Quantum AI

# Benchmarking NISQ and QEC experiments with tensor networks
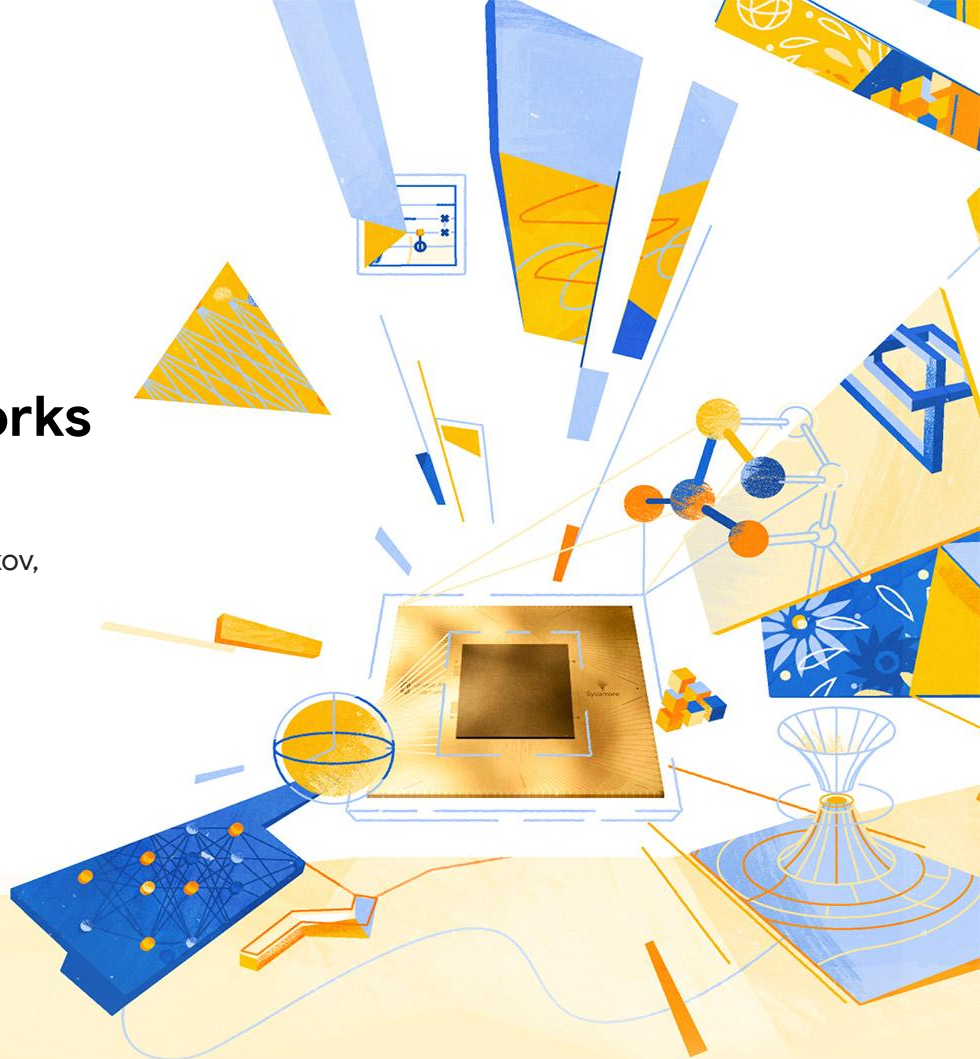
**Benjamin Villalonga**

S. Boixo, S. Mandrà, M. Newman, N. Shutty, K. Kechedzhi, S. Isakov, X. Mi, V. Smelyanskiy, and many others

Workshop on Tensor Networks
IPAM @ UCLA
**Los Angeles 2024**

Google

# Outline

Quantum AI

# Outline

Quantum AI

# The *status quo* of quantum computing experiments

**Useful** or physically motivated **applications** (error mitigation):
- Topological phases of matter, majorana edge modes, non-abelian statistics (Satzinger et al. 2021, Mi et al. 2022, Andersen et al. 2022, …)
- Time crystals (Mi et al. 2022, …)
- Information scrambling quantum systems (Mi et al. 2021)
- Floquet evolution of transverse field Ising model (Y. Kim et al., 2023)
- MERA implementation (Haghshenas et al. 2023)
- Dissipative cooling (Mi et al. 2023)
- Graph problems (Deng et al. 2023)
- Other experiments from Harvard/QuEra, IBM, Quantinuum, USTC, …
- . . .

**NISQ**

**Beyond-classical demonstration** attempts (usually no error mitigation involved):
- Random circuit sampling (Arute et al. 2019, Wu et al., 2021, Zhu et al. 2022, Morvan et al. 2023, Bluvstein et al. 2024)
- Gaussian BosonSampling (Zhong et al. 2020, Zhong et al. 2021, Madsen et al. 2022, Deng et al. 2023)

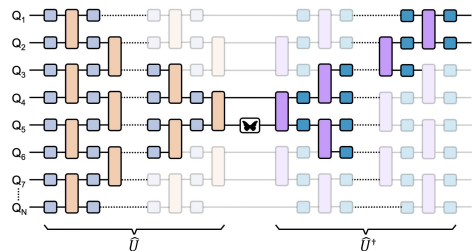Early demonstrations of **quantum error correction**:
- Surface code implementations (Krinner et al. 2022, Zhao et al. 2022)
- Surface code error suppression (Google 2022)
- Other codes (Ofek et al. 2016, Fluhmann et al. 2019, Champagne-Ibarcq et al. 2020, Grimm et al. 2020, Chen et al. 2021, Egan et al. 2021, Ryan-Anderson et al. 2021, Sundaresan et al. 2022)
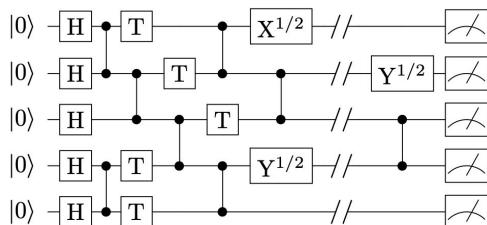
**QEC**

