

A 2D Nearest-Neighbor Quantum Architecture for Factoring

Paul Pham
University of Washington
Quantum Theory Group
Box 352350, Seattle, WA 98195, USA,
ppham@cs.washington.edu,
<http://www.cs.washington.edu/homes/ppham/>

Krysta M. Svore
Microsoft Research
Quantum Architectures and Computation Group
One Microsoft Way, Redmond, WA 98052, USA
ksvore@microsoft.com,
<http://research.microsoft.com/en-us/people/ksvore/>

March 25, 2013

This document responds to comments by Referee 1, which were received on November 30, 2012. These comments are quoted and responded to below. We thank the Referee for the constructive feedback and suggestions.

1 General Comments

In this paper, the authors present in detail a new version of the modular exponentiation component of Shor's algorithm, with attention to the constraints of a 2-D planar graph of moderate degree for qubit connectivity. They use teleportation-based fanout to move qubits around within the machine. The modular addition uses a novel method, depending upon carry-save arithmetic and a small triangular lattice as its unit cell.

Perhaps the most novel part of the arithmetic is the approach to modular multiplication. In a traditional modular multiplication of two numbers x and y , the n^2 bit-wise products $x_i y_j$ are calculated and laid out in a trapezoid shifted to give the correct column (power of two), such that each entry is $2^i 2^j x_i y_j$, which of course can be represented by a single bit in the right place. Then the columns

are added, creating a $2n$ -bit number, which then must be further operated on to perform the modulo N operation. (N is assumed to be an n -bit number with the high-order bit being a one.)

In contrast, in this approach, each of the n^2 entries is a full n -bit number, $(2^i 2^j \bmod m) x_i y_j$. By adding those numbers using their modular circuit, the full $xy \bmod N$ value can be calculated directly. As proposed, this requires $O(n^3)$ bits (qubits) in the register. By combining the n^2 partial results in a log-depth tree structure, the depth for a modular multiplication becomes $O(\log n)$ times the constant depth of their 3-to-2 carry-save modular operation.

Up through section 5, I was prepared to recommend nearly immediate acceptance of the paper. The writing is clear and elegant and the technical work both valuable and polished. In sections 6 and 7, however, I have some doubts about the technical work, and the writing seems a bit more rushed, and the tail of the paper is not yet satisfactory.

Most importantly, it is disappointing that the authors have not produced a more complete estimate of the number of qubits (resources) required, as well as the actual circuit depth. A rough estimate, at least, should be very achievable given the level of detail already developed in the paper. The authors mention this as future work, but without it, the value of the paper is substantially diminished, and it does not seem unreasonable to expect it to appear here.

We thank the referee for the valuable feedback. More detailed numerical upper bounds for circuit resources are provided throughout the paper for each building block and the overall architecture. To do so, we have clarified our architectural model to include the notion of modules, which represent non-contiguous lattices (quantum computers) which can communicate and operate in parallel. This is described in the new Section 2 (Background).

2 Other significant technical comments:

The authors appear not to have noticed that half of the n^2 numbers in their multiplier require only a single bit. As long as $i + j \leq n$, the modulo operation results in the same number, allowing us to avoid using a full n -bit number: $(2^i 2^j \bmod m) x_i y_j = 2^i 2^j x_i y_j$.

Thank you for pointing this out, it is incorporated into the resource estimates for the modular multiplier, in the partial product creation section.

3 Section 2.2

“this ‘consumes’ the cat state” – this sentence is confusing.

We meant that the cat state remains entangled with the original source qubit and its fanned-out, entangled copies. We have since discovered via personal conversation with Aram Harrow and Dan Browne that this statement is no longer true. It is possible to “un-fanout” and therefore disentangle the cat state from the source and target qubits, allowing us to reuse this state. However, we have not been able to confirm the exact circuit at the time of this writing, so we have removed this sentence. Instead, we discuss an alternative method of unfanout, which is alternating teleports and CNOTs in a logarithmic-depth binary tree.

4 Section 3

Choi and Van Meter were not the first to consider 2-D architectures; most of the solid-state proposals and even some of the ion trap proposals worked with a 2-D layout. Kielpinski’s 2002 proposal might have been the first 2-D architecture. Working out exact algorithms on the structures came later, but papers by Kubiakowicz’s group and Chong/Oskin clearly included at least some level of work on the movement of qubits in a planar system, albeit with less attention to the abstraction of logical qubit connectivity.

We intended to say that Choi and Van Meter were the first to construct an adder explicitly optimized for a 2D architecture of qubits. We have reworded this paragraph to accurately reflect their contribution and added appropriate references.

I’m not sure a modular adder is “extended” to do modular exponentiation. “Composed” or “used”, perhaps?

We have reworded this sentence.

Your citation of Gossett has no year.

Corrected.

“all other factoring implementations”? That’s a rather broad characterization. Cleve and Watrous long ago proposed using a parallel reduction tree of multiplications before the QFT. Van Meter and Itoh investigated in detail the tradeoffs in resource consumption for doing this.

We have clarified this sentence by referring to all nearest-neighbor factoring implementations. By “implementations,” we mean a concrete mapping to an architecture, such as those given by [Fowler et al. 2004] and [Kutin 2006]. While we do acknowledge the Cleve and Watrous paper, which gives similar parallel results to those in the Kitaev-Shen-Vyalyi book, both assume arbitrary connectivity of qubits.

