





Vegapunk: Accurate and Fast Decoding for Quantum LDPC Codes with Online Hierarchical Algorithm and Sparse Accelerator

Kaiwen Zhou¹, Liqiang Lu^{1*}, Debin Xiang¹, Chenning Tao¹, Anbang Wu², Jingwen Leng², Fangxin Liu², Mingshuai Chen¹, Jianwei Yin^{1*}

¹Zhejiang University, ²Shanghai Jiao Tong University

MICRO 2025, Seoul, Korea

Outline

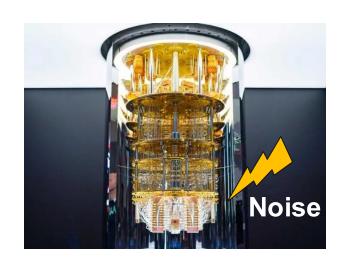


- Background
- Motivation
- Vegapunk Algorithm
- Vegapunk Accelerator
- Experiment

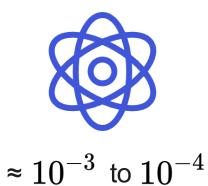
Quantum Computers are Noisy



The most significant barrier for large-scale quantum computing: Errors



Current Physical Error Rate



Required Error Rate

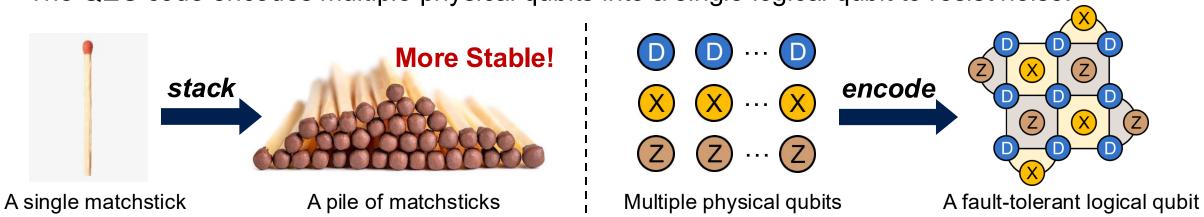


Quantum Error Correction (QEC) is crucial for practical, large-scale quantum computing.

Basics of QEC Code



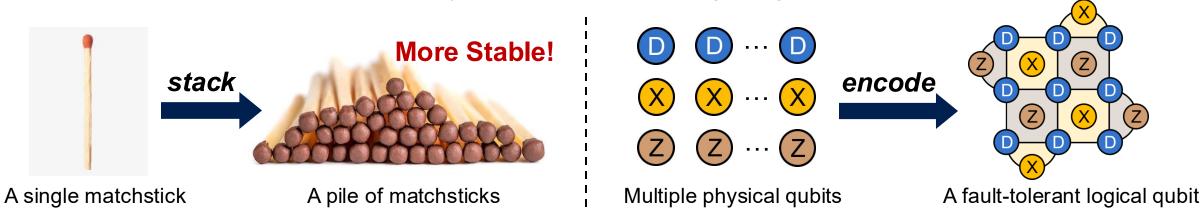
The QEC code encodes multiple physical qubits into a single logical qubit to resist noise.



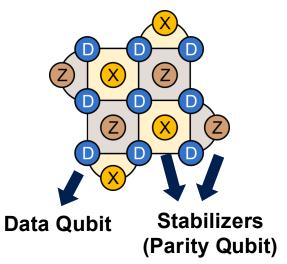
Basics of QEC Code



The QEC code encodes multiple physical qubits into a single logical qubit to resist noise.



Consider surface code [[9, 1, 3]] as an example to illustrate key parameters [[n, k, d]] of QEC codes.



Key Parameters [[n, k, d]]

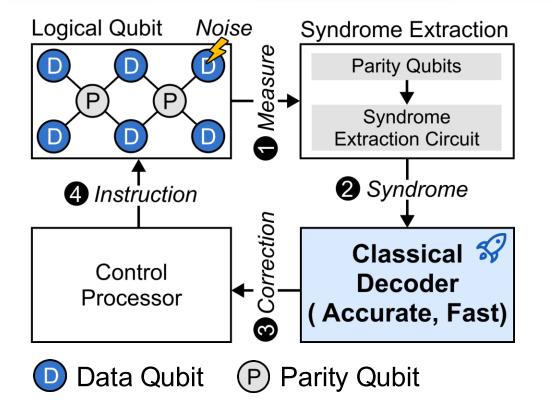
- **n**: the number of data qubits in a logical block.
- **k**: the number of encoded logical qubits.
- **d**: the code distance, quantifying the error tolerance.
- Encoding Rate: k / (number of physical qubits).

Error Detection

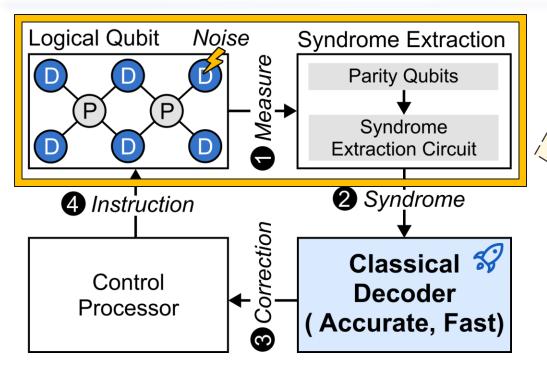
- Bit Flip (X) Errors → Z
 Stabilizers
- Phase Flip (Z) Errors → X
 Stabilizers







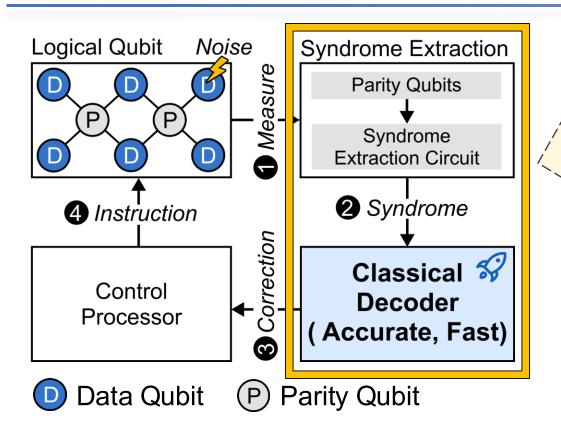




1 The syndrome extraction circuit is executed on parity qubits to **generate the syndrome**.

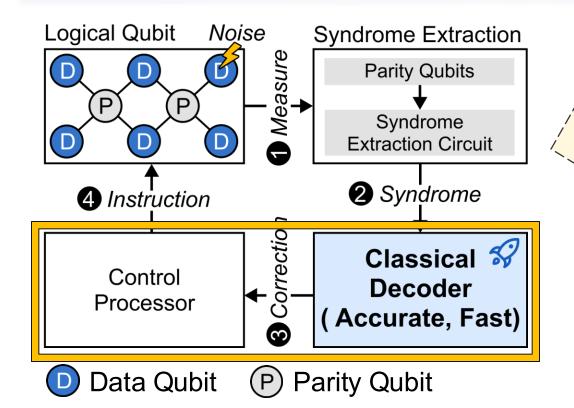
Data Qubit Parity Qubit





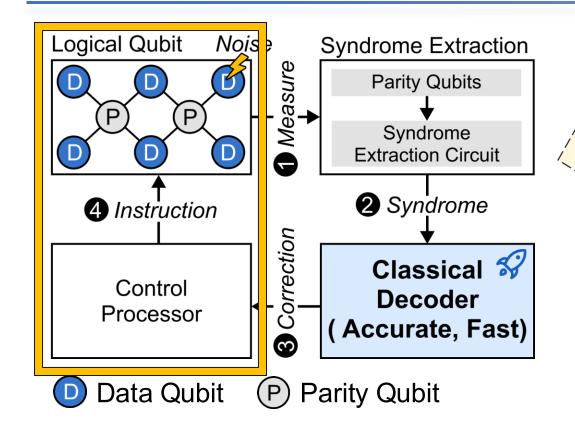
- 1 The syndrome extraction circuit is executed on parity qubits to **generate the syndrome**.
- ② The syndrome is sent to a classical decoder, which identifies error locations on data qubits.





