# The **FL**uid **A**llocation of **S**urface code **Q**ubits model for early fault-tolerant resource estimation
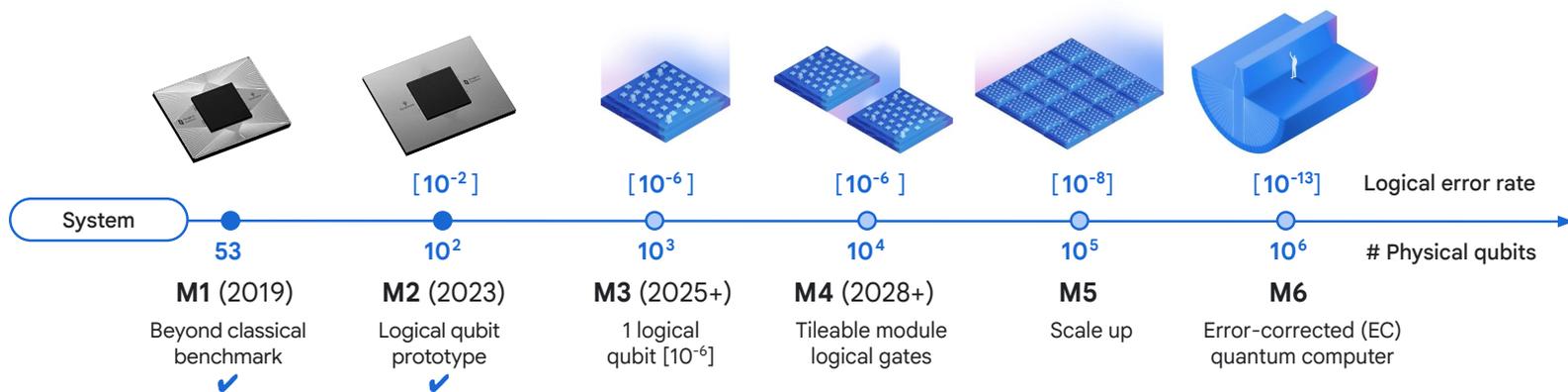
**William J. Huggins**, Tanuj Khattar, Amanda Xu, Matthew Harrigan, Christopher Kang, Guang Hao Low, Dmitri Maslov, Austin Fowler, Nicholas C. Rubin, Ryan Babbush

Google

What problem are we trying to solve?

# Understanding the transition to fault-tolerance



| System | $[10^{-2}]$ | $[10^{-6}]$ | $[10^{-6}]$ | $[10^{-8}]$ | $[10^{-13}]$ | Logical error rate |
|---|---|---|---|---|---|---|
| 53 | $10^2$ | $10^3$ | $10^4$ | $10^5$ | $10^6$ | # Physical qubits |

**M1** (2019)
Beyond classical benchmark ✔

**M2** (2023)
Logical qubit prototype ✔

**M3** (2025+)
1 logical qubit $[10^{-6}]$

**M4** (2028+)
Tileable module logical gates

**M5**
Scale up

**M6**
Error-corrected (EC) quantum computer

Quantum AI

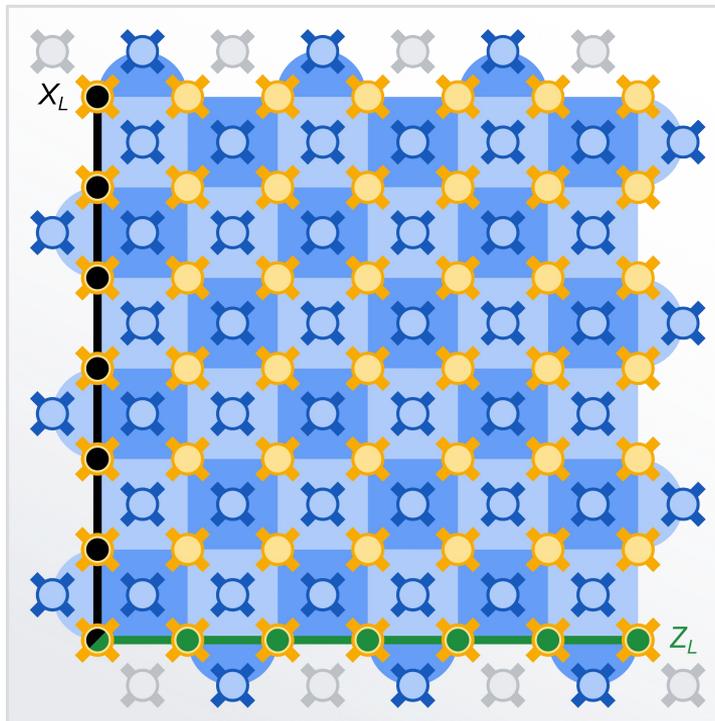# Understanding the transition to fault-tolerance

