.

We focus first on distilling the ancilla state  $|A\rangle = \frac{1}{\sqrt{2}} (|0\rangle + e^{i\frac{\pi}{4}}|1\rangle)$ , which is used to construct the T or  $\pi/8$  phase gate [2, 90, 107]. In the next section, Fig. 14 provides an illustration of why this process is important by showing the fault-tolerant construction of a Toffoli gate in Layer 5 using resources in Layer 4; specifically, ancilla distillation circuits constitute over 90% of the computing effort for a single Toffoli gate. As a result, the analysis in Appendix B 1 contends that these distillation circuits account for the majority of resources in a surface code quantum computer executing Shor’s algorithm. In particular, for every qubit used by the algorithm, approximately 10 qubits are working in the background to generate the necessary distilled ancillas. The ancilla distillation circuit in Fig. 14 shows one level of  $|A\rangle$  distillation, but a lengthy program like Shor’s will typically require two levels (one concatenated on another). Moreover, since perhaps trillions of distilled  $|A\rangle$  ancillas will

Parameter	Symbol	Value
Circuit depth	-	6 clock cycles
Circuit area	$A_{\text{distill}}$	12 logical qubits
Circuit volume	$V( A^{(1)}\rangle)$	72 qubits $\times$ cycles
Factory rate (level $n$ )	$R_{\text{factory}}( A^{(n)}\rangle)$	$A_{\text{factory}}/V( A^{(n)}\rangle)$ ancillas/cycle

TABLE IV. Resource analysis for a distillation factory. These factories are crucial to quantum computers which require ancillas for universal gates. Magic state distillation uses Clifford gates and measurement, so the circuit can be deformed to reduce depth and increase area, or vice versa, while keeping volume approximately constant.

be needed for the algorithm, we create a “distillation factory” [9, 43], which is a dedicated region of the computer that continually produces these states as fast as possible. Speed is important, because ancilla distillation can be the rate-limiting step in quantum circuits [10].

Each  $|A\rangle$  distillation circuit will require 15 lower-level  $|A\rangle$  states, but they are not all used at the same time. For simplicity we will use a clock cycle for each gate equal to the time to implement a logical CNOT (see Section VII for more on this point), so that with initialization and measurement, the distillation circuit requires 6 cycles. By only using  $|A\rangle$  ancillas when they are needed, the circuit can be compacted to require at most 12 logical qubits at a given instant. We characterize the computing effort by a “circuit volume,” which is the product of logical memory space (*i.e.* area of the com-

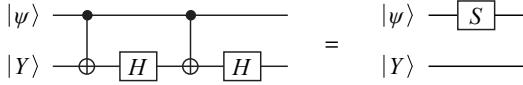


FIG. 12. Circuit decomposition for a logical **S** gate which uses an ancilla  $|Y\rangle$ , but does not consume it. The inverse operation  $S^\dagger$  can be created by running this circuit in reverse.

puter) and time. The circuit volume of  $|A\rangle$  distillation is  $V(|A^{(1)}\rangle) = (12\text{logical qubits}) \times (6\text{ clock cycles}) = 72$ . A two-level distillation will require 16 distillation circuits, or a circuit volume of  $V(|A^{(2)}\rangle) = 1152$ . An efficient distillation factory with area  $A_{\text{factory}}$  will produce on average  $A_{\text{factory}}/V(|A^{(2)}\rangle)$  distilled ancillas per clock cycle. Analysis of this problem in the context of Shor's algorithm is given in Appendix B 1, and Table IV lists a summary of these results.

### C. Logical phase gate without measurement

The **S** gate, or phase gate, is the final component of the Clifford group absent from the fundamental set of fault-tolerant gates in Section V A. Previous implementations of the surface code presented a method for creating this gate by consuming a distilled  $|Y\rangle$  state in a projective measurement-based circuit [2, 90]. However, this approach forces the quantum computer to distill a high-fidelity  $|Y\rangle$  ancilla for each **S** gate, which can be very costly in both fundamental gates and qubits.

We consider an alternative method that uses the  $|Y\rangle$  ancilla without consuming it to make the **S** gate, which was originally presented in Ref. [97]. The circuit uses only four fundamental gates and, unlike the previous technique, is deterministic because measurement is not needed. Since  $|Y\rangle$  ancillas are not consumed, one can distill a handful of such states when a quantum computer is turned on, then preserve them for later use. The circuit in Fig. 12 is equivalent to a simple **S** gate on the control qubit. This is because  $|Y\rangle$  is the  $+i$  eigenstate of the operator  $iY$ , and so the controlled- $iY$  gate will impart a phase  $+i$  only if the control qubit is in the state  $|1\rangle$ , which is identical to the **S** gate. Note also that  $S^\dagger$  can be created by running the circuit in Fig. 12 backwards. This technique allows one to implement the entire Clifford group without measurement in the surface code. Moreover, since **S** gates are used frequently in quantum algorithms, this improved gate construction substantially reduces the complexity of a quantum computer since fewer of the resource-intensive state distillations are necessary.

### D. Approximating arbitrary logical gates

The primary function of the Logical layer is to decompose arbitrary unitary gates from the quantum algorithm

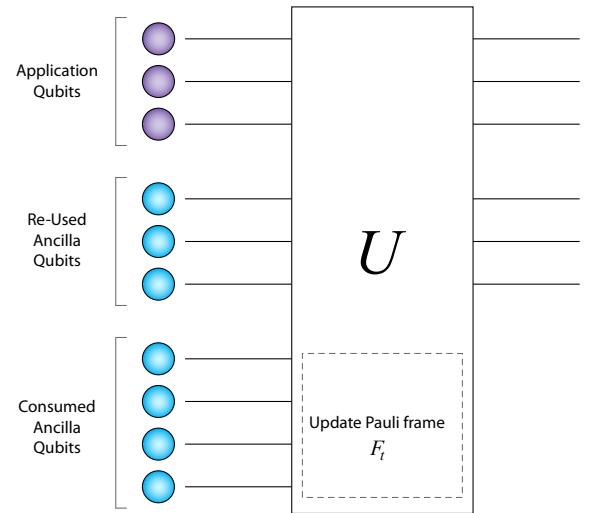


FIG. 13. Color. Constructing an arbitrary gate in the Logical layer. Application qubits are visible to the quantum algorithm, while logical ancilla qubits facilitate universal quantum computation. Some ancillas are re-used, such as  $|Y\rangle$  for **S** gates, while other ancillas are consumed, such as  $|A\rangle$  for **T** gates. Often when an ancilla is consumed in a circuit that uses measurement, the circuit is probabilistic, and the Pauli frame is updated conditional on the measurement result.

into circuits containing fundamental gates available from the QEC layer. The circuits in the Logical layer act on application qubits (used explicitly by the quantum algorithm) and ancilla logical qubits which facilitate universal quantum computation, as shown in Fig. 13. Since arbitrary quantum gates are not available directly, they must be approximated in some fashion, where the total resources required is a function of the approximation accuracy. We cover briefly some of the methods which can be employed to produce such arbitrary gates; a more comprehensive survey of techniques is given in Ref. [106].

When constructing approximations to a unitary operation in the Logical layer, one seeks to implement a quantum circuit that approximates the desired unitary with a minimal overhead in terms of gates and ancillas produced by Layer 3. We denote approximation accuracy as

$$\varepsilon_{\text{approx}} = \sqrt{\frac{d - |\text{tr}(U^\dagger U_{\text{approx}})|}{d}}, \quad (6)$$

where  $U_{\text{approx}}$  is the fault-tolerant quantum circuit that approximates the desired unitary  $U$ , and  $d$  is the dimensionality of these operators [105]. Several techniques exist for approximating arbitrary single-qubit gates, which can be generalized to arbitrary multi-qubit gates:

- Gate approximation sequences, such as those produced by the Solovay-Kitaev algorithm [54, 104] or Fowler's algorithm [105], generate a sequence of gates from the fault-tolerant set (e.g. the set

$\{X, Y, Z, H, S, T\}$ ) that approximates the desired unitary  $U$ . The depth of these sequences scales as  $O(\log^c(\varepsilon_{\text{approx}}))$  with  $c \approx 4$  for Solovay-Kitaev sequences and  $O(\log(\varepsilon_{\text{approx}}))$  for Fowler sequences.

- Phase kickback uses a special ancilla register and a quantum adder to produce fault-tolerant phase rotations [39, 108, 109]. The depth of phase kickback circuits is  $O(\log(\varepsilon_{\text{approx}}))$  or  $O(\log \log(\varepsilon_{\text{approx}}))$  depending on the quantum adder [110–112]. The ancilla register, which is not consumed and can be reused, has size  $m$  qubits to approximate a phase rotation to precision  $\frac{\pi}{2^m}$  radians, which is also  $O(\log(\varepsilon_{\text{approx}}))$ .
- “Teleportation gates” [29] can yield very fast quantum circuits, but typically a special purpose ancilla required for each such gate must be computed in advance, which demands a larger and more complex quantum computer; teleportation gates that increase performance in large-scale quantum computing are used extensively in the architecture of Ref. [10] and in the simulation algorithms analyzed in Ref. [106].

Choosing among these methods depends on the capabilities of the quantum architecture, such as available logical qubits for parallel computation, and on the desired performance characteristics of the computer.

