



Clique: Better Than Worst-Case Decoding for Quantum Error Correction

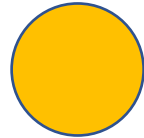
Gokul Subramanian Ravi¹, Jonathan M. Baker¹, Arash Fayyazi², Sophia Fuhui Lin¹,
Ali Javadi-Abhari³, Massoud Pedram², Frederic T. Chong¹

1: UChicago, 2: USC, 3: IBM

*Best of both worlds approach,
combining two schools of QEC
decoding, for 100x–10,000x gains!!*

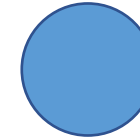
Error correction for fault tolerant quantum systems

Physical qubit

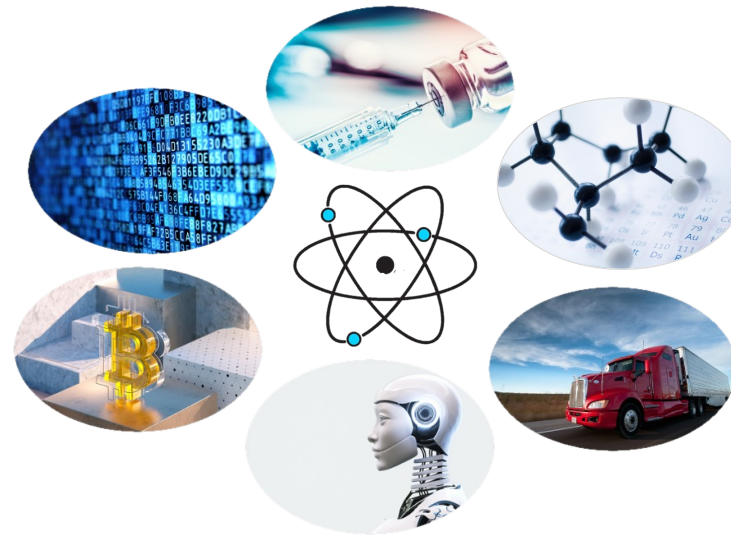


ER = $10^{-2} - 10^{-4}$

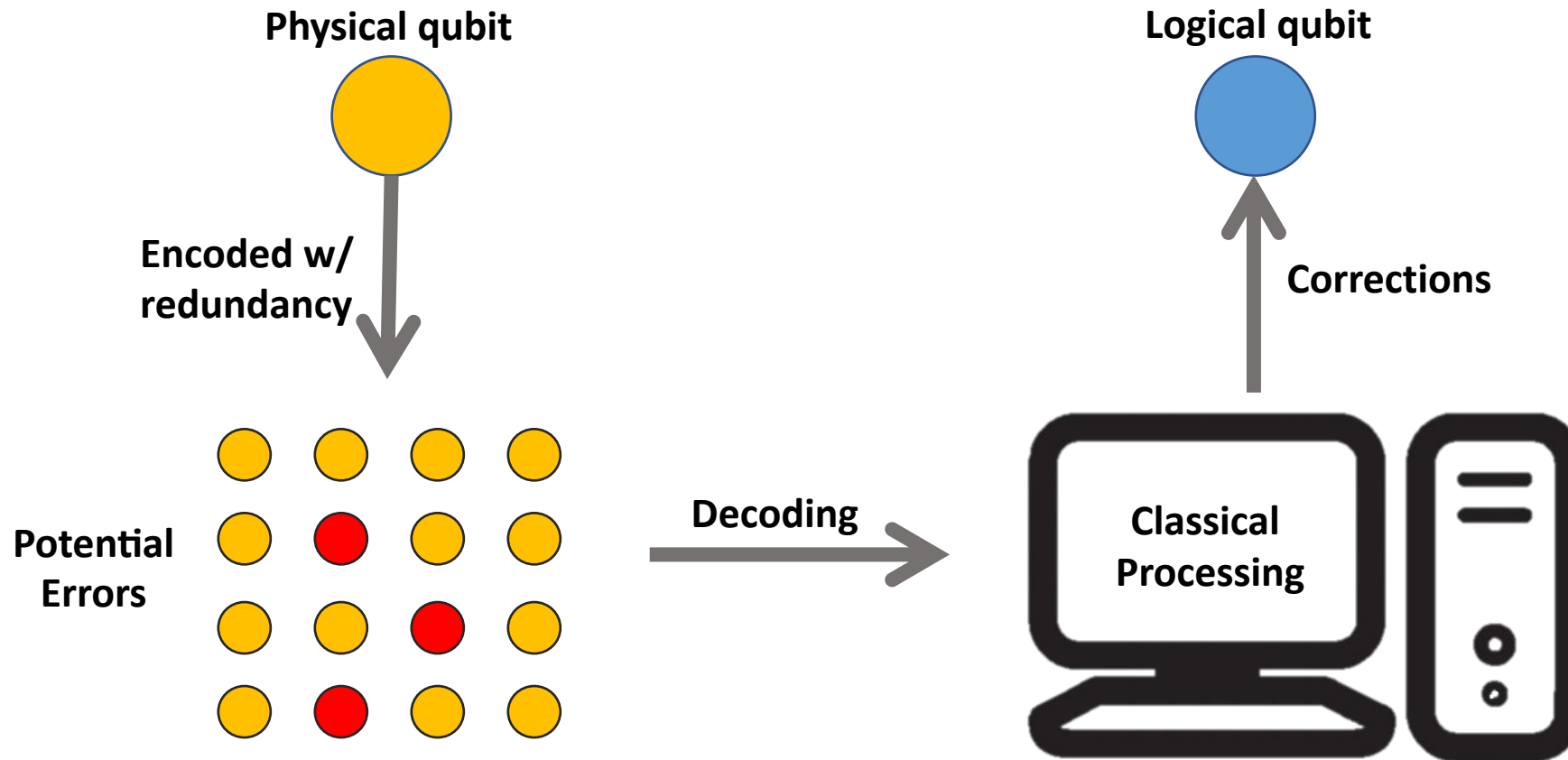
Logical qubit



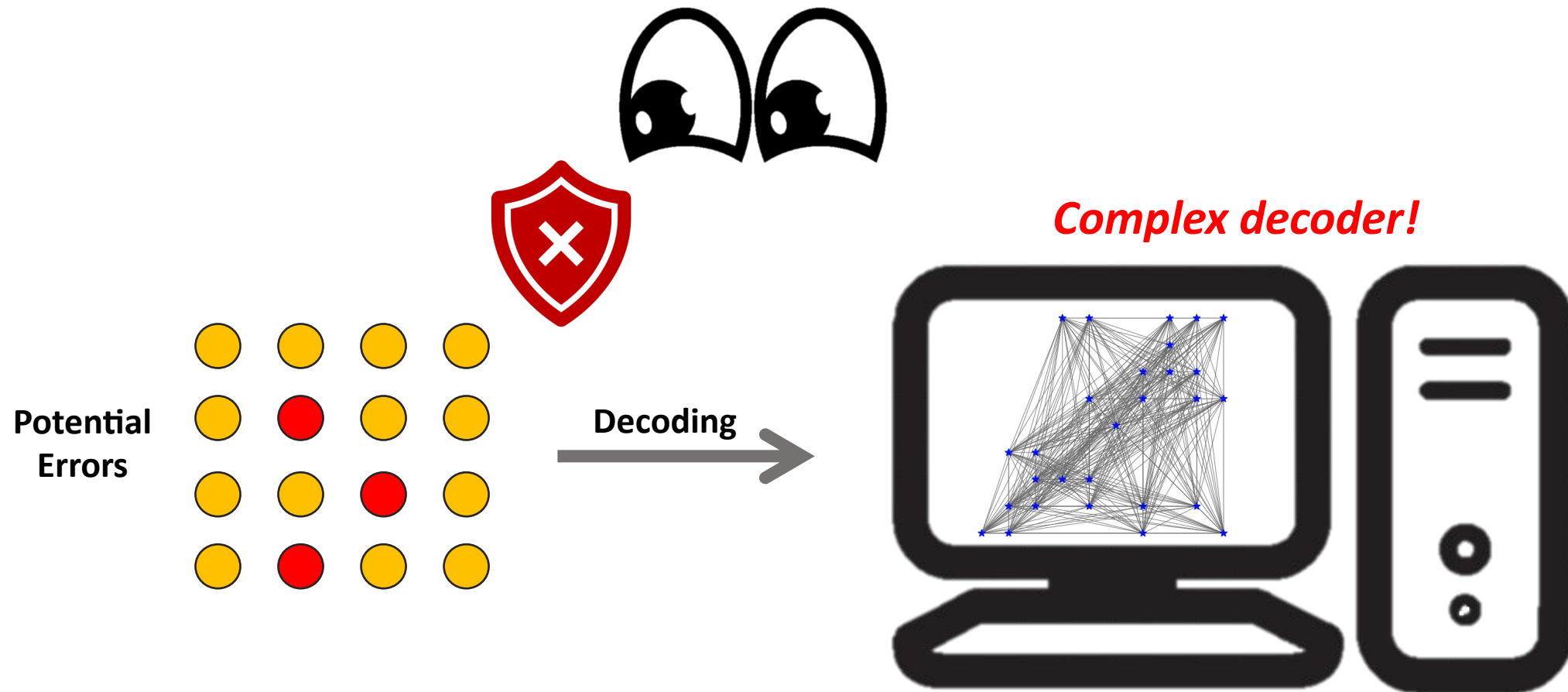
ER = $10^{-6} - 10^{-15}$



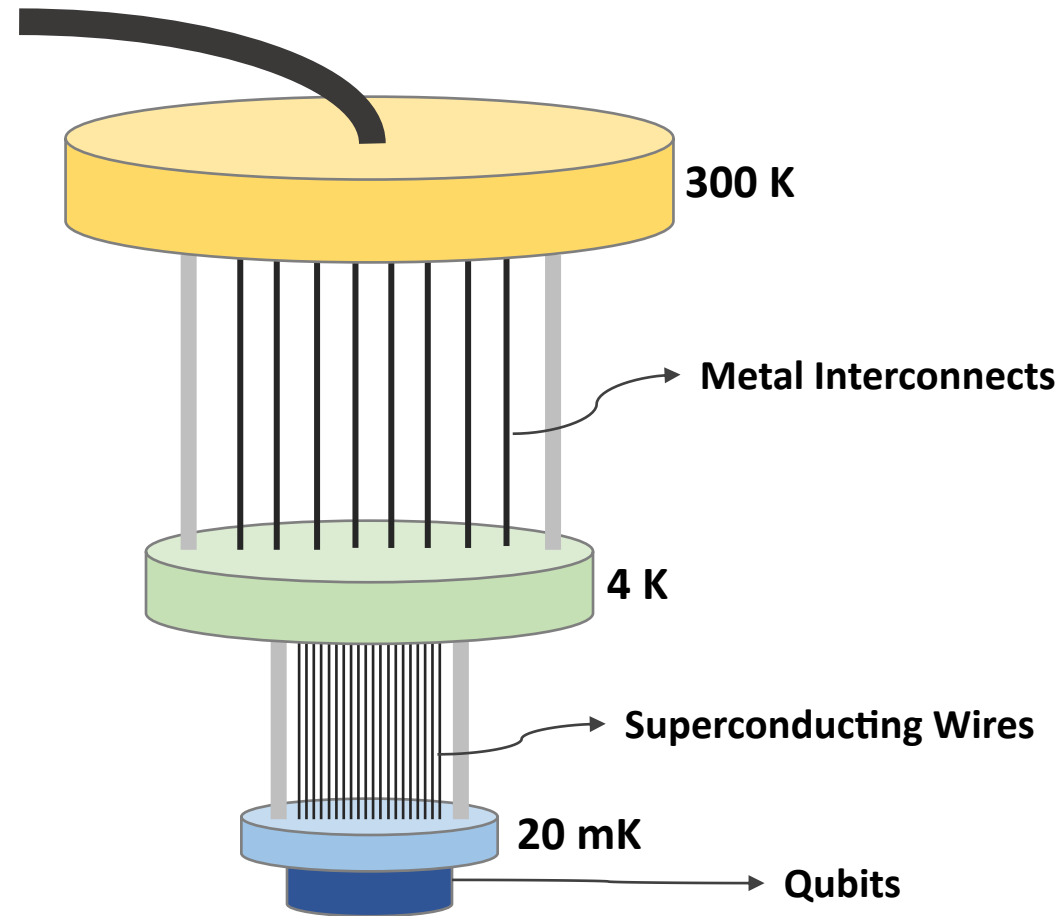
Error correction for fault tolerant quantum systems



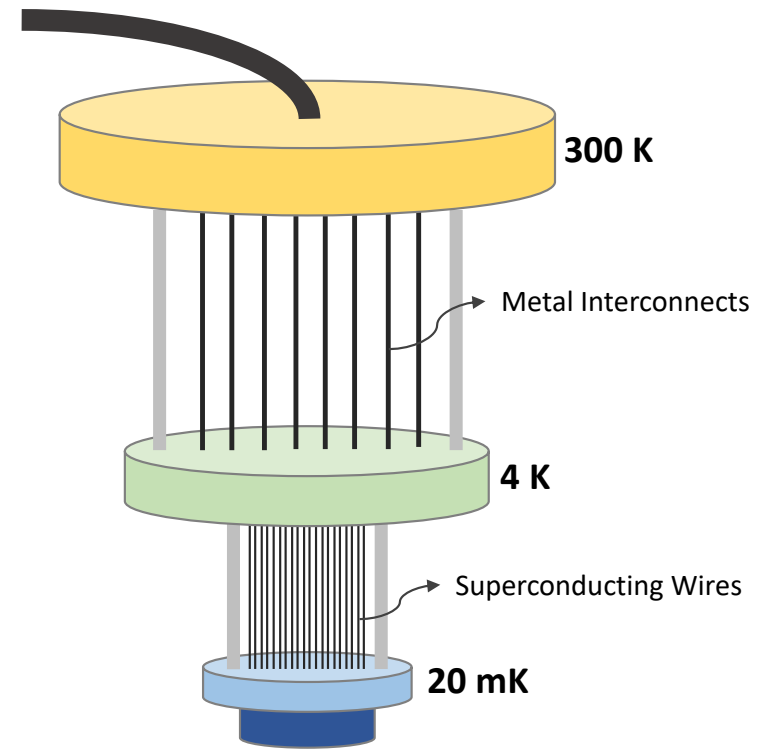
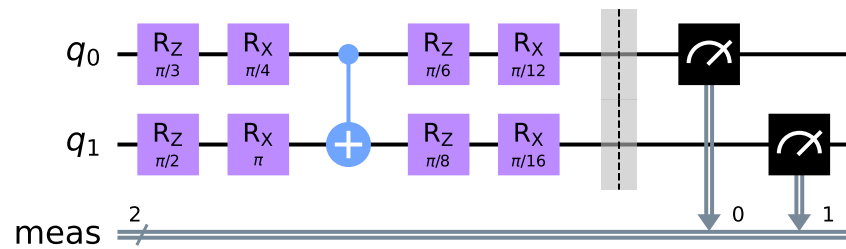
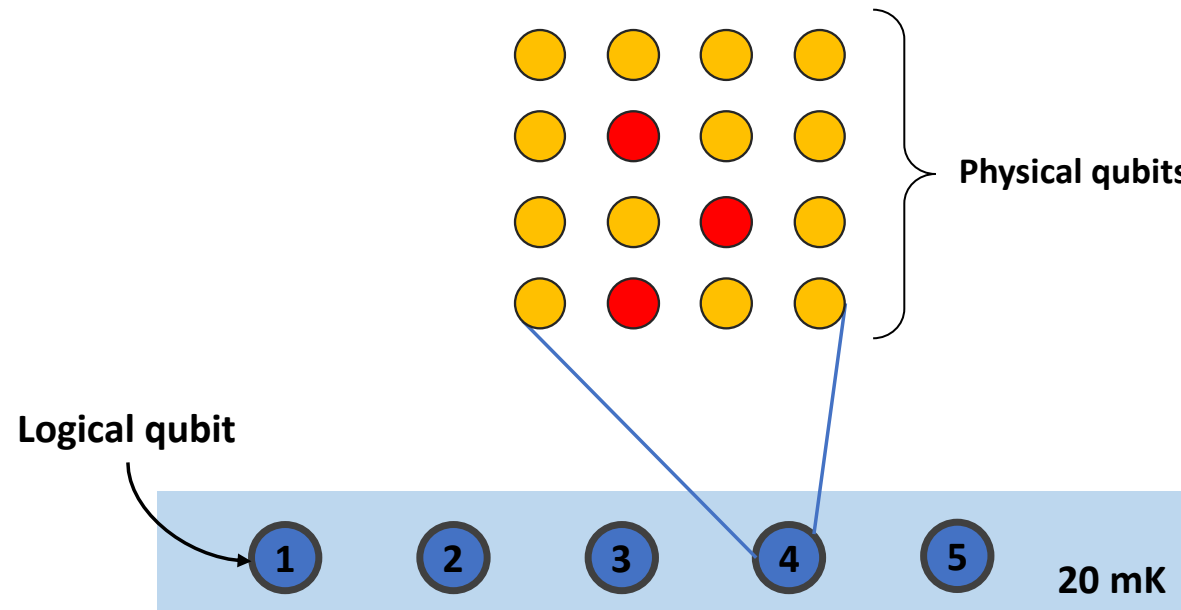
Error correction for fault tolerant quantum systems



Scope: Cryogenic quantum systems



Scope: Cryogenic quantum systems



System-level view: Traditional outside-fridge QEC decoding

Tbps I/O bandwidth → **bandwidth bottleneck!**

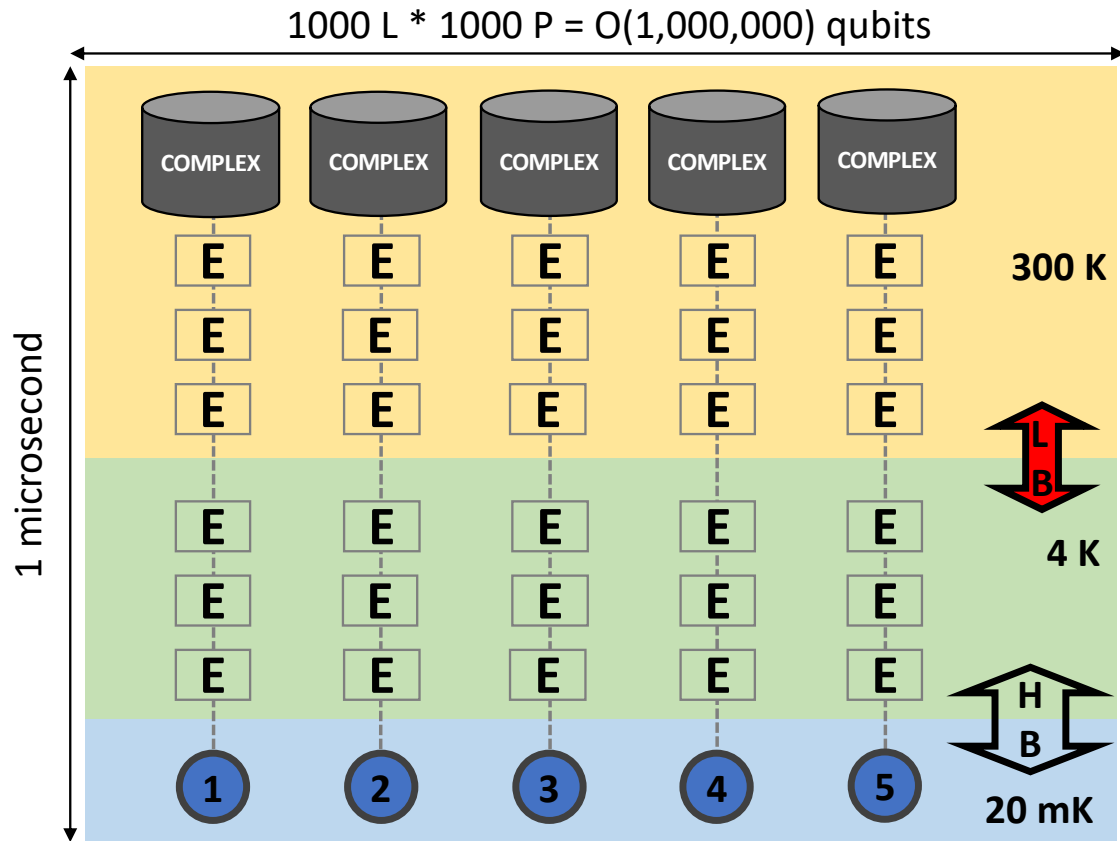
[Fowler, PR-A '12]

..

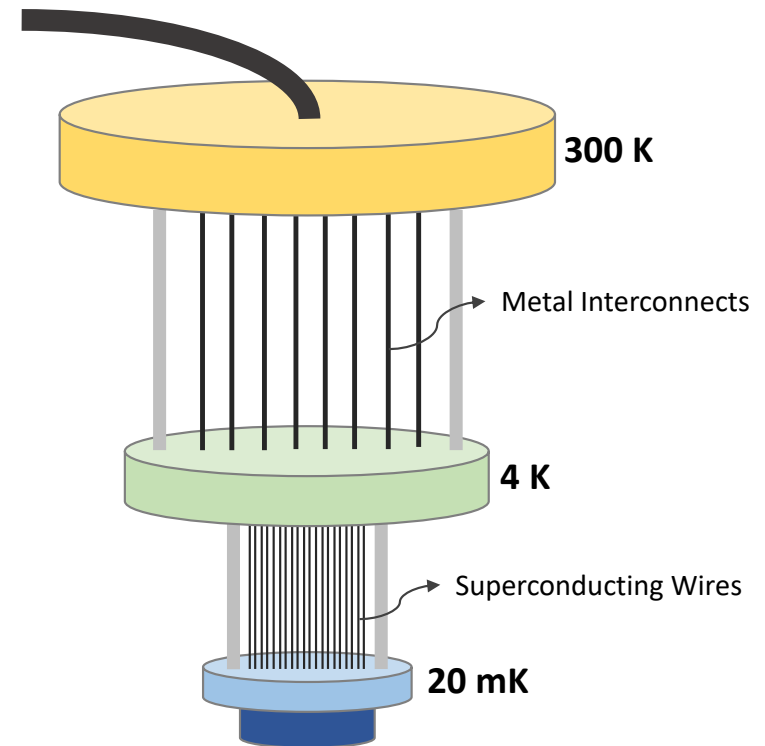
..

[Das, HPCA '22]

Low latency needed for functional correctness!!



E: Error Signatures



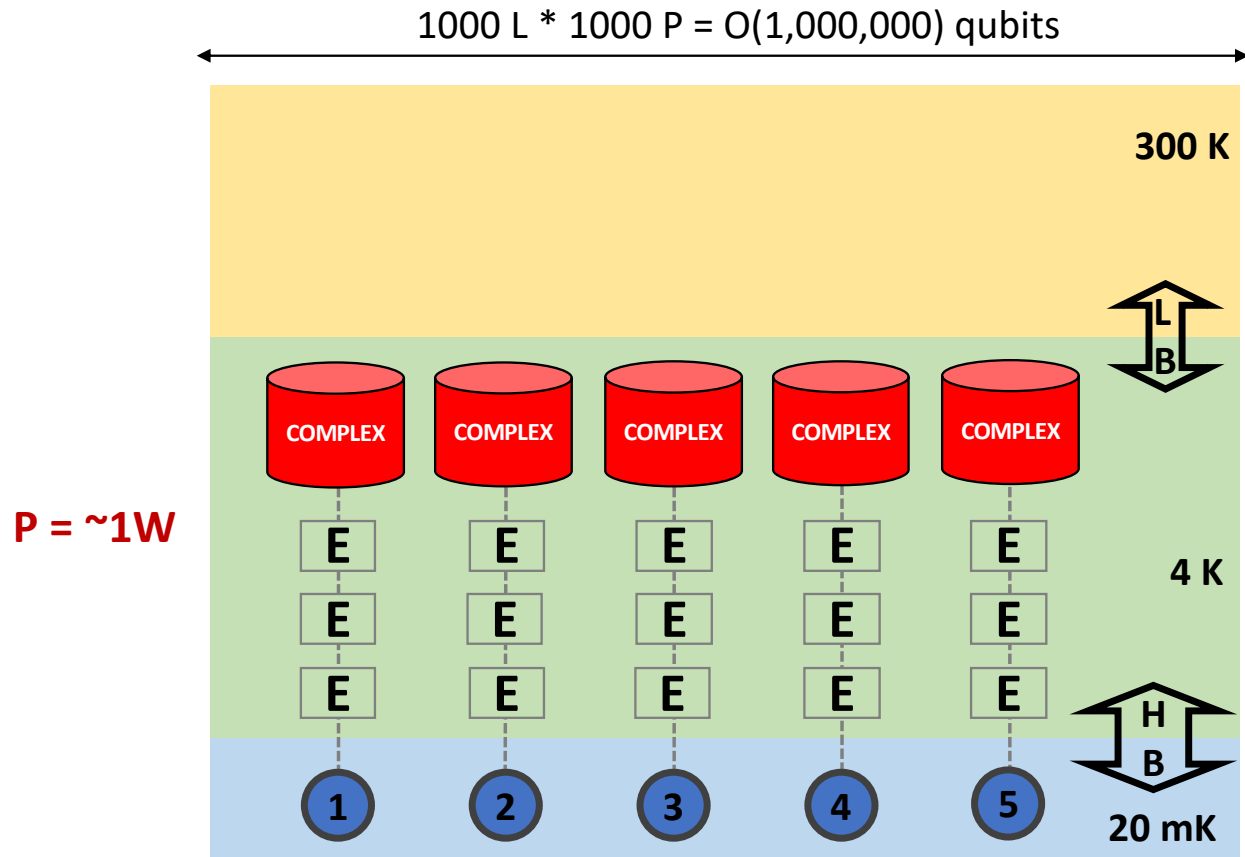
System-level view: Cryogenic inside-fridge QEC decoding

Limited cryogenic power budget (~1W) cryo-resource bottleneck!

