

./sample.bib

On The Cost of Fault-Tolerant Shallow Circuits.

Michael Ben-Or David Ponomarev

April 24, 2025

Abstract

In this work we study the overall depth overhead cost required for constructing fault tolerance circuits. We focus on shallow depth circuits classes, In particular, \mathbf{QAC}_0 , $\mathbf{QNC}_{0,f}$ and \mathbf{QNC}_1 and certain known problem candidates for demonstrating quantum advantage such as factoring [?] and Instantaneous Quantum Polynomial-time [?], [?]. We only give a partial answers, Yet, clues that might pave the way towards a full understanding of the complexity versus fault tolerance trade-off.

1 Introduction.

The question about the feasibility of computation under noise is almost as ancient as the computer science field itself, initialized by Von Neumann [?] at the time that classical computation putted in debuts. Time been pass and the followed works had pointed that not even a polynomial computation in the presence of noise is still reasonable but one can implement a fault tolerance version at a most constant times cost at the circuits depth [?]. Or in asymptotic sense, classical computation in the presence of noise is as exactly hard as computation in ideal environment.

Once again, the feasibility question raised again, this time regarding quantum computing, and while an intensive work has been done, and also succeed to prove that polynomial quantum computation can be made fault tolerance, [?],[?] and even with only constant overhead at the original circuit width [?], the required depth over-head is till not well understood. We stress out that in all the familiar constructions, in construct to Pippenger [?], original constant-depth gates are mapped to asymptotically grow¹ depth gates.

This work address the above, We ask whether a magnitude depth overhead is an unavoidable price that one has to pay. And, in particular, whether an ideal \mathbf{QNC}_1 circuits can be computed in noisy- \mathbf{QNC}_1 circuits. We show how using the ideas presented in [?] and [?] gives almost immediately $\mathbf{QNC}_{0,f} \subset \text{noisy-}\mathbf{QNC}_1$ and that sampling from IQP [?] also can be done in logarithmic depth circuits.

¹Note, that here, classical computation is also counted in the overall depth cost

2 Notations.

In the following we present the notations used in the paper, readers who familiar with the literature of coding theory and quantum fault tolerance might skip sec:classical and sec:quantum and continue directly to sec:decoders which introduces less standard notations.

Once again, the feasibility question raised again, this time regarding quantum computing, and while an intensive work has been done, and also succeed to prove that polynomial quantum computation can be made fault tolerance, [?],[?] and even with only constant overhead at the original circuit width [?], the required depth over-head is till not well understood. We stress out that in all the familiar constructions, in construct to Pippenger [?], original constant-depth gates are mapped to asymptotically grow^a depth gates.

(a) location.

^aNote, that here, classical computation is also counted in the overall depth cost

2.1 Classical and Quantum Circuits.

location (i, j) of C . blue [COMMENT] Add here a figure of classical circuit, that demonstrates locations.

2.2 Classical Coding Theory.

The notation used in this paper follows standard conventions for coding theory. We use n to represent the length of the code, k for the code's dimension, and ρ for its rate. The minimum distance of the code will be denoted as d , and the relative distance, i.e., d/n , as δ . In this paper, n and k will sometimes refer to the number of physical and logical bits. Codes will be denoted by a capital C followed by either a subscript or superscript. When referring to multiple codes, we will use the above parameters as functions. For example, $\rho(C_1)$ represents the rate of the code C_1 .

Square brackets are used to present all these parameters compactly, and we use them as follows: $C = [n, k, d]$ to declare a code with the specified length, dimension, and distance. Any theorem, lemma, or claim that states a statement that is true in the asymptotic sense refers to a family of codes. The parity check matrix of the code will be denoted as H , with the rows of H representing the parity check equations. The generator matrix of the code will be denoted as G , with the rows of G representing the basis of codewords. The syndrome of a received

the right round, and to the second stage a left round. For the whole procedure, we will call a single correction round.

