

## Chapter 16

# Fault-tolerant quantum computers

R. Van Meter et al., “Distributed quantum computation architecture using semiconductor nanophotonics, arXiv:0906.2686v2 [quant-ph] (2009)

N.C. Jones et al., “A layered architecture for quantum computing using quantum dots, arXiv:1010.5022v1 [quant-ph] (2010).

International Journal of Quantum Information  
 © World Scientific Publishing Company

## DISTRIBUTED QUANTUM COMPUTATION ARCHITECTURE USING SEMICONDUCTOR NANOPHOTONICS

RODNEY VAN METER<sup>1,\*</sup>, THADDEUS D. LADD<sup>2,3</sup>, AUSTIN G. FOWLER<sup>4</sup>,  
 and YOSHIHISA YAMAMOTO<sup>2,3</sup>

<sup>1</sup>*Faculty of Environment and Information Studies, Keio University,  
 5322 Endo, Fujisawa, Kanagawa, 252-8520, Japan*

<sup>2</sup>*Edward L. Ginzton Laboratory, Stanford University, Stanford, CA, 94305-4088, USA*

<sup>3</sup>*National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo-to 101-8430, Japan*

<sup>4</sup>*Center for Quantum Computing Technology, University of Melbourne, Victoria 3010, Australia*

\*Email: rdv@sfc.wide.ad.jp

In a large-scale quantum computer, the cost of communications will dominate the performance and resource requirements, place many severe demands on the technology, and constrain the architecture. Unfortunately, fault-tolerant computers based entirely on photons with probabilistic gates, though equipped with “built-in” communication, have very large resource overheads; likewise, computers with reliable probabilistic gates between photons or quantum memories may lack sufficient communication resources in the presence of realistic optical losses. Here, we consider a compromise architecture, in which semiconductor spin qubits are coupled by bright laser pulses through nanophotonic waveguides and cavities using a combination of frequent probabilistic and sparse deterministic entanglement mechanisms. The large photonic resource requirements incurred by the use of probabilistic gates for quantum communication are mitigated in part by the potential high-speed operation of the semiconductor nanophotonic hardware. The system employs topological cluster-state quantum error correction for achieving fault-tolerance. Our results suggest that such an architecture/technology combination has the potential to scale to a system capable of attacking classically intractable computational problems.

*Keywords:* distributed quantum computation; topological fault tolerance; quantum multicomputer; nanophotonics.

### 1. Introduction

Small quantum computers are not easy to build, but are certainly possible. For these, it is sufficient to consider the five basic DiVincenzo criteria<sup>1,2</sup>: ability to add qubits, high-fidelity initialization and measurement, low decoherence, and a universal set of quantum gates. However, these criteria are insufficient for a large-scale quantum computer. DiVincenzo’s added two communications criteria — the ability to convert between stationary and mobile qubit representations, and to faithfully transport the mobile ones from one location to another and convert back to the stationary representation — are also critical, but so is gate speed (“clock rate”), the parallel execution of gates, the necessity for feasible large-scale classical control systems and feed-forward control, and the overriding issues of manufacturing, including the reproducibility of structures that affect key tuning parameters<sup>3,4</sup>. In light of these considerations, the prospects for large-scale quantum computing

2 Van Meter, Ladd, Fowler and Yamamoto

are less certain.

Advances in understanding what constitutes an attractive technology for a quantum computer are married to advances in quantum error correction. These improvements include the theoretical thresholds below which the application of quantum error correction actually *improves* the error rate of the system<sup>5</sup>, increases in the applicability of known classical techniques<sup>6,7,8,9</sup>, understanding of feasible implementation of error correcting codes<sup>10,11,12,13,14,15,16</sup>, design of error suppression techniques suited to particular technologies or error models<sup>17,18,19,20,21,22,23,24</sup>, advances in purification techniques<sup>25,26,27,28,29,30</sup>, and experimental advances toward implementation<sup>31,32,33</sup>. Among the most important, and radical, new ideas in quantum error correction is *topological quantum error correction* (tQEC), for example *surface codes*<sup>34,35,36,37,38</sup>. These codes are attracting attention due to their high error thresholds and their minimal demands on interconnect geometries, but work has just begun on understanding the impact of tQEC on quantum computer architecture, including determining the hardware resources necessary and the performance to be expected<sup>39,40,41,42</sup>.

The effective fault tolerance threshold in tQEC depends critically on the microarchitecture of a system, principally the set of qubits which can be regarded as direct neighbors of each qubit. As connectivity between qubits increases, both the operations required to execute error correction and the opportunities for “crosstalk” as sensitive qubits are directly exchanged decline, allowing the system to more closely approach theoretical limits.

Here, we argue that even for tQEC schemes that require only nearest-neighbor quantum gates in a two-dimensional lattice geometry, communication resources will continue to be critical. We present an architecture sketch in which efficient quantum communication is used to compensate for architecture inhomogeneities, such as physical qubits which must be separated by large effective distances due to hardware constraints, but also due to qubits missing from the lattice due to manufacturing defects. Assuming a homogeneous architecture may be acceptable for small-scale systems, but in order to create a system that will grow to solve practical, real-world problems, distributed computation and a focus on the necessary communications is required. Further, our design explicitly recognizes that not all communications channels are identical; they vary in the fidelity of created entanglement and physical and temporal resources required. This philosophy borrows heavily from established principles in classical computer architecture<sup>43</sup>. Classically, satisfying the demands of data communication is one of the key activities of system architects<sup>44</sup>. Our design process incorporates this philosophy.

No computing system can be designed without first considering its target *workload* and *performance goals*<sup>43,45</sup>. The level of imperfection we allow for quantum operations depends heavily on the application workload of the computer. Our goal is the detailed design (and ultimately implementation) of a large-scale system: more than ten thousand logical qubits capable of running  $10^{11}$  Toffoli gates within a reasonable time (days or at most a few months). For example, such a system could factor a 2,000-bit number using Shor’s algorithm<sup>46</sup>. This choice of scale affects the amount of error in quantum operations that we can tolerate. Steane analyzes the strength of error resilience in a system in terms of  $KQ$ , the product of the number of logical qubits in an application ( $Q$ ) and the depth

(execution time, measured in Toffoli gate times) of the application ( $K$ )<sup>10</sup>. Our goal is to tune the error management system of our computer to achieve a logical error per Toffoli gate executed of  $p_L \ll 1/KQ$ , with  $KQ \sim 10^{15}$ <sup>47</sup>.

Under most realistic technological assumptions, the resources required to reach adequate  $KQ$  values are huge. Nearly all proposed matter qubits are at least microns in size, when control hardware is included. For chip-based systems, a simple counting argument demonstrates that more qubits are required than will fit in a single die, or even a single wafer. This argument forces the implementation to adopt a distributed architecture, and so we require that a useful technology have the ability to entangle qubits between chips<sup>47,48</sup>.

As an example architecture supporting rich communications, we are designing a device based on semiconductor nanophotonics, using the spin of an unpaired electron in a semiconductor quantum dot as our qubit, with two-qubit interactions mediated via cavity QED. We plan to use tQEC to manage run-time, soft faults, and to design the architecture to be inherently tolerant of fabricated and grown defects in most components.

Our overall architecture is a *quantum multicomputer*, a distributed-memory system with a large number of nodes that communicate through a multi-level interconnect. The distributed nature will allow the system to scale, circumventing a number of issues that would otherwise place severe constraints on the maximum size and speed of the system, hence limiting problems for which the system will be suitable.

Within this idiom, many designs will be possible. The work we present here represents a solid step toward a complete design, giving a framework for moving from the overall multicomputer architecture toward detailed node design. We can now begin to estimate the actual hardware resources required, as well as establish goals (such as the necessary gate fidelity and memory lifetimes) for the development of the underlying technology.

Section 2 presents background on the techniques for handling of errors in a quantum computer that we propose to use. Section 3 qualitatively presents our hardware building blocks: semiconductor quantum dots, nanophotonic cavities and waveguides, and the optical schemes for executing gates. Section 4 presents a qualitative description of the resources employed in the complete system. In particular, it describes how some quantum dots, used for communication, are arranged for deterministic quantum logic mediated by coupled cavity modes, while other quantum dots are indirectly coupled via straight, cavity-coupled waveguides for purification-enhanced entanglement creation. Long columns of these basic building blocks span the surface of a chip, and many chips are coupled together to create the complete multicomputer. Preliminary quantitative resource counts appear in section 5.

## 2. Multi-level Error Management

A computer system is subject to both *soft faults* and *hard faults*; in the quantum computing literature, “fault tolerance” refers to soft faults. A soft fault is an error in the operation of a normally reliable component. Soft faults can be further divided into errors on the quantum state (managed through dynamically-executed quantum error correction or purification), and the loss of qubit carrier (e.g., loss of a photon, ion or the electron in a quantum dot,

depending on the qubit technology). Qubit loss may be addressed by using erasure codes, or, in the case of tQEC, through special techniques for rebuilding the lattice state<sup>49</sup>. In this section, we introduce our approach to managing these multiple levels of errors, which will be further developed in the following sections.

### **2.1. Defect Tolerance and Quantum Communication**

Hard faults are either manufactured or “grown” defects (devices that stop working during the operational lifetime of the system). With adequate hardware connectivity, flexible software-based assignment of roles to qubits will add hard fault tolerance, allowing the system to deal with both manufactured and grown defects.

The percentage of devices that work properly is called the *yield*. In our system, most of the components are expected to have high yields, but the quantum dots themselves will likely have low yields, at least in initial fabrication runs and possibly in ultimate devices. These faults occur in part due to the difficulty of growing optically active quantum dots in prescribed locations, but more due to the difficulty of assuring each dot is appropriately charged and tuned near the optical wavelength of the surrounding nanophotonic hardware, to be further discussed in Sec. 3.3.

The presence of hard faults means that the connectivity of the quantum computer begins in a random configuration, which we can determine by device testing. As a result, the architecture will have an inhomogeneous combination of high-fidelity connections where pairs of neighboring qubits are good and low-fidelity connections between more distant qubits. To compensate for the low-fidelity connections, we choose to use *entanglement purification* to bring long-distance entangled-states up to the fidelity we desire for building our complete tQEC lattice. This choice means that the system will naturally use many of the techniques developed for quantum repeaters<sup>50,21,30</sup>, and portions of the system will require similar computation and communication resources, used in a continuous fashion. Details of these procedures are presented in Sec. 4.

### **2.2. Topological Fault Tolerance**

On top of purified states, we employ *topological error correction* (tQEC),<sup>34,35,36,37</sup> in particular the two-dimensional scheme introduced by Raussendorf and Harrington<sup>38,51,52</sup>. In this scheme, the action of the quantum computer is the sequential generation and detection of a cluster state, and error correction proceeds by checking against expected quantum correlations for that state. Logical qubits are defined by deliberately altering these correlations at a pair of boundaries in an effectively three-dimensional lattice of physical qubits. These boundaries may be the extremities of the lattice or holes<sup>a</sup> of various shapes “cut” into the lattice by choosing not to entangle some qubits. The qubits

<sup>a</sup>These holes are commonly called “defects” in the topological computing literature, as they are similar to defects in a crystal; in this paper, we reserve the term “defect” for a qubit that does not function properly, i.e. a manufacturing defect.

in the interior of the lattice have their state tightly constrained, whereas pairs of boundaries are associated with a degree of freedom that is used as the logical qubit.

The simplicity of the gate sequences used to constrain the qubits in the lattice interior and the independence of these gate sequences on the size of the system are directly responsible for tQEC's high threshold error rate of approximately 0.8% for preparation, gate, storage and measurement errors<sup>35,53</sup>, the highest threshold found to date for a system with only nearest neighbor interactions.

In 2-D, we choose to make holes that are squares of side length  $d$ . Logical operators take the form of rings and chains of single-qubit operators — chains connect pairs of holes, rings encircle one of the holes. If we associate  $X_L$  with chains and  $Z_L$  with rings (or vice versa), it can be seen that these operators will always intersect an odd number of times ensuring anticommutation. Braiding holes around one another can implement logical CNOT, as shown in Figure 1.

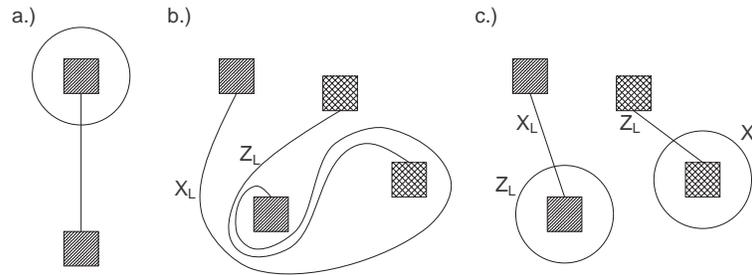


Fig. 1. Logical qubits in topologically error-corrected systems are represented by unentangled “holes” in a high-entangled cluster state on a lattice. The lattice itself is not shown; the squares represent the holes. a.) A single logical qubit is associated with two holes. Logical operators are rings and chains of single qubit operators. b.) Moving holes around one another by changing the error correction circuits on the boundary of holes results in the deformation and ultimately braiding of logical operators. c.) Equivalent form of the braided logical operators after pinching together sections, to form disjoint rings and chains. The mapping of logical operators represents logical CNOT with the left logical qubit as control.

tQEC offers important architectural advantages over other error-suppression schemes, such as concatenated codes. Most importantly, unlike tQEC, many concatenated codes lose much of their effectiveness when long-distance gates are precluded by the underlying technology. In addition, the amount of error correction applied in tQEC can be controlled more finely than with concatenated codes, which have a property that every time an additional level of error correction is used, the number of physical qubits grows by at least an order of magnitude. tQEC's error-protection strength, in contrast, improves incrementally with each additional row and column added to the lattice.

Logical errors are exponentially suppressed by increasing the circumference and separation of holes. This can be inferred directly from Figure 1 — the number of physical qubit errors required to form an unwanted logical operation grows linearly with circumference and separation. The threshold error rate  $p_{th}$  is defined to be the error rate at which increasing the resources devoted to error correction neither increases nor decreases the logical

error — the error rate at which the errors corrected are balanced by the errors introduced by the error correction circuitry. Assuming a hole circumference and separation of  $4d$ , for physical error rates  $p < p_{th}$ , error suppression of order  $O((p/p_{th})^{\alpha d})$  will be observed. The factor  $\alpha$  depends on the details of the error correction circuits. Assuming the error correction circuits do not copy single errors to multiple locations,  $\alpha \sim 2$  as a circumference of  $4d$  implies that a chain of approximately  $2d$  errors can occur before our error correction system will mis-correct the state and give a logical error.

Related tQEC schemes exist in 3-D and 2-D<sup>35,51,52,38</sup>. The 3-D scheme makes use of a 3-D cluster state and the measurement-based approach to computing — all qubits are measured in various bases, and measurement results processed to determine both the bases of future measurements and the final result of the computation. This approach is well-suited to a technology with short-lived qubits (e.g., photons, which are easily lost) or slow measurement. The 2-D scheme requires a 2-D square lattice of qubits that are not easily lost plus fast measurement. Given these two properties, the threshold is slightly higher than the 3-D case and certain operations, such as logical measurement, can be performed more quickly. Barring these minor caveats, the 2-D scheme is a simulation of the 3-D scheme, in which one dimension of the 3-D lattice becomes time.

### 2.3. Logical Gates in Topological Error-Corrected Systems

When making use of topological error correction, only a small number of single logical qubit gates are possible — namely  $X_L$ ,  $Z_L$  and logical initialization and measurement in these bases. Logical initialization and measurement in the  $X_L$  and  $Z_L$  bases can be implemented using initialization and measurement of regions of single qubits encompassing the defects in the  $X$  and  $Z$  bases. The only possible multiple logical qubit gate, logical CNOT, can be implemented by braiding the correct type of defects in a prescribed manner as shown in Figure 1. This set of gates is not universal.

To achieve universality, rotations by  $\pi/2$  and  $\pi/4$  around the  $X_L$  and  $Z_L$  axes can be added to the logical gate set. These gates, however, require the use of specially-prepared  $S$  states where  $|S\rangle = |0\rangle + e^{i\theta}|1\rangle$ ,  $\theta = \pi/2, \pi/4$ . Fault-tolerant creation of the  $S$  states involves use of the concatenated decoding circuits for the 7-qubit Steane code and 15-qubit Reed-Muller code respectively to distill a set of low-fidelity  $S$  states into a single higher-fidelity one. Convergence is rapid — if the input states have average probability of error  $p$ , the output states will have error probabilities of  $7p^3$  and  $35p^3$  respectively<sup>35</sup>.

This implies that for most input error rates, two levels of concatenation will be more than sufficient. Nevertheless, this still represents a large number of logical qubits, implying the need for  $S$  factories throughout the computer and the dedication of most of the qubits in the computer to generate the necessary  $S$  states at a sufficient rate. This will impact the resource counting for our target application, as we discuss in Section 5.

When using an  $S$  state, the actual gate applied will be a random rotation by either  $+\theta$  or  $-\theta$ . Error corrected logical measurement must be used to determine which gate was applied and hence whether a corrective  $2\theta$  gate also needs to be applied. If  $2\theta = \pi/2$ , the correction must be applied before further gates are applied, introducing a temporal gate

ordering. This time ordering prevents arbitrary quantum circuits involving non-Clifford group gates being implemented in constant time.

### 3. Hardware Elements

In considering the hardware in which to implement this architecture, by far the most important pending question is the choice of quantum dot type, which will also determine the semiconductor substrate and operational wavelengths.

#### 3.1. Quantum Dots

The best type of quantum dot to employ remains an open question. Charged, self-assembled InGaAs quantum dots in GaAs are appealing due to their high oscillator strength and near-IR wavelength. These dots have been engineered into cavities in the strong coupling regime<sup>54</sup> and recent experiments have demonstrated complete ultrafast optical control of a single electron spin qubit trapped in the dot<sup>55,56</sup>. However, it is challenging to make high-yield CQED devices from these dots due to their high inhomogeneous broadening and the challenges of site selectivity, although progress continues in designing tunable quantum dots<sup>57,58</sup> in prescribed locations<sup>59</sup>. Sufficient homogeneity for a scalable system, however, may require a more homogeneous kind of quantum dot, such as those defined by a single donor impurity and its associated donor-bound-exciton state. Donor-bound excitons in high quality silicon and GaAs are remarkably homogeneous, both in their optical transitions and in the Larmor frequencies of the bound spin providing the qubit. However, the isolation of single donors in these systems has been challenging. Donor impurities in silicon would seem almost ideal, since isotopic purification can give long spin coherence times<sup>60</sup> and extremely homogeneous optical transitions<sup>61</sup>, but optical control in this system is hindered by silicon's indirect band-gap. A II-VI semiconductor such as ZnSe may provide a nearly ideal compromise – single fluorine impurities in ZnSe have been isolated, shown to have a comparable oscillator strength to quantum dots, and incorporated into microcavities<sup>62</sup>. Recently, sufficient homogeneity has been available to observe interference from photons from independent devices<sup>63</sup>. However, this system comes with its own challenges, such as the less convenient blue emission wavelength. Nitrogen-Vacancy centers in diamond<sup>64,65,66</sup> have also attracted heavy attention recently, but the diamond substrate remains a challenging one for implementing the nanophotonic hardware that supports the quantum computer.

Regardless of the type of quantum dot, there are several common physical features which are to be employed for quantum information processing. The dot has a two-level ground state, provided by the spin of trapped electrons in a global applied magnetic field. This spin provides the physical qubit. The dot also has several optical excited states formed from the addition of an exciton to the dot. One of these excited states forms an optical  $\Lambda$ -system with the two ground states, allowing not only single qubit control via stimulated Raman transitions<sup>67</sup>, but also selective optical phase shifts of dispersive light<sup>68</sup> (to be discussed in Sec. 3.3) or state-selective scattering<sup>69,70,71</sup>. These enable several possible means to achieve entanglement mediated by photons.

### 3.2. Nanophotonics

The quantum dots will be incorporated in small cavities to enhance their interaction with weak optical fields. Cavities may be made from a variety of technologies, including photonic crystal defects and microdisks. Here, we will focus on suspended microdisk cavities.

The small microdisks are in turn coupled to larger waveguides arranged as disks, rings, or straight ridges, which carry qubit-to-qubit communication signals. These waveguides can be ridges topographically raised above the chip surface, or line-defects in photonic crystals. Our present focus is on ridge-type waveguides. Waveguides are well-advanced and relatively low-loss, although it is best to make the waveguides as straight as possible, and to avoid crossing two waveguides in the floor plan. Silicon at telecom wavelengths, for example, makes a good waveguide for our purposes, as it is almost transparent to 1.5  $\mu\text{m}$  light, with a loss of about 0.1 dB/cm. The coherent processing of single photons in on-chip waveguides has recently been well demonstrated for ridge-type silica waveguides<sup>72</sup>.

The “no crossing waveguides” restriction is one of the two key issues driving device layout. The other is the need to route signals to more than one possible destination, for which high-speed, low-loss optical switching is required. Good optical switches are difficult to build: many designs have poor transmission of the desired signals and poor extinction of the undesired ones, and tend to be large and slow. In our architecture, we focus on microdisk-type or microring-type add/drop filters. In suspended silica systems, these switches have been shown to have insertion losses as low as 0.001 dB for the “bus” when the microdisk is off-resonant; optical loss from the bus to the drop port can be as low as 0.3 dB when the system is resonant<sup>73</sup>. On-chip switches in semiconductor platforms do not typically feature such nearly ideal behavior but continue to improve. For example, 40  $\mu\text{m}$  by 12  $\mu\text{m}$  multi-ring add-drop switches with a loss of a few dB were recently demonstrated in a silicon platform<sup>74</sup>.

We need to individually control the resonance of every optical microdisk in the circuit; these microdisks provide the add/drop switches and qubit-hosting cavities. Ultimately, it is the ability to rapidly move these microdisk resonators into and out of near-resonance with the waveguided control light that provides the quantum networking capability. A candidate method for this is to employ the optical nonlinearity of the semiconductor substrate. A strong, below-gap laser beam focused from above onto one of the cavities will shift its index of refraction through a combination of heating, carrier creation, and intrinsic optical nonlinearities<sup>67</sup>. The laser pulses for this may be carried through free space from a micromirror array<sup>75</sup>.

To complete the architecture, we will also need mode-locked lasers for single-qubit control, modulated CW-lasers for quantum non-demolition (QND) measurements as well as deterministic and heralded entanglement gates, and photodiodes to measure the intensity of the control light. Lasers and photodiodes are expensive in both space and manufacturing cost, so an ideal system will be carefully engineered to minimize the number required. Mode-locked lasers with repetition frequency tuned to the Larmor frequency of spin qubits will be used for fast single-qubit rotations<sup>67</sup>. These lasers may be directed by the same micromirror used for switching. More slowly modulated single-frequency lasers will be

used for qubit initialization, measurement, and entanglement operations. These lasers may be incorporated into the chip, or injected via a variety of coupling technologies. The photodiodes are intended to measure intensity of pulses with thousands to millions of photons, rather than single-photon counting, which allows the possibility of fast, on-chip, cavity-enhanced photodiodes; however, off-chip detectors may be more practical depending on the semiconductor employed.

These resources are crucial, as they are needed for every single-qubit measurement and heralded entangling operation. These operations dominate the operation of a cluster-state-based quantum computer. However, these same technologies are evolving rapidly for classical optoelectronic interconnects, and are expected to continue to improve in coming years.

### **3.3. Executing Physical Gates**

Four types of physical gates are employed in this architecture.

The first type of gate is arbitrary single qubit rotations, which may be performed efficiently using picosecond pulses from a semiconductor mode-locked laser with pulse repetition frequency tuned to the qubit's Larmor frequency<sup>67,56</sup>. A cavity is not needed for this operation, and the pulses used are sufficiently far detuned from the qubit and the cavity resonance that the cavity plays little role. The phase and angle of each rotation is determined via switching pulses through fixed delay routes, as described in Ref. 67. The performance of this gate is limited by spurious excitations created in the vicinity of the quantum dot by the pulse<sup>76</sup> and not by optical loss or other architectural considerations.

The next type of gate is the quantum-non-demolition QND measurement of a single qubit. This gate is critical, since the initialization and measurement of every qubit is very frequent in our tQEC architecture, and the QND gate allows both. A QND measurement makes use of the optical microcavity containing the dot, and operates with the cavity well detuned from the dot's optical transitions. In such a configuration, an optical transition to one qubit ground state may present a different effective index of refraction for a cavity mode than the optical transition to the other qubit ground state. This results in a qubit-dependent optical phase shift of a slow optical pulse coupled in and out of the waveguide. This optical pulse may then be mixed with an unshifted pulse from the same laser to accomplish a homodyne measurement of the phase shift. In one variation of this scheme, this phase is detected as a change in the polarization direction of a linearly polarized optical probe beam; this has been demonstrated for quantum dots both with<sup>77</sup> and without<sup>78</sup> a microcavity; larger phase shifts have also been observed in neutral dots in improved photonic crystal cavities<sup>79</sup>. Simulations indicate that pulses with a timescale of about 100 ps may be used for this gate<sup>68,67</sup>.

These first two gate types are single-qubit gates. For generating entanglement between distant qubits, two further gates are employed: a deterministic, nearest-neighbor gate, and a non-deterministic gate for heralded entanglement generation for distant qubits.