# Ingredients

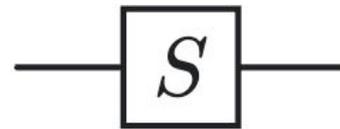# Logical qubits in the surface code



Physical qubit count $= 2 * (d+1)^2$
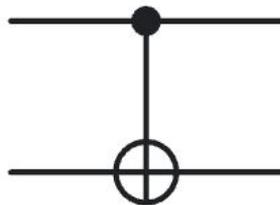
Logical Error Rate $\approx 0.1 * \Lambda^{-(d+1)/2}$

Surface code error correction gives exponentially more time for quadratically more space.
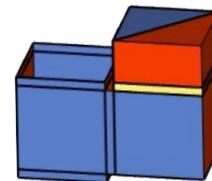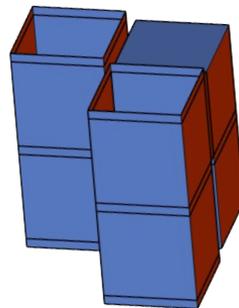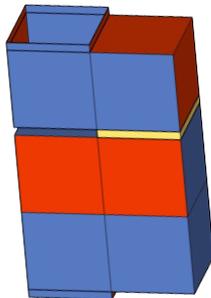
2D local connectivity is sufficient to store information and perform gates using lattice surgery

# Everything* requires ancilla spacetime

Gates in NISQ

Gates in the surface code

# Walking surface codes

With limited qubits, we expect to have $N_{ancilla} \ll N_{data}$

We will want to move ancilla around

Walking surface codes let us move logical qubits

Whole patches can move together



McEwen, M., Bacon, D. & Gidney, C. Relaxing hardware requirements for surface code circuits using time-dynamics. *Quantum* **7**, 1172 (2023).

