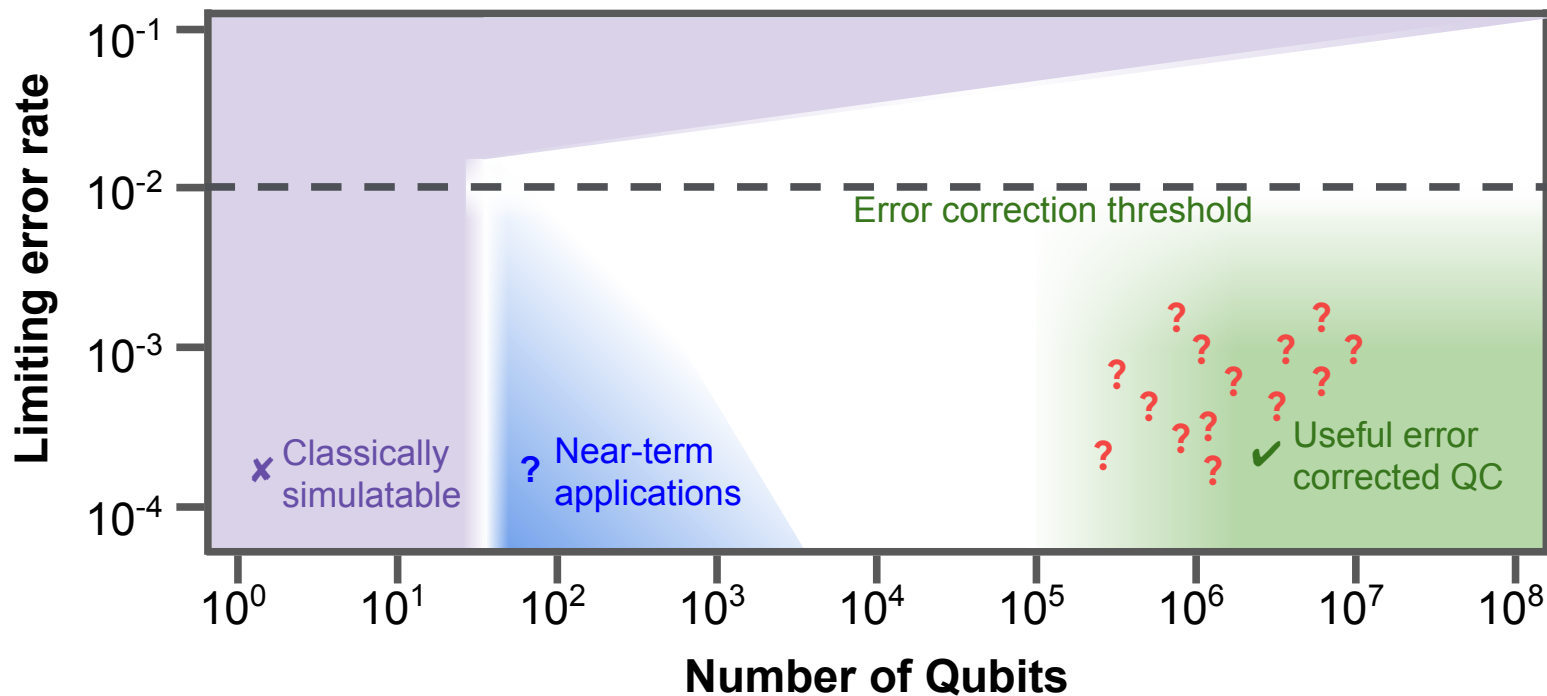


# The **FL**uid **A**llocation of **S**urface code **Q**ubits model for early fault-tolerant resource estimation

William J. Huggins  
WERQSHOP - July 2025

# Understanding the transition to fault-tolerance



# Error mitigation for today's quantum computers...

Time evolution of the 2D Ising model

$$H = J \sum_{\langle ij \rangle} Z_i Z_j + h \sum_i X_i$$

Second-order Trotter

$$N_{2q} \approx 2N^2 \frac{T}{\Delta t}$$

Target  $O(1)$  errors

$$N_{2q} \lesssim p_{2q}^{-1}$$

For example:

$$T = 1?$$

$$\delta T = .04$$

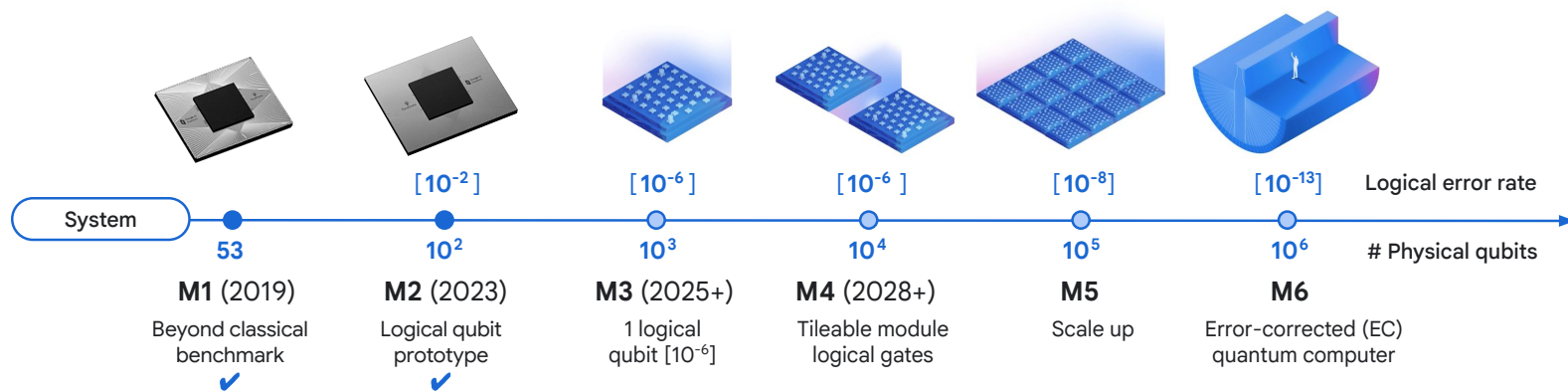
$$N = 10$$

$$p_{2q}^{-1} = .001$$

$$T \lesssim .2$$

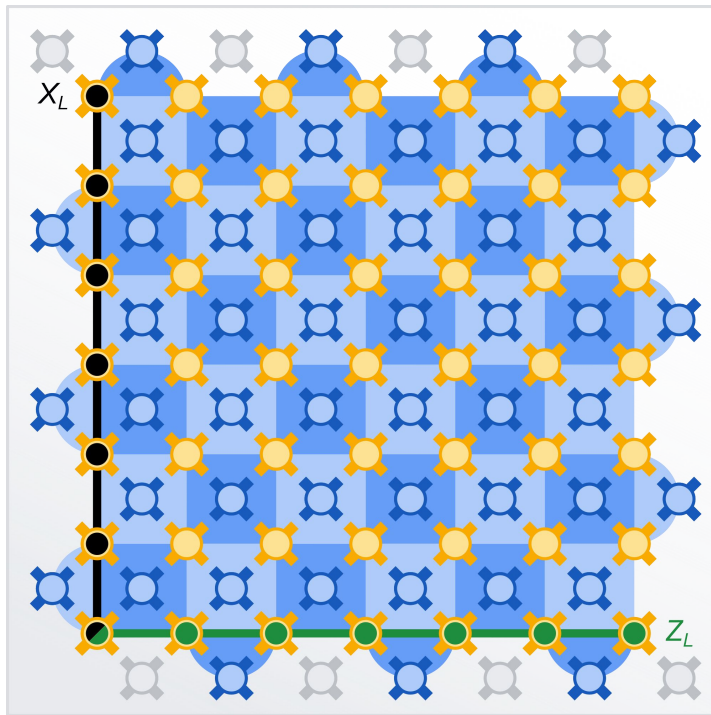


# Understanding the transition to fault-tolerance



# Ingredients

# Surface code logical qubit



$$\text{Physical qubit count} = 2 * (d+1)^2$$

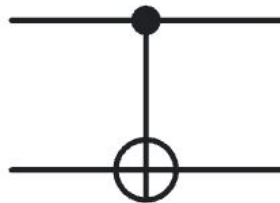
$$\text{Logical Error Rate} \approx 0.1 * \Lambda^{-(d+1)/2}$$

Surface code error correction gives exponentially more time for quadratically more space.

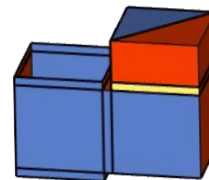
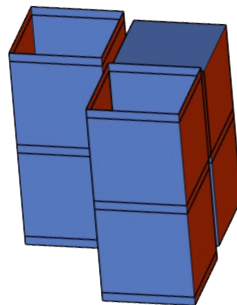
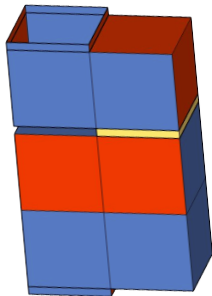
2D local connectivity is sufficient to store information and perform gates using lattice surgery

