

BP + OSD Algorithm for Quantum Error Correction

Lecture Note based on arXiv:2005.07016

Roffe, White, Burton, and Campbell (2020)

Contents

1	Introduction	3
1.1	Overview	3
1.2	Learning Objectives	3
2	Classical Error Correction	4
2.1	Linear Codes	4
2.2	Hamming Weight and Distance	4
2.3	Error Model and Syndrome	5
2.4	The Decoding Problem	5
2.5	Probabilistic Graphical Models	6
2.6	Undirected Probabilistic Graphical Models	6
2.7	Factor Graphs	6
2.8	Comparing Graph Representations	7
3	Belief Propagation Decoder	9
3.1	Introduction and Motivation	9
3.2	The Message Passing Mechanism	9
3.3	BP Algorithm: Step-by-Step	10
3.4	BP Dynamics: Classical vs. Quantum Constraints	17
4	Ordered Statistics Decoding (OSD)	20
4.1	The Key Insight	20
4.2	OSD-0: The Basic Algorithm	20
4.3	Higher-Order OSD (OSD- λ)	24
4.4	Combination Sweep Strategy (OSD-CS)	27
5	The Complete BP+OSD Decoder	30
5.1	Algorithm Flow	30
5.2	Complete Algorithm: BP+OSD-CS	31
5.3	BP Convergence and Performance Guarantees	32
6	Minimum Weight Perfect Matching (MWPM) Decoder	41
6.1	Maximum Likelihood Decoding and MWPM	41
6.2	The Matching Polytope	41
6.3	Dual Formulation and Optimality Conditions	42
6.4	The Blossom Algorithm	43
7	Quantum Error Correction Basics	45
7.1	Qubits and Quantum States	45
7.2	Pauli Operators	45
7.3	Binary Representation of Pauli Errors	46
7.4	CSS Codes	46
7.5	Syndrome Measurement in CSS Codes	47
7.6	Quantum Code Parameters	47

7.7	The Hypergraph Product Construction	48
8	Results and Performance	49
8.1	Error Threshold	49
8.2	Experimental Results	49
8.3	Complexity	49
9	Summary	50
9.1	Key Takeaways	50
9.2	The BP+OSD Recipe	50
10	References	51
	Bibliography	51

1 Introduction

1.1 Overview

This lecture note introduces the **BP+OSD decoder** for quantum error correction:

- **BP** = Belief Propagation (a classical decoding algorithm)
- **OSD** = Ordered Statistics Decoding (a post-processing technique)

Together, BP+OSD provides a general-purpose decoder for **quantum LDPC codes** (Low-Density Parity Check codes).

1.2 Learning Objectives

By the end of this note, you will understand:

1. How classical error correction codes work
2. The Belief Propagation algorithm for decoding
3. Why BP fails for quantum codes (the degeneracy problem)
4. How OSD fixes the degeneracy problem
5. The complete BP+OSD decoding algorithm

2 Classical Error Correction

2.1 Linear Codes

All arithmetic in this note is performed in **binary** (modulo 2):

- $0 + 0 = 0$, $1 + 0 = 0 + 1 = 1$, $1 + 1 = 0$
- This is also written as XOR: $a \oplus b = (a + b) \bmod 2$
- Vectors and matrices use element-wise mod-2 arithmetic

2.2 Hamming Weight and Distance

Definition. The **Hamming weight** of a binary vector \mathbf{v} is the number of 1s it contains $|\mathbf{v}| = \sum_i v_i$. The **Hamming distance** between two vectors \mathbf{u} and \mathbf{v} is the number of positions where they differ: $d(\mathbf{u}, \mathbf{v}) = |\mathbf{u} + \mathbf{v}|$.

For example, for $\mathbf{v} = (1, 0, 1, 1, 0)$: Hamming weight $|\mathbf{v}| = 3$ and for $\mathbf{u} = (1, 1, 0, 1, 0)$ and $\mathbf{v} = (1, 0, 1, 1, 0)$: $\mathbf{u} + \mathbf{v} = (0, 1, 1, 0, 0)$ and $d(\mathbf{u}, \mathbf{v}) = 2$.

Definition. An $[n, k, d]$ **linear code** \mathcal{C} is a set of binary vectors (called **codewords**) where n is the **block length** (number of bits in each codeword), k is the **dimension** (number of information bits encoded), and d is the **minimum distance** (minimum Hamming weight among non-zero codewords). The **rate** of the code is $R = k/n$.

A linear code can be defined by an $m \times n$ **parity check matrix** H . H_{ij} denotes the entry in row i , column j of matrix H . m is the number of rows in H (number of parity checks), n is the number of columns in H (number of bits), and $\text{rank}(H)$ is the number of linearly independent rows. By the rank-nullity theorem: $k = n - \text{rank}(H)$.

For example, The **[3, 1, 3] repetition code** encodes 1 bit into 3 bits by triplication.

Parity check matrix:

$$H = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix} \tag{1}$$

$$\text{Verification: } H \cdot \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \checkmark \text{ and } H \cdot \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \checkmark$$

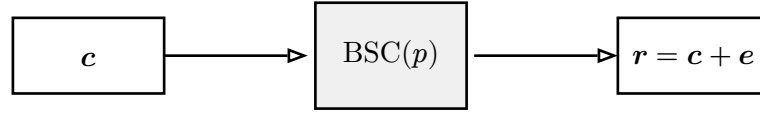
So the codewords are: $\mathcal{C}_H = \{(0, 0, 0), (1, 1, 1)\}$

Parameters: $n = 3$ bits, $k = 3 - 2 = 1$ info bit, $d = 3$ (weight of $(1, 1, 1)$)

2.3 Error Model and Syndrome

Definition. In the **binary symmetric channel** (BSC) with error probability p :

- Each bit is independently flipped with probability p
- Original codeword: \mathbf{c}
- Error pattern: \mathbf{e} (a binary vector, $e_i = 1$ means bit i flipped)
- Received word: $\mathbf{r} = \mathbf{c} + \mathbf{e}$



\mathbf{e} : random error

Figure 1: Binary symmetric channel model

Definition. The **syndrome** of a received word \mathbf{r} is:

$$\mathbf{s} = H \cdot \mathbf{r} \quad (2)$$

Since $H \cdot \mathbf{c} = \mathbf{0}$ for any codeword, we have:

$$\mathbf{s} = H \cdot \mathbf{r} = H \cdot (\mathbf{c} + \mathbf{e}) = H \cdot \mathbf{c} + H \cdot \mathbf{e} = \mathbf{0} + H \cdot \mathbf{e} = H \cdot \mathbf{e} \quad (3)$$

Key Point. The syndrome depends **only on the error**, not on which codeword was sent! This is what makes syndrome-based decoding possible.

2.4 The Decoding Problem

Given: Parity check matrix H and syndrome $\mathbf{s} = H \cdot \mathbf{e}$

Find: The most likely error \mathbf{e}^* that could have produced \mathbf{s}

Definition. Maximum likelihood decoding finds:

$$\mathbf{e}^* = \arg \min_{\mathbf{e}: H \cdot \mathbf{e} = \mathbf{s}} |\mathbf{e}| \quad (4)$$

That is, the minimum Hamming weight error consistent with the syndrome.

2.5 Probabilistic Graphical Models

Before introducing the Belief Propagation algorithm, we need to understand how probabilistic inference problems can be represented as graphs.

Definition. A **probabilistic graphical model** (PGM) is a graph-based representation of a probability distribution. Nodes represent random variables, and edges encode conditional independence relationships. PGMs enable efficient inference algorithms by exploiting the structure of the distribution.

There are two main families of PGMs:

- **Directed graphical models** (Bayesian networks): edges have direction, representing causal relationships
- **Undirected graphical models** (Markov networks): edges are undirected, representing symmetric dependencies

For error correction, we use undirected models because parity constraints are symmetric — no variable “causes” another.

2.6 Undirected Probabilistic Graphical Models

Definition. An **undirected probabilistic graphical model** (also called a **Markov network** or **Markov random field**) represents a joint probability distribution as:

$$P(\mathbf{x}) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{x}_c) \quad (5)$$

where:

- $\mathbf{x} = (x_1, \dots, x_n)$ are random variables
- \mathcal{C} is a set of **cliques** (fully connected subgraphs)
- $\psi_c(\mathbf{x}_c)$ is a **potential function** over variables in clique c
- $Z = \sum_{\mathbf{x}} \prod_c \psi_c(\mathbf{x}_c)$ is the **partition function** (normalization constant)

Key Point. The UAI format mentioned in the getting started guide represents exactly this structure: variables (detectors), cliques (error mechanisms), and potential functions (error probabilities).

For binary error correction with syndrome \mathbf{s} , we want to compute:

$$P(\mathbf{e} \mid \mathbf{s}) \propto \prod_c \psi_c(\mathbf{e}_c) \quad (6)$$

where each potential ψ_c encodes a parity constraint.

2.7 Factor Graphs

To understand the Belief Propagation algorithm, we need the concept of **factor graphs**.

Definition. A **factor graph** is a bipartite graph $G = (V, U, E)$ representing the parity check matrix H . The **data nodes** are set $V = \{v_1, v_2, \dots, v_n\}$ such that each node v_j corresponds to each column of H . A **parity nodes** are set $U = \{u_1, u_2, \dots, u_m\}$ such that

each node u_i corresponds to each row of H . An **edges** $E = \{(v_j, u_i) : H_{ij} = 1\}$ connects v_j to u_i exists if $H_{ij} = 1$.

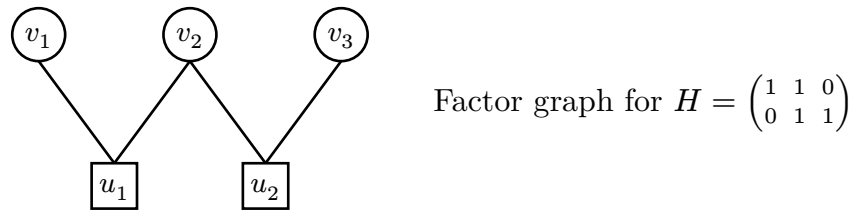


Figure 2: Factor graph for the $[3, 1, 3]$ repetition code with node conventions

The **neighborhoods** of nodes are defined as: $V(u_i) = \{v_j : H_{ij} = 1\}$.

2.8 Comparing Graph Representations

Three related graph representations appear in coding theory and probabilistic inference:

Definition. Comparison of graph representations:

1. Undirected Probabilistic Graphical Model (Markov Network):

- Nodes = random variables
- Edges = direct dependencies between variables
- Cliques = groups of mutually dependent variables
- Represents: $P(\mathbf{x}) = \frac{1}{Z} \prod_c \psi_c(\mathbf{x}_c)$

2. Factor Graph:

- Two types of nodes: variable nodes AND factor nodes
- Bipartite structure: edges only between variables and factors
- Explicitly represents factorization of the distribution
- Represents: $P(\mathbf{x}) = \frac{1}{Z} \prod_i f_i(\mathbf{x}_{N(i)})$

3. Tanner Graph:

- A special case of factor graph for error correction codes
- Variable nodes = bits (columns of H)
- Factor nodes = parity checks (rows of H)
- Represents: parity check matrix H structure

Key Point. Key relationships:

- Factor graphs are a **bipartite refinement** of Markov networks that make the factorization explicit
- Tanner graphs are factor graphs **specialized for linear codes** where factors represent parity constraints
- All three represent the same probability distribution, but factor graphs enable more efficient message-passing algorithms
- The UAI format represents Markov networks (cliques and potentials), while BP operates on the factor graph representation

Property	Markov Network	Factor Graph	Tanner Graph
Node types	Variables only	Variables + Factors	Bits + Checks
Graph structure	General	Bipartite	Bipartite
Edges represent	Dependencies	Factor membership	Parity constraints
Used for	General inference	Message passing	Code decoding
BP efficiency	Less efficient	Efficient	Efficient

Table 1: Comparison of graph representations

Why use factor graphs for BP? The bipartite structure of factor graphs makes message passing natural:

- Messages flow between variables and factors
- Each factor collects evidence from its variables
- Each variable aggregates information from its factors
- No need to handle complex clique structures

For error correction, the Tanner graph (factor graph) representation is ideal because:

- Parity checks are naturally factors (XOR constraints)
- Bits are naturally variables (error indicators)
- The sparse structure (H has few 1s) gives efficient $O(n)$ BP iterations

Definition. An (l, q) -LDPC code is a linear code whose parity check matrix H satisfies:

- Each column has at most l ones (each bit is in at most l checks)
- Each row has at most q ones (each check involves at most q bits)

The matrix H is called **sparse** because l and q are small constants independent of n .