Quantum AI

# Two use cases for tensor networks

## NISQ

Benchmarking experiments with tensor networks
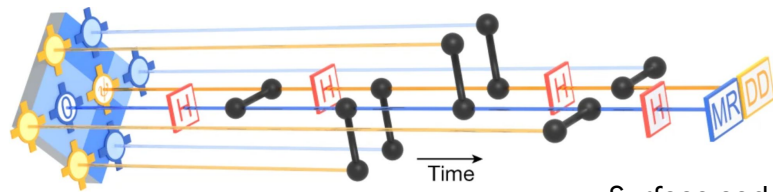
OTOC
(Mi, et al. 2021)

RCS
(Boixo, et al. 2017)

**Exploiting structure:**
compressibility, low
entanglement, ...
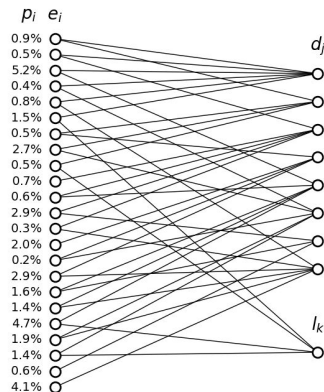
**Worst case:**
brute force contraction

## QEC

Decoding (and benchmarking)  with tensor networks

Time

Surface code
(Google. 2023)

- **Mapping** decoding
  to TN contraction

- **Contract** efficiently

Quantum AI

# Outline

Quantum AI

# A simulation problem (1 / 2)

**Experiment**



- Observable

- Probability amplitude

- Samples

- ...

Characterization:

- Fidelity?
- Is the experiment giving right results?
- What kind of results do we expect?
- ...

Challenging beyond-classical claims:

- What's the classical computational cost?
- What are the hardness guarantees?
- Is the experiment beyond-classical?
- ...

(classical) **Simulation**

Quantum AI

# A simulation problem (2 / 2)

Under special circumstances there are specialized techniques:
- Clifford circuits
- Clifford + T circuits
- Matchgate circuits
- Localized dynamics
- Large noise rates (hinder entanglement formation)
- ...

In the generic case we need brute force

**Circuit** ⟷ **Tensor network**

- Observable
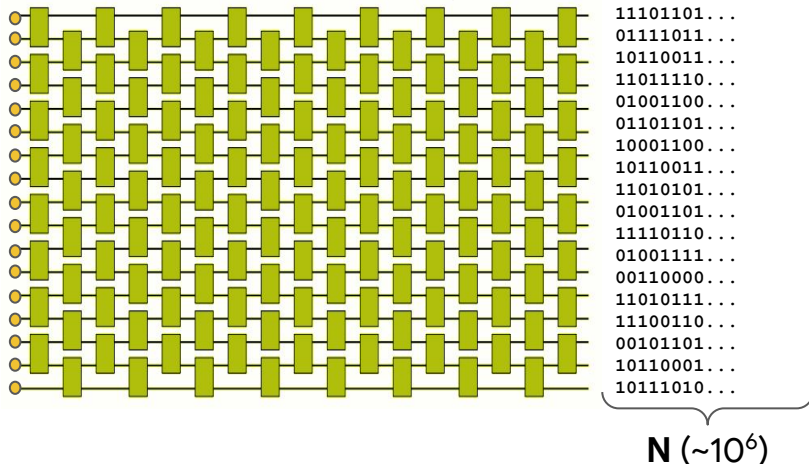- Probability amplitude
- Samples
- ...

As **quantity** or **primitive** for it

Quantum AI

# Making brute force less brute: the RCS case study (1/3)

## Random Circuit Sampling (RCS)



```
11101101...
01111011...
10110011...
11011110...
01001100...
01101101...
10001100...
10110011...
11010101...
01001101...
11110110...
01001111...
00110000...
11010111...
11100110...
00101101...
10110001...
10111010...
```
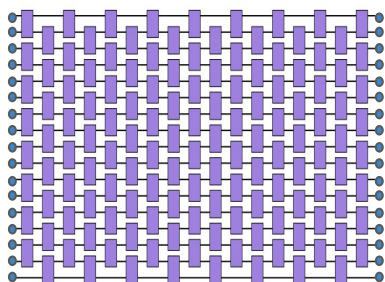
**N** (~$10^6$)

Estimate fidelity **f** (from samples)
**f** > 0 within a few **σ's**?

Can a classical computer perform this task (*in a reasonable amount of time*)?

Fairly strong complexity theory guaranties for the hardness of this task.

(Boixo et al. 2016, Aaronson et al. 2016, Bouland et al. 2019 & 2021, Movassagh et al. 2020, ...)