# Everything\* requires ancilla spacetime

Gates in NISQ



Gates in the surface code



# Rearranging the ancilla with walking codes

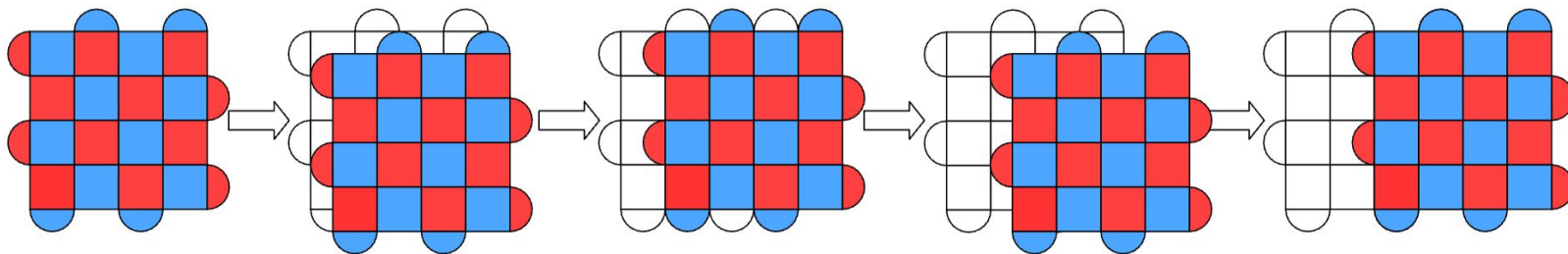
With limited qubits, we expect to have

$$N_{\text{ancilla}} \ll N_{\text{data}}$$

Walking surface codes let us move logical qubits

We will want to move ancilla around

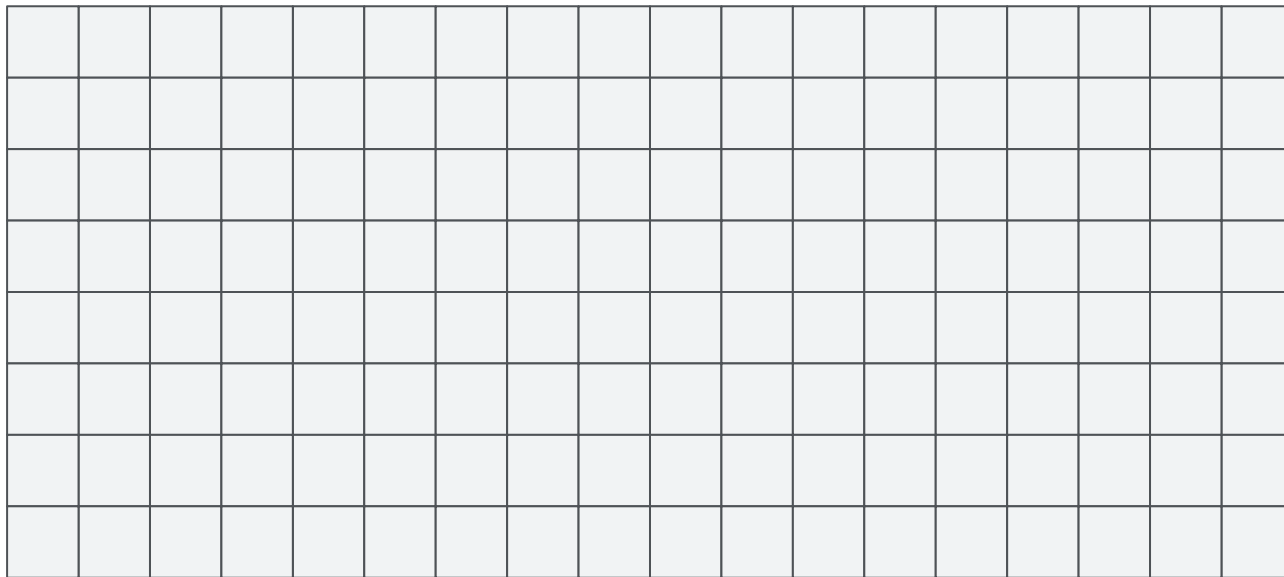
Whole patches can move together



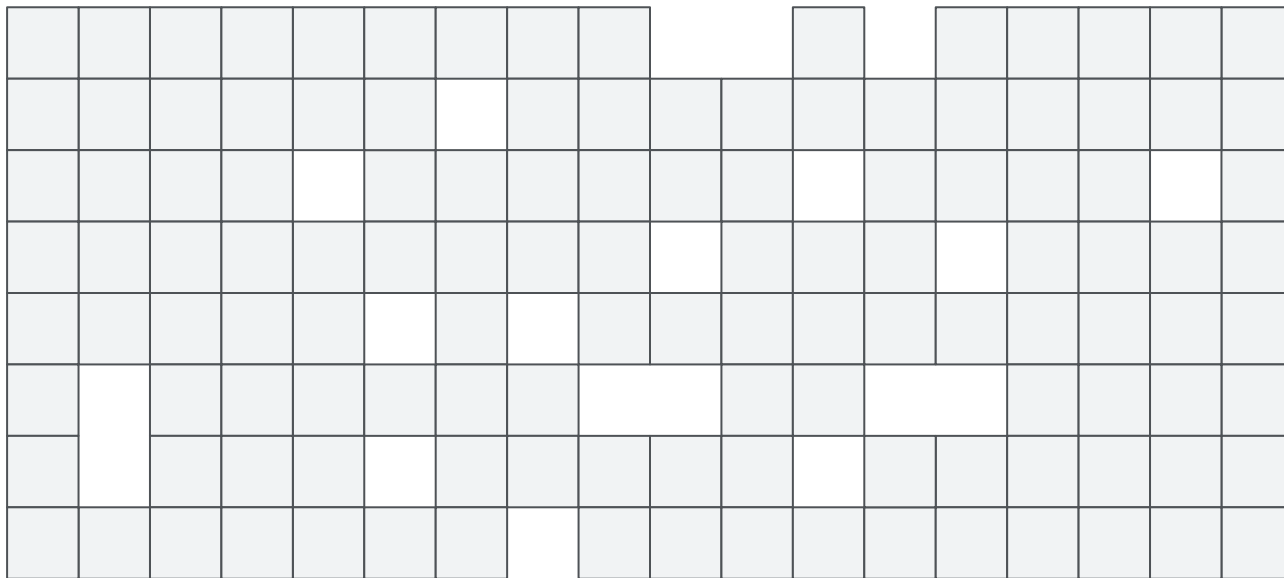


# Walking codes in action

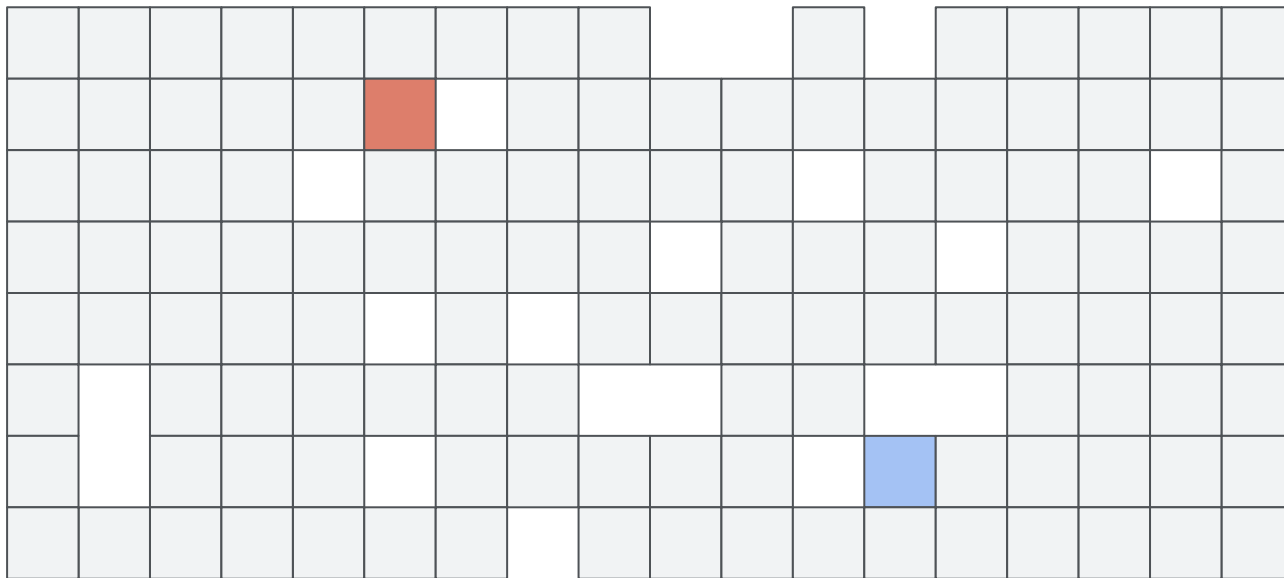
---