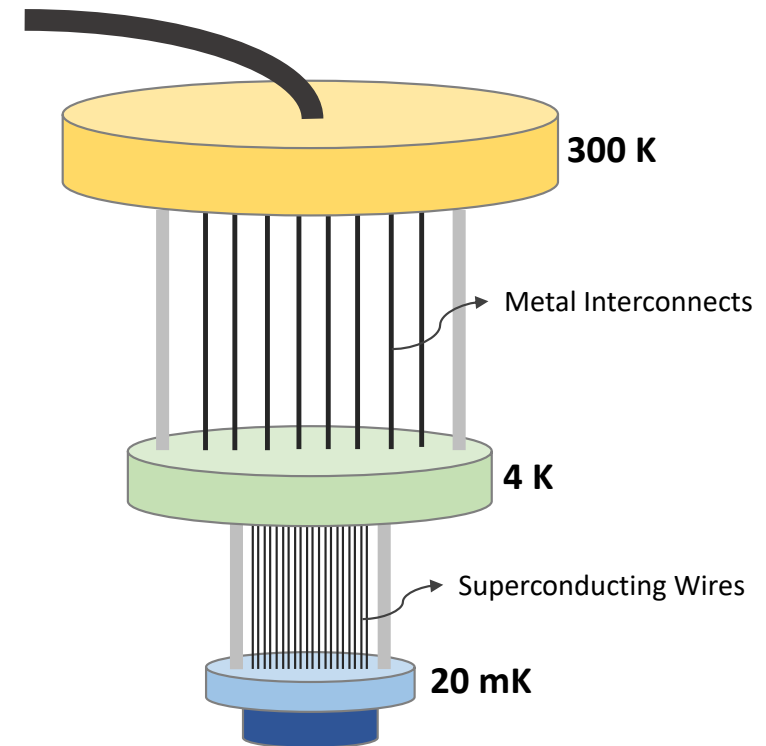
[Holmes, ISCA '20]

[Byun, ISCA '22]

[Ueno, HPCA '22]

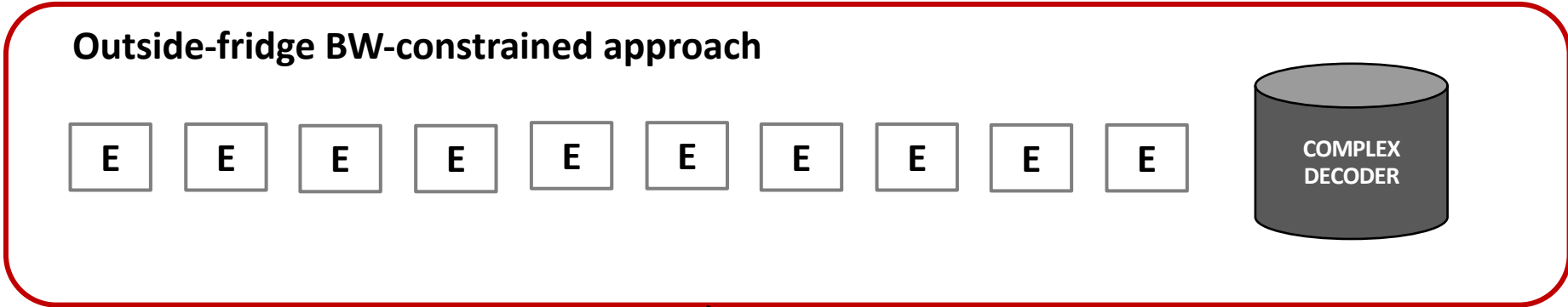


E: Error Signatures

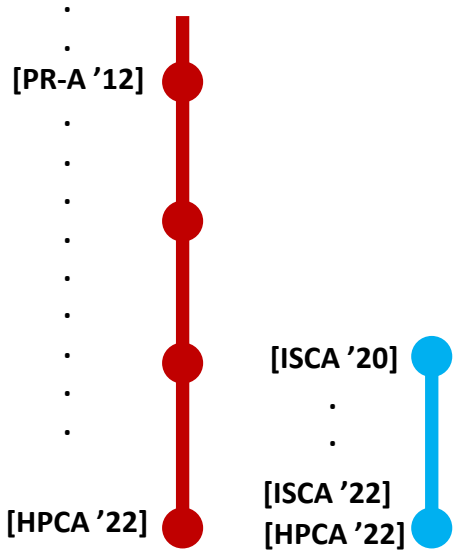
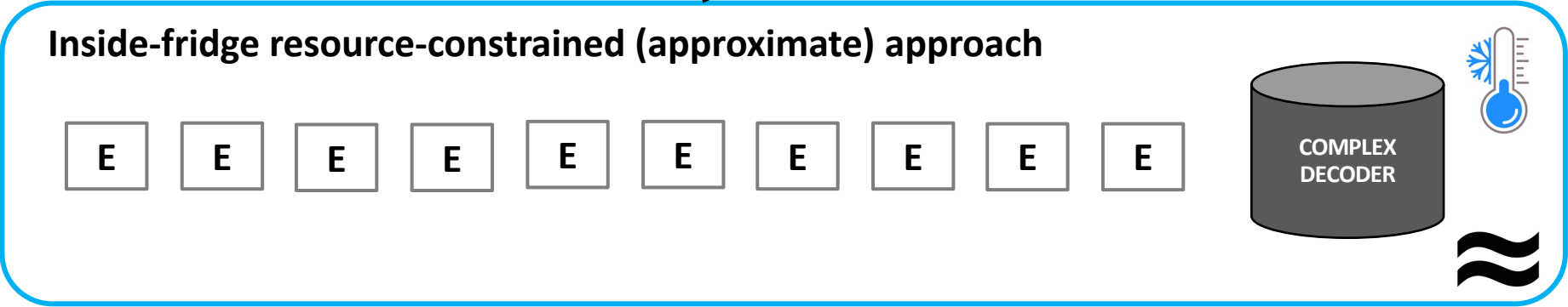


Better than worst case approach to decoding

Key Insight: Not all errors hard to decode → Separate common trivial errors from rare complex errors.

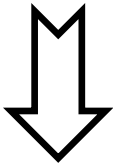


Decode ALL errors (one-size-fits-all)



Better than worst case approach to decoding

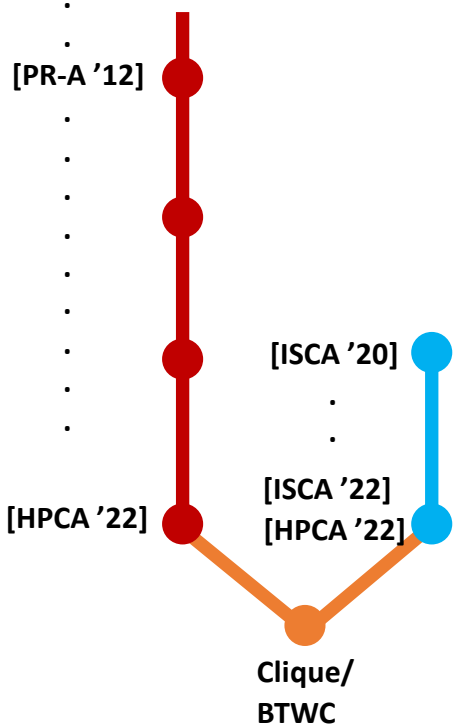
Key Insight: Not all errors hard to decode → Separate common trivial errors from rare complex errors.



Easy to separate!
Most errors are trivial!

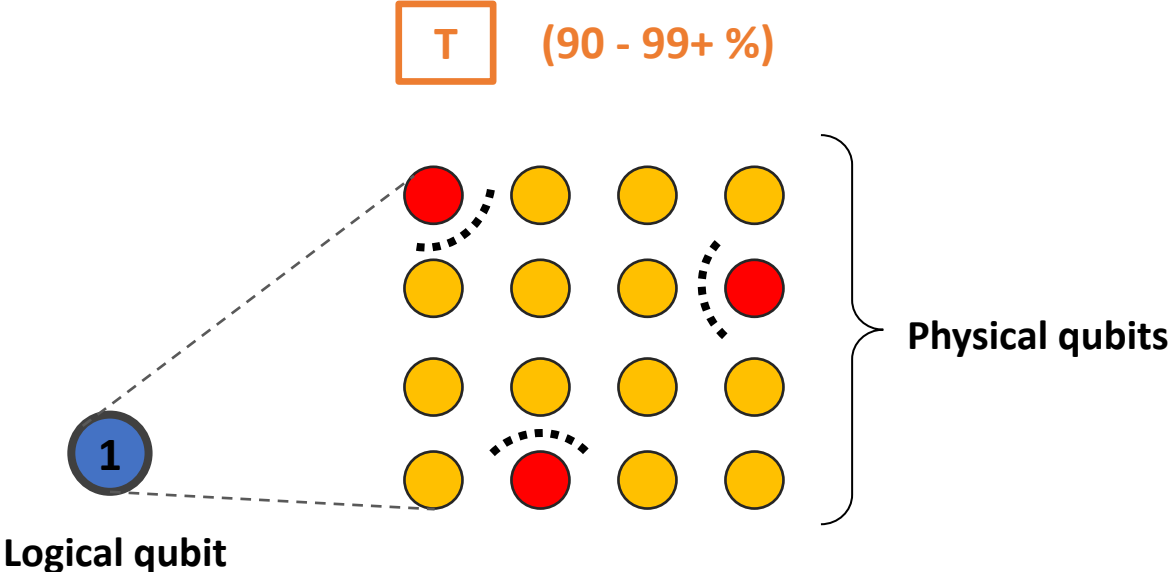


T: Trivial-to-decode
C: Complex-to-decode

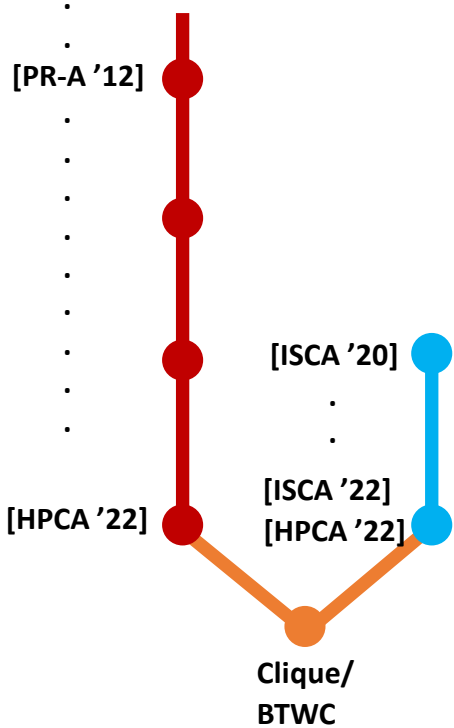


Better than worst case approach to decoding

Isolated errors: Trivial to decode and very common!!



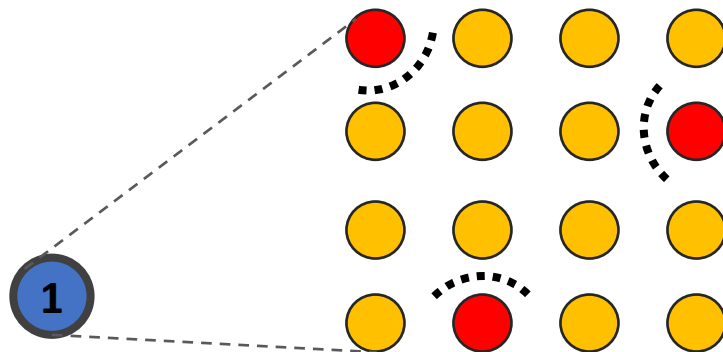
T: Trivial-to-decode
C: Complex-to-decode



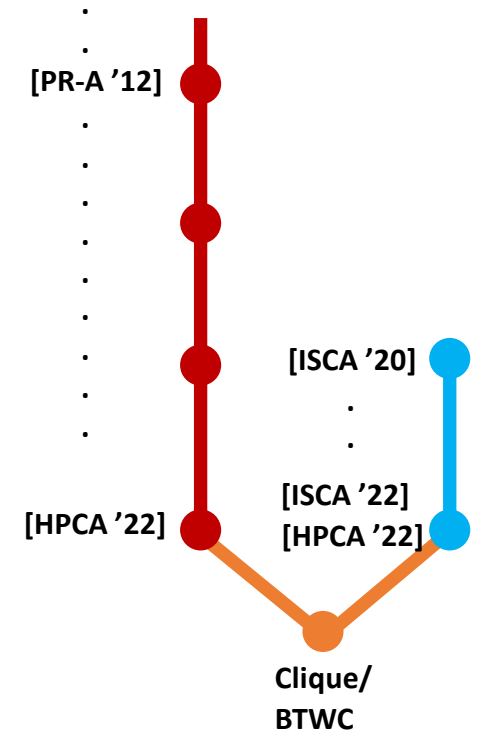
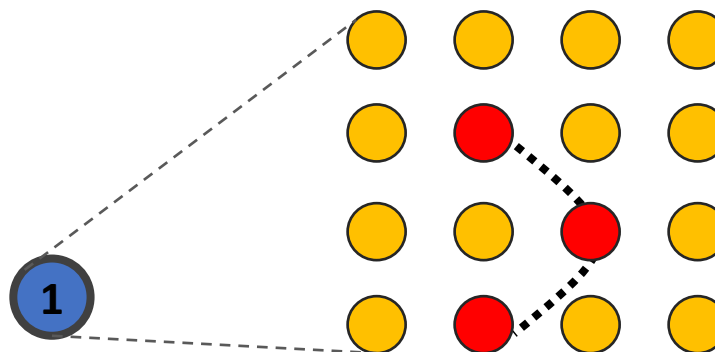
Better than worst case approach to decoding

Error chains: Hard to decode and very rare!!

T (90 - 99+ %)



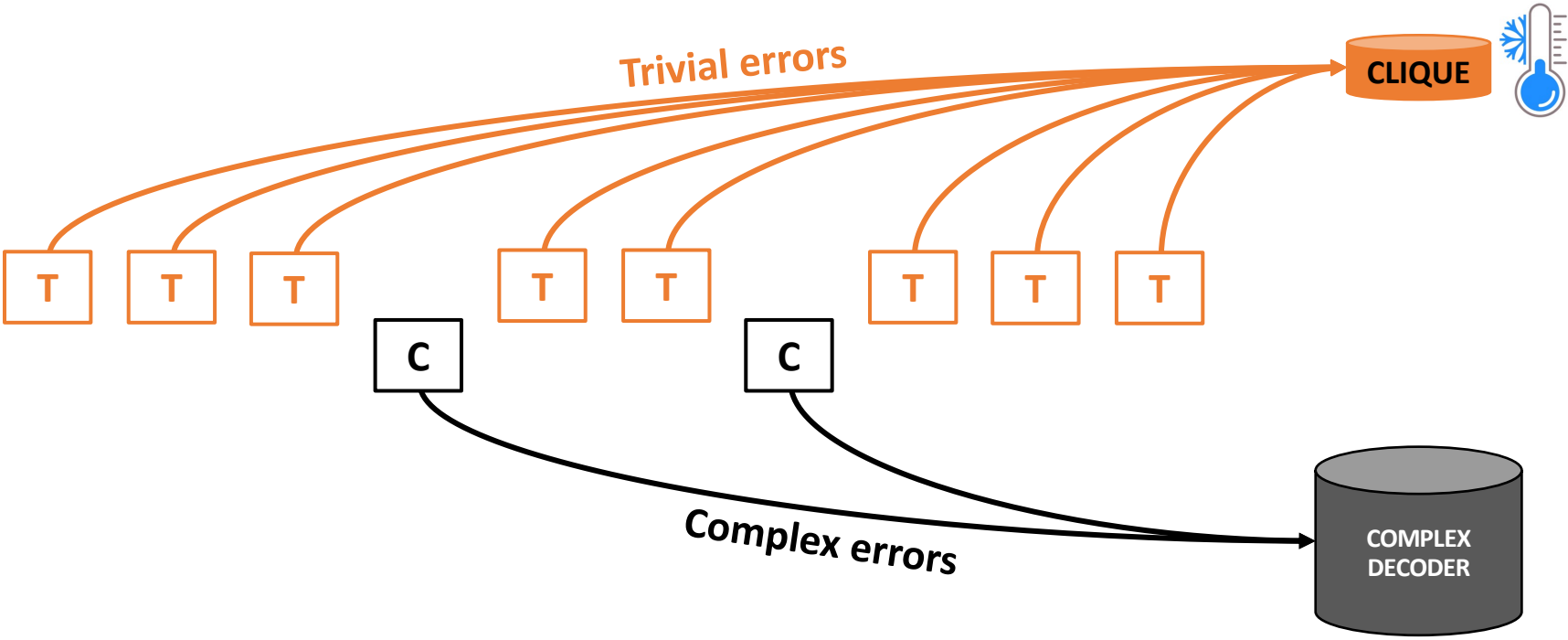
C (<1 - 10 %)



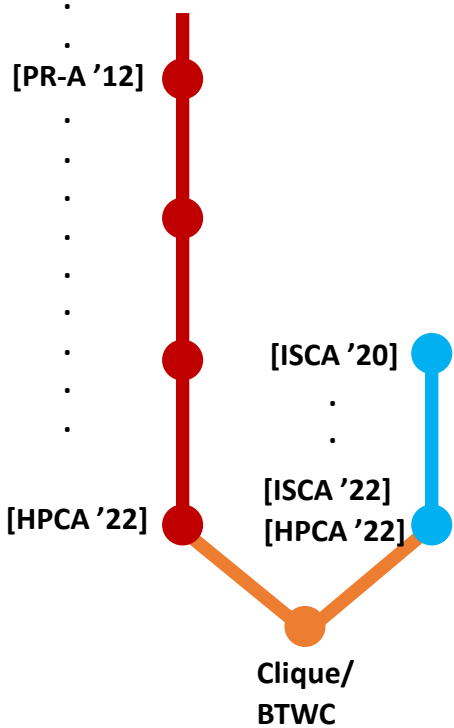
T: Trivial-to-decode
C: Complex-to-decode

Better than worst case approach to decoding

Common trivial errors → simple cryogenic 'Clique' decoder.
Rare complex errors → outside-fridge SOTA complex decoder.

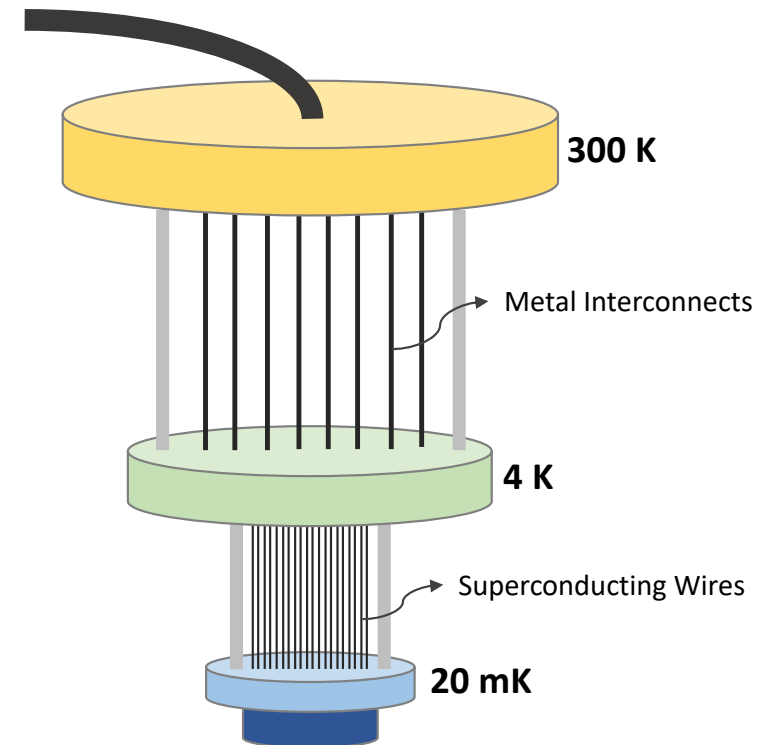
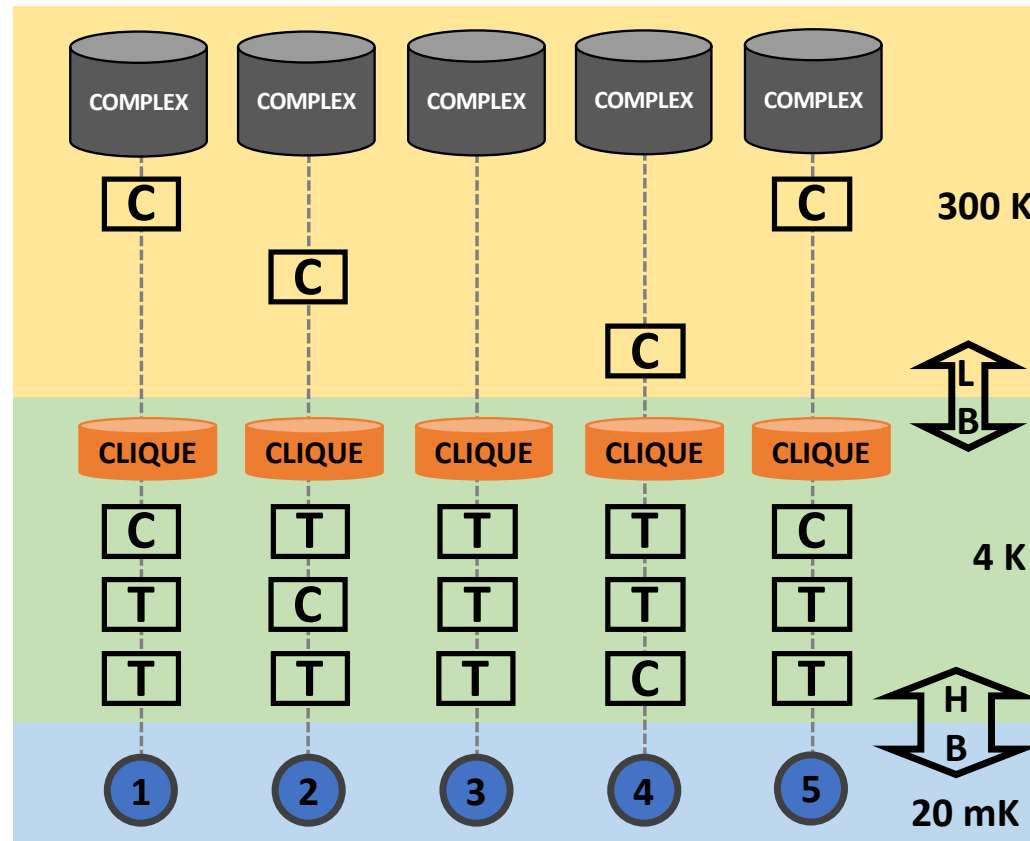