## Tensor network



*Classical adversary*

Sampling algorithm (Markov et al. 2018):

1. Compute **p(x)** for bit strings **x** chosen uniformly at random
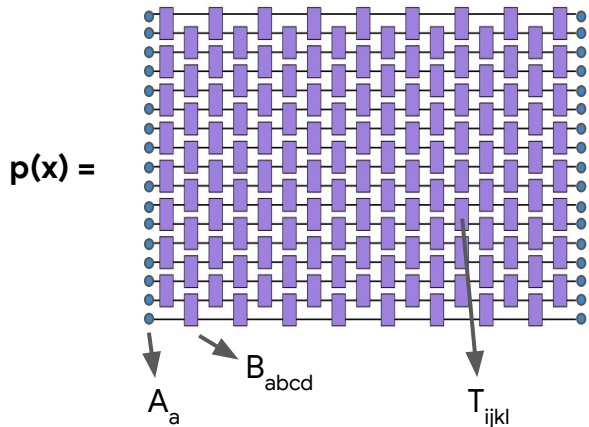2. Accept **x** with probability **p(x)N/M**

Acceptance ratio **1/M ~ 1/10**

Modified rejection sampling: **frugal sampling**

```
11101101...
01111011...
10110011...
11011110...
01001100...
01101101...
10001100...
10110011...
11010101...
01001101...
11110110...
01001111...
00110000...
11010111...
11100110...
00101101...
10110001...
10111010...
```

Quantum AI

p(x) =

$$\sum_{c}\sum_{a}\sum_{f}\sum_{b}\ldots A_a B_{abcd}\ldots T_{ijkl}\ldots$$

$A_a$

$B_{abcd}$

$T_{ijkl}$

Order of contraction dramatically affects computational cost.
Time and memory requirements lower bounded by **treewidth of line graph**
(Markov & Shi 2008)

**Goal:** optimize tensor network contraction ***ordering*** **(O)**. (Gray & Kourtis 2020)
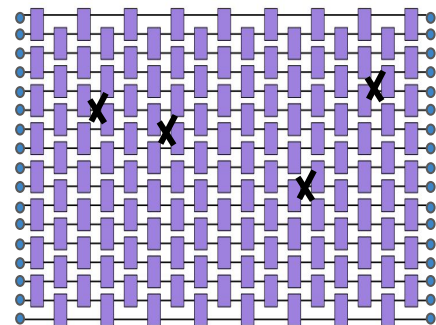
**Memory?**
For current experiments this leads memory requirements ~$10^4\times$
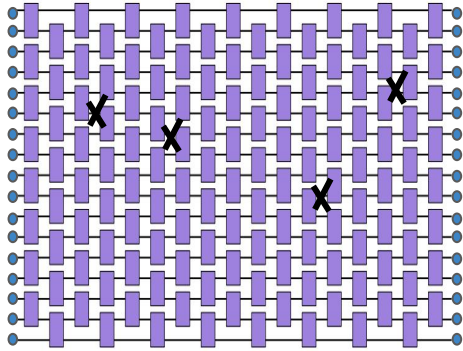the total memory of the largest supercomputer on Earth

**Solution:** project (slice) a set of variables **(S)** and perform sum *a posteriori* (Alibaba, 2018).
✅Alleviates memory requirements *and* parallelizes computation
❌Need to contract an exponential number TNs

Quantum AI

# Making brute force less brute: the RCS case study



**Optimization problem:**

Contraction cost is function of ordering and slices **C(O, S)**.
Memory usage **M(O, S)<M\*** (total memory available).

**Pedantic detail of our encoding:**

We take slices from an ordered set of candidates **P** until **M<M\*** is satisfied, so our cost function is really **C(O, P)**.
Contraction cost is function of ordering and slices **C(O, S)**.

Plethora of "tricks" from the literature can be beneficial in practice:

- **Sparse output:** so millions of amplitudes can be computed with a single contraction (**F. Pan** et al. 2021)
- **Details of hardware gates:** `fSim` gate can be exploited for faster contractions (Google 2019 & **F. Pan** et al. 2021)
- **Memoization:** reuse of intermediate computations across branches of the computation (Kalachev et al. 2021)
- **Experimental fidelity:** low target fidelity speeds up simulation (Markov et al. 2018, Villalonga et al. 2019)

All these accounted for in **C(O, P)**.
Highly optimized evaluation of **C(O, P)**. Current experiments are close to ~1000 two-qubit gates.

Quantum AI

# Example results (1/2)

**Optimized runtimes for RCS experiments:**

| Exp. | 1 amp. FLOPs | 1 million noisy samples | | |
|---|---|---|---|---|
| | | FLOPs | XEB fid. | Time |
| SYC-53 [4] | $6 \times 10^{17}$ | $2 \times 10^{17}$ | $2 \times 10^{-3}$ | 6 s |
| ZCZ-56 [5] | $6 \times 10^{19}$ | $6 \times 10^{19}$ | $6 \times 10^{-4}$ | 20 min |
| ZCZ-60 [6] | $1 \times 10^{21}$ | $1 \times 10^{23}$ | $3 \times 10^{-4}$ | 40 days |
| SYC-70 | $5 \times 10^{23}$ | $6 \times 10^{25}$ | $2 \times 10^{-3}$ | 50 yr |
| SYC-67 | $2 \times 10^{23}$ | $2 \times 10^{37}$ $2 \times 10^{28}$ $2 \times 10^{25}$ | $1 \times 10^{-3}$ | $1 \times 10^{13}$ yr $1 \times 10^{4}$ yr* 12 yr** |