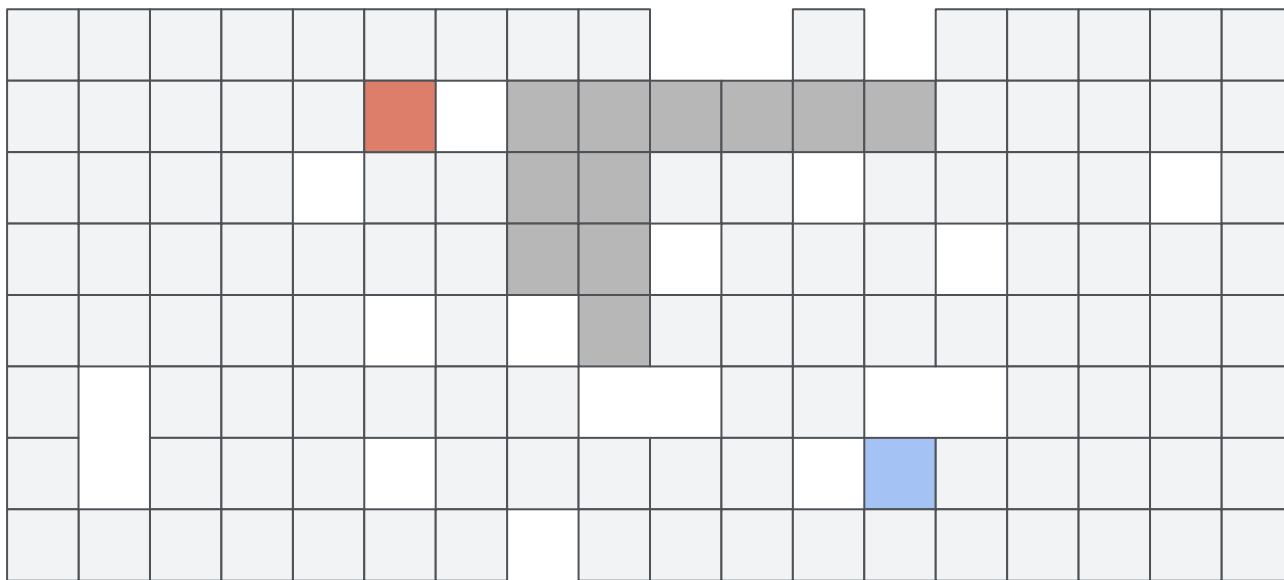
# Walking codes in action



# Walking codes in action



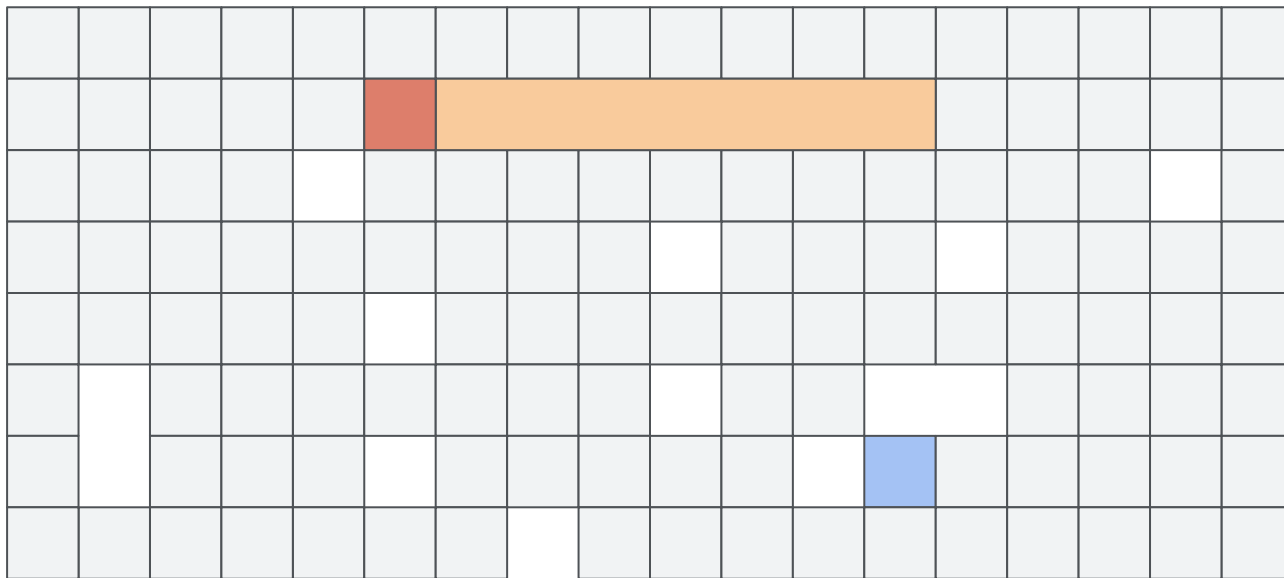
# Walking codes in action



Time 0 → Time 1

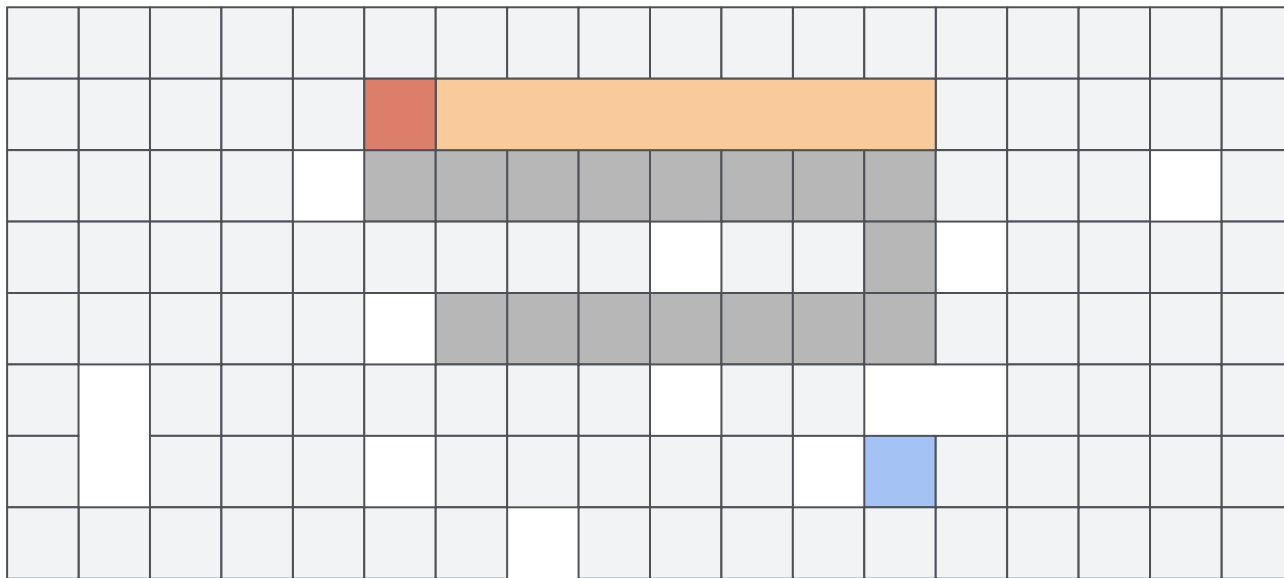


# Walking codes in action



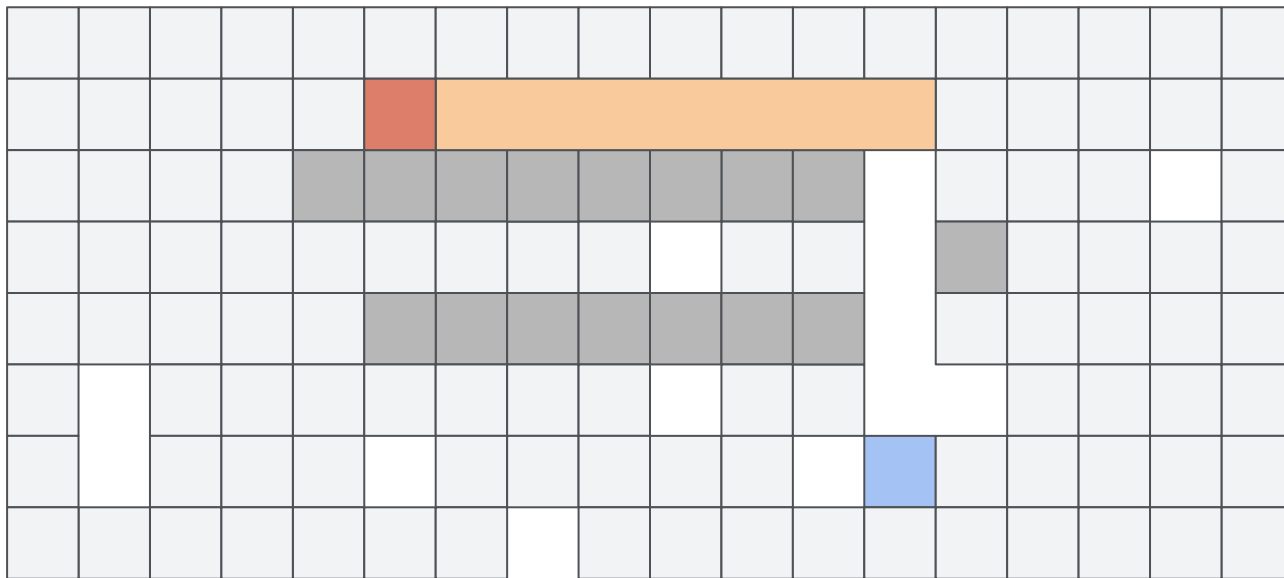
**Time 1 → Time 2**

# Walking codes in action



**Time 1 → Time 2**

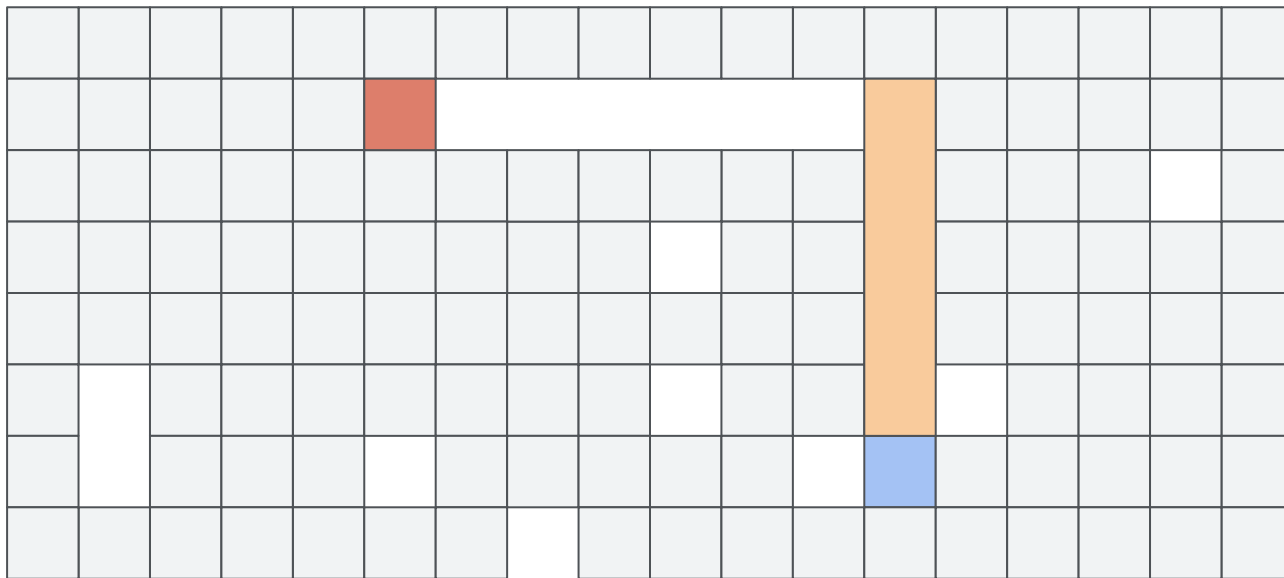
# Walking codes in action



Time 1 → Time 2

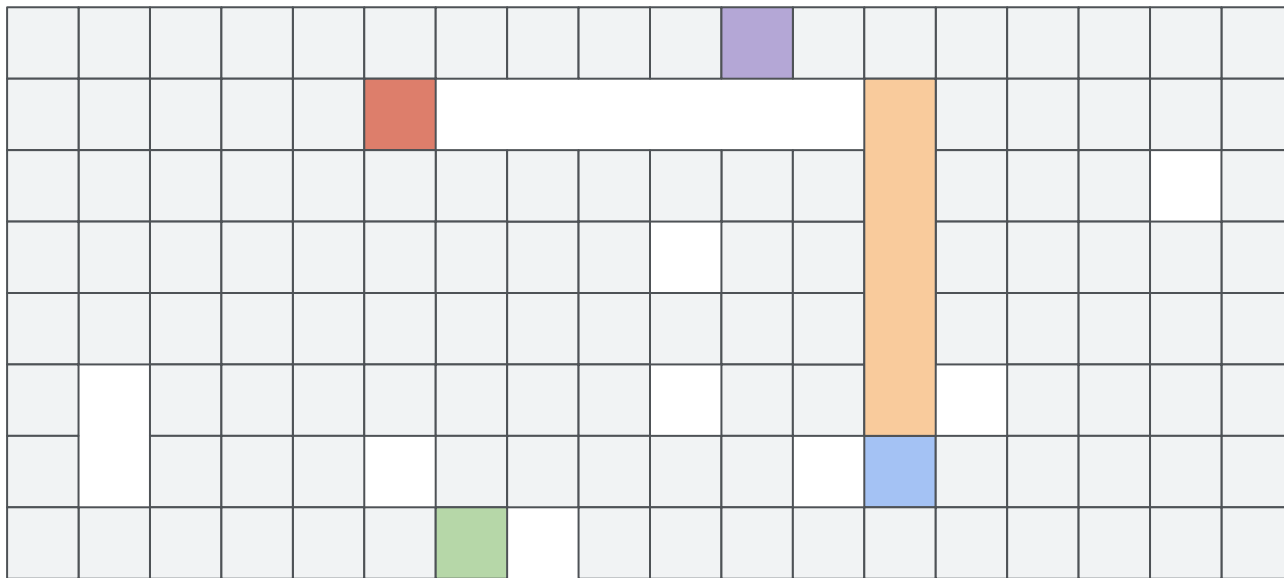


# Walking codes in action



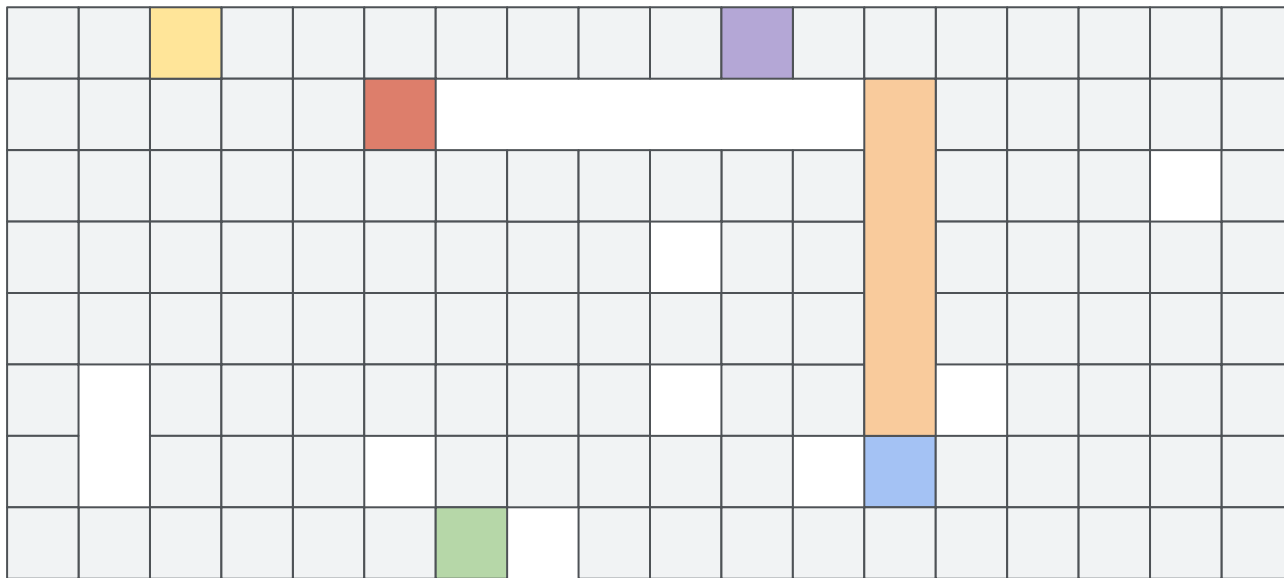