T: Trivial-to-decode
C: Complex-to-decode



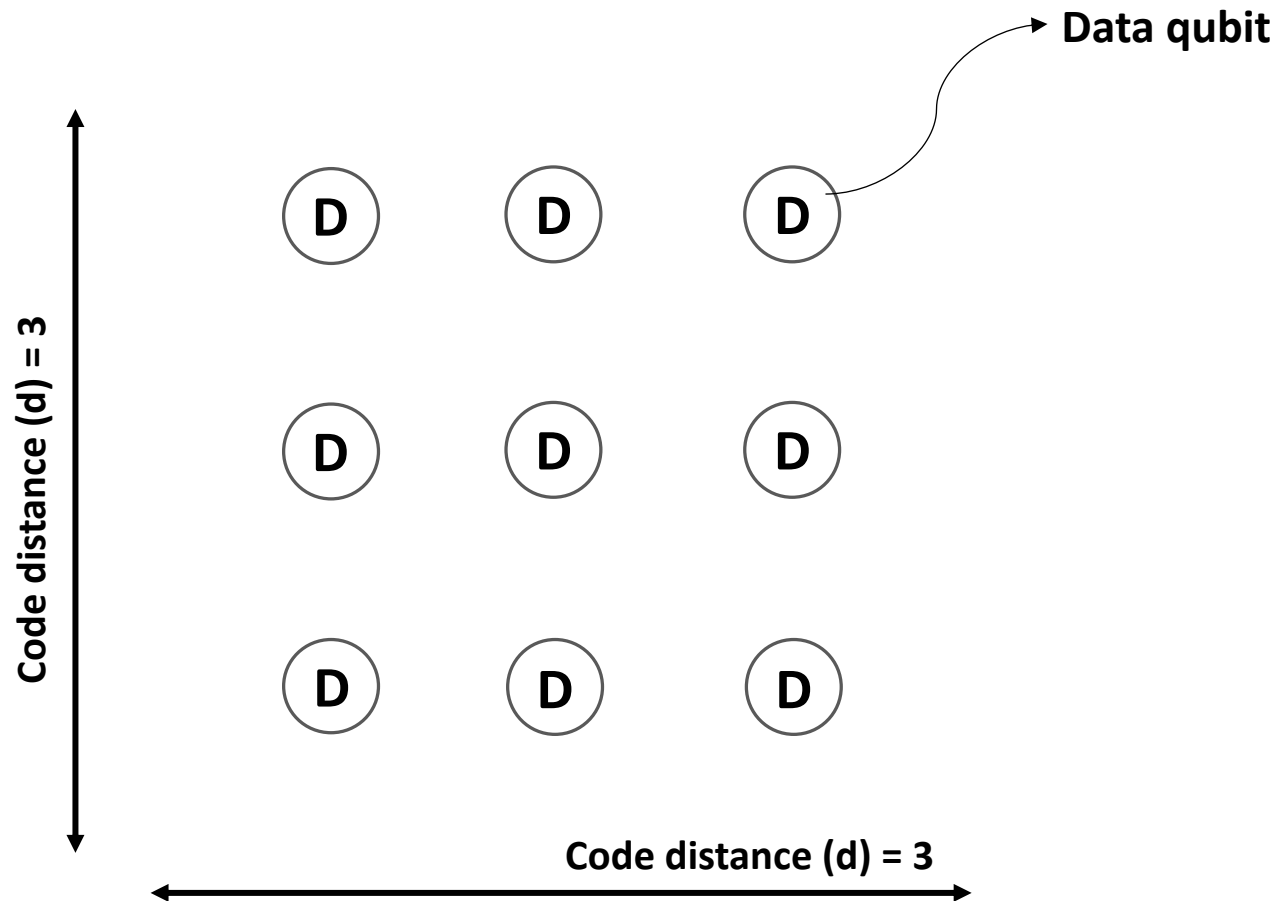
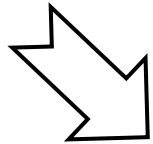
System-level view: Better than worse-case decoding

Reduced outside-fridge decoding → No bandwidth bottleneck!
Reduced inside-fridge decoding HW → No cryo-resource bottleneck!



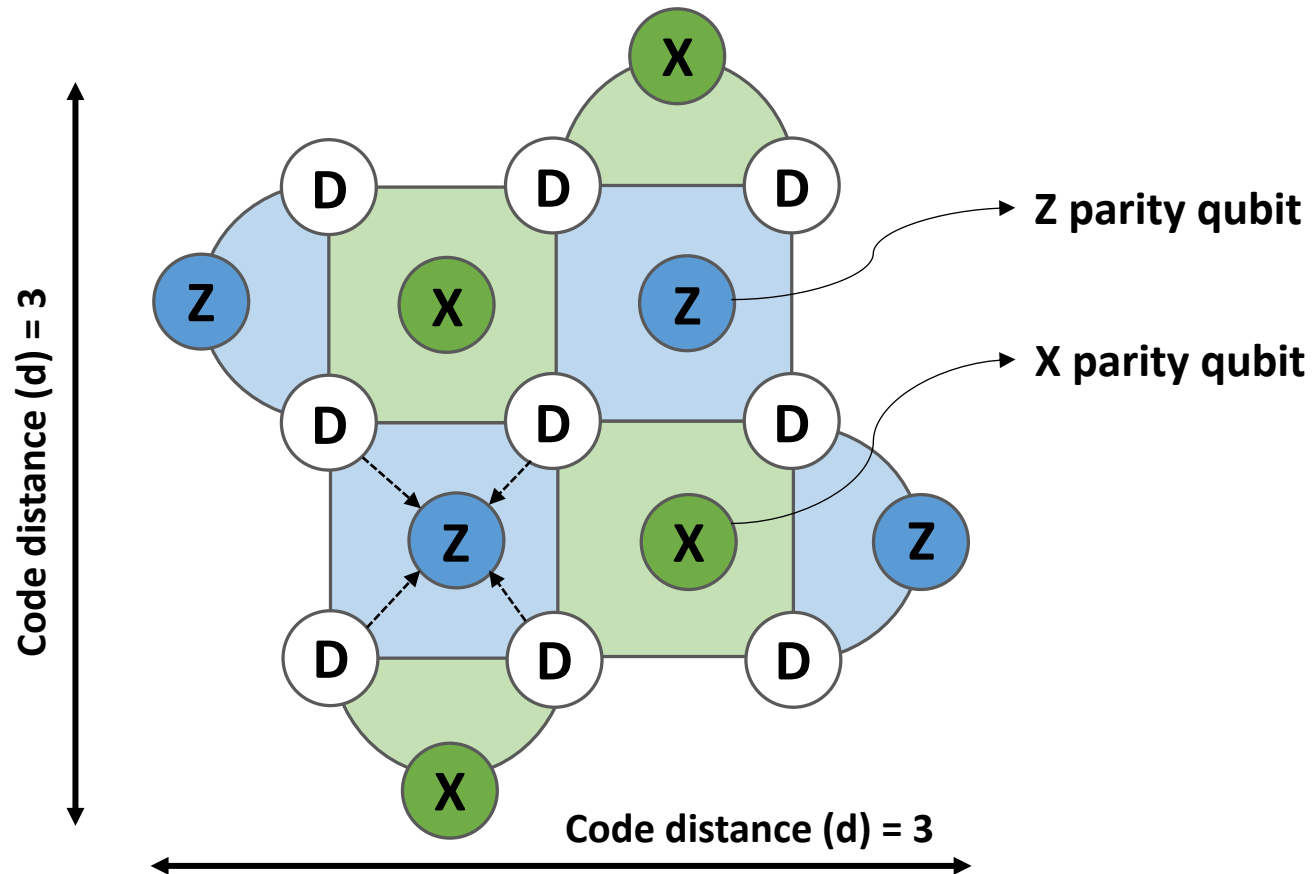
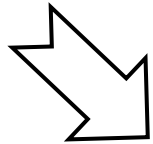
How does QEC work? (Surface code)

1 logical qubit



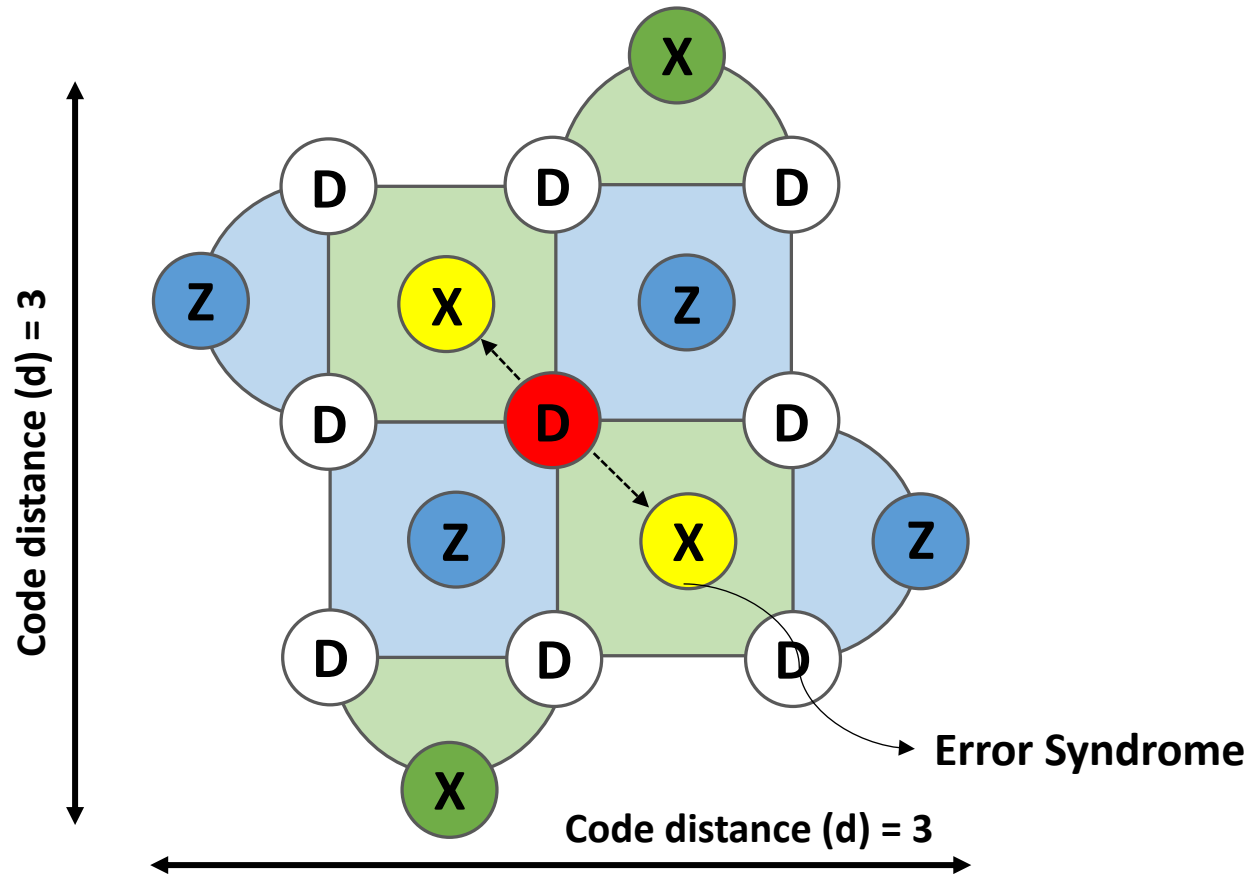
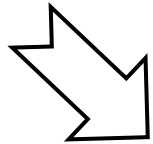
How does QEC work? (Surface code)

1 logical qubit



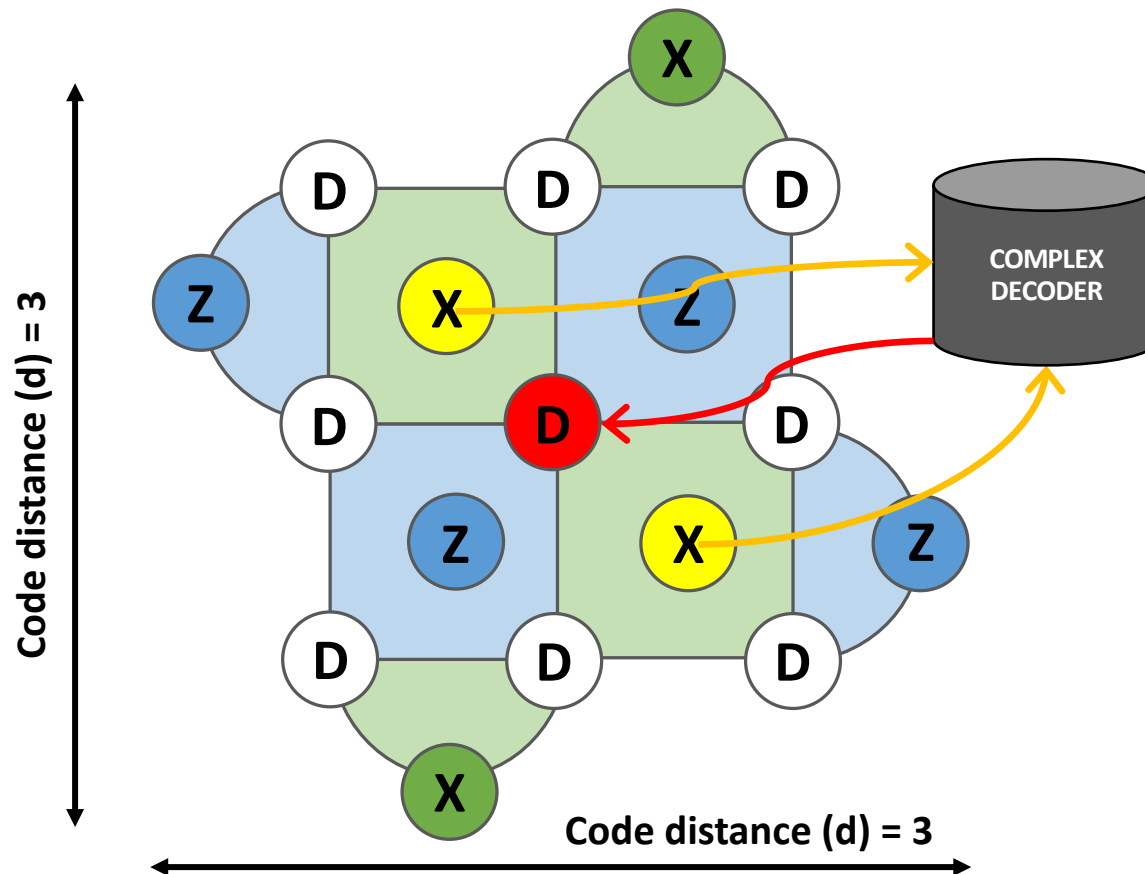
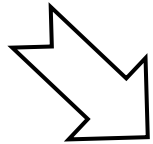
How does QEC work? (Surface code)

1 logical qubit

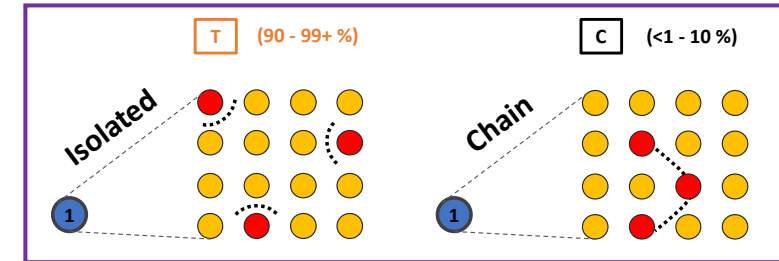
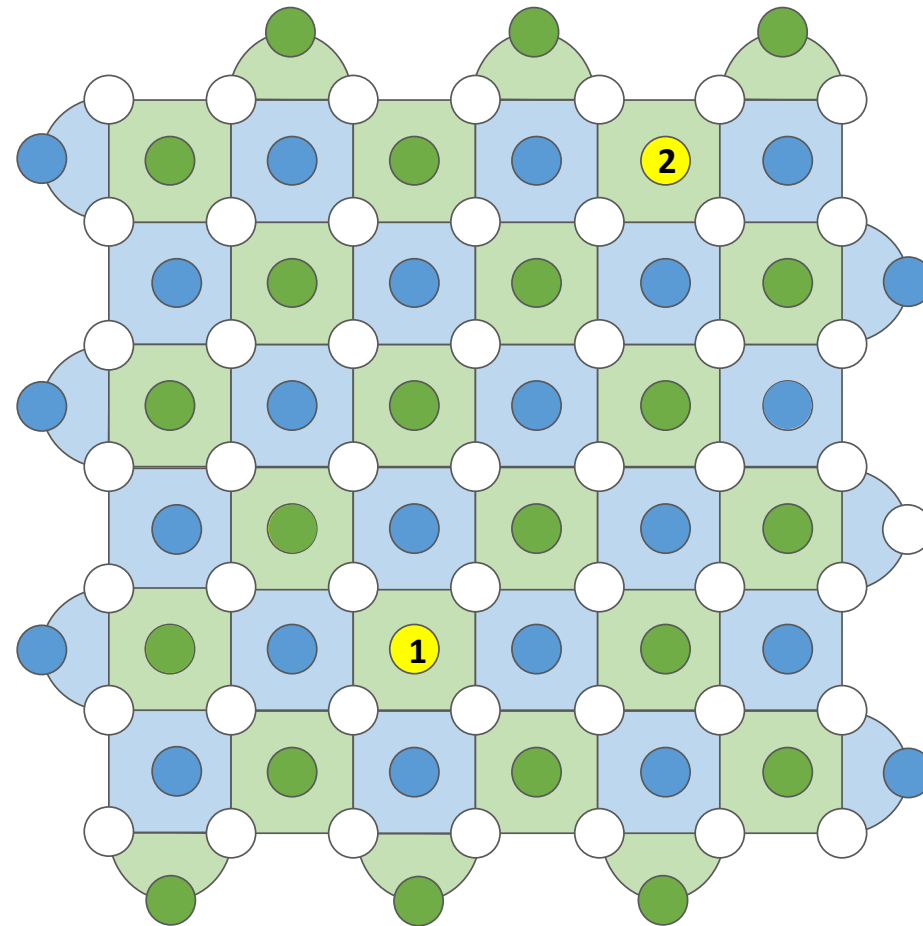


How does QEC work? (Surface code)

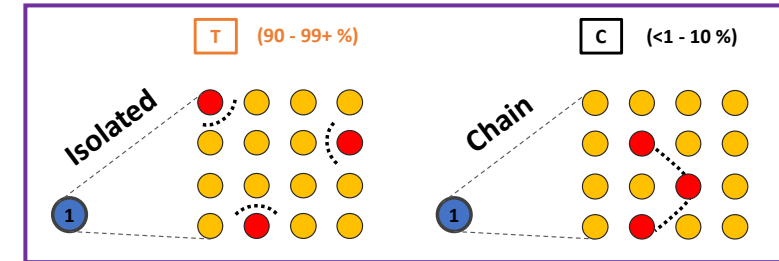
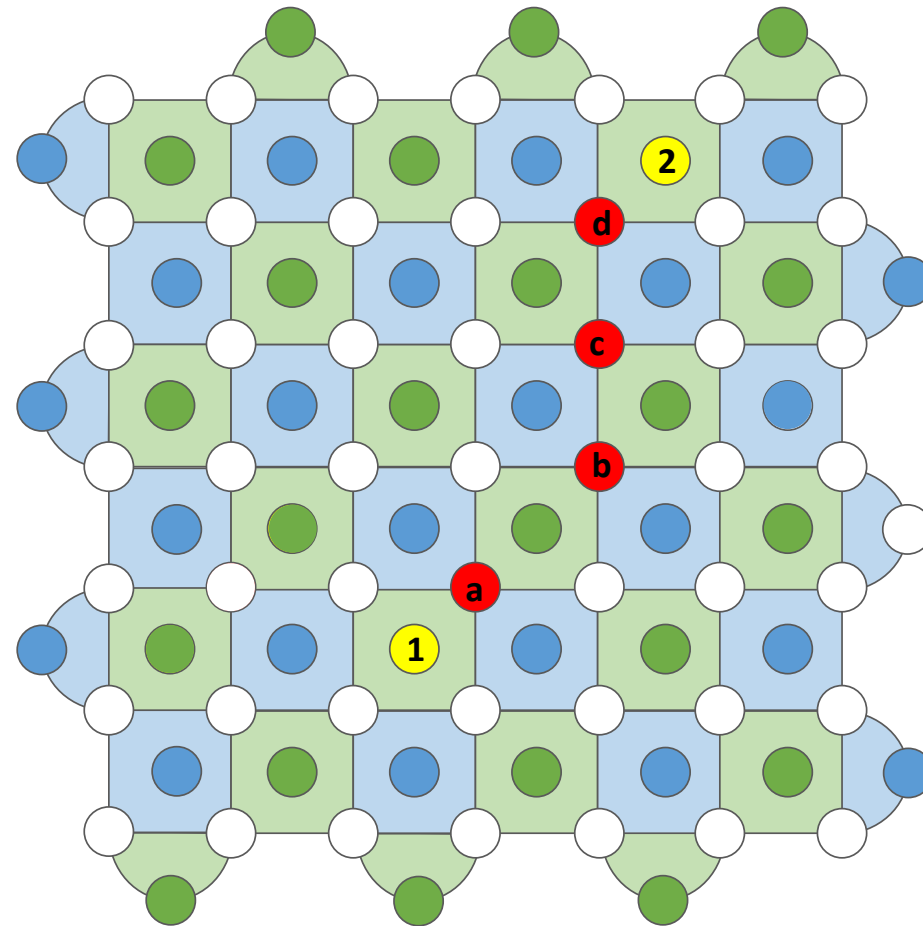
1 logical qubit



Why are error chains hard to decode?

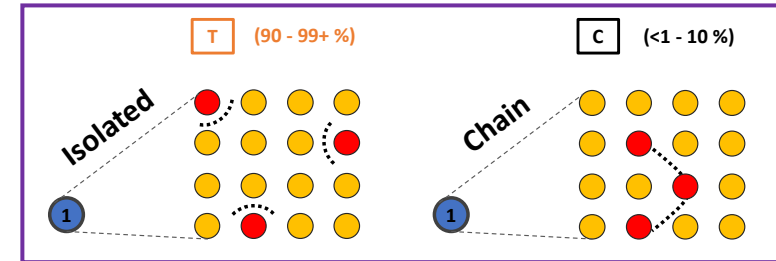
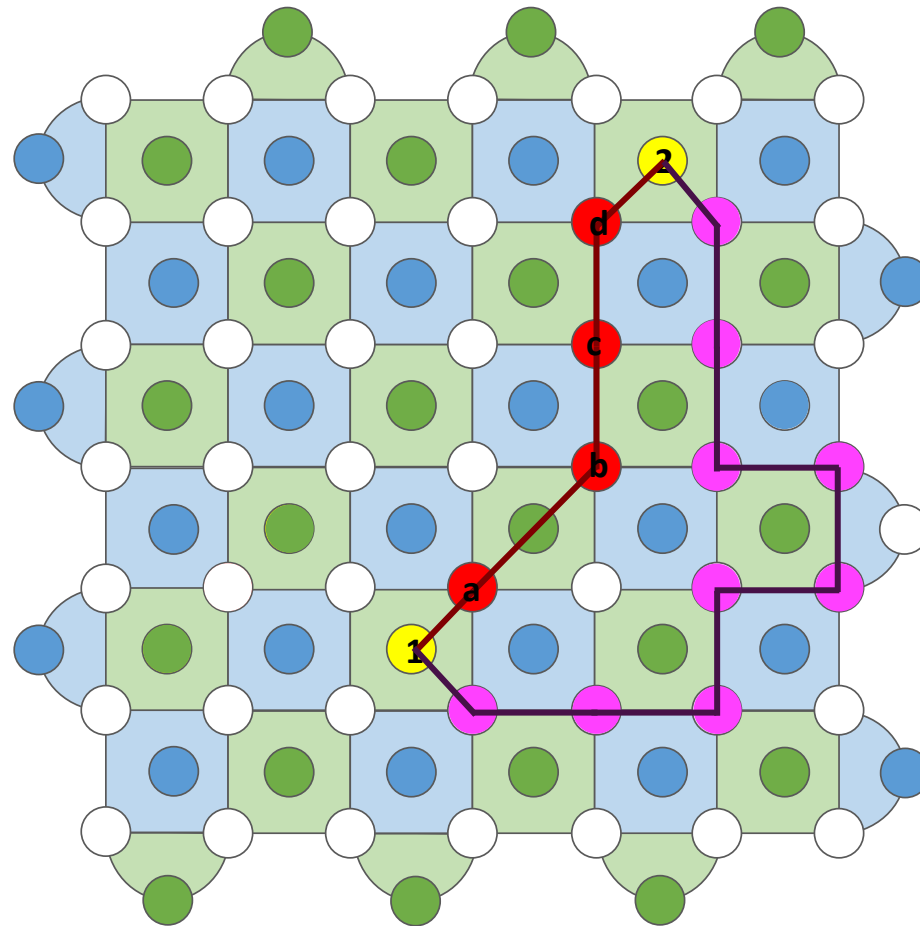


Why are error chains hard to decode?