Google, 2023

Parallelizing over independent GPUs on Frontier

*Assuming distributed contractions over all RAM.
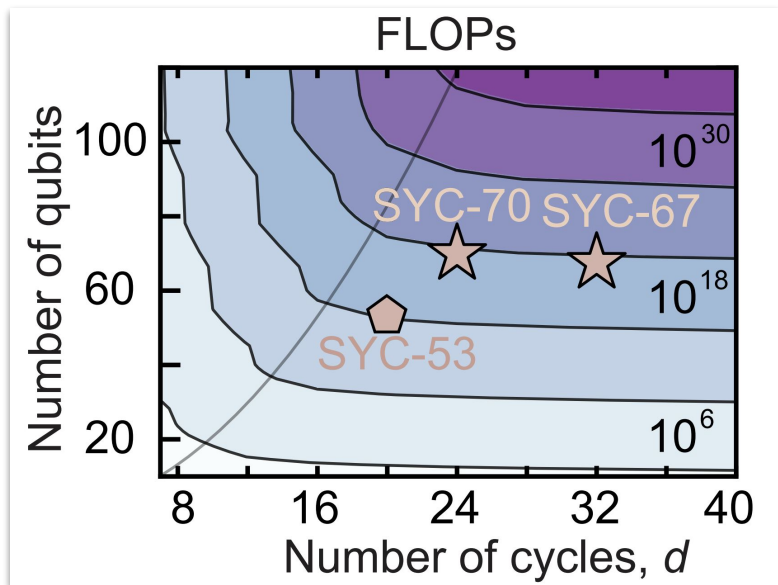
**Assuing distributed contractions using secondary storage.

*&** without inter-node communication times (for stronger adversary against BC claim)

Quantum AI

# Example results (1/2)

**Complexity vs. circuit size:**



Google, 2023

2D architecture (**L × L**) similar to experiment

- At low depth, cost **exp(d × L)**
- At large depth, cost **exp(L × L) = exp(#qubits)**

Quantum AI

# The future of NISQ applications: noise vs. computational volume

Signal (fidelity) decreases exponentially with volume of computation (for generic circuits, ~#two-qubit gates).

Computations are limited to finite sizes, which limits their classical computational cost.

**RCS** experiments beyond classical?

**Strongly established**

**Useful / physical** experiments beyond classical?

**Not yet**

Strongly supported by highly optimized TN contraction results

**Will there be a useful NISQ application before QEC is achievable?**

# Outline

Quantum AI

# The setup: 3 performance contributing factors (1/2)
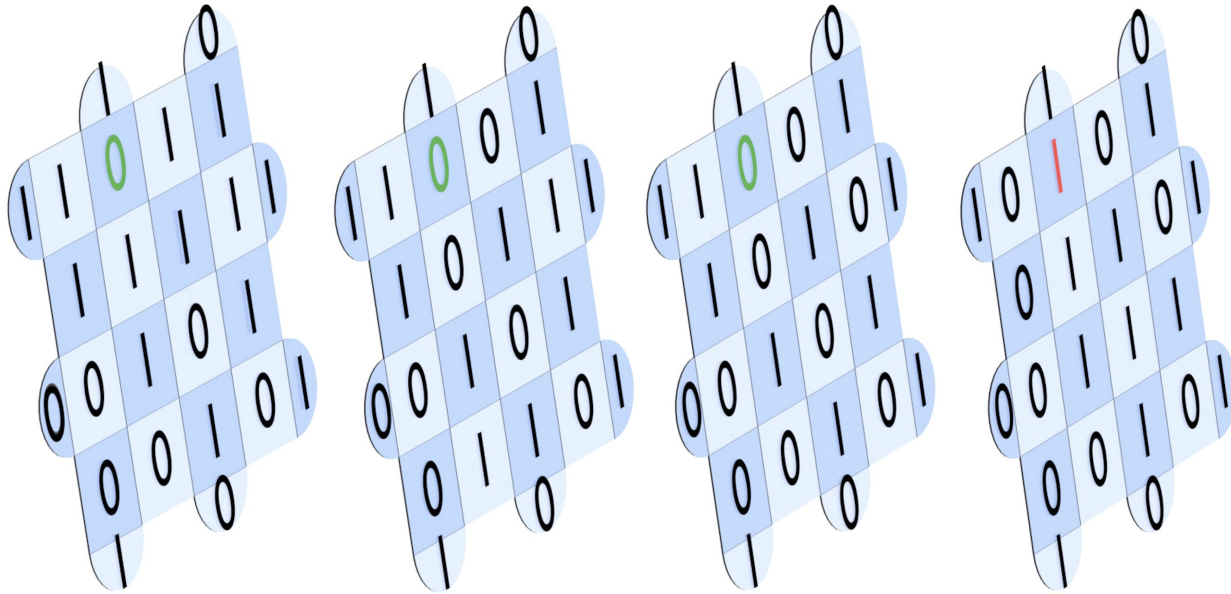


For low enough error rates:
- Encode logical operators (qubit) over many physical qubits
- Early demonstration with **memory experiment:**
  - Initialize system in eigenstate of **X** or **Z**
  - Run several rounds of surface code, each one measuring parity checks (operators)
  - Decode: infer from parity checks whether logical operator has changed value
- Decoding has as input an understanding of physical errors: **error model**

What determines the quality of the experiment (of the logical qubit)?
- Hardware: roughly physical error rates
- Error model
- **Decoder**