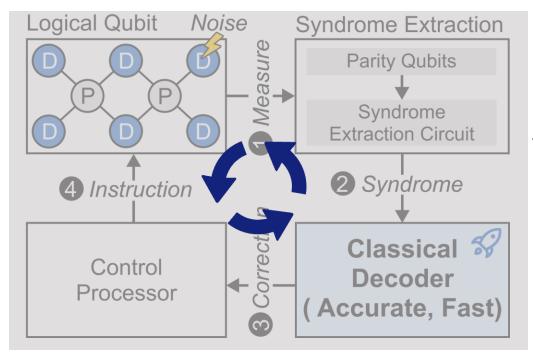
- 1 The syndrome extraction circuit is executed on parity qubits to **generate the syndrome**.
- ② The syndrome is sent to a classical decoder, which identifies error locations on data qubits.
- 3 The decoder **computes a logical correction** and sends it to the control processor.





- 1 The syndrome extraction circuit is executed on parity qubits to **generate the syndrome**.
- ② The syndrome is sent to a classical decoder, which identifies error locations on data qubits.
- 3 The decoder computes a logical correction and sends it to the control processor.
- 4 The control processor receives the correction and sends QEC instructions to correct errors.





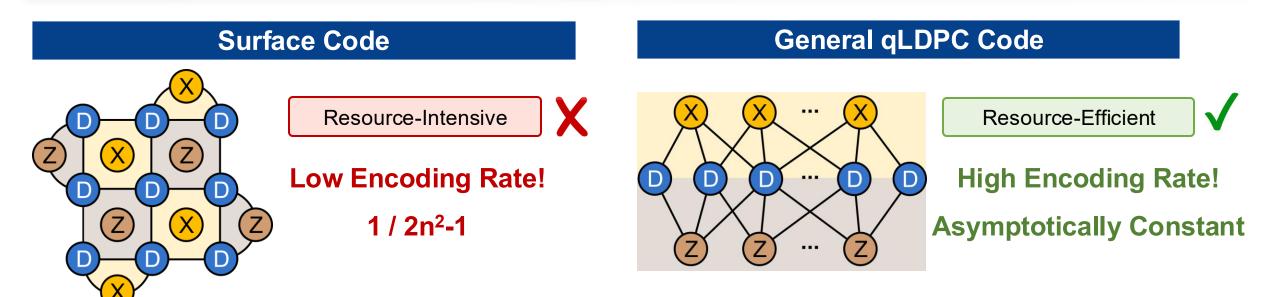
- ① The syndrome extraction circuit is executed on parity qubits to **generate the syndrome.**
- ② The syndrome is sent to a classical decoder, which identifies error locations on data qubits.
- 3 The decoder computes a logical correction and sends it to the control processor.
- The control processor receives the correction and sends QEC instructions to correct errors.

- Data Qubit
- P Parity Qubit
- The QEC process is executed in cycles to protect quantum information.
- Decoder identifies the location and type of errors, according to the syndrome.

Scalable Architectural Solution: qLDPC Code







Higher encoding rate means fewer physical qubits required for the same number of logical qubits.

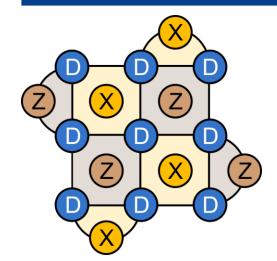
Scalable Architectural Solution: qLDPC Code





Surface Code

General qLDPC Code

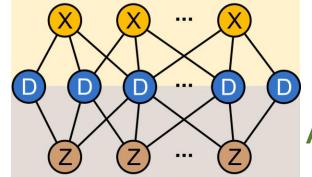


Resource-Intensive



Low Encoding Rate!

1 / 2n²-1



Resource-Efficient



High Encoding Rate!

Asymptotically Constant

- Higher encoding rate means fewer physical qubits required for the same number of logical qubits.
- Let's take an example:
- Surface code [[1452,12,11]] requires 1452 data qubits to encode 12 logical qubits.
- Bivariate Bicycle (BB) code [[144,12,12]], a type of qLDPC codes, achieves the same number of logical qubits using only 144 data qubits!

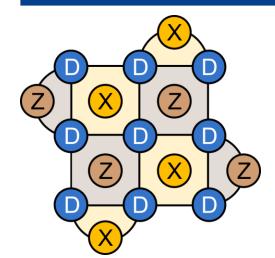
Scalable Architectural Solution: qLDPC Code





Surface Code

General qLDPC Code

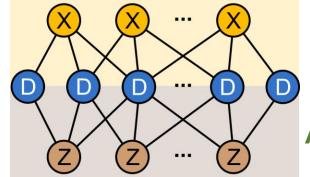


Resource-Intensive



Low Encoding Rate!

1 / 2n²-1



Resource-Efficient



High Encoding Rate!
Asymptotically Constant

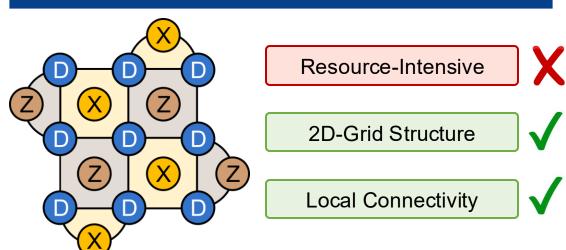
- Higher encoding rate means fewer physical qubits required for the same number of logical qubits.
- Let's take an example:
- Surface code [[1452,12,11]] requires 1452 data qubits to encode 12 logical qubits.

We focus on qLDPC code, since it offers a scalable solution for superconducting platforms.

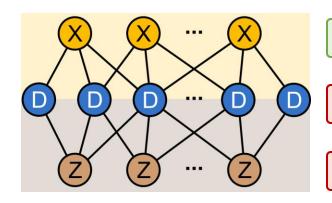
qLDPC Decoding is Not Easy



Surface Code



General qLDPC Code



Resource-Efficient



Hypergraph Structure



Non-Local Connectivity

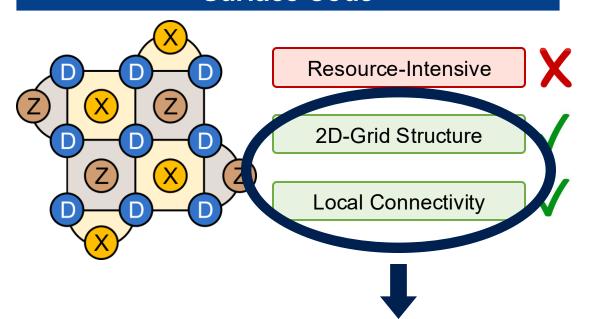




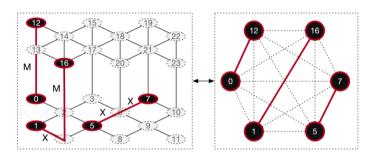
qLDPC Decoding is Not Easy



Surface Code

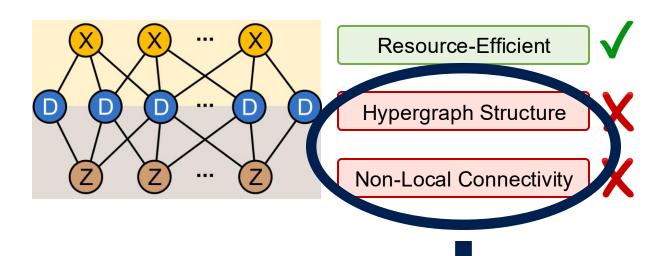


Efficient Decoding (Polynomial-Time)



Minimum Weight **Perfect Matching** (MWPM)

General qLDPC Code



Inefficient Decoding (NP-Hard)

$$e^{\text{ml}} = \arg \max_{e} \{\ln(P(e))\}$$

$$= \arg \max_{e} \left(C - \sum_{j} w_{j} e_{j}\right)$$
s.t. $s = H_{Z}e$,

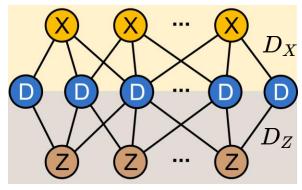
Maximum Likelihood Decoding (MLD) [1]

Decoding Objective and Quantum Degeneracy

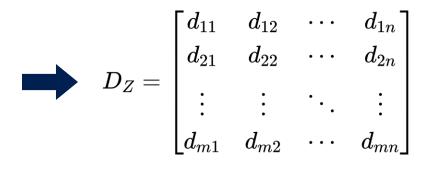




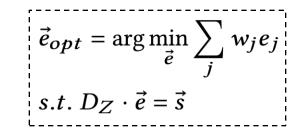
qLDPC decoding can be formulated as a minimum weight decoding problem, which is NP-hard.