# Walking surface codes in action

# Walking surface codes in action

# Walking surface codes in action

# Walking surface codes in action



**Time 0** → Time 1

# Walking surface codes in action



Time 0 → **Time 1**

# Walking surface codes in action



**Time 1** → Time 2

# Walking surface codes in action



**Time 1** → Time 2

# Walking surface codes in action



Time 1 → **Time 2**

# Walking surface codes in action



**Time 2** → Time 3

# Walking surface codes in action



**Time 2** → Time 3

# Walking surface codes in action



**Time 2** → Time 3

# Walking surface codes in action



**Time 2** → Time 3

# Walking surface codes in action



Time 2 → **Time 3**

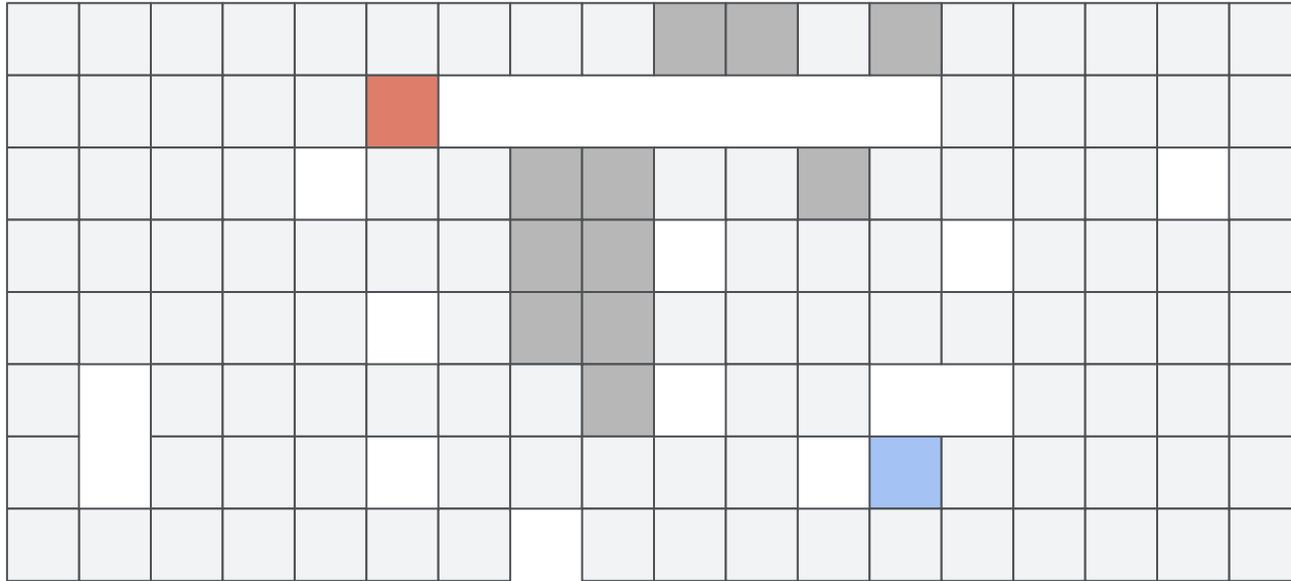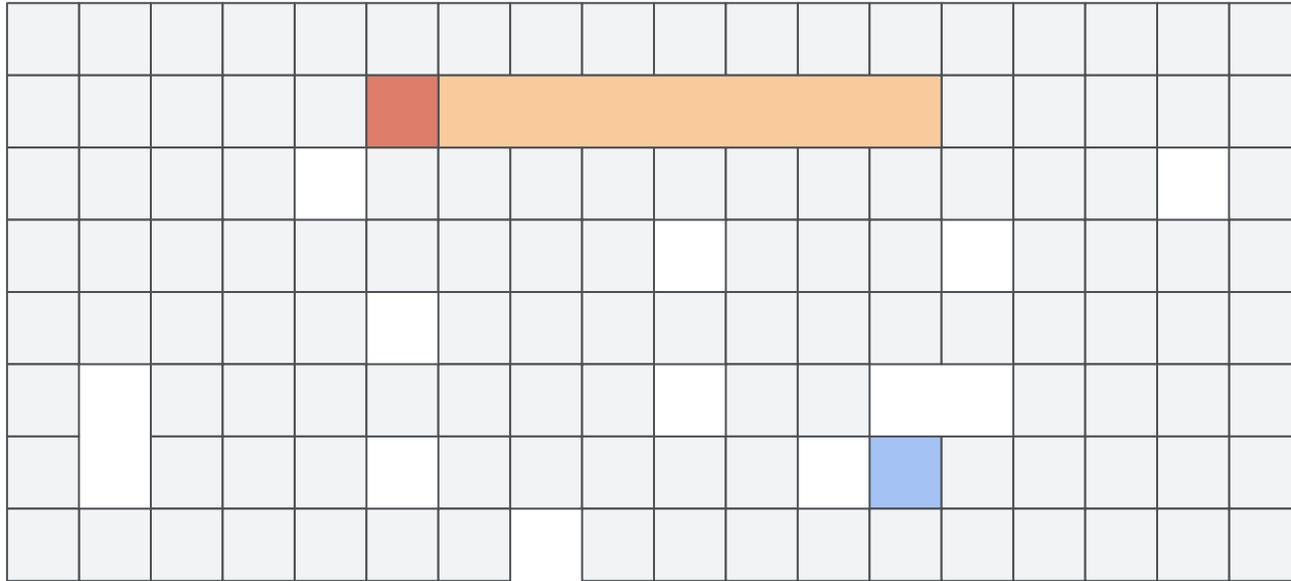# Magic state cultivation



End-to-End Cultivation

Magic state cultivation prepares a resource state we can use to perform a T gate.

The key ideas:
- Check the T state in a low-distance color code
- Rapidly grow to a larger code and postselect
- Low but finite error rate with much lower spacetime volume than previous techniques



Fowler Devitt 2013
"A bridge to lower overhead quantum computation"
$\varepsilon \approx 2e\text{-}12$

Gidney Fowler 2019
"Efficient magic state factories with a catalyzed CCZ→2T transformation"
$\varepsilon \approx 3e\text{-}11$

Fowler Gidney 2018
"Low overhead quantum computation using lattice surgery"
$\varepsilon \approx 9e\text{-}17$

This paper
$\varepsilon \approx 4e\text{-}6$

This paper
$\varepsilon \approx 2e\text{-}9$

# Magic state cultivation



End-to-End Cultivation

Magic state cultivation prepares a resource state we can use to perform a T gate.