Quantum AI

$0 \xleftrightarrow{\text{detector}} 0 \xleftrightarrow{\text{detector}} 0 \xleftrightarrow{\substack{\text{detector} \\ \text{detection} \\ \text{event}}} 1$
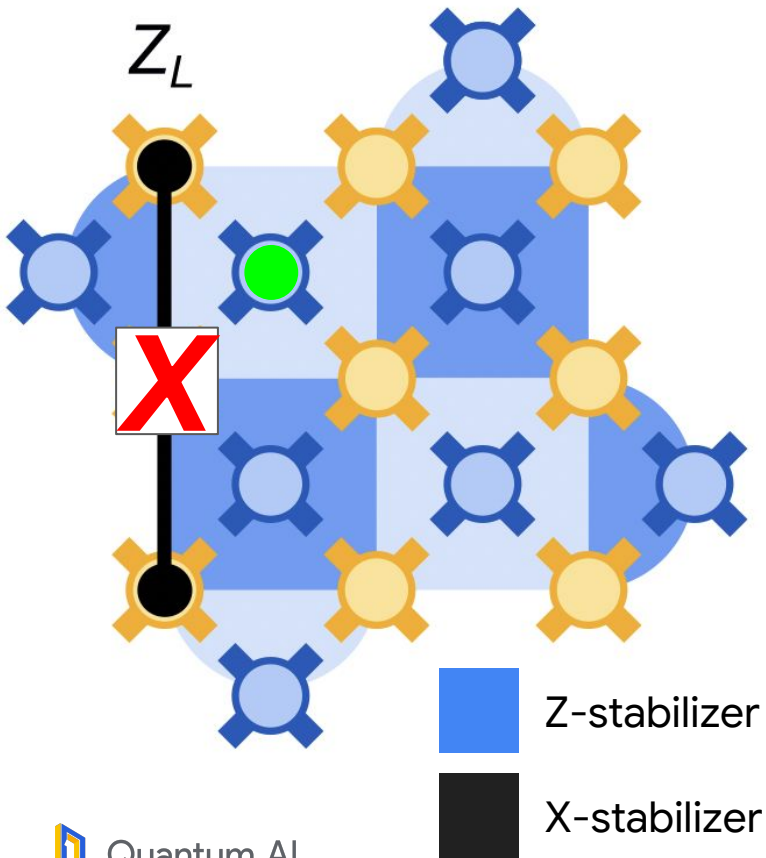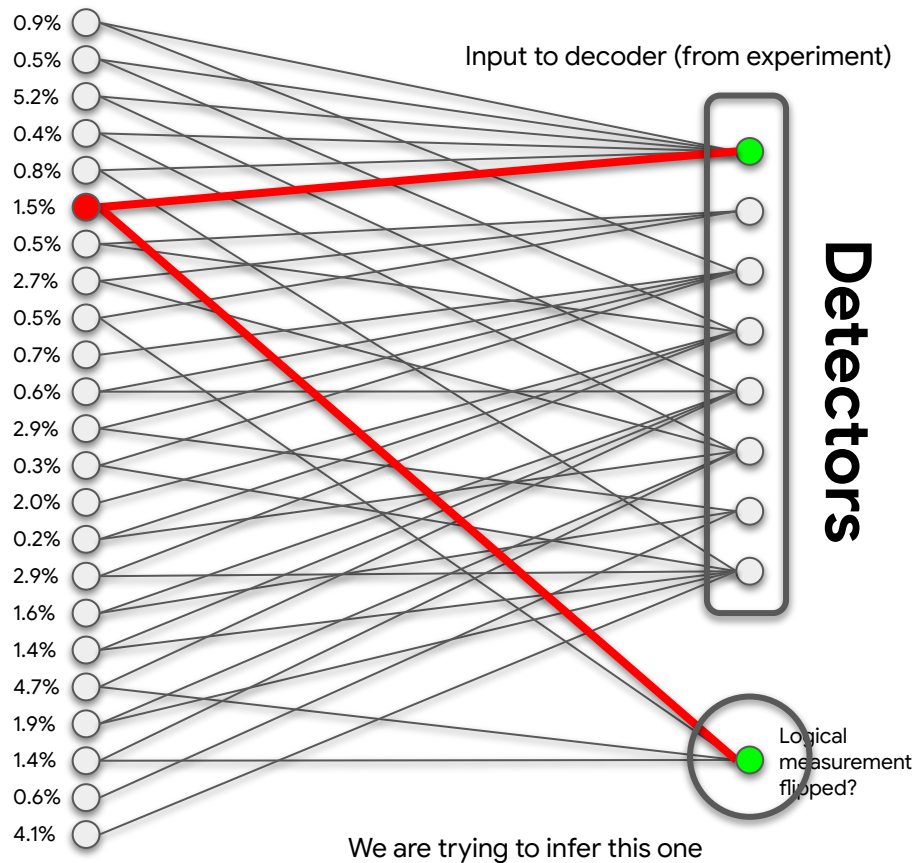
**Detectors** — comparisons of measurements that should agree

**Detection events** — when they don't

From detection events, can we make a guess on the the logical operator flipped/not flipped?

Quantum AI

# The hyper-graph error model

$Z_L$

X

Z-stabilizer

X-stabilizer

Quantum AI

Error mechanisms

0.9%
0.5%
5.2%
0.4%
0.8%
1.5%
0.5%
2.7%
0.5%
0.7%
0.6%
2.9%
0.3%
2.0%
0.2%
2.9%
1.6%
1.4%
4.7%
1.9%
1.4%
0.6%
4.1%