The deterministic, nearest-neighbor gate will be mediated by a common microdisk mode connecting the cavities joining nearby qubits. The phase or amplitude of this cavity

mode may be altered by the state of the qubits with which it interacts, which in turn changes the phase or population of those qubits. The gate is achieved by driving the coupled cavity mode with one or more appropriately modulated optical pulses from a CW laser. The light is allowed to leak out of the cavity and may then be discarded. The amplitude version of such a gate was proposed in 1999 by Imamoglu et al.<sup>80</sup>, and may be viewed as a pair of stimulated Raman transitions for two qubits driven by two CW lasers and their common cavity mode. This gate is known to require high- $Q$  cavities. The phase version of this gate, described in Ref. 81, is an adaptation of the “qubus” gates proposed by Spiller et al. in 2006<sup>82</sup>; more detailed design and simulation of this gate in the present context is in progress<sup>83</sup>.

If such deterministic gates are available, one may naturally ask whether a fully two-dimensional architecture of coupled qubits is more viable than the communication-based architecture we present here. Indeed, if truly reliable cavity QED systems can be developed in the large-scale, deterministic photonic-based gates<sup>84</sup> may enable highly promising single-photon-based architectures for tQEC<sup>85</sup>. However, the devices that will enable deterministic CQED gates in solid-state systems are unlikely to be fully reliable.

In particular, high-fidelity deterministic gates require extremely low optical loss between qubits, and therefore cannot easily survive coupling to straight waveguides or to other elements in the photonic circuit such as switches and fibers. For generating entanglement through these elements, stochastic but heralded entanglement schemes are used, similar to gates in linear optics except with physical quantum memory. Combined with local single-qubit rotations, QND measurements, and deterministic nearest-neighbor gates, this heralded entanglement allows quantum teleportation. Heralded entanglement is the bottleneck resource in quantum wiring. Heralded entanglement gates come in several flavors, but fortunately each type requires the same basic qubit and cavity resource; they vary in the strength of the optical field used and the method of optical detection. Which type to employ depends on the amount of loss between the qubits to be entangled.

For qubits with relatively low loss between them, such as those coupled to a common waveguide without traversing to the drop port of a switch, so-called “hybrid” schemes are attractive<sup>86,68</sup>. In these schemes, the QND measurement discussed above is extended to two qubits, distinguishing odd-parity qubit subspaces from even-parity states. For some detection schemes, such as  $x$ -homodyne detection, this parity gate may be deterministic, up to single-qubit operations which depend on measurement results<sup>87,88</sup>. If such parity gates are available, “repeat-until-success” schemes for quantum computation are very attractive<sup>89</sup>, and have been proposed for use in multicomputer-like distributed systems<sup>90</sup>. However, if weak CQED nonlinearities are employed with lossy waveguides, these detection schemes fail<sup>86,68</sup>. In this case,  $p$ -homodyne detection may still show strong performance, but the parity gate is incomplete. The heralded measurement of an odd-parity state may project qubits into an entangled state with probability  $\simeq 50\%$ , but when this fails no entanglement is present. As in schemes using linear optics, this allows probabilistic quantum logic. With the addition of an extra ancilla qubit, this partial parity-gate may be combined into a probabilistic CNOT gate for entanglement purification.

This scheme is attractive due to its use of relatively bright laser light and near ideal

probability of successful heralding. However, it is strongly subject to loss, as has been discussed previously<sup>68</sup>. More complex measurement schemes may improve the fidelity of such gates at the expense of their probability of heralding a success<sup>91</sup>. For very lossy connections, the number of photons in the optical pulse might be reduced to an average of less than one photon, in which case single-photon scattering schemes<sup>69,70,71</sup> would be employed. These schemes succeed much more infrequently, as they rely on the click of a single photon detector projecting the combined qubit/photon system into one where no photons were lost, a possibility whose probability decreases with loss. Here, we consider only many-photon qubus gates using homodyne detection as discussed in Ref. 68; we compensate for different connections with different loss rates only by changing the intensity of the optical pulses employed, whose optimum varies with loss. The detection scheme remains constant across the architecture.

Although proposals for nonlocal, deterministic gates exist, their performance is always hindered by optical loss. This is an inevitability: if photons are mediating information between qubits, the loss of those photons into the environment inevitably reveals some information about the quantum states of the qubits, causing decoherence. A well-designed photon-mediated architecture should use a hierarchy of photon-mediation schemes to provide high-success-probability gates at low distances and highly loss-tolerant gates at higher distances, and the qubus mechanisms allow some degree of hierarchical tuning without adding extra physical resources.

In the present discussion, we discuss performance entirely in terms of optical loss. Photons may be lost in waveguides, from cavities, from the cavity-waveguide interfaces, and from spontaneous emission. An approximation of the amount of decoherence-causing loss at a quantum-dot-loaded cavity and cavity/waveguide interface, when running hybrid CQED-based gates optimally, is the inverse of the cooperativity factor  $C$ <sup>68</sup>. This factor arises from the ratio of spontaneous emission into a cavity mode (assumed to be overcoupled to the waveguide) to spontaneous emission into other modes. It scales as the quality factor of the cavity divided by its mode volume, so the cavities containing qubits are designed small to maximize this factor. When we discuss qubit-to-qubit optical loss, this loss should be considered as the linear loss in the waveguide connecting the qubits plus about  $C^{-1}$ . Cooperativity factors between self-assembled quantum dots and the whispering gallery modes of suspended microdisks have been shown to approach 100<sup>92,93</sup>, corresponding to a cavity-induced loss limit of 0.04 dB.

#### **4. Architecture: Layout and Operational Basics**

In this section, we qualitatively describe our architecture and its operation. Many of the design decisions described here will be justified numerically in Section 5.

##### **4.1. Architecture Axes**

The basic structural element of our system is one-dimensional: a waveguide with a tangent series of microdisks, each connected to one or more smaller microdisks containing quantum dots, as in Fig. 2. The shared bus nature of a single waveguide offers the advantage

that the qubit at one end can communicate quickly and easily with the qubit at the other end; this long-distance interaction has the potential to accelerate some algorithms and aids in defect tolerance, as we will show below. However, that shared nature makes the bus itself a performance *bottleneck* in the system, as contention for access to the bus and the measurement device forces some actions to be postponed<sup>94</sup>.

This limitation on concurrent operation makes it natural to consider using multiple columns. Columns are connected by teleportation, aided by heralded entanglement and purification. The resulting structure, developed in Figures 2 to 5, is a set of many columns, defined by long, vertical waveguides, interspersed with smaller, circular and oval waveguides, and qubits in cavities tangential to the waveguides. The vertical waveguides are of two types: *logic* waveguides, which are used to execute operations between qubits within one column, and *teleportation* waveguides, which are used to create and purify connections between columns within a single chip or between chips. The small, colored circles represent the smallest microcavities containing quantum-dot qubits. The different colors represent different roles for particular qubits, which we describe in Section 4.2. The teleportation columns do not use the smaller, higher- $Q$  circular waveguides to couple qubits deterministically. Instead, as in Figures 3 and 4, they use larger racetrack-shaped waveguides that can support a larger number of qubits which are only stochastically entangled, called transceiver qubits. The qubits along one racetrack can be used to purify ancilla qubits, allowing us to connect qubits in potentially distant parts of the chip, or to connect to off-chip resources.

The architecture in Fig. 5 is designed to minimize both the length of waveguides and the number of switches traversed by pulses carrying quantum information. Note that signals introduced onto the waveguide snaking through the chip will not be perfectly switched into the detectors, implying some accumulated noise; however, this effect can be mitigated with appropriate detector time binning and sufficiently large microdisk  $Q$ -factors in the switches.

A single node has two axes of growth. The length of a logical waveguide column and the number of columns provide the basic rectangular layout, which will have some flexibility but is ultimately limited by the size of chip that can be practically fabricated, packaged and used. To give a concrete example, if we set the vertical spacing of the red lattice qubits to  $50\ \mu\text{m}$  and the column-to-column spacing to  $100\ \mu\text{m}$ , 100 qubits in each vertical column and 100 columns will result in the active area of the chip being 5 mm by 10 mm.

A third axis of growth is the number of chips that are connected into the overall system – the number of nodes in our multicomputer. In previous work, we have been concerned with the topology and richness of the interconnection network between the nodes of a multicomputer using CSS codes, finding that a linear network is adequate for many purposes<sup>95,94</sup>. The extension of nodes into the serpentine teleportation waveguide in Fig. 5 enables such a linear-network multicomputer, although the additional necessary resources for bridging lossier chip-to-chip connections will not be considered here.

The structures in our architecture are large by modern VLSI standards; the principle fabrication difficulty is accurate creation of the gap between the cavities and the waveguides. That spacing must be 10-100nm, depending on the microdisk and waveguide size

and quality factors<sup>93</sup>. The roughness of the cavity edge is a key fabrication characteristic that determines the quality of the cavity, and ultimately the success of our device.

Although the device architecture and quantum dot technology are not yet fixed, we include images of test-devices fabricated using e-beam lithography following the methodology described in Ref. 93, only to help visualize future devices. Figures 2 and 3 include scanning electron microscope images of a device created in a GaAs wafer containing a layer of self-assembled InAs quantum dots<sup>93</sup>. More scalable fabrication techniques than e-beam lithography must ultimately be developed for scalability; promising routes include nanoimprint lithography<sup>96</sup> and deep sub-wavelength photolithography<sup>97,98,99</sup>.

#### 4.2. Qubit Roles and Basic Circuits

The different colors for the qubit quantum dots in Figure 3 represent different roles within the system. Physically, the cavities are identical, but they are coupled to different waveguides, allowing them to interact directly with different sets of qubits. Within those connectivity constraints, their roles are software-defined and flexible. Finding the correct hardware balance among the separate roles is a key engineering problem. The answer will depend on many parameters of the physical system, including the losses in switches and couplers, and will no doubt change with each successive technological generation.

The red qubits in the figures, in the column vertically placed between the larger circles, are the *lattice* qubits. Those that are functional are assigned an effective  $(x, y)$  position in the 2-D lattice used to implement tQEC. These are subsequently divided into *code* qubits, which are never directly measured, and *syndrome* qubits, which are regularly measured following connections to code qubits in order to maintain the topologically protected surface code. The ideal number and density of syndrome qubits among code qubits depends on the yield. Within a column, all functional nearest neighbor pairs of qubits can be coupled in parallel. Non-nearest-neighbor couplings can only occur sequentially. For very low yields, in which code qubits rarely have nearest-neighbor couplings, only a few syndrome qubits per column are required as the syndrome circuits must largely be implemented sequentially, implying the syndrome qubits can be reused.

The blue qubits, or *transceiver* qubits, are aligned with the racetracks and the long purification waveguides. These qubits are used to create Bell pairs between column groups within the same device, or between devices. Because purification is a very resource-intensive process, the transceiver qubits are numerically the dominant type.

The green qubits, sandwiched between the column of circles and the column of racetracks, are *ancilla* qubits, used to deterministically connect stochastically created entangled states among (blue) transceiver qubits to (red) lattice qubits. The green qubits also play an auxiliary role during the purification of the blue qubits.

The circuit, or program, for executing purification on the blue qubits is shown in Figure 3. The blue qubits have previously been measured and are thus initialized to a known state. Then, qubits in a given teleportation column of Figure 5 are entangled with qubits in either the same column or the one neighbouring it to the right using the heralded entanglement generation technique discussed in Sec. 3.3. Note that waveguide loss prevents the

efficient entangling of qubits in widely separated teleportation columns. In general, a laser pulse is inserted in the teleportation waveguide at a given column, coupled with a qubit in that column, coupled with a second qubit either in that column or the one neighbouring it to its right and then switched out of the teleportation waveguide and measured. This process is repeated in rapid succession, building a pool of low-fidelity entangled pairs, creating the  $|\Psi^+\rangle$  states at the left edge of Figure 3.

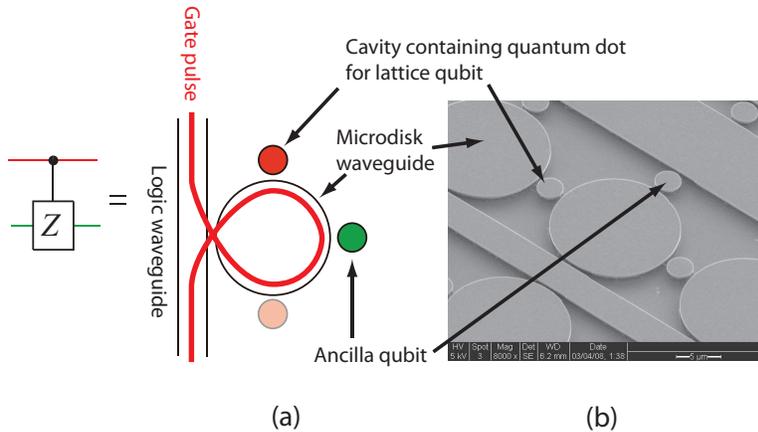


Fig. 2. (a) Layout and pulse path for executing a local, high-fidelity controlled-Z gate. An optical pulse couples from the straight waveguide to the microdisk waveguide; the two qubits of interest are introduced to the logic gate by bringing their cavities into resonance with the optical pulse. (b) Scanning electron micrograph of a non-functional demonstration device, fabricated in GaAs with (unshown) InAs quantum dot layer. The structures are underetched following the methods presented in Ref. 93.

Once the base-level entangled pairs are created, the circuit in Figure 3 is executed within each column, which employs two probabilistic parity gates to achieve the controlled-NOT operations used in entanglement purification. Purification proceeds until entangled state fidelities are considered sufficient for computation. At that time the purified entanglement between blue transceiver qubits is used to make an appropriate entangled (green) ancilla which are connected to the target lattice qubits.

Finally, the high-fidelity Bell pairs are used to create the tQEC lattice, using the clustering circuit shown in Fig. 4.

#### 4.3. Lattice

The most important issue in the generation of a cluster state in our geometry is the physical asymmetry between connections within a column, those with other columns, and those between dies. The hierarchy of connection distances in our system will be characterized in terms of the number of laser pulses and measurements required to achieve entanglement of a particular fidelity.

Entangling two qubits connected to the same circular waveguide is straightforward; we can refer to these as “cavity connected” or “C-connected.” Racetracks are a longer, and

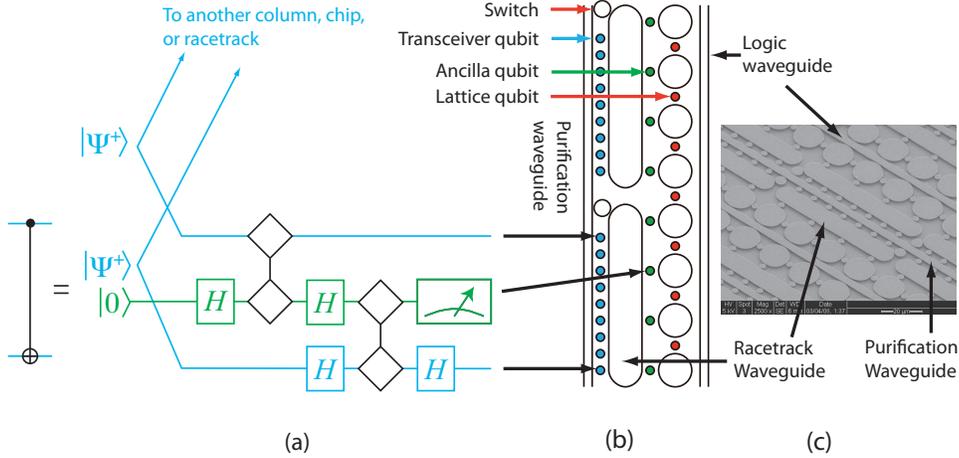


Fig. 3. (a) Partial circuit for executing purification on long-distance Bell pairs. The diamonds represent a probabilistic parity gate which projects two qubits into an odd-parity subspace with probability of approximately 50%. These gates are achieved via pulses routed through the racetrack waveguides via the ring-waveguide labelled “switch”. All measurements are in the  $X$  basis. (b) The basic layout unit is a column of racetrack and circular waveguides sandwiched between the straight purification and logic waveguides. (c) Zoom-out of the same device shown in Fig. 2(b).

slightly lower-fidelity, form of cavity; we refer to two ancillae or two transceiver qubits on the same racetrack as “R-connected”, or racetrack-connected. Two lattice qubits connected through an R-connected Bell pair are said to be indirectly connected, or “I-connected”.

Within a logic column, many deterministic gates on C-connected qubits can be performed without purification, and a high level of parallelism may be employed. The pulses that execute deterministic gates on the logic waveguide couple into the cavities only weakly, and do not need to be measured after the gate, making it possible that the same strong pulse could be used to execute several gates concurrently. If we label the qubits with the pattern  $ABABA\dots$ , we may be able to couple all of the  $AB$  pairs in one entangling time slot, then couple all of the  $BA$  pairs in the second time slot.

The fidelity of W connections is dominated by the efficiency of coupling pulses into and out of cavities, as the loss in the waveguide will be negligible. When connecting two lattice qubits in columns separated by a purification waveguide, we require moderate amounts of purification. The purification ancillae are themselves W-connected; the post-purification lattice connection we refer to as “ $P_W$ -connected”.

Finally, qubits that do not share the same purification waveguide must be connected using a pulse that transits one or more switches. We refer to these physical connections as  $X$  or  $X_{i,j}$  connections, where  $i$  is the number of switches and  $j$  is the number of I/O ports that must be transited. Lattice qubits connected after purification we refer to as  $P_X$ -connected.

The  $P_W$ -connections and  $P_X$ -connections will be most strongly subject to bottlenecks from the limited number of laser pulses and detection events in our architecture, and are therefore the focus of our numeric studies in the next section.

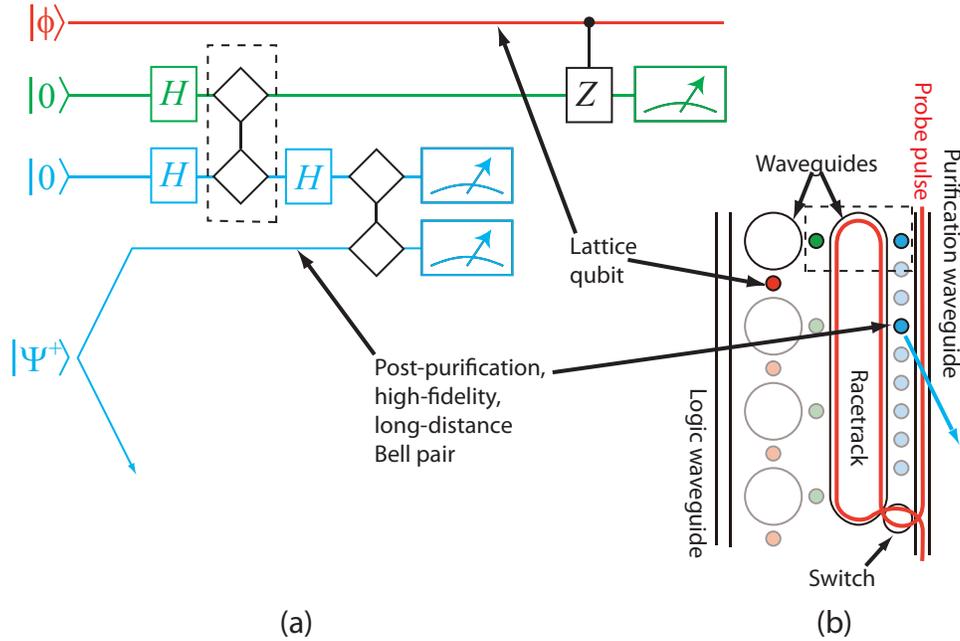
16 *Van Meter, Ladd, Fowler and Yamamoto*


Fig. 4. (a) Partial circuit and (b) qubit/cavity layout and pulse path for executing long-distance clustering operations. This circuit and a matching one elsewhere in the system execute the logical controlled-Z gate between two lattice (red) qubits in a teleported fashion (which we call telegate) by using a high-fidelity Bell pair built on transceiver (blue) qubits. The four qubits used in this circuit are highlighted in the layout. The second transceiver qubit and the ancilla (green) are used as ancillae in this circuit. The diamonds represent probabilistic ( $P \approx 50\%$ ) parity gates on the racetrack-shaped waveguide, between either the two transceiver qubits or the transceiver and the ancilla. The gate in the dashed-line box in (a) is executed by enabling the two qubits in the box in (b). All measurements are in the X basis. The physical CZ gate in the top row is performed using the circuit of Figure 2.

## 5. Resource Estimates

Given a set of technological constraints (pulse rate, error rate, qubit size, maximum die size), a complete architecture will balance a set of tradeoffs to find a sweet spot that efficiently meets the system requirements (application performance, success probability, cost). Minimizing lattice refresh time is the key to both application-level performance and fault tolerance, but demands increased parallelism (hence cost); in our system, this favors a very wide, shallow lattice, which is more difficult to use effectively at the application level. Increasing the number of application qubits increases the parallelism of many applications (including the modular exponentiation that is the bottleneck for Shor's algorithm), but if the space dedicated to the singular factory does not increase proportionally, performance will not improve.

We begin by describing the communication costs and the impact of loss on the lattice refresh cycle time in a generic 2-D multicomputer layout, from which we can calculate the effective logical clock cycle time for executing gates on application qubits. With these concepts in hand, we then propose an architecture, and calculate its prospective performance.

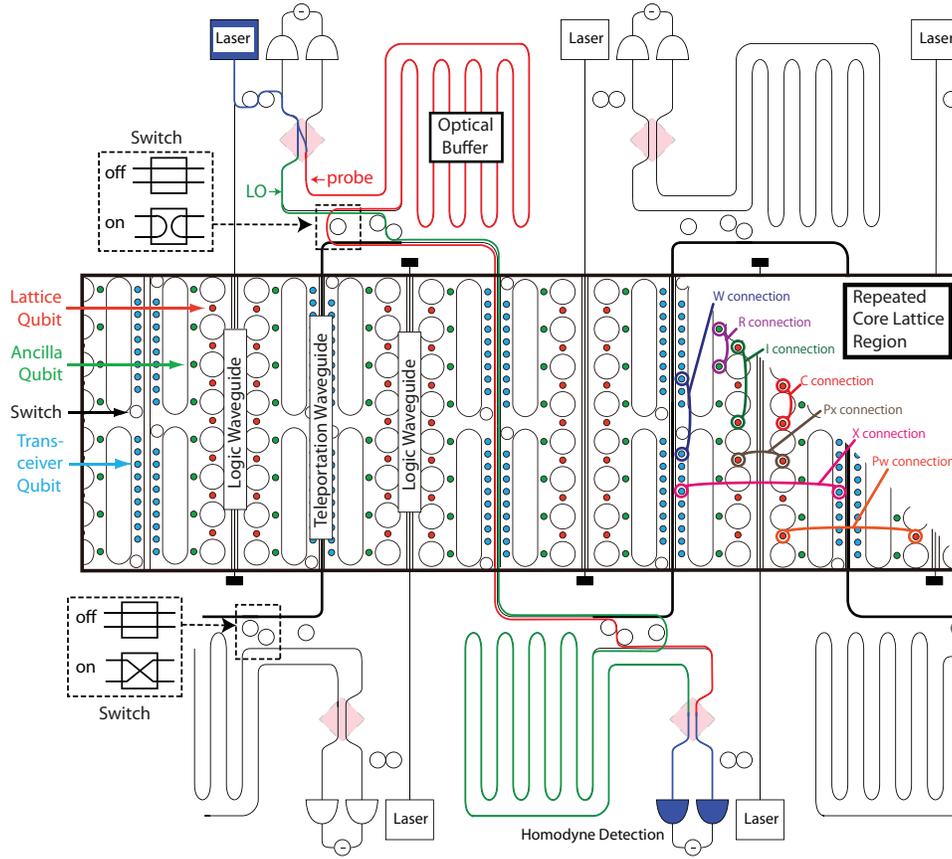


Fig. 5. The nanophotonic quantum multicomputer architecture. Small microdisks containing lattice, ancilla, and transceiver qubits are color-coded while waveguides and microdisk-based add-drop switches are indicated by black lines. This schematic indicates the critical elements of the nanophotonic chip-layout described in the text, but the structures shown are not to-scale. In particular, the modulated CW lasers and detectors shown are the largest elements and are likely to be off-chip. The pink squares indicate the location of beam-splitters defined by evanescently coupled ridge-waveguides, which split a single laser pulse (indicated by a blue line) into probe (red line) and local oscillator (LO, green line) optical pulses. These pulses travel two paths; one is buffered by a serpentine waveguide which delays the probe by several times the pulse width of approximately 100 ps. (The pulse colors are schematic only; these pulses are to be monochromatic.) The probe is switched to follow the LO along the same route through the teleportation waveguides of the core chip, which depend on the qubits to be coupled. Single passes from top-to-bottom, such as the one shown by the red and green lines, enable the similar “W connections” and “Pw connections” between qubits as shown on the right. A U-shaped path (not-shown) would enable the longer-distance “X” and “Px” connections. Lasers directly coupled into waveguides enable C connections and mediate logic within the circular microdisks connecting lattice qubits to ancilla qubits. The rectangular region in the center is repeated many times vertically and horizontally.