Time 2 → Time 3

# Walking codes in action



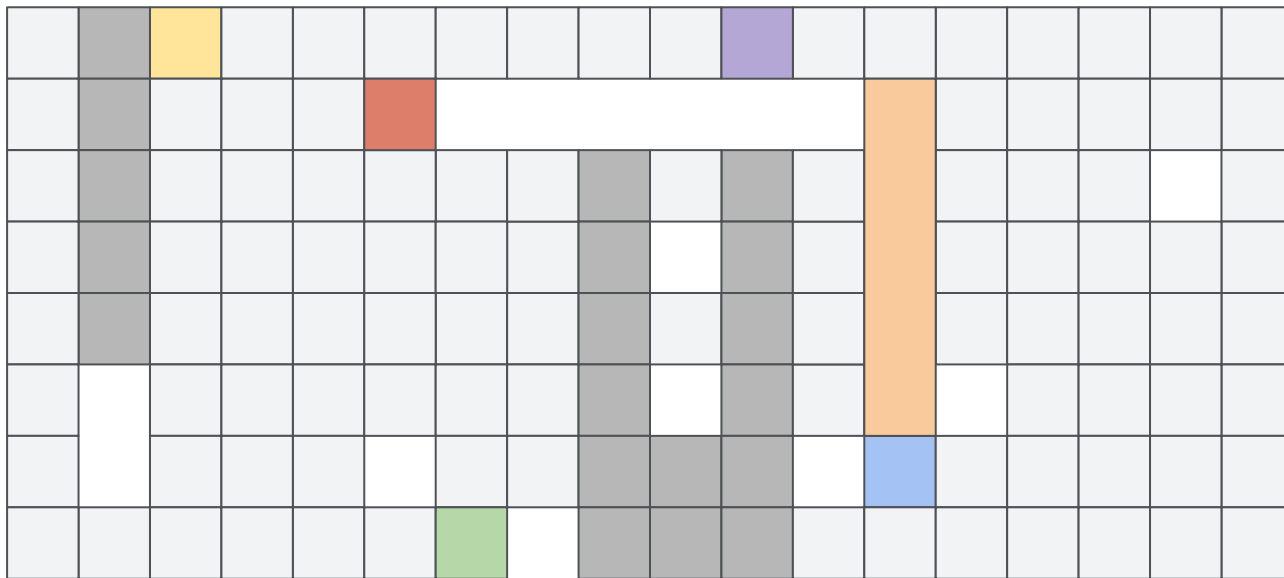
**Time 2 → Time 3**

# Walking codes in action



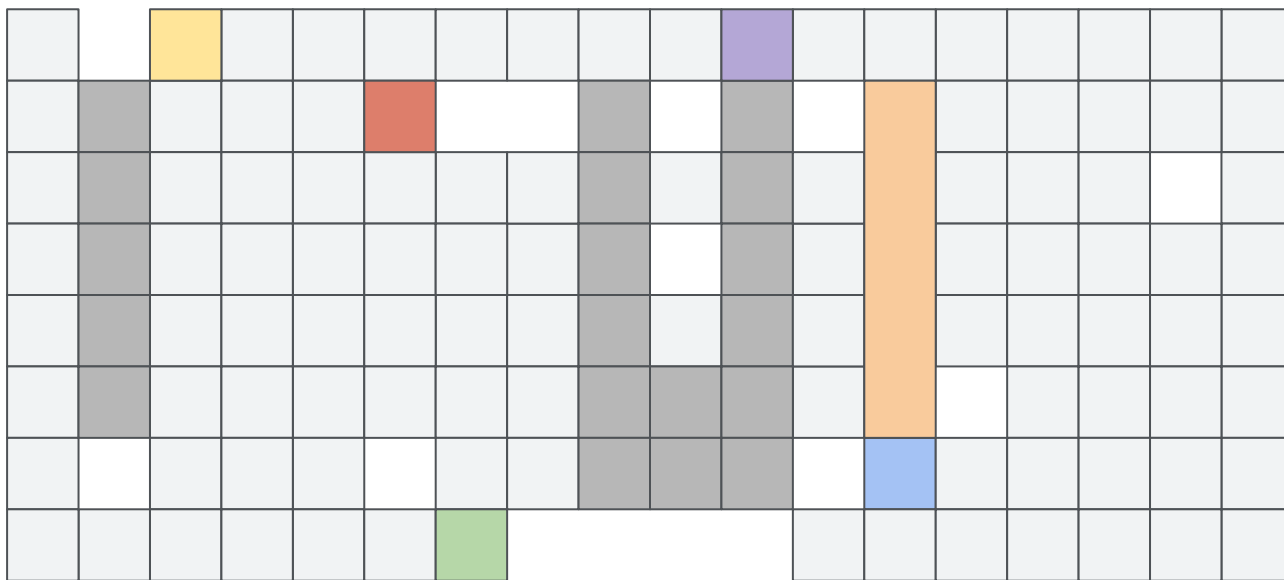
**Time 2 → Time 3**

# Walking codes in action



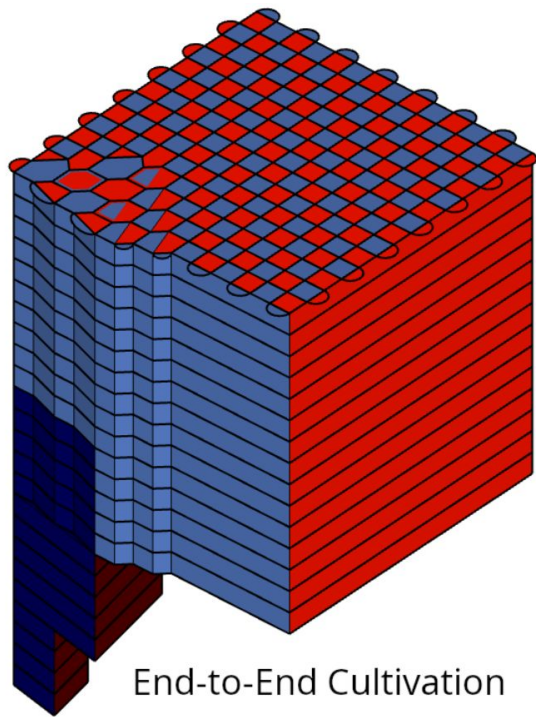
**Time 2 → Time 3**

# Walking codes in action



Time 2 → Time 3

# Magic state cultivation

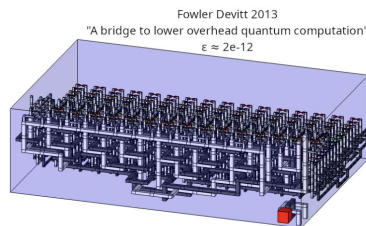


End-to-End Cultivation

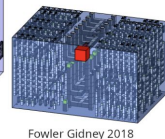
Magic state cultivation prepares a resource state we can use to perform a T gate.