The key ideas:
- Check the T state in a low-distance color code
- Rapidly grow to a larger code and postselect
- Low but finite error rate with much lower spacetime volume than previous techniques

## Magic state cultivation: growing T states as cheap as CNOT gates

Craig Gidney, Noah Shutty, and Cody Jones

Google Quantum AI, California, USA
September 27, 2024

We refine ideas from [KLZ96; JBH16; CN20; Bom+24; GJ23; Gid+23; BH F24] to efficiently prepare good $|T\rangle$ states. We call our construction "magic state cultivation"

# Magic state cultivation

We use simulations to determine the expected spacetime volume given a physical error rate and a target logical error rate

We optimize by varying:
- The amount of postselection
- The choice of color code distance

(Performed using the code and methodology of Gidney *et al.*)

# The model

# Fluid allocation of surface code qubits



A simple cartoon

- Space is horizontal and time is vertical
- Data qubits are orange, fluid ancilla are blue

We make sure we have enough ancilla volume

$$\text{TOTAL ANCILLA VOLUME} = \sum_{g} \text{ANCILLA VOLUME}(g)$$

$$\text{DEPTH} \geq \frac{\text{TOTAL ANCILLA VOLUME}}{\text{NUMBER OF FLUID ANCILLA}}$$

With more fluid ancilla…

- The ancilla volume (blue) is conserved
- The overall spacetime volume is reduced (orange and blue)

# Gate costs in the FLASQ model

| Basic gates | FLASQ ancilla volume | Measurement depth | Notes |
|---|---|---|---|
| $X$ / $Y$ / $Z$ | $0$ | 0 | Implemented in software |
| $X$ / $Z$ basis measurement (or initialization) | $0$ | 0 | |
| $H$ | $7$ | 0 | Includes the cost of a patch rotation |
| $S$ / $S^{\dagger}$ | $5.5$ | 0 | |
| $T$ / $T^{\dagger}$ | $1.5v(p_{phys}, p_{cult}) + t_{react} + 6$ | 1 | Depends on physical error rate, $p_{phys}$, and target logical error rate, $p_{cult}$ |
| Move | $5p(q_1, q_2)$ | 0 | Moves a qubit to an empty patch |
| $CNOT$ / $CZ$ | $5p(q_1, q_2)$ | 0 | |

# The reaction limit

Most operations can be rearranged freely in space and time

Some measurements have to be performed serially

- The choice of measurement basis depends on earlier measurement outcomes
- So the control software has to catch up before the measurement can be made

This is the "reaction limit"

- We often assume 10 microseconds
- We don't expect this to be the limiting factor in early fault-tolerance

Just in case, we ensure that

$$\text{DEPTH} \geq \text{REACTION TIME} \times \text{MEASUREMENT DEPTH}$$

Quantum AI

# Applications

Quantum AI

# Ising model time dynamics

## Time evolution

$$H = -J \sum_{\langle i,j \rangle} Z_i Z_j + g \sum_i X_i$$

## Estimate correlation functions

- Time-evolve after a high-entanglement quench
- Measure the expectation value of $Z_{tot}^2 = \frac{1}{N^2} \sum_{j,k} Z_j Z_k$
- Target an absolute error $\leq 0.01$ with high probability

## Classically challenging

- Canonical benchmark for quantum many-body physics
- Recent tensor network heuristics (Mandra 2025) struggle to perfectly converge and bound errors at an $11 \times 11$ system size

## FLASQ lets us estimate the volume ($V$)

## Phenomenological error model

- Surface code suppression: $\Lambda = 0.01/p_{phys}$
- Logical error rate: $p_{cyc} \approx 0.03 \Lambda^{-(d+1)/2}$
- Magic state error: $p_{mag}$ from cultivation

## Probabilistic error cancellation

- We can mitigate residual errors to get an unbiased estimator
- Sampling overhead: $\Gamma^2 \approx \exp\left(4 p_{cyc} d S_{cliff} + 4 p_{mag} M\right)$

Quantum AI

# Ising model resource estimates



Time to Solution for an 11 × 11 Ising Model

## Time-to-solution

- 20 second-order Trotter steps
- Optimized over choice of code distance and cultivation parameters
- Assumes target standard deviation of $\sigma = 0.0045$

## Cost drivers

- Large number of samples for target precision
- $\approx 7,400$ arbitrary rotations

# Optimal code distance



Optimal Code Distance for Ising Simulations (Conservative FLASQ Estimates)