[H]

<p>(a) Single decoding round</p> $x \in F_2^n$ $\arg \min \{y \in C : y + x \}$ $d(y, C) < v \in V^+ \quad x'_v \leftarrow$ $\arg \min \{y \in C_0 : y + x _v\}$ $v \in V^- \quad x'_v \leftarrow$ $\arg \min \{y \in C_0 : y + x _v\}$ x	<p>every node=[draw,shape=circle]; node (v0) at (0,1) u_1; node (v6) at (0,3) u_2; node foreach in 1,2,3,4,5 (v) at (3,) v; (v0) - (v1) (v0) - (v2) (v0) - (v5) (v6) - (v3) (v6) - (v4) (v6) - (v1);</p> <p>(b) location.</p>
--	---

If the error is at wight less than βn then a single round of the majority reduce the error by at least constant fraction. Denote by $S^{(0)} \subset V^+$ and $T^{(0)} \subset V^-$ the subsets of left and right vertices adjacent to the error. And denote by $T^{(1)} \subset T^{(0)}$ the right vertices such any of them is connect by at least $\frac{1}{2}\delta_0\Delta$ edges to vertices at $S^{(0)}$. Note that that any vertex in $V^-/T^{(1)}$ has on his local view less than $\frac{1}{2}\delta_0\Delta$ faulty bits, So it corrects into his right local view in the first right correction round. Therefore after the right correction round the error is set only on $T^{(1)}$'s neighbourhood, namely at size at most $\Delta|T^{(1)}|$. We will show that this amount is strictly lower by a constant factor than $|e|$.

First, let's use the expansion property (def:mix) for getting an upper bound on

$$T^{(1)} \text{ size: } 1 - \frac{2\delta_0\Delta|T^{(1)}| \leq \Delta \frac{|T^{(1)}||S^{(0)}|}{n} + \lambda\sqrt{|T^{(1)}||S^{(0)}|} \left(\frac{1}{2}\delta_0\Delta - \frac{|S^{(0)}|}{n}\Delta \right) |T^{(1)}| \leq \lambda\sqrt{|T^{(1)}||S^{(0)}|} |T^{(1)}| \leq \left(\frac{1}{2}\delta_0\Delta - \frac{|S^{(0)}|}{n}\Delta \right)^{-2} \lambda^2 |e|$$

Since any left vertex adjoins to at most Δ faulty bits we have that $\Delta|S^{(0)}| \leq |e|$.

Combing with the inequality above we get:

$$\Delta|T^{(1)}| \leq \left(\frac{1}{2}\delta_0\Delta - \frac{|e|}{n} \right)^{-2} \lambda^2 |e| \text{ Hence for } |e|/n \leq \beta = \frac{1}{2}\delta_0\Delta - \sqrt{2\lambda} \text{ it holds that } \Delta|T^{(1)}| \leq \frac{1}{2}|e|. \text{ Namely the error is reduced by half.}$$

2.2.2 Pippenger Construction.

The main insight behind Pippenger's construction is that against non trivial noise there is no point in decoding completing back to the code space since it's likely that in the follow turn the block will absorb a magnificent error. So instead of full decoding, one can perform a partial decoding round, for example, a single round of local majority. Keeping reducing the error will ensure that he distance of the state, with hight probability, will remain close (enough) to the ideal state along an ideal computation.

In the original construction, any bit encoded using the repetition code, notice that the repetition code can be thought as expander code on the bipartite Δ -regular graph where the local code C_0 is the repetition code on Δ bits. Then one can use lemma:reduce to show that a single round of decoding reduces the error by half.

Encode any bit of the original code by n bits using the repetition code. Replace any of the AND, and the OR gate by their transversal gates. In each even turn apply a single round of the majority decoder. Finally at the end apply the decoder.

Once again, the feasibility question raised again, this time regarding quantum computing, and while an intensive work has been done, and also succeed to prove that polynomial quantum computation can be made fault tolerance, [?],[?] and even with only constant overhead at the original circuit width [?], the required depth overhead is till not well understood. We stress out that in all the familiar constructions, in construct to Pippenger [?], original constant-depth gates are mapped to asymptotically grow^a depth gates.

(a) location.