Key Point. LDPC codes are important because their sparse structure enables efficient decoding via Belief Propagation with complexity $O(n)$ per iteration.

3 Belief Propagation Decoder

3.1 Introduction and Motivation

The rediscovery of Low-Density Parity-Check (LDPC) codes in the late 1990s marked a paradigm shift in coding theory, transitioning from algebraic decoding algorithms to probabilistic iterative decoding that approaches the Shannon limit [1]. Central to this revolution is the **Belief Propagation** (BP) algorithm [2], a message-passing protocol that operates on the graphical representation of codes.

Key Point. BP in Modern Communications: BP decoding powers critical communication standards:

- Wi-Fi (IEEE 802.11n/ac/ax)
- Satellite communication (DVB-S2)
- 5G New Radio

While its practical efficacy is undisputed, the mathematical rigor underlying its convergence behavior involves multiple theoretical frameworks: Density Evolution for asymptotic analysis, Bethe Free Energy for variational optimization, and trapping set theory for failure mechanisms.

The convergence of BP is understood through different lenses depending on the regime. In the asymptotic limit of infinite block length, convergence is probabilistic and governed by Density Evolution [3]. In finite-length regimes, convergence is variational, linked to minimization of the Bethe Free Energy [4]. However, combinatorial substructures known as trapping sets can arrest decoding, creating error floors [5].

3.2 The Message Passing Mechanism

Belief Propagation (BP), also called the **sum-product algorithm**, is an iterative message-passing algorithm on the factor graph.

Definition. The goal of BP is to compute, for each bit j , the **marginal probability**:

$$P_1(e_j) = P(e_j = 1 \mid \mathbf{s}), \quad (7)$$

given $\mathbf{s} = H \cdot \mathbf{e}$ is the syndrome of the error \mathbf{e} . This is called a **soft decision** – it tells us how likely each bit is to be flipped.

We use the following notation throughout:

- p as the channel error probability (probability each bit flips)
- $m_{v_j \rightarrow u_i}$ as the message from data node v_j to parity node u_i
- $m_{u_i \rightarrow v_j}$ as the message from parity node u_i to data node v_j
- Messages represent beliefs about whether $e_j = 1$

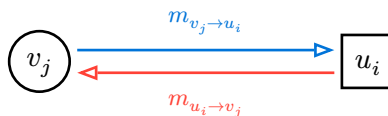


Figure 3: Messages passed between data and parity nodes

The BP-algorithm requires the quantification of how certain we are about the bit being flipped or not. This is done using **log-likelihood ratios** (LLR).

Definition. Instead of probabilities, BP uses **log-likelihood ratios** (LLR) for numerical stability:

$$\text{LLR}(e_j) = \log \frac{P(e_j = 0)}{P(e_j = 1)} \quad (8)$$

- $\text{LLR} > 0$ means $e_j = 0$ is more likely (bit probably correct)
- $\text{LLR} < 0$ means $e_j = 1$ is more likely (bit probably flipped)
- $|\text{LLR}|$ indicates confidence level

For the channel with error probability p , the **channel LLR** is:

$$p_l = \log \frac{1-p}{p} \quad (9)$$

Since $p < 0.5$ in practice, we have $p_l > 0$.

3.3 BP Algorithm: Step-by-Step

Step 1: Initialization

Set the channel LLR:

$$p_l = \log \frac{1-p}{p} \quad (10)$$

Initialize all messages from data nodes to parity nodes with the channel prior:

$$m_{v_j \rightarrow u_i} := p_l \quad \text{for all edges } (v_j, u_i) \quad (11)$$

Why? Before any message passing, the only information we have about each bit is from the **channel itself**. Since each bit flips independently with probability p , the initial belief is simply the channel's prior: “this bit is probably correct” (because $p < 0.5$, so $p_l > 0$).

Code: The implementation adds small epsilon values for numerical stability:

```
# Initialize channel LLRs
channel_llr = torch.log((1 - channel_probs + 1e-10) / (channel_probs + 1e-10))

# Initialize qubit-to-check messages with channel prior
msg_q2c = channel_llr[qubit_edges].unsqueeze(0).expand(batch_size, -1).clone()
msg_c2q = torch.zeros(batch_size, num_edges, device=device)
```

Step 2: Parity-to-Data Messages

Each parity node u_i sends a message to each connected data node v_j :

- **Min-sum form:**

$$m_{u_i \rightarrow v_j} = (-1)^{s_i} \cdot \alpha \cdot \prod_{v'_j \in V(u_i) \setminus v_j} \text{sign}(m_{v'_j \rightarrow u_i}) \cdot \min_{v'_j \in V(u_i) \setminus v_j} |m_{v'_j \rightarrow u_i}| \quad (12)$$

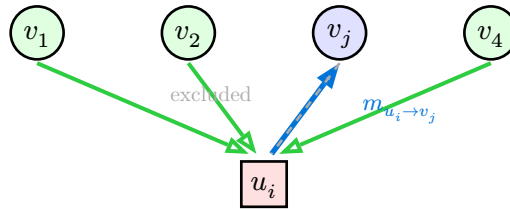
- **Sum-Product form:**

$$m_{u_i \rightarrow v_j} = (-1)^{s_i} \cdot 2 \tanh^{-1} \left(\prod_{v'_j \in V(u_i) \setminus v_j} \tanh \left(\frac{m_{v'_j \rightarrow u_i}}{2} \right) \right) \quad (13)$$

Where:

- s_i = the i -th syndrome bit (given as input, either 0 or 1)
- $V(u_i) \setminus v_j$ = all neighbors of u_i except v_j (defined in Section 3.5)
- $\text{sign}(x) = +1$ if $x \geq 0$, else -1
- $\alpha = 1 - 2^{-t}$ is a **damping factor** at iteration t (helps convergence)

Why? A parity check enforces that XOR of all connected bits equals the syndrome bit s_i . The check node tells v_j : “Based on what I know about the **other** bits, here’s how likely **you** are to be flipped.” If $s_i = 0$, the parity check says “even number of flipped bits.” If the other bits all look correct (positive LLR), then v_j should also be correct. If one other bit looks flipped (negative LLR), then v_j should be correct to maintain even parity. The formula computes this XOR-like logic in LLR form.



Check node collects info from v_1, v_2, v_4
to compute message to v_j (excluding v_j 's own message)

Figure 4: Parity-to-data message: u_i uses info from all neighbors **except** v_j to tell v_j what it should be

Code: The implementation computes signs and magnitudes separately, using a sorting trick to efficiently find the minimum excluding each edge:

```
def _check_to_qubit_minsum(self, msg_q2c, syndromes):
    for c in range(num_checks):
        edges = self.check_to_edges[c]
        incoming = msg_q2c[:, edges] # (batch, degree)

        # Separate signs and magnitudes
        signs = torch.sign(incoming) # (batch, degree)
        mags = torch.abs(incoming) # (batch, degree)

        # Product of all signs
        total_sign = torch.prod(signs, dim=1, keepdim=True)
```

```

# Apply syndrome: flip sign if syndrome is 1
syndrome_sign = 1 - 2 * syndromes[:, c:c+1]
total_sign = total_sign * syndrome_sign

# For each edge, divide out its sign to get product of others
outgoing_signs = total_sign / (signs + 1e-10)

# Min magnitude excluding each edge (second minimum trick)
sorted_mags, _ = torch.sort(mags, dim=1)
min_mag = sorted_mags[:, 0:1]
second_min = sorted_mags[:, 1:2] if sorted_mags.shape[1] > 1 else min_mag

# If edge has the min, use second_min; else use min
is_min = (mags == min_mag)
outgoing_mags = torch.where(is_min, second_min, min_mag)

# Apply scaling factor
scaling = 0.625
msg_c2q[:, edges] = scaling * outgoing_signs * outgoing_mags

```

Key Point. Second Minimum Trick: To compute $\min_{v'_j \neq v_j} |m_{v'_j}|$ efficiently, we sort all magnitudes once. For each edge: if it holds the minimum, use the second minimum; otherwise use the first minimum. This avoids recomputing $O(d)$ minimums for each of d edges.

Sum-Product Code: The sum-product variant uses the tanh identity for exact computation:

```

def _check_to_qubit_sumproduct(self, msg_q2c, syndromes):
    for c in range(num_checks):
        edges = self.check_to_edges[c]
        incoming = msg_q2c[:, edges] # (batch, degree)

        # Compute tanh(LLR/2) - clamp for numerical stability
        half_llr = torch.clamp(incoming / 2, min=-20, max=20)
        tanh_vals = torch.tanh(half_llr) # (batch, degree)

        # Product of all tanh values
        total_prod = torch.prod(tanh_vals, dim=1, keepdim=True)

        # Apply syndrome: flip sign if syndrome is 1
        syndrome_sign = 1 - 2 * syndromes[:, c:c+1]
        total_prod = total_prod * syndrome_sign

        # For each edge, divide out its tanh contribution
        outgoing_prod = total_prod / (tanh_vals + 1e-10)

        # Clamp to valid range for atanh (-1, 1)
        outgoing_prod = torch.clamp(outgoing_prod, min=-1+1e-7, max=1-1e-7)

        # Convert back: 2 * atanh(prod)
        msg_c2q[:, edges] = 2 * torch.atanh(outgoing_prod)

```

Key Point. Min-Sum vs Sum-Product: Min-sum is an **approximation** of sum-product, not an equivalent algorithm. The relationship comes from the identity:

$$2 \tanh^{-1} \left(\prod_i \tanh(x_i/2) \right) \approx \text{sign} \left(\prod_i x_i \right) \cdot \min_i |x_i| \quad (14)$$

This approximation holds because $\tanh(x/2) \approx \text{sign}(x)$ for large $|x|$, and the product of \tanh values is dominated by the smallest magnitude input. The scaling factor $\alpha = 0.625$ in min-sum compensates for the systematic overestimation of confidence that results from this approximation [6].

Algorithm	Formula	Trade-off
Sum-Product	$2 \tanh^{-1}(\prod \tanh(x/2))$	Exact but slower (\tanh/atanh)
Min-Sum	$\alpha \cdot \text{sign}(\prod x) \cdot \min x $	Approximate but faster

Table 2: Comparison of check-to-qubit message algorithms

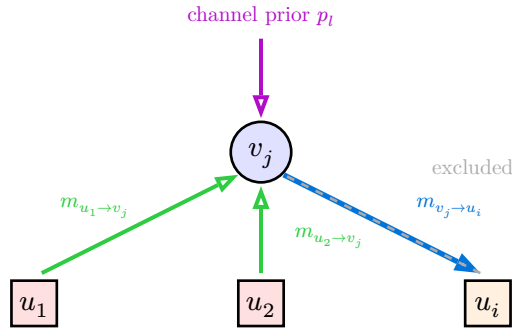
Step 3: Data-to-Parity Messages

Each data node v_j sends a message to each connected parity node u_i :

$$m_{v_j \rightarrow u_i} = p_l + \sum_{u'_i \in U(v_j) \setminus u_i} m_{u'_i \rightarrow v_j} \quad (15)$$

Where $U(v_j) \setminus u_i =$ all parity neighbors of v_j except u_i .

Why? A data node collects evidence from multiple parity checks. Each check provides independent information about whether this bit is flipped. The data node sums up all this evidence (in LLR, multiplication of probabilities becomes addition).