Tanner Graph



Parity Check Matrix



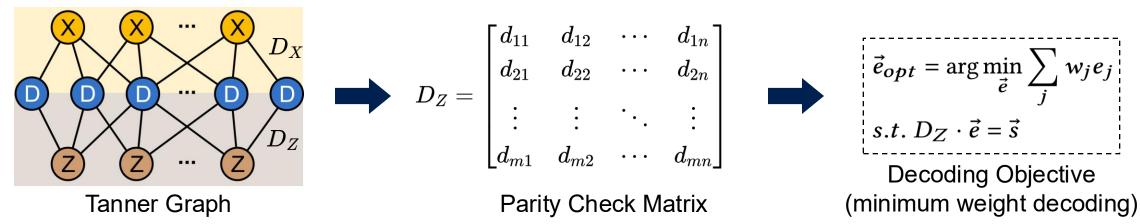
Decoding Objective (minimum weight decoding)

Decoding Objective and Quantum Degeneracy

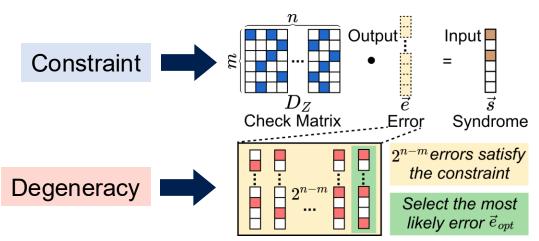




qLDPC decoding can be formulated as a minimum weight decoding problem, which is NP-hard.



Decoding suffers from quantum degeneracy, where multiple errors satisfy the check constraint.

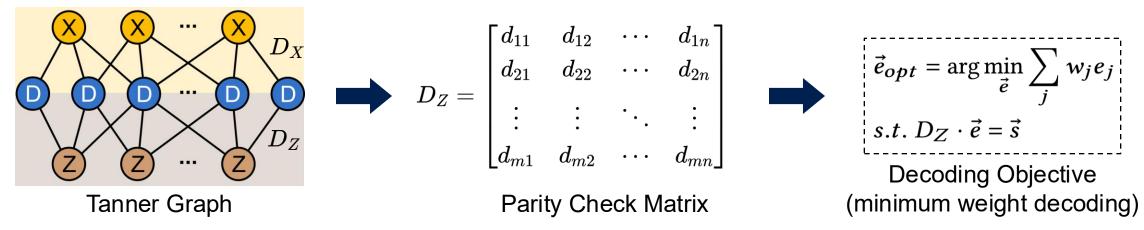


Decoding Objective and Quantum Degeneracy (§

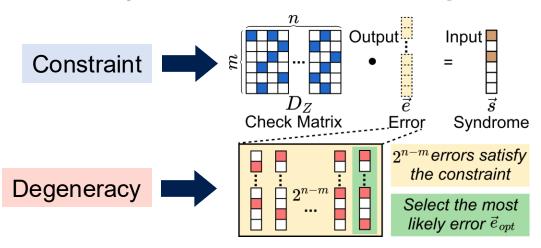


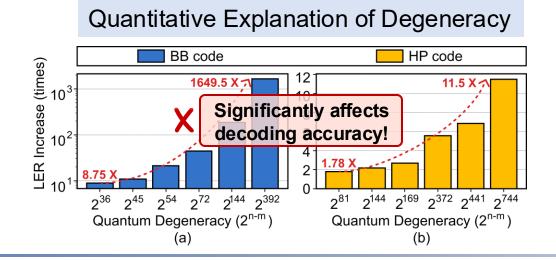


qLDPC decoding can be formulated as a minimum weight decoding problem, which is NP-hard.



Decoding suffers from quantum degeneracy, where multiple errors satisfy the check constraint.





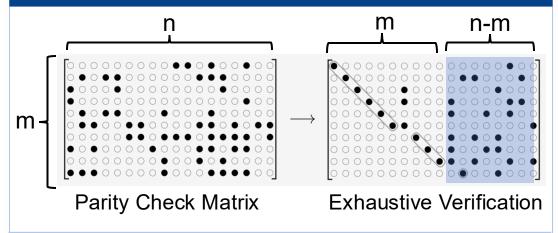




Belief Propagation (BP) $|\mu_{v_2}|$ $[\mu_{v_3}]$ $|\mu_{v_4}|$ $|\mu_{v_5}|$ $|\mu_{v_6}|$ $|\mu_{v_7}|$ Prior Probs. **Data Qubits** v_2 v_3 v_4 v_5 v_7 v_6 **✓** Iterate c_2 Parity Checks Tanner Graph **Syndromes**

- BP iteratively updates probability information by passing messages on a Tanner graph to identify the most likely error pattern.
- Core characteristics: high parallelism and low computational complexity.

BP+Ordered Statistics Decoding (OSD)



- If BP decoding fails, OSD part executes Gaussian elimination, receiving probability information from BP to exhaustively verify possible solutions.
- Core characteristics: high decoding accuracy and high computational complexity.

The vast majority of existing qLDPC decoders are built upon BP/OSD improvements.

Outline



- Background
- Motivation
- Vegapunk Algorithm
- Vegapunk Accelerator
- Experiment

1. Prior Decoders: Accuracy-Latency Tradeoff





Existing qLDPC decoders fail to achieve both accurate and fast decoding.



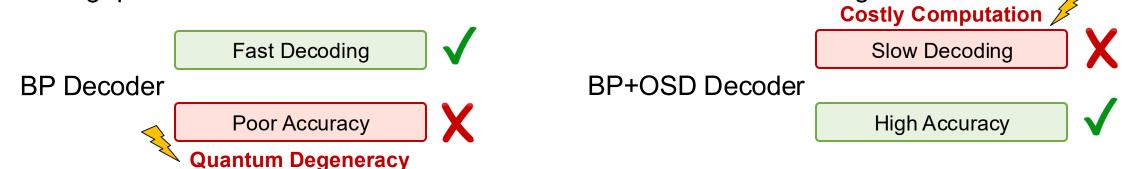


1. Prior Decoders: Accuracy-Latency Tradeoff

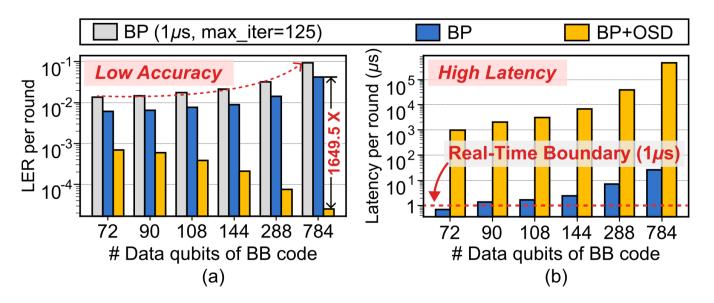




Existing qLDPC decoders fail to achieve both accurate and fast decoding.



Decoding performance profiling of BP and BP+OSD decoders under 0.1% circuit-level noise.