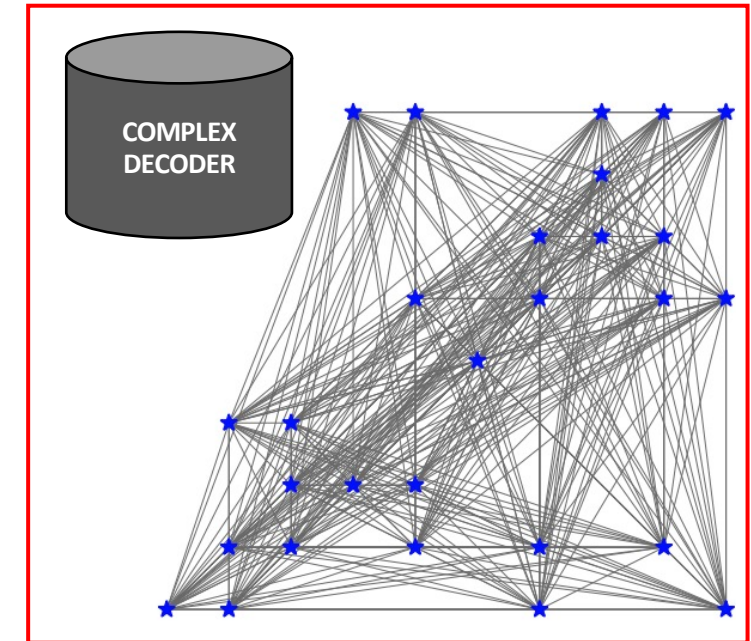
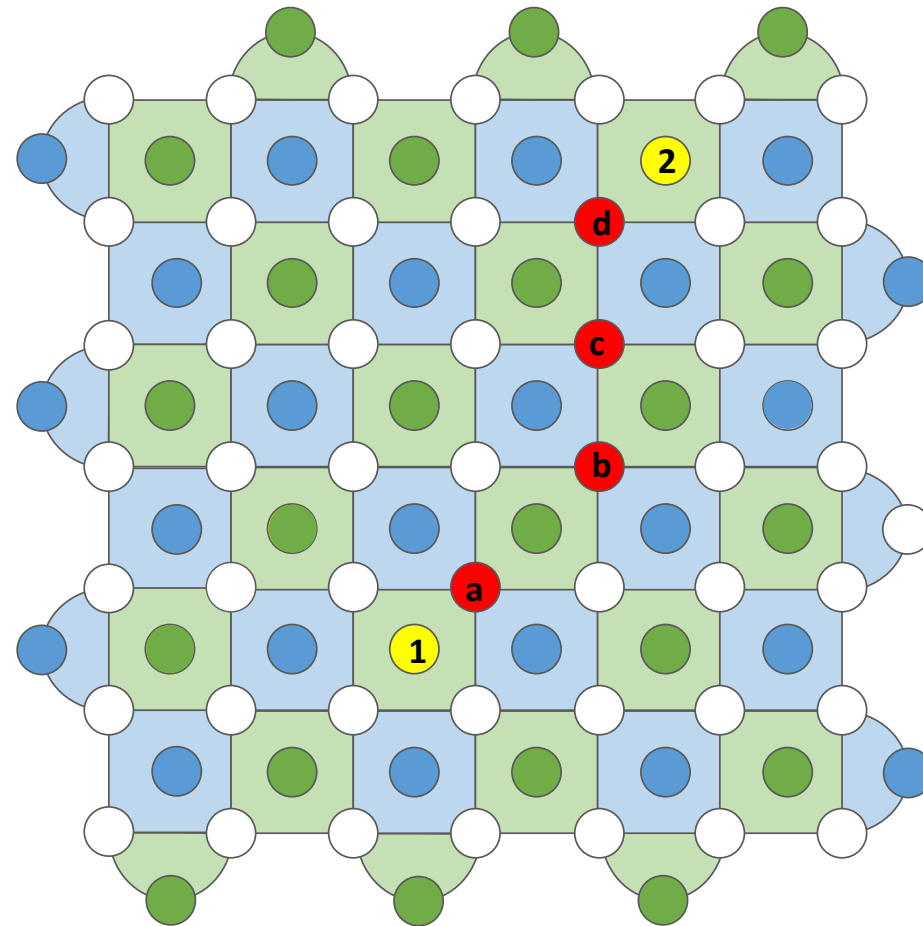
Why are error chains hard to decode?

{ 1 2 } :
● X 9
● X 4



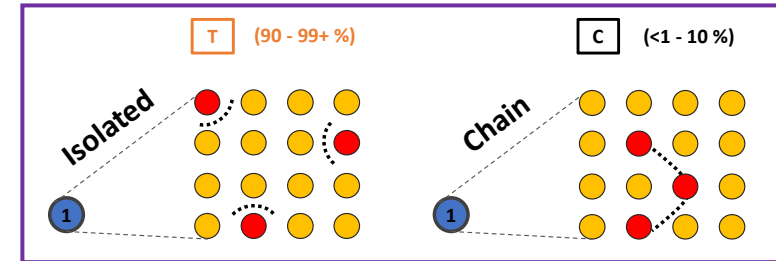
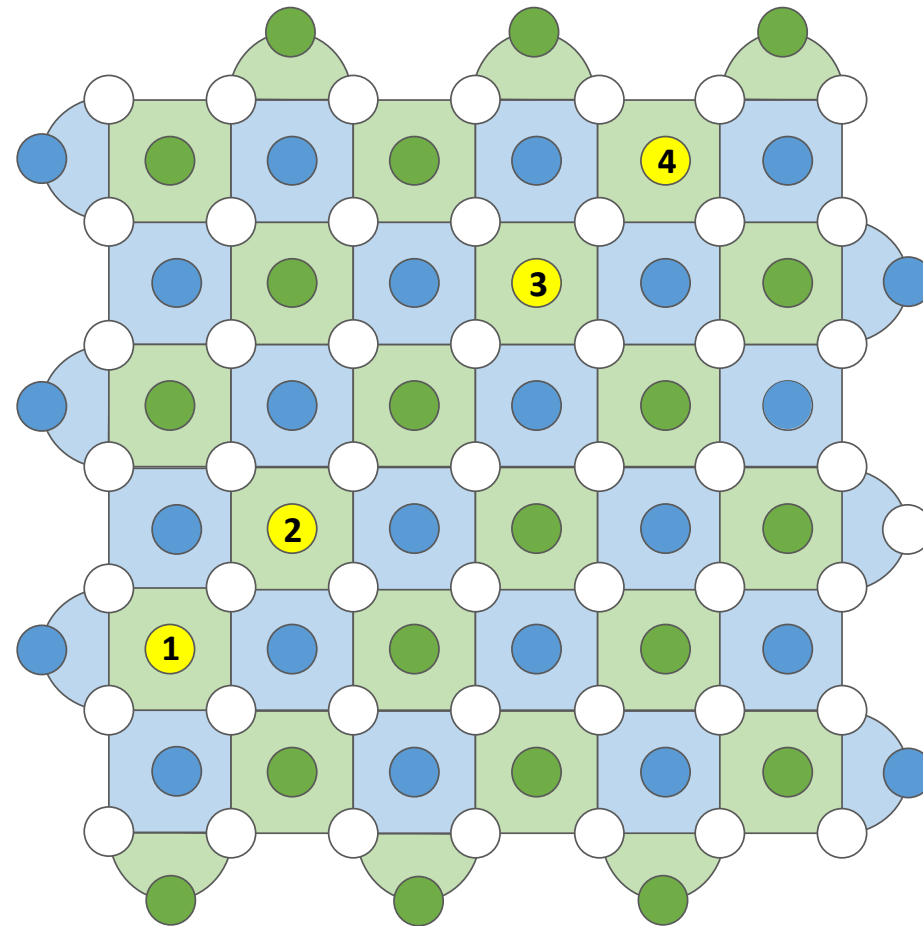
Chained data errors trigger non-local syndromes which are challenging to pair and decode.

Why are error chains hard to decode?

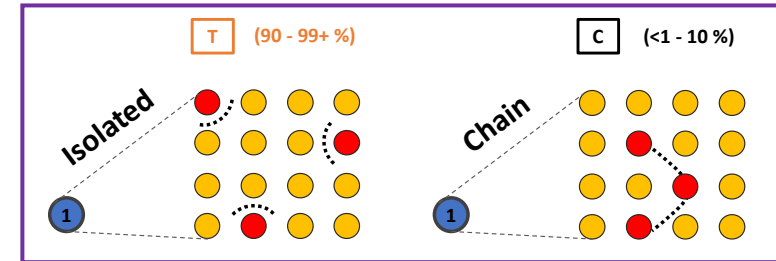
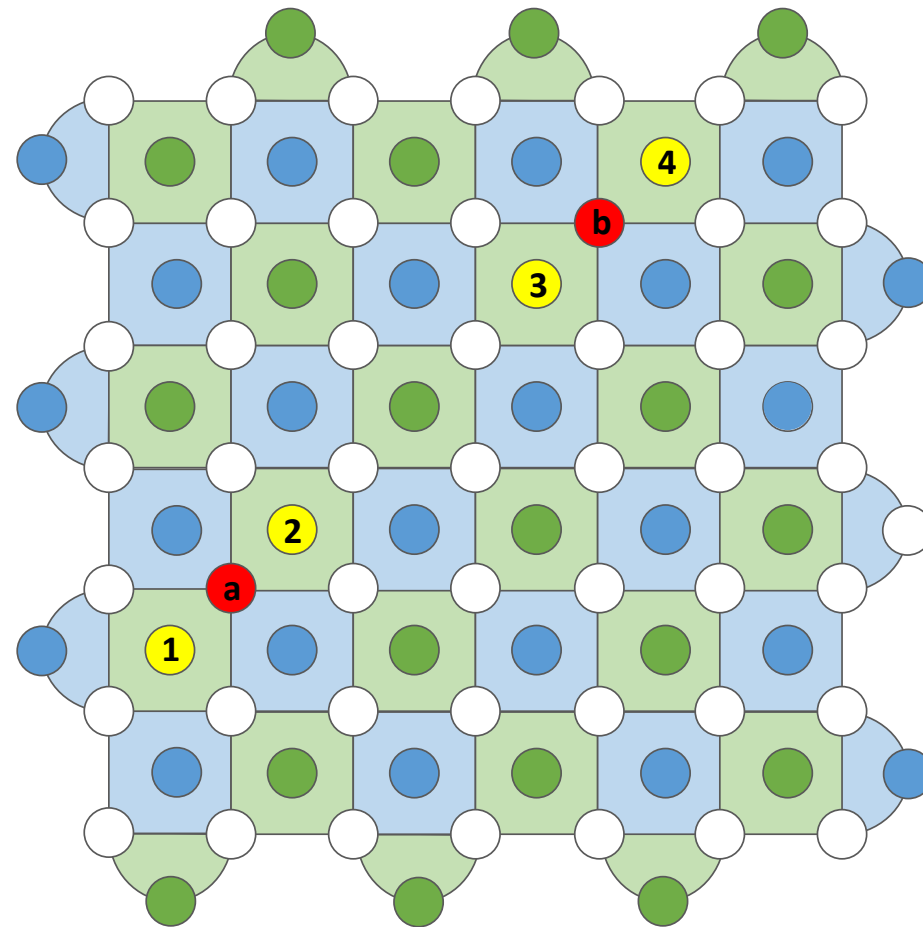


Minimum Weight Perfect Matching

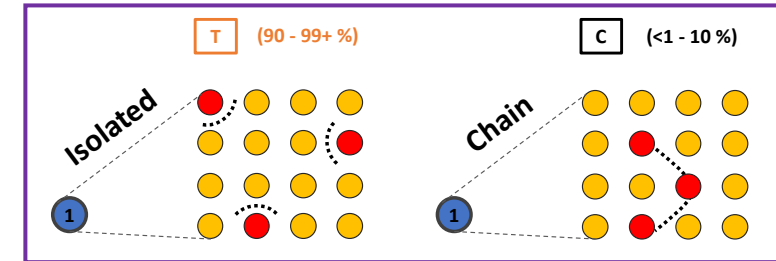
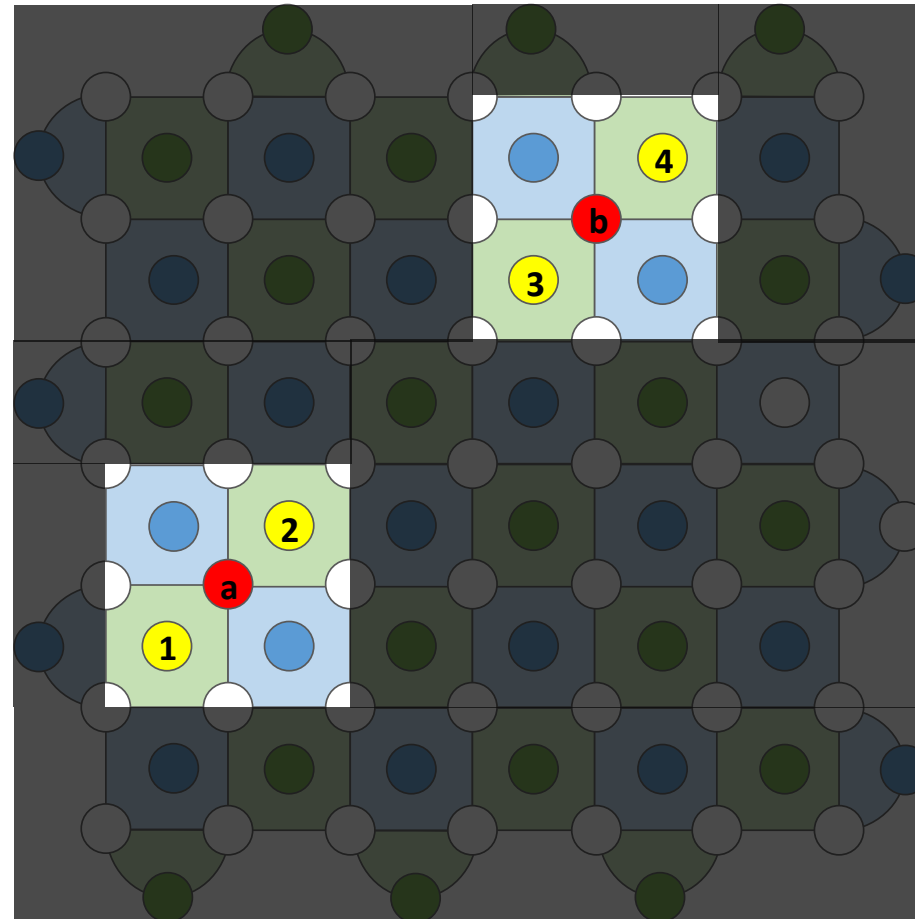
Why are isolated errors trivial to decode?



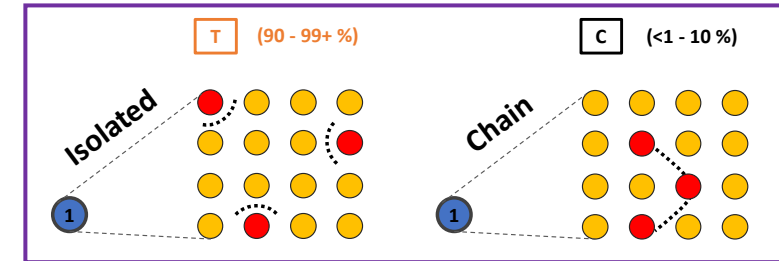
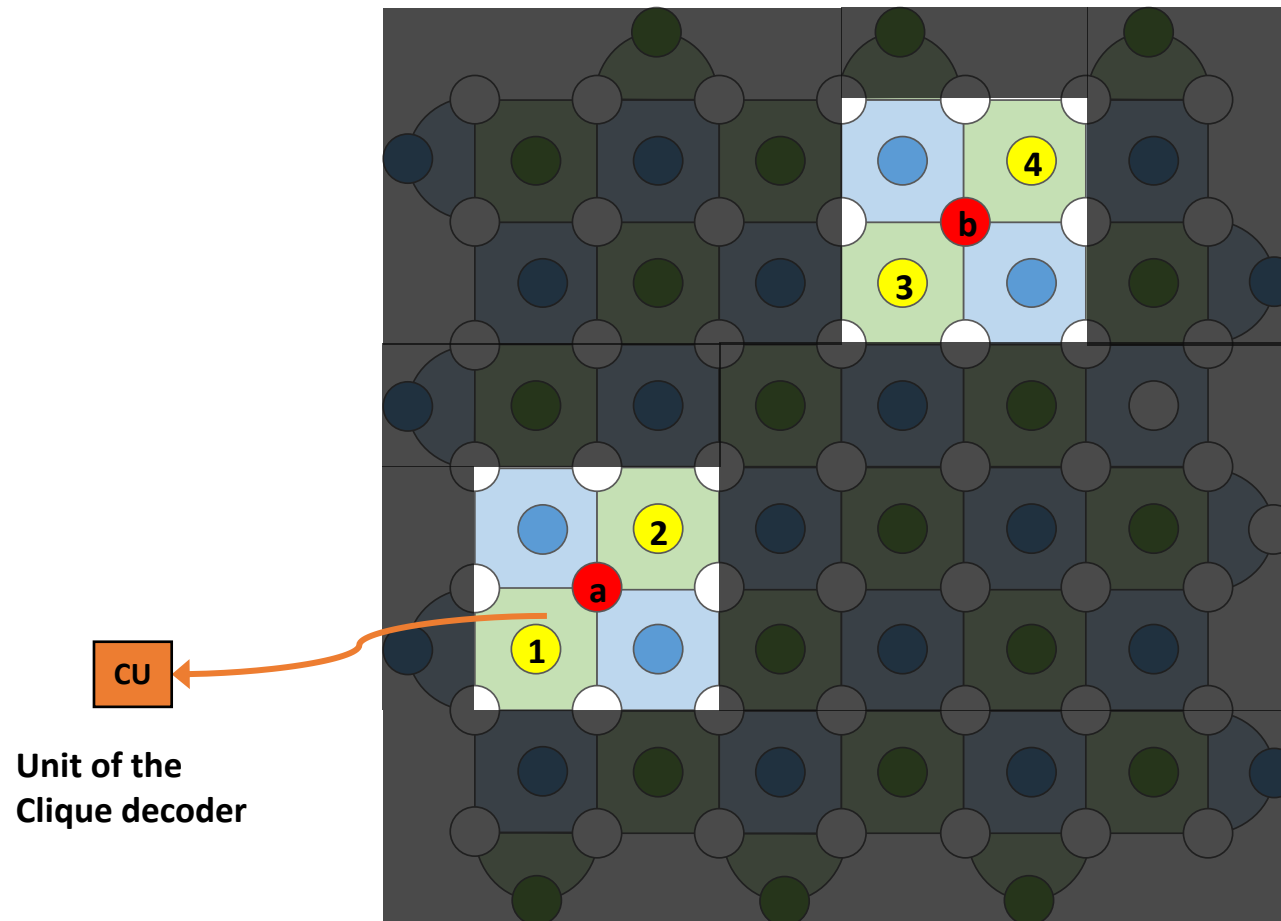
Why are isolated errors trivial to decode?



Why are isolated errors trivial to decode?

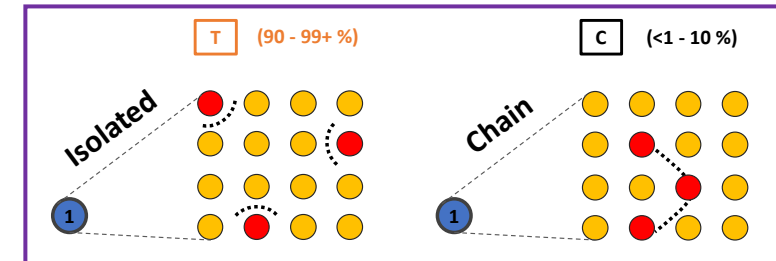
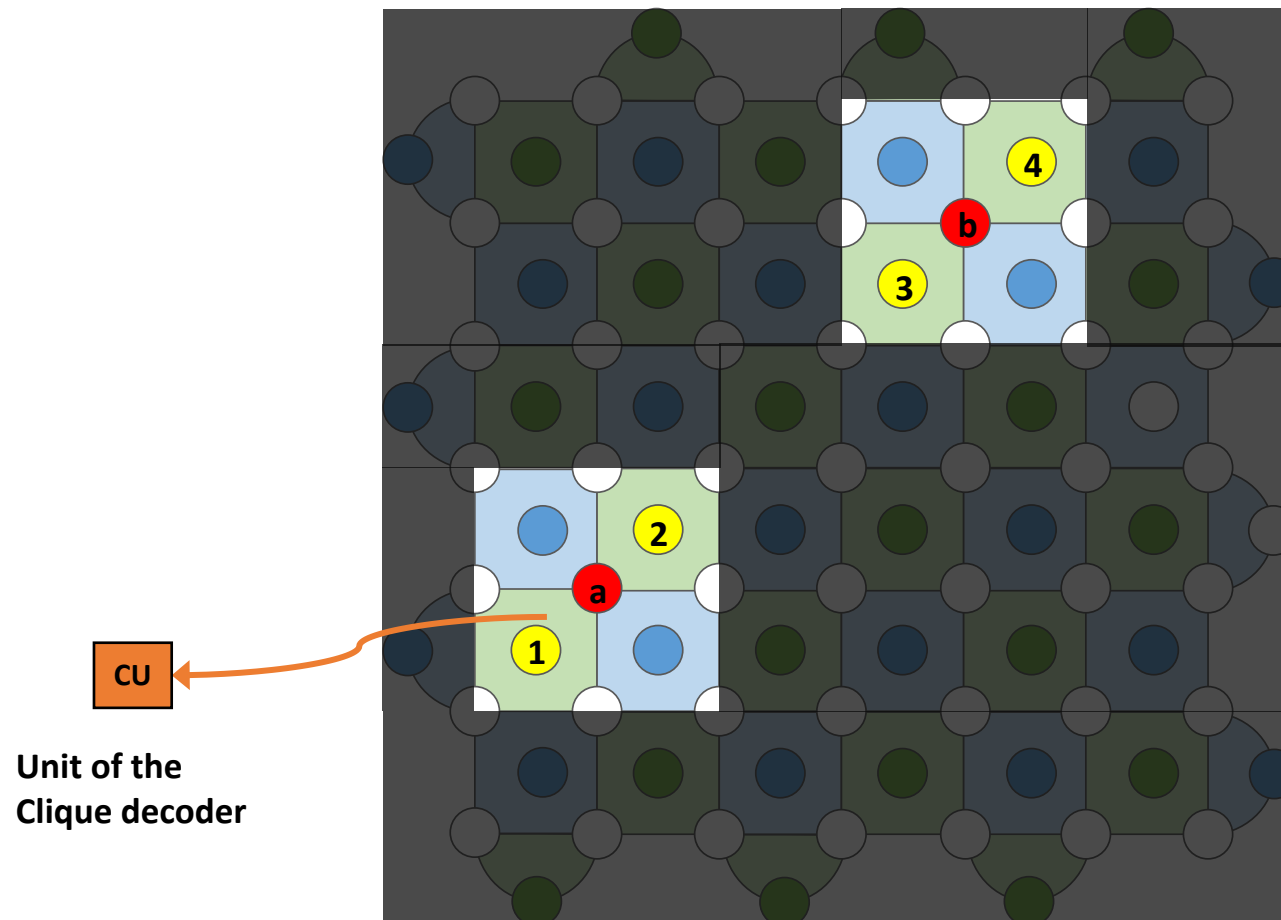


How clique trivially decodes isolated errors



Isolated data errors only trigger locally paired syndromes which are easy to decode.