^aNote, that here, classical computation is also counted in the overall depth cost

2.3 Quantum Codes.

A quantum code over n qubits is an embedding of $\mathcal{H}_2^{\otimes k}$ as a subspace of $\mathcal{H}_2^{\otimes n}$. Similar to classical codes, we will call n and k the physical and logical qubits. The embeddings of states in $\mathcal{H}_2^{\otimes k}$ are called codewords or encoded states. In addition, we will use the term "logical operator" (i.e. logical X_i) to describe an operator that acts on the code space exactly as it would act on the logical space $\mathcal{H}_2^{\otimes k}$ (in our example, turning on and off the encoded state corresponds to the i th qubit exactly as X_i acts as Pauli X on the i th qubit in $\mathcal{H}_2^{\otimes k}$).

We will denote by X and Z the single X and Z Pauli operators, by X_i the application of X on the i th qubit and nothing else (identity) on the rest of the qubits. By $X^{(v)}$ for some $v \in F_2^n$, we mean the operator composed by applying X on each of the qubits whose index is a non-trivial coordinate of v and identity elsewhere. In a similar fashion, we define $Z^{(v)}$. When the context is clear, we will allow ourselves to omit the brackets, i.e. Z^v . The weight of a Pauli operator is the number of coordinates on which the operator acts non-trivially. Recall that the set of Pauli $+I$ spans all the Hermitian matrices. We say that the Pauli weight of an operator is the maximal weight of a Pauli in its Pauli decomposition. For example, consider the operator $A = IXX + ZII$, the weight of A is 2.

The distance of a quantum code is the minimal weight of an operator that takes one codeword to another. We use the standard bracket notation to describe quantum states and in addition, we define for a vector space $A \subset F_2^n$ the notation A to represent the uniform superposition of all the vectors belonging to that space, namely: $A = \frac{1}{\sqrt{|A|}} \sum_{x \in A} x$. We define in the same way the notation to hold for affine spaces, $x + A$. We will use \propto to denote a quantum states up to normalization factor, for example $\psi \propto 0 + 1$ means that $\psi = \frac{1}{\sqrt{2}}(0 + 1)$. A CSS code is a quantum code defined by a pair of classical codes C_X and C_Z ,

satisfying $C_Z^\perp \subset C_X$, such that any codeword of it has the form $x + C_Z^\perp$, where $x \in C_X$. We will use Q to refer to a CSS code in general and use C_X/C_Z^\perp to refer to the vectors associated with the X -generators or the encoded states in the computational basis. In the same way, C_Z/C_X^\perp refers to the Q in the phase basis. We will say that a CSS code Q is a LDPC if C_X and C_Z are both LDPC codes. Our construction uses the classical Tanner code [?], the expander codes [?], and Hyperproduct code (quantum expanders) [?], [?], [?]. We will not describe these constructions and refer the reader to those papers for further information.

2.4 Decoders.

We denote by C_g the good qLDPC code [?] [?] [?], and by C_{ft} the concatenation code presented at [?] (ft stands for fault tolerance). For a code C_y , we use Φ_y, E_y, D_y to denote the channel maps circuits into the their matched circuits compute in the code space, the encoder, and the decoder, respectively. We use Φ_U to denote the 'Bell'-state storing the gate U . We say that a state ψ is at a distance d from a quantum code C if there exists an operator U that sends ψ into C such that U is spanned on Paulis with a degree of at most d . Sometimes, when the code being used is clear from the context, we will say that a block B of qubits has absorbed at most d noise if the state encoded on B is at a distance of at most d from that code.

3 Todo:

1. Move to encoding each qubit by logarithmic width (instead of chunks) the reason is that the gate teleportation becomes complicated when it applied over higher dimension.
2. Then showing for 2-qubit gates set that is indeed works.
3. Treating separately to noise observed in two qubits gates.

4 Fault tolerance Toffoli.

blue [COMMENT] In that section the \cdot operation is the pair wise product (pair wise AND).