# Validating against a hand compilation



| | 4 rotation synthesis units | | | |
|---|---|---|---|---|

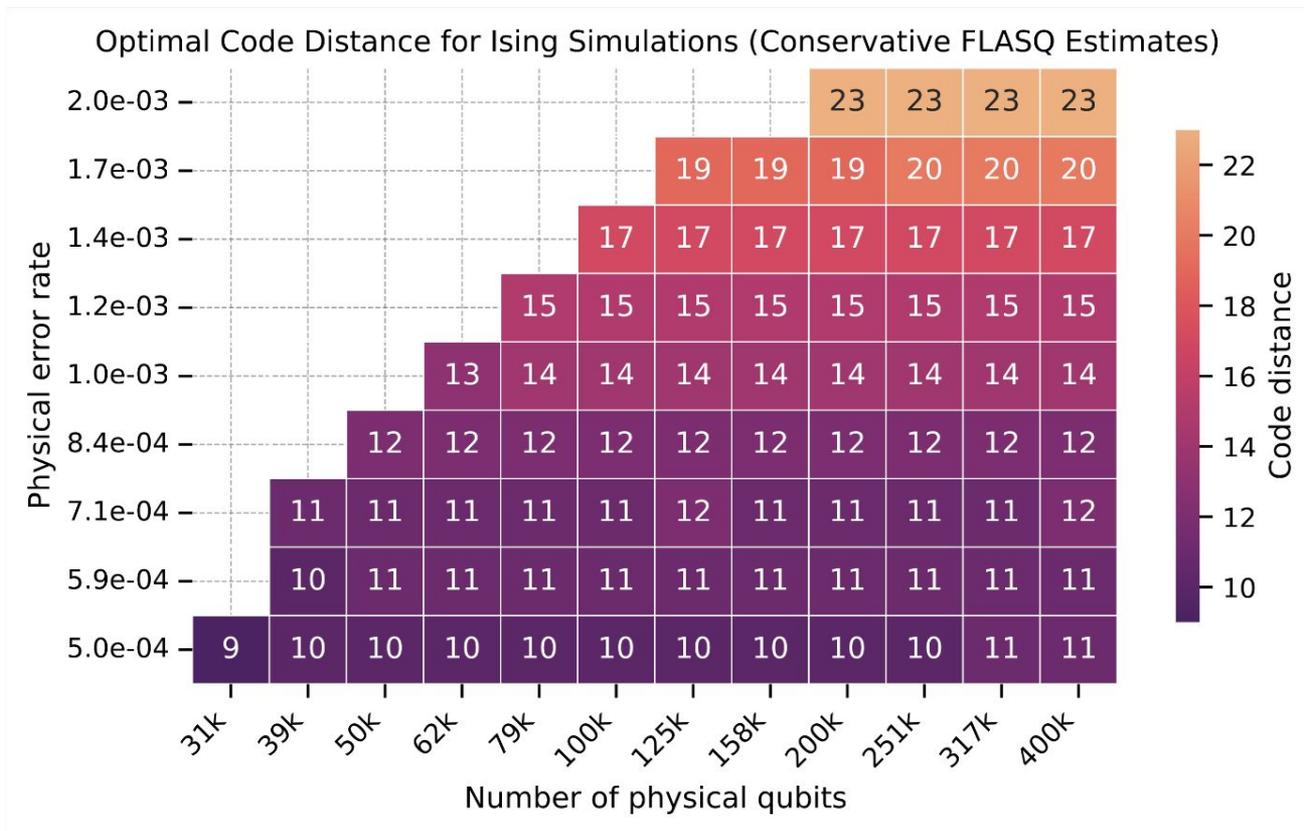| | | | | | |
|---|---|---|---|---|---|
| 1 | | | | | 104 |
| 2 | 19 | 36 | 53 | 70 | 87 105 |
| 3 | 20 | 37 | 54 | 71 | 88 106 |
| 4 | 21 | 38 | 55 | 72 | 89 107 |
| 5 | 22 | 39 | 56 | 73 | 90 108 |
| 6 | 23 | 40 | 57 | 74 | 91 109 |
| 7 | 24 | 41 | 58 | 75 | 92 110 |
| 8 | 25 | 42 | 59 | 76 | 93 111 |
| 9 | 26 | 43 | 60 | 77 | 94 112 |
| 10 | 27 | 44 | 61 | 78 | 95 113 |
| 11 | 28 | 45 | 62 | 79 | 96 114 |
| 12 | 29 | 46 | 63 | 80 | 97 115 |
| 13 | 30 | 47 | 64 | 81 | 98 116 |
| 14 | 31 | 48 | 65 | 82 | 99 117 |
| 15 | 32 | 49 | 66 | 83 | 100 118 |
| 16 | 33 | 50 | 67 | 84 | 101 119 |
| 17 | 34 | 51 | 68 | 85 | 102 120 |
| 18 | 35 | 52 | 69 | 86 | 103 121 |