The key ideas:

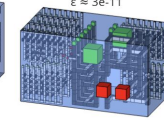
- Check the T state in a low-distance color code
- Rapidly grow to a larger code and postselect
- Low but finite error rate with much lower spacetime volume than previous techniques



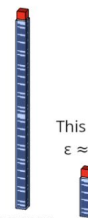
Fowler Devitt 2013  
"A bridge to lower overhead quantum computation"  
 $\epsilon \approx 2e-12$



Fowler Gidney 2018  
"Low overhead quantum computation  
using lattice surgery"  
 $\epsilon \approx 9e-17$



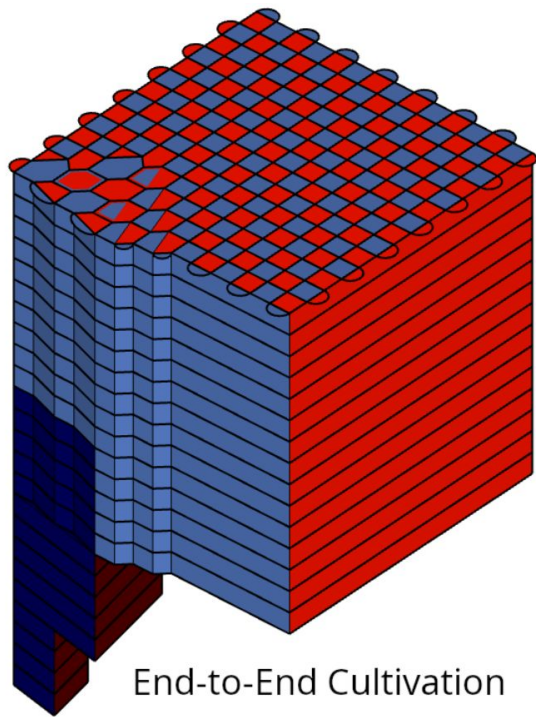
Gidney Fowler 2019  
"Efficient magic state factories with a  
catalyzed CCZ  $\rightarrow$  2T transformation"  
 $\epsilon \approx 3e-11$



This paper  
 $\epsilon \approx 2e-9$

This paper  
 $\epsilon \approx 4e-6$

# 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; HIF24] to efficiently prepare good  $|T\rangle$  states. We call our construction “magic state cultivation”

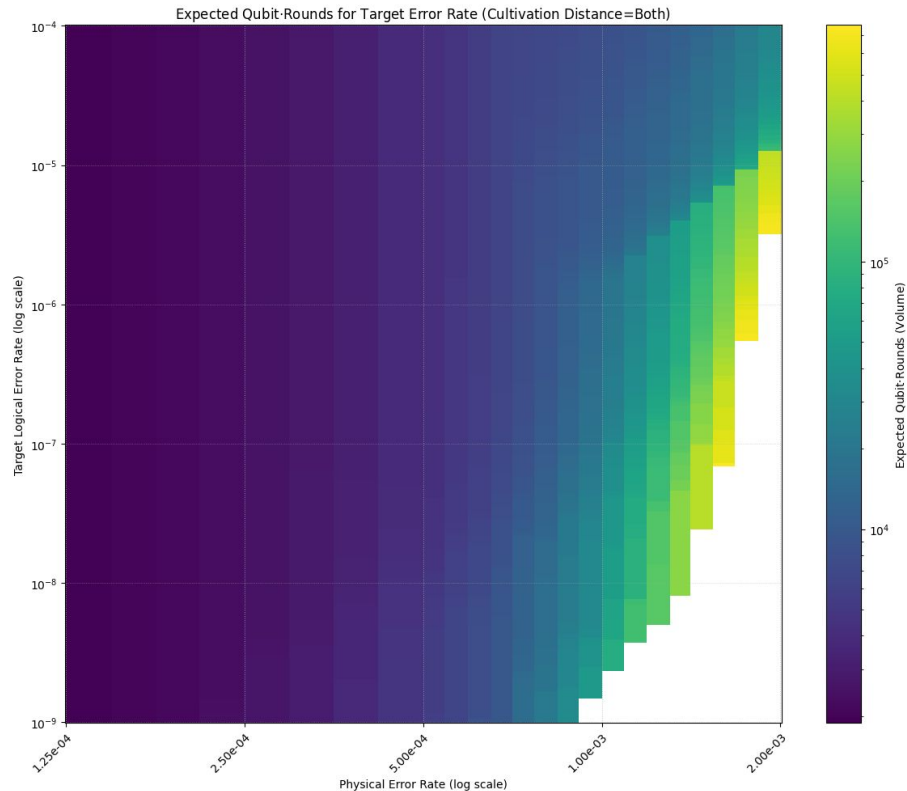
# Magic state cultivation

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

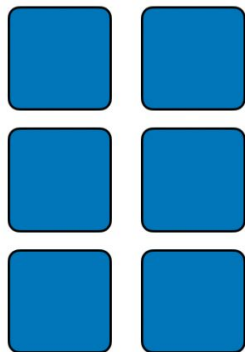
(Performed using the code and methodology of *Gidney et al.*)





# The model

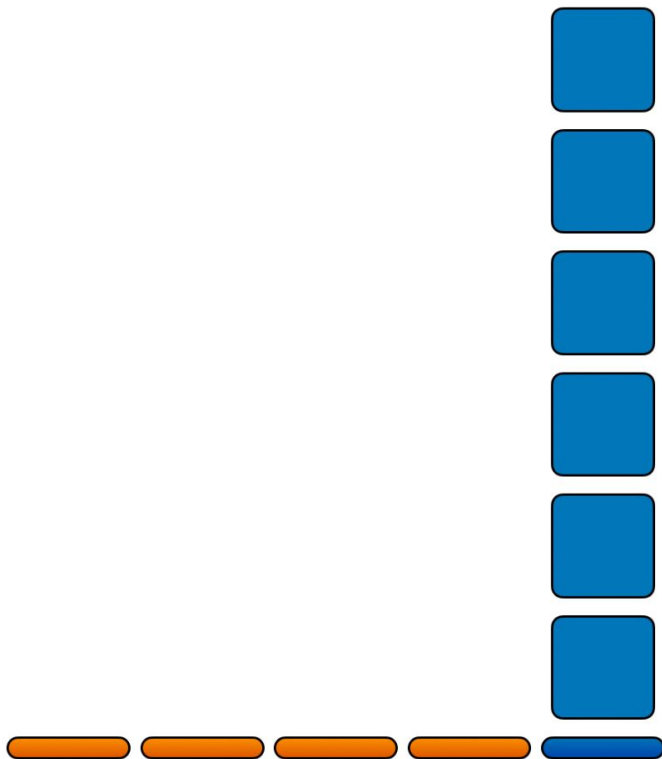
# Fluid Allocation of Surface Code Qubits



Each gate gets a cost that represents the extra spacetime volume it requires.

$$\text{COST} = \sum_g c_g$$

# FLuid Allocation of Surface Code Qubits



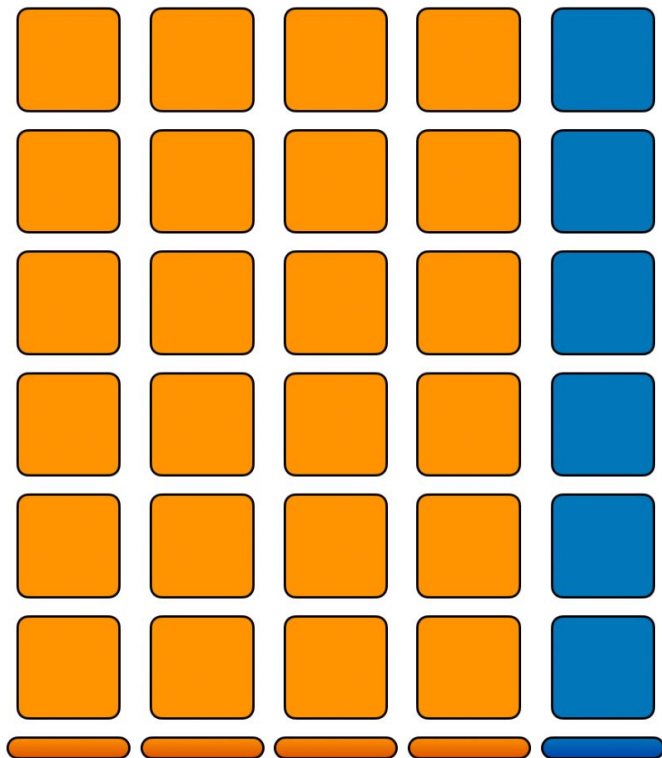
Each gate gets a cost that represents the extra spacetime volume it requires.

$$\text{COST} = \sum_g c_g$$

We make sure we have enough extra spacetime for all the gates.

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

# FLuid Allocation of Surface Code Qubits



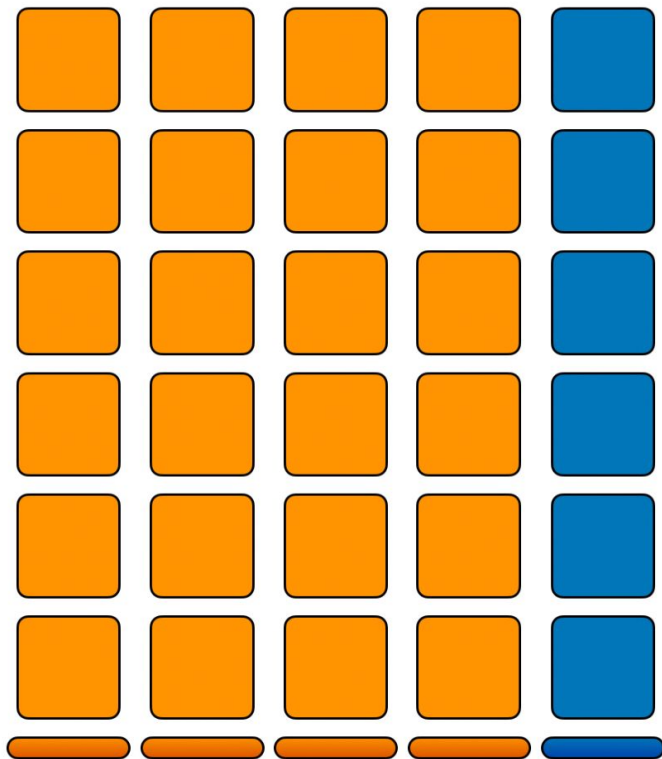
Each gate gets a cost that represents the extra spacetime volume it requires.

$$\text{COST} = \sum_g c_g$$

We make sure we have enough extra spacetime for all the gates.