Assume that $\bar{0}, \bar{1} \in C_X$ and that they belong to two different cosets of C_X/C_Z^\perp . Let $x, y \in \{\bar{0}, \bar{1}\}$.

$$\sum_{z, z', w \in C_Z^\perp} z z' w \sum_{z, z', w \in C_Z^\perp} z z' w + z \cdot z' \sum_{z, z', w \in C_Z^\perp} z + x z' + y w + z \cdot z' \sum_{z, z', w \in C_Z^\perp} z + x z' + y x \cdot y + x \cdot z' + y \cdot z + \quad (1)$$

Since $x, y \in \{\bar{0}, \bar{1}\}$ we have that $x \cdot z'$ equals to either z' or $\bar{0}$. Hence $\sum_{w \in C_Z^\perp} \xi + x \cdot z + w = \sum_{w \in C_Z^\perp} \xi + w$. So the idea is the following, suppose that one has to compute Toffoli at time t over the registers R_1, R_2, R_3 . First, at time 0, he initialize a logical zero C_Z^\perp in each register, then he compute pairwise Toffoli R_1, R_2 into R_3 . That gives the ket $\sum_{z, z', w \in C_Z^\perp} z \cdot z' + w$, immediately afterwards encode R_3 again into a good quantum code. Denote by τ the time required for decoding R_3

back, at time $t - \tau$ start to decode R_3 . Eventually at time t compute again the transversal Toffoli, by equation:toff we gets the desired.

By similar arguments exhibited at claim:noisepa one can show that the errors behaves according to a Pauli noise channel. blue [COMMENT] That is not correct, since the concatenation construction assumes that all the registers initialized to physical zeros in the begging of the computation.

4.1 Another Idea, $z \cdot z'$ cann't contribute too mach.

Clearly we have that $|z \cdot z'| \leq |z|, |z'|$ therefore we have that $\Pr_{z, z' \in C_Z^\perp} [|z \cdot z'| \geq t] \leq \Pr_{z \in C_Z^\perp} [|z| \geq t]$. Now assume that the tanner code by which the code defined is bipartite graph and denote by z_+, z_- the grouping of the z 's generators supported on the even and the odd vertices of the graph. By triangle inequality $|z| = |z_+ + z_-| \leq |z_+| + |z_-|$, So if $|z| > t$ then at least one of $|z_-|, |z_+|$ is greater than $t/2$. Hence via the union bound: $\Pr_{z \in C_Z^\perp} [|z|] \leq \Pr_{z \in C_Z^\perp} [\bigcup_{i \in \pm} |z_i| \geq t/2] \leq \sum_{i \in \pm} \Pr_{z \in C_Z^\perp} [|z_i| \geq t/2]$

Since any two positive (negative) generators are disjoint we have that $|z_+|$ is a sum of the independent random variables each stands for the weight contributed by a positive vertex. Let us denote by V^+, V^- the positive and the negative vertices and for each vertex $v \in V$ we will denote by z_v the bits of z restricted to v edges. So $|z_\pm| = \sum_{v \in V^\pm} |z_v|$. For simplicity assume that $|V^+| = |V^-| = n/2$ and that $\mathbf{E}_{z \in C_A \otimes C_B} [|z|] = \mu$. Then we can use concentration inequality to have: $\Pr_{z \in C_Z^\perp} [|z|] \leq \sum_{i \in \pm} \Pr_{z \in C_Z^\perp} [\sum_{v \in V^i} |z_v| \geq t/2] \leq 2e^{-(\mu - \frac{t}{2})n}$ Thus if $\mu - \gamma \geq O(1)$ (from claim:error) then with high probability the Toffoli is computed up to reducible error.

4.2 Using Polynomials Codes.

Consider the CSS code above F_q defined by the stabilizers $C_Z^\perp = \{x^{2i} : i \geq \frac{1}{2}d - c\}$ and $C_X^\perp = \{x^{2i+1} : i \geq \frac{1}{2}d - c\}$ for some $c = O(1)$. Clearly the X -stabilizers commute with the Z -stabilizers.