20 rows

8 columns

## Static layout

- We couple the data qubits (orange) to the rotation synthesis areas (green)
- We use walking surface codes to shift the red access hallway

## Rotation synthesis gadget

- Each $2 \times 2$ region can implement an arbitrary rotation
- They consume one T state every four logical timesteps

# Validating against a hand compilation



## Static layout

- We couple the data qubits (orange) to the rotation synthesis areas (green)
- We use walking surface codes to shift the red access hallway

## Rotation synthesis gadget

- Each $2 \times 2$ region can implement an arbitrary rotation
- They consume one T state every four logical timesteps

## Comparison with FLASQ estimates

| Logical timesteps (FLASQ) | Logical timesteps (by hand) | Ratio (FLASQ / hand) |
|---|---|---|
| 55995 | 73810 | 0.76 |

Probing the crossover between NISQ and fault-tolerance

# Two modes of operation for a quantum processor

## NISQ Mode

- Treat the device as a large, noisy processor
- Run many parallel copies of the simulation using all available space

## Fault-Tolerant Mode

- Allocate all physical qubits to a single, error-corrected simulation
- (Does not account for the possibility of Heisenberg-limited measurement)

## Simple NISQ error model

- Single-qubit depolarizing noise after every two-qubit gate
- Allows for arbitrary rotations in a single timestep
- Mitigated using Probabilistic Error Cancellation (PEC)

## Phenomenological FT error model

- Rotations compiled into Clifford + T gate set
- Same as other Ising simulation (single-qubit X and Z errors)
- Use PEC to mitigate residual errors

# Moderately deep circuits (20 Trotter Steps)



Log₁₀ Runtime Ratio (Fault-Tolerant / NISQ) for Ising Simulation

## Reading the heatmap

- We plot the log-ratio of the runtimes: $\log_{10}(T_{FT}/T_{NISQ})$
- **Red regions:** NISQ is faster
- **Blue regions:** Fault-Tolerant is faster

## The tradeoff

- NISQ wins at low error rates where the exponential overhead is not too large
- The FT mode is preferable at higher error rates (provided there is enough space)

# Deeper circuits (40 Trotter Steps)



Log$_{10}$ Runtime Ratio (Fault-Tolerant / NISQ) for Ising Simulation (40 Trotter Steps)

## Reading the heatmap

- We plot the log-ratio of the runtimes: $\log_{10}(T_{FT}/T_{NISQ})$
- **Red regions:** NISQ is faster
- **Blue regions:** Fault-Tolerant is faster

## Doubling the depth shifts the balance

- The NISQ calculations hit an exponential wall
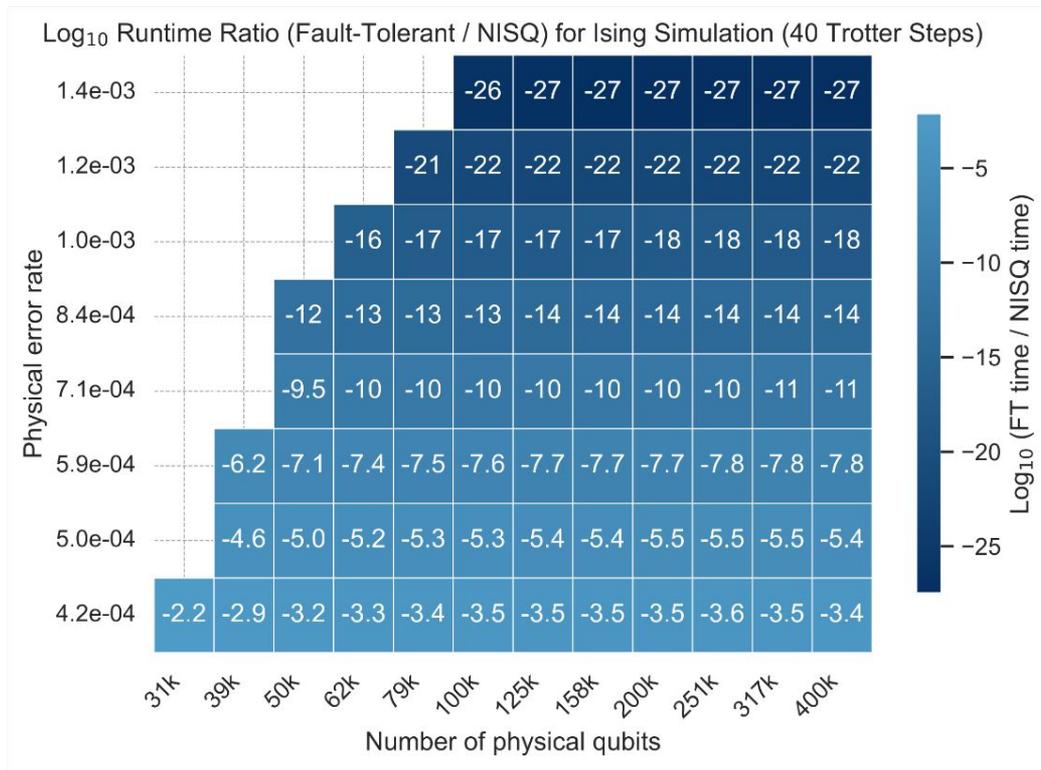- The fault-tolerant cost increases only marginally