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

# FLuid Allocation of Surface Code Qubits



Each gate gets a cost that represents the extra spacetime volume it requires.

$$\text{COST} = \sum_g c_g$$

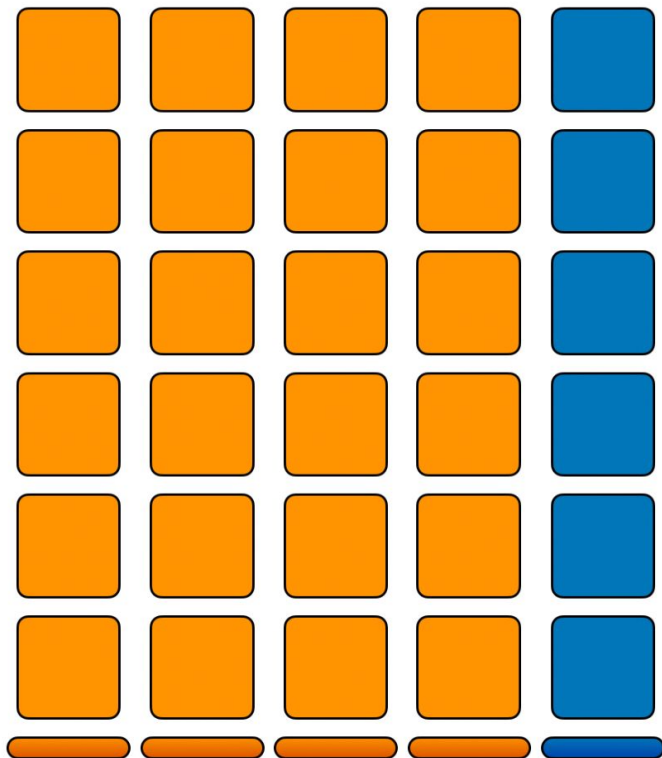
We make sure we have enough extra spacetime for all the gates.

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

(We also make sure the depth is large enough to account for measurements that have to happen sequentially.)

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

# FLuid Allocation of Surface Code Qubits



Each gate gets a cost that represents the extra spacetime volume it requires.

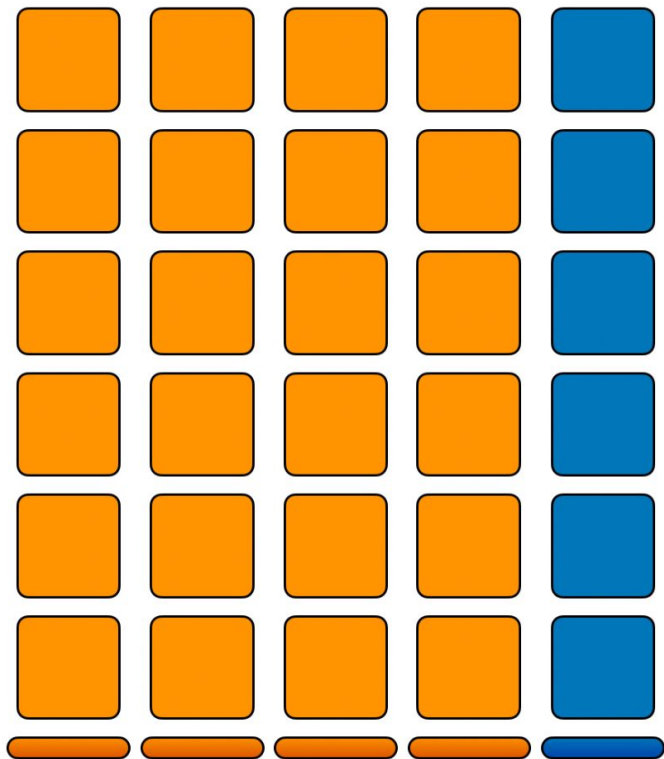
$$\text{COST} = \sum_g c_g$$

We make sure we have enough extra spacetime for all the gates.

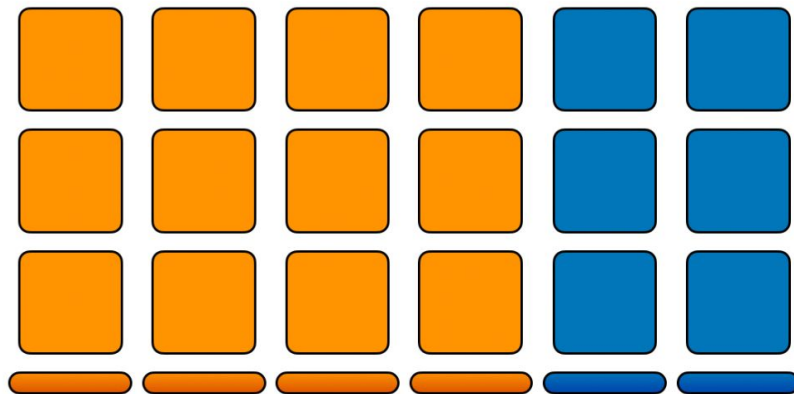
$$\text{DEPTH} \geq \frac{\text{COST}}{\text{NUMBER OF FLUID ANCILLA}}$$



# Fluid Allocation of Surface Code Qubits



The depth shrinks as we increase the number of fluid ancilla qubits



# FLuid Allocation of Surface Code Qubits

Each gate gets a cost that represents the extra spacetime volume it requires.

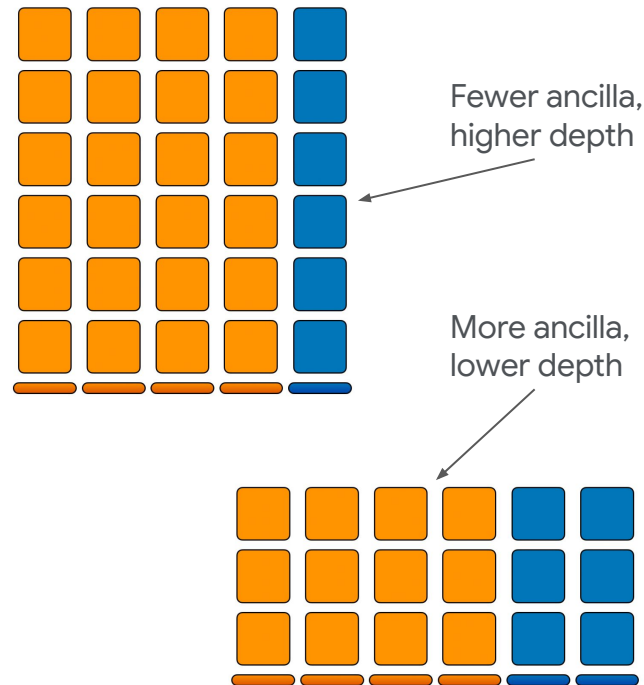
$$\text{COST} = \sum_g c_g$$

We make sure we have enough extra spacetime for all the gates.

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

(We also make sure the depth is large enough to account for measurements that have to happen sequentially.)

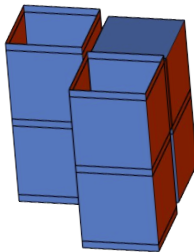
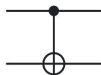
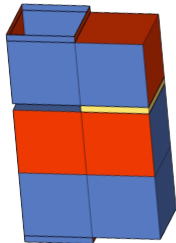
$$\text{DEPTH} \geq \text{REACTION TIME} \cdot \text{MEASUREMENT DEPTH}$$





# Gate costs in the FLASQ model

Basic gates	Conservative cost	Optimistic cost	Notes
$X / Y / Z$	0	—	Implemented in software
$X / Z$ basis measurement (or initialization)	0	—	
$H$	3	1.5	Patch rotation each time (patch rotation half the time)
$S / S^\dagger$	1.5	1	Gate performed statically (performed in motion half the time)
$T / T^\dagger$	$1.5v(p_{phys}, p_{cult}) + 2$	$v(p_{phys}, p_{cult}) + 2$	Buffer to account for packing constraints (no buffer) Depends on physical error rate ( $p_{phys}$ ) and target logical error rate ( $p_{cult}$ )
$CNOT / CZ$	$2p(q_1, q_2) + 3$	$p(q_1, q_2)$	Buffer to account for extra routing (no buffer)



Cost in units of  $d^3$ , represents volume required above and beyond memory.

$p(q_1, \dots, q_k)$  denotes the minimum length (measured in the manhattan distance) required to connect the qubits.

The volume required to cultivate a magic state depends on the underlying physical error rate  $p_{phys}$  and the target logical error rate  $p_{cult}$ .

# Gate costs in the FLASQ model

Basic gates	Conservative cost	Optimistic cost	Notes
$X / Y / Z$	0	—	Implemented in software
$X / Z$ basis measurement (or initialization)	0	—	
$H$	3	1.5	Patch rotation each time (patch rotation half the time)
$S / S^\dagger$	1.5	1	Gate performed statically (performed in motion half the time)
$T / T^\dagger$	$1.5v(p_{phys}, p_{cult}) + 2$	$v(p_{phys}, p_{cult}) + 2$	Buffer to account for packing constraints (no buffer) Depends on physical error rate ( $p_{phys}$ ) and target logical error rate ( $p_{cult}$ )
$CNOT / CZ$	$2p(q_1, q_2) + 3$	$p(q_1, q_2)$	Buffer to account for extra routing (no buffer)
Additional primitive gates	Cost		Notes
$T_X = HTH / T_X^\dagger$	$\text{Cost}(T)$	—	Can be implemented directly using cultivation
$R_Z(\theta)$	$\approx (4.86 + .53 \log_2(\epsilon_{rot})) \text{Cost}(T) + 2 \text{Cost}(S) + 2 \text{Cost}(H)$	$\approx (4.86 + .53 \log_2(\epsilon_{rot})) \text{Cost}(T) + \frac{7}{8} \text{Cost}(S) + \frac{3}{8} \text{Cost}(H)$	Depends on target rotation synthesis error ( $\epsilon_{rot}$ , in the diamond norm)
Move	$p(q_1, q_2)$	—	Swaps the state of two qubits, one of which must be in the $ 0\rangle$ state
Toffoli	$4 \text{Cost}(T) + 2p(q_1, q_2, q_3) + 10$	$4 \text{Cost}(T) + p(q_1, q_2, q_3) + 10$	Buffer to account for extra routing (no buffer)
$k$ -target $CNOT$	$2p(q_0, \dots, q_k) + 3$	$p(q_0, \dots, q_k)$	Buffer to account for extra routing (no buffer)
$k$ -parity-controlled $X$ gate	$2p(q_0, \dots, q_k) + 3$	$p(q_0, \dots, q_k)$	Buffer to account for extra routing (no buffer) Controlled based on parity of multiple control qubits
$Y$ basis measurement (or initialization)	.5	—	

Cost in units of  $d^3$ , represents volume required above and beyond memory.