Data node sums: channel prior + messages from u_1, u_2
to send to u_i (excluding u_i 's own message)

Figure 5: Data-to-parity message: v_j combines channel prior with info from other checks

Intuition: Why exclude u_i ? To avoid **echo effects**. If v_j included the message it previously received from u_i , that information would bounce back, creating a feedback loop. On a **tree-structured graph**, this exclusion ensures each piece of evidence is counted exactly once, making BP exact. On graphs with cycles, this is an approximation.

Code: The implementation efficiently computes this by computing the total sum and subtracting each edge's contribution:

```
def _qubit_to_check(self, msg_c2q, channel_llr):
    for q in range(num_qubits):
        edges = self.qubit_to_edges[q]
        incoming = msg_c2q[:, edges] # (batch, degree)

        # Sum of all incoming messages plus channel prior
        total_sum = incoming.sum(dim=1, keepdim=True) + channel_llr[q]

        # For each edge, subtract its own contribution
        msg_q2c[:, edges] = total_sum - incoming
```

Key Point. Why subtract? Instead of computing $n - 1$ sums for each of n edges (expensive), we compute one total sum and subtract each term. This reduces complexity from $O(d^2)$ to $O(d)$ per qubit.

Step 4: Compute Soft Decisions

For each bit j , compute the total belief (sum of all evidence):

$$P_1(e_j) = p_l + \sum_{u_i \in U(v_j)} m_{u_i \rightarrow v_j} \quad (16)$$

Why? Unlike Step 3, here we include **all** incoming messages (no exclusion). This is the final belief about bit j , combining the channel prior with evidence from **every** connected parity check. The result is the log-posterior probability ratio.

Code: The implementation sums all incoming messages (no exclusion) and converts to probability:

```
def _compute_marginals(self, msg_c2q, channel_llr):
    for q in range(num_qubits):
        edges = self.qubit_to_edges[q]

        # Total LLR = channel + sum of ALL incoming
        total_llr = channel_llr[q] + msg_c2q[:, edges].sum(dim=1)

        # Convert LLR to probability: P(1) = sigmoid(-LLR)
        marginals[:, q] = torch.sigmoid(-total_llr)
```

Key Point. Sigmoid conversion: The relationship between LLR and probability is:

$$\text{LLR} = \log \frac{P(0)}{P(1)} \Rightarrow P(1) = \frac{1}{1 + e^{\text{LLR}}} = \sigma(-\text{LLR}) \quad (17)$$

PyTorch's `torch.sigmoid(-total_llr)` computes this efficiently.

Step 5: Make Hard Decisions

Convert soft decisions to a binary estimate:

$$e_j^{\text{BP}} = \begin{cases} 1 & \text{if } P_1(e_j) < 0 \quad (\text{more likely flipped}) \\ 0 & \text{otherwise} \quad (\text{more likely correct}) \end{cases} \quad (18)$$

This gives us the BP estimate $e^{\text{BP}} = (e_1^{\text{BP}}, e_2^{\text{BP}}, \dots, e_n^{\text{BP}})$.

Why? The sign of LLR directly tells us the most likely value: $P_1 > 0$ means $P(e_j = 0) > P(e_j = 1)$, so the bit is probably correct. $P_1 < 0$ means the bit is probably flipped.

Step 6: Check Convergence

Verify if the estimate satisfies the syndrome equation:

$$H \cdot e^{\text{BP}} = s \quad ? \quad (19)$$

- **If yes:** BP has **converged**. Return e^{BP} and soft decisions P_1 .
- **If no:** Go back to Step 2 and repeat.
- **If max iterations reached:** BP has **failed to converge**.

Why iterate? On graphs with cycles, a single pass doesn't propagate information globally. Each iteration allows beliefs to travel further through the graph. Eventually, if the error is correctable, the hard decisions will satisfy all parity checks.

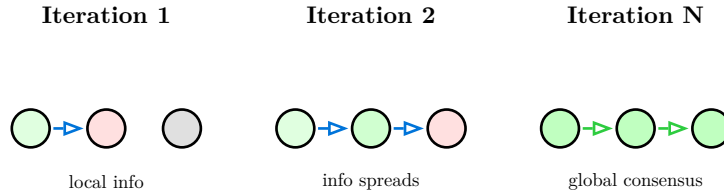


Figure 6: Information propagates further with each iteration until convergence

3.3.1 Damping for Convergence

Theory: Damping prevents oscillation by mixing old and new messages:

$$m^{(t+1)} = \gamma \cdot m^{(t)} + (1 - \gamma) \cdot m_{\text{new}} \quad (20)$$

where $\gamma \in [0, 1)$ is the damping factor.

Code: Applied after computing new check-to-qubit messages:

```
# Main iteration loop
for _ in range(max_iter):
    # Compute new check-to-qubit messages
    msg_c2q_new = self._check_to_qubit_minsum(msg_q2c, syndromes)

    # Damping: blend old and new messages
    msg_c2q = damping * msg_c2q + (1 - damping) * msg_c2q_new

    # Update qubit-to-check messages
    msg_q2c = self._qubit_to_check(msg_c2q, channel_llr)
```

Damping Value	Behavior	Use Case
$\gamma = 0$	No damping (full update)	Simple graphs, fast convergence
$\gamma = 0.2$	Light damping	Typical default
$\gamma = 0.5$	Strong damping	Graphs with short cycles

Table 3: Effect of damping factor on BP convergence

3.3.2 Summary: Complete BP Iteration

Putting it all together, one BP iteration consists of:



$$m_{u \rightarrow v} = \text{minsum/sumproduct} \quad m' = \gamma m + (1 - \gamma)m_{\text{new}} \quad m_{v \rightarrow u} = p_l + \sum m - m_{\text{in}}$$

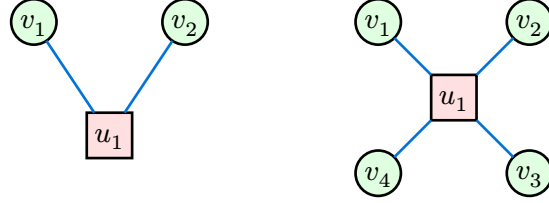
Figure 7: One BP iteration: message updates with damping

3.4 BP Dynamics: Classical vs. Quantum Constraints

To understand why BP fails on quantum codes, we compare its behavior on a minimal classical graph (2 bits) versus a minimal quantum stabilizer (4 bits).

3.4.1 Problem Setup: 2-Bit vs. 4-Bit Parity

Consider two scenarios with channel error probability $p = 0.1$ and observed syndrome $s = 1$ (odd parity): (i). A check node u_1 connected to 2 variables v_1, v_2 (e.g., a simple parity check). (ii). A check node u_1 connected to 4 variables v_1, \dots, v_4 (e.g., a surface code plaquette).



Case A: 2-Bit Check

Case B: 4-Bit Plaquette

Figure 8: Comparison of message passing geometry

Both cases start with the same channel LLR:

$$\text{LLR}_{\text{channel}} = \ln\left(\frac{1-p}{p}\right) = \ln(9) \approx 2.197. \quad (21)$$

The check node sends a message to v_1 based on evidence from **all other** neighbors. Formula: $m_{u \rightarrow v} = (-1)^s \cdot 2 \tanh^{-1}\left(\prod_{v' \in N(u) \setminus v} \tanh(m_{v' \rightarrow u}/2)\right)$

Case A (2-Bit): Strong Correction v_1 has only 1 neighbor (v_2). The product contains a single term $\tanh(1.1) \approx 0.8$.

$$m_{u_1 \rightarrow v_1} = -2 \tanh^{-1}(0.8) \approx -2.197 \quad (22)$$

The check says: “My other neighbor is definitely correct, so **you** must be wrong.” The message magnitude **matches** the channel prior but flips the sign.

Case B (4-Bit): Signal Dilution v_1 has 3 neighbors (v_2, v_3, v_4). The product accumulates uncertainty: $0.8^3 \approx 0.512$.

$$m_{u_1 \rightarrow v_1} = -2 \tanh^{-1}(0.512) \approx -1.13 \quad (23)$$

The check says: “I see an error, but it could be any of you 4. I suspect you, but weakly.” The penalty (-1.13) is **weaker** than the channel prior (+2.197).

3.4.2 Step 3: Variable Update and Failure

We now compute the updated belief for v_1 :

Metric	Case A (2-Bit)	Case B (4-Bit)
Prior LLR	+2.197	+2.197
Check Msg	-2.197	-1.13
Net LLR	0	+1.067
Prob($e_1 = 1$)	50%	$\approx 26\%$
Hard Decision	$e_1 = 0$ or 1	$e_1 = 0$

Why Quantum BP Fails:

1. **Case A (Classical):** The LLR drops to 0. BP correctly identifies “maximum uncertainty.” It knows it cannot distinguish between e_1 and e_2 .
2. **Case B (Quantum):** The LLR remains positive (+1.067). BP remains “confident” that the bit is correct.
 - The hard decision outputs $e = (0, 0, 0, 0)$.
 - **Parity Check:** $0 + 0 + 0 + 0 = 0 \neq 1$.
 - **Result:** BP fails to find a valid solution because the “blame” is diluted across too many qubits.

In a full surface code, this effect compounds. Multiple conflicting checks reduce the LLR toward 0, leading to a state of “Split Belief” where the decoder is paralyzed by symmetry.

Key Point. Summary of the Failure Mode:

- **Ambiguity Dilution:** In high-weight stabilizers (like 4-bit plaquettes), the check node message is too weak to overturn the channel prior.
- **Invalid Hard Decisions:** Unlike the 2-bit case where uncertainty is explicit ($L = 0$), the 4-bit case results in false confidence ($L > 0$) that violates the parity constraint.

3.4.3 The Degeneracy Problem

The Degeneracy Problem In quantum codes, **multiple distinct error patterns** can be physically equivalent (differing only by a stabilizer). BP, which is designed to find a **unique** correct error, fails to distinguish between these equivalent patterns, leading to decoding failure.

Definition. Two errors e_1 and e_2 are **degenerate** if they produce the same syndrome s and their difference forms a stabilizer:

$$H \cdot e_1 = H \cdot e_2 = s \quad \text{and} \quad e_1 + e_2 \in \text{Stabilizers} \quad (24)$$

Unlike classical codes where distinct low-weight errors are rare, quantum stabilizer codes are **constructed** from local degeneracies (the stabilizers themselves).

For example, consider the Toric code where qubits reside on the edges of a lattice. A stabilizer B_p acts on the 4 qubits surrounding a plaquette. The weight-2 error occurs on the **left** edges is degenerate with the weight-2 error on the **right** two edges.

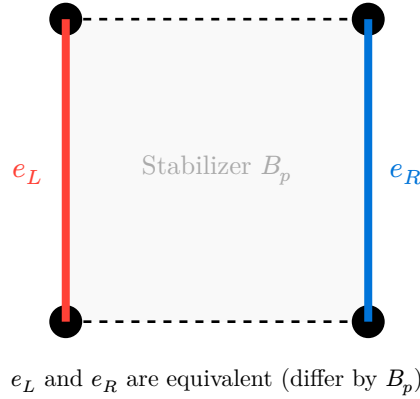


Figure 9: Symmetric degeneracy in a Toric code plaquette

When BP runs on this symmetric structure, it encounters a fatal ambiguity.

1. **Perfect Symmetry:** The graph structure for e_L is identical to the graph structure for e_R .
2. **Message Stalling:** BP receives equal evidence for the left-side error and the right-side error.
3. **Marginal Probability 0.5:** The algorithm converges to a state where every qubit in the loop has $P(e_i) \approx 0.5$.

Definition. The Hard Decision Failure: When $P(e_i) \approx 0.5$, the log-likelihood ratio is ≈ 0 .

- If we threshold **below** 0.5: The decoder outputs $e = \mathbf{0}$ (no error). This fails to satisfy the syndrome.
- If we threshold **above** 0.5: The decoder outputs $e = e_L + e_R$. This is the stabilizer itself (a closed loop), which effectively applies a logical identity but fails to correct the actual open string error.

This is why BP alone typically exhibits **zero threshold** on the Toric code: as the code size increases, the number of these symmetric loops increases, and BP gets confused by all of them simultaneously.