Detectors
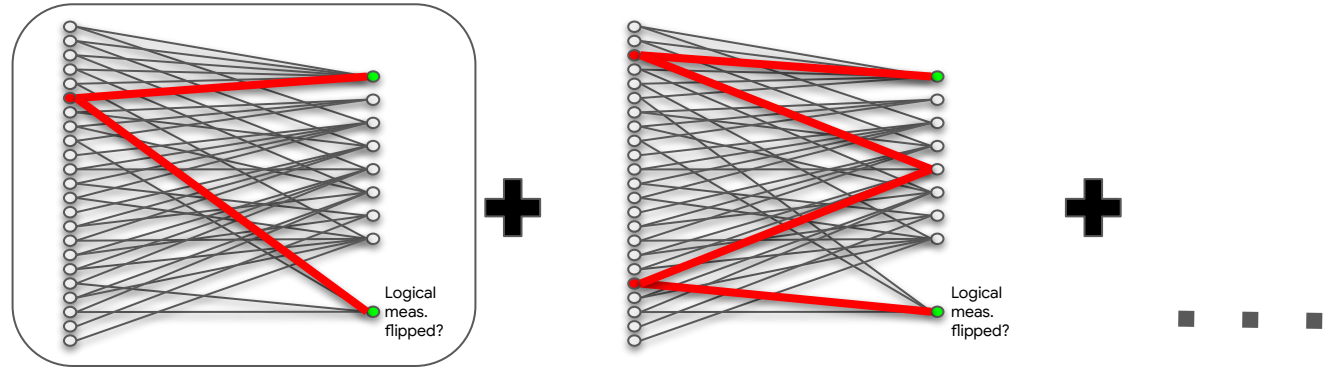
Input to decoder (from experiment)
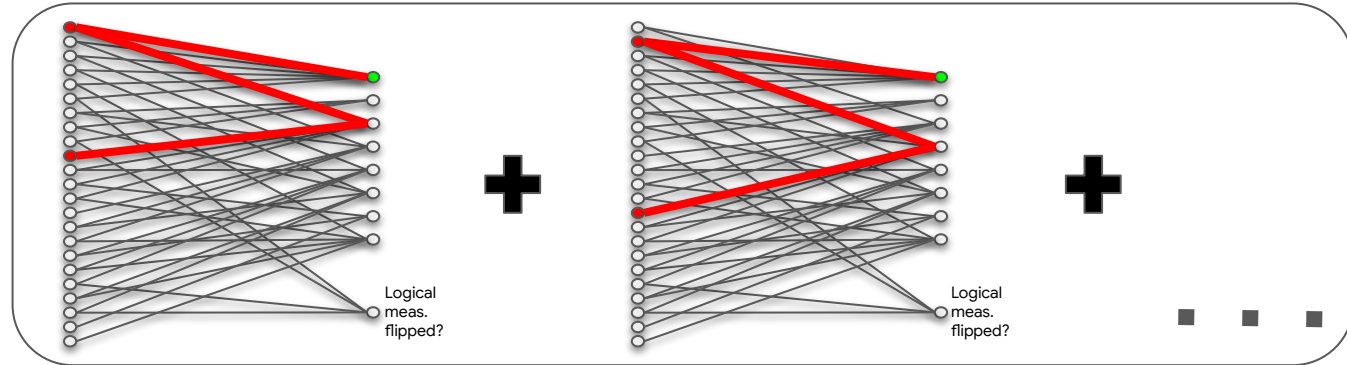
Logical measurement flipped?

We are trying to infer this one

# The (maximum-likelihood) decoding problem

Most likely error...

... might not belong to the most likely **set** of errors.



Logical meas. flipped?

Logical meas. flipped?

Logical meas. flipped?

Logical meas. flipped?

Maximum likelihood decoding = *optimal* decoding

Quantum AI

Error model ($p_i$+ graph)

$p_i$  $e_i$

Hidden ($e_i$)

Experiment ($d_j$)

$d_j$

$l_k$

Trying to infer ($l_k$)
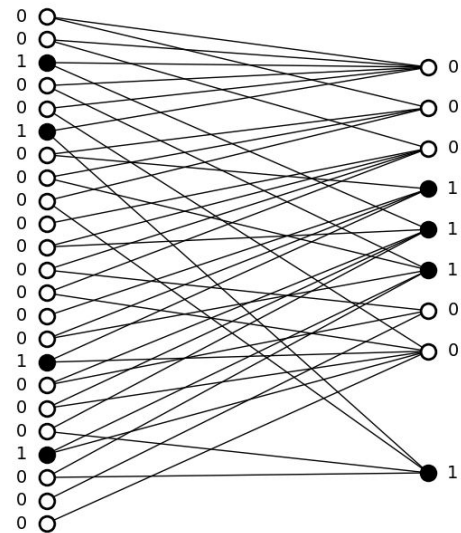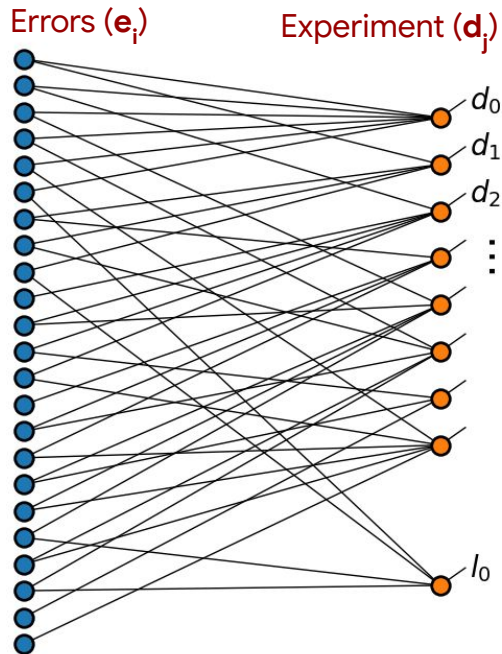
$$\Pr(\vec{e}) = \prod_i p_i^{e_i} \cdot (1 - p_i)^{1-e_i}$$

$$\vec{d} = \vec{f}(\vec{e}) \text{ and } \vec{l} = \vec{g}(\vec{e})$$

$$L\left(\vec{l}\,|\vec{d}\right) \propto \sum_{\vec{e}:\left[\vec{f}(\vec{e})=\vec{d}\right]\wedge\left[\vec{g}(\vec{e})=\vec{l}\right]} \Pr(\vec{e})$$

Quantum AI

# A tensor network ML decoder for all hyper-graph error models

(Initial proponent: Bravyi et al. 2014)