### 5.1. Communications and Lattice Refresh

Figure 6 shows the residual infidelity and the cost in teleportation waveguide pulses as a function of the loss in the probe beam from qubit to qubit through the waveguides. Purification is performed using only Bell pairs of symmetric fidelities, and is run until final

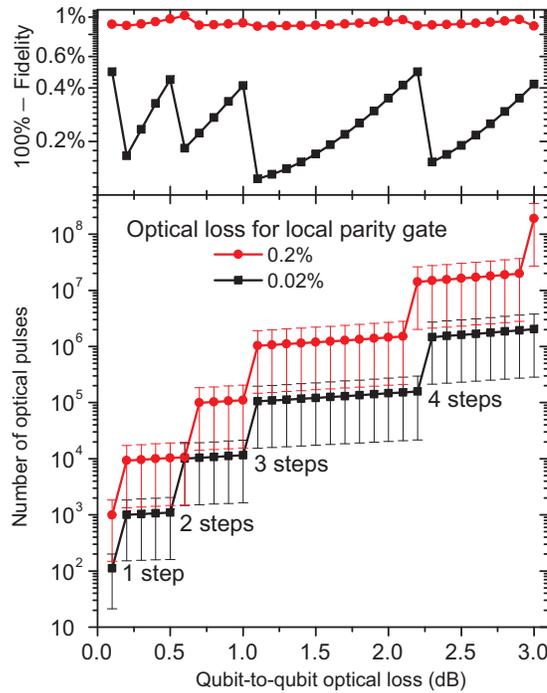


Fig. 6. For qubus connections, impact of signal loss on the final fidelity achievable using symmetric purification. Error bars represent the RMS of the number of pulses, which is close to the average number; the distribution is strongly Poisson-like.

fidelity saturates or until fidelity is better than 99.5%. The two curves represent two values of round-trip loss in the racetrack waveguides used for local parity gates; with local loss of 0.2%, we cannot achieve a final fidelity above the threshold for tQEC. Thus, we establish an engineering goal of 0.02% loss or better.

The values in Fig. 6 are calculated by generating a Markov probability matrix for the protocol of symmetric purification<sup>100</sup>, where each matrix transition requires the generation and detection of an optical pulse in the teleportation waveguide. Probabilities and fidelities for each step are found using the formalism presented in Ref. 68. Many of these transitions are deterministic, but some are not due to the probability of parity gates failing or the purification protocol failing. Exponentiation of this matrix allows the direct calculation of the probability of completing the protocol in a given number of steps, allowing calculation of the probability density function for completion of purification vs. number of optical pulses. These probability distributions are strongly Poissonian. They are used to calculate the average and root-mean-square number of pulses plotted in Fig. 6.

This Markov analysis is useful for estimating performance, but overestimates the required spatial and temporal resources considerably. The strictly symmetric purification

routine assumed here makes less than ideal use of qubit memory; alternative resource management strategies can lead to order-of-magnitude improvements in speed without a comparable increase in size, as considered, for example, in Ref. 30. Also, the calculation we have performed assumes that when parity gates fail in the circuit shown in Fig. 3(a), the entire procedure fails and entangled pairs must be regenerated and repurified. In fact, if one parity gate succeeds and the other fails, then one Bell pair preserves some of its entanglement and may be kept, possibly with a Pauli correction, for subsequent purification rounds. Optimizing the purification procedure to account for such possibilities is difficult to do analytically; Monte Carlo simulations such as those in Ref. 30 may estimate the worth of these strategies, but we leave such simulations for future work.

With the proper layout, we can connect multiple chips into a two-dimensional structure. With  $V$  rows of  $H$  chips each, and a chip that consists of  $C$  columns each containing  $R$  rows of lattice qubits, we have a physical structure capable of supporting an  $HC \times VR$  lattice. In such a multicomputer, entangling pulses may be destined for another qubit in the same column in the same chip, another qubit in the same column but the chip below, or in the neighboring column to the left or right. With multiple possible destinations, switching is naturally required; we can arrange the switching so that vertical connections are  $X_{1,1}$  connections and horizontal ones are  $X_{2,1}$  connections. Assessing the scalability of such a system and establishing guidelines for configuring the system depend on understanding these connections.

Table 1 lists the costs for the lattice building operations on such a switched multicomputer architecture. We compare two logical lattices, a direct-mapped  $HC \times VR$  logical lattice and a sub-lattice-organized  $HCs \times VR/s$  logical lattice in which each physical column is used as a small  $R/s \times s$  lattice<sup>b</sup>. The physical yield affects the probability that two neighboring lattice qubits and their shared ancilla are good, and hence the probability that a  $C$  connection can be used. Additionally, for low yields ( $y < 0.8$ ), we assign only a few qubits per column as tQEC syndrome qubits, forcing all lattice cycle operations to use  $P_W$ -connected gates.

Table 1. Number and types of connections per physical waveguide for lattice-building for an  $H \times V$  multicomputer with  $C \times R$  lattice qubits per node and  $HC$  total laser input ports and lattice sub-factor  $s$ . Expressions assume  $R \bmod s = 0$ .  $R_f = Ry_e = Ry_p(1 - (1 - y_p)^2)$ , the functional number of qubits in a column.

Connection type	100% yield	physical yield $y_p$
$C$	$2V(R - s)$	$n_C = 2V(R_f - s)y_p^2$ (for $y_p \geq 0.8$ ) or 0 ( $y_p < 0.8$ )
$P_W$	$V(2R - R/s)$	$n_W = V(2R_f - R_f/s) + 2V(R_f - s) - n_C$
$V$ neighbor ( $P_X(X_{1,1})$ )	$2s(V - 1)$	$n_{X1} = 2s(V - 1)$
$H$ neighbor ( $P_X(X_{2,1})$ )	$VR/s$	$n_{X2} = VR_f/s$

We observe several qualitative facts about this architecture:

<sup>b</sup>The table assumes that  $R \bmod s = 0$ . Although that is not a requirement, the expressions are more complex for  $R \bmod s \neq 0$ ; without careful structuring, potentially as many as half of the  $P_W$  connections may become  $P_X$  for  $X_{1,1}$ .

- The lattice cycle time is constant as  $H$  increases, but the number of lasers and measurement devices must increase proportionally.
- To first order, the lattice cycle time scales linearly with  $VR$ , but second-order effects will likely make it worse than linear.
- The number of  $X_{2,1}$  connections favors a sub-lattice with a large  $s$ , but the minimum size of the logical lattice limits  $s$ ; we require  $14d \leq VR/s$ .
- Increasing lattice cycle time hurts fidelity due to memory degradation.
- Increasing lattice cycle time hurts application performance.

The total lattice refresh cycle time is  $t_{lat} = t_{pulse}p_{lat}$ , where  $p_{lat}$  is the number of pulse time steps in the complete cycle. The final, logical clock rate for application gates depends on both the refresh cycle and the temporal extent of the lattice holes as they move through the system to execute logical gates. We can visualize the movement of the holes through the temporal dimension as “pipes” routed in a pseudo-3-D space. To maintain the same  $4d$  perimeter and spacing about the hole as it extends into the temporal dimension, each hole movement will also have to extend for  $5d$  lattice refresh cycles. We have used  $d = 14$  as the length of one side of each square hole. The temporal spacing must be  $4d = 56$ , implying that the fastest rate at which hole braiding can occur is  $5d = 70$  lattice refresh cycles.

In our architecture, the logical clock rate is  $\Omega(d^2)$ . The number of refresh cycles per logical gate is  $\Theta(d)$ . The refresh time itself is  $\Omega(R) = \Omega(d)$ ; because we must choose  $R \propto d$ , the number of pulses grows at least linearly in  $d$ . As the columns lengthen, fidelity falls and the number of pulses per cycle grows, creating a positive feedback in  $d$  and cycle time.

## 5.2. Proposed Architecture and Performance

Table 2 summarizes our initial strawman architecture, depicted in Fig. 5. To factor an  $n$ -bit number using Shor’s algorithm, we would like to have  $6n$  logical qubits. Having established a goal of factoring a 2,048-bit number, we need 12,288 logical qubits.

Ultimately, the execution of application algorithms in tQEC requires, as at the physical level, two components: communication and computation. Logical communication consists of routing the pipes through the pseudo-3-D lattice. These pipes can route through the space with only a fixed temporal extent, allowing the equivalent of “long distance” gates in the circuit model. They do, however, consume space in the lattice, creating a direct tradeoff between the physical size of the system and the time consumed. Additionally, the shape of the logical lattice determines how efficiently logical qubits can be placed and routed. We assign 25% of the logical qubit space for wiring and hole movement space.

Computation, for many algorithms, will be dominated by Toffoli gates; as some of the operations are probabilistic, an average of over ten  $S$  and  $T$  states are required for each. Shor’s algorithm requires some  $40n^3$  Toffoli gates:  $5n^2$  adder calls<sup>102</sup> (after optimizations to modulo arithmetic and one level of indirection in the arithmetic<sup>103</sup>), each requiring  $10n$  Toffoli gates<sup>104</sup>. The total of  $40n^3 = 3.2 \times 10^{11}$  Toffoli gates require over  $10^{12}$   $S$  states. Again, a direct tradeoff can be made between space and time, as the  $S$  states can be

Table 2. Summary of our proposed serpentine, add-drop filter architecture.  $M = 2^{20} \sim 10^6$ .

<i>System Hardware</i>	
Chip lattice, $C \times R$	$128 \times 770$
Multicomputer setup, $H \times V$	$65536 \times 1$
Physical lattice size (in qubits)	$8M \times 770 = 6.46 \times 10^9$
Laser ports	4M
Measurement devices	16M
Purification/entanglement pulse rate	10 GHz
Switch type	add-drop filter
Required physical yield	$y_p = 40\%$
Effective yield for lattice qubits	$y_e = y_p(1 - (1 - y_p)^2) = 25.6\%$
Functional column height	$R_f = Ry_e = 196$
Required local optical loss	0.02%
Required adjusted gate error rate	$p_{err} \leq p_{thresh}/4 \sim 0.2\%$
Required memory coherence time	$t_{mem} \geq 1000t_{lat} = 49 \text{ msec}$
<i>Communication Costs</i>	
$W, P_W$ connection	0.1dB, $p_W = 111$ pulses
$X_{0,0}, P_X$ conn. (neighboring column)	0.4dB, $p_X = 1068$ pulses
<i>Lattice Operations</i>	
Sub-lattice factor $s$	1
Logical lattice	$8M \times 196$
Pulses per lattice cycle (avg.)	$p_{lat} \sim n_W p_W + n_{X2} p_X = 4.9 \times 10^5$
Lattice cycle time	$t_{lat} = p_{lat} t_{pulse} = 49 \mu\text{sec}$
<i>Logical Qubit Operations</i>	
Hole separation constant	$d = 14$
Lattice area per qubit (at rest, loosely packed)	$14d \times 9d = 196 \times 126 = 24696$
Lattice area per qubit (at rest, tightly packed)	$10d \times 5d = 140 \times 70 = 9800$
Hole movement time	$t_{move} = 5dt_{lat} = 3.41 \text{ msec}$
Hole braiding time	$t_{braid} = 5dt_{lat} = 3.41 \text{ msec}$
Toffoli gate construction	Nielsen & Chuang <sup>101</sup> , p. 182
Finished $ S\rangle$ states per Toffoli gate (avg.)	11.5
Total braidings of $ S\rangle$ states per Toffoli	1795
Toffoli gate time $t_{tof}$	$\sim 14t_{braid} = 48 \text{ msec}$
<i>Application Operations</i>	
Maximum capacity, in logical qubits	119836
Number of application logical qubits	$6n = 12288$
$ S\rangle$ factory space	77589
“wiring” space	$25\% = 29959$
<i>Shor</i>	
Length of number to be factored	$n = 2048$
Adder	Carry-lookahead
Adder time	$t_{add} = 4 \log_2 n t_{tof} = 2.1 \text{ seconds}$
Modulo & indirect arithmetic	$w = 2, p = 11, \sim 5 \times$ faster than basic VBE <sup>102,103</sup>
Number of adder calls	$n_{add} = 4n^2 = 1.68 \times 10^7$
Number of adders executed in parallel	1
Number of Toffoli gates	$n_{tof} = 40n^3 = 3.2 \times 10^{11}$
Time to execute algorithm only	$3.5 \times 10^7$ seconds (409 days)
Time to create singular states	$2.7 \times 10^7$ seconds (314 days)
Final execution time	409 days

built in parallel. For our system and this size of problem, rough balance is achieved with about 65% of the logical qubits dedicated to the  $|S\rangle$  factory.

The multicomputer organization is wide and shallow, to minimize refresh cycle time.

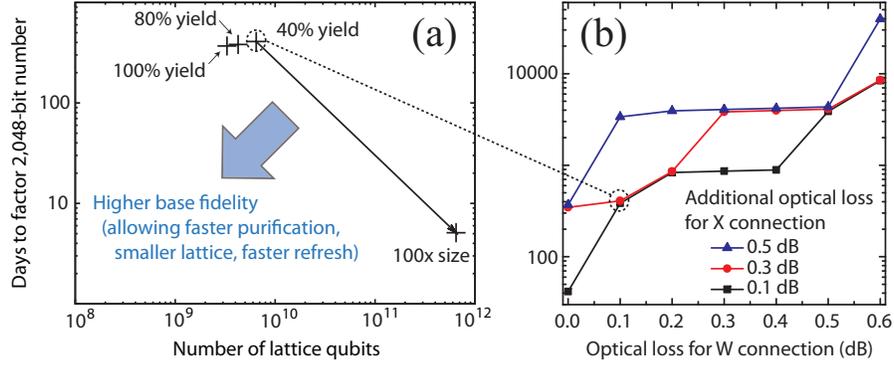


Fig. 7. Factoring time for 2,048-bit number using Shor's factoring algorithm. a) Our baseline proposal, with 40% yield, 0.1dB W connections and 0.4dB X connections, can be improved by increasing the size and application-level parallelism of the system. Improving yield above 40% reduces necessary resources only moderately, but raising the fidelity of the base-level entangled pairs has a major impact on both system size and performance. b) Achieving low-loss connections is critical to performance.

Once we have decided to limit  $V$  to 1, the detailed chip layout simplifies, allowing the serpentine waveguide shown in Fig. 5. In this architecture, W connections are high fidelity, there are no  $V$  neighbors ( $X_{1,1}$  connections), and connections to neighboring columns need not leave the chip except at chip boundaries. The  $n_{X2}$  from Table 1 is still  $VR_f/s$ , but physical connections are X connections with a loss of only about 0.4dB. The vertical height of a single chip will only accommodate enough cavities for a direct-mapped lattice,  $s = 1$ .

Figure 7a shows the execution time for our proposed system. A 2048-bit number should be factorable in just over 400 days, if the technological characteristics in Table 2 can be met. The system is large, requiring more than six billion lattice qubits and several times that total number when ancillae and transceivers are included. At the application level, much more parallelism is available if a larger system is built. A system one hundred times larger would factor the number in about five days.

Figure 7b shows execution time as a function of the loss in our two key connection types, the intra-column W connections and the inter-column X connections. Minimizing the additional loss incurred in inter-column travel helps hold execution time within reasonable bounds.

Reaching toward the desirable lower left corner of Fig. 7a requires improving the base-level entanglement fidelity or reducing the number of pulses used to purify Bell pairs. Our system is fairly robust to yield. Below 40% it is difficult to build a system capable of running tQEC, but above that level, increasing yield has only minor effects on temporal and spatial resources. This gives a clear message: pursue fidelity and quality of components at the expense of yield.

## 6. Discussion

Our design focuses on the communications within a quantum computer, building on a natural hierarchy of connectivity ranging from direct coupling of neighbors on one physical axis of our chip through medium-fidelity, waveguide-based purification coupling on the other axis, to distant, switched connections requiring substantial purification. Thus, while we refer to our design as a quantum multicomputer with each node consisting of a single chip, it is more accurate to regard the connections between qubits as occurring on a set of levels rather than a simple internal/external distinction. Founded on quantum dots connected via cavity QED and nanophotonic waveguides and using topological error correction, this proposal represents progress toward a practical quantum computer architecture. The physical technologies are maturing rapidly, and tQEC offers both operational flexibility and a high threshold on realistic architectures such as ours.

While the overall architecture (multicomputer) and the system building blocks (tQEC, purification circuits, etc.) have been established, much work remains to be done. The most important pending decision is the actual choice of semiconductor and quantum dot type. The cavity  $Q$  and memory lifetime, which dramatically affect our ability to build and maintain the lattice cluster state, will be critical factors in this decision. The yield of functional qubits will ultimately drive the types of experiments that are feasible.

With the decision of semiconductor and the key technical parameters in hand, it will become possible to more quantitatively analyze the mid-level design choices of node size, layout tradeoffs, and the numbers of required lasers and photodiodes. The control system for managing the qubits and cavity coupling will be a large engineering effort involving optics, electronic circuits, and possibly micromechanical elements. Finally, application algorithms need to be implemented and optimized and run-time systems deployed, which will require the creation of large software tool suites.

One of our goals in this work is to establish target values for experimental parameters that must be achieved for such a large system to work. For the chip design and system configuration we present here, we estimate that the yield of functional quantum dots must be at least 40%, the local optical loss must be better than 0.02%, the adjusted gate error rate better than 0.2%, and the memory coherence time about 50 milliseconds or more. The exact values of these goals depend on the architecture, system scale, and application; the entire system is summarized in Table 2.

As a final comment, the physical resources demanded by this architecture are daunting. Other architectures for quantum computers are comparably daunting. The current work is intended in large part to reveal the scope of the problem. With realistic resources such as lossy waveguides, finite-yield qubits, and finite chip-sizes, the added overhead for error correction makes quantum computers very expensive by current standards. We must rely on engineering advancements to improve nanophotonic and quantum dot devices as well as VLSI-like manufacturing capabilities to realize a quantum computer with a realistic cost. Indeed, our current understanding of how to make very large quantum computers is often likened to classical computers before VLSI techniques were developed. The successful technologies enabling practical approaches to building large computers are likely yet to be

discovered, but architectures such as the one we have presented and the defect-tolerant, communication-oriented design principles we have used are expected to provide the guiding context for these new technologies.

### Acknowledgments

This work was supported by NSF, with partial support by MEXT and NICT. We acknowledge the support of the Australian Research Council, the Australian Government, and the US National Security Agency (NSA) and the Army Research Office (ARO) under contract number W911NF-08-1-0527. The authors thank Shinichi Koseki for fabricating and photographing the test structure and Shota Nagayama for help with the figures. We thank Jim Harrington, Robert Raussendorf, Ray Beausoliel, Kae Nemoto, Bill Munro, and the QIS groups at HP Labs and NII, for many useful technical discussions. We also would like to thank Skype, Ltd. for providing the classical networking software that enabled the tri-continental writing of this manuscript.

### References

1. D.P. DiVincenzo. The physical implementation of quantum computation. *Fortschritte der Physik*, 48(9-11):771–783, 2000.
2. David P. DiVincenzo. Quantum Computation. *Science*, 270(5234):255–261, 1995.
3. Timothy P. Spiller, William J. Munro, Sean D. Barrett, and Pieter Kok. An introduction to quantum information processing: applications and realisations. *Contemporary Physics*, 46:406, 2005.
4. Rodney Van Meter and Mark Oskin. Architectural implications of quantum computing technologies. *ACM Journal of Emerging Technologies in Computing Systems*, 2(1):31–63, January 2006.
5. Dorit Aharonov and Michael Ben-Or. Fault-tolerant quantum computation with constant error rate. <http://arXiv.org/quant-ph/9906129>, June 1999. extended version of STOC 1997 paper.
6. Todd Brun, Igor Devetak, and Min-Hsiu Hsieh. Correcting quantum errors with entanglement. *Science*, 314(5798):436–439, 2006.
7. Dave Bacon and Andrea Casaccino. Quantum error correcting subsystem codes from two classical linear codes. [quant-ph/0610088](http://arXiv.org/quant-ph/0610088), October 2006.
8. D. Bacon. Operator quantum error-correcting subsystems for self-correcting quantum memories. *Physical Review A*, 73(1):12340, 2006.
9. D.J.C. MacKay, G. Mitchison, and P.L. McFadden. Sparse-graph codes for quantum error correction. *IEEE Transactions on Information Theory*, 50(10):2315, 2004.
10. Andrew M. Steane. Overhead and noise threshold of fault-tolerant quantum error correction. *Physical Review A*, 68:042322, 2003.
11. Andrew M. Steane. Quantum computer architecture for fast entropy extraction. *Quantum Information and Computation*, 2(4):297–306, 2002. <http://arxiv.org/quant-ph/0203047>.
12. Simon J. Devitt, Austin G. Fowler, and Lloyd C. Hollenberg. Simulations of Shor’s algorithm with implications to scaling and quantum error correction. <http://arXiv.org/quant-ph/0408081>, August 2004.
13. Dean Copley, 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*, pages 65–74, 2003.

14. T. Szkopek, P. O. Boykin, H. Fan, V. P. Roychowdhury, E. Yablonovitch, G. Simms, M. Gyure, and B. Fong. Threshold error penalty for fault-tolerant quantum computation with nearest neighbor communication. *IEEE Trans. on Nanotech.*, 5:42, 2006.
15. Darshan D. Thaker, Tzvetan Metodi, Andrew Cross, Isaac Chuang, and Frederic T. Chong. CQLA: Matching density to exploitable parallelism in quantum computing. In *Computer Architecture News, Proc. 33rd Annual International Symposium on Computer Architecture*. ACM, June 2006.
16. Mark G. Whitney, Nemanja Isailovic, Yatish Patel, and John Kubiawicz. A fault tolerant, area efficient architecture for Shor's factoring algorithm. In *Proc. 36th Annual International Symposium on Computer Architecture*, June 2009.
17. E. Collin, G. Ithier, A. Aassime, P. Joyez, D. Vion, and D. Esteve. NMR-like control of a quantum bit superconducting circuit. *Physical Review Letters*, 93:157005, October 2004.
18. Lieven M.K. Vandersypen and Isaac Chuang. NMR techniques for quantum computation and control. *Rev. Modern Phys.*, 76:1037, 2004.
19. D. A. Lidar, I. L. Chuang, and K. B. Whaley. Decoherence-free subspaces for quantum computation. *Physical Review Letters*, 81(12):2594–2597, September 1998.
20. Daniel A. Lidar and K. Birgitta Whaley. *Irreversible Quantum Dynamics*, volume 622 of *Lecture Notes in Physics*, chapter Decoherence-Free Subspaces and Subsystems. Springer, 2003.
21. W. Dür and H.J. Briegel. Entanglement purification and quantum error correction. *Rep. Prog. Phys.*, 70:1381–1424, 2007.
22. Z. W. E. Evans, A. M. Stephens, J. H. Cole, and L. C. L. Hollenberg. Error correction optimisation in the presence of X/Z asymmetry, 2007.
23. Austin G. Fowler, Charles D. Hill, and Lloyd C. L. Hollenberg. Quantum error correction on linear nearest neighbor qubit arrays. *Physical Review A*, 69:042314, 2004.
24. A. Yu. Kitaev. Quantum computations: algorithms and error correction. *Russian Math. Surveys*, 52(6):1191–1249, 1997.
25. C.H. Bennett, G. Brassard, S. Popescu, B. Schumacher, J.A. Smolin, and W.K. Wootters. Purification of noisy entanglement and faithful teleportation via noisy channels. *Physical Review Letters*, 76(5):722–725, 1996.
26. J.I. Cirac, A. Ekert, S.F. Huelga, and C. Macchiavello. Distributed quantum computation over noisy channels. *Physical Review A*, 59:4249, 1999.
27. J. Dehaene, M. Van den Nest, B. De Moor, and F. Verstraete. Local permutations of products of Bell states and entanglement distillation. *Physical Review A*, 67(2):22310, 2003.
28. C. Kruszynska, A. Miyake, H.J. Briegel, and W. Dür. Entanglement purification protocols for all graph states. *Physical Review A*, 74(5):52316, 2006.
29. E.N. Maneva and J.A. Smolin. Improved two-party and multi-party purification protocols. *Contemporary Mathematics Series*, 305:203–212, 2000.
30. Rodney Van Meter, Thaddeus D. Ladd, W. J. Munro, and Kae Nemoto. System design for a long-line quantum repeater. *IEEE/ACM Transactions on Networking*, 17(3):1002–1013, June 2009.
31. E. Knill, R. Laflamme, R. Martinez, and C. Negrevergne. Benchmarking quantum computers: the five-qubit error correcting code. *Physical Review Letters*, 86(25):5811–5814, June 2001.
32. J. Chiaverini, D. Leibfried, T. Schaetz, M. D. Barrett, R. B. Blakestad, J. Britton, W. M. Itano, J. D. Jost, E. Knill, C. Langer, R. Ozeri, and D. J. Wineland. Realization of quantum error correction. *Nature*, 432:602–605, 2004.
33. T.B. Pittman, B.C. Jacobs, and J.D. Franson. Demonstration of quantum error correction using linear optics. *Physical Review A*, page 052332, May 2005.
34. Eric Dennis, Alexei Kitaev, Andrew Landahl, and John Preskill. Topological quantum memory. *J. Math. Phys.*, 43:4452–4505, 2002.
35. Robert Raussendorf, Jim Harrington, and Kovid Goyal. Topological fault-tolerance in cluster

- state quantum computation. *New Journal of Physics*, 9:199, 2007.
36. A.Y. Kitaev. Fault-tolerant quantum computation by anyons. *Annals of Physics*, 303(1):2–30, 2003.
  37. M.H. Freedman, A. Kitaev, M.J. Larsen, and Z. Wang. Topological quantum computation. *American Mathematical Society*, 40(1):31–38, October 2002.
  38. Robert Raussendorf and Jim Harrington. Fault-tolerant quantum computation with high threshold in two dimensions. *Physical Review Letters*, 98:190504, 2007.
  39. Simon J. Devitt, Austin G. Fowler, Todd Tilma, William J. Munro, and Kae Nemoto. Classical processing requirements for a topological quantum computing system, 2009.
  40. S. J. Devitt, A. G. Fowler, A. M. Stephens, A. D. Greentree, L. C. L. Hollenberg, W. J. Munro, and Kae Nemoto. Architectural design for a topological cluster state quantum computer. *arXiv:0808.1782*, 2008.
  41. D. P. DiVincenzo. Fault tolerant architectures for superconducting qubits. *arXiv:0905.4839*, 2009.
  42. René Stock and Daniel F. V. James. Scalable, high-speed measurement-based quantum computer using trapped ions. *Physical Review Letters*, 102(17):170501, 2009.
  43. John L. Hennessy and David A. Patterson. *Computer Architecture: A Quantitative Approach*. Morgan Kaufman, 4th edition, 2006.
  44. William James Dally and Brian Towles. *Principles and Practices of Interconnection Networks*. Elsevier, 2004.
  45. Rodney Van Meter, Kohei M. Itoh, and Thaddeus D. Ladd. Architecture-dependent execution time of Shor’s algorithm. In *Proc. Int. Symp. on Mesoscopic Superconductivity and Spintronics (MS+S2006)*, February 2006.
  46. Peter W. Shor. Algorithms for quantum computation: Discrete logarithms and factoring. In *Proc. 35th Symposium on Foundations of Computer Science*, pages 124–134, Los Alamitos, CA, 1994. IEEE Computer Society Press.
  47. Rodney Doyle Van Meter III. *Architecture of a Quantum Multicomputer Optimized for Shor’s Factoring Algorithm*. PhD thesis, Keio University, 2006. available as arXiv:quant-ph/0607065.
  48. Jeffrey Yepez. Type-II quantum computers. *International Journal of Modern Physics C*, 12(9):1273–1284, 2001.
  49. Thomas M. Stace, Sean D. Barrett, and Andrew C. Doherty. Thresholds for topological codes in the presence of loss. *Physical Review Letters*, 102(20):200501, 2009.
  50. H.-J. Briegel, W. Dür, J.I. Cirac, and P. Zoller. Quantum repeaters: the role of imperfect local operations in quantum communication. *Physical Review Letters*, 81:5932–5935, 1998.
  51. A.G. Fowler, A.M. Stephens, and P. Groszkowski. High threshold universal quantum computation on the surface code. *Arxiv preprint arXiv:0803.0272*, 2008.
  52. Austin G. Fowler and Kovid Goyal. Topological cluster state quantum computing, 2008.
  53. D.S. Wang, A.G. Fowler, A.M. Stephens, and L.C.L. Hollenberg. Threshold error rates for the toric and surface codes. *Arxiv preprint arXiv:0905.0531*, 2009.
  54. J.P. Reithmaier, G. Sek, A. Löffler, C. Hofmann, S. Kuhn, S. Reitzenstein, L.V. Keldysh, V.D. Kulakovskii, T.L. Reinecke, and A. Forchel. Strong coupling in a single quantum dot–semiconductor microcavity system. *Nature*, 432(7014):197–200, 2004.
  55. J. Berezovsky, M. H. Mikkelsen, N. G. Stoltz, L. A. Coldren, and D. D. Awschalom. Picosecond coherent optical manipulation of a single electron spin in a quantum dot. *Science*, 320:349, 2008.
  56. D. Press, T. D. Ladd, B. Y. Zhang, and Y. Yamamoto. Complete quantum control of a single quantum dot spin using ultrafast optical pulses. *Nature*, 456:218–221, 2008.
  57. C. Kistner, T. Heindel, C. Schneider, A. Rahimi-Iman, S. Reitzenstein, S. Hofling, and A. Forchel. Demonstration of strong coupling via electro-optical tuning in high-quality QD-micropillar systems. *Optics Express*, 16(19):15006, 2008.