How clique trivially decodes isolated errors

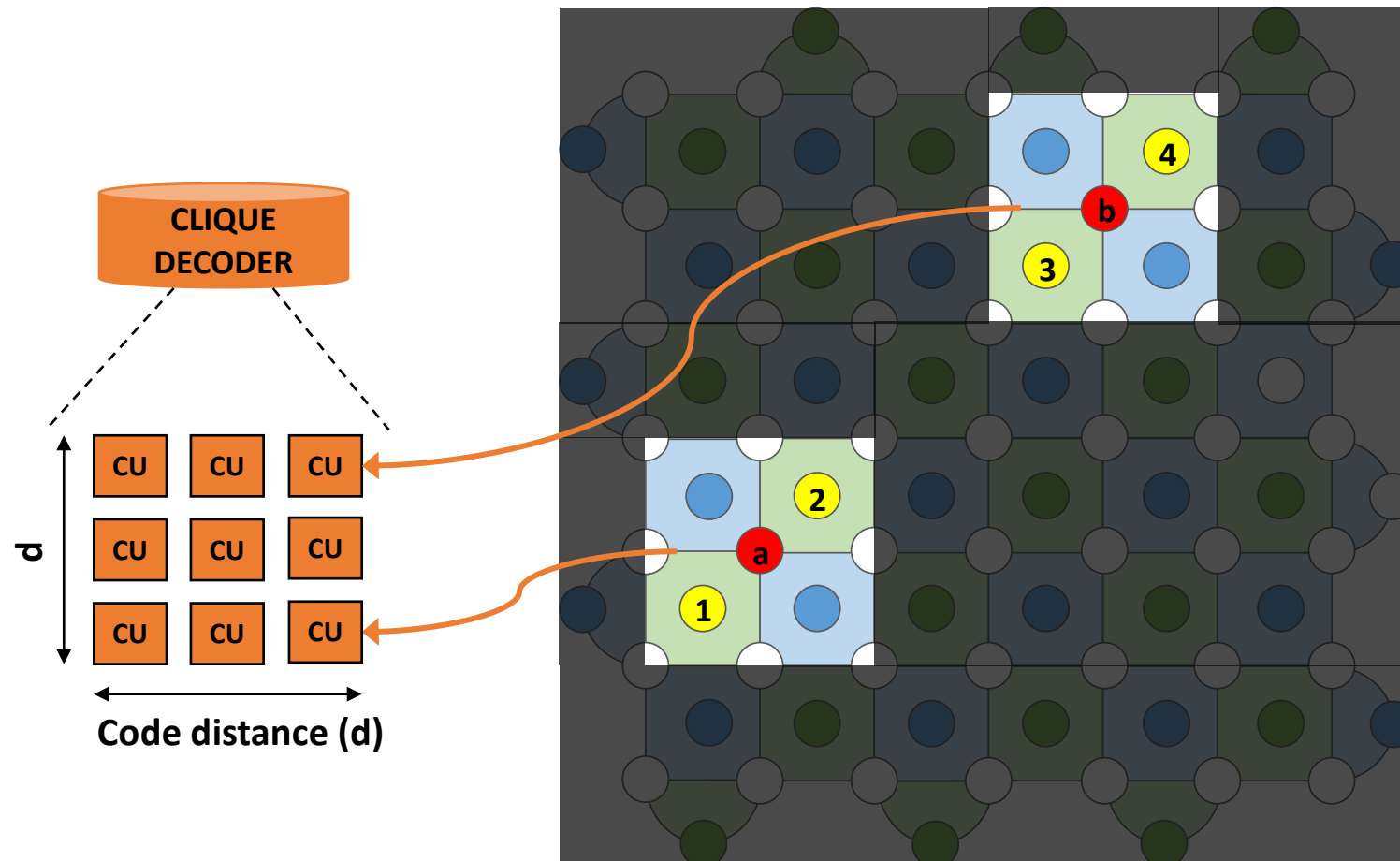


Additional subtleties to detecting if all syndromes only correspond to isolated data errors!

Additional logic required to handle syndrome measurement errors!

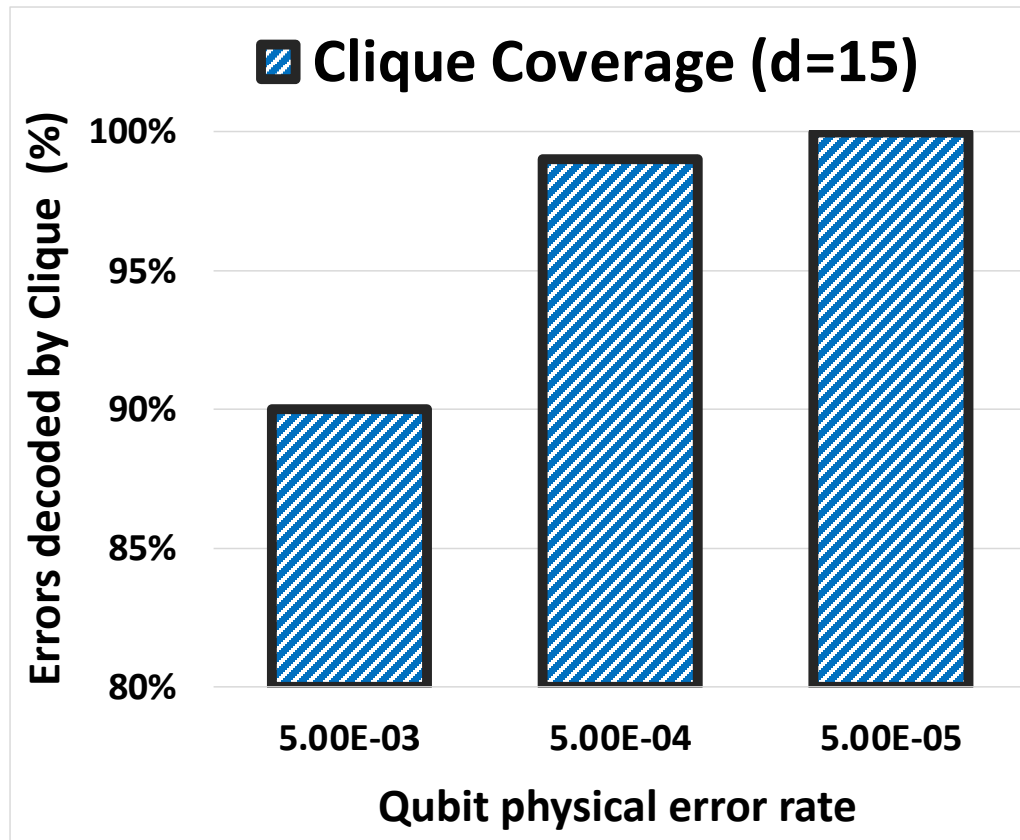
Clique decoder architecture

CU: ~10 combinational gates.
Clique decoder: d^2 CUs.
Linear Clique scaling wrt. physical qubits.



Quantitative benefits: Fridge I/O bandwidth reduction

90 - 100% of decodes handled trivially by Clique, largely eliminating outside-fridge decoding.

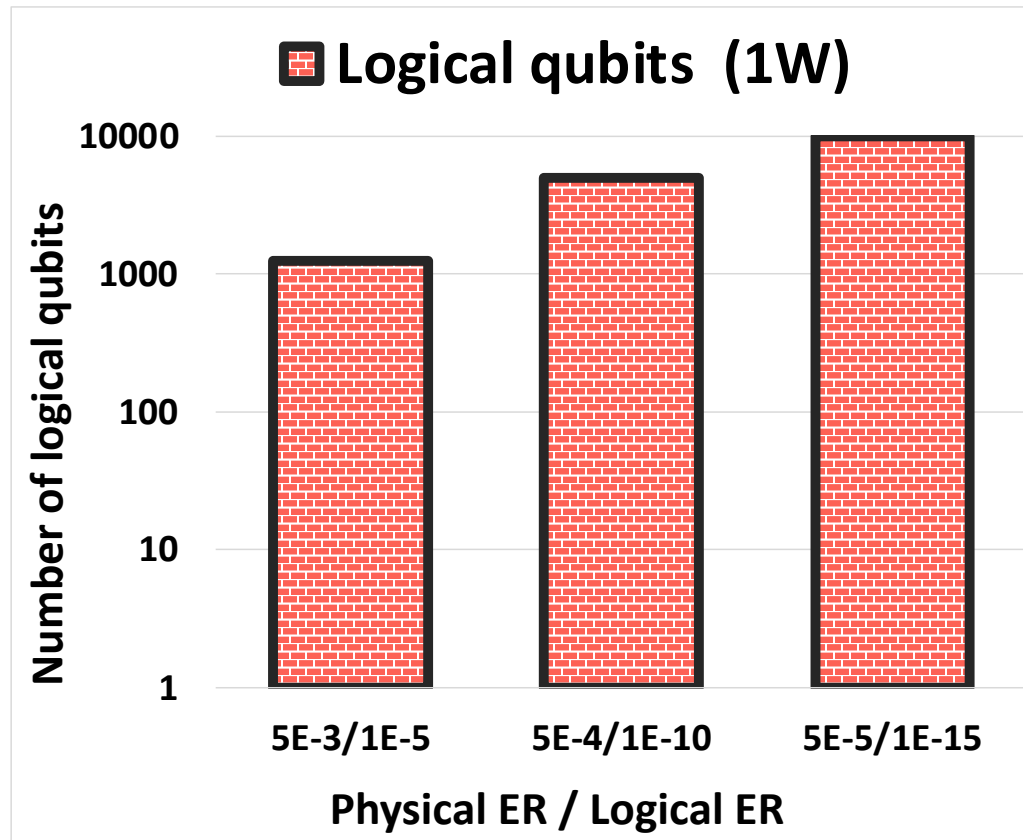


Comparison to AFS compression [Das, HPCA '22]:
Clique BW reduction is 10-10,000x greater than AFS which is an entirely off-chip decoding scheme but employs data compression on error I/O data.

Quantitative benefits: Cryo-resource requirement

Clique supports 2.5M physical qubits at 1W power → 1000s of logical qubits.

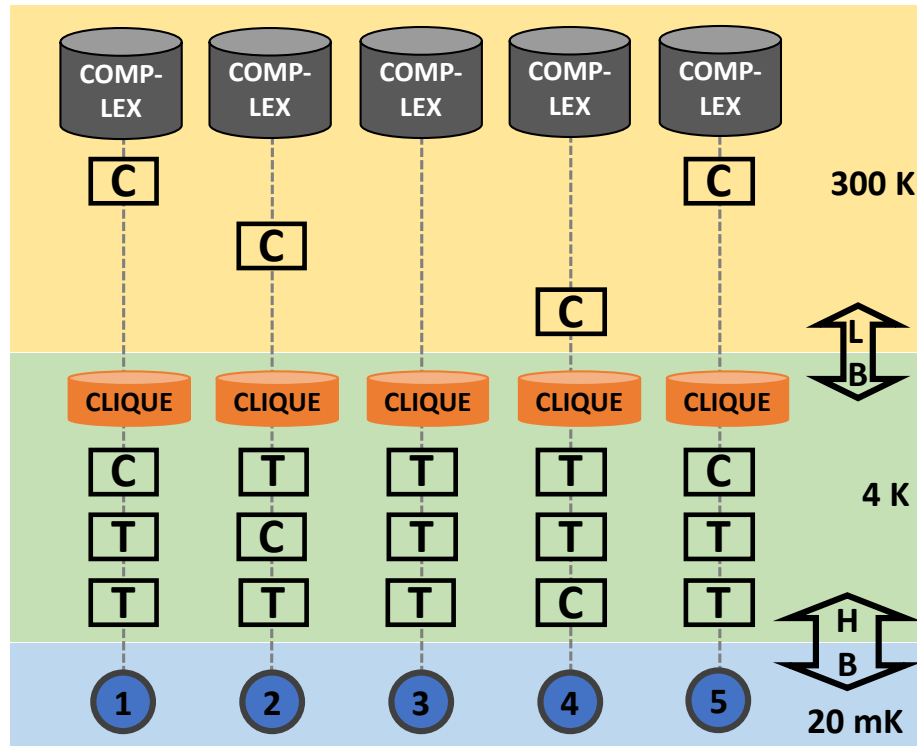
2.5M physical qubits!



Comparison to NISQ+ [Holmes, ISCA '20]:

At $d=9$, Clique requires 25-80x lower on-chip resources compared to NISQ+, an approximate fully cryogenic decoder. Greater benefits at higher code distances.

Key Takeaways



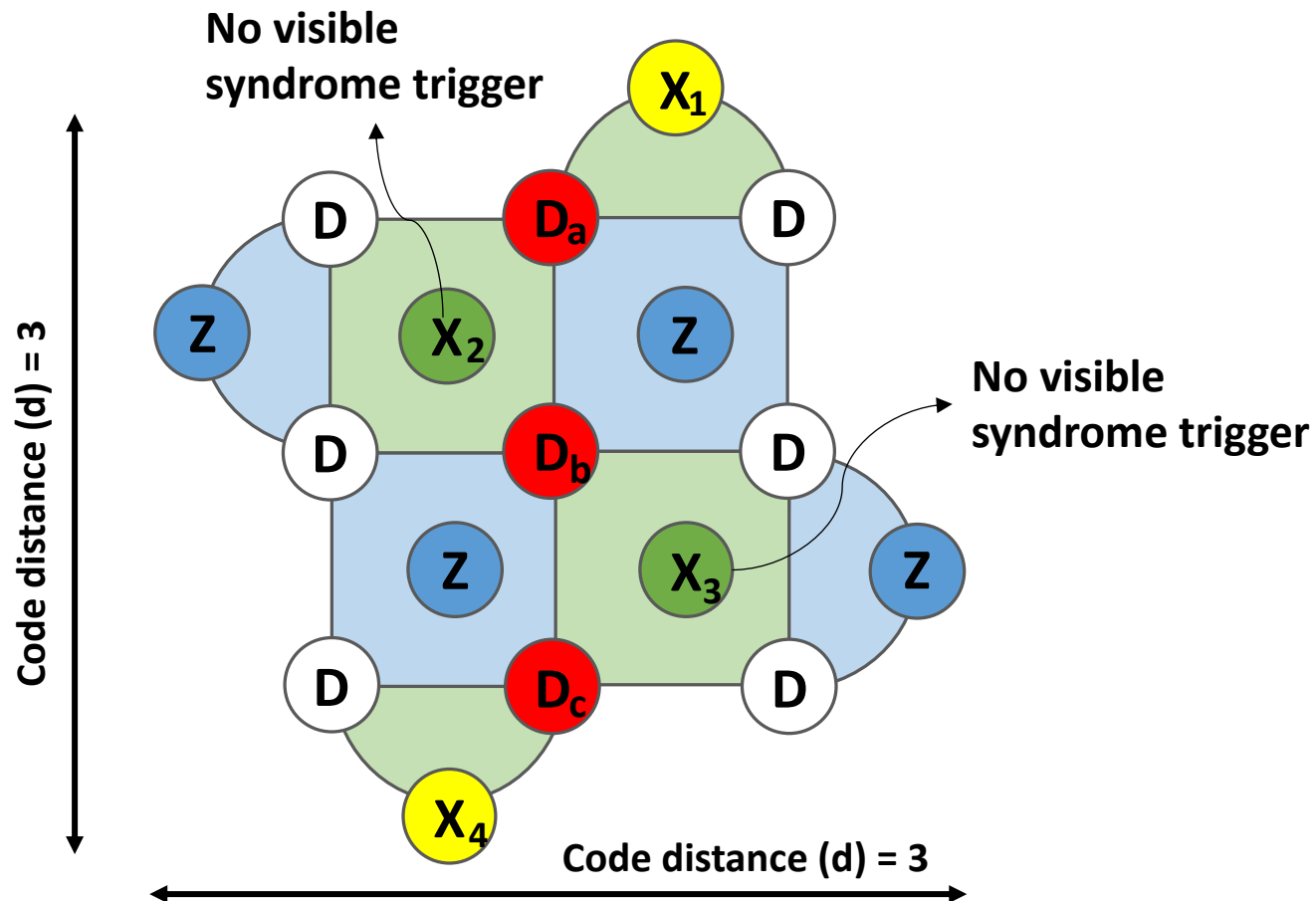
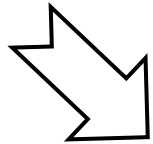
1. QEC decoding suffers severe bottlenecks: bandwidth, area, power, thermal.
2. BTWC approach: common trivial errors can be handled separately from rare complex errors
3. Clique: A lightweight cryogenic decoder for accurately decoding and correcting common-case trivial errors.
4. High fridge I/O bandwidth reduction and low cryo-resource requirement (2-4 orders of magnitude benefits over SOTA).

Thank you!

gravi@uchicago.edu

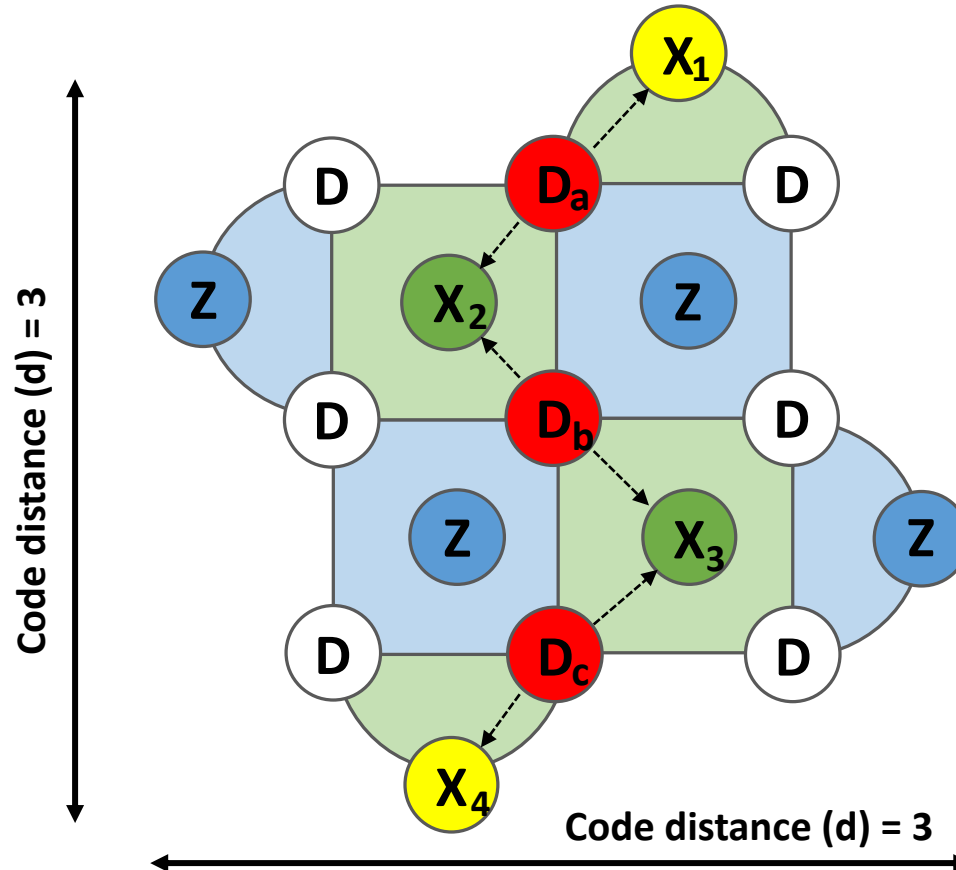
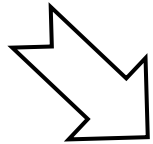
How does QEC work? (Surface code)

1 logical qubit



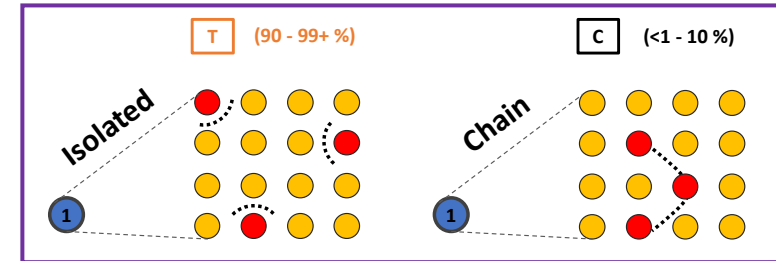
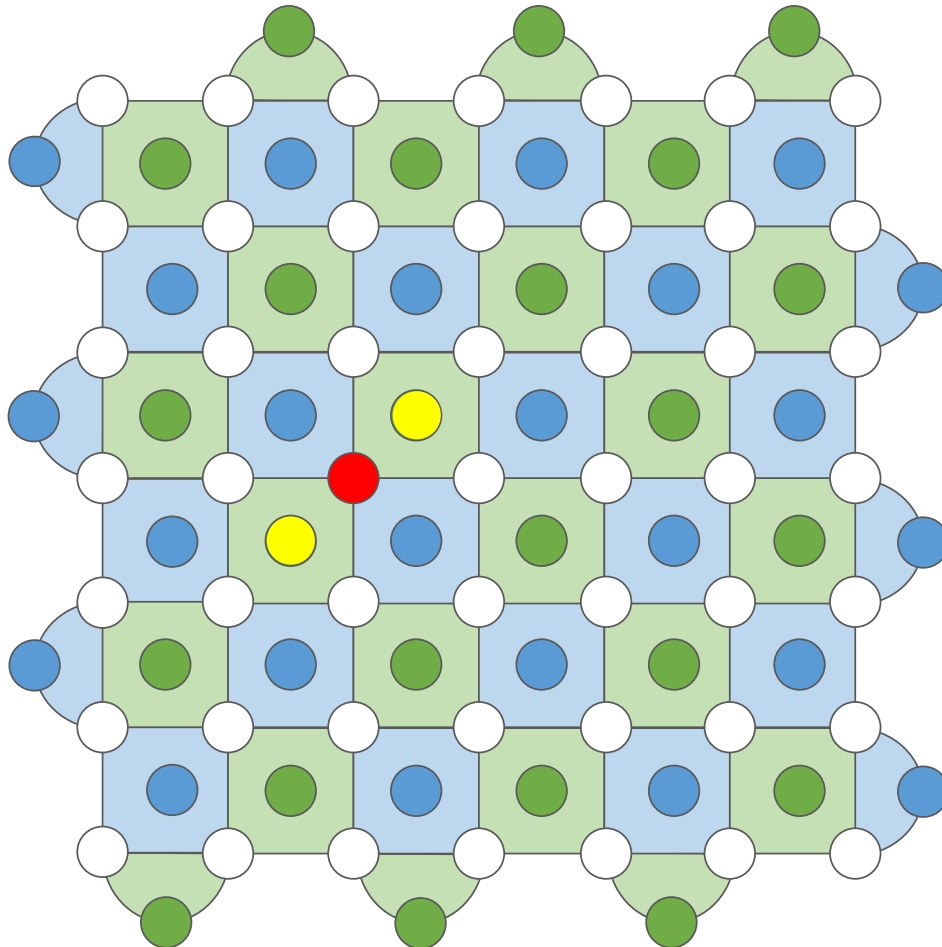
How does QEC work? (Surface code)

1 logical qubit



Error chains generate less syndrome information.

Why are isolated errors much more common than error chains?