Observations:

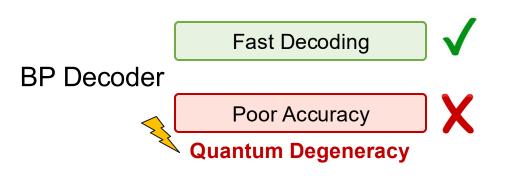
- BP's accuracy is significantly lower than BP+OSD's, rendering it unacceptable.
- Neither BP nor BP+OSD can meet the real-time decoding requirements (1µs) when scaled to large-scale qLDPC codes.

1. Prior Decoders: Accuracy-Latency Tradeoff





Existing qLDPC decoders fail to achieve both accurate and fast decoding.
 Costly Computation



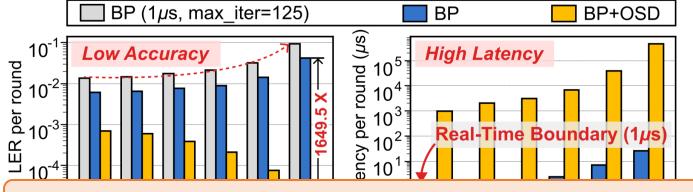
BP+OSD Decoder

Slow Decoding

High Accuracy



Decoding performance profiling of BP and BP+OSD decoders under 0.1% circuit-level noise.



Observations:

- BP's accuracy is significantly lower than BP+OSD's, rendering it unacceptable.
- Neither BP nor BP+OSD can meet the

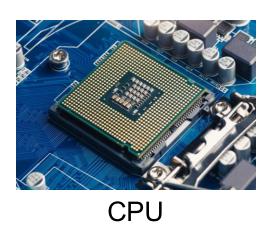


Goal: Develop an accurate qLDPC decoder that can process each round of syndrome in 1µs.

2. Low Latency: Why Dedicated Accelerator?



• General-purpose processors (e.g., CPUs or GPUs) are inefficient at handling the fine-grained and bit-level operations required in decoding.







Decoding Process

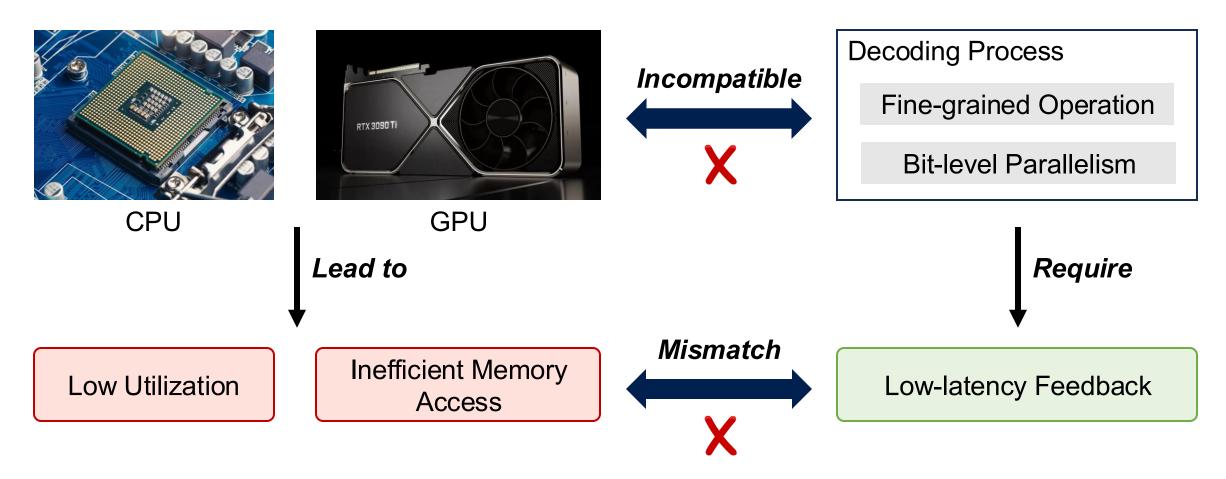
Fine-grained Operation

Bit-level Parallelism

2. Low Latency: Why Dedicated Accelerator?



 General-purpose processors (e.g., CPUs or GPUs) are inefficient at handling the fine-grained and bit-level operations required in decoding.



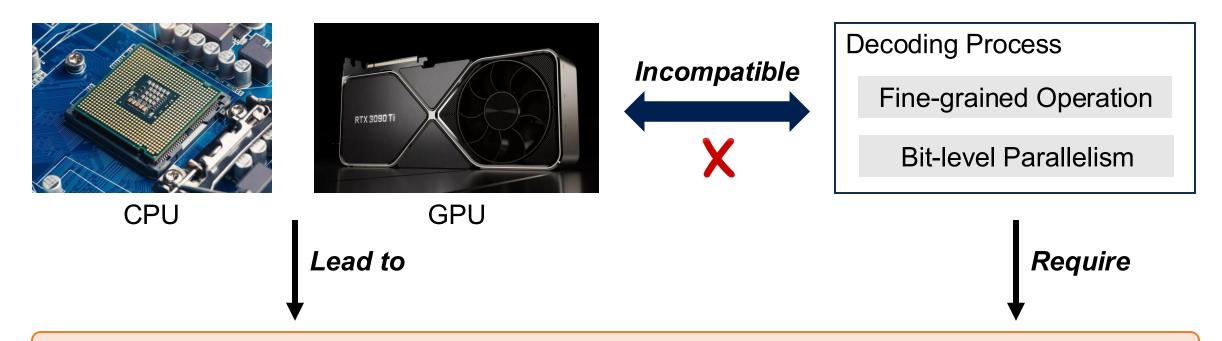
2. Low Latency: Why Dedicated Accelerator?





27

 General-purpose processors (e.g., CPUs or GPUs) are inefficient at handling the fine-grained and bit-level operations required in decoding.





We choose to design a dedicated accelerator for qLDPC decoding based on FPGAs.

Outline

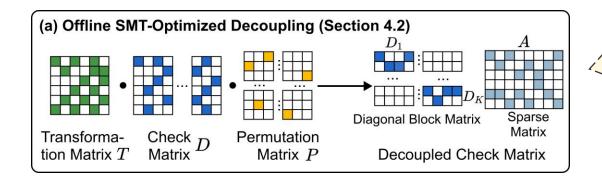


- Background
- Motivation
- Vegapunk Algorithm
- Vegapunk Accelerator
- Experiment



Workflow of Vegapunk Algorithm





To Address **Low Accuracy**:

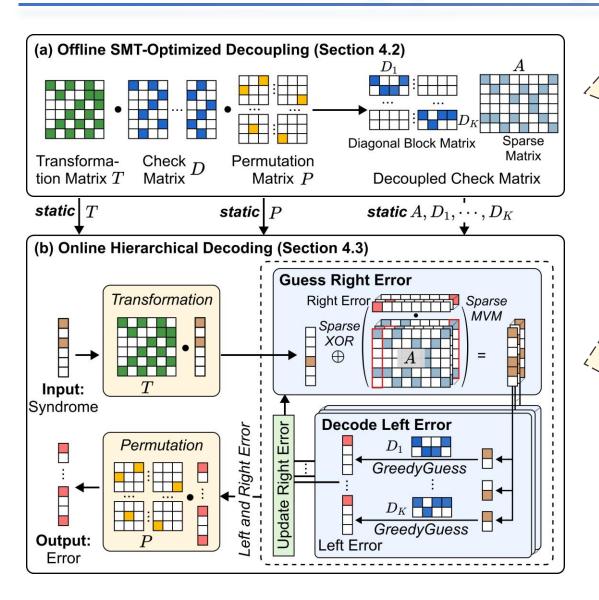
- **Insight:** Degeneracy arises from the check matrix having far more columns than rows.
- **Solution:** Decouple the original check matrix into smaller sub-matrices, to balance the number of columns and rows.

29 kaiwenzhou@zju.edu.cn Kaiwen Zhou



Workflow of Vegapunk Algorithm





To Address **Low Accuracy**:

- **Insight:** Degeneracy arises from the check matrix having far more columns than rows.
- **Solution:** Decouple the original check matrix into smaller sub-matrices, to balance the number of columns and rows.