## VI. LAYER 5: APPLICATION

The Application layer is where quantum algorithms are executed. The efforts of Layers 1 through 4 have produced a computing substrate that supplies any arbitrary gate needed. The Application layer is therefore not concerned with the implementation details of the quantum computer—it is an ideal quantum programming environment. We do not introduce any new algorithmic methods here, but rather we are interested in how to accurately estimate the quantum computing resources required for a target application. This analysis can indicate the feasibility of a proposed quantum computer design, which is a worthwhile consideration when evaluating the long-term prospects of a quantum computing research program.

A quantum engineer could start here in Layer 5 with a specific application in mind and work down the layers to determine the system design necessary to achieve desired functionality. We take this approach for QuDOS by examining two interesting quantum algorithms: Shor’s factoring algorithm and simulation of quantum chemistry. A rigorous system design is beyond the scope of the present work, but we consider the computing resources required for each application in sufficient detail that one may gauge the engineering effort necessary to design a quantum computer based on QuDOS technology.

### A. Elements of the Application Layer

The Application layer is composed of *application* qubits and gates that act on the qubits. Application qubits are logical qubits used explicitly by a quantum algorithm (see Fig. 13); as discussed in Section V, many logical qubits are also used to distill ancilla states necessary to produce a universal set of gates, but these distillation logical qubits are not visible to the algorithm in Layer 5. When an analysis of a quantum algorithm quotes a number of qubits without reference to fault-tolerant error correction, often this means the number of application qubits [99, 113–115]. Similarly, Application-layer gates are equivalent in most respects to logical gates; the distinction is made according to what resources are visible to the algorithm or deliberately hidden in the machinery of the Logical layer, which affords some discretion to the computer designer.

A quantum algorithm could request any arbitrary gate in Layer 5, but not all quantum gates are equal in terms of resource costs. We saw in Section VB that distilling  $|A\rangle$  ancillas, which are needed for T gates, is a very expensive process. For example, Fig. 14 shows how Layers 4 and 5 coordinate to produce an Application-layer Toffoli gate, illustrating the extent to which ancilla distillation consumes resources in the computer. When ancilla preparation is included, T gates can account for over 90% of the circuit complexity in a fault-tolerant quantum algorithm (*cf.* Ref. [10] as well). For this reason, we count resources for applications in terms of Toffoli gates. This is a natural choice, because the level of ancilla distillation, number of virtual qubits, *etc.* depend on the choice of hardware, error correction, and many other design-specific parameters; by comparison, number of Toffoli gates is machine-independent since this quantity depends only on the algorithm (much like the number of application qubits mentioned above). To determine error correction or hardware resources for a given algorithm, one can take the Layer 5 resource estimates and work down through Layers 4 to 1, which is an example of modularity in this architecture framework. Using the analysis in the preceding sections, an Application-layer Toffoli gate in QuDOS has an execution time of 930  $\mu\text{s}$  (31 logical gate cycles including the S gate circuits, discussed in Section VII).

### B. Shor’s algorithm

Perhaps the most well-known application of quantum computers is Shor’s algorithm, which decomposes an integer into its prime factors [98]. Solving the factoring problem efficiently would compromise the RSA cryptosystem [116]. Because of the prominence of Shor’s algorithm in the field of large-scale, fault-tolerant quantum computing, we analyze the resources required to factor a number of size typical for RSA.

A common key length for RSA public-key cryptogra-

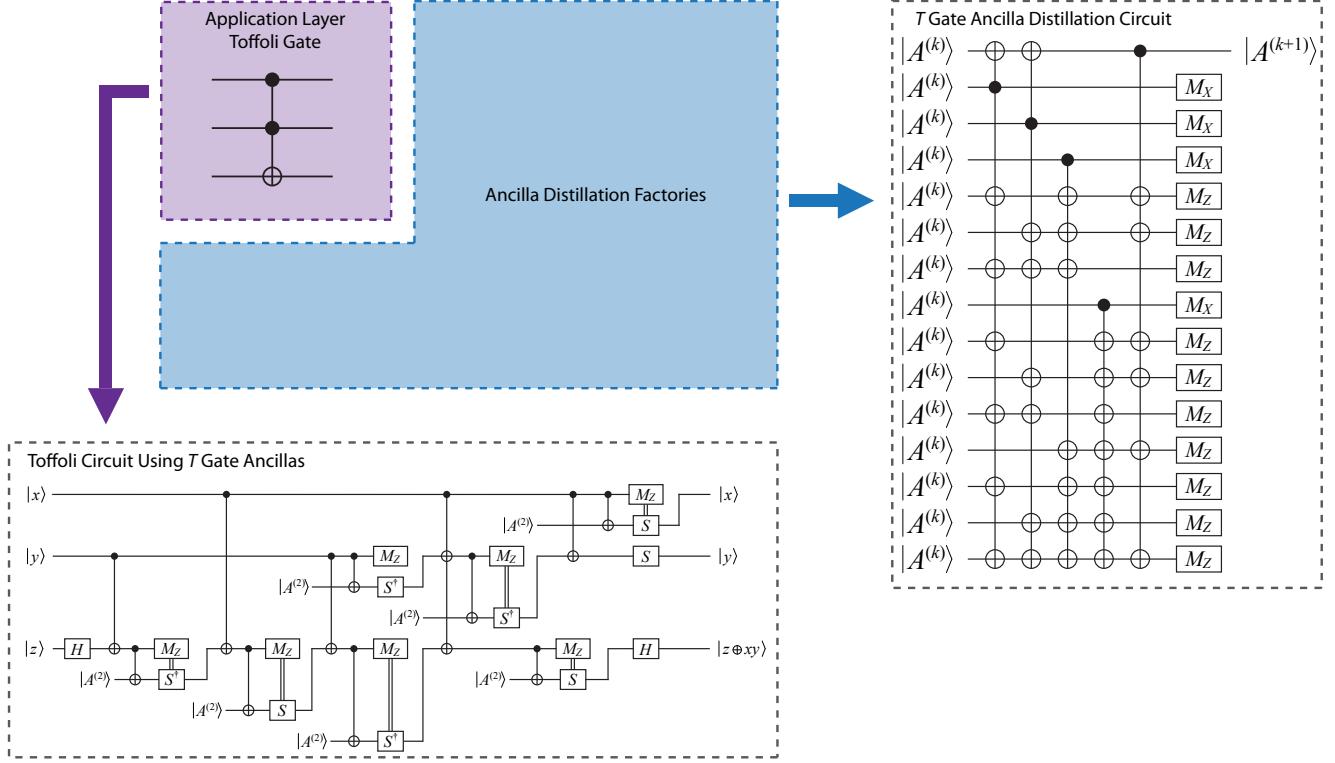


FIG. 14. Color. A Toffoli gate ( $|x, y, z\rangle \rightarrow |x, y, z \oplus xy\rangle$ ) at the Application layer is constructed with assistance from the Logical layer, using the decomposition in Ref. [54]. There are only three application qubits, but substantially more logical qubits are needed for distillation circuits in Layer 4. The  $|A^{(2)}\rangle$  ancillas are the result of two levels of distillation ( $|A^{(0)}\rangle$  is an injected state) on the ancilla required for T gates. Note that each time an ancilla is used with measurement, the Pauli frame may need to be updated. The ancilla-based circuit for S gates (see Fig. 12) is not shown here, for clarity.

phy is 1024 bits. Factoring a number this large is not trivial, even on a quantum computer, as the following analysis shows. Fig. 15 shows the expected run time on QuDOS for one iteration of Shor’s algorithm versus key length in bits for two different quantum computers: one where system size increases with the problem size, and one where the system size is limited to  $10^5$  logical qubits (including application qubits). For the fixed-size quantum computer, the runtime begins to grow faster than the minimal circuit depth when factoring numbers 2048 bits and higher. Fixing the machine size highlights the importance of the ancilla distillation factories. For this instance of Shor’s algorithm, about 90% of the machine should be devoted to distillation; if insufficient resources are devoted to distillation, performance of the factoring algorithm plummets. For example, the 4096-bit factorization devotes  $\sim 75\%$  of the machine to distillation, but about  $3\times$  as many factories would be needed to achieve maximum execution speed in the lower trace in Fig. 15. Many design parameters in an implementation of Shor’s algorithm can be tuned as desired, and we collect the details of our analysis in Appendix B 1. We should also mention here that Shor’s algorithm is probabilistic, so a few iterations may be required [98].

### C. Quantum simulation

Quantum computers were inspired by the problem that simulating quantum systems on a classical computer is fundamentally difficult. Feynman postulated that one quantum system could simulate another much more efficiently than a classical processor, and he proposed a quantum processor to perform this task [117]. Quantum simulation is one of the few known quantum algorithms that solves a useful problem believed to be intractable on classical computers, so we analyze the resource requirements for quantum simulation in the quantum architecture we propose.