**Errors ($e_i$)**      **Experiment ($d_j$)**

$$\mathrm{Pr}(\vec{e}) = \prod_i p_i^{e_i} \cdot (1 - p_i)^{1 - e_i}$$

$$\vec{d} = \vec{f}(\vec{e}) \text{ and } \vec{l} = \vec{g}(\vec{e})$$

$$L\left(\vec{l}|\vec{d}\right) \propto \sum_{\vec{e}:[\vec{f}(\vec{e}) = \vec{d}] \wedge [\vec{g}(\vec{e}) = \vec{l}]} \mathrm{Pr}(\vec{e})$$

$$L\left(l_0|\vec{d}\right) =$$



$$= \begin{cases} p_i & \text{if } \alpha_0 = \alpha_1 = \ldots = 1 \\ 1 - p_i & \text{if } \alpha_0 = \alpha_1 = \ldots = 0 \\ 0 & \text{otherwise} \end{cases}$$

Propagates error to:
- Detectors
- Logical operator(s)

$$= \begin{cases} 1 & \text{if } \alpha_0 + \alpha_1 + \ldots \text{ even} \\ 0 & \text{if } \alpha_0 + \alpha_1 + \ldots \text{ odd} \end{cases}$$
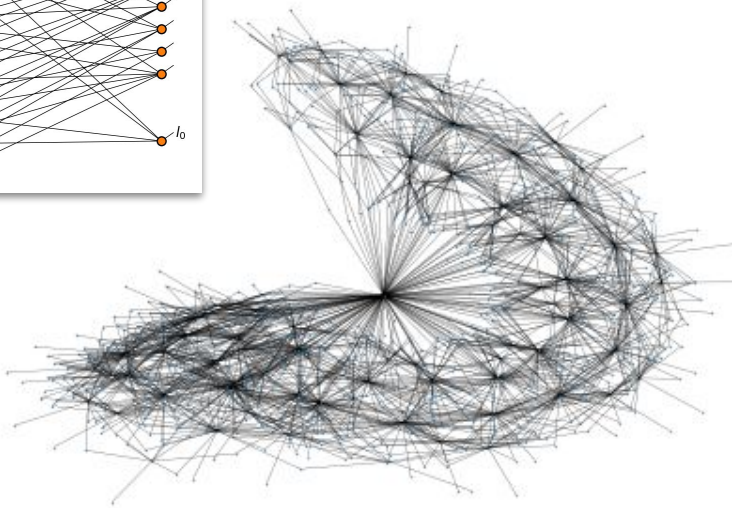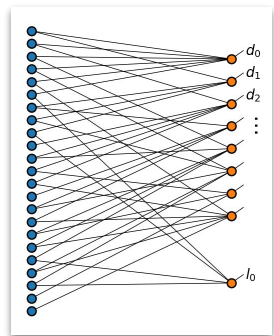
Two ways of seeing it:
- Enforces right parity with $d_j$ and $l_k$
- Kills error configurations that violate constraints

(Piveteau et al. 2023 also uses error hyper-graph as starting point)
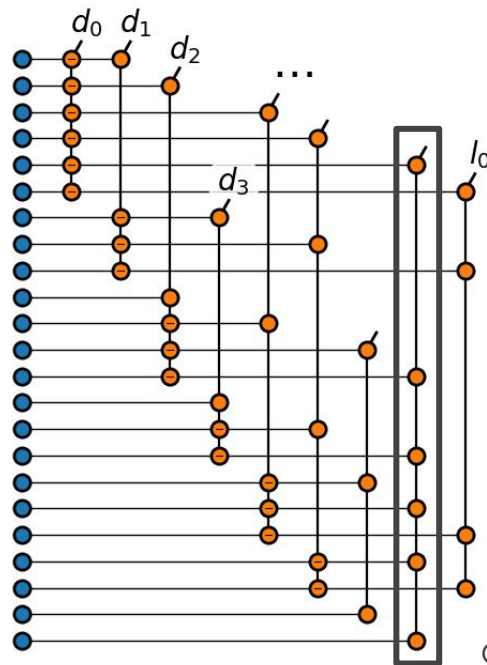
Decoding:

$$L\left(l_0 = 0|\vec{d}\right) \geq L\left(l_0 = 1|\vec{d}\right)$$
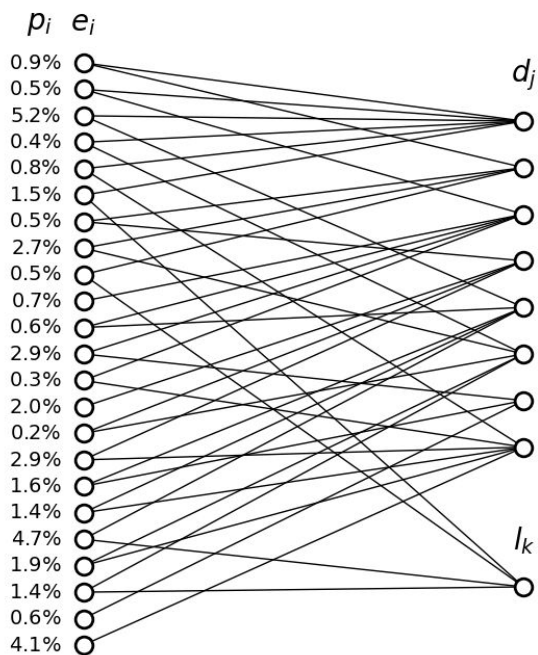
**Exact contraction not scalable**

**Approximate contraction?**

$$\beta = \begin{cases} 1 & \text{if } (\alpha_0 + \alpha_1 + \dots \text{ even}) \\ & \wedge (\alpha_x = \beta) \\ 0 & \text{otherwise} \end{cases}$$