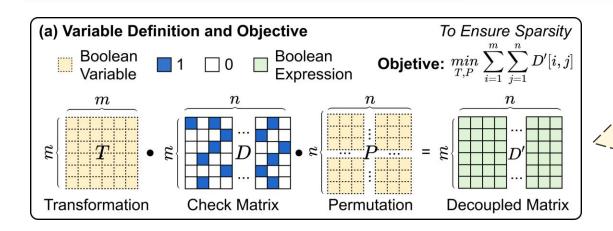
To Address **High Latency**:

- **Insight:** The probability of an error pattern exponentially decreases with the number of '1's it contains.
- **Solution:** A highly parallel greedy decoding algorithm that incrementally increases the number of '1's in error pattern, evaluates their likelihood at each step, and incorporates an early stopping mechanism.



Offline SMT-Optimized Decoupling



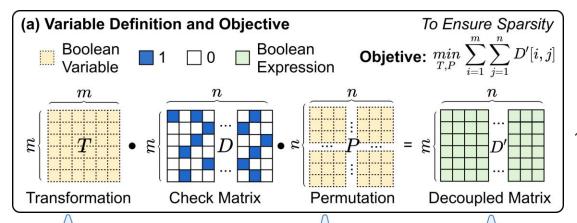


Fundamental Idea: Transform the Check Matrix

- **Left-multiplying** it with a transformation matrix *T*.
- **Right-multiplying** it with a permutation matrix *P*.

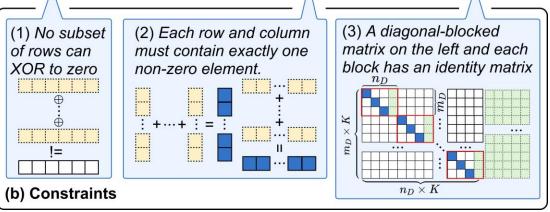
Offline SMT-Optimized Decoupling





Fundamental Idea: Transform the Check Matrix

- Left-multiplying it with a transformation matrix T.
- Right-multiplying it with a permutation matrix P.



Fundamental Idea: Reformulate into an SMT Problem

- By setting constraints, the decoupled matrix consists
 of a diagonal-block matrix and a sparse matrix.
- The objective of our SMT solver is to maximize the sparsity of the decoupled matrix.



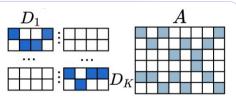
Permutation Matrix P



Transformation Matrix *T*



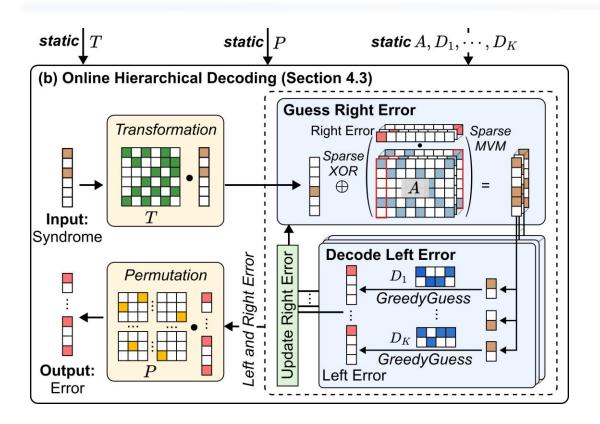
Decoupled Matrix *D*′





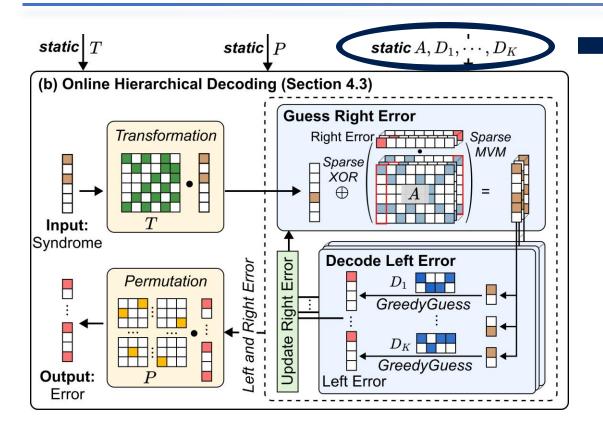




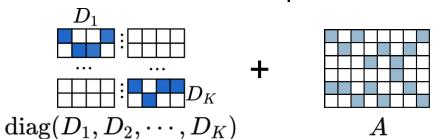




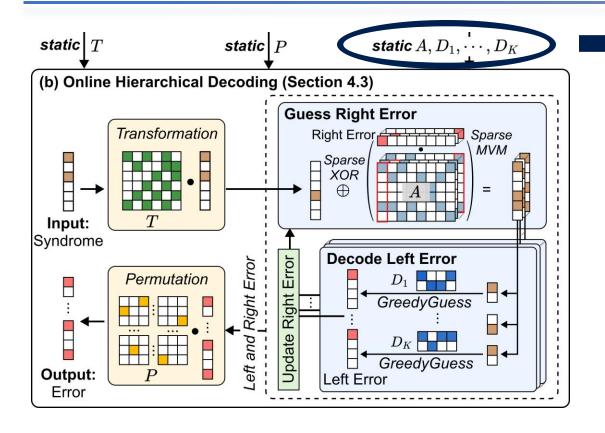




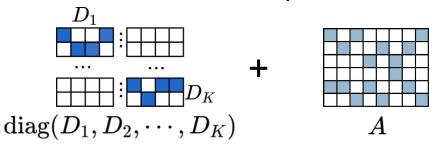
The check matrix is decomposed into two parts:



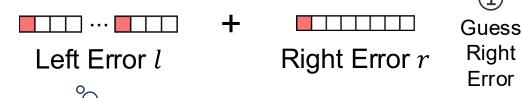




The check matrix is decomposed into two parts:

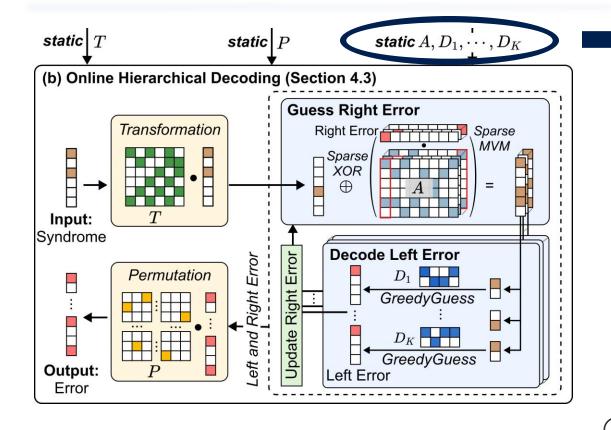


Therefore, we propose to partition the full error into *left error* and *right error* accordingly:

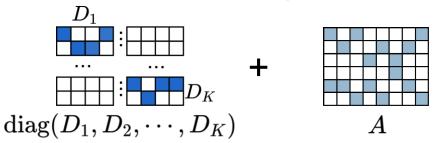




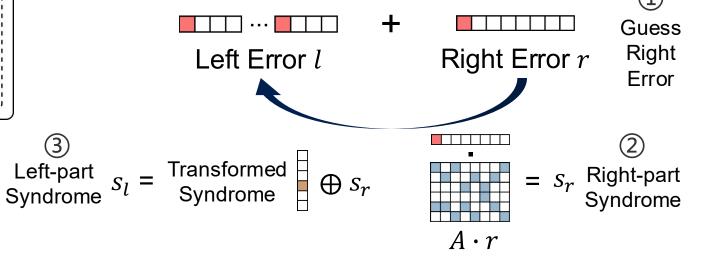




The check matrix is decomposed into two parts:



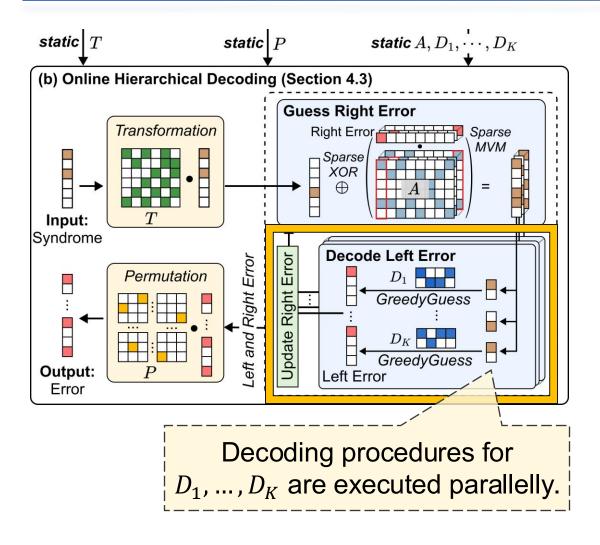
Therefore, we propose to partition the full error into *left error* and *right error* accordingly:





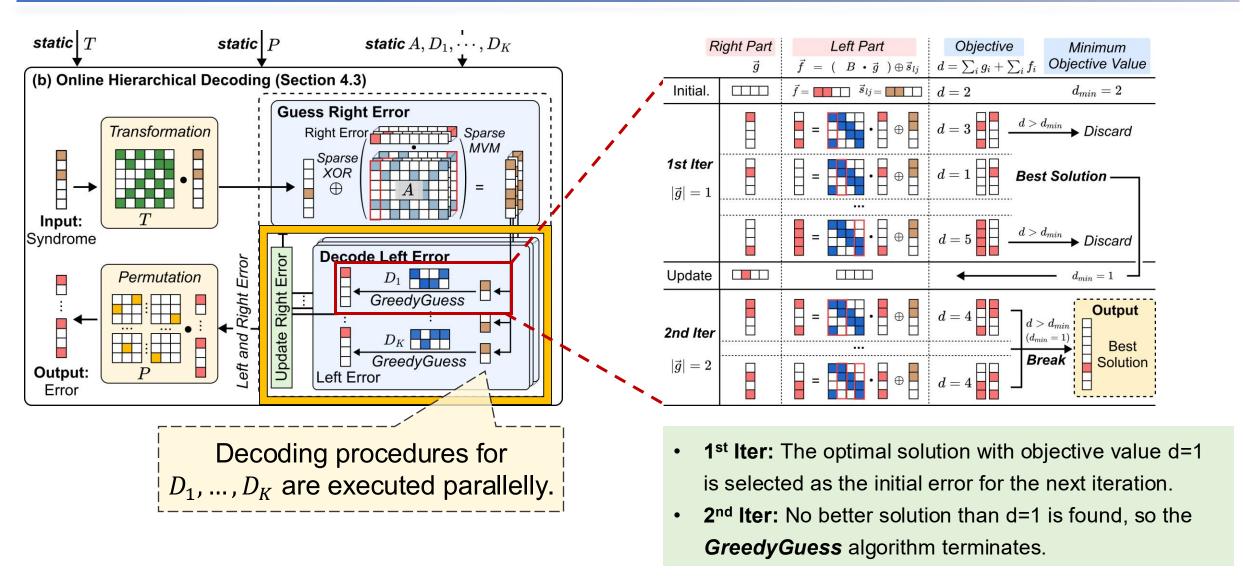
Online Hierarchical Decoding





Online Hierarchical Decoding





Outline

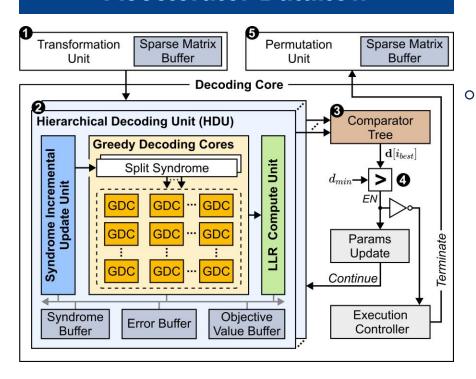


- Background
- Motivation
- Vegapunk Algorithm
- Vegapunk Accelerator
- Experiment

Accelerator Architecture



Accelerator Dataflow



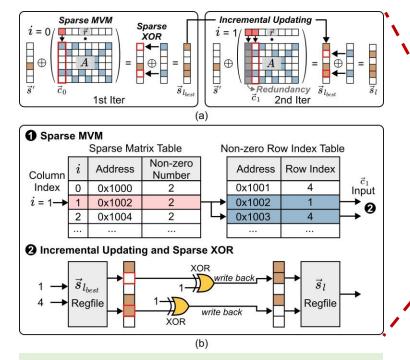
Customized design based on the online hierarchical algorithm.

- The accelerator consists of a decoding core, a transformation unit, and a permutation unit.
- The execution controller determines whether the decoding should continue.

Accelerator Architecture

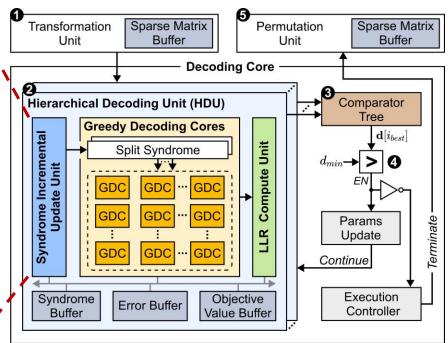


Syndrome Update Unit



- Sparse MVM and sparse XOR operations are employed.
- Incremental updating greatly reduces redundant computations during decoding.

Accelerator Dataflow



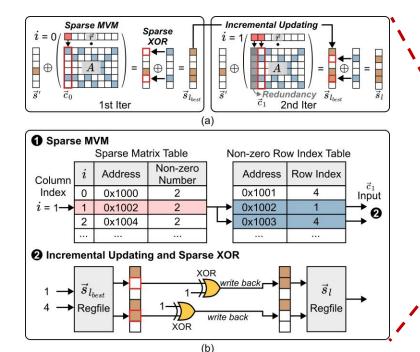
- The accelerator consists of a decoding core, a transformation unit, and a permutation unit.
- The execution controller determines whether the decoding should continue.

Customized design based on the online hierarchical algorithm.

Accelerator Architecture

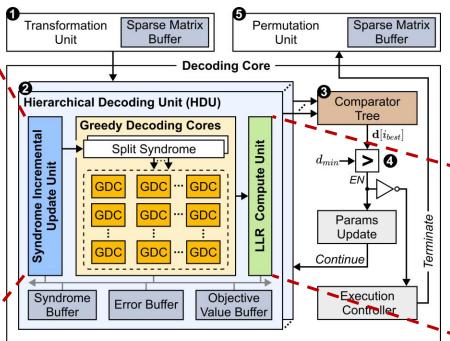


Syndrome Update Unit



- Sparse MVM and sparse XOR operations are employed.
- Incremental updating greatly reduces redundant computations during decoding.

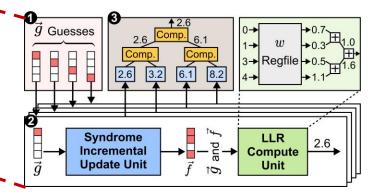
Accelerator Dataflow



- The accelerator consists of a decoding core, a transformation unit, and a permutation unit.
- The execution controller determines whether the decoding should continue.

Customized design based on the online hierarchical algorithm.

LLR Compute Unit



- The adder tree replaces floatingpoint multiplications for fast objective value computation.
- The comparison tree **efficiently identifies** the minimum value.

Outline



- Background
- Motivation
- Vegapunk Algorithm
- Vegapunk Accelerator
- Experiment

Evaluation Setup



Baselines:

• BP^[1], BP+OSD^[2] (CS(7) version), BP+LSD^[3], and BPGD^[4].

■ Benchmarks:

• (1) Bivariate Bicycle (BB) codes; (2) Hypergraph Product (HP) codes.

■ Implementation:

Vegapunk: High-Level Synthesis (HLS) C++, Xilinx Alveo U50 FPGA at 250 MHz.