As far as I am aware, no one has worked out the details of Draper’s transform adder taking into account the need to do Solovay-Kitaev decomposition. This may add a very large factor to the execution time.

Agreed, we are not aware of anyone working out compilation of the Draper transform adder to a universal gate set, such as the Clifford group plus $\pi/8$ gates. It’s been shown that compilation will require $O(\log(1/\epsilon))$ overhead, where ϵ is the required precision. Throughout, we choose to avoid use of the QFT in our circuit due to the compilation required for each small rotation gate. We have added a sentence regarding this compilation cost.

Do you think the Zalka approximate multiplier approach actually works?

It is likely that something similar to the Zalka approximate multiplier actually works in practice for the majority of input values, if not the exact implementation described in the Zalka and Kutin papers. However, a rigorous theoretical argument, or empirical verification by simulation, has not yet appeared in the literature.

While it’s okay to include a ”forward reference” to a forthcoming paper of your own that carries more detail, you can’t ask us to ”refer to” it!

Agreed. Citation has been removed.

I don’t think the BKP ”exact” circuit guarantees a result, does it? Shor’s own original algorithm only probabilistically gives the correct answer, and that probability is a matter of some debate in the literature. Papers by Fowler (2004), Miquel (1996), Garcia-Mata (2007) and others provide different estimates.

Correct. This was a misleading statement and has since been corrected.

The journal style will ultimately dictate this, but when the bibliographic labels in the text are alphabetic, the references are usually ordered alphabetically.

We have updated the style.

5 Section 4

It would be worth pointing out that qubit 0 is the low-order qubit.

Corrected.

”At the level of bits, a CSA...” this sentence is awkward, reword.

Corrected.

Fig. 4 shows the layout, but it doesn't exactly match the circuit of Fig. 3. Having the exact circuit to accompany Fig. 4 would be useful.

The layout in Fig. 4 matches Fig. 3 if you look at the outgoing qubit values in Fig. 3. We have clarified this in the text.

6 Section 5

It feels somewhat like the phrasing on constant depth is a bit misleading in this section. Please reread and make sure it is easy for the reader to follow your claims.

Can you please provide specific examples of misleading phrasing?

Proof of Lemma 1: "O.ur" typo.

Corrected.

A little attention to classical versus quantum addends in this section might help the reader.

Unfortunately, due to time constraints, we were unable to incorporate this feedback into our revisions.

Fig. 5: labeling the lines themselves as "Layer 1", "Layer 2", etc. might help the reader.

This is a great improvement and has been incorporated.

"at no layer generates" – "no layer generates"?

Corrected.

Fig. 6: Swap the left and right ends of this figure to make it correspond to Fig. 5 as closely as possible, and point out this correspondence by also labeling the "layers" and the time axis here. "FANOUT RAIL" – of which variable(s)?

Unfortunately, due to time constraints, we were not able to incorporate this suggestion. The fanout rails are of the only qubits connected to them, namely, v_4 , u_4 , and v_5 , respectively.

My estimate here is that spatial resources are $29n$ qubits, temporal resources are $12n$ Toffoli gates. Is that about right?

This is approximately correct, and we have since added a table with a more accurate upper bound of the circuit resources.

7 Section 6

If you are worried that “quantum number” will mislead some readers, would “quantum integer” be better?

This is a great suggestion, and we have changed all mentions of “quantum number” to “quantum integer.”

Bottom of 6.1 is the place to expand the discussion to incorporate the above on the size of partial products and their impact on resource utilization.

A new section on partial product creation has been added, with more detailed discussion of the resources involved.

Your white tiles are hard to see on my printout.

The tile colors have since been removed, since the tiles must always be kept around during a modular multiplier’s operation. Therefore, the notion of colors and active status is no longer relevant.

Important: Your second and third rules under “black tiles” appear to conflict. Reword to be precise.

Thanks for this observation, please see previous response.

Bottom of Sec. 6: more discussion of the size and number of addends and resource consumption is needed.

More detailed numerical estimates are added in a new section at the end of Section 6. However, we think the number and size of addends for both serial and parallel multiplication are self-explanatory.

Fig. 10: In earlier figures, arrows were used only to indicate addend motion via teleportation, correct? A different symbol to indicate actual multiplication would make the figure clearer. Otherwise, asserting that an arrow itself is “log n depth” will confuse the reader. The tree approach used in this figure is not original, and in fact dates to the 2000 paper by Cleve and Watrous, at least.

We have clarified the use of arrows in the parallel modexp figure, as well as the modular multiple addition figure, where the arrows have a different meaning.

We do not claim that this tree structure is original, merely that it is now possible to do the necessary communication (teleportation and fan-out) in constant depth. That is, the tree structure before was a theoretical construction on an abstract architecture, and now based on the previous nearest-neighbor implementations of modular arithmetic, correspond to an actual physical tree structure if this architecture were to be fabricated and observed to run over time.

Sections 7, 8 and 9 should be extended with more discussion, and actual resource consumption figures. Some of the papers cited in Fig. 11 (which ought to be a table, not a figure) provide detailed estimates on resources (depth and width), as do Beckman et al. (1996, uncited here, but should be) and Van Meter and Itoh. Both of those latter papers offer several configurations that may make direct comparison tricky.

The requested discussion has been added with resource consumption, and the tables with the label “Figure” have been corrected. The citation to Beckman has been added.

Overall, as noted, this paper will be a valuable contribution to the literature once these issues are addressed.

We thank the referee for the valuable feedback.