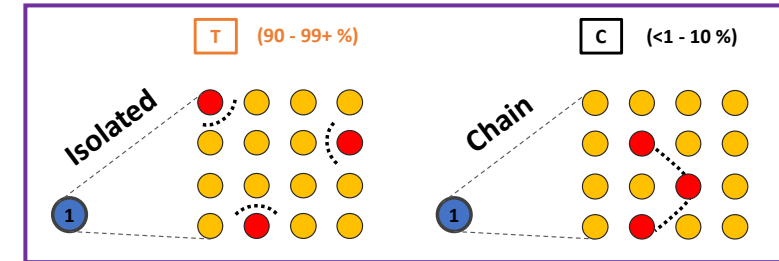
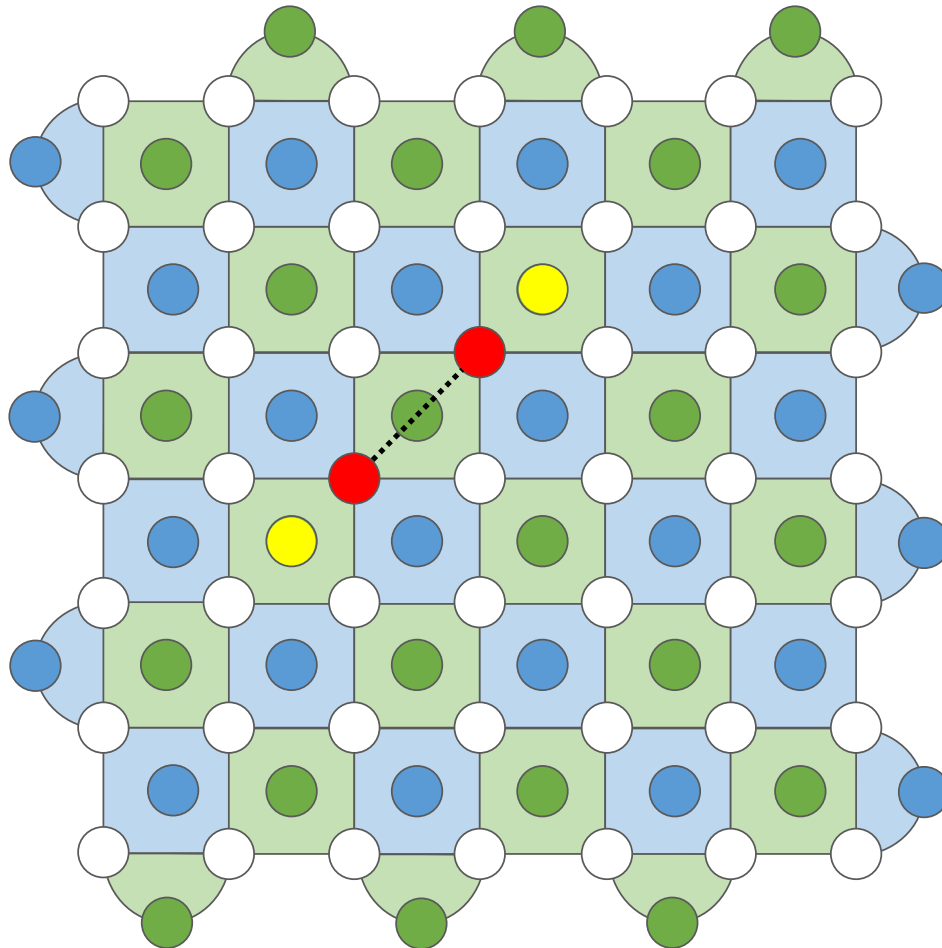


**1 logical qubit encoded in
49 physical data qubits ($d=7$)**

PER= 10^{-3} (1 in 1000), $N = 49$

P (1 error in block) = $N * \text{PER} = 4.9\%$

Why are isolated errors much more common than error chains?



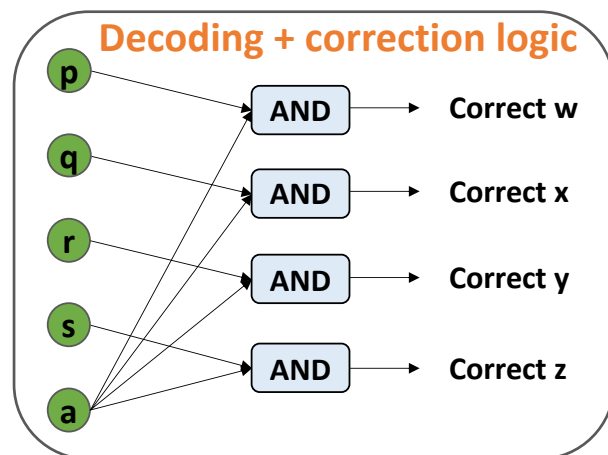
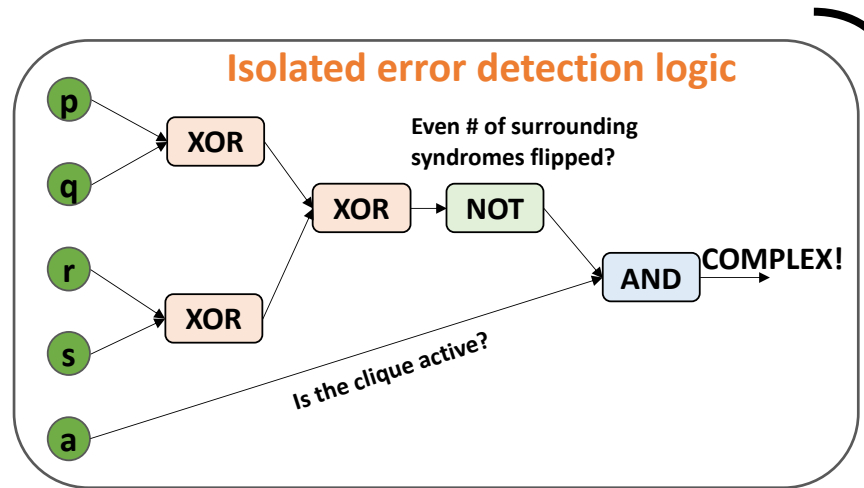
1 logical qubit encoded in
49 physical data qubits ($d=7$)

PER= 10^{-3} (1 in 1000), $N = 49$

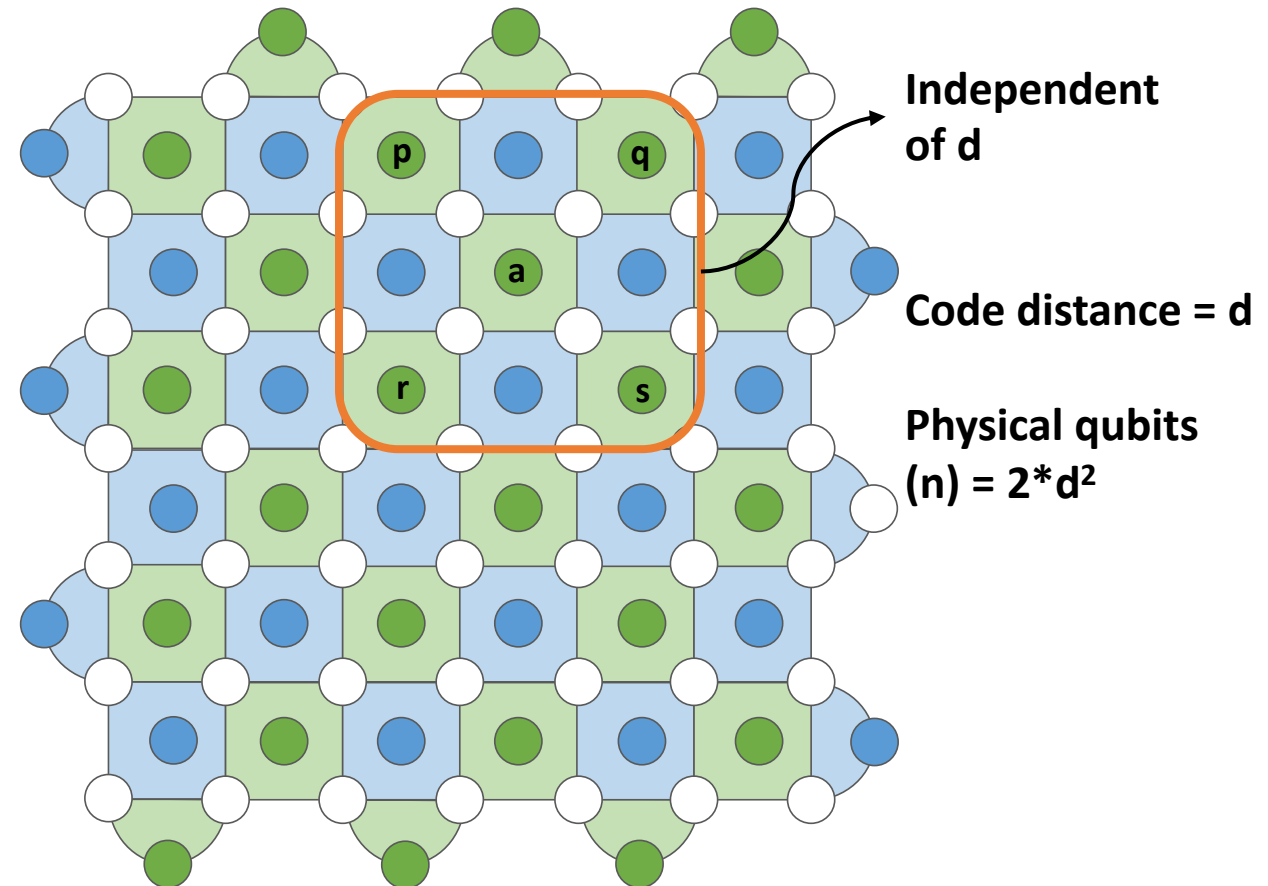
P (2 adjacent errors anywhere in block)
= $6 * N * \text{PER}^2 = 0.03\%$
(160x less likely than the isolated case)

Clique decoder hardware design

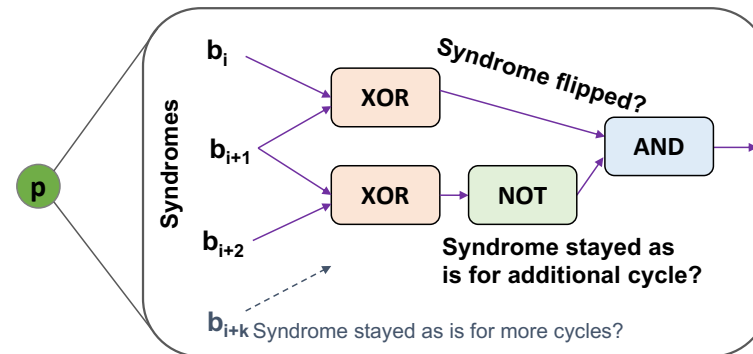
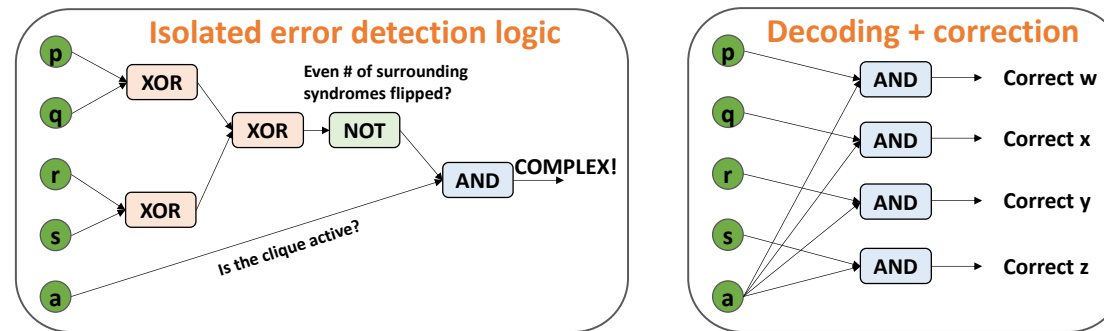
Lightweight hardware suited to cryo-domain: < 10 combinational logic gates per clique unit.
Total Clique decoder cost scales linearly in the number of qubits.



$\times O(n)$

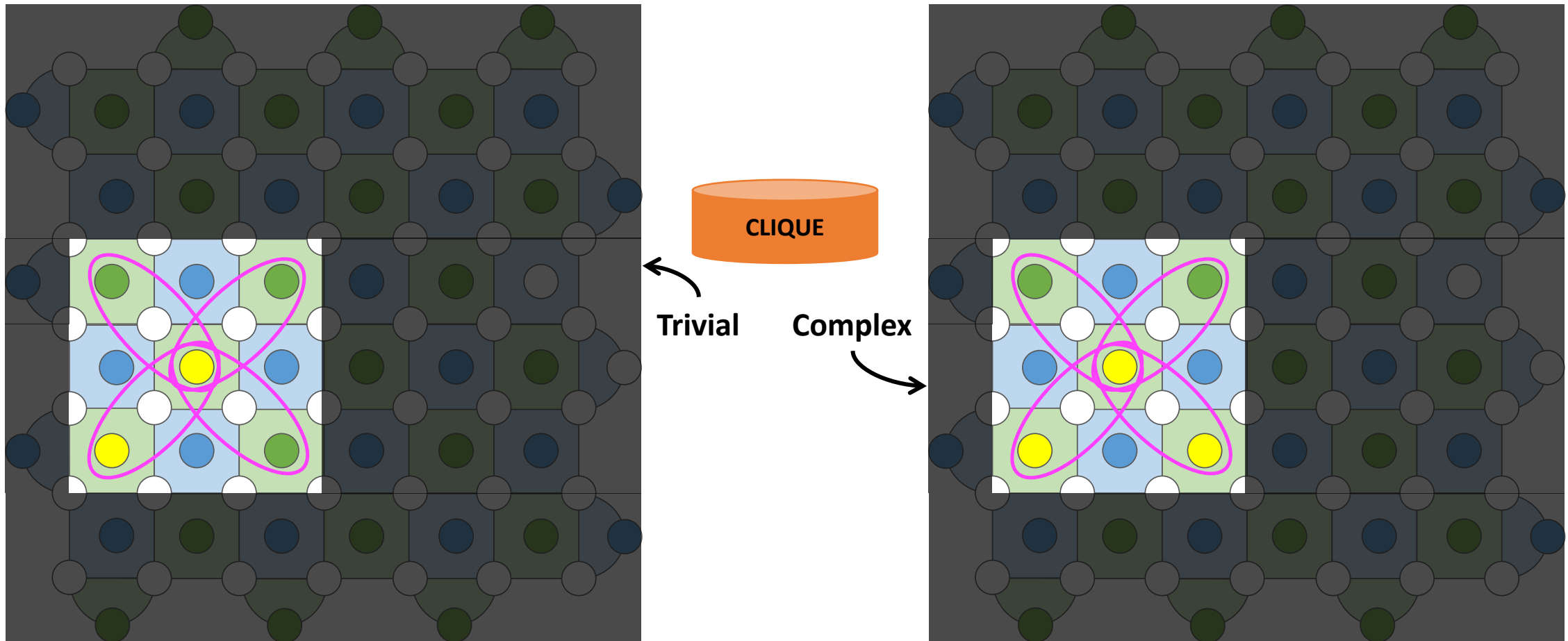


Clique decoder hardware design



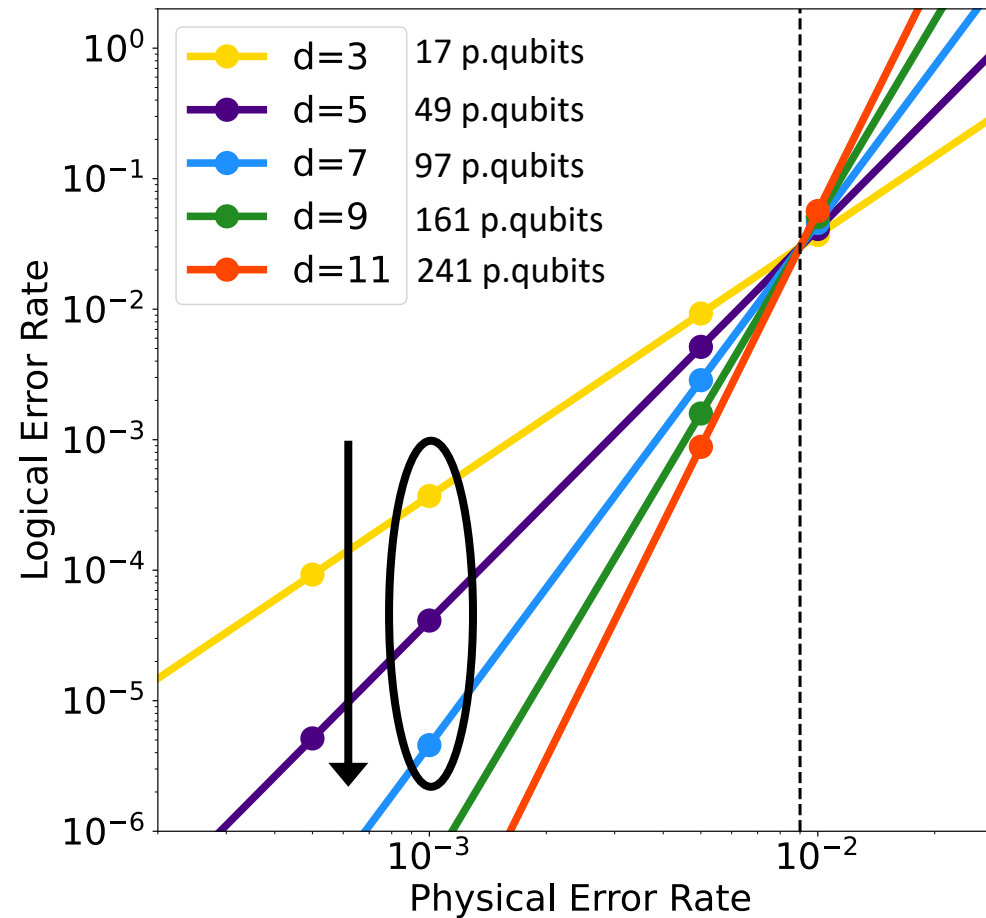
How to trivially detect isolated errors?

Isolated error litmus test: If the center of a clique is set, and if an odd number of neighbor syndromes are set, the clique can be trivially decoded.

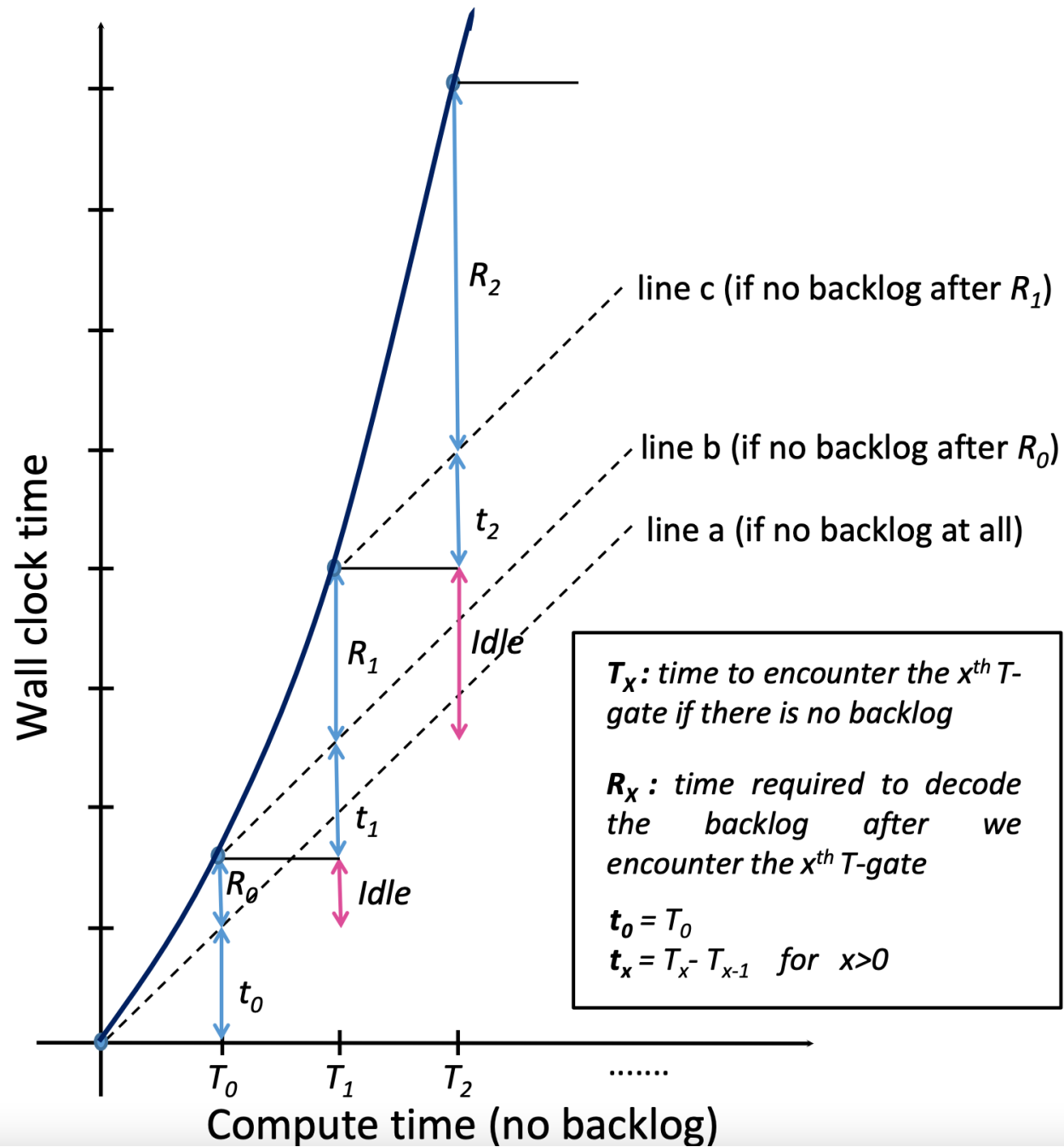


How QEC works: Surface code

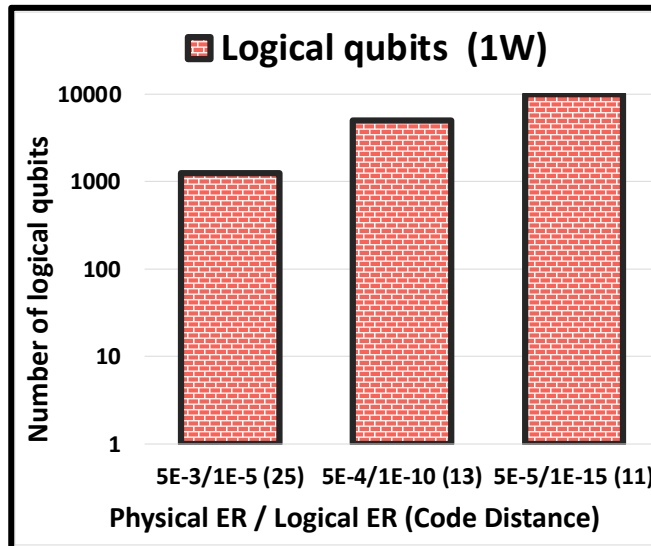
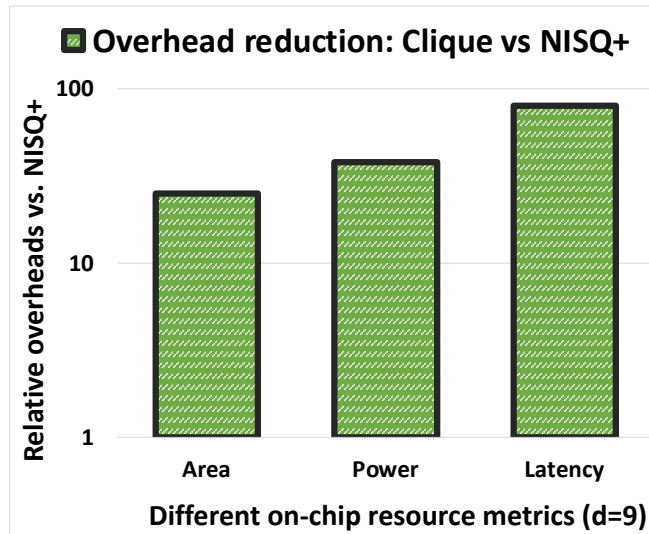
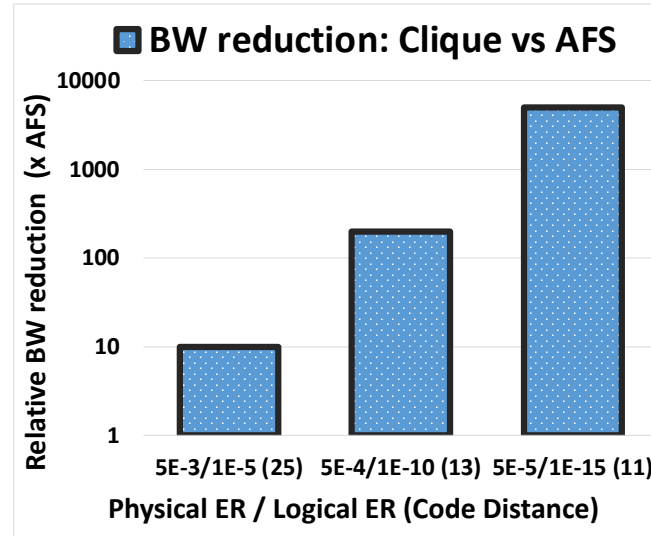
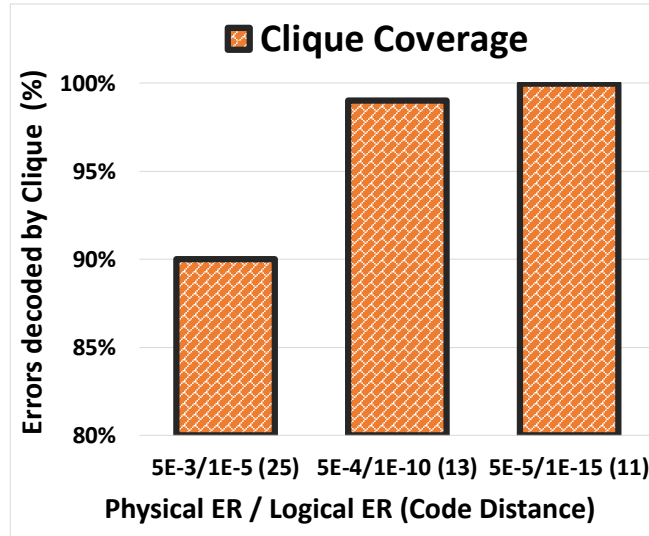
Better logical qubits possible with higher code distance, but with increased overheads.