We specifically consider fault-tolerant quantum simulation. Other methods of simulation are under investigation [118–120], but they lie outside the scope of this work. The particular example we select is simulating the Schrödinger equation for time-independent Hamiltonians in first-quantized form, where each Hamiltonian represents the electron/nuclear configuration in a molecule [100, 121]. An application of such a simulation is to determine ground- and excited-state energy levels in a molecule. We select first-quantized instead of second-quantized form for better resource scaling at large

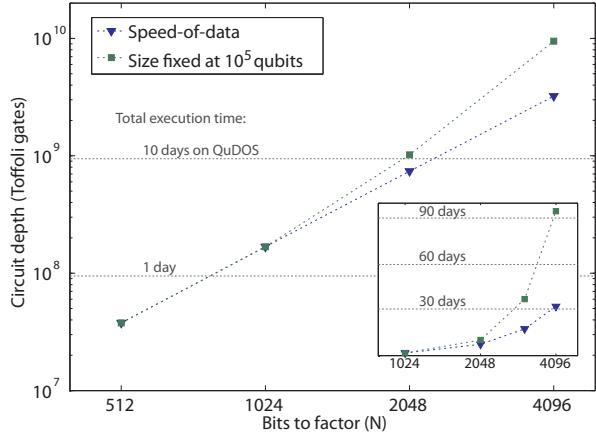


FIG. 15. Color. Execution time for Shor’s algorithm, using the same circuit implementation as Ref. [43]. The vertical axis shows circuit depth, in terms of Toffoli gates, and the plot is labeled with estimated runtime on the QuDOS architecture. The blue trace is a quantum computer whose size in logical qubits scales as necessary to compute at the speed of data (no delays). The green trace is a machine with  $10^5$  logical qubits, which experiences rapidly increasing delays as problem size increases beyond 2048 bits since insufficient resources are available to distill ancillas for T gates, a necessary component of Shor’s algorithm. The inset shows the same data on a linear vertical scale, illustrating when the quantum computer experiences delays for lack of enough qubits.

problem sizes [122].

Fig. 16 shows the time necessary to execute the simulation algorithm for determining an energy eigenstate on the QuDOS computer as a function of the size of the simulation problem, expressed in number of electrons and nuclei. First-quantized form stores the position-basis information for an electron wavefunction in a quantum register, and the complete Hamiltonian is a function of one- and two-body interactions between these registers, so this method does not depend on the particular molecular structure or arrangement; hence, the method is very general. Note that the calculation time scales linearly in problem size, as opposed to the exponential scaling seen in classical methods. The precision of the simulation scales with the number of time steps simulated [99], and this example uses  $2^{10}$  time steps for a maximum precision of about 3 significant figures. Details of this simulation algorithm can be found in Appendix B 2.

#### D. Large-scale quantum computing

The factoring algorithm and quantum simulation represent interesting applications of large-scale quantum computing, and for each the computing resources required of a layered architecture based on QuDOS are listed in Table V. The algorithms are comparable in total resource costs, as reflected by the fact that these two example problems require similar degrees of error correc-

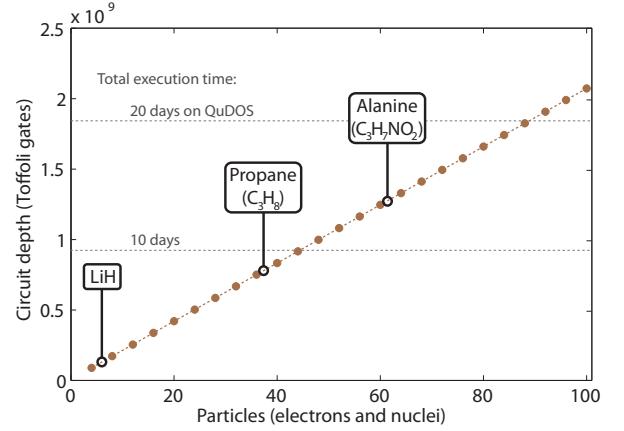


FIG. 16. Color. Execution time for simulation of a molecular Hamiltonian in first-quantized form, as a function of problem size. The horizontal axis is number of particles being simulated, and the plot is labeled with some interesting examples from chemistry. The vertical axis is circuit depth in Toffoli gates, and the plot is labeled with estimated runtime on QuDOS. Each simulation uses 12-bit spatial precision in the wavefunction and  $2^{10}$  time steps for 10-bit precision in readout, or at most  $\sim 3$  significant figures. The linear scaling in algorithm runtime versus problem size is due to two-body potential energy calculations, which constitute the majority of the quantum circuit. The number of potential energy calculations increases quadratically with problem size, but through parallel computation they require linear execution time, as described in Appendix B 2.

tion (hence very similar  $KQ$  product). The simulation algorithm is more compact than Shor’s, requiring in particular fewer logical qubits for distillation, which reflects the fact that this algorithm performs fewer arithmetic operations in parallel. However, Shor’s algorithm has a shorter execution time in this analysis. Both algorithms can be accelerated through parallelism if the quantum computer has more logical qubits to work with [101, 106].

## VII. TIMING CONSIDERATIONS

Precise timing and sequencing of operations are crucial to making an architecture efficient. In the framework we present here, an upper layer in the architecture depends on processes in the layer beneath, so that logical gate time is dictated by QEC operations, and so forth. This system of dependence of operation times is depicted for QuDOS in Fig. 17. The horizontal axis is a logarithmic scale in the time to execute an operation at a particular layer, while the arrows indicate fundamental dependence of one operation on other operations in lower layers.

Examining Fig. 17, we see that the timescales increase as one goes to higher layers. This is because a higher layer must often issue multiple commands to layers below. Using QuDOS as an example, the Virtual layer must construct a virtual 1-qubit gate from a sequence of spin-

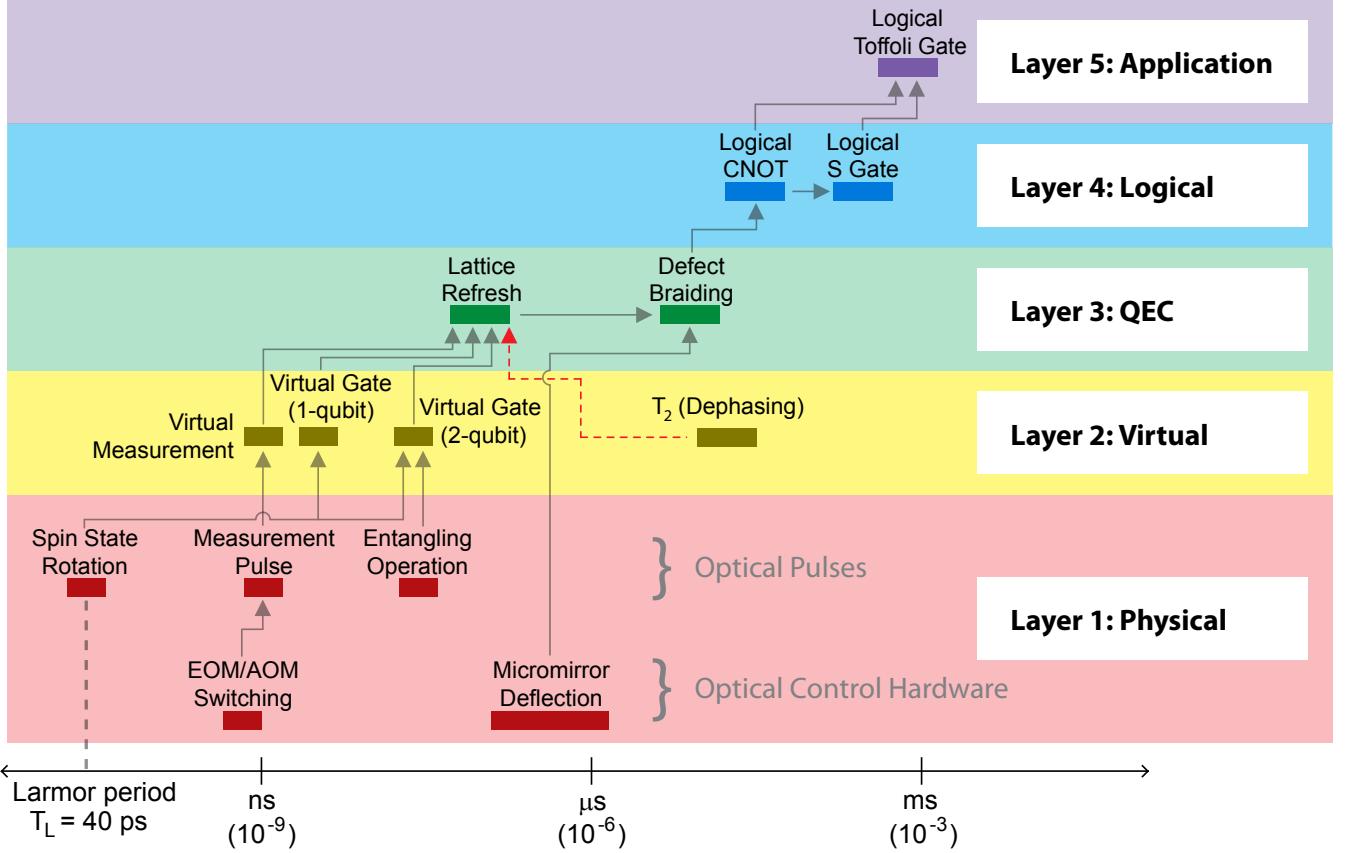


FIG. 17. Color. Relative timescales for critical operations in QuDOS within each layer. Each bar indicates the approximate timescale of an operation, and the width indicates that some operation times may vary with improvements in technology. The arrows indicate dependence of higher operations on lower layers. The red arrow signifies that the surface code lattice refresh must be 2–3 orders of magnitude faster than the dephasing time in order for error correction to function. The Application layer is represented here with a Toffoli gate, which is a common building block of quantum algorithms. Complete algorithm runtimes can vary significantly, depending on both the capabilities of the quantum computer and the specific way each algorithm is compiled, such as to what extent calculations are performed in parallel. The optical control hardware is discussed in Appendix A.

state rotations. This process includes the duration of the laser pulses and the delays between pulses, which all add together for the total duration of the virtual gate. A crucial point shown in Fig. 17 is that the time to implement a logical quantum gate can be orders of magnitude greater than the duration of each individual physical process, such as a laser pulse. This increase in operation time is an important consideration for quantum computer designs which rely on comparatively slower physical processes. At the same time, a quantum computer with a subset of very fast control mechanisms is limited by the slowest essential gate process, as QuDOS can only operate as fast as the 2-qubit entangling gate in Layer 1 permits. For large-scale quantum computing, the speed of logical, error-corrected operations is the crucial figure of merit.