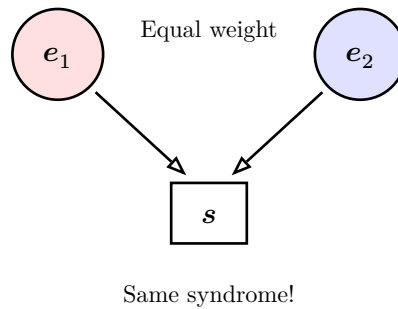


Figure 10: Two errors with the same syndrome cause BP to fail

When BP encounters degenerate errors of equal weight:

1. BP assigns high probability to **both** solutions e_1 and e_2
2. The beliefs “split” between the two solutions
3. BP outputs $e^{\text{BP}} \approx e_1 + e_2$
4. Check: $H \cdot e^{\text{BP}} = H \cdot (e_1 + e_2) = s + s = \mathbf{0} \neq s$

5. BP fails to converge!

Key Point. For the Toric code, degeneracy is so prevalent that **BP alone shows no threshold** — increasing code distance makes performance worse, not better!

4 Ordered Statistics Decoding (OSD)

4.1 The Key Insight

The parity check matrix H (size $m \times n$ with $n > m$) has more columns than rows and cannot be directly inverted. However, we can select a subset of $r = \text{rank}(H)$ linearly independent columns to form an invertible $m \times r$ submatrix.

Definition. For an $m \times n$ matrix H with $\text{rank}(H) = r$:

- **Basis set** $[S]$: indices of r linearly independent columns
- **Remainder set** $[T]$: indices of the remaining $k' = n - r$ columns
- $H_{[S]}$: the $m \times r$ submatrix of columns in $[S]$ (this is invertible!)
- $H_{[T]}$: the $m \times k'$ submatrix of columns in $[T]$

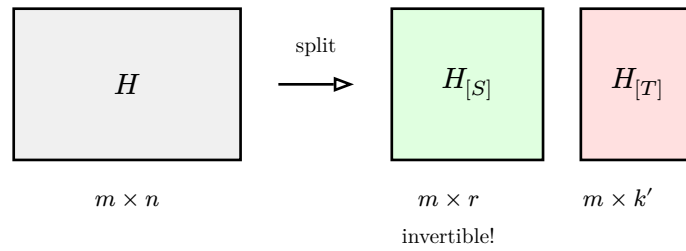


Figure 11: Splitting H into basis and remainder parts

OSD then resolves split beliefs and **forcing a unique solution**. It first choose the basis selection $[S]$ through BP soft decisions guide. and then calculate the matrix inversion on $H_{[S]}$ eliminates ambiguity.

4.2 OSD-0: The Basic Algorithm

Ordered Statistics Decoding (OSD) was introduced by Fossorier and Lin as a soft-decision decoding algorithm for linear block codes that approaches maximum-likelihood performance [7]. The algorithm was later extended with computationally efficient variants [8]. For quantum LDPC codes, the BP+OSD combination was shown to be remarkably effective by Panteleev and Kalachev, and further developed by Roffe et al. [9].

Definition. **OSD-0** finds a solution by:

1. Choosing a “good” basis $[S]$ using BP soft decisions P_1
2. Solving for the basis bits via matrix inversion
3. Setting all remainder bits to zero

Algorithm 2: OSD-0

Input:

- Parity matrix H (size $m \times n$, rank r)
- Syndrome \mathbf{s}
- BP soft decisions $P_1(e_1), \dots, P_1(e_n)$ (from Algorithm 1)

Steps:

1. **Rank bits by probability:** Sort bit indices by P_1 values: most-likely-flipped first.
Result: ordered list $[O_{\text{BP}}] = (j_1, j_2, \dots, j_n)$
2. **Reorder columns:** $H_{[O_{\text{BP}}]}$ = matrix H with columns reordered by $[O_{\text{BP}}]$
3. **Select basis:** Scan left-to-right, select first r linearly independent columns. Basis indices: $[S]$. Remainder indices: $[T]$ (size $k' = n - r$)
4. **Solve on basis:** $\mathbf{e}_{[S]} = H_{[S]}^{-1} \cdot \mathbf{s}$
5. **Set remainder to zero:** $\mathbf{e}_{[T]} = \mathbf{0}$ (zero vector of length k')
6. **Remap to original ordering:** Combine $(\mathbf{e}_{[S]}, \mathbf{e}_{[T]})$ and undo the permutation

Output: $\mathbf{e}^{\text{OSD-0}}$ satisfying $H \cdot \mathbf{e}^{\text{OSD-0}} = \mathbf{s}$

Figure 12: OSD-0 algorithm

The GPU-accelerated implementation in `batch_osd.py` realizes OSD-0 through Gaussian elimination rather than explicit matrix inversion. Below is the mapping between the theoretical definition and the code logic.

4.2.1 Step 1: Choosing a “Good” Basis (Sorting)

The definition requires selecting basis columns $[S]$ based on BP soft decisions. For quantum error correction, the implementation sorts by **probability descending** rather than reliability $|p - 0.5|$.

```
# Sort by probability descending (highest probability first)
# High-probability errors become pivots; low-probability errors become free variables
sorted_indices = np.argsort(probs)[::-1]
```

Key Point. Why probability-based sorting for quantum codes?

Traditional OSD uses reliability $|p - 0.5|$, but this fails for quantum codes where:

- Most qubits have $p \approx 0$ (no error) \rightarrow reliability ≈ 0.5
- Identified errors have $p \approx 1 \rightarrow$ reliability ≈ 0.5

Both have similar reliability despite being very different! Sorting by probability directly places likely errors in $[S]$ (pivots) and unlikely errors in $[T]$ (free variables set to 0).

4.2.2 Step 2: Solving for Basis Bits (RREF)

Instead of computing $H_{[S]}^{-1}$ explicitly, the code computes the **Reduced Row Echelon Form** (RREF) of the augmented matrix $[H_{\text{sorted}} \mid \mathbf{s}]$. This is numerically stable and solves the linear system in one pass. For example, consider a parity check matrix H with 3 checks and 6 variables, and syndrome $\mathbf{s} = (1, 0, 1)^T$. Assume BP has already sorted columns by probability (column 0 = highest probability error). Then the initial augmented matrix $[H_{\text{sorted}} \mid \mathbf{s}]$ writes:

$$\left(\begin{array}{cccccc|c} 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 1 \end{array} \right) \quad (25)$$

Iteration 1: Process column 0

- Find pivot: Row 0 has a 1 in column 0 ✓
- Eliminate: Row 1 also has 1 in column 0, so XOR row 1 with row 0

$$\left(\begin{array}{cccccc|c} 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 & 0 & 1 \end{array} \right) \quad \text{pivot_cols} = [0] \quad (26)$$

Iteration 2: Process column 1

- Find pivot: Row 1 has a 1 in column 1 ✓
- Eliminate: Row 2 also has 1 in column 1, so XOR row 2 with row 1

$$\left(\begin{array}{cccccc|c} 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{array} \right) \quad \text{pivot_cols} = [0, 1] \quad (27)$$

Iteration 3: Process column 2

- Find pivot: No rows have a 1 in column 2 below the current pivot row ×
- Skip this column (it becomes a free variable)

Iteration 4: Process column 3

- Find pivot: Row 2 has a 1 in column 3 ✓
- Eliminate: Rows 0 and 1 have 1s in column 3, XOR both with row 2

$$\left(\begin{array}{cccccc|c} 1 & 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{array} \right) \quad \text{pivot_cols} = [0, 1, 3] \quad (28)$$

Final Result: Pivot columns $[S] = \{0, 1, 3\}$ (basis variables) and free columns $[T] = \{2, 4, 5\}$ (remainder variables).

In the code, the H matrix is stored as a 2D ‘int8’ array to fully utilized the GPU’s integer arithmetic capabilities. The fuction ‘compute_rref’ implements Gaussian elimination over GF(2):

```
def _get_rref_cached(self, sorted_indices: np.ndarray, syndrome: np.ndarray):
    # Reorder columns by sorted indices
    H_sorted = self.H[:, sorted_indices]
    # Build augmented matrix [H_sorted | s]
    augmented = np.hstack([H_sorted, syndrome.reshape(-1, 1)]).astype(np.int8)
    # Compute RREF in-place
    pivot_cols = self._compute_rref(augmented)
    return augmented, pivot_cols

def _compute_rref(self, M: np.ndarray) -> List[int]:
    m, n = M.shape
    pivot_row = 0
    pivot_cols = []

    for col in range(n - 1): # Don't pivot on syndrome column
        if pivot_row >= m:
```

```

        break
    # Find a row with 1 in this column
    candidates = np.where(M[pivot_row:, col] == 1)[0]
    if len(candidates) == 0:
        continue # No pivot in this column

    # Swap to bring pivot to current row
    swap_r = candidates[0] + pivot_row
    if swap_r != pivot_row:
        M[[pivot_row, swap_r]] = M[[swap_r, pivot_row]]

    pivot_cols.append(col)

    # Eliminate all other 1s in this column (XOR in GF(2))
    rows_to_xor = np.where(M[:, col] == 1)[0]
    rows_to_xor = rows_to_xor[rows_to_xor != pivot_row]
    if len(rows_to_xor) > 0:
        M[rows_to_xor, :] ^= M[pivot_row, :]

    pivot_row += 1

return pivot_cols

```

- `pivot_cols`: The column indices where pivots were found. These form the basis set $[S]$.
- After RREF, the basis submatrix has an identity-like structure, and the syndrome column contains the solution values.

4.2.3 Step 3: Reading the OSD-0 solution:

Continue from the previous example, the result can be read from the pivot columns and syndrome column s , as illustrated by the following steps:

- Set free variables to zero: $e_2 = e_4 = e_5 = 0$
- Read pivot values from the transformed syndrome column:
 - Row 0: pivot at column 0, syndrome value = 1 $\rightarrow e_0 = 1$
 - Row 1: pivot at column 1, syndrome value = 1 $\rightarrow e_1 = 1$
 - Row 2: pivot at column 3, syndrome value = 0 $\rightarrow e_3 = 0$
- Solution in sorted order then must be $e_{\text{sorted}} = (1, 1, 0, 0, 0, 0)$, and can be verified through:

$$H_{\text{sorted}} \cdot e_{\text{sorted}} = \begin{pmatrix} 1 & 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 & 0 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \oplus 0 \\ 1 \oplus 1 \\ 0 \oplus 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = s \quad \checkmark \quad (29)$$

In the code, the solution is extracted by initializing the solution vector to zeros and updating **only** the pivot positions using the transformed syndrome.

```

# OSD-0 Solution: Initialize all bits to 0 (remainder bits stay 0)
solution_base = np.zeros(self.num_errors, dtype=np.int8)

# Build pivot-to-row mapping and extract solution
pivot_row_map = {}
for r in range(augmented.shape[0]):
    row_pivots = np.where(augmented[r, :self.num_errors] == 1)[0]
    if len(row_pivots) > 0:

```

```

col = row_pivots[0]
if col in pivot_cols:
    pivot_row_map[col] = r
    # Pivot bit = transformed syndrome value
    solution_base[col] = augmented[r, -1]

```

- `solution_base = np.zeros(...)`: Ensures $e_{[T]} = \mathbf{0}$ (OSD-0 constraint).
- `augmented[r, -1]`: The transformed syndrome value. Since the basis submatrix is now identity-like, this directly gives $e_{[S]}$.

4.2.4 Step 4: Inverse Mapping (Unsort)

Finally, the solution is mapped back to the original column ordering:

```

# Remap from sorted order back to original order
estimated_errors = np.zeros(self.num_errors, dtype=int)
estimated_errors[sorted_indices] = final_solution_sorted
return estimated_errors

```

Theoretical Step	Code Realization (<code>batch_osd.py</code>)
1. Sort by soft decisions	<code>np.argsort(probs)[::-1]</code> (probability descending)
2. Select basis $[S]$	<code>pivot_cols</code> from <code>_compute_rref</code>
3. Matrix inversion	RREF transforms basis to identity structure
4. Solve $e_{[S]}$	<code>solution_base[col] = augmented[r, -1]</code>
5. Set $e_{[T]} = \mathbf{0}$	<code>np.zeros(...)</code> initialization
6. Unsort	<code>estimated_errors[sorted_indices] = solution</code>

Table 4: Mapping OSD-0 theory to `batch_osd.py` implementation

4.3 Higher-Order OSD (OSD- λ)

OSD-0 assumes the remainder error bits are zero ($e_{[T]} = \mathbf{0}$). While this provides a valid solution, it forces all “correction” work onto the basis bits $[S]$, which may result in a high-weight (improbable) error pattern.