Graph structure is **adjacency matrix** of error hyper-graph

Quantum AI

$$L\left(l_0|\vec{d}\right) =$$

error: **(1-p$_i$ , p)**

parity

parity & copy

copy & copy

$$\alpha_0 \quad \square \quad \alpha_1 \; = \; \begin{cases} 1 & \text{if } (\alpha_0 = \alpha_1) \;\wedge\; (\beta_0 = \beta_1) \\ 0 & \text{otherwise} \end{cases}$$

Quantum AI

# A tensor network ML decoder for all hyper-graph error models (5/5)



$p_i$ $e_i$

0.9%
0.5%
5.2%
0.4%
0.8%
1.5%
0.5%
2.7%
0.5%
0.7%
0.6%
2.9%
0.3%
2.0%
0.2%
2.9%
1.6%
1.4%
4.7%
1.9%
1.4%
0.6%
4.1%

$d_j$

$l_k$

$d_0$ $d_1$ $d_2$ $d_3$ ...

$l_0$

error: **(1-p$_i$ , p)**

$\alpha_0$
$\alpha_1$

parity

$\alpha_1$ $\alpha_2$ $\alpha_0$

parity & copy

$\alpha_0$ $\alpha_1$ $\alpha_x$ $\beta$

copy & copy

$\beta_0$ $\alpha_0$ $\alpha_1$ $\beta_1$

**Approximate contraction:**
MPS evolution with finite $\chi$
(left to right)

Decoding:

$$L\left(l_0 = 0 | \vec{d}\right) \geq L\left(l_0 = 1 | \vec{d}\right)$$

**?**

Quantum AI

# Results (1/2)

Milestone experiment on error suppression using the surface code
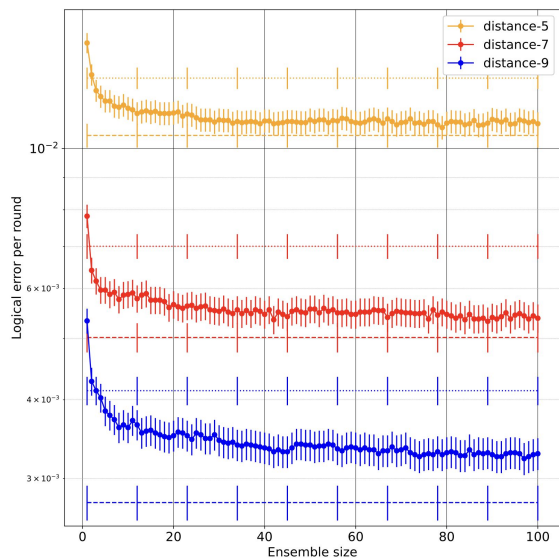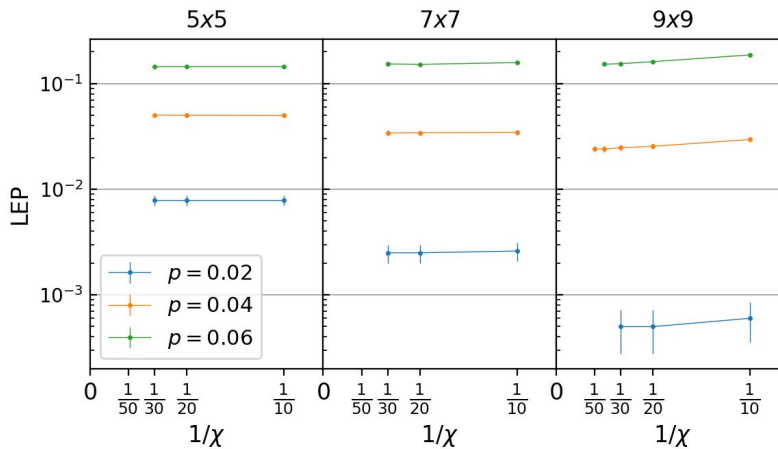Google, 2023 - *Nature* 614, no. 7949 (2023): 676-681

# Results (2/2)

Benchmarking performance of faster / scalable decoders
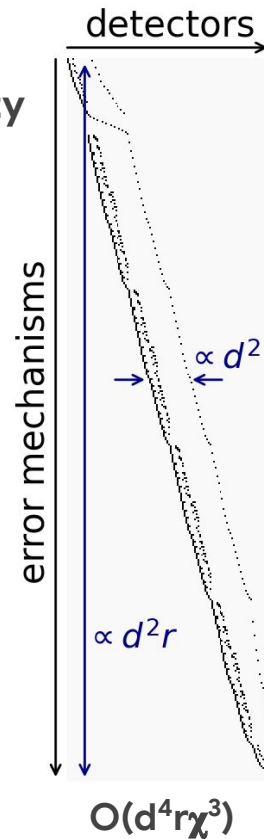N. Shutty, M. Newman, **BV**, 2024 - *arXiv:2401.12434* (2024)

**Benchmark of Harmony**



**Convergence**



**Complexity**



$O(d^4 r \chi^3)$

# Outline

Quantum AI

# Conclusion

Two applications of tensor networks to experimental quantum computing:

- Highly-optimized tensor network contraction for benchmarking NISQ experiments:
  - Strong evidence for RCS being beyond classical
  - Insightful method to challenge useful beyond-classical claims
- Decoding for QEC:
  - Decode *arbitrary* error hyper-graph codes
  - Benchmark experimental hardware and error model quality
  - Benchmark performance of fast, scalable decoders

**References**

Latest RCS paper: Google, *arXiv:2304.11119* (2023)

Surface code error suppression: Google, *Nature* 614, no. 7949 (2023): 676-681

Harmony decoding: Noah Shutty, Michael Newman, and Benjamin Villalonga, *arXiv:2401.12434* (2024)

Quantum AI