Fig. 17 also highlights the fact that different control operations in the computer occur on substantially different timescales; achieving synchronization of these processes is an important function for a quantum computer

architecture. To facilitate this process, each layer in the architecture has an internal “clock frequency,” which is characteristic of the timescale of operations in that layer. These clock cycle times for each layer in QuDOS are listed in Table VI, along with the operations which define them. Even within the same layer, some processes may take different lengths of time to execute, so setting a clock cycle synchronizes these operations. Accordingly, as one layer builds on operations in a lower layer, the two layers are naturally synchronized.

Synchronization alone is not sufficient for a quantum computer to function. Consider again the control cycle in Fig. 2. Extracting and processing the error syndrome must be executed on timescales of the same order as the duration of a logical gate, or else errors will accumulate faster than they can be detected. This function is performed by classical circuitry, but the required computing effort may not be trivial. Fast quantum operations can be a burden when error correction requires complex (classical) calculations, as is the case for the sur-

Computing Resource		Shor's Algorithm (1024-bit)	Molecular Simulation (alanine)
Layer 5	Application qubits	6144	6650
	Circuit depth (Toffoli)	$1.68 \times 10^8$	$1.27 \times 10^9$
Layer 4	Log. distillation qubits	66564	15860
	Logical clock cycles	$5.21 \times 10^9$	$3.94 \times 10^{10}$
Layer 3	Code distance	31	31
	Error per lattice cycle	$2.58 \times 10^{-20}$	$2.58 \times 10^{-20}$
Layer 2	Virtual qubits	$4.54 \times 10^8$	$1.40 \times 10^8$
	Error per virtual gate	$1.00 \times 10^{-3}$	$1.00 \times 10^{-3}$
Layer 1	Quantum dots (area on chip)	$4.54 \times 10^8$ ( $4.54 \text{ cm}^2$ )	$1.40 \times 10^8$ ( $1.40 \text{ cm}^2$ )
<b>Execution time (est.)</b>		<b>1.81 days</b>	<b>13.7 days</b>

TABLE V. Summary of the computing resources in a layered architecture based on the QuDOS platform, for Shor's algorithm factoring a 1024-bit number (same implementation as Ref. [43]) and the ground state simulation of the molecule alanine ( $\text{C}_3\text{H}_7\text{NO}_2$ ) using first-quantized representation. Further details about the algorithms are provided in Appendix B.

Layer	Clock Cycle	Limiting Operation
(1) Physical	8 ns	Laser repetition frequency
(2) Virtual	32 ns	Virtual 1-qubit gate
(3) QEC	256 ns	Lattice refresh (syndrome circuit)
(4) Logical	30 $\mu$ s	Logical CNOT

TABLE VI. Clock cycle times for Layers 1 to 4 in our analysis of QuDOS. The cycle time in each layer is determined by a fundamental control operation. Many operations possess some flexibility that would permit tradeoffs in execution time and system size, and better methods may be discovered.

face code. Devitt *et al.* [37] and Fowler *et al.* [96, 123] examined this problem, finding that the processing requirements for surface code error correction are not trivial; performing these calculations “live” where the results may be needed within *e.g.* 10  $\mu$ s could be one of the more important problems for engineering a quantum computer. Still, the recent progress in this area suggests that some combination of improved algorithm software and custom hardware can achieve the necessary performance [123].

### VIII. DISCUSSION

We have presented a layered framework for a quantum computer architecture. The layered framework has two major strengths: it is *modular*, and it facilitates *fault tolerance*. The layered nature of the architecture hints at modularity, but the defining characteristic of the layers we have chosen is encapsulation. Each of the layers has a unique and important purpose, and that layer bundles the related operations to fulfill this purpose. Additionally, each layer plays the role of resource manager, since

often many operations in a lower layer are combined in a higher layer. Since technologies in quantum computing will evolve over time, layers may need replacement in the future, and encapsulation makes integration of new processes a more straightforward task.

Fault tolerance is at present the biggest challenge for quantum computers, and the organization of layers is deliberately chosen to serve this need. Arguably, Layers 1 and 5 define any quantum computer, but the layers in between are devoted exclusively to creating fault tolerance in an intelligent fashion. Layer 2 uses simple control to mitigate systematic errors, so this layer is positioned close to the Physical layer where techniques like dynamical decoupling and decoherence-free subspaces are most effective. Layer 3 hosts quantum error correction (QEC), which is essential for large-scale circuit-model quantum computing on any hardware, such as executing Shor's algorithm on a 1024-bit number. There is a significant interplay between Layers 2 and 3, because Layer 2 enhances the effectiveness of Layer 3. Finally, Layer 4 fills the gaps in the gate set provided by Layer 3 to form any desired unitary operation to arbitrary accuracy, thereby providing a complete substrate for universal quantum computation in Layer 5.

QuDOS, a specific hardware platform we introduce here, demonstrates the power of the layered architecture concept, but it also highlights a promising set of technologies for quantum computing, which are particularly noteworthy for the fast timescales of quantum operations, the high degree of integration possible with solid state fabrication, and the adoption of several mature technologies from other fields of engineering. The execution times for fundamental quantum operations are discussed in Section II G, but the importance of these fast processes becomes clear in Fig. 17, where the overhead resulting from virtual gates in Layer 2, QEC in Layer 3, and gate constructions in Layer 4 increases the time to implement quantum gates from nanoseconds in the Physical layer to milliseconds in the Application layer, or six orders of magnitude. In this context, a quantum computer needs very fast physical operations.

One of our principal objectives is to better understand the resources required to construct a quantum computer that solves a problem intractable for classical computers. Common figures of merit for evaluating quantum computing technology are gate fidelity, operation time, and qubit coherence time. This investigation goes further to show how connectivity and classical control performance are also crucial. Designing a quantum computer requires viewing the system as a whole, such that tradeoffs and compatibility between component choices must be addressed. A holistic picture is equally important for comparing different quantum computing technologies, such as ion traps or superconducting circuits. This work illustrates how to approach the complete challenge of designing a quantum computer, so that one can adapt these techniques to develop architectures for other quantum computing technologies we have not considered

here. By doing so, differing system proposals can be compared within a common framework, which gives aspiring quantum engineers a common language for determining the best quantum computing technology for a desired application.

## ACKNOWLEDGMENTS

This work was supported by the National Science Foundation CCF-0829694, the Univ. of Tokyo Special Coordination Funds for Promoting Science and Technology, NICT, and the Japan Society for the Promotion of Science (JSPS) through its “Funding Program for World-Leading Innovative R&D on Science and Technology (FIRST Program).” NCJ was supported by the National Science Foundation Graduate Fellowship. AGF acknowledges support from the Australian Research Council, the Australian Government, and the US National Security Agency (NSA) and the Army Research Office (ARO) under contract W911NF-08-1-0527.

## Appendix A: Parallel Control of Laser Pulses in QuDOS

The QuDOS design depends on applying millions (or billions) of laser control pulses in parallel. For completeness, we outline here a method for achieving this level of control, but a detailed analysis of the engineering problem lies outside the scope of this work. Imagine that the 2D array of quantum dots in a cavity is an image plane, like a projector screen. The challenge is to create a precisely controlled optical pattern on this screen. We need two key elements for this scheme to work: the ability to modulate laser pulses to each point on the screen (*i.e.* quantum dot), and the ability to focus laser signals near the diffraction limit. Similar concepts were presented in Ref. [5].

To solve the first problem, we propose to use an array of MEMS mirrors. This technology was developed for high-definition projectors and optical switches for telecommunications [124, 125], but the same devices are being adapted for use in quantum information processing [126, 127]. Since MEMS mirrors are based on the same fabrication techniques as integrated circuits, a controllable mirror array with millions of units has been demonstrated commercially [124], and even larger arrays may be possible.

The 2-qubit gate mechanism in Section IID requires the quantum dots to be in relatively close proximity ( $1 \mu\text{m}$ ), which is close to the wavelength of the laser light (920 nm). Therefore, any optical patterns will have to compete with diffraction. This is a familiar problem in photolithography, so we propose to adopt the method of phase-shift masking [128, 129] from this field. In essence, a “mask” is defined by a transparent plate patterned in such a way that light passing through the mask receives a

phase shift that is a function of position within the image plane. This technique creates interference of light coming from different directions in such a manner that one can produce diffraction-limited patterns on the quantum dot array.

The MEMS mirrors and phase-shift masks work together as follows. Operations on the surface code follow a highly regular pattern of virtual gates; more specifically, the surface code can be constructed by building a cluster state and performing measurement on selected virtual qubits to create defects in the lattice [90]. The cluster state operations are decomposed into a sequence of laser pulses, which are in turn created by appropriately-designed phase-shift masks. Separately, a pattern of measurement pulses for each virtual qubit are modulated by a MEMS array. All of these optical signals are multiplexed together and sent to the quantum dot array, which is depicted in Fig. 18.