■ Noise Model:

- BB code: circuit-level noise with initial errors, gate errors, measurement errors, and reset errors.
- HP codes: phenomenological noise with bit-flip errors and measurement errors.

Evaluation Metrics:

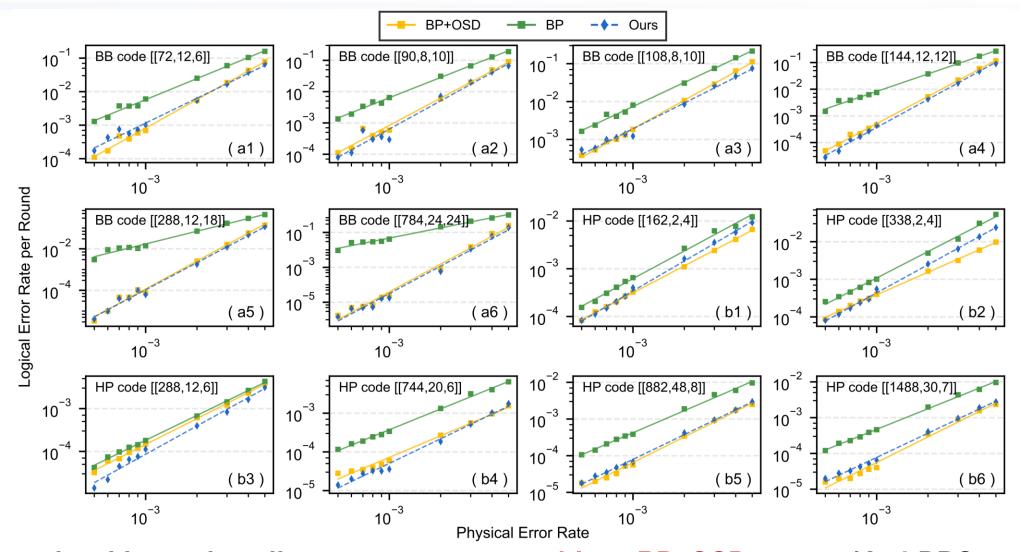
- (1) Logical Error Rate (LER) per round p_L ; (2) latency; (3) accuracy threshold p_T .
- We calculate accuracy threshold p_T by $\ln p_L = k \ln p + (1 k) \ln p_T$, where p is physical error rate and k is constant parameter. A higher accuracy threshold indicates greater tolerance to physical error rates.

[1] Pearl, Judea. "Reverend Bayes on inference engines: A distributed hierarchical approach." 2022. [3] Hillmann, et al. "Localized statistics decoding: A parallel decoding algorithm for qLDPC codes." 2024. [2] Fossorier, et al. "Soft-decision decoding of linear block codes based on ordered statistics." 2002. [4] Yao, Hanwen, et al. "Belief propagation decoding of qLDPC codes with guided decimation." 2024.



Logical Error Rate (LER)





Vegapunk achieves decoding accuracy comparable to BP+OSD across 12 qLDPC codes.



Decoupled Result and Decoding Performance



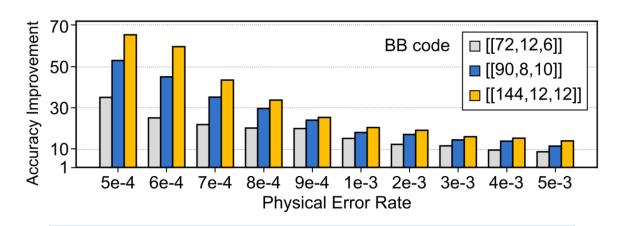


Code Type	Code Notation	Check Matrix Shape	Decoupled Matrices		Accuracy Threshold (%)			Latency per Round		
			A shape (Spars.*)	D _i shape (Spars.*)	BP	BP+OSD- CS(7)	Vegapunk	BP	BP+OSD- CS(7)	Vegapunk
BB Codes	[[72,12,6]]	[36,360]	[36,108] (6)	[6,42] (3)	0.020	0.112	0.091	694ns	0.98ms	720ns
	[[90,8,10]]	[45,450]	[45,135] (6)	[9,63] (3)	0.020	0.110	0.125	1104ns	2.02ms	732ns
	[[108,8,10]]	[54,540]	[54,162] (6)	[6,42] (3)	0.016	0.065	0.060	1400ns	3.03ms	732ns
	[[144,12,12]]	[72,720]	[72,216] (6)	[6,42] (3)	0.017	0.132	0.149	2387ns	6.77ms	732ns
	[[288,12,18]]	[144,1440]	[144,432] (6)	[12,84] (3)	0.006	0.186	0.196	7009ns	38.6ms	780ns
	[[784,24,24]]	[392,3920]	[392,1176] (6)	[28,196] (3)	0.002	0.213	0.227	25881ns	449ms	840ns
Average	-	•	-		0.013	0.136	0.141	6412ns	83.4ms	756ns
HP Codes	[[162,2,4]]	[81,243]	[81,81] (2)	[9,18] (2)	0.167	0.377	0.268	151ns	0.18ms	264ns
	[[338,2,4]]	[169,507]	[169,169] (2)	[13,26] (2)	0.094	0.261	0.175	282ns	0.43ms	276ns
	[[288,12,6]]	[144,432]	[144,144] (4)	[12,24] (4)	0.622	0.674	0.777	72ns	0.19ms	432ns
	[[744,20,6]]	[372,1116]	[372,372] (4)	[31,62] (4)	0.350	2.176	1.455	578ns	1.93ms	480ns
	[[882,48,8]]	[441,1323]	[441,441] (5)	[63,126] (3)	0.237	0.798	0.768	735ns	6.66ms	526ns
	[[1488,30,7]]	[744,2232]	[744,744] (4)	[31,62] (4)	0.228	0.819	0.775	1875ns	11.7ms	492ns
Average	-	-	-	-	0.283	0.851	0.703	616ns	3.52ms	412ns

Ablation Study



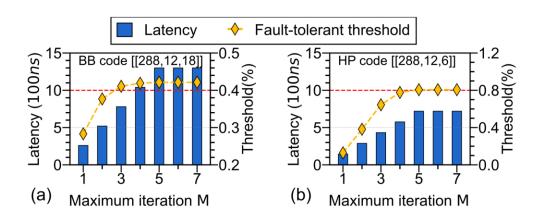
Offline Decoupling Strategy



Result: On three BB codes, the offline decoupling strategy yields accuracy improvements of 17.9×, 26.1×, and 31.1×, respectively.

Analysis: The decoupling strategy partitions the check matrix into smaller submatrices, which not only mitigates quantum degeneracy but also reduces the search space for greedy decoding.

Maximum Iteration M



Result: As M increases, both decoding latency and accuracy threshold improve, with **M=3 offering the best trade-off**.

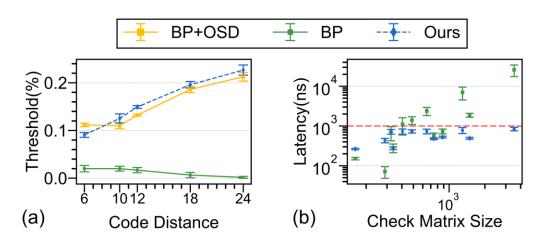
Analysis: A larger M expands the greedy decoding's search space, enhancing accuracy to some extent but also increasing latency, with diminishing returns.



Scalability and Comparison with SOTA



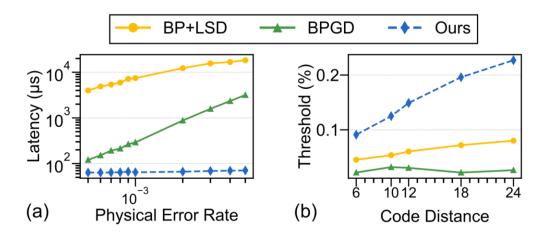
Scalability



Result: both Vegapunk and BP+OSD exhibit increasing thresholds as the code distance increases. In contrast, BP shows a decreasing trend.