$p(q_1, \dots, q_k)$  denotes the minimum length (measured in the manhattan distance) required to connect the qubits.

The volume required to cultivate a magic state depends on the underlying physical error rate  $p_{phys}$  and the target logical error rate  $p_{cult}$ .

# Applications

# Error mitigation overhead with probabilistic error cancellation

$$\text{NUMBER OF SAMPLES} \propto \Gamma^2 \frac{\|O\|}{\epsilon^2}$$

We can reduce the resource requirements by combining error correction and error mitigation.

$$p_{cyc} = .1\Lambda^{(d+1)/2}, \quad \Lambda = \frac{.01}{p_{phys}}$$

Here we use probabilistic error cancellation and a simple noise model.

$$V = N_{active} \cdot \text{DEPTH} + \text{COST}$$

Our FLASQ model lets us estimate how many opportunities we have for errors.

$$\Gamma^2 \approx (1 - p_{cyc})^{-4dV}$$

We can convert this into an estimate of the error mitigation overhead

# Resource estimates for lattice models

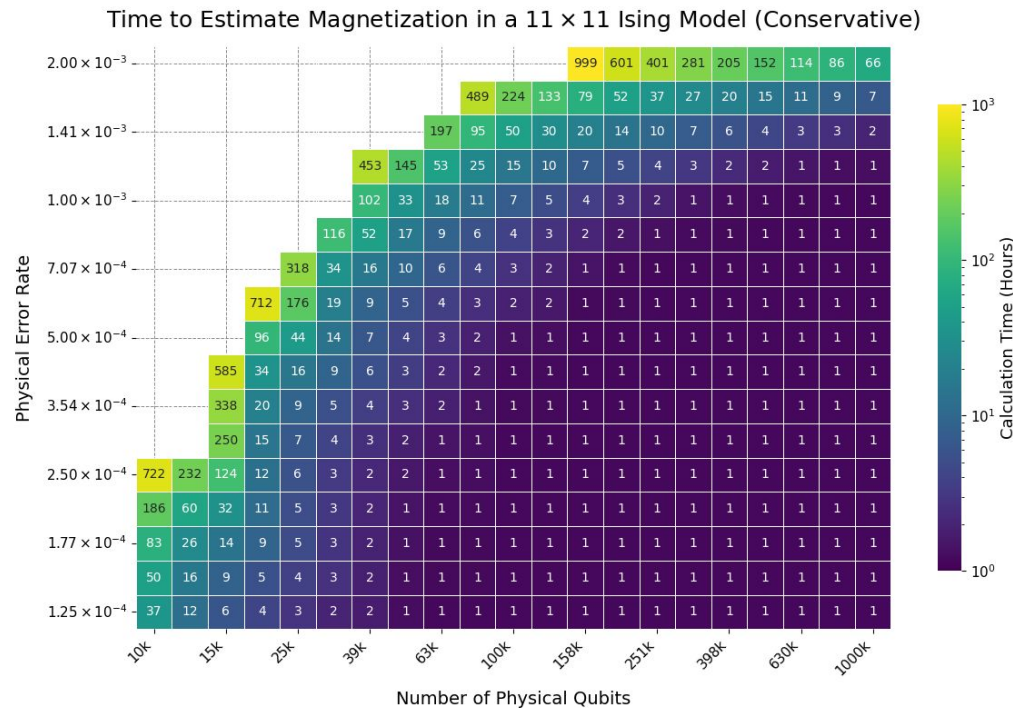
- 2D Ising model

$$H = J \sum_{\langle ij \rangle} Z_i Z_j + h \sum_i X_i$$

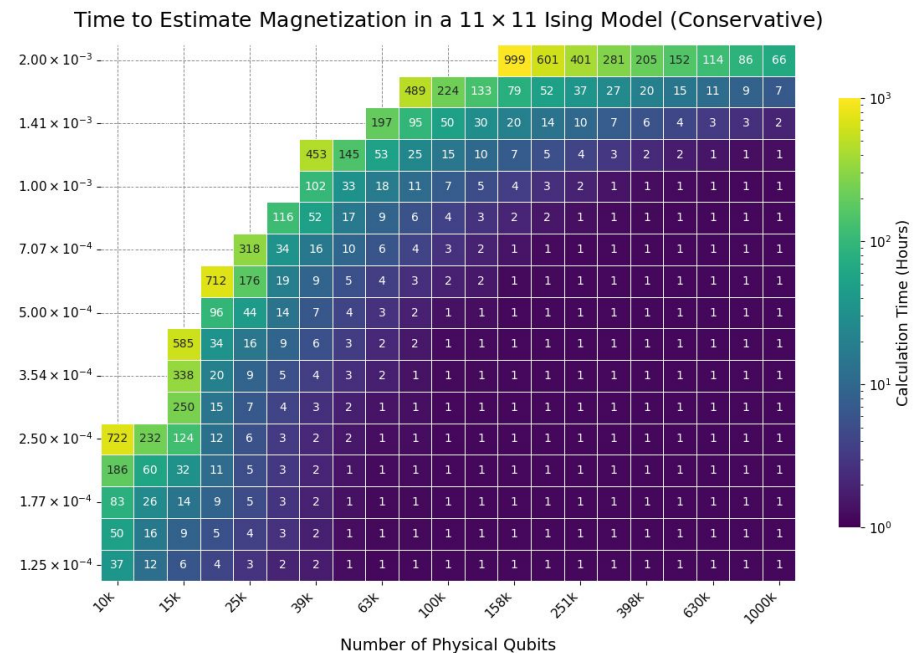
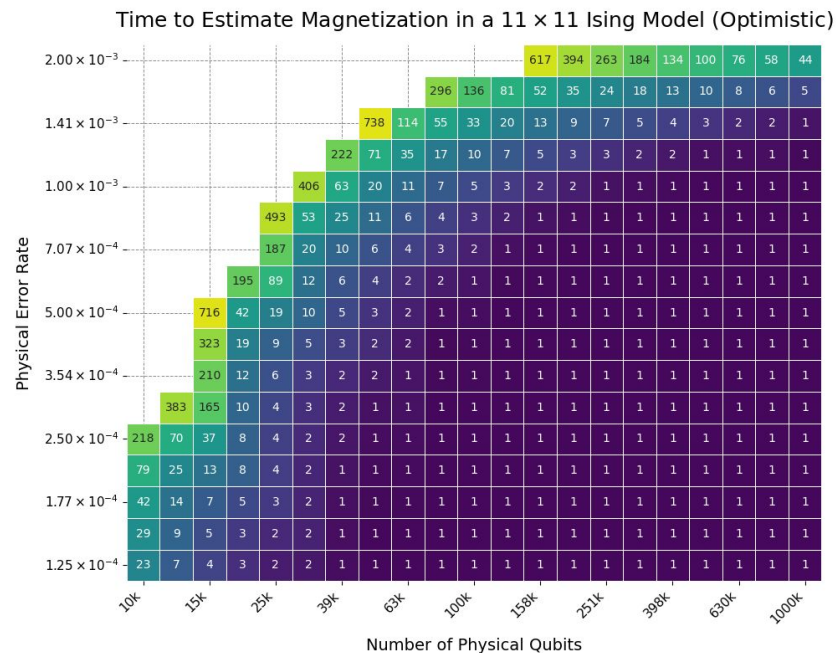
- Time-evolve after a quench and measure magnetization

- Probing difficult regime

- Near phase transition
- $T = 1$ ,  $dt=0.05$ , 2nd-order Trotter
- .002 standard deviation (comparable with disagreement between classical methods)<sup>[1]</sup>



# Resource estimates for lattice models



# Synthesizing rotations is expensive

Even the best methods for synthesizing a single-qubit Z rotation require many T gates<sup>[2]</sup>:

$$T(R_Z) = 4.86 + 0.53 \log_2 \epsilon^{-1}$$

This is the dominant cost for many applications.

Can we reduce it?

Hamming weight phasing lets us implement K identical  $R_Z$  rotations in parallel using only  $\log_2(K)$  rotations and  $\approx 4K$  additional T gates.

It works by calculating the Hamming weight into an ancilla register and applying rotations to each of the ancilla.

Standard compilations use additional ancilla to minimize the number of T gates.

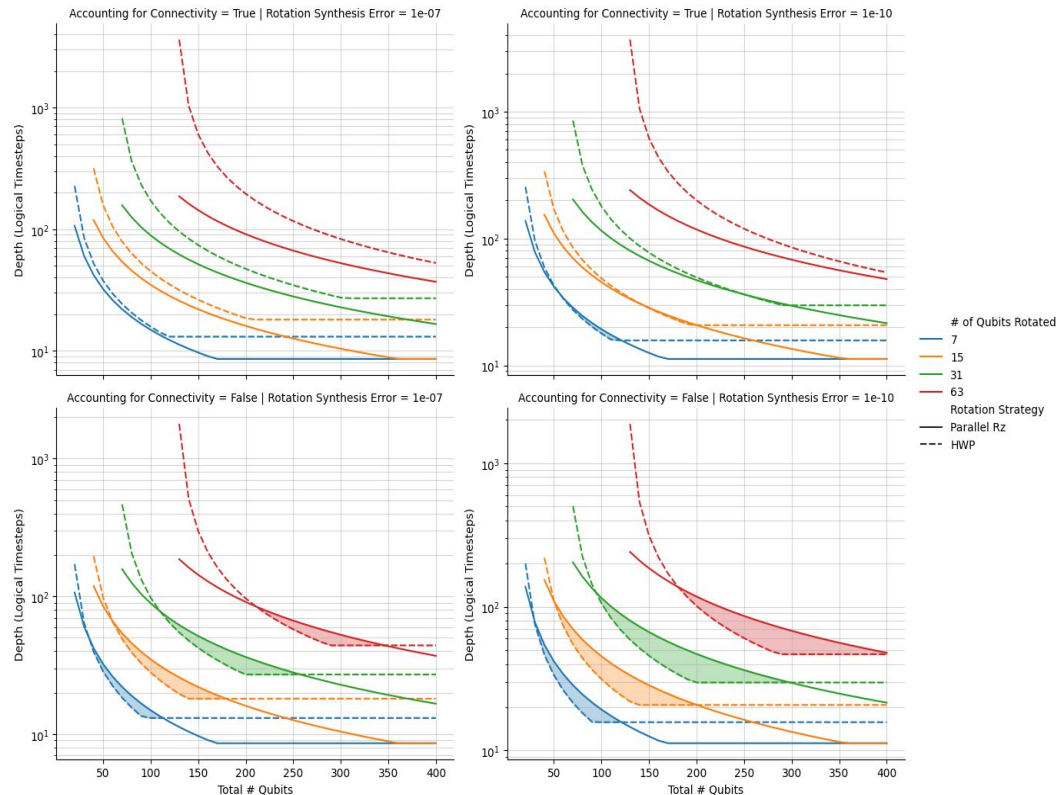
# Lowering the cost with Hamming weight phasing?