58. I. Fushman, D. Englund, A. Faraon, N. Stoltz, P. Petroff, and J. Vuckovic. Controlled phase shifts with a single quantum dot. *Science*, 320(5877):769–772, 2008.
59. C. Schneider, M. Strauß, T. Sünner, A. Huggenberger, D. Wiener, S. Reitzenstein, M. Kamp, S. Höfling, and A. Forchel. Lithographic alignment to site-controlled quantum dots for device integration. *Appl. Phys. Lett.*, 92(18):183101, 2008.
60. A.M. Tyryshkin, S. A. Lyon, A. V. Astashkin, and A. M. Raitisimring. Electron spin-relaxation times of phosphorous donors in silicon. *Phys. Rev. B*, 68:193207, 2003.
61. A. Yang, M. Steger, D. Karaiskaj, M. L. W. Thewalt, M. Cardona, K. M. Itoh, H. Riemann, N. V. Abrosimov, M. F. Churbanov, A. V. Gusev, A. D. Bulanov, A. K. Kaliteevskii, O. N. Godisov, P. Becker, H.-J. Pohl, J. W. Ager III, and E. E. Haller. Optical detection and ionization of donors in specific electronic and nuclear spin states. *Phys. Rev. Lett.*, 97:227401, 2006.
62. A. Pawlis, M. Panfilova, D. J. As, K. Lischka, K. Sanaka, T. D. Ladd, and Y. Yamamoto. Lasing of donor-bound excitons in ZnSe microdisks. *Phys. Rev. B*, 77:153304, 2008.
63. K. Sanaka, A. Pawlis, T. D. Ladd, K. Lischka, and Y. Yamamoto. Indistinguishable photons from independent semiconductor nanostructures. *Phys. Rev. Lett.*, 2009. in press.
64. M. V. G. Dutt, L. Childress, L. Jiang, E. Togan, J. Maze, F. Jelezko, A. S. Zibrov, P. R. Hemmer, and M. D. Lukin. Quantum register based on individual electronic and nuclear spin qubits in diamond. *Science*, 316(5829):1312–1316, 2007.
65. C. Santori, Ph. Tamarat, P. Neumann, J. Wrachtrup, D. Fattal, R. G. Beausoleil, J. Rabeau, P. Olivero, A. D. Greentree, S. Praver, F. Jelezko, and P. Hemmer. Coherent population trapping of single spins in diamond under optical excitation. *Phys. Rev. Lett.*, 97(24):247401, 2006.
66. Balasubramanian, G. et al. Ultralong spin coherence time in isotopically engineered diamond. *Nature Mater.*, 8:383–387, 2009.
67. S.M. Clark, K.M.C. Fu, T.D. Ladd, and Y. Yamamoto. Quantum computers based on electron spins controlled by ultra-fast, off-resonant, single optical pulses. *Physical Review Letters*, 99:040501, 2007.
68. T. D. Ladd, P. van Loock, K. Nemoto, W. J. Munro, and Y. Yamamoto. Hybrid quantum repeater based on dispersive CQED interactions between matter qubits and bright coherent light. *New J. Phys.*, 8:184, 2006.
69. C. Cabrillo, J. I. Cirac, P. García-Fernández, and P. Zoller. Creation of entangled states of distant atoms by interference. *Phys. Rev. A*, 59(2):1025, 1999.
70. L. Childress, J.M. Taylor, A.S. Sørensen, and M.D. Lukin. Fault-tolerant quantum repeaters with minimal physical resources and implementations based on single-photon emitters. *Physical Review A*, 72(5):52330, 2005.
71. E. Waks and J. Vuckovic. Dipole induced transparency in drop filter cavity-waveguide systems. *Phys. Rev. Lett.*, 96:153601, 2006.
72. A. Politi, M. J. Cryan, J. G. Rarity, S. Y. Yu, and J. L. O’Brien. Silica-on-silicon waveguide quantum circuits. *Science*, 320(5876):646, 2008.
73. H. Rokhsari and K. J. Vahala. Ultralow loss, high  $Q$ , four port resonant couplers for quantum optics and photonics. *Phys. Rev. Lett.*, 92(25):253905, Jun 2004.
74. Yurii Vlasov, William M. J. Green, and Fengnian Xia. High-throughput silicon nanophotonic wavelength-insensitive switch for on-chip optical networks. *Nature Photonics*, March 2008. doi:10.1038/nphoton.2008.31.
75. J. Kim et al. System design for large-scale ion trap quantum information processor. *Quantum Information and Computation*, 5(7):515–537, 2005.
76. S. M. Clark, K.-M. C. Fu, Q. Zhang, T D. Ladd, C. Stanley, and Y. Yamamoto. Ultrafast optical spin echo for electron spins in semiconductors. *Phys. Rev. Lett.*, 2009. in press.
77. 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, 2006.

28 *Van Meter, Ladd, Fowler and Yamamoto*

78. M. Atature, J. Dreiser, A. Badolato, and A. Imamoglu. Observation of Faraday rotation from a single confined spin. *Nat. Phys.*, 3:101, 2007.
79. I. Fushman, D. Englund, A. Faraon, N. Stoltz, P. Petroff, and J. Vuckovic. Controlled phase shifts with a single quantum dot. *Science*, 320(5877):769, 2008.
80. A. Imamoglu, D. D. Awschalom, G. Burkard, D.P Divincenzo, D. Loss, M. Shermin, and A. Small. Quantum information processing using quantum dot spins and cavity QED. *Phys. Rev. Lett.*, 83:4204, 1999.
81. Shi-Biao Zheng. Unconventional geometric quantum phase gates with a cavity QED system. *Phys. Rev. A*, 70(5):052320, Nov 2004.
82. T. P. Spiller, K. Nemoto, S. L. Braunstein, W. J. Munro, P. van Loock, and G. J Milburn. Quantum computation by communication. *New J. Phys.*, 8:30, 2006.
83. Thaddeus D. Ladd and Yoshihisa Yamamoto. in preparation.
84. L. M. Duan and H. J. Kimble. *Phys. Rev. Lett.*, 92:127902, 2004.
85. A. M. Stephens, Z. W. E. Evans, S. J. Devitt, A. D. Greentree, A. G. Fowler, W. J. Munro, J. L. O'Brien, K. Nemoto, and L. C. L. Hollenberg. Deterministic optical quantum computer using photonic modules. *Physical Review A*, 78(3), 2008.
86. P. van Loock, T. D. Ladd, K. Sanaka, F. Yamaguchi, K. Nemoto, W. J. Munro, and Y. Yamamoto. Hybrid quantum repeater using bright coherent light. *Phys. Rev. Lett.*, 96:240501, 2006.
87. S. D. Barrett, Pieter Kok, Kae Nemoto, R. G. Beausoleil, W. J. Munro, and T. P. Spiller. Symmetry analyzer for nondestructive bell-state detection using weak nonlinearities. *Phys. Rev. A*, 71(6):060302, Jun 2005.
88. W.J. Munro, K. Nemoto, and T.P. Spiller. Weak nonlinearities: a new route to optical quantum computation. *New Journal of Physics*, 7:137, May 2005.
89. Sean D. Barrett and Pieter Kok. Efficient high-fidelity quantum computation using matter qubits and linear optics. *Phys. Rev. A*, 71(6):060310, Jun 2005.
90. Yuan Liang Lim, Sean D. Barrett, Almut Beige, Pieter Kok, and Leong Chuan Kwek. Repeat-Until-Success quantum computing using stationary and flying qubits. *Physical Review Letters*, 95(3):30505, 2005.
91. P. van Loock, N. Lutkenhaus, W. J. Munro, and K. Nemoto. Quantum repeaters using coherent-state communication. *Physical Review A*, 78(6), 2008.
92. E. Peter, P. Senellart, D. Martrou, A. Lemaître, J. Hours, J. M. Gérard, and J. Bloch. Exciton-photon strong-coupling regime for a single quantum dot embedded in a microcavity. *Phys. Rev. Lett.*, 95:067401, 2005.
93. Shinichi Koseki, Bingyang Zhang, Kristiaan De Greve, and Yoshihisa Yamamoto. Monolithic integration of quantum dot containing microdisk microcavities coupled to air-suspended waveguides. *Applied Physics Letters*, 94:051110, February 2009.
94. Rodney Van Meter, W. J. Munro, Kae Nemoto, and Kohei M. Itoh. Arithmetic on a distributed-memory quantum multicomputer. *ACM Journal of Emerging Technologies in Computing Systems*, 3(4):17, January 2008.
95. Rodney Van Meter, Kae Nemoto, and William J. Munro. Communication links for distributed quantum computation. *IEEE Transactions on Computers*, 56(12):1643–1653, December 2007.
96. Stephen Y. Chou, Peter R. Krauss, and Preston J. Renstrom. Imprint lithography with 25-nanometer resolution. *Science*, 272(5258):85–87, 1996.
97. Linjie Li, Rafael R. Gattass, Erez Gershgoren, Hana Hwang, and John T. Fourkas. Achieving  $\lambda/20$  resolution by one-color initiation and deactivation of polymerization. *Science*, 324(5929):910–913, 2009.
98. Trisha L. Andrew, Hsin-Yu Tsai, and Rajesh Menon. Confining light to deep subwavelength dimensions to enable optical nanopatterning. *Science*, 324(5929):917–921, 2009.
99. Timothy F. Scott, Benjamin A. Kowalski, Amy C. Sullivan, Christopher N. Bowman, and

- Robert R. McLeod. Two-color single-photon photoinitiation and photoinhibition for subdiffraction photolithography. *Science*, 324(5929):913–917, 2009.
100. W. Dür, H.-J. Briegel, J. I. Cirac, and P. Zoller. Quantum repeaters based on entanglement purification. *Physical Review A*, 59(1):169–181, Jan 1999.
  101. Michael A. Nielsen and Isaac L. Chuang. *Quantum Computation and Quantum Information*. Cambridge University Press, 2000.
  102. Vlatko Vedral, Adriano Barenco, and Artur Ekert. Quantum networks for elementary arithmetic operations. *Phys. Rev. A*, 54:147–153, 1996. <http://arXiv.org/quant-ph/9511018>.
  103. Rodney Van Meter and Kohei M. Itoh. Fast quantum modular exponentiation. *Physical Review A*, 71(5):052320, May 2005.
  104. Thomas G. Draper, Samuel A. Kutin, Eric M. Rains, and Krysta M. Svore. A logarithmic-depth quantum carry-lookahead adder. *Quantum Information and Computation*, 6(4&5):351–369, July 2006.

# A Layered Architecture for Quantum Computing Using Quantum Dots

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 Computer 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

\* Corresponding author:

E-mail: ncjones@stanford.edu

**Abstract.** We address the challenge of designing a quantum computer architecture with a layered framework that is modular and facilitates fault-tolerance. The framework is flexible and could be used for analysis and comparison of differing quantum computer designs. Using this framework, we develop a complete, layered architecture for quantum computing with optically controlled quantum dots, showing how a myriad of technologies must operate synchronously to achieve fault-tolerance. Our design deliberately takes advantage of the large possibilities for integration afforded by semiconductor fabrication. Quantum information is stored in the electron spin states of a charged quantum dot controlled by ultrafast optical pulses. Optical control makes this system very fast, scalable to large problem sizes, and extensible to quantum communication or distributed architectures. The design of this quantum computer centers on error correction in the form of a topological surface code, which requires only local and nearest-neighbor gates. We analyze several important issues of the surface code that are relevant to an architecture, such as resource accounting and the use of Pauli frames. Furthermore, we investigate the performance of this system and find that Shor's factoring algorithm for a 2048-bit number can be executed in approximately one week.

PACS numbers: 03.67.Pp, 03.67.Lx, 85.35.Be, 73.21.La, 85.40.Hp

## Contents

<b>1</b>	<b>Introduction to the Layered Architecture</b>	<b>3</b>
1.1	Prior Work on Quantum Computer Architecture . . . . .	3
1.2	Layered Framework . . . . .	4
1.3	What is a qubit? . . . . .	5
1.4	Two Ways to Protect Quantum Information . . . . .	5
1.5	Communication Protocols . . . . .	7
1.6	Interaction between Layers . . . . .	7

<i>CONTENTS</i>	2
<b>2 Layer 1: Physical</b>	<b>7</b>
2.1 Components of the Physical Layer . . . . .	7
2.1.1 Electron Spin within a Quantum Dot . . . . .	8
2.1.2 Planar DBR Microcavity . . . . .	9
2.1.3 Ultrafast Optical Pulses for Spin-State Rotation . . . . .	9
2.1.4 Spin Entangling Operation . . . . .	9
2.1.5 Optical Patterns for Control Pulses . . . . .	11
2.1.6 Quantum Non-Demolition (QND) Measurement . . . . .	13
2.1.7 Detector Array . . . . .	14
2.1.8 Static Decoherence (Memory Errors) . . . . .	14
2.1.9 Dynamic Coherent and Incoherent Errors . . . . .	15
2.2 Layer 1 Performance . . . . .	15
2.2.1 Manipulating the Spin-basis Bloch Sphere . . . . .	15
2.2.2 Entangling Operation . . . . .	16
2.2.3 MEMS Micromirror Switching . . . . .	17
<b>3 Layer 2: Virtualization</b>	<b>17</b>
3.1 Components of the Virtualization Layer . . . . .	17
3.1.1 Virtual Qubit . . . . .	17
3.1.2 Virtual Gates . . . . .	18
3.2 Layer 2 Performance . . . . .	18
3.2.1 Virtual Qubit . . . . .	18
3.2.2 Virtual Gates . . . . .	21
<b>4 Layer 3: Quantum Error Correction</b>	<b>21</b>
4.1 Components of the QEC Layer . . . . .	22
4.1.1 Architecture and the Surface Code . . . . .	23
4.1.2 Pauli Frames . . . . .	23
4.1.3 Measurement and Detector Arrays . . . . .	24
4.2 Layer 3 Performance . . . . .	24
4.2.1 Size of the Surface Code . . . . .	24
4.2.2 Surface Code Operations in Time . . . . .	26
4.2.3 Local Error Correction Processing . . . . .	27
4.2.4 Pauli Frames in Action . . . . .	27
<b>5 Layer 4: Logical</b>	<b>28</b>
5.1 Functions of the Logical Layer . . . . .	28
5.2 Layer 4 Performance . . . . .	29
5.2.1 Singular State Distillation . . . . .	29
5.2.2 Construction of a Toffoli Gate . . . . .	31
5.2.3 Summary of Logical Layer Performance . . . . .	31
<b>6 Application Layer</b>	<b>31</b>
<b>7 Timing Considerations</b>	<b>32</b>
<b>8 Discussion</b>	<b>34</b>

## 1. Introduction to the Layered Architecture

A computer architecture defines and organizes the components of a system, their roles, and the interfaces between them. In computer systems, the architecture of a system determines its performance, the difficulty of implementation, and its flexibility. A good architecture exposes the strengths of its underlying technologies while avoiding unnecessary dependence on a specific technology, allowing independent evolution over time, and occasionally wholesale replacement of components or subsystems. Developing a flexible framework is particularly important for the nascent field of quantum computing, where relatively little work on architecture has been performed. This problem is important since an architecture provides structure, not only for the quantum computer itself but also for the designers — organizing the system design can also serve to organize the conceptual and logistical problems of engineering a computer.

Here, we propose a layered architecture for quantum computing which is both *modular* and *fault-tolerant*. The objective is to develop a framework for building up a quantum computer from individual components, while also providing a means to compare different approaches to quantum computing, such as nitrogen-vacancy centers in diamond, quantum dots, trapped ions, or atoms in optical lattices [1]. This architecture has many universal aspects applicable to different physical hardware, but to make this discussion concrete, we introduce a new quantum computer architecture based on **Quantum Dots with Optically-controlled Spins**, or QuDOS. The organizing principles of the architecture are explained as this specific quantum computer implementation is developed step-by-step.

### 1.1. Prior Work on Quantum Computer Architecture

Many different quantum computing technologies are under experimental investigation [1]. Since DiVincenzo introduced his fundamental criteria for a viable quantum computing technology [2] and Steane emphasized the difficulty of designing systems capable of running quantum error correction (QEC) adequately [3,4], several groups of researchers have outlined various additional taxonomies addressing the architectural needs of large-scale systems [5,6]. For many technologies, small-scale interconnects have been proposed, 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 researchers. In particular, the issue of heterogeneity in systems has received relatively little attention.

Kielpinski *et al.* proposed a scalable ion trap technology utilizing separate memory and computing areas [7]. 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 [8,9]. Oskin *et al.* [10], building on the Kane solid-state NMR technology [11], 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. Duan and Monroe proposed the use of photonic qubits to distribute entanglement between ions located in distant traps [12], and such photonic channels could be utilized to realize a modular, scalable distributed

quantum computer [13]. Fowler *et al.* [14] 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.* [15, 16] have investigated automated layout and optimization of circuit designs specifically for ion trap architectures, and Isailovic *et al.* [17, 18] have studied interconnection and data throughput issues in similar ion trap systems.

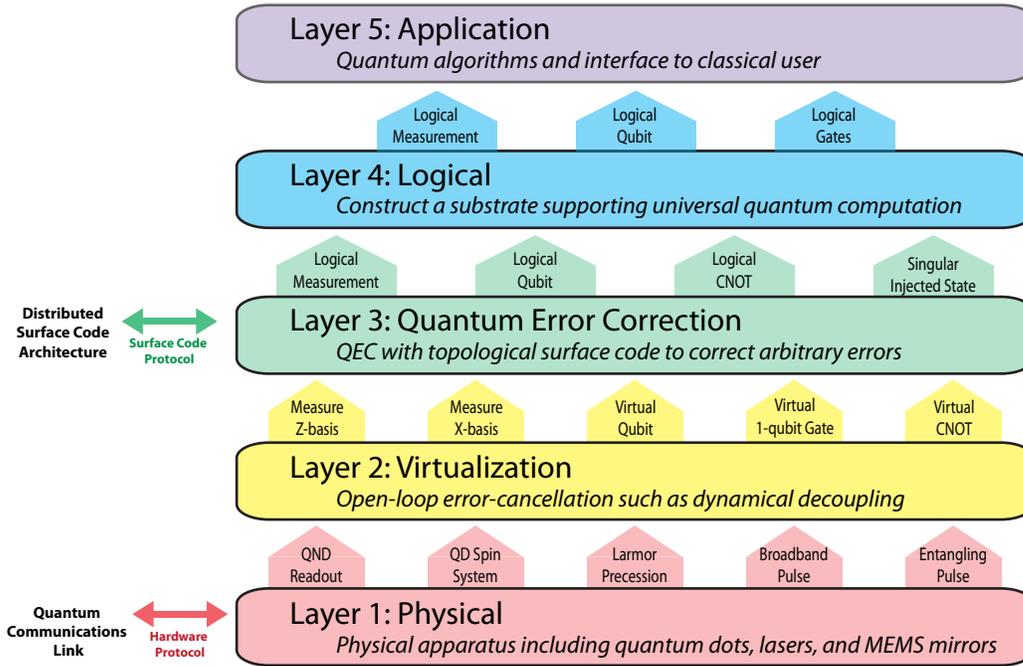
With the recent advances in the operation of the topological codes and its desirable characteristics of a high *practical* threshold and need for only nearest-neighbor interactions, research effort has shifted toward architectures capable of building and maintaining large two- and three-dimensional cluster states [19, 20].

The abstract framework of a quantum multicomputer [21] recognizes that large-scale systems demand heterogeneous interconnects; in most quantum computing technologies, it may not be possible to build monolithic systems that contain, couple, and control billions of physical qubits. This architectural framework was extended in a recent paper designed around nanophotonic coupling of electron spin quantum dots that explicitly uses *multiple* levels of interconnect with varying coupling fidelities (hence, purification requirements), as well as the ability to operate with a very low yield of functional devices [22]. 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 architecture proposed in this paper.

### 1.2. Layered Framework

A good architecture must have a simple structure while also efficiently managing the complex array of resources in a quantum computer. Our architecture 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. The utility of this scheme is that the many functions in a quantum computer are organized in a manner that aids understanding for the designer and translates directly into an effective control structure for the device itself. 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 which organizes the operations in the architecture. Figure 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, 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 underpinning the quantum computer. The layers between (Virtualization, 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* qubits and quantum gates at the Application layer.



**Figure 1.** Layered control stack which forms the framework of a quantum computer architecture. Vertical arrows indicate services provided to a higher layer. Arrows on the left margin indicate that communication protocols are necessary to connect multiple quantum information devices, but this topic is outside the scope of this work.

### 1.3. What is a qubit?

The fundamental unit of information in a quantum computer is the qubit. For our purposes, we reserve the terminology “qubit” for an isolated system where the quantum information is closed under  $SU(2)$  algebra [23]. A simple example is a two-level system (TLS), such as the spin of an electron. The reason for this distinction is to reserve “qubit” for its meaning as an information unit, not a physical system. As we will show by example in section 2.1.1, the underlying physical system may be more complex than a TLS, but these details are hidden by layers of abstraction in the architecture. The “qubit” first appears as an output of Layer 2 (Virtualization, section 3), where physical processes are organized into quantum information in the form of “virtual qubits.” Layer 3 (Quantum Error Correction, section 4) constructs a “logical qubit” from many virtual qubits and gates. These objects and the processes which create them are explained in detail in subsequent sections.