The configuration of defects in the surface code changes more slowly than one cycle of the syndrome extraction circuit. Because the defect boundaries must all be separated by the code distance  $d$ , the pattern of defects can be rearranged every  $d/4$  lattice steps. Therefore, the MEMS mirrors, which control where measurements are made, can be rearranged every  $2 \mu\text{s}$  in QuDOS, which is compatible with current technology [130]. Still, one has to account for the time required to reposition a set of mirrors. Two sets of mirrors are used in an alternating sequence: one is repositioning while the other is actively used, as shown in the top of Fig. 18. Electro-optic modulators can quickly multiplex laser pulses between the two mirror arrays.

## Appendix B: Application Layer Details

We provide here a brief summary of the computational complexity for Shor’s algorithm and quantum simulation in first-quantized form. This analysis produces the resource estimates in Section VI.

### 1. Shor’s algorithm

We adopt the same implementation of Shor’s algorithm given in Van Meter *et al.* [43]. In order to determine the performance of Shor’s algorithm at Layer 5, we must look at how efficiently Layer 4 prepares Toffoli gates. Let us suppose that the quantum computer has capacity for  $10^5$  logical qubits; in general, one can interchange logical capacity and algorithm execution time. To factor an  $N$ -bit number, approximately  $6N$  application qubits are used by the algorithm itself, with the remainder of the logical qubits used to produce the crucial  $|A\rangle$  ancillas. Implementations with fewer application qubits are possible [113, 115], but the performance of such circuits is dramatically slower, especially if one is restricted to a limited set of gates. As shown in Section VB, one

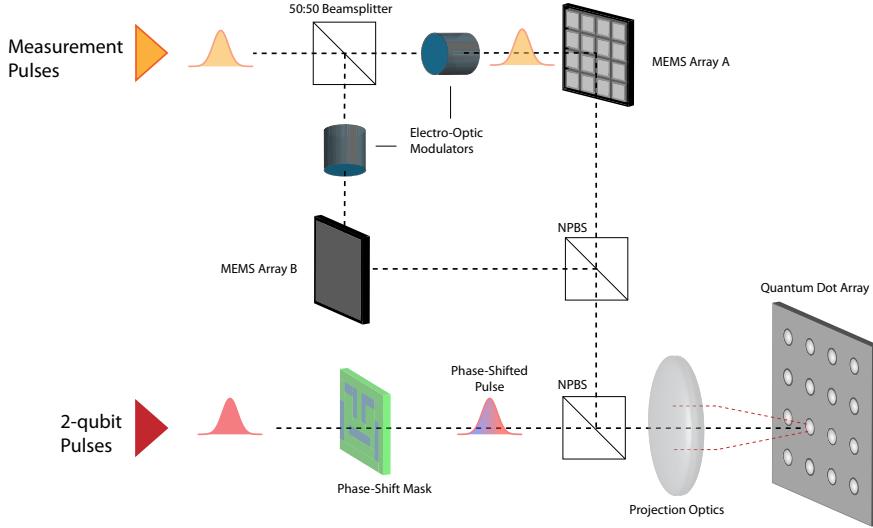


FIG. 18. Color. Optical setup in QuDOS which controls many quantum dots in parallel. Phase-shift masks are used to produce diffraction-limited optical patterns which entangle the physical qubits (bottom). Computation is achieved by measuring the resulting cluster-state; the measurement pattern is controlled by MEMS mirrors (top). Two MEMS mirror arrays are multiplexed with electro-optic modulators, so that one can reposition mirrors while the other is “active.”

Bits to Factor	Ancilla Factory Cross-Section (Logical Qubits)	Distillation Rate ( $ A\rangle$ per cycle)	Consumption Rate (Max) ( $ A\rangle$ per cycle)
512	96928	84.1	32.1
1024	93856	81.5	57.8
2048	87712	76.1	105.1
4096	75424	65.5	192.7
8192	50848	44.1	355.7
16384	1696	1.5	660.6

TABLE VII. Generation rates and maximal consumption rates for a  $10^5$ -qubit quantum computer running Shor’s factoring algorithm. When the speed-of-data consumption rate is higher than the distillation rate, Shor’s algorithm experiences delays.

Operator	Maximum Memory Size (Logical Qubits)	Circuit Depth (Layer 4 Clock Cycles)
Kinetic Energy	$334 \times B$	$1.55 \times 10^5$
Potential Energy	$369 \times B$	$6.26 \times 10^5 \times B$
QFT	$272 \times B$	$2.57 \times 10^4$

TABLE VIII. Resource requirements for the operators in first-quantized molecular simulation with  $B$  particles and 12-bit spatial precision including ancilla distillation.

round of  $|A\rangle$  distillation requires a volume of computing resources with cross-section 12 logical qubits and time 6 CNOT cycles. We arrange the excess ( $10^5 - 6N$ ) qubits in “factories” which distill ancillas as fast as possible.

As before, we define a Logical layer clock cycle as 1 CNOT gate. We express the rate at which the factory generates ancillas by mean number of ancillas pro-

duced per clock cycle. Two levels of distillation will require 16 distillation circuits (15 at the first level, 1 at the second level), which uses a circuit volume of  $V_{\text{distill}} = 16 \times (12 \log_2 \text{qubits}) \times (6 \text{ clock cycles})$ . For a given cross-section area  $A_{\text{factory}}$  of the quantum computer devoted to distillation, the maximal rate of ancilla production is given by  $A_{\text{factory}}/V_{\text{distill}}$ . Calculations of these values are given in Table VII.

We need to determine whether the quantum computer can run as fast as the circuit depth in Layer 5, or whether the distillation of  $|A\rangle$  states limits performance. Using a construction like Fig. 14, the depth of the Toffoli gate is 31 clock cycles, where each S gate requires 4 cycles as shown in Fig. 12, and the circuit requires 7 distilled  $|A\rangle$  ancillas. The circuit uses the carry-lookahead adder construction in Ref. [112], which requires  $\sim 10N$  Toffoli gates in total with a circuit depth of ( $\sim 4 \log_2 N$ )  $t_{\text{Toffoli}}$ , or  $\sim 124 \log_2 N$  cycles. Using these figures, the maximal consumption rate of ancillas can be calculated, as shown in Table VII. As the size of the number to be factored increases, a fixed-size quantum computer is at some point unable to generate enough ancillas to run the algorithm at maximum speed; when this happens, execution time is limited by the distillation process. One can make a crude estimate from Table VII that an efficient quantum computer for Shor’s algorithm must devote 90% of its resources to distillation. By similar arguments, a “minimal” size quantum computer that holds just the algorithm qubits and distills one  $|A\rangle$  ancilla at a time will be very slow.

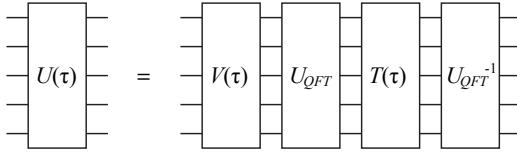


FIG. 19. Circuit representation for one iteration of the Hamiltonian propagator in first-quantized form. The QFT is performed on the wavefunction, transforming between position-basis and momentum-basis.

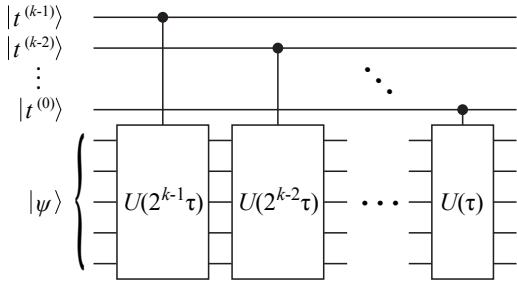


FIG. 20. The time-evolution of the Hamiltonian is produced by iterating the system propagator over many time steps. After evolution, a quantum Fourier transform of the time vector transforms the system into the energy eigenbasis, allowing readout of an energy eigenvalue.

## 2. Quantum simulation

We utilize the method in Ref. [100] to perform simulation in first-quantized form. Each electron wavefunction is represented on a 3-dimensional Cartesian grid with 12 bits of precision in each dimension, which requires a quantum register of 36 qubits per particle. We

elect to use a different set of adders and multipliers than Ref. [100], opting instead for simple ripple-carry adders which suffice for 12-bit precision [111]. First, the potential energy operator is calculated in the position-basis. We transform the wavefunction representation from position-basis to momentum-basis with the quantum Fourier transform (QFT), allowing efficient evaluation of the kinetic energy operator. The inverse QFT transforms our system back to position-basis. The quantum circuit representation of the system propagator  $\mathcal{U}$  is depicted in Fig. 19.

The resource requirements for each of the kinetic ( $\mathcal{T}$ ), potential ( $\mathcal{V}$ ), and QFT operators are summarized in Table VIII. The parameters in Table VIII were derived assuming parallel calculation of commuting operator terms; for example, the Coulomb interaction between particles  $\alpha$  and  $\beta$  can be calculated simultaneously as  $\gamma$  and  $\delta$ , because these terms in the Hamiltonian commute and the circuits are disjoint [106]. Moreover, we have used the preceding analysis in Appendix B 1 to include in these figures the size of ancilla factories, which is  $\sim 260B$  logical qubits in order to simulate a system of  $B$  particles. For this parallel simulation algorithm, ancilla production consumes about 70% of the quantum computer.

The circuit construction in Fig. 19 is just one iteration of the system propagator. Estimating an energy eigenvalue requires simulation of the system at discrete time steps, so the propagator is repeated many times [99], as shown in Fig. 20. After evolving the propagator along these time steps, the system is transformed to the energy eigenbasis via a QFT operation on the time vector  $|t\rangle$  [54]. The precision in the final answer is limited by the number of bits in  $|t\rangle$ , so for this analysis we assume the system is evolved for  $2^{10}$  time steps, which offers at most  $\sim 3$  decimal digits of precision.