5 The Noise Model

Informally classical noisy circuits describe the running computation of circuits when the bits have probabilities to flip. As exactly to the classical case, in noisy quantum circuits qubits have probabilities to fault. We formalise the noise model by defining a channel $\mathcal{N} : \mathcal{C} \rightarrow \mathcal{D}(\mathcal{C})$ that given an ideal circuit induce distribution over circuits. For example, one can consider a Pauli channel, which after each gate of the original circuit, either do nothing with probability $1 - p$ or, with probability $1 - p$ impose uniformly one of the Pauli operators X, Z, Y . Formally: Pauli channel $\mathcal{N} : \mathcal{C} \rightarrow \mathcal{D}(\mathcal{C})$ defined to give on input $C \in \mathcal{C}$ the distribution over circuits \tilde{C} where any even location of $(i, 2j)$ of \tilde{C} equals to the (i, j) location of C , and any odd location $(i, 2j + 1)$ of \tilde{C} is the density operator $(1 - p)I + \frac{p}{3}(X + Y + Z)$.

The Pauli channel is characterized by exhibits an independent noise on the qubits, Yet for most of the fault tolerance construction a much more weaker property is required to be assumed. We say that a channel is a local stochastic noise

channel if the probability to error to be occur is exponentially decays at the number of qubits the error supports.

An error channel $\mathcal{N} : \mathcal{C} \rightarrow \mathcal{D}(\mathcal{C})$ will be said to be a local stochastic noise channel if there exists a constant c such the probability to a fault to be applied on locations (I, j) , where I is a subset of qubits, is less than c^{-n} .

Another important property of a noise model which we consider in this work is the accessibility to fresh qubits, also known as resets gate. When having an access to fresh qubits one can assume that in any time in the computation there are qubits at the 0 states. Usually those qubits are used to measured the syndrome relative to an error correction code. It was proven that without an access to fresh qubits quantum circuits cannot last than logarithmic depth without mixing into a fully mixed state, meaning to be turned into complete garbage [?]. That result also holds for a classical noisy computation.

An error channel $\mathcal{N} : \mathcal{C} \rightarrow \mathcal{D}(\mathcal{C})$ will be said to has a fresh qubits access if location (i, j) in an output gate \tilde{C} has a non zero probability to exhibits a fault if there is a $j' < j$ such a location (i, j') such that on the input circuit C , at location (i, j') a non identity gate is posed.

We close this section by formalize the noisy-QNC₁ class. We denote by noisy-QNC₁ the class of decision problems solvable by logarithmic-depth quantum circuits, subjected to a local stochastic noise, with bounded probability of error. We mention that in [?], it was proved how a fault tolerance circuit with an access to fresh qubits, at logarithmic depth, can be converted to a log depth circuits without a fresh qubits access at the cost which is at most polynomial in wide. Meaning that Proving that QNC₁ \subset noisy-QNC₁ implies also that QNC₁ can be computed, in the presence of noise without an access to fresh qubits.

6 Fault Tolerance (With Resets gates) at Linear Depth.

There exists a value $p_{th} \in (0, 1)$ such that if $p < p_{th}$, then any quantum circuit C with a depth of D and a width of W can be computed by a p -noisy circuit C' , which allows for resets. The depth of C' is at most $\max \{O(D), O(\log(WD))\}$.

6.1 Initializing Magic for Teleportation gates and encodes ancillaries.

The Protocol:

1. Initialization of zeros: The qubits are divided into blocks of size $|B|$. Each block is encoded in C_g using $D_{ft}\Phi_{ft}[E_g]0^{|B|}$.
2. Initialization of Magic for Teleportation gates: The gates in the original circuit are encoded in C_g using $D_{ft}\Phi_{ft}[E_g]\Phi_U$.
3. Gate teleportation: Each gate in the original circuit is replaced by a gate teleportation.
4. Error reduction: After the initialization step, at each time tick, each block runs a single round of error reduction.

[From [?]] Assuming that an error $|e| \leq \gamma n$, i.e e is supported on less than γn bits, then a single correction round reduce e to an error e' such that $|e'| <$

$\nu|e|$. The gate $D_{ft}\Phi_{ft}[E_g]$ initializes states encoded in C_g subject to a $3p$ -noise channel. Clearly, with high probability, $\Phi_{ft}[E_g]$ successfully encodes into $C_{ft} \circ C_g$, let's say with probability $1 - \frac{1}{\text{poly}(n)}$. Denote by E_i and D_i the encoder and decoder at the i th level of the concatenation construction. Consider the decoder under \mathcal{N} action: $P_2D_1P_2D_2, \dots, P_{i-1}D_iP_i$, by the fault-tolerance construction, a logical error at the i th stage occurs with probability p^{2^i} . Therefore, by the union bound, the probability that in one of the steps the circuit absorbs an error that is not corrected is less than $p + p^2 + p^4 + \dots < 2p$. Hence, any decoded qubit absorbs noise with probability less than $2p$.