Comparing FLASQ to previous FT estimates

# Reduced resource estimates compared to prior work



Ising Model Runtime Comparison: FLASQ (with PEC) vs Beverland et al. (2022)

## Previous state-of-the-art

- Beverland et al. (2022) analyzed a $10 \times 10$ Ising model (4th-order Trotter)
- Relied purely on QEC + traditional Magic State Distillation
- Strict target error budget of $\epsilon = 0.001$ (No error mitigation)

## FLASQ predicts much lower runtimes

- Same error budget, but PEC allows for smaller code distances
- Cultivation and error mitigation together reduce the cost of non-Clifford gates
- **More than an order of magnitude reduction in space and time costs at realistic physical error rates** ($p_{phys} = 10^{-3}$)

# Case Study: Hamming Weight Phasing

# The cost of synthesizing parallel rotations

## The Problem

- Synthesizing many arbitrary rotations can dominate the cost of an algorithm, even with cultivation
- Many applications (including the Ising model simulations) require executing many identical $R_Z(\theta)$ rotations in parallel

## One approach: Hamming weight phasing (HWP)

- Calculate the Hamming weight of the $N$ target qubits into an ancilla register of size $\approx \log_2(N)$
- Apply scaled rotations only to the register: $R_Z(\theta)$ to bit 1, $R_Z(2\theta)$ to bit 2, etc.
- Uncompute the Hamming weight
- T-count drops to $\approx 4N + \mathcal{O}(\log N \cdot \log \epsilon^{-1})$

# The cost of synthesizing parallel rotations

## The Problem

- Synthesizing many arbitrary rotations can dominate the cost of an algorithm, even with cultivation
- Many applications (including the Ising model simulations) require executing many identical $R_Z(\theta)$ rotations in parallel

## One approach: Hamming weight phasing (HWP)

- Calculate the Hamming weight of the $N$ target qubits into an ancilla register of size $\approx \log_2(N)$
- Apply scaled rotations only to the register: $R_Z(\theta)$ to bit 1, $R_Z(2\theta)$ to bit 2, etc.
- Uncompute the Hamming weight
- T-count drops to $\approx 4N + \mathcal{O}(\log N \cdot \log \epsilon^{-1})$

## Theoretical T-counts for precision $\epsilon = 10^{-5}$

| Qubits to Rotate | HWP T-Count | Parallel $R_Z$ T-Count | T-Count Ratio ($R_Z$/HWP) |
|---|---|---|---|
| 15 | $\approx 112$ | $\approx 270$ | $\approx 2.4\times$ |
| 43 | $\approx 252$ | $\approx 774$ | $\approx 3.1\times$ |
| 121 | $\approx 576$ | $\approx 2{,}178$ | $\approx 3.8\times$ |