- 
- [1] T.D. Ladd, F. Jelezko, R. Laflamme, Y. Nakamura, C. Monroe, and J.L. O'Brien, "Quantum computers," *Nature* **464**, 45–53 (2010).
- [2] Austin G. Fowler, Ashley M. Stephens, and Peter Groszkowski, "High-threshold universal quantum computation on the surface code," *Phys. Rev. A* **80**, 052312 (2009).
- [3] David P. DiVincenzo, "The physical implementation of quantum computation," *Fortschritte der Physik* **48**, 771–783 (2000).
- [4] Andrew M. Steane, "Quantum computer architecture for fast entropy extraction," *Quantum Info. Comput.* **2**, 171–183 (2002).
- [5] Andrew M. Steane, "How to build a 300 bit, 1 Gig operation quantum computer," *Quantum Info. Comput.* **7**, 297–306 (2007).
- [6] Timothy P. Spiller, William J. Munro, Sean D. Barrett, and Pieter Kok, "An introduction to quantum information processing: applications and realizations," *Contemporary Physics* **46**, 407–436 (2005).
- [7] Rodney Van Meter and Mark Oskin, "Architectural implications of quantum computing technologies," *ACM Journal of Emerging Technologies in Computing Systems* **2**, 31–63 (2006).
- [8] J.M. Taylor, H.-A. Engel, W. Dür, A. Yacoby, C. M. Marcus, P. Zoller, and M. D. Lukin, "Fault-tolerant architecture for quantum computation using electrically controlled semiconductor spins," *Nature Physics* **1**, 177–183 (2005).
- [9] A. M. Steane, "Space, time, parallelism and noise requirements for reliable quantum computing," *Fortschritte der Physik* **46**, 443–457 (1998).
- [10] N. Isailovic, M. Whitney, Y. Patel, and J. Kubiatowicz, "Running a quantum circuit at the speed of data," in *35th International Symposium on Computer Architecture, 2008 (ISCA '08)* (2008).
- [11] Tzvetan S. Metodi, Darshan D. Thaker, Andrew W. Cross, Frederic T. Chong, and Isaac L. Chuang, "A quantum logic array microarchitecture: Scalable quantum data movement and computation," in *Pro-*

- ceedings of the 38th International Symposium on Microarchitecture MICRO-38 (2005) pp. 305–318, *Preprint arXiv:quant-ph/0509051v1*.
- [12] D. Kielpinski, C. Monroe, and D. Wineland, “Architecture for a large-scale ion-trap quantum computer,” *Nature* **417**, 709–711 (2002).
- [13] Dean Copsey, Mark Oskin, Tzvetan Metodiev, Frederic T. Chong, Isaac Chuang, and John Kubiatowicz, “The effect of communication costs in solid-state quantum computing architectures,” in *Proceedings of the Fifteenth Annual ACM Symposium on Parallel Algorithms and Architectures (SPAA '03)* (ACM, New York, NY, USA, 2003) pp. 65–74.
- [14] K.M. Svore, A.V. Aho, A.W. Cross, I. Chuang, and I.L. Markov, “A layered software architecture for quantum computing design tools,” *Computer* **39**, 74–83 (2006).
- [15] M. Oskin, F.T. Chong, I.L. Chuang, and J. Kubiatowicz, “Building quantum wires: the long and the short of it,” in *30th International Symposium on Computer Architecture, 2003 (ISCA '03)* (2003) pp. 374–385.
- [16] B.E. Kane, “A silicon-based nuclear spin quantum computer,” *Nature* **393**, 133–137 (1998).
- [17] Matteo Mariantoni, H. Wang, T. Yamamoto, M. Neeley, Radoslaw C. Bialczak, Y. Chen, M. Lenander, Erik Lucero, A.D. O’Connell, D. Sank, M. Weides, J. Wenner, Y. Yin, J. Zhao, A.N. Korotkov, A.N. Cleland, and John M. Martinis, “Implementing the Quantum von Neumann Architecture with Superconducting Circuits,” *Science* **334**, 61–65 (2011).
- [18] J.I. Cirac, P. Zoller, H.J. Kimble, and H. Mabuchi, “Quantum state transfer and entanglement distribution among distant nodes in a quantum network,” *Phys. Rev. Lett.* **78**, 3221–3224 (1997).
- [19] S. J. van Enk, H.J. Kimble, J.I. Cirac, and P. Zoller, “Quantum communication with dark photons,” *Phys. Rev. A* **59**, 2659–2664 (1999).
- [20] A.M. Steane and D.M. Lucas, “Quantum computing with trapped ions, atoms and light,” *Fortschritte der Physik* **48**, 839–858 (2000).
- [21] L.-M. Duan, M.D. Lukin, J.I. Cirac, and P. Zoller, “Long-distance quantum communication with atomic ensembles and linear optics,” *Nature* **414**, 413–418 (2001).
- [22] R. Van Meter, K. Nemoto, and W.J. Munro, “Communication links for distributed quantum computation,” *Computers, IEEE Transactions on* **56**, 1643–1653 (2007).
- [23] L.-M. Duan and C. Monroe, “Colloquium: Quantum networks with trapped ions,” *Rev. Mod. Phys.* **82**, 1209–1224 (2010).
- [24] Jungsang Kim and Changsoon Kim, “Integrated optical approach to trapped ion quantum computation,” *Quantum Info. Comput.* **9**, 181–202 (2009).
- [25] Austin G. Fowler, William F. Thompson, Zhizhong Yan, Ashley M. Stephens, B.L.T. Plourde, and Frank K. Wilhelm, “Long-range coupling and scalable architecture for superconducting flux qubits,” *Phys. Rev. B* **76**, 174507 (2007).
- [26] M. Whitney, N. Isailovic, Y. Patel, and J. Kubiatowicz, “Automated generation of layout and control for quantum circuits,” in *Proceedings of the 4th International Conference on Computing Frontiers* (ACM Press New York, NY, USA, 2007) pp. 83–94.
- [27] M. Whitney, N. Isailovic, Y. Patel, and J. Kubiatowicz, “A fault tolerant, area efficient architecture for Shor’s factoring algorithm,” in *36th International Symposium on Computer Architecture, 2009 (ISCA '09)* (2009).
- [28] N. Isailovic, Y. Patel, M. Whitney, and J. Kubiatowicz, “Interconnection networks for scalable quantum computers,” in *33rd International Symposium on Computer Architecture, 2006 (ISCA '06)* (2006) pp. 366–377.
- [29] Daniel Gottesman and Isaac L. Chuang, “Demonstrating the viability of universal quantum computation using teleportation and single-qubit operations,” *Nature* **402**, 390–393 (1999).
- [30] Jeremy Levy, “Quantum-information processing with ferroelectrically coupled quantum dots,” *Phys. Rev. A* **64**, 052306 (2001).
- [31] Austin G. Fowler, Simon J. Devitt, and Lloyd C. L. Hollenberg, “Implementation of Shor’s Algorithm on a Linear Nearest Neighbour Qubit Array,” *Quantum Info. Comput.* **4**, 237–251 (2004).
- [32] Panos Aliferis and Andrew W. Cross, “Subsystem Fault Tolerance with the Bacon-Shor Code,” *Phys. Rev. Lett.* **98**, 220502 (2007).
- [33] James E. Levy, Anand Ganti, Cynthia A. Phillips, Benjamin R. Hamlet, Andrew J. Landahl, Thomas M. Gurrieri, Robert D. Carr, and Malcolm S. Carroll, “The impact of classical electronics constraints on a solid-state logical qubit memory,” (2009), *Preprint arXiv:0904.0003v1*.
- [34] James E. Levy, Malcolm S. Carroll, Anand Ganti, Cynthia A. Phillips, Andrew J. Landahl, Thomas M. Gurrieri, Robert D. Carr, Harold L. Stalford, and Erik Nielsen, “Implications of electronics constraints for solid-state quantum error correction and quantum circuit failure probability,” *New J. Phys.* **13**, 083021 (2011).
- [35] Yaakov S. Weinstein, C. Stephen Hellberg, and Jeremy Levy, “Quantum-dot cluster-state computing with encoded qubits,” *Phys. Rev. A* **72**, 020304 (2005).
- [36] René Stock and Daniel F. V. James, “Scalable, high-speed measurement-based quantum computer using trapped ions,” *Phys. Rev. Lett.* **102**, 170501 (2009).
- [37] Simon J. Devitt, Austin G. Fowler, Todd Tilma, W. J. Munro, and Kae Nemoto, “Classical processing requirements for a topological quantum computing system,” *International Journal of Quantum Information* **8**, 121–147 (2010).
- [38] Simon J Devitt, Ashley M Stephens, William J Munro, and Kae Nemoto, “Integration of highly probabilistic sources into optical quantum architectures: perpetual quantum computation,” *New J. Phys.* **13**, 095001 (2011).
- [39] Alexei Yu. Kitaev, Alexander H. Shen, and Mikhail N. Vyalyi, *Classical and Quantum Computation*, 1st ed. (American Mathematical Society, 2002).
- [40] Mark Oskin, Frederic T. Chong, and Isaac L. Chuang, “A practical architecture for reliable quantum computers,” *IEEE Computer* **35**, 79–87 (2002).
- [41] Simon J. Devitt, William J. Munro, and Kae Nemoto, “High performance quantum computing,” *Progress in Informatics* **8**, 1–7 (2011), *Preprint arXiv:0810.2444v1*.
- [42] Rodney Van Meter, Kae Nemoto, W. J. Munro, and Kohei M. Itoh, “Distributed arithmetic on a quantum multicomputer,” *SIGARCH Comput. Archit. News* **34**, 354–365 (2006).
- [43] Rodney Van Meter, Thaddeus D. Ladd, Austin G.