Definition. Higher-order OSD improves this by testing non-zero configurations for the remainder bits $e_{[T]}$. For any chosen hypothesis $e_{[T]}$, the corresponding basis bits $e_{[S]}$ are uniquely determined to satisfy the syndrome:

$$e_{[S]} = H_{[S]}^{-1} \cdot (s + H_{[T]} \cdot e_{[T]}) \quad (30)$$

It is straightforward to show that the constructed error $e = (e_{[S]}, e_{[T]})$ always satisfies the parity check equation $H \cdot e = s$ and the OSD-0 is a special case of OSD- λ when $\lambda = 0$.

$$\begin{aligned}
H \cdot e &= (H_{[S]} \ H_{[T]}) \cdot \begin{pmatrix} e_{[S]} \\ e_{[T]} \end{pmatrix} = H_{[S]} \cdot e_{[S]} + H_{[T]} \cdot e_{[T]} \\
&= H_{[S]} \cdot \left[H_{[S]}^{-1} \cdot (s + H_{[T]} \cdot e_{[T]}) \right] + H_{[T]} \cdot e_{[T]} \\
&= I \cdot (s + H_{[T]} \cdot e_{[T]}) + H_{[T]} \cdot e_{[T]} \\
&= s + H_{[T]} \cdot e_{[T]} + H_{[T]} \cdot e_{[T]} = s + \mathbf{0} = s
\end{aligned} \quad (31)$$

Then the problem change to find the minimum soft-weight solution for the remainder bits $e_{[T]}$. A naive way to do this is to implement an exhaustive search (OSD-E) testing on all $2^{k'}$ patterns. This guarantees finding the minimum weight solution. Unfortunately, the remainder set $[T]$ has size $k' = n - r$, which is exponentially large in the code parameters $n - r$. To make this feasible, we restrict the search to the search depth λ , i.e., the **most suspicious** λ bits in $[T]$ (those with highest error probability among free variables) and accelerate this search process by using the GPU. The `batch_osd.py` implementation accelerates OSD-E by evaluating all 2^λ candidates in parallel on GPU. Here is how it works:

Step 1: Identify Search Columns

Select the λ free variables with highest error probability (most suspicious):

```
# Get free columns (not pivots)
all_cols = set(range(self.num_errors))
free_cols = sorted(list(all_cols - set(pivot_cols)))

# Sort free columns by probability (highest first = most suspicious)
free_cols_with_prob = [(col, probs[sorted_indices[col]]) for col in free_cols]
free_cols_with_prob.sort(key=lambda x: -x[1])

# Select top osd_order free variables for search
search_cols = [col for col, _ in free_cols_with_prob[:osd_order]]
```

Step 2: Generate All 2^λ Candidates

```
# Exhaustive: Generate all  $2^k$  combinations using bit manipulation
num_candidates = 1 << len(search_cols) #  $2^k$ 
candidates_np = np.array([
    [(i >> j) & 1 for j in range(len(search_cols))]
    for i in range(num_candidates)
], dtype=np.int8)
```

Step 3: Parallel Evaluation on GPU

All candidates are evaluated simultaneously using batched matrix operations:

```
def _evaluate_candidates_gpu(self, candidates, augmented, search_cols, probs_sorted,
                             pivot_cols):
    num_candidates = candidates.shape[0]

    # Transfer to GPU
    M_subset = torch.from_numpy(augmented[:, search_cols]).float().to(self.device)
    syndrome_col = torch.from_numpy(augmented[:, -1]).float().to(self.device)

    # Compute modified syndromes for ALL candidates in parallel
    # target_syndrome = (s + M @ e_T) % 2
    target_syndromes = (syndrome_col.unsqueeze(0) + candidates.float() @ M_subset.T)
    % 2

    # Initialize solution matrix (num_candidates x n)
    cand_solutions = torch.zeros(num_candidates, self.num_errors, device=self.device)

    # Set free variable values from candidates
    cand_solutions[:, search_cols] = candidates.float()

    # Solve for pivot variables using the RREF structure
    for r in range(augmented.shape[0]):
```

```

row_pivots = torch.where(augmented_torch[r, :] == 1)[0]
if len(row_pivots) > 0:
    pivot_c = row_pivots[0].item()
    if pivot_c in pivot_cols:
        # Pivot value = modified syndrome for this row
        cand_solutions[:, pivot_c] = target_syndromes[:, r]

# Compute soft-weighted costs and return best solution
costs = self._compute_soft_weight_gpu(cand_solutions, probs_sorted)
best_idx = torch.argmax(costs)
return cand_solutions[best_idx]

```

Step 4: Soft-Weight Cost Function

In our code, the OSD uses the **soft-weighted cost** based on log-probabilities [9]:

Definition. Soft-Weighted Cost (Log-Probability Weight). For an error pattern $e = (e_1, \dots, e_n)$ with bit-wise error probabilities $p_i = P(e_i = 1)$, the soft-weighted cost is:

$$W_{\text{soft}(e)} = \sum_{i:e_i=1} (-\log p_i) = -\sum_{i=1}^n e_i \cdot \log p_i \quad (32)$$

Lower cost indicates a more probable error pattern.

There actually are several other cost functions that can be used for OSD, such as the Hamming weight, Euclidean distance, and LLR-based weight, as listed in the following table:

Cost Function	Formula	Properties
Hamming Weight [10]	$W_{H(e)} = \sum_i e_i$	Counts flipped bits; ignores probabilities
Soft Weight (Log-Prob) [9]	$W_{\text{soft}(e)} = -\sum_i e_i \log p_i$	Weights by $-\log p_i$; approximates ML
Euclidean Distance [11]	$d_E^2 = \sum_i (r_i - c_i)^2$	For AWGN channels with continuous signals
LLR-Based Weight [12]	$W_{\text{LLR}(e)} = \sum_i e_i L_i $	Uses log-likelihood ratios $L_i = \log\left(\frac{p_i}{1-p_i}\right)$

Table 5: Comparison of cost functions for selecting the best error pattern

Then why soft weight is preferred for BP+OSD? This is because the soft weight approximates the Maximum-Likelihood Decoding (ML) objective [13]. The ML decoder selects the error pattern e^* that maximizes the posterior probability

$$e^* = \arg \max_e P(e \mid \text{syndrome}). \quad (33)$$

Taking the logarithm, which is a monotonic transformation, and suppose the error e_i are independent, each with probability $P(e_i = 1) = p_i$, this becomes:

$$e^* = \arg \max_e \sum_i [e_i \log p_i + (1 - e_i) \log(1 - p_i)] \quad (34)$$

For sparse errors where most $e_i = 0$, minimizing $W_{\text{soft}(e)} = -\sum_i e_i \log p_i$ closely approximates the ML objective.

Key Point. Question: Why not use the Hamming weight or Euclidean distance as the cost function?

In the code, the soft-weight cost function is implemented as follows:

```
def _compute_soft_weight_gpu(self, solutions, probs):
    # Clip to avoid log(0)
    probs_clipped = torch.clamp(probs, 1e-10, 1 - 1e-10)
    # Log-probability weights: -log(p) penalizes flipping low-probability bits
    log_weights = -torch.log(probs_clipped)
    # Total cost = sum of weights for flipped bits
    costs = (solutions * log_weights).sum(dim=1)
    return costs
```

4.4 Combination Sweep Strategy (OSD-CS)

To allow for a larger search depth (e.g., $\lambda \approx 50 - 100$) without exponential cost, we use the **combination sweep** strategy, first proposed for reducing error floors in classical LDPC codes and adapted for quantum codes by Roffe et al. [9].

Definition. OSD-CS assumes the true error pattern on the remainder bits is **sparse**. Instead of checking **all** 2^λ patterns on λ most suspicious bits, it only checks those with low Hamming weight ($w = 0, 1, 2$). Exhausted on those strings only take $C_\lambda^0 + C_\lambda^1 + C_\lambda^2 = 1 + \lambda + \frac{\lambda(\lambda-1)}{2}$ candidates.

Algorithm Steps:

1. **Sort:** Select the λ most suspicious positions in $[T]$ (highest probability among free variables).
2. **Sweep:** Generate candidate vectors $\mathbf{e}_{[T]}$ with:
 - **Weight 0:** The zero vector (equivalent to OSD-0).
 - **Weight 1:** All single-bit flips among the λ bits.
 - **Weight 2:** All pairs of bit flips among the λ bits.
3. **Select:** Calculate $\mathbf{e}_{[S]}$ for each candidate, compute the soft-weighted cost, and pick the best one.

Method	Complexity	Use Case
OSD-0	$O(1)$	Fastest, baseline performance
OSD-E (Exhaustive)	$O(2^\lambda)$	Optimal for small λ (≤ 15)
OSD-CS (Comb. Sweep)	$O(\lambda^2)$	Near-optimal for large λ (≈ 60)

Table 6: Comparison of OSD Search Strategies

Why OSD-CS works for Quantum Codes? In the low-error regime relevant for QEC, it is statistically very unlikely that the optimal solution requires flipping 3+ bits in the specific subset of “uncertain” remainder bits. Checking only weights 0, 1, and 2 captures the vast majority of likely error configurations while reducing complexity from exponential to polynomial (quadratic) [9].

Definition. Combination sweep is a greedy search testing configurations by likelihood:

1. **Sort remainder bits:** Order bits in $[T]$ by error probability (most likely first)

2. **Test weight-0:** The zero vector (OSD-0 baseline)
3. **Test weight-1:** Set each single bit in $e_{[T]}$ to 1 (all k possibilities)
4. **Test weight-2:** Set each pair among the first λ bits to 1

Keep the minimum soft-weight solution found.

Recall that the **binomial coefficient** $\binom{\lambda}{2} = \frac{\lambda(\lambda-1)}{2}$ counts ways to choose 2 items from λ .

Total configurations: $1 + k + \binom{\lambda}{2}$

With $\lambda = 60$: approximately $1 + k + 1770$ configurations (vs 2^k for exhaustive search!)

4.4.1 OSD-CS Implementation in `batch_osd.py`

The GPU-accelerated implementation realizes OSD-CS by explicitly generating sparse error patterns (weights 0, 1, and 2) instead of iterating through all binary combinations.

Step 1: Generating Sparse Candidates

The `_generate_osd_cs_candidates` method generates candidate vectors $e_{[T]}$ with structured loops:

```
def _generate_osd_cs_candidates(self, k: int, osd_order: int) -> np.ndarray:
    """Generate OSD-CS (Combination Sweep) candidate strings."""
    candidates = []

    # Weight 0: Zero vector (OSD-0 baseline)
    candidates.append(np.zeros(k, dtype=np.int8))

    # Weight 1: Single-bit flips (k candidates)
    for i in range(k):
        candidate = np.zeros(k, dtype=np.int8)
        candidate[i] = 1
        candidates.append(candidate)

    # Weight 2: Two-bit flips (limited to osd_order)
    for i in range(min(osd_order, k)):
        for j in range(i + 1, min(osd_order, k)):
            candidate = np.zeros(k, dtype=np.int8)
            candidate[i] = 1
            candidate[j] = 1
            candidates.append(candidate)

    return np.array(candidates, dtype=np.int8)
```

- `np.zeros(k)`: The “baseline” hypothesis (OSD-0 solution).
- `range(k)` loop: Adds k candidates, each with a single bit flip at index i .
- Nested `range(limit)` loop: Adds $\binom{\lambda}{2}$ candidates, representing pairs of flips at indices (i, j) .