### 1.4. Two Ways to Protect Quantum Information

Quantum information is fragile, so the most important role of a fault-tolerant quantum computing architecture is to protect quantum information from errors caused by

both coupling to the environment and imperfect control operations. There are two fundamental approaches to protecting qubits used in the layered architecture. The first technique addresses systematic errors, which are correlated in time. When considering decoherence of a qubit, the environment exerts an unknown coupling to the qubit, but perhaps the noise power spectral density of the bath has a characteristic coherence time longer than the timescales of control operations. This situation can be addressed by dynamical decoupling (DD) [24, 25], which is a class of control techniques for reducing qubit decoherence caused by an environment. Similarly, control pulses may have a repeatable bias (such as laser intensity fluctuations), so that the same error is consistent between pulses at different times. Much like DD, there are sequences of pulses known as compensation sequences [26, 27] which reduce control errors by having multiple faulty pulses combine to create a more accurate quantum gate. This collection of techniques resides in Layer 2, the Virtualization layer (see section 3). The purpose of Layer 2 is to take raw physical processes and shape them into the abstract components of quantum information — qubits and quantum gates — which is why the qubit first appears in Layer 2.

The second important method for protecting quantum information is quantum error correction (QEC) [28]. Much like its classical analogue, QEC encodes quantum information in an error-correcting code, which is characterized by the ability to identify and correct arbitrary errors in the fundamental qubits and gates (provided their probability of occurrence is below a certain threshold). The whole of Layer 3 is devoted to QEC, which in this investigation is a topological surface code [29]. In general, other QEC schemes can be incorporated into this architecture. Errors manifest as a “syndrome” found by projective measurement operations in a “syndrome extraction” circuit in the quantum computer. QEC fails when the most likely pattern of errors corresponding to a syndrome is not correct, which happens when error rates are too high (above threshold). The hallmark of QEC is its ability to correct arbitrary errors.

The distinction between Layers 2 and 3 is subtle but important. Qualitatively, it would seem that Layer 2 is open-loop control since the sequence of control operations does not depend on the state of the system, while Layer 3 quantum error correction incorporates feedback by measuring the system and changing future operations conditioned on the measurements. However, the methodology of QEC has advanced so that this is no longer accurate; in particular, the accumulated errors can be handled by Pauli frames (see sections 4.1.2 and 4.2.4). In this manner, “error correction” consists of post-processing measurement results, so feedback into control operations does not occur. Still, while not traditional closed loop control, the quantum measurements do alter the quantum data; in particular they project drifts in the continuous Hilbert space of the virtual qubits into discrete substates which may be analyzed via digital error correction techniques. Another possibility is separating error mitigation techniques by local and non-local control operations. Dynamical decoupling typically uses only local gates, but there is a related concept known as the decoherence-free subspace (DFS) [30] which requires non-local gates (*e.g.* coupling multiple electron spins). A DFS encodes a qubit into a system with many more degrees of freedom; the particular encoding exploits a symmetry in the system so that the subspace spanned by the qubit is invariant to some non-unitary coupling to the environment (decoherence). This behavior is reminiscent of QEC, but there is a crucial distinction that firmly separates DD, DFS, and compensation sequences in Layer 2 from QEC in Layer 3. Layer 2 techniques do not extract information about the state of the quantum computer, whereas Layer 3 does. The information gathered by Layer

3 is not the state of the quantum information being protected, but rather the likely pattern of errors which have occurred. Therefore, Layer 2 does not monitor the state of the system and never uses projective measurement. Layer 3 monitors the system for errors, and accomplishes this with measurement. Nevertheless, Layers 2 and 3 are not redundant; the importance of each, along with their synergy, is discussed in sections 3 and 4.

### *1.5. Communication Protocols*

Just as modern digital computers are frequently networked together, quantum computers may need to share a quantum information connection. To communicate quantum information between two devices, as in quantum repeaters [31] or in a distributed architecture [22], an appropriate communication protocol for Layer 1 (Physical) must be devised. These two quantum computers could be of wholly different technologies (say ion trap vs. quantum dot) if a practicable protocol exists. Moreover, a distributed surface code architecture [22] would also require a communication protocol [32] in Layer 3. The location of these protocols in the layered framework is noted in Figure 1, but this topic is considered outside the scope of the present work.

### *1.6. Interaction between Layers*

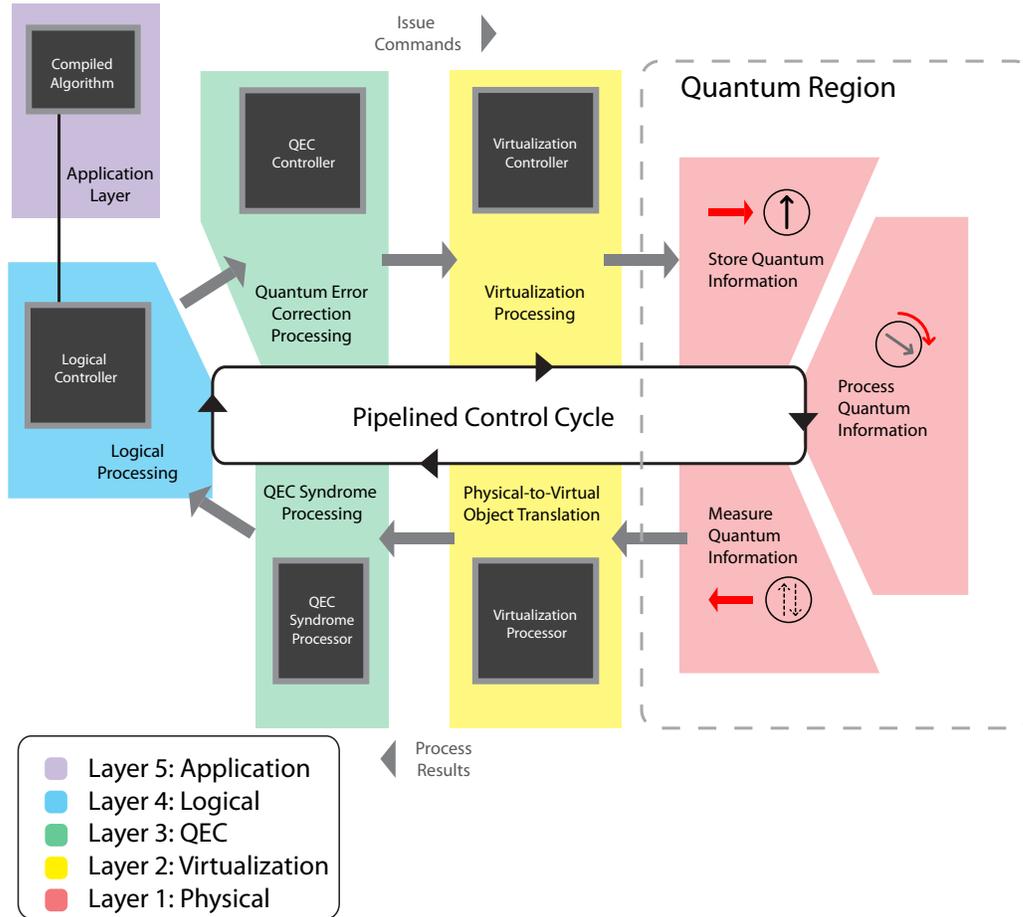
For the quantum computer to function efficiently, each layer must issue instructions to layers below in a tightly defined sequence. However, a robust system must also be able to handle errors caused by faulty devices. To satisfy both criteria, a control loop must handle operations at all layers simultaneously while also processing syndrome measurement to correct errors which occur. A prototype for this control loop is shown in Figure 2.

The primary control cycle defines the behavior of the quantum computer in this architecture since all operations must interact with this loop. As discussed later, timing is critically important, so this cycle does not simply issue a single command and wait for the result before proceeding — pipelining is essential [33]. Moreover, Figure 2 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 to the underlying design of the quantum computer.

## **2. Layer 1: Physical**

The physical layer is the foundation of the quantum computer. All truly quantum effects happen here, with higher layers building complicated operations from sequences of processes performed at the physical layer. As a result, the physical layer exists solely to provide services to layers above, and no decision- or branching-based controls run here, as occurs in the upper layers. Implementing a quantum computer architecture begins at Layer 1, where basic hardware for storing and manipulating quantum information is constructed. We illustrate this process with a quantum computer based on the optical control of charged quantum dots known as QuDOS.

### *2.1. Components of the Physical Layer*



**Figure 2.** Primary control cycle of the quantum computer. Whereas the control stack in Figure 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.

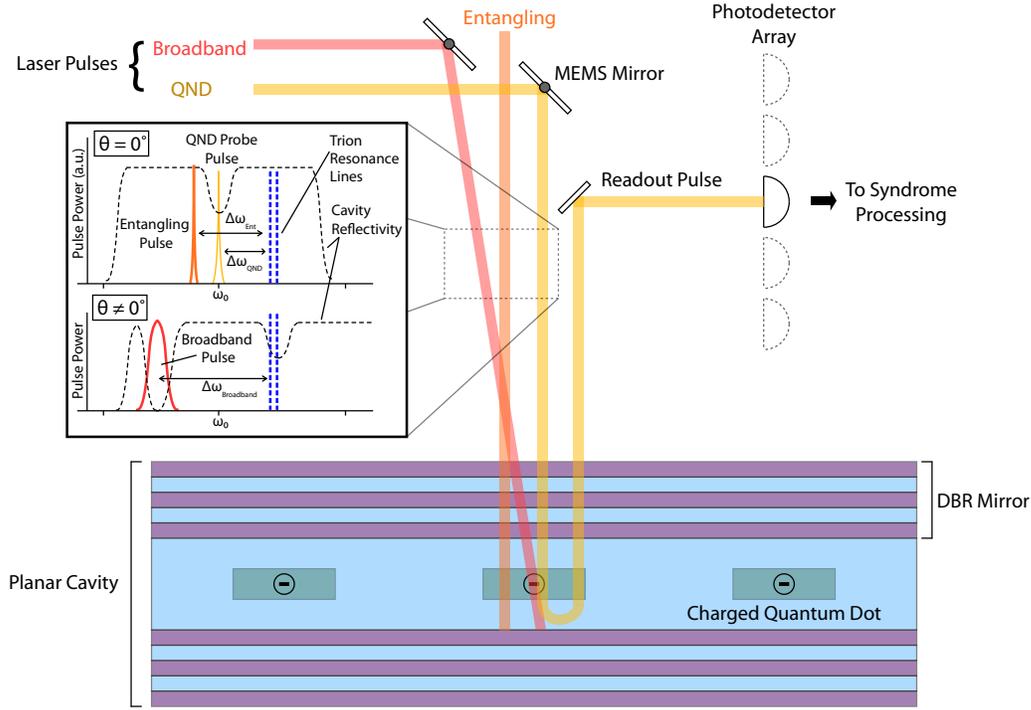
*2.1.1. Electron Spin within a Quantum Dot* A quantum computer must have the ability to store information between processing steps. The information carrier in QuDOS is the spin of an electron bound within an InGaAs self-assembled quantum dot (QD) surrounded by GaAs substrate [34–39]. These QDs can be 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 [40], 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 2.1.3. Second, there is continuous phase evolution between the two spin states, which in conjunction with optical pulses provides complete unitary control of the electron spin vector.

*2.1.2. Planar DBR Microcavity* Accessing the quantum properties of a single electron spin system requires an enhanced interaction with light, and so an optical microcavity is necessary. To facilitate the two-dimensional array of the surface code detailed in Layer 3, this microcavity must be planar in design, and so the cavity is constructed from two distributed Bragg reflector (DBR) mirrors stacked vertically with a  $\lambda/2$  cavity layer in between. This cavity is grown by molecular beam epitaxy (MBE). The QDs are embedded at the center of this cavity to maximize interaction with antinodes of the cavity field modes. Figure 3 illustrates quantum dots arranged at the center of a planar cavity. Using MBE, high-quality ( $Q > 10^5$ ) microcavities can be grown with alternating layers of GaAs/AlAs [41].

*2.1.3. Ultrafast Optical Pulses for Spin-State Rotation* The ability to perform fast manipulations of the quantum states stored in a quantum computer is essential for performing operations faster than decoherence processes can corrupt them, as well as for ensuring a fast overall algorithm execution time [6]. In QuDOS, ultrafast optical pulses centered 900–950 nm rotate the spin vector of an electron within a QD [42, 43]. By virtue of being short in duration, these pulses are broad in frequency, facilitating stimulated Raman transitions between the spin levels through excited-state trion levels. Therefore, the complete dynamics of the state rotation depends on a four-level system (consisting of the two metastable spin ground states and two excited trion states). Other control pulses in QuDOS (see sections 2.1.4 and 2.1.6) require a high-Q microcavity which has a narrow transmission window at the cavity resonance; the bandwidth of the broadband pulses is significantly larger than the transmission bandwidth of the cavity resonance. As a result, the broadband pulse cannot be sent directly into the microcavity. This problem is circumvented by sending the broadband pulses at angled (rather than normal) incidence. The cavity response is shifted to higher frequencies, so that a red-detuned pulse can enter the cavity at the first minimum in the cavity reflectivity as a function of frequency. Alternatively, one could send red-detuned pulses at normal incidence, sacrificing the majority of each pulse which is reflected; this approach is only viable if significantly more optical power is available. Figure 3 shows the three laser pulses used in QuDOS, as well as their power spectrums.

*2.1.4. Spin Entangling Operation* The construction of a practical, scalable two-qubit gate in a quantum dot architecture remains the most challenging element of the hardware. In quantum dots with transverse confinement provided by electrostatic gates, electronic manipulation of the electron wavefunction allows control over the exchange interaction, providing fast ( $\sim 100$  ps) quantum gates, as proposed some time ago [45] and demonstrated in numerous experiments [46]. Employing such gates for a hybrid system with both optical and electrical control is certainly possible, but requires further development of the optical control of electrically defined quantum dots [47]. Entanglement of directly tunnel-coupled vertically stacked InAs quantum dots has also been demonstrated [48], but the scalability of this coupling mechanism



**Figure 3.** Schematic diagram showing the primary components of the Physical Layer. The inset image shows the spectrum for the various laser pulses that implement different quantum operations. The quantum non-demolition (QND) measurement and entangling pulses are modulated continuous-wave laser pulses, which are narrow in frequency bandwidth; these pulses are sent at normal incidence. The broadband pulses which rotate the electron spin state are angled relative to normal incidence, which shifts the cavity reflectivity response to higher frequencies. This enables the red-detuned pulse to enter the cavity at the first dip in the reflectivity. Each of the different laser pulses has a detuning relative to the trion (excited state) resonance frequency. The entangling pulse is also detuned from cavity resonance in a manner prescribed in Ref. [44].

is uncertain. An exotic but promising possibility includes optically inducing longer-range, exciton-mediated exchange interactions [49, 50].

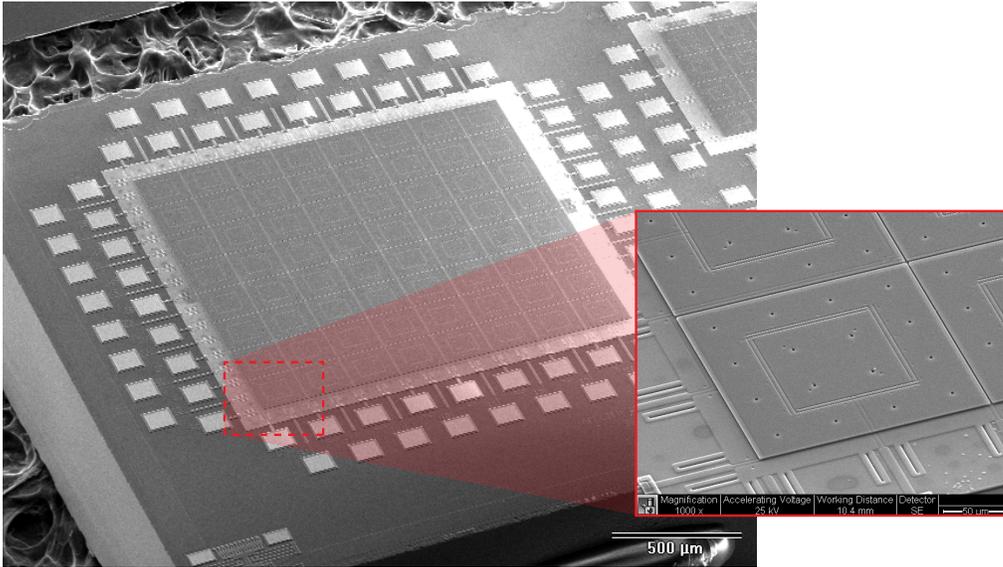
A fast, all-optically controlled two-qubit gate would certainly be attractive, and early proposals [51] identified the importance of employing the nonlinearities of cavity QED. Ref. [51] suggests the application of two lasers for both single-qubit and two-qubit control; more recent developments have indicated that both single-qubit gates [42, 52] and two-qubit gates [44] can be accomplished using only a single optical pulse. However, the demands on the optical microcavity system are challenging. 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  $V$ . Large values of this can be achieved via cavities with strong transverse confinement, such as the microdisk cavities proposed in Ref. [51]. This arrangement

poses challenges for scalability. An obvious modification is to couple these cavities with waveguides; an architecture employing this approach was discussed in Ref. [22], and in this case substantial additional physical resources are needed to mitigate optical losses at the cavity-waveguide interfaces. For the present architecture, 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 [34]. While existing cooperativity factors achieved this way are not estimated to be high enough to produce quantum gates which meet the fault-tolerant threshold alone, advanced control techniques and multi-spin encodings (such as for “virtual qubits”; see sections 3.1.1 and 3.2.1) may enable this technology to function with acceptable error rates.

Whether a microplanar microcavity arrangement will offer sufficient nonlinearity for an all-optically controlled dot architecture or whether electronically controlled gates will be employed will depend on forthcoming experimental developments. However, for the purposes of the present architecture, both gates are short range, a constraint handled by the surface code quantum error correction we employ (see section 4) which demands only nearest-neighbor interactions. Short range gates can severely limit the efficacy of other quantum error correction schemes [53]. Long range couplings, for example to form bridges over optically inactive regions of a single chip or for chip-to-chip connections, will likely employ a variety of different quantum optical techniques which sacrifice speed for tolerance to optical loss; for a discussion, see Refs. [22] and [54]. Incorporating such long-distance links into the present architecture must occur at both the Physical and QEC layers; the inclusion of such interconnections into QuDOS is the subject of future work.

*2.1.5. Optical Patterns for Control Pulses* Parallelism in operations is essential for QuDOS to function efficiently. The quantum operations are driven by laser pulses sent into the planar microcavity (see Figure 3). However, one cannot use a single laser for each quantum dot since such a system would not feasibly scale to a fault-tolerant quantum computer. Instead, this architecture uses a finite set of lasers; each laser can illuminate the entire array of quantum dots, which in QuDOS is estimated to be about  $10^9$  QDs in size (see section 5.2). To implement desired gates, one must allow this laser light to reach the target quantum dots, while blocking it from interacting with the other QDs. Since the QDs are arranged in a square lattice, one can think of the presence or absence of a light pulse at each as a pixel in an image, and the overall image forms an optical pattern across the surface of the planar cavity. The problem of multiplexing control signals to individual QDs is therefore solved by controlling an optical pattern which illuminates the QD array. Two major challenges must be addressed: (1) the architecture demands control over very many QDs, and (2) the spacing of the QDs is approximately  $1\ \mu\text{m}$  apart, which approaches the diffraction limit of the light for control operations ( $\sim 900\ \text{nm}$ ). The following sections address each of these concerns individually.

*MEMS Micromirrors* Optical control of the QD array requires more than simply generating the light pulses with lasers. These pulses must interact with the correct target quantum dots in time and location. We estimate the number of virtual qubits (and hence QDs) needed in this architecture to be on the order of  $10^9$  (see section



**Figure 4.** SEM image of an array of micromirrors. This image is an experimental sample of vertically actuated mirrors, whereas the MEMS array in QuDOS would require smaller size mirrors with tilting actuation at very high speeds. The inset image shows a magnified view of a single mirror.

5.2). Multiplexing a single laser to millions of targets is a daunting task, especially integrated into a single system. To achieve this, MEMS micromirrors are fabricated in a two-dimensional array so that each mirror serves as an optical modulator for its corresponding QD [55, 56]. Figure 4 shows an example of  $8 \times 8$  array of micromirrors fabricated on a silicon substrate, where each mirror acting as a pixel can move vertically to induce piston motion. A much larger array ( $\sim 10^6$  pixels) of tilting micromirrors is commercially available for use in projection displays today [57]. The commercial digital mirror devices (DMD) feature a switching time of  $\sim 5 \mu\text{s}$  with individual pixel size as small as  $\sim 10 \mu\text{m}$  [58]. With further device optimization, a switching time of  $\leq 1 \mu\text{s}$  is feasible [59]. Using this technology, the laser light falling on each pixel can be turned “on” or “off” by reflecting the light towards or away from the QD. The light pattern pointed towards the QD can be imaged onto the QD array using imaging optics and phase-shift masking.

*Phase-shift Masking* Controlling the individual quantum dots in this computer requires focusing light in a complex pattern with a resolution close to the diffraction limit of the light being used. The quantum dots are spaced  $1 \mu\text{m}$  apart, while the control pulses have wavelength 900–950 nm. Designing such a system would be a formidable challenge, but fortunately it has been achieved already in a mature industry: photolithography for the fabrication of integrated circuits. A typical approach in photolithography is to design an optical mask such that light shined through the mask creates a desired optical profile on an image plane parallel to the mask. Among the various types of masks for manipulating light, a recent technique

is the use of phase-shift masking [60, 61]. Rather than blocking or passing light (as in opaque masks), the phase shift mask is transparent everywhere, but the mask consists of regions which impart different phase shifts to the light which passes through. These phase shifts cause interference patterns in the light on an image plane after the mask. Although the optical pulses are broadband compared to monochromatic laser light, the bandwidth is sufficiently narrow that interference patterns are preserved. Development of these phase-shift masks is complex but routine for photolithography, so QuDOS can leverage a well-studied engineering problem to control the optical pattern of light striking the quantum dot array. The phase-shift masks will create an interference pattern which focuses the laser beam to certain target QDs, while the MEMS mirrors can modulate whether the light pulse is sent to a group of QDs.

*2.1.6. Quantum Non-Demolition (QND) Measurement* The essential measurement operation in QuDOS consists of an optical pulse which uses dispersive quantum non-demolition (QND) readout 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. External photodetectors (section 2.1.7) measure the phase-shift, thereby inducing measurement of the electron spin.

The physics of such a QND measurement has favorable engineering consequences. The fact that the probe pulse is off-resonant means that inhomogeneity among various quantum dots can be tolerated to a higher degree than is true in schemes involving resonant pulses. The technique is simple, and does not require additional “ancilla” quantum dots (or other structures) to be fabricated. Finally, the probe pulse can have a relatively high photon count, resulting in less stringent detector requirements. Several results in recent years have demonstrated the promise of this mechanism for measurement: multi-shot experiments by Berezovsky *et al.* [62] and Atature *et al.* [63] have measured spin-dependent phase shifts in charged quantum dots, and Fushman *et al.* [64] observed a large phase shift induced by a neutral quantum dot in a photonic crystal cavity.

There are however several challenges related to this scheme. First, the measurement must be “single shot” — after just one probe pulse is applied, the measurement result via photodetection is correct with high probability. For a sufficiently large detuning, the ability to complete a single-shot QND measurement depends on the cooperativity factor  $C$  of the cavity. For a sufficiently large detuning  $\Delta$ , the ability to complete a single-shot QND measurement depends only weakly on  $\Delta$ ; the critical factor is the cooperativity factor  $C$ . The phase shift in the probe pulse,  $\theta$ , scales as the detuning, as well as with  $C$ . However, the probability for a photon to create a trion state, which decays by spontaneous emission, scales as  $C/\Delta^2$  [65]. For an input probe pulse that is in a coherent state, the number of photons required to resolve a phase shift  $\theta$  scales as  $1/\theta^2 \propto \Delta^2/C^2$ , indicating that the probability of spontaneous emission during a single-shot measurement (which would spoil its QND character, i.e. introduce measurement error) scales as  $1/C$ , independent of  $\Delta$ . Consequently, a large cooperativity factor of the cavity (e.g.  $C \sim 10^3$ ) may allow single-shot dispersive QND measurements to be carried out with low measurement error.

A second challenge with this measurement scheme involves the selection rules of the quantum dot in a magnetic field. The ultrafast single qubit rotation scheme [42] requires a  $\Lambda$ -system for the electron spin, but this is only available when the static magnetic field is applied perpendicular to the semiconductor growth axis (Voigt

geometry). Thus far, the multi-shot demonstrations of dispersive QND readout have only been performed with parallel orientation [63] or low strength magnetic fields [62]. The demonstration of single-shot QND measurement in Voigt geometry, with high magnetic fields, is still possible, but the phase shifts are likely to be smaller. Specifically, there are two active  $\Lambda$ -systems which provide competing phase shifts, so only the difference between the two phase shifts can be detected. The actual value of the achievable phase shift will depend on the trion energy structure, which in turn depends on quantum dot growth parameters such as strain and dot ellipticity.

*2.1.7. Detector Array* To enable parallel operations in QuDOS, measurement must also be performed in a parallel fashion by an array of photodetectors. A measurement light pulse reflected from the cavity is directed to an integrated grid of CMOS imagers, which is an alternative imaging technology to the more common CCD [66]. Additional control circuitry for Layer 3 is also embedded in this array. The surface code quantum error correction implemented in Layer 3 (see section 4.1.3) must process the outcomes of measurement on syndrome qubits to determine correction operations for errors that have occurred. To make this error analysis step fast and efficient, the necessary processors for error correction are integrated with the array of photodetectors.