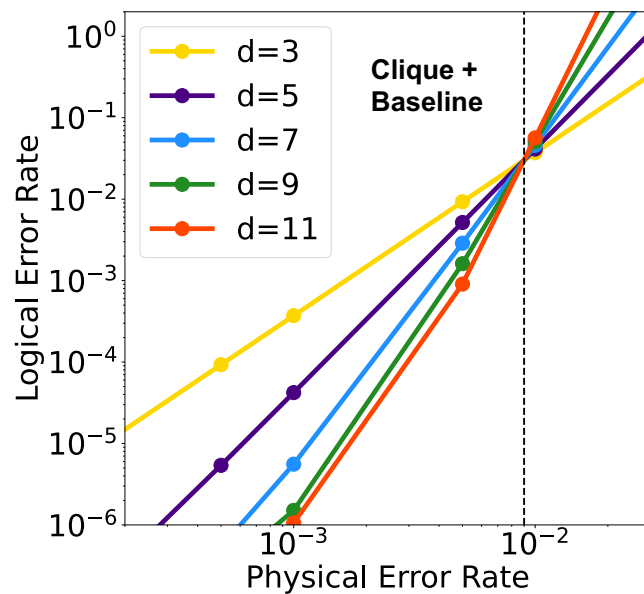
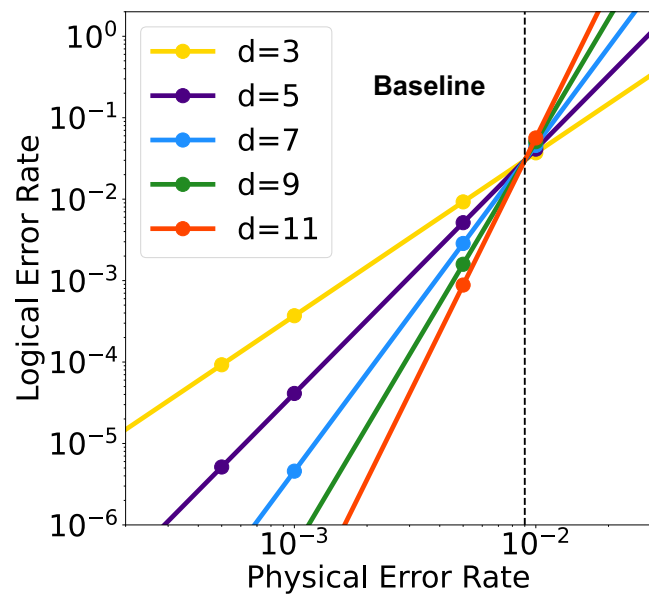
**Logical error decreases
with code distance**



What are the quantitative benefits?

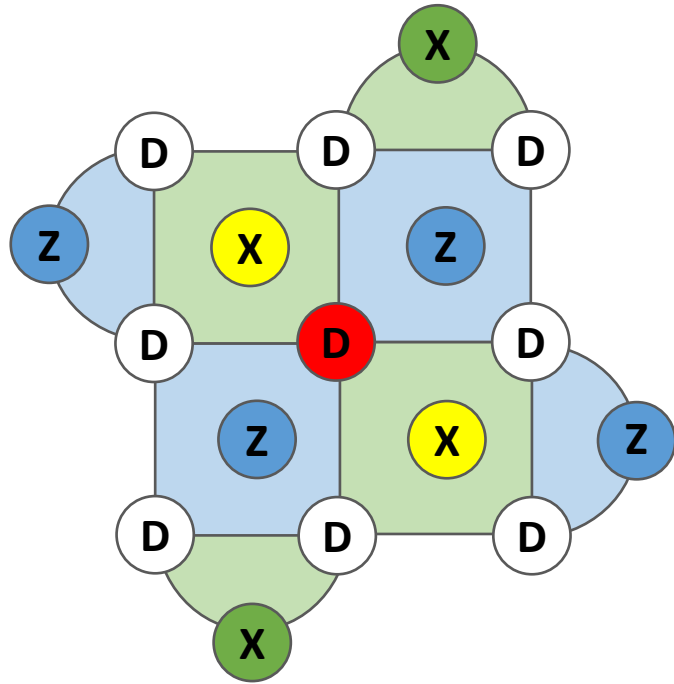


Logical error rate

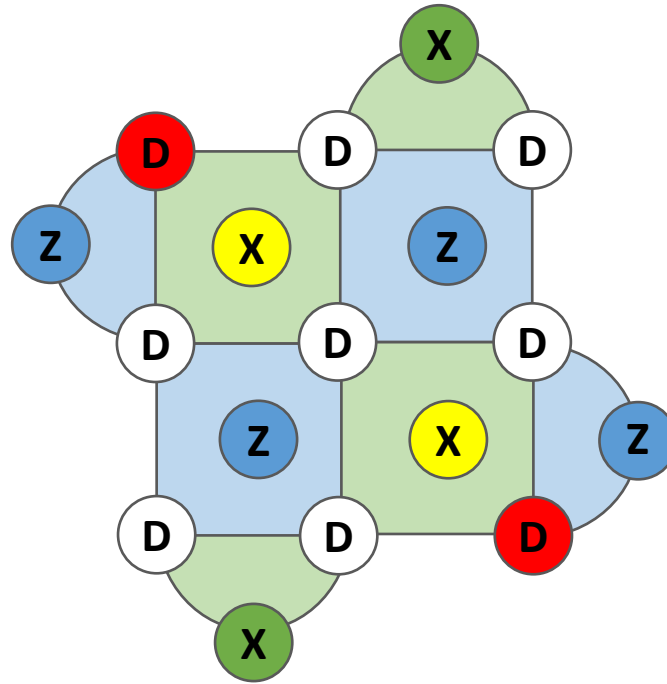


How does QEC decoding work?

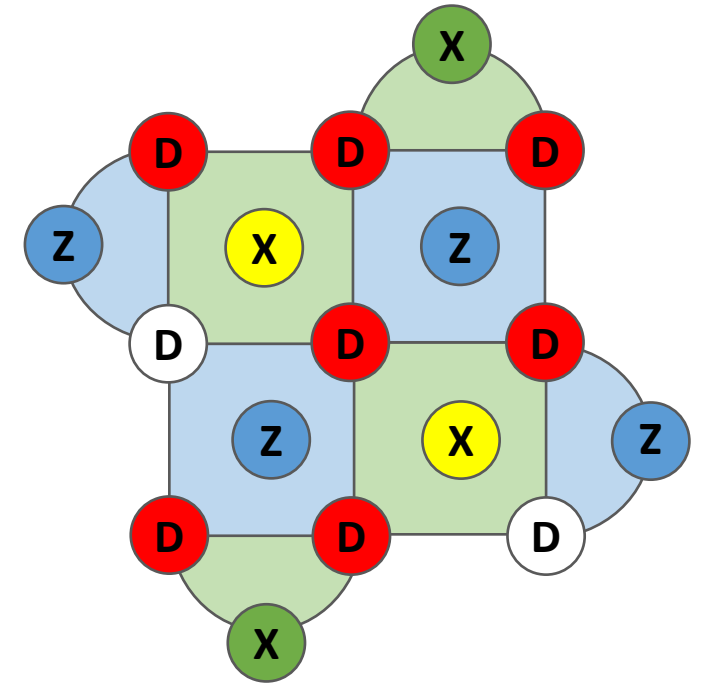
Likelihood of 1 data qubit error: $P = 10^{-2}$



Likelihood of 2 errors: $P^2 = 10^{-4}$



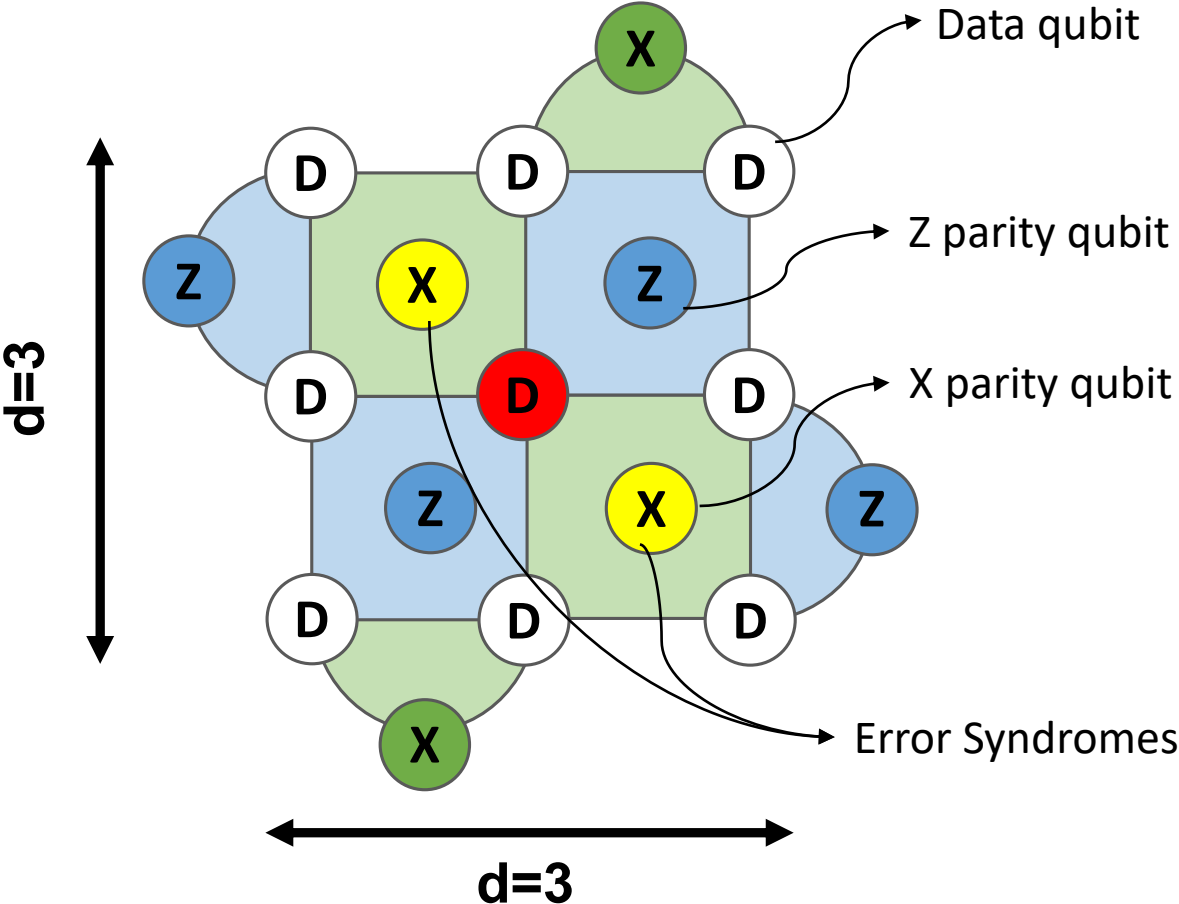
Likelihood of 7 errors: $P^7 = 10^{-14}$



The decoder returns a solution that is most likely: a decoding that produces the lowest number of data errors that satisfies the error syndrome pattern.

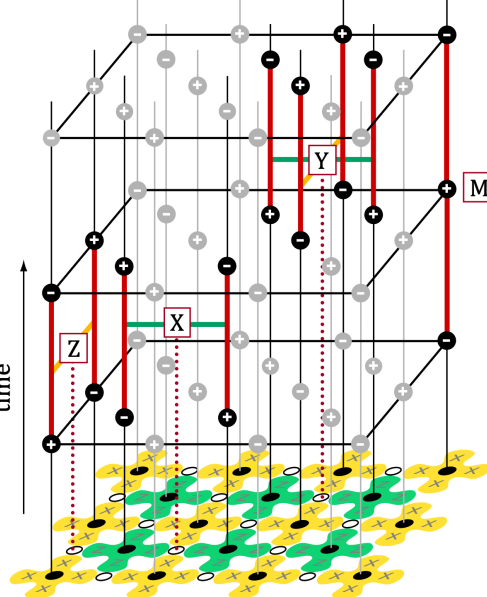
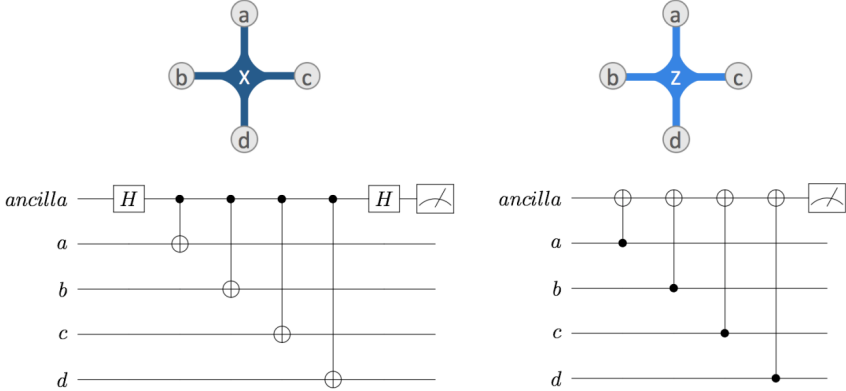
Background: Surface Codes

1 logical qubit w/ rotated surface code



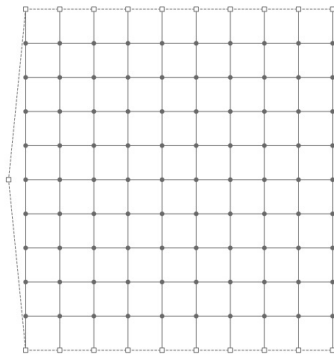
Support for high physical error rates

NISQ+ [Holmes2020]

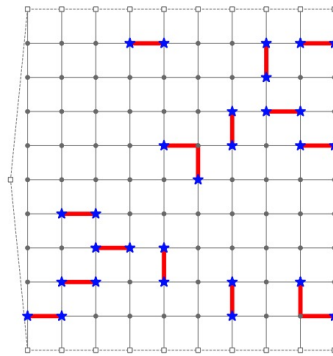


Surface Codes... [Fowler2012]

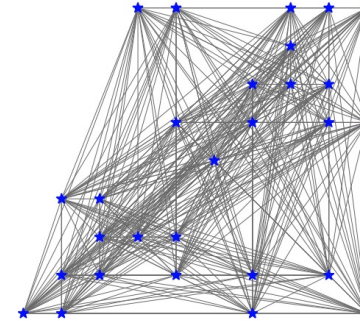
Decoding is complex at large code distances



(a) Matching graph

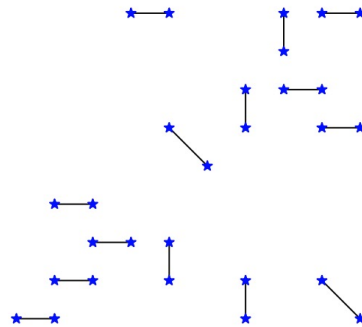


(b) Error

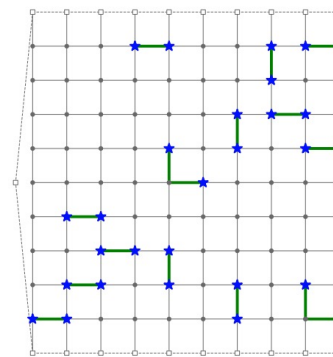


(c) Syndrome graph

1) Increases with code distance.
2) Multiplied by number of logical qubits.



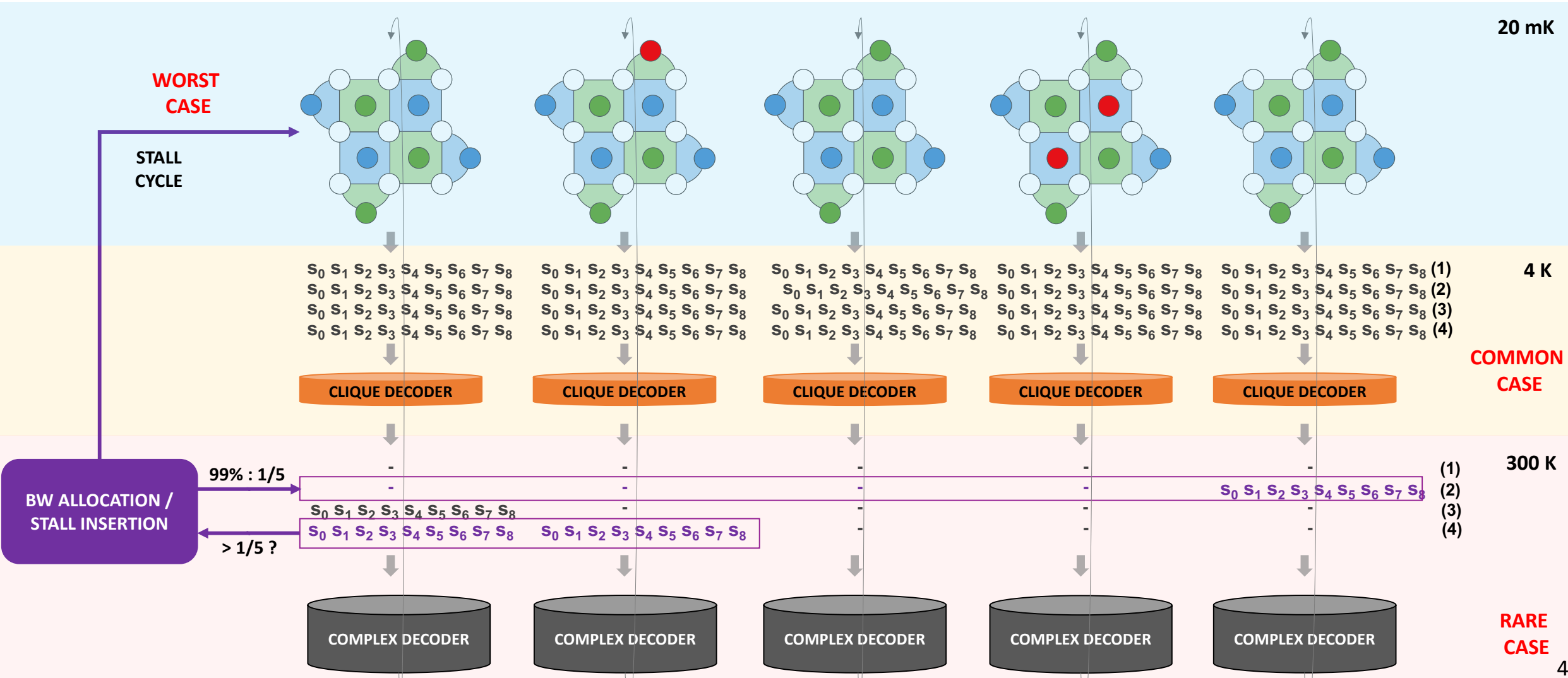
(d) Minimum-weight perfect matching



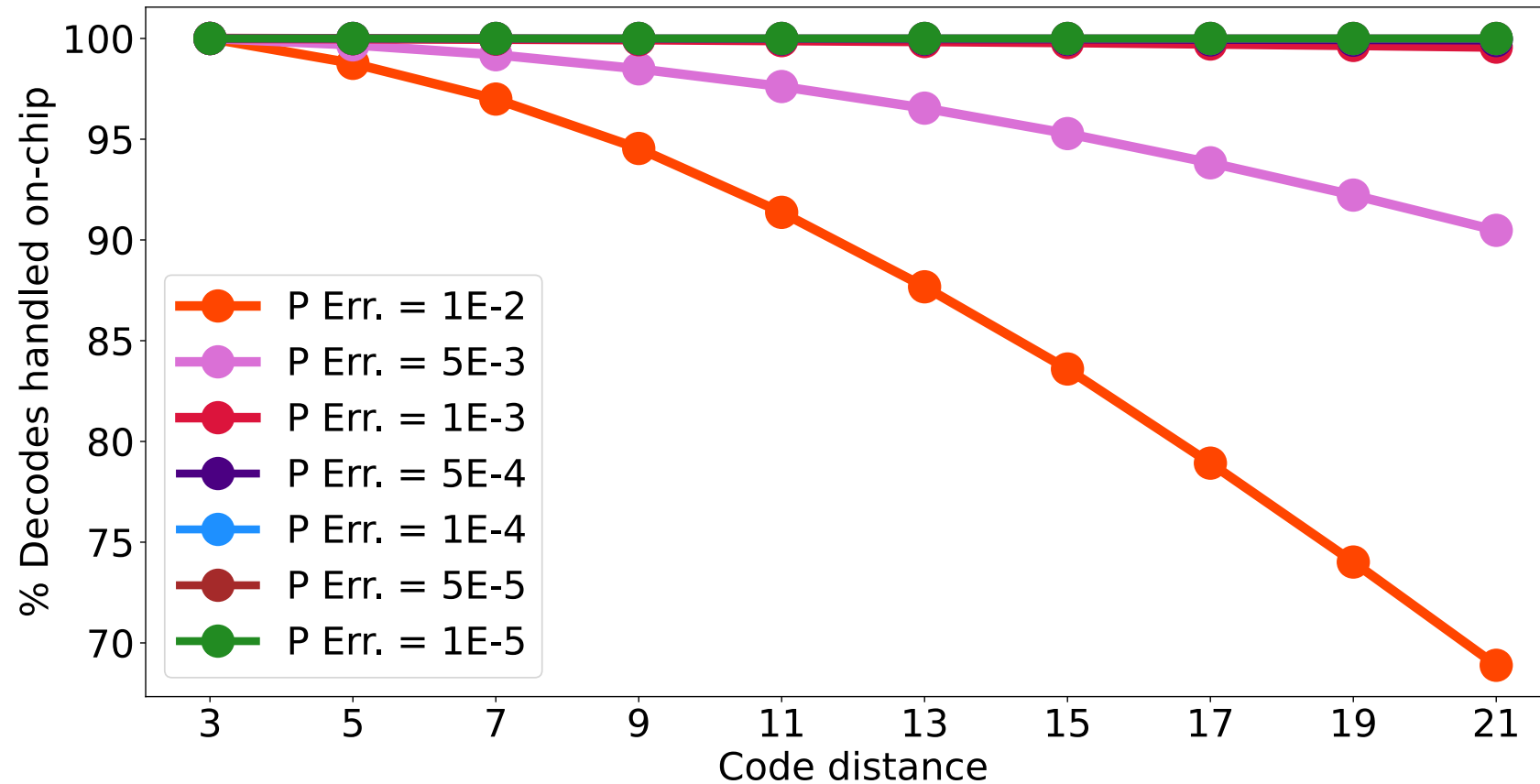
(e) Correction

pyMatching
[Higgott2021]