Step 2: Integration into the Solve Loop

The `solve` method switches between OSD-E and OSD-CS based on the `osd_method` parameter:

```
# Generate candidates based on method
if osd_method == 'combination_sweep':
    # OSD-CS:  $O(\lambda^2)$  sparse candidates
    candidates_np = self._generate_osd_cs_candidates(len(search_cols), osd_order)
```

```

else:
    # Exhaustive:  $O(2^k)$  all combinations
    num_candidates = 1 << len(search_cols)
    candidates_np = np.array([[ (i > j) & 1 for j in range(len(search_cols)) ]
                               for i in range(num_candidates)], dtype=np.int8)

# Transfer to GPU and evaluate all candidates in parallel
candidates = torch.from_numpy(candidates_np).to(self.device)
best_solution_sorted = self._evaluate_candidates_gpu(
    candidates, augmented, search_cols, probs_sorted, pivot_cols
)

```

Both methods use the same GPU evaluation function—only the candidate generation differs.

Complexity Comparison

Method	Candidates	Example ($\lambda = 15$)
OSD-E	2^λ	32768
OSD-CS	$1 + k + \binom{\lambda}{2}$	$\approx 1 + k + 105$

Table 7: Number of candidates for OSD-E vs OSD-CS

Weight Class	Code Realization
Weight 0	<code>candidates.append(np.zeros(k))</code>
Weight 1	<code>for i in range(k): candidate[i] = 1</code>
Weight 2	<code>for i in range(limit): for j in range(i+1, limit): ...</code>

Table 8: Mapping OSD-CS theory to `batch_osd.py` implementation

5 The Complete BP+OSD Decoder

5.1 Algorithm Flow

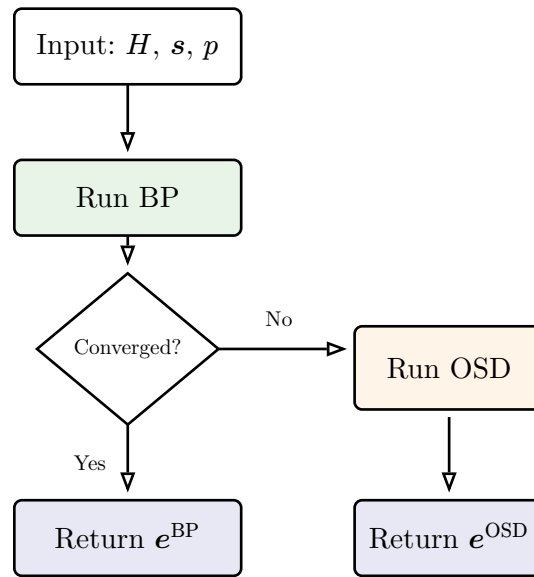


Figure 13: BP+OSD decoder flowchart

Key Point.

- If BP succeeds (converges): use BP result — fast!
- If BP fails: use OSD to resolve degeneracy — always gives valid answer

5.2 Complete Algorithm: BP+OSD-CS

Algorithm 3: BP+OSD-CS Decoder

Input: Parity matrix H ($m \times n$, rank r), syndrome s , error prob p , depth $\lambda=60$
Output: Error estimate e satisfying $H \cdot e = s$

```
function BP_OSD_CS(H, s, p,  $\lambda$ ):
    // ===== STAGE 1: Run Belief Propagation (Algorithm 1) =====
    (converged, e_BP, P_1) = BP(H, s, p)
    if converged:
        return e_BP

    // ===== STAGE 2: OSD-0 (Algorithm 2) =====
    [0_BP] = argsort(P_1)           // Sort: most likely flipped first
    H_sorted = H[:, 0_BP]          // Reorder columns
    [S] = first r linearly independent columns of H_sorted
    [T] = remaining  $k' = n - r$  columns

    e_[S] = H_[S]^(-1)  $\times$  s        // Solve on basis
    e_[T] = zeros( $k'$ )              // Set remainder to zero
    best = (e_[S], e_[T])
    best_wt = hamming_weight(best)

    // ===== STAGE 3: Combination Sweep =====
    // Weight-1 search: try flipping each remainder bit
    for i = 0 to  $k'-1$ :
        e_[T] = zeros( $k'$ ); e_[T][i] = 1
        e_[S] = H_[S]^(-1)  $\times$  (s + H_[T]  $\times$  e_[T])
        if hamming_weight((e_[S], e_[T])) < best_wt:
            best = (e_[S], e_[T])
            best_wt = hamming_weight(best)

    // Weight-2 search: try flipping pairs in first  $\lambda$  bits
    for i = 0 to min( $\lambda$ ,  $k'$ )-1:
        for j = i+1 to min( $\lambda$ ,  $k'$ )-1:
            e_[T] = zeros( $k'$ ); e_[T][i] = 1; e_[T][j] = 1
            e_[S] = H_[S]^(-1)  $\times$  (s + H_[T]  $\times$  e_[T])
            if hamming_weight((e_[S], e_[T])) < best_wt:
                best = (e_[S], e_[T])
                best_wt = hamming_weight(best)

    return inverse_permute(best, 0_BP) // Remap to original ordering
```

Listing 1: Complete BP+OSD-CS algorithm

5.3 BP Convergence and Performance Guarantees

Theorem (BP Convergence on Trees). [2], [14] If the factor graph $G = (V, U, E)$ is a **tree** (contains no cycles), then BP converges to the **exact** marginal probabilities $P(e_j = 1 \mid \mathbf{s})$ in at most d iterations, where d is the diameter of the tree (maximum distance between any two nodes).

Proof. We prove exactness by induction on the tree structure, using the factorization property of graphical models.

Factorization on Trees: For a tree-structured factor graph, the joint probability distribution factors as:

$$P(\mathbf{x}) = \frac{1}{Z} \prod_{a \in \mathcal{F}} \psi_a(\mathbf{x}_{\mathcal{N}(a)}) \quad (35)$$

where \mathcal{F} is the set of factors, $\mathcal{N}(a)$ are neighbors of factor a , and Z is the partition function.

Key Property: On a tree, removing any node v separates the graph into disjoint connected components (subtrees). By the global Markov property, variables in different subtrees are conditionally independent given v .

Base Case (Leaf Nodes): Consider a leaf variable node v with single neighbor (factor) a . The message $\mu_{v \rightarrow a}(x_v)$ depends only on the local evidence $P(y_v \mid x_v)$. Since there are no other dependencies, this message is exact at iteration 1.

Inductive Step: Assume messages from all nodes at distance $> k$ from root are exact. Consider node u at distance k with neighbors $\mathcal{N}(u) = \{a_1, \dots, a_m\}$.

For message $\mu_{u \rightarrow a_i}(x_u)$, the BP update is:

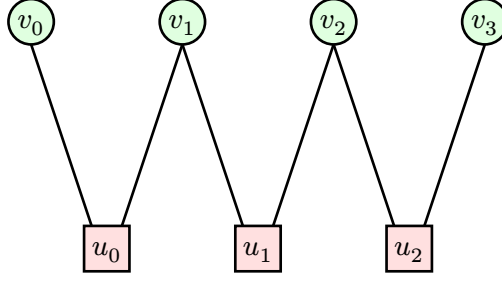
$$\mu_{u \rightarrow a_i}(x_u) \propto P(y_u \mid x_u) \prod_{a_j \in \mathcal{N}(u) \setminus \{a_i\}} \mu_{a_j \rightarrow u}(x_u) \quad (36)$$

By the separation property, removing u creates m independent subtrees rooted at $\{a_1, \dots, a_m\}$. By the inductive hypothesis, messages from these subtrees are exact marginals of their respective subtrees. Since subtrees are conditionally independent given u , the product of messages equals the joint probability of all subtree configurations, making $\mu_{u \rightarrow a_i}(x_u)$ exact.

Termination: After d iterations (where d is the tree diameter), messages have propagated from all leaves to all nodes. Each node's belief $b_{v(x_v)} \propto P(y_v \mid x_v) \prod_{a \in \mathcal{N}(v)} \mu_{a \rightarrow v}(x_v)$ equals the exact marginal $P(x_v \mid \mathbf{y})$ by the factorization property.

Therefore, BP computes exact marginals on trees in d iterations. \square

Example: Consider the $[7, 4, 3]$ Hamming code with tree-structured factor graph:



Tree structure: diameter $d = 4$, BP converges in 4 iterations

Figure 14: Tree-structured code where BP gives exact solution

For this tree with syndrome $\mathbf{s} = (1, 0, 0)$ and $p = 0.1$:

- BP converges in $d = 4$ iterations
- Output: $\mathbf{e}^{\text{BP}} = (1, 0, 0, 0, 0, 0, 0)$ (single bit flip at position 0)
- This is the **exact** maximum likelihood solution

Theorem (BP Performance on Graphs with Cycles). [15], [16] For an (l, q) -LDPC code with factor graph of **girth** g (minimum cycle length), BP provides the following guarantees:

1. **Local optimality:** If the true error \mathbf{e}^* has Hamming weight $|\mathbf{e}^*| < g/2$, then BP converges to \mathbf{e}^* with high probability (for sufficiently small p).
2. **Approximation bound:** For codes with girth $g \geq 6$ and maximum degree $\Delta = \max(l, q)$, if BP converges, the output \mathbf{e}^{BP} satisfies:

$$|\mathbf{e}^{\text{BP}}| \leq (1 + \varepsilon(g, \Delta)) \cdot |\mathbf{e}^*| \quad (37)$$

where $\varepsilon(g, \Delta) \rightarrow 0$ as $g \rightarrow \infty$ for fixed Δ .

3. **Iteration complexity:** BP requires $O(g)$ iterations to propagate information across the shortest cycle.

Proof. Part 1 (Local optimality): Consider an error \mathbf{e}^* with $|\mathbf{e}^*| < g/2$. In the factor graph, the neighborhood of radius $r = \lfloor g/2 \rfloor - 1$ around any error bit is a tree (no cycles within distance r). Within this tree neighborhood:

- BP computes exact marginals (by Theorem 1)
- The error bits are separated by distance $\geq g/2$
- No interference between error regions

Therefore, BP correctly identifies each error bit independently, giving $\mathbf{e}^{\text{BP}} = \mathbf{e}^*$.

Part 2 (Approximation bound): For $|\mathbf{e}^*| \geq g/2$, cycles create dependencies. The approximation error comes from:

- **Double-counting:** Evidence circulates through cycles
- **Correlation:** Nearby error bits are not independent

For girth g , the correlation decays exponentially with distance. The number of length- g cycles through a node is bounded by Δ^g . Using the correlation decay lemma for loopy belief propagation, the relative error in log-likelihood ratios is:

$$\varepsilon(g, \Delta) \leq C \cdot \Delta^{2-g/2} \quad (38)$$

for some constant C . This translates to the weight approximation bound.

Part 3 (Iteration complexity): Information propagates one edge per iteration. To detect a cycle of length g , messages must travel distance g , requiring $O(g)$ iterations. \square

Key Point. Practical implications:

- Codes with large girth g (e.g., $g \geq 8$) allow BP to correct more errors
- Random LDPC codes typically have $g = O(\log n)$, giving good BP performance
- Structured codes (e.g., Toric code with $g = 4$) have small girth, leading to BP failures
- The degeneracy problem in quantum codes compounds the cycle problem, making OSD necessary

5.3.1 Density Evolution Framework

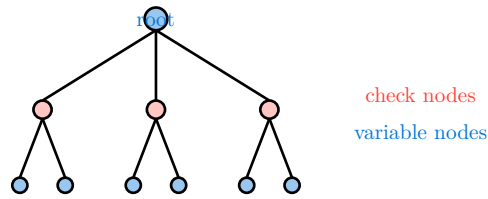
We now develop the rigorous theoretical foundations that explain **when** and **why** BP converges in different regimes, drawing from asymptotic analysis via density evolution [3], variational optimization through statistical physics [4], and combinatorial failure modes [5].

For infinite-length random LDPC codes, convergence is analyzed through the **density evolution** method [15], which tracks the probability distributions of messages rather than individual message values.