Thus, overall, we can bound the probability of a single qubit being faulty by: $\Pr[fault] = \Pr[fault|\Phi_{ft}[E_g]] \cdot \Pr[\Phi_{ft}[E_g]] + \Pr[fault|\overline{\Phi_{ft}[E_g]}] \cdot \Pr[\overline{\Phi_{ft}[E_g]}]$
 $\leq \Pr[fault|\Phi_{ft}[E_g]] + \Pr[\overline{\Phi_{ft}[E_g]}] \leq 2p + \frac{1}{\text{poly}(n)} \leq 3p$

In our construction, we use the concatenation code to encode blocks of length $\log(n)$. Therefore, any $\text{poly}(n)$ in the above should be replaced by $\log(n)$. However, this does not affect anything since the inequality does not depend on n .

With a probability $1 - \frac{WD}{|B|} \cdot D2e^{-2|B|(\beta-p)}$, the total amount of noise absorbed in a block at any given time t , is less than γn . Consider the i th block, denoted by B_i . By applying Hoeffding's inequality, we have that the probability that more than $\beta|B|$ qubits are flipped at time t is less than $2e^{-2|B|(\beta-p)}$. By using the union bound over all blocks at all time locations, we can conclude that with probability $1 - \frac{WD}{|B|} \cdot D2e^{-2|B|(\beta-p)}$, the noise absorbed in a block is less than $|\beta|B$ for the entire computation.

Let X_t denote the support size of the error over B_i at time t . Using claim:error, we can bound the total amount of error absorbed by a block until time t as follows: $X_t \leq \nu \cdot (X_{t-1} + \beta|B|) \leq \nu(\gamma + \beta)|B| \leq \gamma|B|$

The total depth of the circuit is $O(D) + O(\log^c |B|)$. The gate for encoding $|B|$ -length blocks in C_g is a Clifford gate and can therefore be computed in $O(\log |B|)$ depth. The encoding of the magic/bell states is done by first computing them in the logical space (un-encoded qubits) and then encode them using the encoder. Hence, the fault-tolerant version of both initializing ancillaries and magic states/bell states costs $O((\log |B|) \cdot \log^c(|B| \log |B|))^2$ depth [?]. Backing into C_g from C_{ft} by decoding the concatenation code takes exactly as long as the encoding, namely $O((\log |B|) \cdot \log^c(|B| \log |B|))$.

Then, using the bell measurements, any of the logical gates takes $O(1)$ depth. Since we only perform a single round of error correction, the remaining computation until the last decoding stage takes at most constant time of the original depth. Finally, we pay $O(\log |B|)$ for complete decoding. Summing all, we get: $O(\log |B| \cdot \log^c(|B| \log |B|)) + O(D) + O(\log |B|)$
 $= O(D) + O(\log^c |B|)$

Assuming that W is polynomial in D , taking the block length to be $|B| = \log((W \cdot D)^c)$, as shown in claim:prob, results in a linear fault tolerance construction with a success probability of $1 - \frac{1}{\log^{c^2}(W \cdot D)}$. This means that the fault tolerance version of circuits in QNC_1 has a logarithmic depth. Additionally, using the construction in [?] produces a polynomial fault tolerance circuit in the reversible gates setting. blue [COMMENT] We missed the fact that it requires non trivial classical computation to compute what gate should be applied after

²The width of the original circuit is $|B|^2$ so the number of locations is $|B|^2 \cdot \log |B|$

the gate teleportation (i.e UPU^\dagger).

7 Does $\text{NC}_1 \subset \text{noisy-QNC}_1$?

8 Does Factoring $\subset \text{noisy-QNC}_1$?

$$\begin{aligned} D(n) &= \Theta(\log n) + D(\sqrt{n}) \\ \Rightarrow D(n) &= \Theta(\log n) \end{aligned}$$