Proposal: Better than worst case decoding for QEC

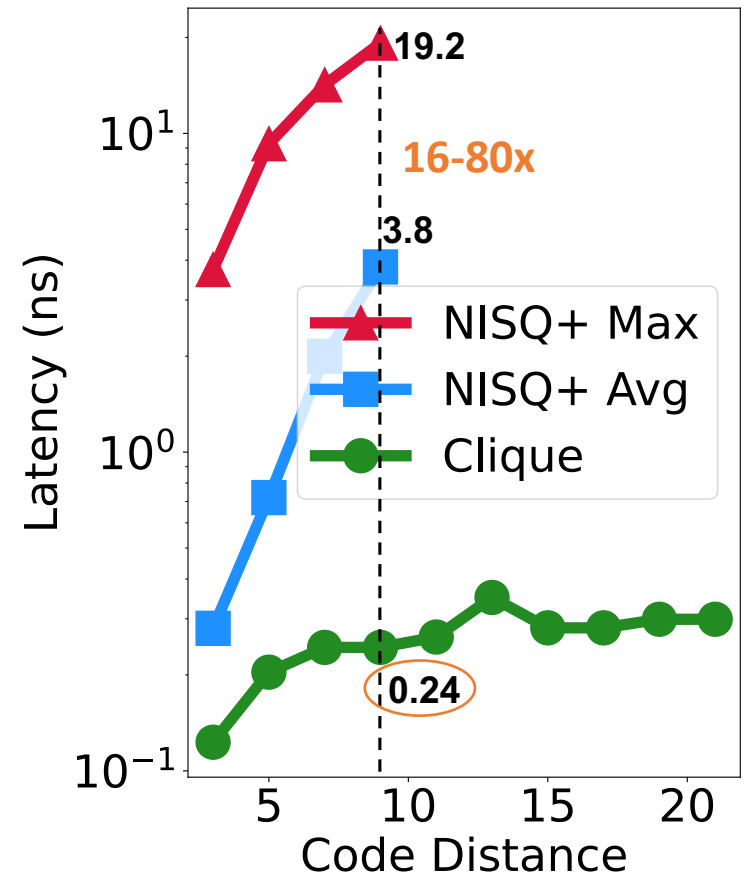
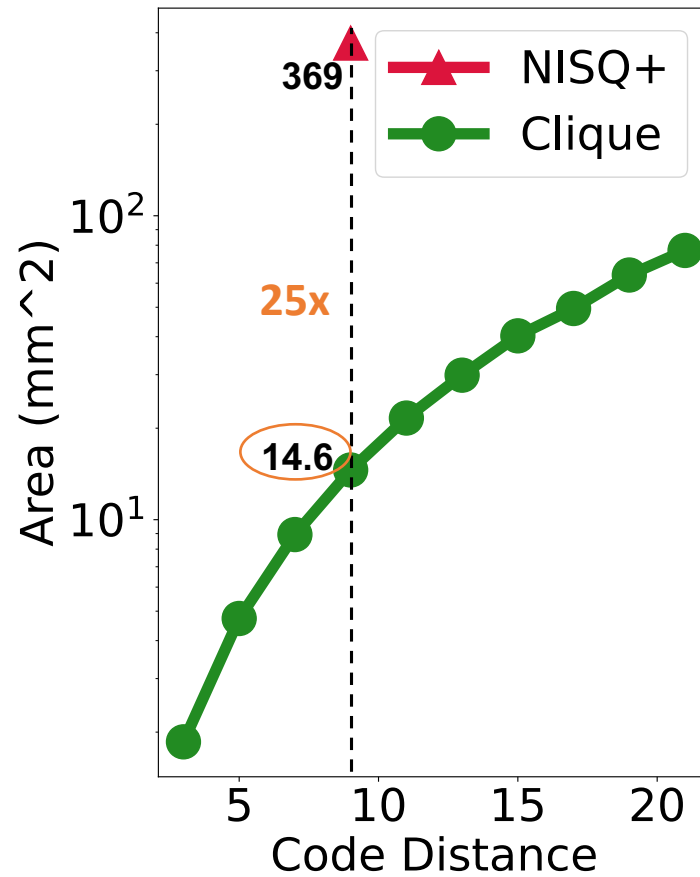
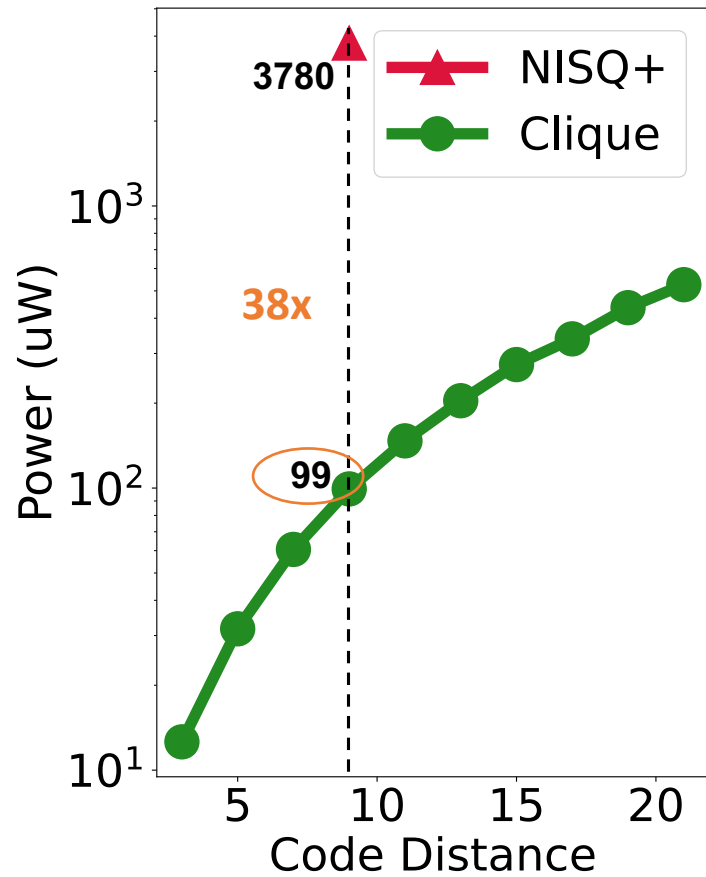


Results: Clique Decoder Coverage

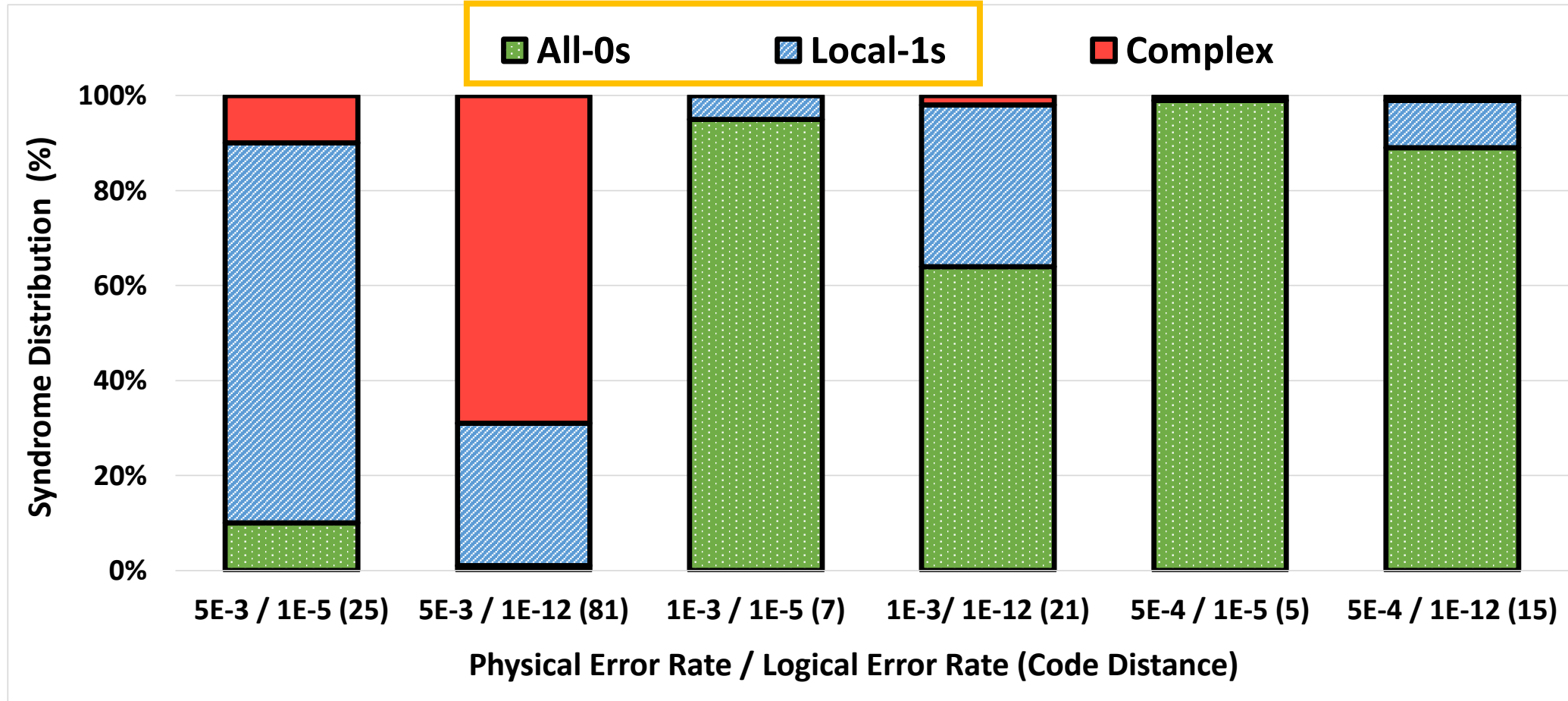


10-1000x greater bandwidth reduction compared to AFS which is entirely off-chip decoding but employs data compression on the syndrome data that is be sent off chip

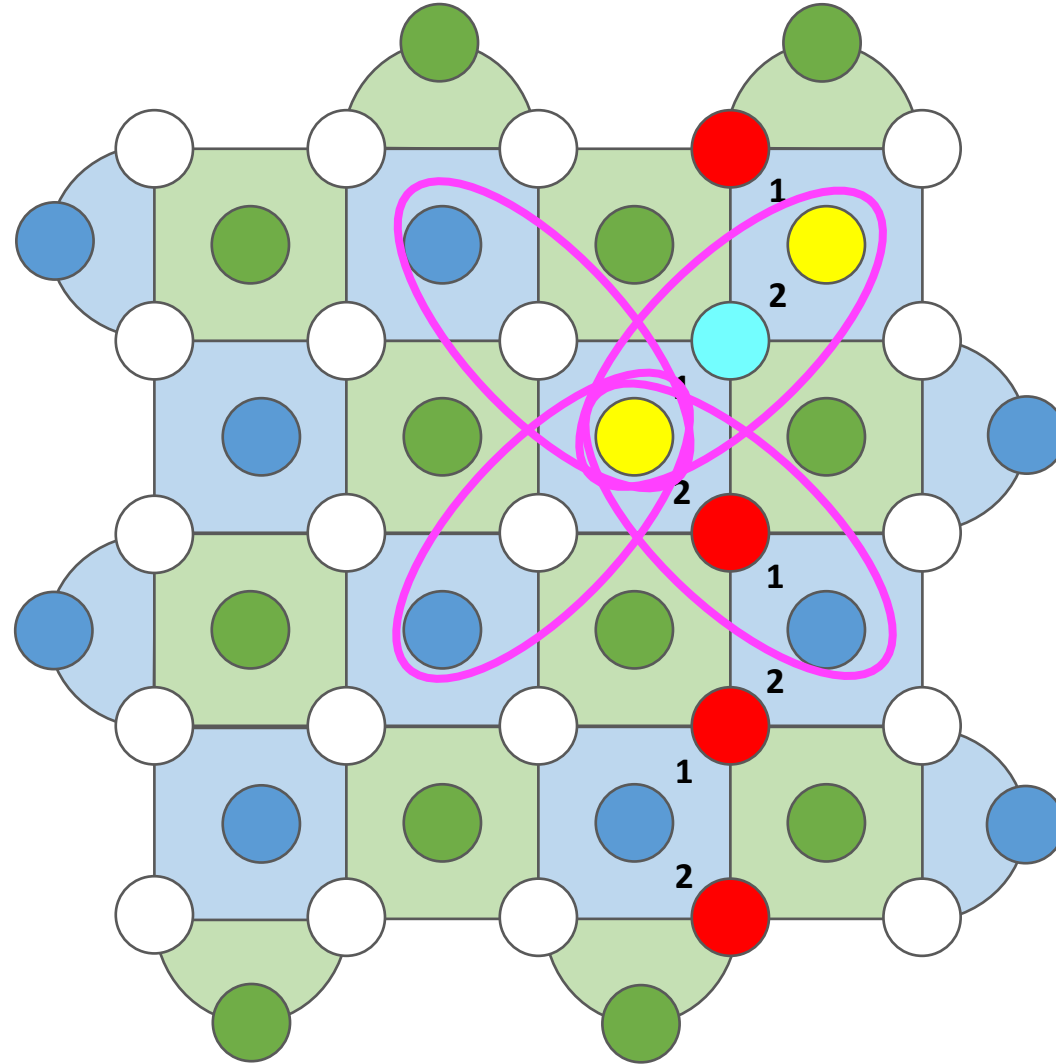
Results: Overheads compared to NISQ+



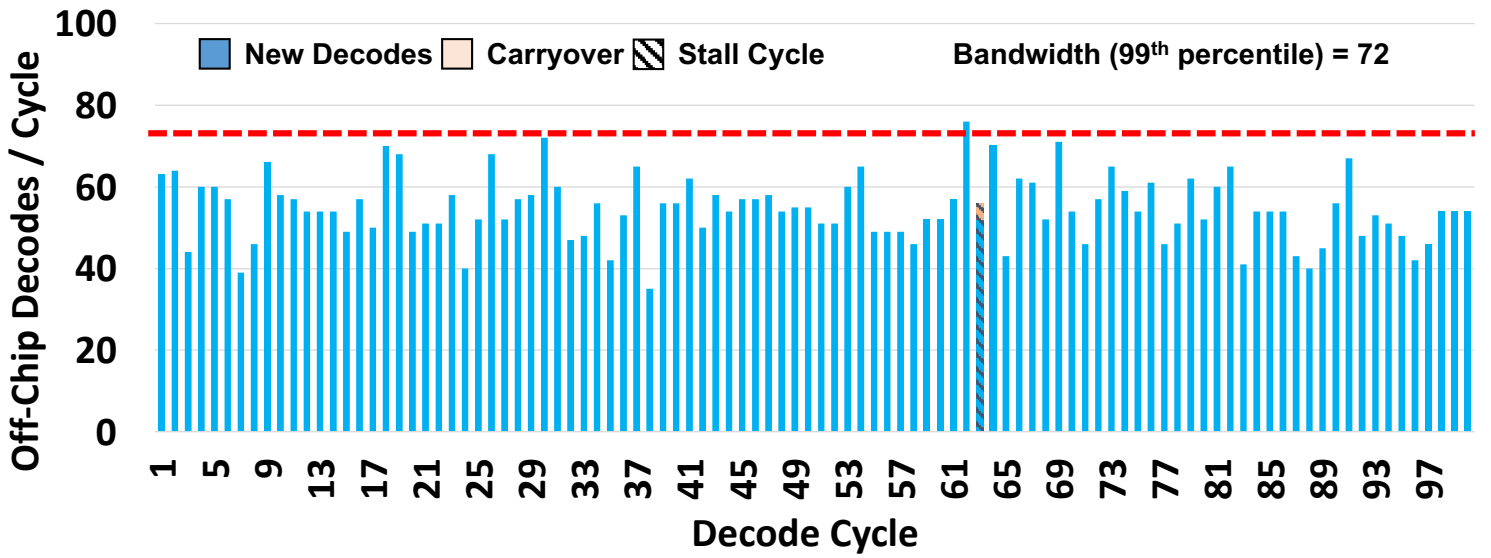
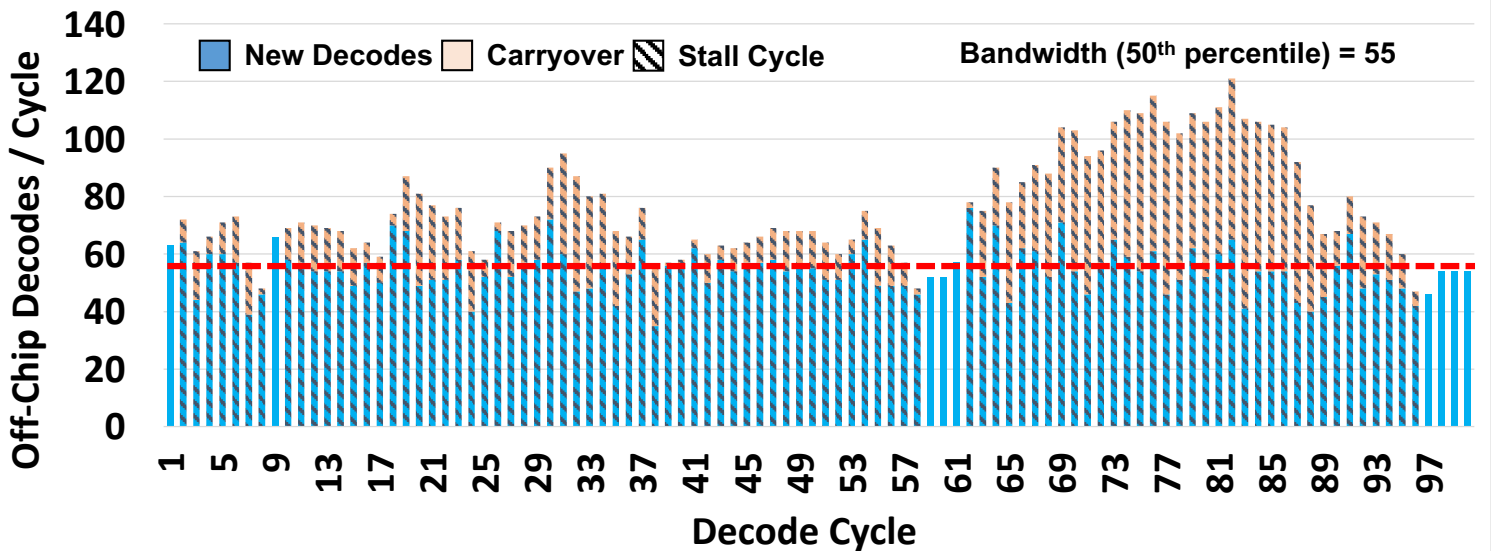
Observation: Error distribution vs Error rates



Logical errors (both Clique and complex decoder)

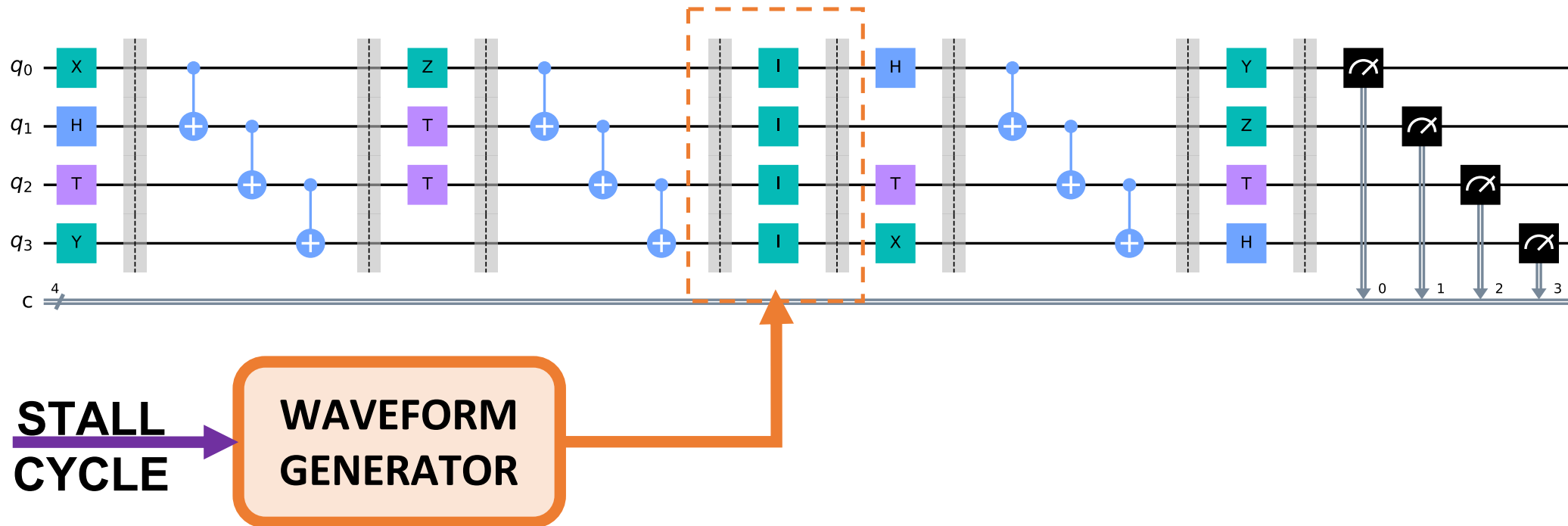


Statistical Off-chip Bandwidth Allocation

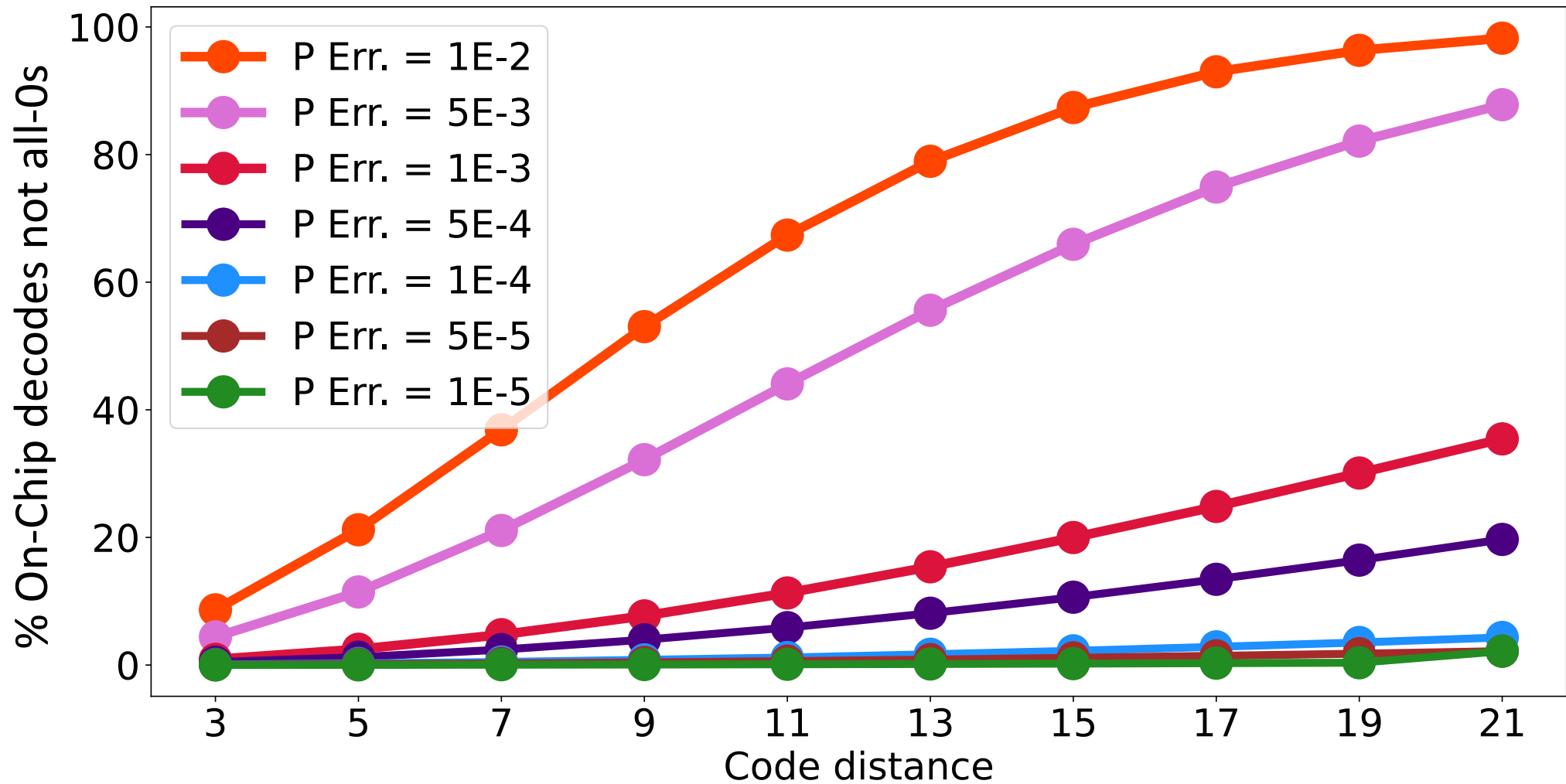


Errors need to be resolved every cycle*

Idle cycle insertion on stall cycle



Results: Clique Decoder Coverage – not all 0s



Results: Bandwidth Allocation vs Stalling trade-offs