Definition. Cycle-Free Horizon: For a random LDPC code with block length $n \rightarrow \infty$, the **computation tree** of depth l rooted at any edge is the subgraph containing all nodes reachable within l hops. The cycle-free horizon property states:

$$\lim_{n \rightarrow \infty} \mathbb{P}(\text{cycle in depth-}l \text{ tree}) = 0 \quad (39)$$

This means that for any fixed number of iterations l , the local neighborhood appears tree-like with probability approaching 1 as $n \rightarrow \infty$.



Computation tree of depth $l = 2$: no cycles

Figure 15: Locally tree-like structure in large random graphs

Key Point. The cycle-free horizon is the mathematical justification for applying tree-based convergence proofs to loopy graphs in the asymptotic limit. It explains why BP performs well on long random LDPC codes despite the presence of cycles.

Definition. Concentration Theorem: Let Z be a performance metric (e.g., bit error rate) of BP after l iterations on a code randomly drawn from ensemble $\mathcal{C}(n, \lambda, \rho)$, where $\lambda(x)$ and $\rho(x)$ are the variable and check node degree distributions. For any $\varepsilon > 0$:

$$\mathbb{P}(|Z - \mathbb{E}[Z]| > \varepsilon) \leq e^{-\beta n \varepsilon^2} \quad (40)$$

where $\beta > 0$ depends on the ensemble parameters.

Interpretation: As $n \rightarrow \infty$, almost all codes in the ensemble perform identically to the ensemble average. Individual code performance concentrates around the mean with exponentially small deviation probability.

Key Point. Concentration visualization: The performance metric Z concentrates exponentially around its ensemble average $\mathbb{E}[Z]$. For large block length n , the probability of deviation greater than ε decays as $e^{-\beta n \varepsilon^2}$, meaning almost all codes perform identically to the average.

The proof of the Concentration Theorem uses martingale theory:

Proof. Proof sketch via Doob's Martingale:

1. **Martingale Construction:** View code selection as revealing edges sequentially. Define $Z_i = \mathbb{E}[Z \mid \text{first } i \text{ edges revealed}]$. This forms a Doob martingale: $\mathbb{E}[Z_{i+1} \mid Z_0, \dots, Z_i] = Z_i$.
2. **Bounded Differences:** In a sparse graph with maximum degree Δ , changing a single edge affects at most $O(\Delta^l)$ messages after l iterations. Since Δ is constant and l is fixed, the change in Z is bounded: $|Z_i - Z_{i-1}| \leq c/n$ for some constant c .
3. **Azuma-Hoeffding Inequality:** For a martingale with bounded differences $|Z_i - Z_{i-1}| \leq c_i$:

$$\mathbb{P}(|Z_m - Z_0| > \varepsilon) \leq 2 \exp\left(-\frac{\varepsilon^2}{2 \sum_{i=1}^m c_i^2}\right) \quad (41)$$

4. **Application:** With $m = O(n)$ edges and $c_i = O(1/n)$, we have $\sum c_i^2 = O(1/n)$, giving:

$$\mathbb{P}(|Z - \mathbb{E}[Z]| > \varepsilon) \leq 2 \exp\left(-\frac{\varepsilon^2 n}{2C}\right) = e^{-\beta n \varepsilon^2} \quad (42)$$

where $\beta = 1/(2C)$. □

Theorem (Threshold Theorem). For a code ensemble with degree distributions $\lambda(x), \rho(x)$ and a symmetric channel with noise parameter σ (e.g., standard deviation for AWGN), there exists a unique **threshold** σ^* such that:

1. If $\sigma < \sigma^*$ (low noise): As $l \rightarrow \infty$, the probability of decoding error $P_e^{(l)} \rightarrow 0$
2. If $\sigma > \sigma^*$ (high noise): $P_e^{(l)}$ remains bounded away from zero

The threshold is determined by the fixed points of the density evolution recursion:

$$P_{l+1} = \Phi(P_l, \sigma) \quad (43)$$

where Φ is the density update operator combining variable and check node operations.

Key Point. Threshold phenomenon: There exists a sharp transition at σ^* . Below this threshold (low noise), BP converges to zero error as iterations increase. Above threshold (high noise), errors persist. This sharp phase transition is characteristic of random LDPC ensembles.

Key Point. Why the threshold exists: The density evolution operator Φ has two competing fixed points:

- **All-correct fixed point:** Messages concentrate at $\pm\infty$ (high confidence)
- **Error fixed point:** Messages remain near zero (low confidence)

Below threshold, the all-correct fixed point is stable and attracts all trajectories. Above threshold, the error fixed point becomes stable, trapping the decoder.

5.3.2 Variational Perspective: Bethe Free Energy

The density evolution framework applies to infinite-length codes. For finite loopy graphs, we need a different lens: **statistical physics** [4]. This reveals that BP is actually performing **variational optimization** of an energy function.

Definition. Bethe Free Energy: For a factor graph with variables $\mathbf{x} = (x_1, \dots, x_n)$ and factors ψ_a , let $b_{i(x_i)}$ be the **belief** (pseudo-marginal) at variable i and $b_{a(x_a)}$ be the belief at factor a . The Bethe Free Energy is:

$$F_{\text{Bethe}(b)} = \sum_a \sum_{\mathbf{x}_a} b_{a(\mathbf{x}_a)} E_{a(\mathbf{x}_a)} - H_{\text{Bethe}(b)} \quad (44)$$

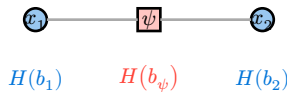
where the **Bethe entropy** approximates the true entropy using local entropies:

$$H_{\text{Bethe}} = \sum_a H(b_a) + \sum_i (1 - d_i) H(b_i) \quad (45)$$

Here d_i is the degree of variable i , and $H(b) = -\sum_x b(x) \log b(x)$ is the Shannon entropy.

Constraints: Beliefs must be normalized and **marginally consistent**:

$$\sum_{\mathbf{x}_a \setminus x_i} b_{a(\mathbf{x}_a)} = b_{i(x_i)} \quad \text{for all } i \in a \quad (46)$$



$$\text{Bethe entropy: } H_{\text{Bethe}} = H(b_\psi) + (1 - 2)H(b_1) + (1 - 2)H(b_2)$$

($d_1 = d_2 = 2$ for this graph)

Figure 16: Bethe entropy decomposes global entropy into local terms

Key Point. Intuition: The Bethe approximation treats each factor independently, summing local entropies. The $(1 - d_i)$ correction prevents double-counting: a variable connected to d_i factors appears in d_i factor entropies, so we subtract $(d_i - 1)$ copies of its individual entropy.

Theorem (Yedidia-Freeman-Weiss). [4] A set of beliefs $\{b_i, b_a\}$ is a **fixed point** of the Sum-Product BP algorithm if and only if it is a **stationary point** (critical point) of the Bethe Free Energy $F_{\text{Bethe}(b)}$ subject to normalization and marginalization constraints.

Equivalently: BP performs coordinate descent on the Bethe Free Energy. Each message update corresponds to minimizing F_{Bethe} with respect to one edge's belief.

Proof. Proof sketch via Lagrangian:

1. **Constrained optimization:** Form the Lagrangian:

$$\mathcal{L} = F_{\text{Bethe}(b)} + \sum_{i,a} \sum_{x_i} \lambda_{ia}(x_i) \left(b_{i(x_i)} - \sum_{x_a \setminus x_i} b_{a(x_a)} \right) + \text{normalization terms} \quad (47)$$

2. **Stationarity conditions:** Taking $\partial \mathcal{L} / \partial b_a = 0$ and $\partial \mathcal{L} / \partial b_i = 0$:

$$b_{a(x_a)} \propto \psi_{a(x_a)} \prod_{i \in a} \exp(\lambda_{ia}(x_i)) \quad (48)$$

$$b_{i(x_i)} \propto \prod_{a \in i} \exp(\lambda_{ia}(x_i)) \quad (49)$$

3. **Message identification:** Define messages $\mu_{i \rightarrow a}(x_i) = \exp(\lambda_{ia}(x_i))$. Substituting and enforcing marginalization constraints yields exactly the BP update equations:

$$\mu_{i \rightarrow a}(x_i) \propto P(y_i | x_i) \prod_{a' \in i \setminus a} \mu_{a' \rightarrow i}(x_i) \quad (50)$$

$$\mu_{a \rightarrow i}(x_i) \propto \sum_{x_a \setminus x_i} \psi_{a(x_a)} \prod_{i' \in a \setminus i} \mu_{i' \rightarrow a}(x_{i'}) \quad (51)$$

□

Key Point. Energy landscape interpretation: BP performs gradient descent on the Bethe Free Energy landscape. On trees, there's a single global minimum (correct solution). On loopy graphs, local minima can trap the decoder, corresponding to incorrect fixed points. The contour lines represent energy levels, with BP trajectories flowing toward minima.

Key Point. Implications for convergence:

- **On trees:** Bethe approximation is exact ($F_{\text{Bethe}} = F_{\text{Gibbs}}$), so BP finds the global minimum
- **On loopy graphs:** F_{Bethe} is an approximation. BP finds a local minimum, which may not be the true posterior

- **Stable fixed points** correspond to local minima of F_{Bethe}
- **Unstable fixed points** (saddle points) cause oscillations

This explains why BP can converge to incorrect solutions: it gets trapped in local minima created by graph cycles.

5.3.3 Sufficient Conditions for Convergence

While density evolution guarantees asymptotic convergence and Bethe theory explains fixed points, neither provides **guarantees** for specific finite loopy graphs. We now present rigorous sufficient conditions [17].

Definition. Dobrushin's Influence Matrix: For a graphical model, the influence C_{ij} measures the maximum change in the marginal distribution of variable i caused by fixing variable j :

$$C_{ij} = \sup_{x_j, x_{j'}} \|P(x_i | x_j) - P(x_i | x_{j'})\|_{\text{TV}} \quad (52)$$

where $\|\cdot\|_{\text{TV}}$ is the total variation distance.

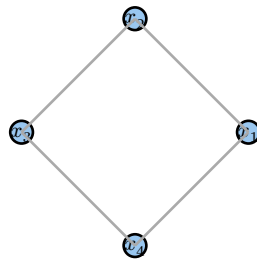
The **Dobrushin interdependence matrix** \mathbf{C} has entries C_{ij} for $i \neq j$ and $C_{ii} = 0$.

Theorem (Dobrushin's Uniqueness Condition). If the Dobrushin matrix satisfies:

$$\|\mathbf{C}\|_{\infty} = \max_i \sum_{j \neq i} C_{ij} < 1 \quad (53)$$

then:

1. The Gibbs measure has a unique fixed point
2. BP converges exponentially fast to this fixed point from any initialization
3. The convergence rate is $\lambda = \|\mathbf{C}\|_{\infty}$



Example: 4-cycle with weak coupling
 $\|\mathbf{C}\|_{\infty} = \max_i \sum_j C_{ij} = 2 \cdot 0.3 = 0.6 < 1 \checkmark$

Figure 17: Dobrushin condition: information dissipates through the graph

Key Point. Limitation for LDPC codes: Error correction codes are designed to **propagate** information over long distances. Parity checks impose hard constraints (infinite coupling strength). Therefore, useful LDPC codes typically **violate** Dobrushin's condition.

While sufficient, Dobrushin's condition is far from necessary. It applies mainly to high-noise regimes where correlations are weak.

Theorem (Contraction Mapping Convergence). View BP as a mapping $T : \mathcal{M} \rightarrow \mathcal{M}$ on the space of messages. If T is a **contraction** under some metric d :

$$d(T(\mathbf{m}), T(\mathbf{m}')) \leq \lambda \cdot d(\mathbf{m}, \mathbf{m}') \quad (54)$$

with Lipschitz constant $\lambda < 1$, then:

1. BP has a unique fixed point \mathbf{m}^*
2. BP converges geometrically: $d(\mathbf{m}^{(t)}, \mathbf{m}^*) \leq \lambda^t d(\mathbf{m}^{(0)}, \mathbf{m}^*)$

Proof. Proof: Direct application of the Banach Fixed Point Theorem. The contraction property ensures:

- **Uniqueness:** If \mathbf{m}^* and \mathbf{m}'^* are both fixed points, then:

$$d(\mathbf{m}^*, \mathbf{m}'^*) = d(T(\mathbf{m}^*), T(\mathbf{m}'^*)) \leq \lambda \cdot d(\mathbf{m}^*, \mathbf{m}'^*) \quad (55)$$

Since $\lambda < 1$, this implies $d(\mathbf{m}^*, \mathbf{m}'^*) = 0$, so $\mathbf{m}^* = \mathbf{m}'^*$.