## If we only count T gates...

- HWP looks very favorable
- But this is missing two important factors

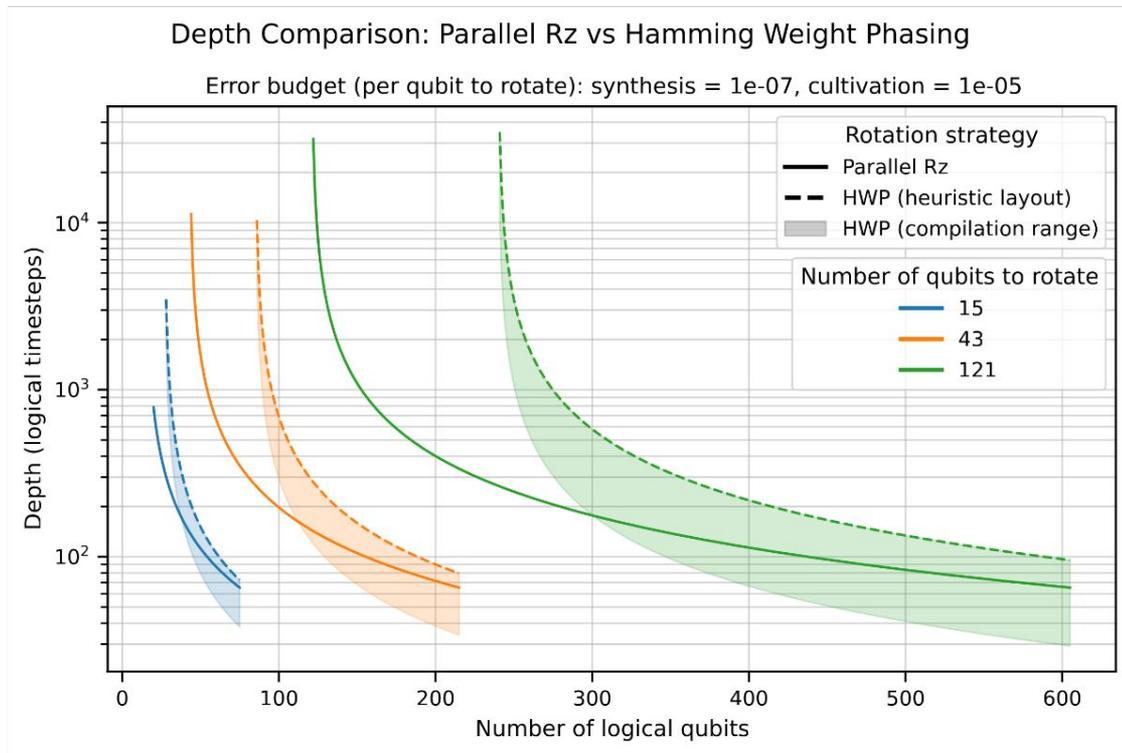# The other costs of Hamming Weight Phasing

Ancilla qubit overhead

- Computing the Hamming weight requires a sophisticated arithmetic circuit
- Prior work minimizes the T count but uses $\approx N$ extra "algorithmic" ancilla qubits
- These ancilla are not available to mediate other operations

Clifford and routing complexity

- The arithmetic circuits require a significant number of Clifford operations
- Compiling the multi-qubit Clifford and non-Clifford operations to a 2D grid adds more overhead

Evaluating the impact of these tradeoffs requires a detailed cost model

Quantum AI

# Evaluating HWP with the FLASQ model



Depth Comparison: Parallel Rz vs Hamming Weight Phasing

Error budget (per qubit to rotate): synthesis = 1e-07, cultivation = 1e-05

**Comparing depth (in terms of logical timesteps)**

- **Solid lines:** Naive Parallel $R_Z$ approach
- **Dotted lines:** HWP with an explicit (but not optimal) layout on a 2D grid
- **Shaded region:** HWP neglecting some fraction of the routing overhead

**Takeaways**

- The other costs of HWP signficantly reduce its potential advantage
- A careful analysis, with an explicit layout, would be required to benefit at all (even with abundant space)

Quantum AI

arxiv.org/abs/2511.08508

Quantum AI

# FLuid Allocation of Surface code Qubits (FLASQ)

Goal: estimate the resources required to implement a quantum circuit ...

- ... In a two-dimensional surface code architecture
- ... With an (currently non-existent) compilation stack

Key assumption: Operations can be freely rearranged in space and time

- Their "extra" ancilla spacetime volume is conserved
- This is justified by the inherent flexibility of lattice surgery and the use of walking surface codes
- We separately account for the reaction depth

This assumption makes it easy to estimate the overall spacetime volume

- This enables estimates of error rates, wall-clock times, etc
- Can be done programmatically (ask me for the code!)
- We can use this ability to balance tradeoffs

But it also ignores some important factors

Quantum AI

# Going forward

### Early fault-tolerance is a moving target

- Magic state cultivation
- Developments in quantum error correction
- Careful combinations of existing techniques can have a large impact

### Just counting T gates may not be sufficient

- If you are dominated by small angle rotations it might be okay
- But the cost of routing and Clifford operations can be large
- Even determining this requires a good model!

### Compilation for early fault-tolerance is a big project

- FLASQ doesn't actually solve the problem, it just guesses the cost
- This work is necessary and may lead to surprising cost reductions

Quantum AI

Thank you

Quantum AI