- Fowler, and Yoshihisa Yamamoto, "Distributed quantum computation architecture using semiconductor nanophotonics," *International Journal of Quantum Information* **8**, 295–323 (2010), *Preprint arXiv:quant-ph/0906.2686v2*.
- [44] John Preskill, "Fault-tolerant quantum computation," (1997), *Preprint arXiv:quant-ph/9712048*.
- [45] John Paul Shen and Mikko H. Lipasti, *Modern Processor Design: Fundamentals of Superscalar Processors* (McGraw-Hill Higher Education, 2005).
- [46] Gunnar Björk, Stanley Pau, Joseph Jacobson, and Yoshihisa Yamamoto, "Wannier exciton superradiance in a quantum-well microcavity," *Phys. Rev. B* **50**, 17336–17348 (1994).
- [47] A. Imamoğlu, D.D. Awschalom, G. Burkard, D.P. DiVincenzo, D. Loss, M. Sherwin, and A. Small, "Quantum Information Processing Using Quantum Dot Spins and Cavity QED," *Phys. Rev. Lett.* **83**, 4204–4207 (1999).
- [48] N.H. Bonadeo, Gang Chen, D. Gammon, and D.G. Steel, "Single quantum dot nonlinear optical spectroscopy," *Physica Status Solidi B* **221**, 5–18 (2000).
- [49] J.R. Guest, T.H. Stievater, Xiaoqin Li, Jun Cheng, D.G. Steel, D. Gammon, D.S. Katzer, D. Park, C. Ell, A. Thränhardt, G. Khitrova, and H.M. Gibbs, "Measurement of optical absorption by a single quantum dot exciton," *Phys. Rev. B* **65**, 241310 (2002).
- [50] J. Hours, P. Senellart, E. Peter, A. Cavanna, and J. Bloch, "Exciton radiative lifetime controlled by the lateral confinement energy in a single quantum dot," *Phys. Rev. B* **71**, 161306 (2005).
- [51] Y. Yamamoto, T.D. Ladd, D. Press, S. Clark, K. Sanaka, C. Santori, D. Fattal, K.-M. Fu, S. Höfling, S. Reitzenstein, and A. Forchel, "Optically controlled semiconductor spin qubits for quantum information processing," *Physica Scripta* **2009**, 014010 (2009).
- [52] M. Bayer, G. Ortner, O. Stern, A. Kuther, A. A. Gorbunov, A. Forchel, P. Hawrylak, S. Fafard, K. Hinzer, T. L. Reinecke, S. N. Walck, J. P. Reithmaier, F. Klopf, and F. Schäfer, "Fine structure of neutral and charged excitons in self-assembled In(Ga)As/(Al)GaAs quantum dots," *Phys. Rev. B* **65**, 195315 (2002).
- [53] S. Reitzenstein, C. Hofmann, A. Gorbunov, M. Strauß, S.H. Kwon, C. Schneider, A. Löffler, S. Höfling, M. Kamp, and A. Forchel, "AlAs/GaAs micropillar cavities with quality factors exceeding 150.000," *Appl. Phys. Lett.* **90**, 251109 (2007).
- [54] Michael A. Nielsen and Isaac L. Chuang, *Quantum Computation and Quantum Information*, 1st ed. (Cambridge University Press, 2000).
- [55] David Press, Thaddeus D. Ladd, Bingyang Zhang, and Yoshihisa Yamamoto, "Complete quantum control of a single quantum dot spin using ultrafast optical pulses," *Nature* **456**, 218–221 (2008).
- [56] Kristiaan De Greve, Peter L. McMahon, N. Cody Jones, *et al.*, *In preparation*.
- [57] Sophia E. Economou, L.J. Sham, Yanwen Wu, and D.G. Steel, "Proposal for optical U(1) rotations of electron spin trapped in a quantum dot," *Phys. Rev. B* **74**, 205415 (2006).
- [58] Susan M. Clark, Kai-Mei C. Fu, Thaddeus D. Ladd, and Yoshihisa Yamamoto, "Quantum computers based on electron spins controlled by ultrafast off-resonant single optical pulses," *Phys. Rev. Lett.* **99**, 040501 (2007).
- [59] T. D. Ladd and Y. Yamamoto, "Simple quantum logic gate with quantum dot cavity QED systems," *Phys. Rev. B* **84**, 235307 (2011).
- [60] G. F. Quinteiro, J. Fernández-Rossier, and C. Piermarocchi, "Long-range spin-qubit interaction mediated by microcavity polaritons," *Phys. Rev. Lett.* **97**, 097401 (2006).
- [61] J. Berezovsky, M.H. Mikkelsen, O. Gywat, N.G. Stoltz, L.A. Coldren, and D.D. Awschalom, "Nondestructive optical measurements of a single electron spin in a quantum dot," *Science* **314**, 1916–1920 (2006).
- [62] Mete Atatüre, Jan Dreiser, Antonio Badolato, and Atac Imamoğlu, "Observation of Faraday rotation from a single confined spin," *Nature Physics* **3**, 101–106 (2007).
- [63] Ilya Fushman, Dirk Englund, Andrei Faraon, Nick Stoltz, Pierre Petroff, and Jelena Vuckovic, "Controlled Phase Shifts with a Single Quantum Dot," *Science* **320**, 769–772 (2008).
- [64] A.B. Young, R. Oulton, C.Y. Hu, A.C.T. Thijssen, C. Schneider, S. Reitzenstein, M. Kamp, S. Höfling, L. Worschech, A. Forchel, and J.G. Rarity, "Quantum-dot-induced phase shift in a pillar microcavity," *Phys. Rev. A* **84**, 011803 (2011).
- [65] David Press, Kristiaan De Greve, Peter L. McMahon, Thaddeus D. Ladd, Benedikt Friess, Christian Schneider, Martin Kamp, Sven Höfling, Alfred Forchel, and Yoshihisa Yamamoto, "Ultrafast optical spin echo in a single quantum dot," *Nature Photonics* **4**, 367–370 (2010).
- [66] Panos Aliferis and John Preskill, "Fault-tolerant quantum computation against biased noise," *Phys. Rev. A* **78**, 052331 (2008).
- [67] D.A. Lidar, I.L. Chuang, and K.B. Whaley, "Decoherence-free subspaces for quantum computation," *Phys. Rev. Lett.* **81**, 2594–2597 (1998).
- [68] Daniel A. Lidar and K. Birgitta Whaley, "Irreversible quantum dynamics," (Springer, Berlin, 2003) Chap. Decoherence-Free Subspaces and Subsystems, pp. 83–120, *Preprint arXiv:quant-ph/0301032*.
- [69] Matthew Grace, Constantin Brif, Herschel Rabitz, Ian Walmsley, Robert Kosut, and Daniel Lidar, "Encoding a qubit into multilevel subspaces," *New J. Phys.* **8**, 35 (2006).
- [70] E.L. Hahn, "Spin echoes," *Phys. Rev.* **80**, 580–594 (1950).
- [71] Lorenza Viola, Emanuel Knill, and Seth Lloyd, "Dynamical decoupling of open quantum systems," *Phys. Rev. Lett.* **82**, 2417–2421 (1999).
- [72] K. Khodjasteh and D. A. Lidar, "Fault-tolerant quantum dynamical decoupling," *Phys. Rev. Lett.* **95**, 180501 (2005).
- [73] Michael J. Biercuk, Hermann Uys, Aaron P. VanDevender, Nobuyasu Shiga, Wayne M. Itano, and John J. Bollinger, "Optimized dynamical decoupling in a model quantum memory," *Nature* **458**, 996–1000 (2009).
- [74] Lorenza Viola and Emanuel Knill, "Robust dynamical decoupling of quantum systems with bounded controls," *Phys. Rev. Lett.* **90**, 037901 (2003).
- [75] Hui Khoon Ng, Daniel A. Lidar, and John Preskill, "Combining dynamical decoupling with fault-tolerant quantum computation," *Phys. Rev. A* **84**, 012305 (2011).
- [76] H.Y. Carr and E.M. Purcell, "Effects of diffusion on free precession in nuclear magnetic resonance experiments,"