The requirements for the detectors will partially be determined by how large the phase shift of the dispersive measurement pulse can be made, which is influenced by the cavity and quantum dot parameters. Fundamentally, however, the detectors are used to make a decision on whether an impinging pulse contains an average photon number below or above a certain threshold. Single photon detection capability is thus not important. It may be possible to compensate for poor detector quantum efficiency by increasing the measurement probe power, although higher quantum efficiency is preferable, since the probability of measurement error scales with probe power. The gain curve of the detector and amplifier circuitry is important, since it may be necessary to distinguish between a 1 nanosecond pulse containing, for example,  $10^5$  photons on average, and a pulse with just 1% more photons on average; the detector must not saturate near the used probe pulse power, and must have sufficient gain at those powers that the small difference in average photon number can be discerned with high probability.

For the expected phase shift, there appears to be no fundamental limitation to engineering a CMOS imaging array that is sufficiently sensitive, fast (GHz operation frequency), and large (one pixel per quantum dot to be simultaneously measured) to meet the requirements for QND readout. Current state-of-the-art CMOS image sensors are not yet advanced enough, especially with respect to speed (frame rate), but rapid progress is being made, driven by commercial requirements in a wide variety of applications [67].

*2.1.8. Static Decoherence (Memory Errors)* The continuous phase evolution discussed in section 2.1.1 would not pose a problem if it was constant — it could be mitigated by synchronizing pulse arrival times to the Larmor period [52]. However, the nuclei in the vicinity of the quantum dot electron also have nonzero spin, so they interact with the electron by the hyperfine interaction. This creates an effective magnetic field with random orientation and bounded magnitude acting on the electron. The effect is that the phase evolution between the spin levels is different for each quantum dot electron, and difficult to determine. However, the nuclear spins are stable

on timescales much longer than the electrons, so that the perturbation to the electron spin is effectively an unknown constant, within the electron  $T_2 \sim 1 \mu\text{s}$  timescale. This phenomenon is responsible for the ensemble dephasing time ( $T_2^* \sim 1 \text{ ns}$  [68]). This error source obscures the system designer’s ability to track the phase evolution of the spin vector on the Bloch sphere, but it is not fatal, as this problem can be addressed by Layer 2 techniques (see section 3). We note also that in some semiconductors, isotopic purification (removing any atoms with nonzero nuclear spin) can improve dephasing times by an order of magnitude [69], but this approach is not possible for QuDOS since there are no stable zero-spin nuclear isotopes of In, Ga or As.

*2.1.9. Dynamic Coherent and Incoherent Errors* Coherent errors, such as deviation in the axis of rotation or angle of rotation on the Bloch sphere, preserve state population in the two-level spin system (which later serves as the foundation of the virtual qubit). In section 3.2, we illustrate techniques in Layer 2 for addressing systematic coherent errors which occur in the Physical layer.

Incoherent processes involve coupling to modes outside of the two-level spin system, which is problematic because this leads to an irrecoverable loss of quantum information. The broadband pulses induce virtual transitions between the metastable spin levels and the excited trion levels. However, there is a possibility that a real excitation of a trion can occur, such as if the detuning from resonance is too small or a phonon interacts with the system to cause actual absorption of a photon (and generation of a trion). The exciton lifetime in the GaAs system is  $T_{1,X} = 1 \text{ ns}$ , so when the trion decays by spontaneous emission, the state of the two-level spin system underpinning our virtual qubit is effectively measured without knowledge of the result, leading to complete depolarization of the qubit.

## *2.2. Layer 1 Performance*

Operational performance is critical when designing a computing system. Execution time and accuracy are particularly important for quantum computers, where expected logical operation speed may play a role in deciding which physical system one chooses. Performance in the physical layer depends on the timing and duration of optical pulses which manipulate the charged QD spin system. Table 1 lists some of the critical parameters for the physical processes required: broadband pulses and precession in the magnetic field are needed for “raw” 1-qubit gate pulses, QND pulses provide spin state measurement, and the entangling operation is the basis of the 2-qubit gate.

*2.2.1. Manipulating the Spin-basis Bloch Sphere* QuDOS must be able to control the charged quantum dot spin system by rotating the spin state represented by a Bloch vector around two orthogonal axes on the Bloch sphere. This requirement arises because we use this system to construct the “virtual qubit” in Layer 2 (see section 3.1.1), and two separate axes of control are necessary for arbitrary  $\text{SU}(2)$  operations.

The first way to control the spin state is through the static magnetic field. The spin states are split in energy, so the relative phase between the two levels precesses at the Larmor frequency of about 25 GHz in a 7 T field. This can be viewed as a continuous rotation around the Z-axis on the Bloch sphere.

Stimulated Raman transitions produced by ultrafast broadband optical pulses incident on the quantum dot coherently rotate population between the spin states, which in the idealized Bloch sphere can be interpreted as X-axis rotations. However,

Operation	Mechanism	Duration	Notes
Spin phase precession (Z-axis)	Magnetic field splitting of spin energy levels	$T_{\text{Larmor}} = 40 \text{ ps}$	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	$\tau_{\text{pulse}} = 14 \text{ ps}$	Red-detuned from spin ground state-trion transitions.
Entangling Operation	Nonlinear phase shift of spin states via coupling to a common cavity mode	$\tau_{\text{entangle}} = 100 \text{ ns}$	CW laser signal modulated by an electro-optic modulator (EOM).
QND Measurement	Dispersive phase-shift of light reflected from planar cavity	$\tau_{\text{QND}} = 1 \text{ ns}$	CW laser signal modulated by an EOM.

**Table 1.** 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.

there are some important non-ideal effects which must be addressed. A perfect X-axis rotation is only possible in the limit of an infinitely fast pulse, where the spin vector precession (due to the magnetic field) on the Bloch sphere during the optical pulse goes to zero. As we see in Table 1, the Larmor frequency is comparable to the broadband pulse duration. Even very fast pulses ( $< 1 \text{ ps}$ ) will still incur significant error, which lowers gate fidelity and increases the burden on Layer 2 to produce gates with error below the threshold of the surface code.

An alternative approach is to tune the broadband pulse amplitude so that the angular velocity of rotation around the X-axis is equal to the velocity of rotation around the Z-axis due to the magnetic field. One can verify that the resulting rotation by an angle  $\pi$  is equivalent to applying a Hadamard gate, which in conjunction with arbitrary Z-rotations from the magnetic field is sufficient to produce any  $\text{SU}(2)$  gate. Moreover, unlike very fast pulses, this operation can in principle produce high-fidelity state rotation.

*2.2.2. Entangling Operation* Universal quantum computation requires a gate which generates entanglement. The surface code requires the virtual gate **Controlled-NOT** (CNOT) which must come from Layer 2. To produce this gate, Layer 1 must provide a mechanism for generating entanglement between the spins of the charged quantum dots. The method proposed for this architecture (discussed in section 2.1.4) couples two electron spins with a common cavity optical mode. Simulation of this process indicates that it requires a modulated continuous-wave laser pulse about 10–100 ns in duration [44]. The effect of this pulse on the virtual qubits formed by two neighboring charged quantum dots is to induce a non-linear phase shift dependent on the state of

the electron spins. This is equivalent to a **Controlled-Z** gate at Layer 2, which can be transformed into the **CNOT** gate with single-qubit operations (in this case, virtual Hadamard gates).

*2.2.3. MEMS Micromirror Switching* Switching delay for the MEMS mirrors depends on their mechanical properties, which are limited by the fabrication processes and actuation requirements. We anticipate that the switching delay can be reduced to  $\leq 1 \mu\text{s}$ , but that alone is insufficient for controlling the crucial optical signals in QuDOS. Most of the optical signals can be applied in a repeated pattern across the entire array of QDs, but measurement and the associated single-qubit gates to change the basis of measurement must be multiplexed to the appropriate quantum dots. Therefore, this one particular set of operations requires two MEMS mirror arrays (designated “*A*” and “*B*” for simplicity) which both couple into a single beamsplitter. Electro-optic modulators (EOMs) control whether laser light signals reach the mirror arrays and reflect to the beamsplitter, which in turn directs light to the QDs. The EOMs alternate which mirror array is “on”, by one transmitting light while the other blocks. When *A* is on, the mirror array is static and a certain measurement pattern is projected onto the QDs for every measurement light pulse. Meanwhile, *B* is re-positioning its mirrors for the next measurement pattern. When *B* is ready, the EOMs switch states, and now *B* is “on” while *A* re-positions its mirrors for the next measurement pattern. By using *alternating MEMS arrays*, we can overcome the slow switching speed of the MEMS mirror technology.

### 3. Layer 2: Virtualization

Quantum information systems are very sensitive to imperfections in their environment and control, which manifest as errors in the stored information. These errors can be systematic or random. Layer 2 sharply reduces systematic errors since this can be accomplished without measuring the system state, which is inherently faster and simpler than error-correcting methods which extract information about errors. Quantum error correction is implemented in Layer 3 to correct general errors, but doing so requires syndrome extraction circuits which implement non-local 2-qubit gates and operate at longer timescales. The purpose of Layer 2 is to reduce the error rate in virtual qubits and gates to the levels sufficient for Layer 3 to function.

#### 3.1. Components of the Virtualization Layer

*3.1.1. Virtual Qubit* The virtual qubit is an abstraction of the underlying physical system. It approximates an ideal qubit as a two-level system whose state is constant until purposefully manipulated. However, the virtual qubit is modeled with real decoherence. In QuDOS, the virtual qubit is created from the two metastable spin states of an electron confined to a QD. The raw physical system has dephasing time  $T_2^* \approx 1 \text{ ns}$  [68] caused by an inhomogeneous distribution of nuclear spins in the environment of the electron. This dephasing time is insufficient for the Layer 3 operations, and so this system must be augmented with dynamical decoupling (DD) techniques [24, 25], which extend the dephasing time of the virtual qubit into the microsecond regime (see section 3.2.1). Additionally, the electron spin vector precesses with the Larmor frequency about the Z-axis on the Bloch sphere, whereas the virtual

qubit is static. This abstraction is achieved by appropriately timing measurement and control optical pulses, as discussed in sections 3.2.1 and 3.2.2.

Measurement of the virtual qubit is achieved by the QND measurement of the spin state from Layer 1 (see section 2.1.6). In principle multiple measurements could be performed in Layer 1 in order to increase measurement fidelity, but this architecture uses single-shot readout for the sake of speed.

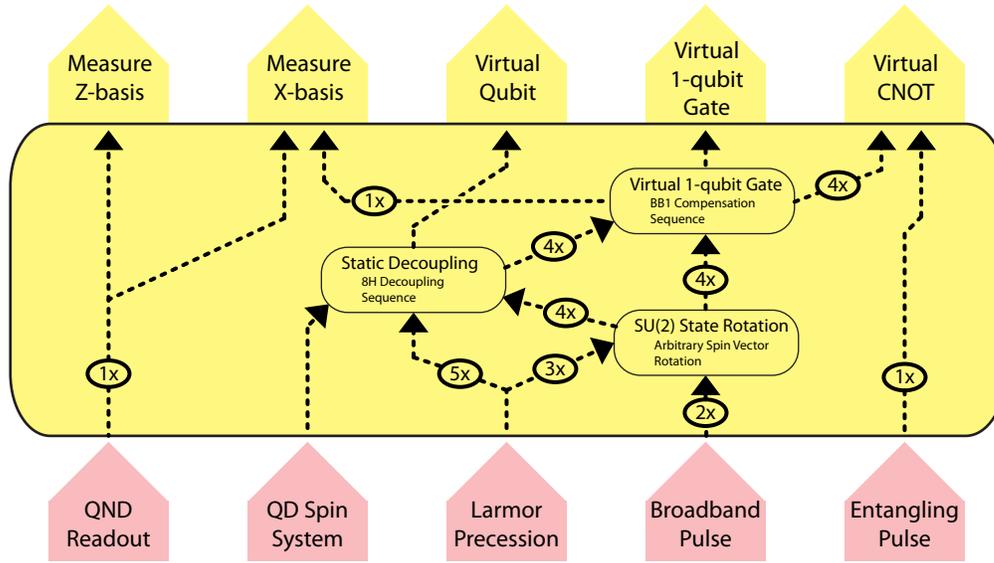
*3.1.2. Virtual Gates* Quantum operations must be implemented by physical hardware which is ultimately faulty to some extent. Many errors are *systematic*, so that they are repeatable, even if they are unknown to the quantum computer designer. In Layer 2, virtual gates manipulate the state of the virtual qubit by combining fundamental control operations in Layer 1 in a manner which creates destructive interference of control errors. Virtual gates must suppress systematic errors as much as possible in order to satisfy the demands of the error correction system (Layer 3). For example, 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. This section explains the theory behind the virtual gate, while section 3.2.2 illustrates how a simple virtual gate scheme is developed in QuDOS.

Efficient schemes exist for eliminating systematic errors. Compensation sequences can correct repeatable (but perhaps unknown) errors in the state rotation operations [26, 27]. This condition is often true since errors are frequently due to imperfections in the Layer 1 processes, such as laser intensity fluctuations over long timescales or the coupling strength of the electron to the optical field (caused by fabrication imperfections). Since these errors are systematic over the timescales of 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. Moreover, many compensation sequences are quite general, so that error-reduction works without knowledge of the type or magnitude of error.

### *3.2. Layer 2 Performance*

The Virtualization layer can sharply reduce systematic device errors, but not random errors; therefore, Layer 2 must operate fast enough to permit Layer 3 to correct the remaining random errors. As discussed in section 2.1.8, memory errors accrue over time regardless of what operation is being executed. Layer 2 mitigates memory and control errors, but if the virtualization operations require too long to execute, the residual error will be above the threshold of the surface code, and Layer 3 cannot function. Figure 5 gives a broad view of Layer 2 in QuDOS, and the following subsections give a detailed analysis of its performance.

*3.2.1. Virtual Qubit* Constructing the virtual qubit requires Layer 2 to conceal the complexity of controlling the QD spin state. The QD electron resides in a strong magnetic environment. The magnetic field induces a splitting of the spin energy levels, causing a time-dependent phase rotation of the spin states. Since the spin levels form the basis of the qubit, this is equivalent to a continuous rotation of the Bloch vector around the Z-axis. Therefore, control pulses must be accurately timed so that they perform the desired operation.

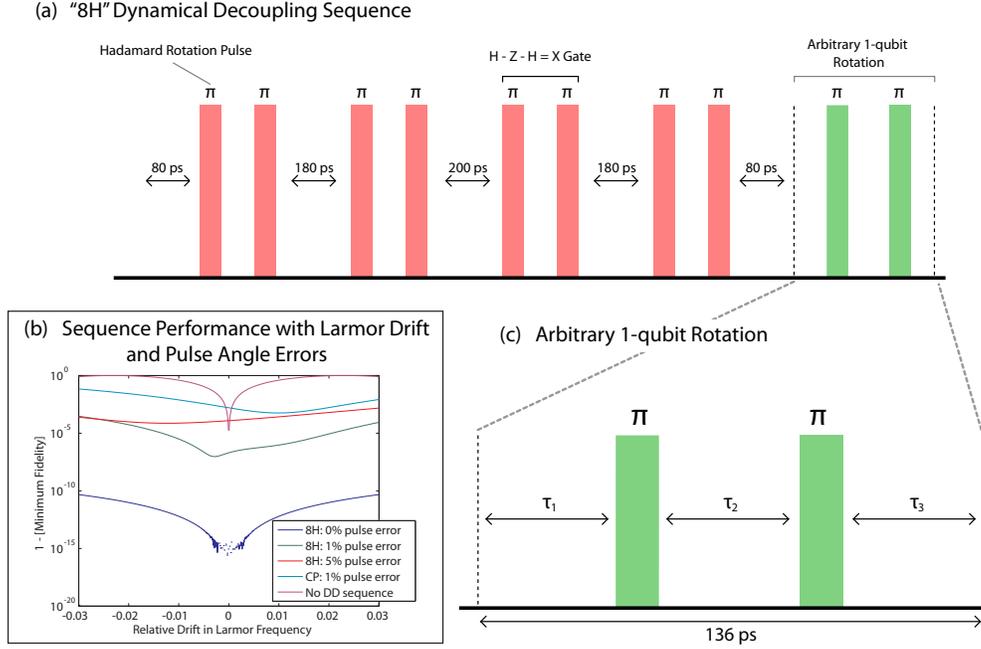


**Figure 5.** The mechanics of the virtualization layer. The outputs of Layer 1 are combined in controlled sequences to produce virtual qubits and gates. Arrows indicate how the output of one process is used by another process. The circled numbers indicate the quantity of a certain resource being used.

Control of the QD spin is complicated by the inhomogeneous nuclear environment which causes the Z-axis rotation to proceed at a somewhat uncertain angular frequency. This problem is mitigated by a dynamical decoupling sequence, so that the system is decoupled from environmental noise and brought into a precisely controlled reference frame at a predictable time. The sequence in Figure 6 illustrates such a decoupling sequence, appropriate for use in this architecture. Although longer sequences may consist of more pulses, to minimize execution time, we have chosen a sequence of eight Hadamard pulses. Instead of using a more common sequence like Carr-Purcell (CP) [70, 71] or Uhrig dynamical decoupling (UDD) [72], the sequence in Figure 6 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.

The virtual qubit is formed by hiding the details of the inhomogeneous phase angular velocity with the DD sequences. Control and readout pulses are timed to arrive exactly when the DD sequence brings the QD spin state back into focus, so that above Layer 2 the virtual qubit appears to be a static quantum memory.

The definition of the virtual qubit is the subspace spanned by the QD electron spin states, which coincides with the measurement process in Layer 1. The measurement pulse and readout projects the electron into one of the spin states. Measurement of the virtual qubit requires that the DD sequence be halted, because decoupling interferes with measurement. Since the measurement pulse is in the Z-basis, rotations around the Z-axis (from the magnetic environment) do not affect the outcome. Neglecting



**Figure 6.** The 8-pulse Hadamard (8H) sequence used in this architecture can eliminate both memory and control errors to first-order. (a) Each bar represents a pulse rotating the Bloch vector by angle  $\pi$  around the axis  $\frac{1}{\sqrt{2}}(\hat{X} + \hat{Z})$ , which is equivalent to the Hadamard gate. The red bars indicate pulses with fixed arrival times which perform dynamical decoupling. The arrival time of the green bars is varied to produce a desired virtual gate. (b) Simulation of the 8H sequence shows that good performance is possible even with both pulse angle errors and drift in the Larmor frequency at a particular quantum dot. In experiments, the Larmor frequency can drift by about  $\pm 2\%$ , which is consistent with the result  $T_2^* \approx 1$  ns [68]. We determine later that based on the threshold of quantum error correcting codes, the error in a virtual gate should be less than  $10^{-3}$ . From simulation, we see that even 5% pulse error in the 8H sequence will reach this performance, whereas 1% pulse error in the CP sequence is insufficient. Not using any dynamical decoupling leads to unacceptable error rates in this system because of the Larmor frequency drift. The simulation runs into a numerical accuracy limit at  $\sim 10^{-15}$ . (c) Construction of an arbitrary 1-qubit rotation requires 136 ps. The delays  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$  are varied to produce an arbitrary gate. The pulses are Hadamard gates.

DD during measurement is acceptable because the longitudinal ( $T_1$ ) relaxation time is very long compared with the measurement pulse duration [73,74].

*3.2.2. Virtual Gates* Virtual gates manipulate the state of the virtual qubit, but they must use Layer 1 to do so. However, the Layer 1 processes have errors which must be suppressed to form successful virtual gates, which are defined as having error rates tolerable for the error correction in Layer 3 to function. Fortunately, control errors are often systematic, and efficient techniques exist for canceling such errors. This quantum computer uses compensation sequences [26,27], in which a series of pulses accomplishes the desired gate in such a manner that pulse errors interfere destructively and cancel to first order.

First, let us construct a “raw” gate, which is a sequence of pulses that accomplishes arbitrary SU(2) rotation of the virtual qubit, but without any error suppression. Figure 6(c) provides a template for constructing any such sequence, which requires at most 136 ps in QuDOS. This operation time is determined by the Larmor precession in Layer 1; the quantity is exact since any errors are predominantly systematic and therefore corrected by Layer 2. To deterministically apply the gate we desire, this sequence must coincide with the common reference frame produced by the DD sequence above. For convenience, we select the delays in Table 3 such that the DD sequence combined with an arbitrary SU(2) gate requires 1 ns. One of the simplest compensation sequences (BB1 [26]) requires 4 arbitrary gates; hence the virtual gate (with error cancelation) requires 4 ns.

The virtual 2-qubit gate is accomplished by constructing a CNOT gate from the entangling operation in Layer 1 and the virtual 1-qubit gates. The DD sequence is designed so that the Ising-like interaction ( $\sigma_z \otimes \sigma_z$ ) component of controlled-phase rotation is allowed while any 1-qubit phase rotations are suppressed. As a result, the CNOT gate can be created by performing 1-qubit virtual gates before and after the entangling operation. However, errors in the entangling operation are due to spontaneous emission [44], which compensation sequences cannot correct; as a result, error in the virtual CNOT gate must be suppressed by design of Layer 1 processes as much as possible.

#### 4. Layer 3: Quantum Error Correction

Error correction schemes remove entropy from an information system. In contrast to Layer 2, quantum error correction schemes [28,75–79] such as the surface code [80] can correct arbitrary errors in the underlying quantum information, assuming the probability of such errors is bounded below a certain threshold [81,82]. This process of information protection is achieved by continually consuming ancilla states prepared to extract entropy from the quantum computer (via syndrome measurement). Layer 3 of this architecture framework is devoted to quantum error correction (QEC), which is vitally important to the successful operation of the quantum computer.

Layers 2 and 3 are not redundant — they are synergistic. The Virtualization layer cannot correct arbitrary errors, and so a large-scale quantum computer will require QEC. However, Layer 2 can mitigate some errors with significantly less effort than would be required in Layer 3, because the Virtualization layer does not extract any information from the system. In this manner, Layer 2 makes Layer 3 more efficient. If errors rates are high, QEC alone may not function at all, and the techniques in the Virtualization layer are essential. The only scenario in which Layer 2 is unnecessary

Operation	Label	Error Cancellation	Composition	Max Duration
Hadamard Rotation	$H$	No	$\frac{1}{2\sqrt{2}}T_{\text{Larmor}}$	14 ps
Z Rotation	$R_Z(\theta)$	No	$\frac{1}{2}T_{\text{Larmor}}$	20 ps
X Rotation	$R_X(\theta)$	No	$H \cdot R_Z(\theta) \cdot H$	48 ps
1-qubit Gate	1Q	No	$R_Z(\theta_1) \cdot H \cdot R_Z(\theta_2) \cdot H \cdot R_Z(\theta_3)$	88 ps
Dynamical Decoupling with 1Q	DD+1Q	Yes ( $T_2$ )	Delay(80 ps) · $R_X(\pi)$ · Delay(180 ps) · $R_X(\pi)$ · Delay(200 ps) · $R_X(\pi)$ · Delay(180 ps) · $R_X(\pi)$ · Delay(80 ps) · 1Q	1 ns
Virtual 1-qubit Gate	Virtual1Q	Yes ( $T_2$ and Gate Error)	DD+1Q <sub>1</sub> · DD+1Q <sub>2</sub> · DD+1Q <sub>3</sub> · DD+1Q <sub>4</sub>	4 ns
Controlled-NOT	CNOT	Partial ( $T_2$ and 1-qubit only)	VirtualH <sup>⊗2</sup> · Controlled-Z · VirtualH <sup>⊗2</sup>	100 ns
Z-basis Measurement/Initialization	MZ/IZ	-	$\tau_{\text{QND}}$	1 ns
X-basis Measurement	MX	-	VirtualH · MZ	5 ns
X-basis Initialization	IX	-	MZ · VirtualH	5 ns

Table 2. Virtualization Layer Operations

(and perhaps harmful) is if the errors in the Physical layer are completely uncorrelated, in which case Layer 2 control techniques yield no benefit. In QuDOS, systematic errors dominate (as witnessed by the fact that  $T_2^*$  is at least three orders of magnitude shorter than  $T_2$  [68]), and so Layer 2 is critical to this architecture.

Within QuDOS, our specific architecture implementation, the surface code is the crucial means to provide logical qubits and gates with the exceptionally low error demanded of a large-scale quantum algorithm such as Shor's factoring algorithm. We will not review the entirety of the surface code here, but instead refer the interested reader to several key works in the field [29, 83, 84]. This section is devoted to the important architecture-related matters of surface code QEC, such as resource requirements in terms of Layer 2 outputs (virtual gates and qubits) and time to implement logical operations.

#### 4.1. Components of the QEC Layer

The QEC layer uses error correction to provide fault-tolerant logical qubits, logical gates, and logical measurement to Layer 4. We explain the salient aspects of the surface code, the error correction scheme in QuDOS, but in general the processes in Layer 3 can vary significantly between different forms of QEC. The surface code provides the ability to correct arbitrary errors with quantum error correction [29, 83, 84]. Virtual

qubits in a broad 2-dimensional array are encoded into a single surface code via single-qubit operations and nearest-neighbor (CNOT) gates. Logical qubits are produced by forming “defects” in the surface code. A defect is a rectangular connected region of virtual qubits in the lattice which have been measured, so that the resulting surface code lattice has an  $SU(2)$  subspace of freedom, equivalent to a qubit. Ref. [29] gives an overview of the steps needed to construct the surface code. For practical matters (explained in Ref. [84]), a logical qubit is constructed from two defects. In contrast to Layer 2, the surface code gathers information on the system state by periodically measuring an error syndrome and using this knowledge to correct errors in post-processing. The probability of an undetected error decreases exponentially as a function of the “distance” [85] of the code, so that logical qubits and gates with arbitrarily low error are possible with a sufficiently large code. However, the virtual qubits and gates must have error rates below the threshold of the surface code (1.4% [86]), so that often error-reduction techniques in Layer 2 are necessary for Layer 3 to function. The error rate in virtual qubits and gates needs to be about an order of magnitude below the threshold, or approximately  $10^{-3}$ , for the surface code to be manageable in size.