We compare parallel rotations synthesis with Hamming weight phasing.

- Top: We lay the circuits out on a 2D grid using a heuristic.
- Bottom: We neglect the cost of long-range interactions.

We shade the region where Hamming weight phasing yields an improvement in the circuit depth.

Comparing Parallel Rotation Synthesis with Hamming Weight Phasing (Conservative FLASQ Model)





Thank you

# Magic state cultivation simulations (expected number of retries)

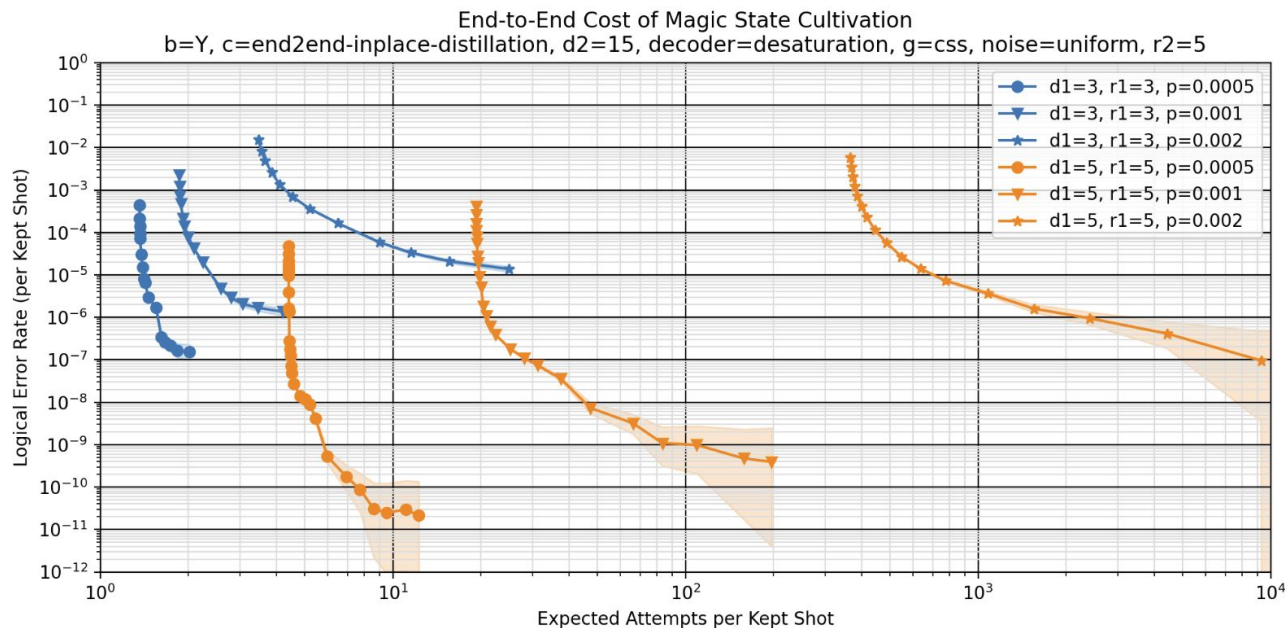
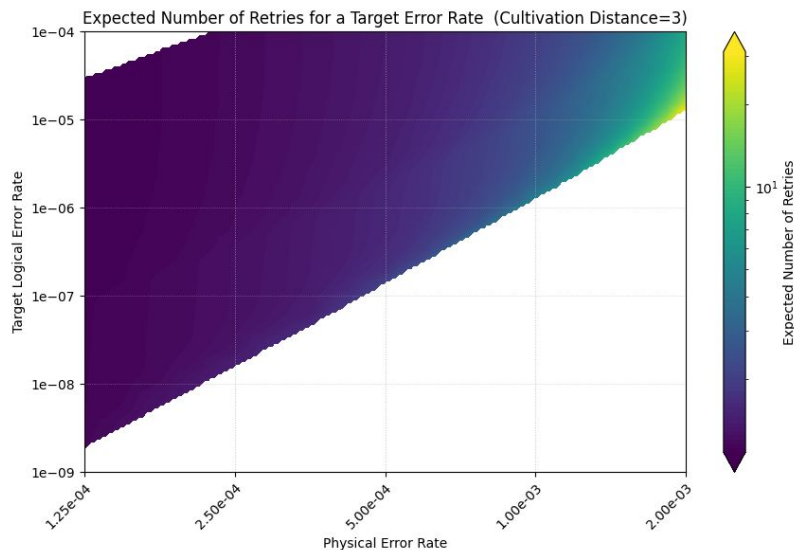
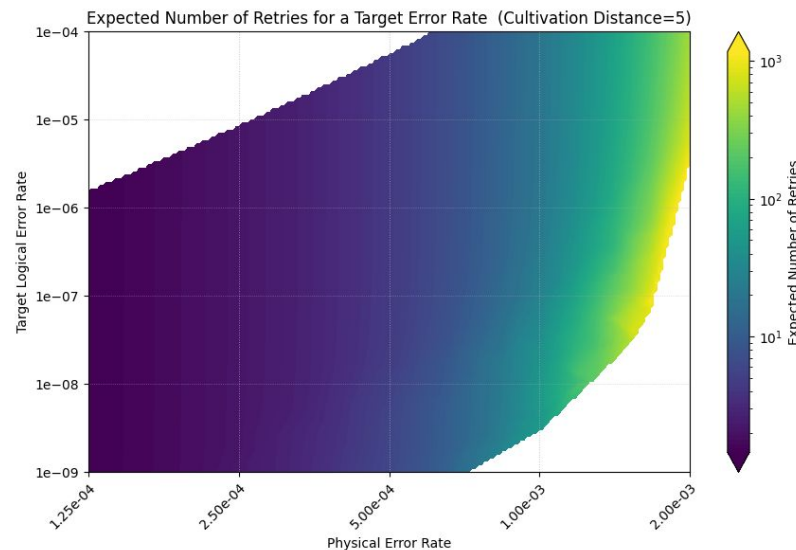


Figure 14: **End to end simulation of magic state cultivation.** Curves show the effect of varying the complementary gap cutoff for rejecting shots, decreasing logical error rate at the cost of increasing expected attempts. [Click here to open an example end-to-end circuit in Crumple.](#) Shaded regions indicate error rate hypotheses with a likelihood within a factor of 1000 of the max likelihood hypothesis, given the sampled data.

# Magic state cultivation simulations (expected number of retries)



- 21 different error rates
- 5 billion distance 3
- 50 billion distance 5



- Same methodology as cultivation paper
- Trimmed data with high uncertainty

# Magic state cultivation simulations (expected spacetime volume)

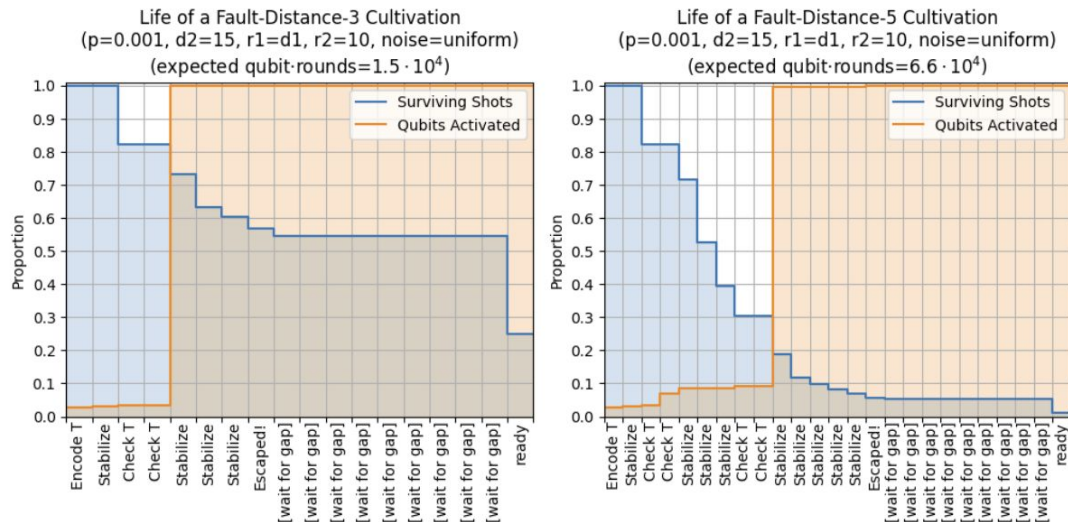


Figure 15: **Growth and survival during magic state cultivation.** Each column is a labelled circuit cycle, which performs roughly ten layers of one/two qubit gates and ends with one or two layers of dissipative gates. The proportion of activated qubits starts small during the injection stage, gradually increases during the cultivation stage, and then jumps for the escape stage. Survival rates drop continuously due to postselected detectors failing, then hold steady for 10 cycles while waiting for the complementary gap to be computed. Expected qubit · rounds is computed by integrating the product of the survival rate times the qubit proportion, then multiplying by the total number of qubits and dividing by the final survival rate. This is slightly optimistic, due to ignoring factors like packing efficiency.