- Phys. Rev. **94**, 630–638 (1954).
- [77] U. Haeberlen and J.S. Waugh, “Coherent averaging effects in magnetic resonance,” Phys. Rev. **175**, 453–467 (1968).
- [78] Götz S. Uhrig, “Keeping a quantum bit alive by optimized  $\pi$ -pulse sequences,” Phys. Rev. Lett. **98**, 100504 (2007).
- [79] Kenneth R. Brown, Aram W. Harrow, and Isaac L. Chuang, “Arbitrarily accurate composite pulse sequences,” Phys. Rev. A **70** (2004).
- [80] Y. Tomita, J.T. Merrill, and K.R. Brown, “Multi-qubit compensation sequences,” New J. Phys. **12**, 015002 (2010).
- [81] Kaveh Khodjasteh and Lorenza Viola, “Dynamically error-corrected gates for universal quantum computation,” Phys. Rev. Lett. **102**, 080501 (2009).
- [82] P. Cappellaro, J. Emerson, N. Boulant, C. Ramanathan, S. Lloyd, and D. G. Cory, “Entanglement assisted metrology,” Phys. Rev. Lett. **94**, 020502 (2005).
- [83] J. M. Elzerman, R. Hanson, L. H. Willems van Beveren, B. Witkamp L. M. K. Vandersypen, and L. P. Kouwenhoven, “Single-shot read-out of an individual electron spin in a quantum dot,” Nature **430**, 431–435 (2004).
- [84] Miro Kroutvar, Yann Ducommun, Dominik Heiss, Max Bichler, Dieter Schuh, Gerhard Abstreiter, and Jonathan J. Finley, “Optically programmable electron spin memory using semiconductor quantum dots,” Nature **432**, 81–84 (2004).
- [85] D. Aharonov and M. Ben-Or, “Fault-tolerant quantum computation with constant error,” in *Proceedings of the Twenty-Ninth Annual ACM Symposium on Theory of Computing (STOC '97)* (ACM, New York, NY, USA, 1997) pp. 176–188.
- [86] John Preskill, “Reliable quantum computers,” P. Roy. Soc. A-Math. Phy. **454**, 385–410 (1998), *Preprint* arXiv:quant-ph/9705031.
- [87] Simon J. Devitt, Kae Nemoto, and William J. Munro, “The idiots guide to quantum error correction,” (2009), *Preprint* arXiv:0905.2794.
- [88] D. Gottesman, *Stabilizer Codes and Quantum Error Correction*, Ph.D. thesis, California Institute of Technology, Pasadena, CA (1997).
- [89] Robert Raussendorf and Jim Harrington, “Fault-tolerant quantum computation with high threshold in two dimensions,” Phys. Rev. Lett. **98**, 190504 (2007).
- [90] R Raussendorf, J Harrington, and K Goyal, “Topological fault-tolerance in cluster state quantum computation,” New J. Phys. **9**, 199 (2007).
- [91] E. Knill, “Quantum computing with realistically noisy devices,” Nature **434**, 39–44 (2005).
- [92] Panos Aliferis, Daniel Gottesman, and John Preskill, “Accuracy threshold for postselected quantum computation,” Quantum Info. Comput. **8**, 181–244 (2008).
- [93] Dave Bacon, “Operator quantum error-correcting subsystems for self-correcting quantum memories,” Phys. Rev. A **73**, 012340 (2006).
- [94] Austin G. Fowler, David S. Wang, and Lloyd C. L. Hollenberg, “Surface code quantum error correction incorporating accurate error propagation,” Quantum Info. Comput. **11**, 8 (2011).
- [95] 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 (2011).
- [96] Austin G. Fowler, Adam C. Whiteside, and Lloyd C. L. Hollenberg, “Towards practical classical processing for the surface code,” (2011), *Preprint* arXiv:1110.5133.
- [97] Panos Aliferis, *Level Reduction and the Quantum Threshold Theorem*, Ph.D. thesis, California Institute of Technology, Pasadena, CA (2007).
- [98] Peter W Shor, “Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer,” SIAM J. Comput. **26**, 1484–1509 (1997).
- [99] Alán Aspuru-Guzik, Anthony D. Dutoi, Peter J. Love, and Martin Head-Gordon, “Simulated quantum computation of molecular energies,” Science **309**, 1704–1707 (2005).
- [100] Ivan Kassal, Stephen P. Jordan, Peter J. Love, Masoud Mohseni, and Alán Aspuru-Guzik, “Polynomial-time quantum algorithm for the simulation of chemical dynamics,” P. Natl. Acad. Sci. USA **105**, 18681–18686 (2008).
- [101] Rodney Van Meter and Kohei M. Itoh, “Fast quantum modular exponentiation,” Phys. Rev. A **71**, 052320 (2005).
- [102] David P. DiVincenzo and Panos Aliferis, “Effective fault-tolerant quantum computation with slow measurements,” Phys. Rev. Lett. **98**, 020501 (2007).
- [103] Simon Anders and Hans J. Briegel, “Fast simulation of stabilizer circuits using a graph-state representation,” Phys. Rev. A **73**, 022334 (2006).
- [104] Christopher M. Dawson and Michael A. Nielsen, “The Solovay-Kitaev Algorithm,” Quantum Info. Comput. **6**, 81 (2006).
- [105] Austin G. Fowler, “Constructing arbitrary Steane code single logical qubit fault-tolerant gates,” Quantum Info. Comput. **11**, 867–873 (2011).
- [106] N. Cody Jones, James D. Whitfield, Peter L. McMahon, Man-Hong Yung, Rodney Van Meter, Alán Aspuru-Guzik, and Yoshihisa Yamamoto, “Simulating chemistry efficiently on fault-tolerant quantum computers,” (2012), *Preprint* arXiv:1204.0567.
- [107] Sergey Bravyi and Alexei Kitaev, “Universal quantum computation with ideal Clifford gates and noisy ancillas,” Phys. Rev. A **71**, 022316 (2005).
- [108] A. Yu. Kitaev, “Quantum measurements and the Abelian Stabilizer Problem,” (1995), *Preprint* arXiv:quant-ph/9511026v1.
- [109] R. Cleve, A. Ekert, C. Macchiavello, and M. Mosca, “Quantum algorithms revisited,” P. Roy. Soc. A-Math. Phy. **454**, 339–354 (1998).
- [110] Vlatko Vedral, Adriano Barenco, and Artur Ekert, “Quantum networks for elementary arithmetic operations,” Phys. Rev. A **54**, 147–153 (1996).
- [111] Steven A. Cuccaro, Thomas G. Draper, Samuel A. Kutin, and David Petrie Moulton, “A new quantum ripple-carry addition circuit,” (2008), *Preprint* arXiv:quant-ph/0410184.
- [112] Thomas G. Draper, Samuel A. Kutin, Eric M. Rains, and Krysta M. Svore, “A logarithmic-depth quantum carry-lookahead adder,” Quantum Info. Comput. **6**, 351–369 (2006).
- [113] Stéphane Beauregard, “Circuit for Shor’s algorithm using  $2n+3$  qubits,” Quantum Info. Comput. **3**, 175–185 (2003).
- [114] Christof Zalka, “Shor’s algorithm with fewer pure qubits,” (2006), *Preprint* arXiv:quant-ph/0601097.
- [115] Yasuhiro Takahashi and Noboru Kunihiro, “A quantum circuit for Shor’s factoring algorithm using  $2n +$

- 2 qubits," *Quantum Info. Comput.* **6**, 184–192 (2006).
- [116] R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Commun. ACM* **21**, 120–126 (1978).
- [117] Richard Feynman, "Simulating physics with computers," *International Journal of Theoretical Physics* **21**, 467–488 (1982).
- [118] Iulia Buluta and Franco Nori, "Quantum simulators," *Science* **326**, 108–111 (2009).
- [119] Julio T. Barreiro, Markus Müller, Philipp Schindler, Daniel Nigg, Thomas Monz, Michael Chwalla, Markus Hennrich, Christian F. Roos, Peter Zoller, and Rainer Blatt, "An open-system quantum simulator with trapped ions," *Nature* **470**, 486–491 (2011).
- [120] J.D. Biamonte, V. Bergholm, J.D. Whitfield, J. Fitzsimons, and A. Aspuru-Guzik, "Adiabatic quantum simulators," *AIP Advances* **1**, 022126 (2011).
- [121] C. Zalka, "Simulating quantum systems on a quantum computer," *P. Roy. Soc. A-Math. Phy.* **454**, 313–322 (1998), *Preprint arXiv:quant-ph/9603026v2*.
- [122] Ivan Kassal, James D. Whitfield, Alejandro Perdomo-Ortiz, Man-Hong Yung, and Alán Aspuru-Guzik, "Simulating Chemistry Using Quantum Computers," *Annu. Rev. Phys. Chem.* **62**, 185–207 (2011).
- [123] Austin G. Fowler, Adam C. Whiteside, and Lloyd C. L. Hollenberg, "Towards practical classical processing for the surface code: timing analysis," (2012), *Preprint arXiv:1202.5602*.
- [124] Texas Instruments, (2011), <http://www.dlp.com>.
- [125] J. Kim *et al.*, "1100x1100 port MEMS-based optical crossconnect with 4-dB maximum loss," *IEEE Photonics Technology Letters* **15**, 1537–1539 (2003).
- [126] Changsoon Kim, C. Knoernschild, Bin Liu, and Jungsang Kim, "Design and characterization of MEMS micromirrors for ion-trap quantum computation," *IEEE Journal of Selected Topics in Quantum Electronics* **13**, 322–329 (2007).
- [127] Caleb Knoernschild, Changsoon Kim, Bin Liu, Felix P. Lu, and Jungsang Kim, "MEMS-based optical beam steering system for quantum information processing in two-dimensional atomic systems," *Optics Letters* **33**, 273–275 (2008).
- [128] Y. Liu and A. Zakhor, "Binary and phase shifting mask design for optical lithography," *IEEE T. Semiconduct. M.* **5**, 138–152 (1992).
- [129] Joanna Aizenberg, John A. Rogers, Kateri E. Paul, and George M. Whitesides, "Imaging profiles of light intensity in the near field: Applications to phase-shift photolithography," *Appl. Opt.* **37**, 2145–2152 (1998).
- [130] G.N. Nielson, R.H. Olsson, P.R. Resnick, and O.B. Spahn, "High-speed MEMS micromirror switching," *Conference on Lasers and Electro-Optics, 2007 (CLEO 2007)*, 1–2 (2007).