*4.1.1. Architecture and the Surface Code* In contrast to some other QEC schemes, the surface code has some key advantages for architecture. In particular, the surface code requires only local and nearest-neighbor gates between qubits in a square lattice. Within this architecture framework, the necessary Layer 2 components for the surface code to function are the injection of single-qubit states needed for non-Clifford gates, a two-dimensional array of qubits with nearest-neighbor coupling (CNOT), and measurement in the X and Z bases [29, 87]. The two-dimensional arrangement with nearest-neighbor CNOT gates is most readily achieved in QuDOS with a physical 2D array of quantum dots, each supporting a virtual qubit. Although the single-qubit Pauli rotations are needed to form a complete set for universal quantum computation, we may neglect these in the present context by simply maintaining a continually-changing Pauli frame in a classical computer and modifying the final measurement results of the quantum computation [88].

*4.1.2. Pauli Frames* A Pauli frame [88, 89] is a simple and efficient classical computing technique to track the result of applying a series of Pauli gates ( $X$ ,  $Y$ , or  $Z$ ) to single qubits. The Gottesman-Knill Theorem implies that tracking Pauli gates can be done efficiently on a classical computer [90]. 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, 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 correction operation *would be applied*, and continue with computation. As stated above, this is permitted for the case of Pauli gates. When a measurement is finally made on a qubit, the result is modified based on the corresponding Pauli gate which should have been applied earlier. This stored Pauli gate is called the Pauli frame [88, 89], 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 quantum computer operations proceeds normally, with the only change being how the final measurement of that qubit is interpreted.

We emphasize that the Pauli frame is a *classical object* stored in the digital circuitry which 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 lattice. 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), so the Pauli frame only tracks one of four values —  $X$ ,  $Y$ ,  $Z$ , or  $I$  — for each virtual qubit in the lattice.

*4.1.3. Measurement and Detector Arrays* As we saw in Layers 1 and 2, the grid of QDs facilitates the nearest neighbor entangling operations for a virtual CNOT gate. Measurement is also done in an array fashion, with a corresponding lattice of photodetectors. This detector array also functions in Layer 3 since the measurement results must be processed at the surface code level. Rather than sending the multitude of measurement results to a separate location (and incur the delays and communication bottlenecks), the surface code error syndrome processors are co-located on-chip with the detectors. This is permitted because we can designate some defects in the surface code as stationary while also never needing to measure the virtual qubits there, so that there is a “shadow” on the detector array. We can use this space for digital logic to process measurement results rather than unused photodetectors. The action of these processors is discussed in sections 4.2.3 and 4.2.4.

## 4.2. Layer 3 Performance

The primary purpose of Layer 3 is to produce logical qubits and gates with arbitrarily low error from the faulty virtual qubits and gates. For this reason, accuracy is the primary performance figure of Layer 3. Assuming 100% yield and independent error sources, any desired logical (Layer 4) accuracy can be achieved, provided that: (a) the Layer 2 operations have error below the Layer 3 threshold (1.4% for the surface code [86]); and (b), the quantum computer has sufficient space in terms of virtual qubits to host a QEC code as large as necessary. For the purposes of QuDOS in this investigation, we assume that both requirements are achievable and analyze the resources needed to realize such a surface code quantum computer. However, we will show that requirement (b) can be very demanding since realistic hardware will have error rates which require very large surface codes. After establishing the space requirements of a fault-tolerant quantum computer architecture (with a specified arbitrary accuracy), we then analyze the time needed to execute Layer 3 operations, which will ultimately determine the logical “clock speed” or operation frequency of the quantum computer in Layer 4. Figure 7 provides a schematic view of the processes inside Layer 3.

*4.2.1. Size of the Surface Code* Quantum error correction schemes generate protected codespaces within a larger Hilbert space formed from an ensemble of qubits. The tradeoff is that instead of a single qubit, the quantum computer now requires many

Parameter	Symbol	Value
Threshold error per virtual gate	$\varepsilon_{\text{thresh}}$	$1.4 \times 10^{-2}$
Error per virtual gate	$\varepsilon_{\text{V}}$	$1 \times 10^{-3}$
Logical circuit depth (in lattice refresh cycles)	$K$	$3.4 \times 10^{11}$
Number of logical qubits (“Shor”, section 5.2)	$Q$	12288
Error per lattice refresh cycle	$\varepsilon_{\text{L}}$	$2.7 \times 10^{-18}$
Surface code distance	$d$	27
Virtual qubits per logical qubit	VQ/LQ	4830

**Table 3.** Parameters Determining the Size of the Surface Code in QuDOS

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 ( $\varepsilon_{\text{V}}$ ), which is an input to Layer 3 from Layer 2
- threshold error per virtual gate of the surface code ( $\varepsilon_{\text{thresh}}$ )
- distance ( $d$ ) of this instance of the surface code
- error per logical gate ( $\varepsilon_{\text{L}}$ ), which is upper-bounded by the performance requirements of the quantum algorithm in Layer 4

To determine  $\varepsilon_{\text{L}}$ , 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

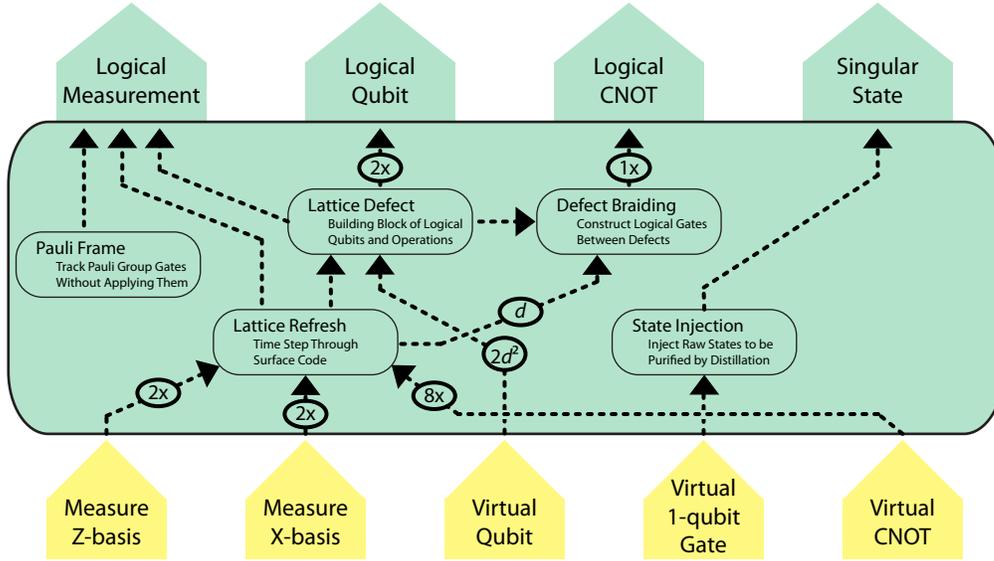
$$P_{\text{fail}} = 1 - (1 - \varepsilon_{\text{L}})^{KQ} \approx KQ\varepsilon_{\text{L}} \quad (1)$$

for small  $\varepsilon_{\text{L}}$ . Therefore, we demand that  $\varepsilon_{\text{L}} \ll 1/KQ$ . Given these quantities, the average error per logical gate in the surface code may be closely approximated [85] by

$$\varepsilon_{\text{L}} \approx C \left( \frac{\varepsilon_{\text{V}}}{\varepsilon_{\text{thresh}}} \right)^{\lfloor \frac{d+1}{2} \rfloor} \quad (2)$$

where  $C$  is a constant determined by the implementation of the surface code. The data in Ref. [85] suggests  $C \approx 3 \times 10^{-2}$ . Therefore, given a known  $\varepsilon_{\text{V}}$ ,  $\varepsilon_{\text{thresh}} = 1.4 \times 10^{-2}$ , and  $C \approx 0.03$ , one can determine the necessary distance  $d$  such that the probability of failure of an entire quantum algorithm is sufficiently small. Table 3 provides an example of these calculations for the QuDOS quantum computer. Error per virtual qubit ( $\varepsilon_{\text{V}}$ ) is also assumed, and the  $K$  and  $Q$  values are for Shor’s algorithm factoring a 2048-bit integer (see section 5.2). We have assumed  $\varepsilon_{\text{L}} \leq 10^{-2}/KQ$ , so that the success probability of the quantum algorithm is greater than 99%.

Determining the necessary distance for the surface code allows one to compute the minimum number of virtual qubits needed to produce one logical qubit, which consists of two defects separated from each other and any other defects or boundaries by the distance of the code. Table 3 calculates this number for QuDOS, but we emphasize that a complete surface code quantum computer will need additional virtual qubits to facilitate movement of defects (braiding) and the distillation of singular qubits needed for non-Clifford logical gates. As a result, the total number of virtual qubits from Layer 2 — and therefore the number of quantum dots from Layer 1 — is larger than simply the product of [virtual qubits per logical qubit]  $\times$  [logical qubits].



**Figure 7.** Process translation in Layer 3. 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. The circled number is the quantity of the corresponding resource which is used.

Accounting for these additional virtual qubits is crucial to accurately estimating the resource requirements for QuDOS. More generally, the quantity of these additional qubits depends significantly on the algorithm one is implementing, since the number of singular qubits is related to the types of logical gates one must implement. A total accounting of the virtual qubits in the surface code is given in section 5.2.

*4.2.2. Surface Code Operations in Time* The fundamental time step in Layer 3 is one unit along the simulated time axis of the topological cluster state [87], which is the lattice refresh cycle of the surface code. Within this architecture, the virtual gates needed for this process are performed in parallel across the entire array of virtual qubits. This parallelism is a fundamental strength of the architecture, because the lattice refresh time can be very fast as shown in Table 4. Moreover, refresh time is independent of system size since all operations proceed in parallel; by contrast, in the architecture in Ref. [22], lattice refresh time depends on the size of one axis of the lattice.

Logical qubits in the surface code are defects in the lattice [29], and logical operations involve braiding these defects through simulated time. Figure 7 illustrates how functions of the surface code are constructed from Layer 2 services. For error-correction purposes, the speed of a braiding operation is constrained by the distance of the code, so that if the minimum spatial separation of defects is 27 virtual qubits (as in Table 3), the time to perform braiding must also take at least 27 lattice refresh cycles. This is because the error chains in the surface code can span both spatial and

Operation	Label	Composition	Max Duration
Lattice Refresh with <i>alternating MEMS arrays</i>	LatticeRefresh	$2 \times (\text{IZ} \cdot 4 \times \text{CNOT} \cdot \text{MZ} \cdot \text{IX} \cdot 4 \times \text{CNOT} \cdot \text{MX})$	1.61 $\mu\text{s}$
Defect Braiding	DefectBraid	$27 \times \text{LatticeRefresh}$	43.5 $\mu\text{s}$
Logical CNOT	LogicalCNOT	DefectBraid	43.5 $\mu\text{s}$

**Table 4.** Layer 3: Surface Code Operations

temporal dimensions, with units along the time axis defined by the lattice refresh time. As such, any defects must be separated by distance  $d$  in time and space within the surface code. Table 4 shows the time required for several Layer 3 processes including the construction of logical gates such as CNOT.

*4.2.3. Local Error Correction Processing* The error syndrome decoding process in Layer 3 requires the location of error chain endpoints and the subsequent matching of these endpoints into a minimum-weight set of error chains [29, 85, 87]. Section 4.1.3 suggested a method for integrating local surface code processors into the photodetector array. Devitt *et al.* describe how to split the syndrome decoding problem into smaller manageable chunks [20], an approach which is supported by the use of local error correction processors. Nevertheless, minimum-weight matching can require a significant number of calculations, so even special-purpose processors may require a substantial amount of time to complete this task. If latency and classical processing cause significant delay of the availability of the logical measurement result, future logical operations depending on the result can be delayed with logical identity gates. This is computationally reasonable as a linear increase of the size of the logical qubit results in only polynomial increase of the classical processing time but an exponential increase of the logical qubit lifetime, implying arbitrary delay can be handled with only logarithmic overhead. As a result, Layer 3 must signal to Layer 4 when error syndrome processing requires more time, so that Layer 4 inserts the necessary logical identity gates into the sequence of logical operations.

*4.2.4. Pauli Frames in Action* In this architecture, Pauli frames are dynamic objects, just like a virtual qubit. However, unlike virtual qubits, they are entirely *classical* objects, since they carry two bits of classical information for each qubit to which they apply. For this reason, they exist in the digital circuitry associated with error syndrome processing, since these same processors determine what the Pauli frame should be.

The manner in which a Pauli frame is implemented could be compared to a classical parity mask. Imagine there is a string of data bits, and one has determined where in this string some bits were flipped by errors. This error correction information is stored in a second bit string (parity mask), which consists of 1s where bit-flip errors occurred, and 0s elsewhere. The recovery operation is then the bitwise XOR of the two strings. Returning to our quantum computer, Layer 3 will continually record the results of the syndrome measurement step. Before any operation requiring logical measurement (such as a non-Clifford gate), any errors which have occurred must be identified. A parity mask for the entire surface code is created, with an entry for each

virtual qubit. The minimum weight matching algorithm processes the accumulated syndrome information and pinpoints the locations of  $X$  errors; the parity mask is then updated by flipping the existing entry corresponding to each  $X$  error. Note that after this update process, it is possible that a virtual qubit will have experienced two  $X$  errors at different times, which cancel each other. The parity mask appropriately has a 0 entry in this event. Complementary to the identification of  $X$  errors, a similar procedure is performed for  $Z$  errors in a second parity mask. The combination of the  $X$  and  $Z$  parity masks is the Pauli frame for Layer 3.

The Pauli frame comes into action whenever a logical measurement is made on a pair of defects in the surface code. The individual virtual qubit measurements are modified based on the Pauli frame in Layer 3. If the measurement basis applied at each virtual qubit commutes with the corresponding Pauli frame entry, the measurement result is unchanged. If the measurement anti-commutes with current Pauli frame entry at this virtual qubit, then the measurement result is flipped. This action is comparable to the bitwise XOR mentioned above for a classical bit string. Note that the presence of both an  $X$  and  $Z$  error in the Pauli frame is tantamount to a  $Y$  error, as global phase is irrelevant to measurement outcome. After the adjusted measurement result is reported, the corresponding Pauli frame entry is reset to  $I$  (0 entry for both  $X$  and  $Z$  masks). We emphasize that the error *syndrome* in Layer 3 is not identical to the Pauli frame. Using the minimum weight matching algorithm, Pauli frames (corresponding to the identified locations of errors) are determined based on the syndrome, which is the location of error chain endpoints without explicit knowledge of the error chains themselves.

## 5. Layer 4: Logical

The Logical layer transforms the outputs of the QEC layer into a complete substrate for quantum computing which is used by the Application layer. The QEC layer provides logical qubits and a limited set of logical gates; however, the Application layer may request any arbitrary quantum gate, and it is the task of Layer 4 to create this gate. A specific implementation of the Logical layer depends on what services Layer 3 provides. We develop Layer 4 in the context of using the surface code in Layer 3, which provides logical qubits, logical CNOT, and injected singular states. In another quantum computer where the QEC layer provides different outputs, a different set of processes in the Logical layer may be needed.

### 5.1. Functions of the Logical Layer

The function of the logical layer is to provide the logical qubits and gates needed for the quantum algorithm in the Application layer. The surface code produces logical qubits and gates with arbitrarily high accuracy. However, the only fault-tolerant gates provided by the surface code are the Pauli 1-qubit gates (trivially performed by updating the Layer 4 Pauli frame), initialization and measurement in the  $X$  and  $Z$  bases, the CNOT gate and the identity gate. Rotations about the  $X$  and  $Z$  Bloch sphere axes can be achieved given ancilla states of the form  $\frac{1}{\sqrt{2}}(|0\rangle + e^{i\theta}|1\rangle)$ , which can be created using non-fault-tolerant techniques. For  $\theta = \pi/2, \pi/4$ , fault-tolerant state distillation circuits can be used to obtain arbitrarily high fidelity ancilla states enabling arbitrarily high fidelity rotations of these angles. Similar techniques can be used to create ancilla states enabling Toffoli to be implemented. By the Solovay-Kitaev

theorem [91], these gates are sufficient to efficiently approximate arbitrary single-qubit logical unitary gates.

*Logical Pauli Frame* Just as in Layer 3, it is unnecessary to implement logical Pauli gates. Instead, a second Pauli frame exists in Layer 4 which functions exactly like its counterpart in Layer 3 (see sections 4.1.2 and 4.2.4). Whenever a logical Pauli gate would be applied, the corresponding entry in the Layer 4 Pauli frame is parity flipped instead. However, the performance requirements of this Pauli frame are not as strict as the one in Layer 3, so the Logical Pauli frame can exist in software which controls the Logical layer, instead of dedicated hardware as is necessary for the Pauli frame in Layer 3. This can be seen in Table 4, where the fastest rate one would need to apply Logical Pauli gates is after each round of defect braidings, or once every 56.4  $\mu\text{s}$ , and the number of entries in the Layer 4 Pauli frame corresponds to logical qubits, which is 12288. Conversely, the Pauli frame in Layer 3 must be updated every lattice refresh cycle (1.61  $\mu\text{s}$ ), and the number of entries is the number of virtual qubits in the surface code:  $9.04 \times 10^8$ , as calculated below in section 5.2. Conventional computers can handle the workload of the Layer 4 Pauli frame in software, but the workload of the Layer 3 Pauli frame demands the custom-designed processors described in sections 4.1.3 and 4.2.3.

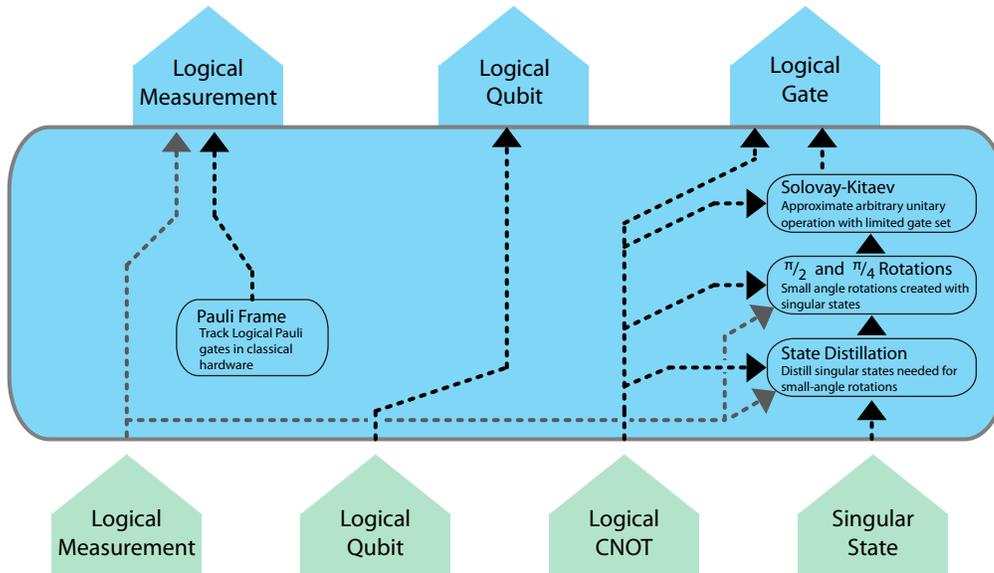
## 5.2. Layer 4 Performance

The Logical layer supports the Application layer, and so we analyze the resources in Layer 4 for the specific purpose of executing Shor’s algorithm in Layer 5. The number of logical qubits for Shor’s algorithm depends on the particular implementation of the algorithm [6, 92–94]. The algorithm adopted for this quantum computer architecture scales as  $\sim 6N$ , where  $N$  is the number of bits in the number to be factored [95], so that factoring a 2048-bit number requires approximately 12288 logical qubits. However, auxiliary logical qubits are also needed for state distillation (explained in section 5.2.1) required for the Toffoli gates used in the modular exponentiation step of Shor’s algorithm (see also section 5.2.2). To factor a 2048-bit number, approximately 90000 logical qubits is sufficient [22]; fewer qubits can be used at the expense of time, since the Toffoli gate operations would be delayed until the injected states are distilled. Additionally, “wiring space” is added for the surface code to enable defects to move and braid when necessary, which is estimated as a 25% overhead in the size of the surface code. With these figures in place, we can estimate the resources required for this quantum computer, as shown in Table 5.

*5.2.1. Singular State Distillation* Singular qubit states described in section 5.1 are necessary to produce arbitrary logical gates in Layer 4. These singular qubits can be produced by magic state distillation [29, 96]. This process consumes a great deal of resources in the quantum computer since many logical qubits are used for distillation. If states are injected with an approximate error of 0.1%, then distilling a  $|Y\rangle$  state ( $\theta = \pi/2$ ) requires two levels of distillation, or at least 49 injected states, to produce one logical  $|Y\rangle$  qubit with infidelity (error) of  $2.4 \times 10^{-24}$ . Using one level of distillation is insufficient because the resulting error in the qubit is  $7.0 \times 10^{-9}$ , which is much greater than the logical error rate ( $\varepsilon_L = 9.8 \times 10^{-19}$ , see section 4.2). Similarly, distilling an  $|A\rangle$  state ( $\theta = \pi/4$ ) requires at least two levels of distillation, or at least 225 injected  $|A\rangle$  states, to produce one logical  $|A\rangle$  qubit with infidelity (error) of

Resource	Label	Composition	QuDOS Quantity
Application Logical Qubits	AppQubits	$6 \times [\textit{bit size of number to be factored}]$	12288
State Distillation Qubits	DistQubits	(Determined by algorithm)	78000
Size of the Surface Code in Virtual Qubits		$1.25 \times \text{VQ/LQ} \times (\text{AppQubits} + \text{DistQubits})$	$9.04 \times 10^8$

**Table 5.** Layer 4 Resources in Terms of Logical Qubits and the Corresponding Size of the Surface Code in Virtual Qubits



**Figure 8.** 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 the Solovay-Kitaev algorithm.

$1.5 \times 10^{-21}$ . Consequently, distillation must be performed continuously in parallel with other logical operations to ensure that these purified states are available on demand. Ref. [84] discusses the resource cost and error scaling of this process in more detail. It is noteworthy that distillation is probabilistic, but the probability of success is high for high-fidelity injected states. We assume the injected states are formed from an initialized virtual qubit and one virtual gate. Since the production of these qubits is critical to the performance of a surface code quantum computer, the injection operations in Layers 1 and 2 should be optimized so that distillation converges to a high-fidelity logical qubit quickly.

Operation	Quantity in QuDOS
Toffoli Gate Time	600 $\mu$ s
ModExp Circuit Depth (Toffoli Gates)	$7.38 \times 10^8$
Minimum Execution Time for Shor's Algorithm	6 days

**Table 6.** Layer 5 Performance for Shor's Algorithm factoring a 2048-bit number.

*5.2.2. Construction of a Toffoli Gate* A construction of the Toffoli gate is given in Ref. [23] on p. 182. The small-angle phase and  $\pi/8$  gates require singular qubits, as described in Ref. [29]. We will assume here that these singular qubits are distilled as needed, but accurately accounting for the braiding operations and volume of the surface code required for singular state production is an active area of research. Hence we account for the time required to implement a Toffoli gate in terms of the number of braiding steps (the depth of this circuit) assuming the singular states are available. There is one braiding for each gate, giving a total of 13 braidings. This implies a minimum Toffoli gate time of 600  $\mu$ s.

*5.2.3. Summary of Logical Layer Performance* We do not calculate the time required to implement any arbitrary gate within the Logical layer, but the prescription is straightforward:

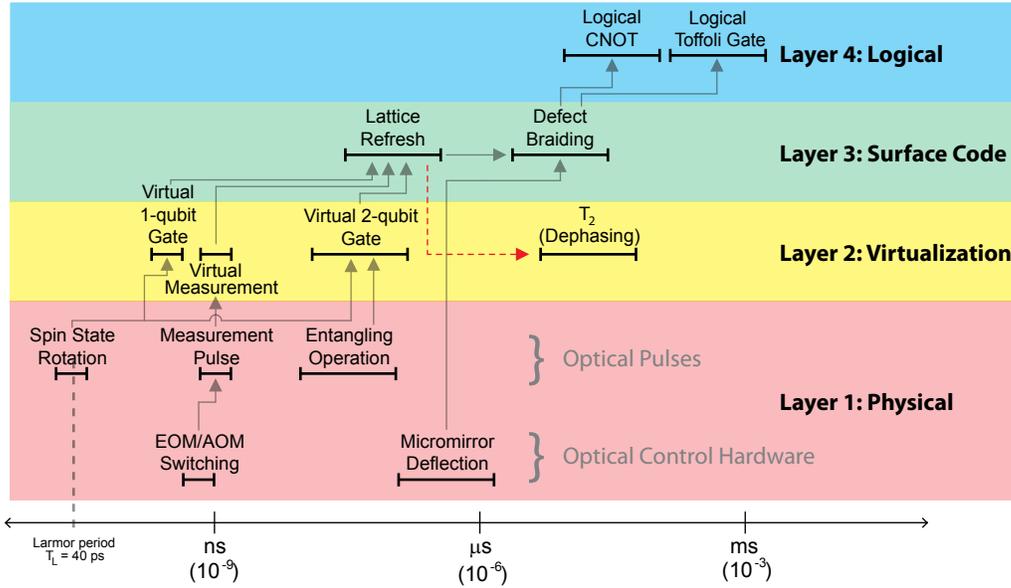
- Specify each logical gate and its accuracy tolerance.
- Use the Solovay-Kitaev algorithm [91] to determine a sequence of available gates from the surface code or state distillation which accurately approximates the desired logical gate.
- Decompose this sequence of gates into the necessary braiding operations, and calculate the total time required.