Analysis: Vegapunk and BP+OSD can fully exploit the error correction capacity of qLDPC codes, while BP is unable to achieve this due to its low decoding accuracy.

Comparison with BP+LSD and BPGD



Result: Vegapunk achieves average speedups of 147.6× and 13.9×, and average accuracy threshold improvements of 2.53× and 7.11×.

Analysis: Vegapunk employs a greedy decoding strategy and effectively leverages the sparsity in decoupled check matrices. The higher accuracy stems from its offline decoupling strategy.

Conclusion



■ Problem:

Accurate and fast decoder for qLDPC codes in the era of fault-tolerant quantum computing.

■ Motivation:

Existing qLDPC decoders face two challenges: (1) low decoding accuracy and (2) high latency.

■ Solution: Vegapunk

• (1) Offline SMT-optimized decoupling; (2) online hierarchical decoding; (3) customized accelerator.

■ Results:

- Vegapunk achieves decoding accuracy on par with BP+OSD.
- Across 12 types of qLDPC codes, Vegapunk consistently achieves decoding latency below 1 µs.



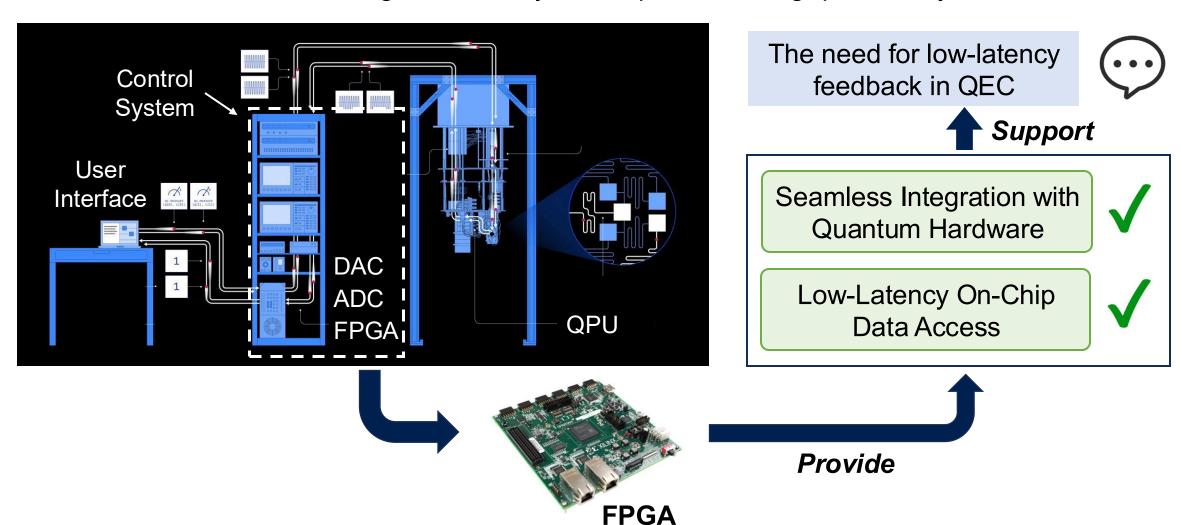
Backup Materials

2. Low Latency: Why Dedicated Accelerator?



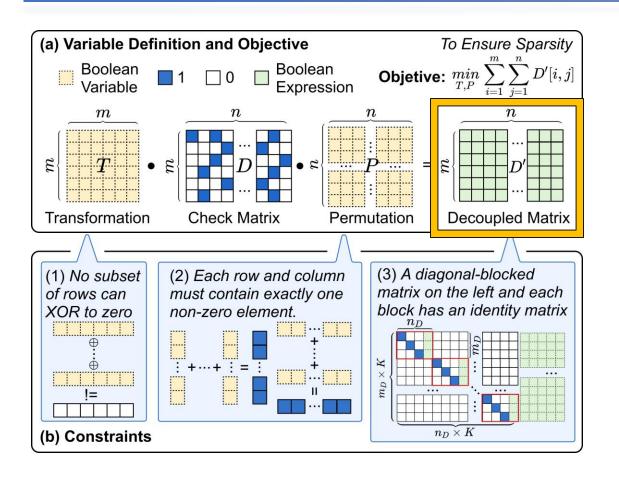


FPGA-based accelerators integrate naturally with superconducting quantum systems.



Offline SMT-Optimized Decoupling



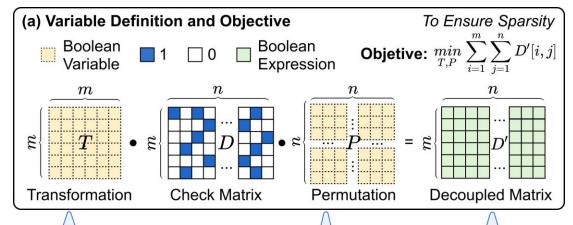


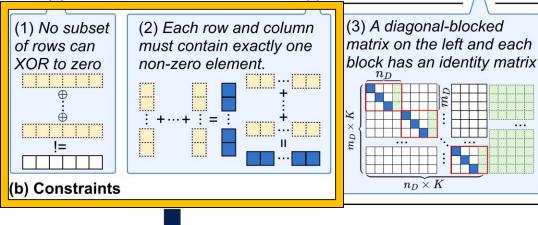
Mathematical Formulation

Variables:
$$D'[i,j] = \sum_{q=1}^{n} \sum_{k=1}^{m} T[i,k]D[k,q]P[q,j]$$

Offline SMT-Optimized Decoupling







Example

Mathematical Formulation

Variables:
$$D'[i,j] = \sum_{q=1}^{n} \sum_{k=1}^{m} T[i,k]D[k,q]P[q,j]$$

Constraints:

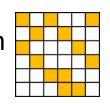
• Transformation matrix *T*:

$$\sum_{j=1}^{m} \bigoplus_{i \in \Delta} T[i, j] > 0, \ \forall \Delta \subseteq \{1, \dots, m\} / \emptyset$$

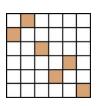
Permutation matrix P:

$$\sum_{j=1}^{n} P[i,j] = 1, \ \forall i \in \{1,\cdots,n\} \ , \ \sum_{i=1}^{n} P[i,j] = 1, \ \forall j \in \{1,\cdots,n\}$$

Transformation matrix *T*

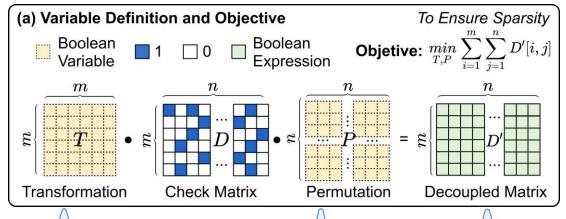


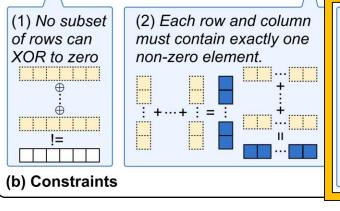
Permutation matrix *P*

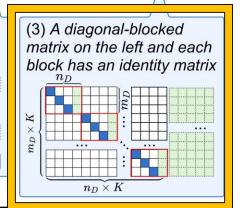


Offline SMT-Optimized Decoupling









Example

Mathematical Formulation

Constraints:

Decoupled matrix D':

$$m_D \cdot K = m, m \le n_D \cdot K \le n$$

 $D'[i \cdot m_D + t, j \cdot n_D + k] = 0, \forall i, j (i \neq j) \in \{1, \dots, K\}$ $\forall t \in \{1, \dots, m_D\}, \forall k \in \{1, \dots, n_D\}$

 $D'[i \cdot m_D + t, i \cdot n_D + t] = 1, \forall i \in \{1, \dots, K\}, \forall t \in \{1, \dots, m_D\}$ $D'[i \cdot m_D + t, i \cdot n_D + k] = 0, \forall i \in \{1, \dots, K\} \ \forall t, k(t \neq k) \in \{1, \dots, m_D\}$