For Shor's algorithm, the modular exponentiation subroutine (**ModExp**) is the bottleneck to performance. The **ModExp** process depends principally on the Toffoli gate, so we use this figure to estimate the run-time of Shor's algorithm in section 6.

## 6. Application Layer

The Application layer hosts the quantum algorithm, such as Shor's algorithm [97], that a classical user wishes to execute. Logical gates constructed in Layer 4 are performed on the logical qubits provided by the QEC layer, and the end result is communicated to the classical user. The Application layer is completely unaware of the underlying hardware, since it interfaces only with Layer 4. Since the lower layers have provided all the resources for quantum computing, the figures of merit in Layer 5 are the number of available qubits and the speed of logical operations, which implies the time required to implement a certain quantum algorithm. The size of QuDOS in terms of both virtual and logical qubits was given in Table 5, and the run-time for Shor's algorithm factoring a 2048-bit number is given in Table 6.

We must note however that Table 6 gives a minimum execution time, which can be slowed by some processes we have not fully analyzed here. If the syndrome decoding process in Layer 3 takes longer than expected, non-Clifford operations, such as gates which require singular states, will be delayed. Additionally, we have not yet accurately



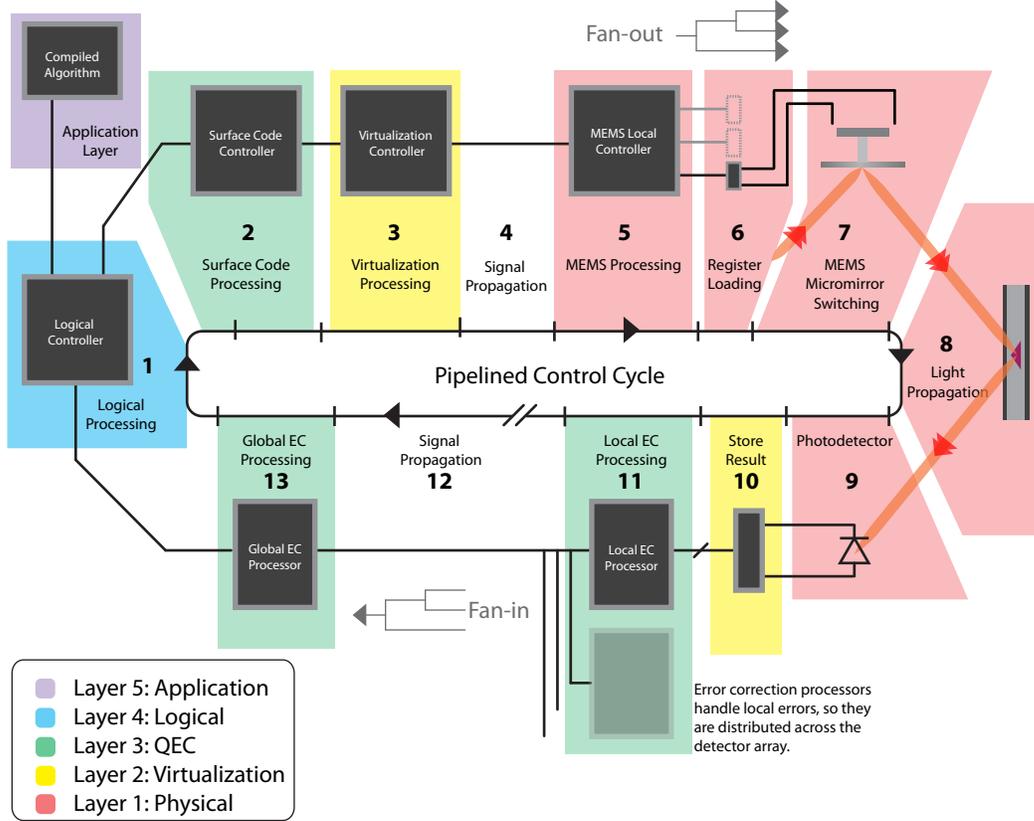
**Figure 9.** Relative timescales for critical operations in QuDOS within each layer. The arrows indicate dependence of higher operations on lower layers. The red arrow signifies that the surface code lattice refresh must be much faster than the dephasing time in order for error correction to function. The Application layer is not included here since quantum algorithms can vary widely in the time they take to implement, while this figure is concerned with the fundamental operations in a quantum computer.

simulated the braiding operations for the singular state distillation in Layer 4. If this procedure requires more time or space in the surface code than estimated above, then the overall algorithm run-time will suffer. Finally, connectivity is a related issue. We assumed in section 5.2 a 25% overhead, but perhaps more is necessary or long-range interactions in the surface code will become bottlenecks. These matters are the subject of future work.

### 7. 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 surface code operations, and so forth. This system of dependence of operation times is depicted in Figure 9. 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 Figure 9, one can 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. For example, the virtualization layer must construct a virtual 1-qubit



**Figure 10.** Primary control cycle of QuDOS. Though presented as a sequence of steps, each layer issues multiple commands to the layer below, and operations are pipelined. Each of the numbered steps represents a process whose time duration must be precisely determined in order for the architecture to function efficiently. Steps 4 and 12 indicate significant delays in the propagation of a control or readout signal.

gate from a sequence of spin-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. Figure 10 shows the QuDOS control loop, which is a more detailed rendition of Figure 2. Here one can see how the essential operations in the architecture interact. In particular, the error syndrome decoding step in Layer 3 is separated into two components, with local processors (described in section 4.2.3) handling small subsections of the surface code while a global processor integrates the results of the local processors into one consistent error pattern. In particular, the global processor corrects any logical measurements based on the Layer 3 Pauli frame (sections 4.1.2 and 4.2.4). The processors in steps 1, 2, and 3 of Figure 10 coordinate the operations in the corresponding layers of the architecture.

The switching delay of  $1 \mu\text{s}$  means that one set of MEMS mirrors cannot multiplex all of the laser pulses in this architecture. Witness in Figure 9 that virtual 1-qubit gates and measurement must operate much faster than this delay permits. Therefore,

two micromirror arrays are needed for the crucial measurement operations in Layer 1 and the associated single-qubit virtual gates which change the basis of measurement. As explained in section 2.2.3, one mirror array is used to multiplex light signals while the second is re-positioning. When the second mirror is in place, the electro-optic modulators switch so that the second array multiplexes light while the first re-positions. In this manner, the relatively slow switching delay of the MEMS mirrors can be circumvented, at the expense of losing some optical power.

## 8. 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. 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 fault-tolerance in an intelligent fashion. Layer 2 uses simple control without monitoring qubit states 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 2048-bit number. There is a significant interplay between Layers 2 and 3, because Layer 2 enhances the effectiveness of Layer 3, which is discussed further in section 4. 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 operation times for fundamental quantum gates are discussed in section 2, but the importance of these fast processes becomes clear in Figure 9, where the overhead of virtual gates in Layer 2 and QEC in Layer 3 increases the time to implement quantum gates from nanoseconds in the Physical layer to milliseconds in the Logical layer, or six orders of magnitude. In this context, a quantum computer needs very fast physical operations.

Much like their classical counterparts, quantum computers need to scale to large size in order to solve useful problems. In fact, when error correction is included, a quantum computer requires sizable classical processing as well. As with the integrated circuit (IC) industry, integrated fabrication is one of the best methods to solve this dilemma. QuDOS is particularly well-suited to device integration since the quantum memory resides in charged quantum dots, which are formed by solid-state fabrication methods inherited from the semiconductor field. Likewise, the MEMS

mirrors and photodetector arrays are also designed to be fabricated with very large-scale integration (VLSI). This approach is proven and robust, and the state of the art is continually being advanced for reasons outside of quantum computing.

The QuDOS architecture has been designed to take advantage of existing mature technologies wherever possible. Doing so allows the architecture to benefit from advances in technology due to other fields of research. For example, the projection optics and phase-shift masking were developed for photolithography in the IC industry. MEMS micromirrors were invented to multiplex light signals in high-definition televisions and projectors. Utilizing these established technologies reduces the burden on quantum computing engineers.

An area which requires further development is the optical engineering in this system. The development of phase-shift masks is routine in photolithography, but it is still very computationally demanding. Moreover, adapting this technique to the current system is complicated by the beamsplitters and MEMS mirrors along the optical paths, as well as the spectral width of the pulses and how they interact with the quantum dots in the planar cavity. Additionally, focusing the laser light signals to the quantum dot array may require a projection optics system, which we have not studied here. Similarly, while integrating logic with photodetectors is feasible, determining the complexity of custom circuits for the syndrome processing step in the surface code is an area of future work. Moreover, integrating digital logic with high-gain photodetectors needed in the context of QuDOS may also present challenges we have not analyzed here.

One of our principle objectives is to better understand the resources required to construct a quantum computer which 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 optical lattices. 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.

- [1] T. D. Ladd, F. Jelezko, R. Laffamme, Y. Nakamura, C. Monroe, and J. L. O'Brien. Quantum computers. *Nature*, 464:45–53, 2010.
- [2] David P. DiVincenzo. The physical implementation of quantum computation. *Fortschritte der Physik*, 48(9-11):771–783, September 2000.
- [3] Andrew M. Steane. Quantum computer architecture for fast entropy extraction. *Preprint* arXiv:quant-ph/0203047., 2002.
- [4] Andrew M. Steane. How to build a 300 bit, 1 Giga-operation quantum computer. *Preprint* arXiv:quant-ph/0412165., 2004.
- [5] Timothy P. Spiller, William J. Munro, Sean D. Barrett, and Pieter Kok. An introduction to quantum information processing: applications and realizations. *Contemporary Physics*, 46(6):407–436, November 2005.
- [6] Rodney Van Meter and Mark Oskin. Architectural implications of quantum computing technologies. *ACM Journal of Emerging Technologies in Computing Systems*, 2(1):31–63, Jan 2006.
- [7] D. Kielpinski, C. Monroe, and D. Wineland. Architecture for a large-scale ion-trap quantum computer. *Nature*, 417:709–711, June 2002.
- [8] Dean Cosey, Mark Oskin, Tzvetan Metodiev, Frederic T. Chong, Isaac Chuang, and John Kubiawicz. 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)*, pages 65–74, New York, NY, USA, 2003. ACM.
- [9] 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(1):74–83, Jan. 2006.
- [10] M. Oskin, F.T. Chong, I.L. Chuang, and J. Kubiawicz. Building quantum wires: the long and the short of it. *30th International Symposium on Computer Architecture, 2003 (ISCA'03)*, pages 374 – 385, June 2003.
- [11] B. E. Kane. A silicon-based nuclear spin quantum computer. *Nature*, 393:133–137, May 1998.
- [12] L.-M. Duan and C. Monroe. Colloquium: Quantum networks with trapped ions. *Rev. Mod. Phys.*, 82(2):1209–1224, Apr 2010.
- [13] Jungsang Kim and Changsoo Kim. Integrated optical approach to trapped ion quantum computation. *Quantum Information and Computation*, 9:181–202, 2009.
- [14] 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(17):174507, Nov 2007.
- [15] M. Whitney, N. Isailovic, Y. Patel, and J. Kubiawicz. Automated generation of layout and control for quantum circuits. *Proceedings of the 4th International Conference on Computing Frontiers*, pages 83–94, 2007.
- [16] M. Whitney, Y. Isailovic, N. and Patel, and J. Kubiawicz. A fault tolerant, area efficient architecture for shor's factoring algorithm. *36th International Symposium on Computer Architecture, 2009 (ISCA'09)*, 2009.
- [17] N. Isailovic, Y. Patel, M. Whitney, and J. Kubiawicz. Interconnection networks for scalable quantum computers. *33rd International Symposium on Computer Architecture, 2006 (ISCA'06)*, pages 366–377, 2006.
- [18] N. Isailovic, M. Whitney, Y. Patel, and J. Kubiawicz. Running a quantum circuit at the speed of data. *35th International Symposium on Computer Architecture, 2008 (ISCA'08)*, 2008.
- [19] René Stock and Daniel F. V. James. Scalable, high-speed measurement-based quantum computer using trapped ions. *Phys. Rev. Lett.*, 102(17):170501, Apr 2009.
- [20] 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(1–2):121–147, 2010.
- [21] Rodney Van Meter, Kae Nemoto, W. J. Munro, and Kohei M. Itoh. Distributed arithmetic on a quantum multicomputer. *SIGARCH Comput. Archit. News*, 34(2):354–365, 2006.
- [22] 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 available as arXiv:quant-ph/0906.2686v2.
- [23] Michael A. Nielsen and Isaac L. Chuang. *Quantum Computation and Quantum Information*. Cambridge University Press, 1 edition, October 2000.
- [24] Lorenza Viola and Emanuel Knill. Robust dynamical decoupling of quantum systems with bounded controls. *Phys. Rev. Lett.*, 90(3):037901, Jan 2003.
- [25] Hui Khoon Ng, Daniel A. Lidar, and John Preskill. Combining dynamical decoupling with fault-tolerant quantum computation. *Preprint* arXiv:quant-ph/0911.3202., 2009.

- [26] Kenneth R. Brown, Aram W. Harrow, and Isaac L. Chuang. Arbitrarily accurate composite pulse sequences. *Phys. Rev. A*, 70(052318), 2004.
- [27] Y Tomita, J T Merrill, and K R Brown. Multi-qubit compensation sequences. *New Journal of Physics*, 12(1):015002, 2010.
- [28] John Preskill. Fault-tolerant quantum computation. *Preprint* arXiv:quant-ph/9712048., 1997.
- [29] Austin G. Fowler, Ashley M. Stephens, and Peter Groszkowski. High-threshold universal quantum computation on the surface code. *Phys. Rev. A*, 80(5):052312, Nov 2009.
- [30] Daniel A. Lidar and K. Birgitta Whaley. Decoherence-free subspaces and subsystems. *Preprint* arXiv:quant-ph/0301032., 2003.
- [31] W. Dür, H.-J. Briegel, J. I. Cirac, and P. Zoller. Quantum repeaters based on entanglement purification. *Phys. Rev. A*, 59(1):169–181, Jan 1999.
- [32] Austin G. Fowler, David S. Wang, Charles D. Hill, Thaddeus D. Ladd, Rodney Van Meter, and Lloyd C. L. Hollenberg. Surface code quantum communication. *Phys. Rev. Lett.*, 104(18):180503, May 2010.
- [33] John Paul Shen and Mikko H. Lipasti. *Modern processor design: fundamentals of superscalar processors*. McGraw-Hill Higher Education, 2005.
- [34] Gunnar Björk, Stanley Pau, Joseph Jacobson, and Yoshihisa Yamamoto. Wannier exciton superradiance in a quantum-well microcavity. *Phys. Rev. B*, 50(23):17336–17348, Dec 1994.
- [35] A. Imamoglu, 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(20):4204–4207, 1999.
- [36] N H Bonadeo, Gang Chen, D Gammon, and D G Steel. Single quantum dot nonlinear optical spectroscopy. *Physica Status Solidi B*, 221(1):5–18, 2000.
- [37] 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(24):241310, Jun 2002.
- [38] 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(16):161306, Apr 2005.
- [39] Y Yamamoto, T D Ladd, D Press, S Clark, K Sanaka, C Santori, D Fattal, K M Fu, S Hfing, S Reitzenstein, and A Forchel. Optically controlled semiconductor spin qubits for quantum information processing. *Physica Scripta*, 2009(T137):014010, 2009.
- [40] 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. Klopff, and F. Schäfer. Fine structure of neutral and charged excitons in self-assembled in(ga)as/(al)gaas quantum dots. *Phys. Rev. B*, 65(19):195315, May 2002.
- [41] 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. *Applied Physics Letters*, 90:251109, 2007.
- [42] 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.
- [43] J Berezovsky, M H Mikkelsen, N G Stoltz, L A Coldren, and D D Awschalom. Picosecond coherent optical manipulation of a single electron spin in a quantum dot. *Science*, 320(5874):349–352, 2008.
- [44] T. D. Ladd and Y. Yamamoto. A simple, robust, and scalable quantum logic gate with quantum dot cavity QED systems. *Preprint* arXiv:quant-ph/0910.4988v1., 2009.
- [45] D. Loss and D. P. DiVincenzo. Quantum computation with quantum dots. *Phys. Rev. A*, 57:120–126, 1998.
- [46] R. Hanson, L. P. Kouwenhoven, J. R. Petta, S. Tarucha, and L. M. K. Vandersypen. Spins in few-electron quantum dots. *Rev. Mod. Phys.*, 79(4):1217–1265, 2007.
- [47] Makoto Kuwahara, Takeshi Kutsuwa, Keiji Ono, and Hideo Kosaka. Single charge detection of an electron created by a photon in a g-factor engineered quantum dot. *Applied Physics Letters*, 96(16):163107, 2010.
- [48] Danny Kim, Samuel G. Carter, Alex Greulich, Allan S. Backer, and Daniel Gammon. Ultrafast optical control of entanglement between two quantum dot spins. *Preprint* arXiv:quant-ph/1007.3733., 2010.
- [49] C. Piermarocchi, P. Chen, L. J. Sham, and D. G. Steel. Optical RKKY interaction between charged semiconductor quantum dots. *Phys. Rev. Lett.*, 89:167402, 2002.
- [50] G. F. Quinteiro, J. Fernandez-Rossier, and C. Piermarocchi. Long-range spin-qubit interaction mediated by microcavity polaritons. *Phys. Rev. Lett.*, 97(9):097401–4, 2006.

- [51] A. Imamoglu, D.D. Awschalom, G. Burkard, D.P. DiVincenzo, D. Loss, M. Shermin, and A. Small. Quantum information processing using quantum dot spins and cavity QED. *Phys. Rev. Lett.*, 83:4204, 1999.
- [52] 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.
- [53] T. Szkopek, P.O. Boykin, Heng Fan, V.P. Roychowdhury, E. Yablonovitch, G. Simms, M. Gyure, and B. Fong. Threshold error penalty for fault-tolerant quantum computation with nearest neighbor communication. *Nanotechnology, IEEE Transactions on*, 5(1):42 – 49, jan. 2006.
- [54] T D Ladd *et al*. “High-Speed Quantum Computer with Semiconductor Spins” in *Semiconductor Quantum Bits*, eds. F. Henneberger and O. Benson, 453. Singapore: Pan Stanford Publishing, 2009.
- [55] 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(2):322–329, 2007.
- [56] Caleb Knoernschild, Changsoon Kim, Bin Liu, Felix P. Lu, and Jungsang Kim. Mem-based optical beam steering system for quantum information processing in two-dimensional atomic systems. *Optics Letters*, 33(3):273–275, 2008.
- [57] Texas Instruments. <http://www.dlp.com>.
- [58] D. Dudley and C. Dunn. Dlp technology — not just for projectors and tvs. *Photonik International*, pages 98–101, 2006.
- [59] 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)*, pages 1–2, may. 2007.
- [60] Y. Liu and A. Zakhor. Binary and phase shifting mask design for optical lithography. *Semiconductor Manufacturing, IEEE Transactions on*, 5(2):138 –152, may. 1992.
- [61] 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(11):2145–2152, 1998.
- [62] 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(5807):1916–1920, 2006.
- [63] Mete Atatüre, Jan Dreiser, Antonio Badolato, and Atac Imamoglu. Observation of faraday rotation from a single confined spin. *Nature Physics*, 3:101–106, 2007.
- [64] Ilya Fushman, Dirk Englund, Andrei Faraon, Nick Stoltz, Pierre Petroff, and Jelena Vuckovic. Controlled Phase Shifts with a Single Quantum Dot. *Science*, 320(5877):769–772, 2008.
- [65] Romain Long, Tilo Steinmetz, Peter Hommelhoff, Wolfgang Hnsel, Theodor W. Hnsch, and Jakob Reichel. Magnetic microchip traps and single-atom detection. *Philosophical Transactions: Mathematical, Physical and Engineering Sciences*, 361(1808):pp. 1375–1389, 2003.
- [66] Jehyuk Rhee and Youngjoong Joo. Wide dynamic range cmos image sensor with pixel level adc. *Electronics Letters*, 39(4):360 – 361, feb. 2003.
- [67] Abbas El Gamal and Helmy Eltoukhy. CMOS Image Sensors. *IEEE Circuits and Devices*, pages 6–20, May/June 2005.
- [68] 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.
- [69] Eisuke Abe, Kohei M. Itoh, Junichi Isoya, and Satoshi Yamasaki. Electron-spin phase relaxation of phosphorus donors in nuclear-spin-enriched silicon. *Phys. Rev. B*, 70(3):033204, Jul 2004.
- [70] H. Y. Carr and E. M. Purcell. Effects of diffusion on free precession in nuclear magnetic resonance experiments. *Phys. Rev.*, 94(3):630–638, May 1954.
- [71] U. Haeberlen and J. S. Waugh. Coherent averaging effects in magnetic resonance. *Phys. Rev.*, 175(2):453–467, Nov 1968.
- [72] Götz S. Uhrig. Keeping a quantum bit alive by optimized  $\pi$ -pulse sequences. *Phys. Rev. Lett.*, 98(10):100504, Mar 2007.
- [73] 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, May 2004.
- [74] 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, September 2004.

- [75] Peter W. Shor. Scheme for reducing decoherence in quantum computer memory. *Phys. Rev. A*, 52(4):R2493–R2496, Oct 1995.
- [76] A. M. Steane. Error correcting codes in quantum theory. *Phys. Rev. Lett.*, 77(5):793–797, Jul 1996.
- [77] A. R. Calderbank and Peter W. Shor. Good quantum error-correcting codes exist. *Phys. Rev. A*, 54(2):1098–1105, Aug 1996.
- [78] D. Gottesman. *Stabilizer Codes and Quantum Error Correction*. PhD thesis, California Institute of Technology, Pasadena, CA, 1997.
- [79] Alexei Kitaev. Fault-tolerant quantum computation by anyons. *Preprint* arXiv:quant-ph/9707021., 1997.
- [80] Sergey B. Bravyi and Alexei Yu. Kitaev. Quantum codes on a lattice with boundary. *Preprint* arXiv:quant-ph/9811052., 1998.
- [81] Emanuel Knill, Raymond Laflamme, and Wojciech Zurek. Accuracy threshold for quantum computation. *Preprint* arXiv:quant-ph/9610011., 1996.
- [82] 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)*, pages 176–188, New York, NY, USA, 1997. ACM.
- [83] Robert Raussendorf and Jim Harrington. Fault-tolerant quantum computation with high threshold in two dimensions. *Phys. Rev. Lett.*, 98(19):190504, May 2007.
- [84] R Raussendorf, J Harrington, and K Goyal. Topological fault-tolerance in cluster state quantum computation. *New Journal of Physics*, 9(6):199, 2007.
- [85] Austin G. Fowler, David S. Wang, and Lloyd C. L. Hollenberg. Surface code quantum error correction incorporating accurate error propagation. *Preprint* arXiv:1004.0255v1., 2010.
- [86] David S. Wang, Austin G. Fowler, and Lloyd C. L. Hollenberg. Quantum computing with nearest neighbor interactions and error rates over 1%. *Preprint* arXiv:quant-ph/1009.3686., 2010.
- [87] Simon J Devitt, Austin G Fowler, Ashley M Stephens, Andrew D Greentree, Lloyd C L Hollenberg, William J Munro, and Kae Nemoto. Architectural design for a topological cluster state quantum computer. *New Journal of Physics*, 11(8):083032, 2009.
- [88] E. Knill. Quantum computing with realistically noisy devices. *Nature*, 434:39–44, 2005.
- [89] David P. DiVincenzo and Panos Aliferis. Effective fault-tolerant quantum computation with slow measurements. *Phys. Rev. Lett.*, 98(2):020501, Jan 2007.
- [90] Simon Anders and Hans J. Briegel. Fast simulation of stabilizer circuits using a graph-state representation. *Phys. Rev. A*, 73(2):022334, Feb 2006.
- [91] Christopher M. Dawson and Michael A. Nielsen. The solovay-kitaev algorithm. *Quantum Inf. Comput.*, 6:81, 2006.
- [92] Austin G. Fowler, Simon J. Devitt, and Lloyd C. L. Hollenberg. Implementation of shor’s algorithm on a linear nearest neighbour qubit array. *Quantum Information and Computation*, 4:237–251, 2004.
- [93] Christof Zalka. Shor’s algorithm with fewer pure qubits. *Preprint* arXiv:quant-ph/0601097., 2006.
- [94] Samuel A. Kutin. Shor’s algorithm on a nearest-neighbor machine. *Preprint* arXiv:quant-ph/0609001., 2006.
- [95] Rodney Van Meter and Kohei M. Itoh. Fast quantum modular exponentiation. *Phys. Rev. A*, 71(5):052320, May 2005.
- [96] Sergey Bravyi and Alexei Kitaev. Universal quantum computation with ideal clifford gates and noisy ancillas. *Phys. Rev. A*, 71(2):022316, Feb 2005.
- [97] Peter W Shor. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. *SIAM J. Comput.*, 26(5):1484–1509, 1997.