- **Convergence:** For any initialization $\mathbf{m}^{(0)}$:

$$d(\mathbf{m}^{(t+1)}, \mathbf{m}^*) = d(T(\mathbf{m}^{(t)}), T(\mathbf{m}^*)) \leq \lambda \cdot d(\mathbf{m}^{(t)}, \mathbf{m}^*) \quad (56)$$

Iterating gives $d(\mathbf{m}^{(t)}, \mathbf{m}^*) \leq \lambda^t d(\mathbf{m}^{(0)}, \mathbf{m}^*)$. □

Key Point. Spectral radius condition: For binary pairwise models, the contraction constant can be computed from the **spectral radius** of the interaction matrix:

$$\rho(\mathbf{A}) < 1, \quad \text{where } A_{ij} = \tanh|J_{ij}| \quad (57)$$

This is sharper than Dobrushin's condition (which corresponds to the L_∞ norm of \mathbf{A}).

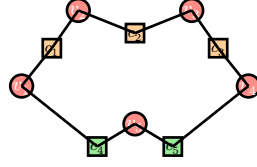
5.3.4 Failure Mechanisms: Trapping Sets

The previous sections explain when BP converges. We now characterize when and why it **fails** [5]. In the high-SNR regime, BP can get trapped in incorrect fixed points due to specific graph substructures.

Definition. (a,b) Absorbing Set: A subset $\mathcal{D} \subseteq V$ of a variable nodes is an (a, b) absorbing set if:

1. The induced subgraph contains exactly b **odd-degree** check nodes (unsatisfied checks)
2. Every variable node $v \in \mathcal{D}$ has **strictly more** even-degree neighbors than odd-degree neighbors in the induced subgraph

Interpretation: If the variables in \mathcal{D} are in error, each receives more “confirming” messages (from satisfied checks) than “correcting” messages (from unsatisfied checks), causing the decoder to stabilize in the error state.



5 error variables (red)
 3 odd-degree checks (orange)
 Even-degree checks (green)

Canonical (5,3) absorbing set: each variable has ≥ 2 even neighbors

Figure 18: The (5,3) absorbing set: a stable error configuration

Key Point. Why BP gets trapped:

1. **Majority vote:** Each variable node performs a weighted majority vote of its check neighbors
2. **Satisfied checks dominate:** In an absorbing set, satisfied checks (even degree) outnumber unsatisfied checks (odd degree) for each variable
3. **Reinforcement loop:** Satisfied checks send messages that **confirm** the error state, while unsatisfied checks send weak correction signals
4. **Stable fixed point:** The configuration becomes a local minimum of the Bethe Free Energy

This is the primary cause of **error floors** in LDPC codes: at high SNR, rare noise patterns that activate absorbing sets dominate the error probability.

Theorem (Absorbing Sets and Error Floors). For an LDPC code with minimum absorbing set size (a_{\min}, b_{\min}) , the error floor is dominated by:

$$P_{\text{error}} \approx \binom{n}{a_{\min}} \cdot p^{a_{\min}} \cdot (1-p)^{n-a_{\min}} \cdot P_{\text{trap}} \quad (58)$$

where P_{trap} is the probability that BP fails to correct the absorbing set configuration.

Implication: Error floor height is determined by the **size** and **multiplicity** of small absorbing sets. Code design focuses on eliminating small absorbing sets.

6 Minimum Weight Perfect Matching (MWPM) Decoder

6.1 Maximum Likelihood Decoding and MWPM

Maximum Likelihood Decoding (MLD) seeks the most probable error pattern \mathbf{e} given syndrome \mathbf{s} and error probabilities $p(\mathbf{e})$.

Definition. Maximum Likelihood Decoding Problem: Given parity check matrix $\mathbf{H} \in \mathbb{F}_2^{m \times n}$, syndrome $\mathbf{s} \in \mathbb{F}_2^m$, and error weights $w_i = \ln\left(\frac{1-p_i}{p_i}\right)$, find:

$$\min_{\mathbf{c} \in \mathbb{F}_2^n} \sum_{i \in [n]} w_i c_i \quad \text{subject to} \quad \mathbf{H}\mathbf{c} = \mathbf{s} \quad (59)$$

For certain code structures, MLD can be efficiently reduced to a graph matching problem.

Theorem (MLD to MWPM Reduction). If every column of \mathbf{H} has at most 2 non-zero elements (each error triggers at most 2 detectors), then MLD can be deterministically reduced to Minimum Weight Perfect Matching with boundaries in polynomial time.

Definition. Detector Graph: Given $\mathbf{H} \in \mathbb{F}_2^{m \times n}$, construct graph $G = (V, E)$ where:

- Vertices: $V = [m] \cup \{0\}$ (detectors plus boundary vertex)
- Edges: For each column i of \mathbf{H} :
 - If column i has weight 2 (triggers detectors x_1, x_2): edge (x_1, x_2) with weight w_i
 - If column i has weight 1 (triggers detector x): edge $(x, 0)$ with weight w_i

Proof. Reduction procedure:

1. **Graph construction:** Build detector graph G from \mathbf{H} as defined above. The boundary operator $\partial : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^{m+1}$ maps edge vectors to vertex vectors, corresponding to the parity check matrix.
2. **Syndrome to boundary:** Given syndrome $\mathbf{s} \in \mathbb{F}_2^m$, identify the set $D \subseteq V$ of vertices with non-zero syndrome values. This becomes the boundary condition for the matching problem.
3. **Shortest path computation:** For all pairs (u, v) where $u, v \in D \cup \{0\}$, compute shortest paths using Dijkstra's algorithm. This requires $O(|D|^2)$ shortest path computations, constructing a complete weighted graph on $D \cup \{0\}$.
4. **MWPM with boundary:** Solve MWPM on the complete graph with boundary vertex $\{0\}$. The solution gives edges whose boundary equals D , which corresponds to the minimum weight error pattern satisfying $\mathbf{H}\mathbf{c} = \mathbf{s}$.

Since each step is polynomial time, MLD reduces to MWPM in polynomial time. \square

6.2 The Matching Polytope

Definition. Weighted Perfect Matching: Given weighted graph $G = (V, E, W)$ where $W = \{w_e \in \mathbb{R} \mid e \in E\}$:

- A **matching** $M \subseteq E$ has no two edges sharing a vertex
- A **perfect matching** covers every vertex in V

- The **weight** of matching M is $\sum_{e \in M} w_e$

The integer programming formulation uses indicator variables $x_e \in \{0, 1\}$:

$$\min_{\mathbf{x}} \sum_{e \in E} w_e x_e \quad \text{subject to} \quad \sum_{e \in \delta(\{v\})} x_e = 1 \quad \forall v \in V, \quad x_e \in \{0, 1\} \quad (60)$$

where $\delta(\{v\})$ denotes edges incident to vertex v .

Theorem (Matching Polytope Characterization). Define the odd set family $\mathcal{O}(G) = \{U \subseteq V : |U| \text{ is odd and } \geq 3\}$. Let:

$$\begin{aligned} P_1(G) &= \text{conv} \left\{ \mathbf{x} : x_e \in \{0, 1\}, \sum_{e \in \delta(\{v\})} x_e = 1 \quad \forall v \in V \right\} \\ P_2(G) &= \left\{ \mathbf{x} : x_e \geq 0, \sum_{e \in \delta(\{v\})} x_e = 1 \quad \forall v \in V, \sum_{e \in \delta(U)} x_e \geq 1 \quad \forall U \in \mathcal{O}(G) \right\} \end{aligned} \quad (61)$$

If edge weights are rational, then $P_1(G) = P_2(G)$.

Implication: The integer program can be relaxed to a linear program by replacing $x_e \in \{0, 1\}$ with $x_e \geq 0$ and adding the **blossom constraints** $\sum_{e \in \delta(U)} x_e \geq 1$ for all odd sets U .

Key Point. Why blossom constraints matter: An odd set U cannot have a perfect matching using only internal edges (odd number of vertices). Therefore, at least one edge must connect to the outside: $\sum_{e \in \delta(U)} x_e \geq 1$. This constraint is necessary and sufficient for the convex hull to equal the integer hull.

6.3 Dual Formulation and Optimality Conditions

Definition. MWPM Dual Problem: The dual of the MWPM linear program is:

$$\max_{\mathbf{y}} \sum_{v \in V} y_v + \sum_{O \in \mathcal{O}(G)} y_O \quad (62)$$

subject to:

$$\begin{aligned} \lambda_e = w_e - (y_{v_1} + y_{v_2}) - \sum_{O: e \in \delta(O)} y_O &\geq 0 \quad \forall e \in E \\ y_O &\geq 0 \quad \forall O \in \mathcal{O}(G) \end{aligned} \quad (63)$$

where λ_e is the **slack** of edge e .

Theorem (KKT Complementary Slackness). Primal solution \mathbf{x} and dual solution $(\mathbf{y}, \{y_O\})$ are optimal if and only if:

1. **Primal feasibility:** $\sum_{e \in \delta(\{v\})} x_e = 1$, $\sum_{e \in \delta(U)} x_e \geq 1$, $x_e \geq 0$
2. **Dual feasibility:** $\lambda_e \geq 0$, $y_O \geq 0$
3. **Complementary slackness:**
 - $\lambda_e x_e = 0$ (tight edges are in matching)

$$\bullet \quad y_O \left(\sum_{e \in \delta(O)} x_e - 1 \right) = 0 \text{ (tight odd sets have positive dual)}$$

6.4 The Blossom Algorithm

The Blossom algorithm, developed by Edmonds (1965), solves MWPM by maintaining primal and dual feasibility while growing alternating trees.

Definition. Alternating structures:

- **M-alternating walk:** Path (v_0, v_1, \dots, v_t) where edges alternate between M and $E \setminus M$
- **M-augmenting path:** M-alternating walk with both endpoints unmatched
- **M-blossom:** Odd-length cycle in an M-alternating walk where edges alternate in/out of M

Key Point. Algorithm overview:

1. **Initialization:** Start with empty matching $M = \emptyset$, dual variables $y_v = 0$
2. **Main loop:** While M is not perfect:
 - **Search:** Find M-alternating walks from unmatched vertices
 - **Augment:** If M-augmenting path found, flip edges along path (add unmatched edges, remove matched edges)
 - **Shrink:** If M-blossom found, contract it to a single vertex, update dual variables
 - **Grow:** If no path/blossom found, increase dual variables to make new edges tight, add to search tree
 - **Expand:** When blossom dual variable reaches zero, uncontract it
3. **Termination:** When all vertices are matched, return M

Theorem (Blossom Algorithm Correctness and Complexity). The Blossom algorithm:

1. Maintains primal feasibility, dual feasibility, and complementary slackness throughout
2. Terminates with an optimal MWPM
3. Runs in $O(|V|^3)$ time with careful implementation

Iteration bounds:

- Augmentations: at most $|V|/2$
- Contractions: at most $2|V|$
- Expansions: at most $2|V|$
- Edge additions: at most $3|V|$
- Total: $O(|V|^2)$ iterations, each taking $O(|V|)$ time

Key Point. Why MWPM for quantum codes:

- Surface codes and other topological codes have parity check matrices where each error triggers at most 2 stabilizers
- This structure allows efficient MLD via MWPM
- Blossom algorithm provides polynomial-time optimal decoding
- Practical implementations achieve near-optimal thresholds (10-11% for surface codes)

- Contrast with BP: MWPM finds global optimum but is slower; BP is faster but can get trapped in local minima

7 Quantum Error Correction Basics

7.1 Qubits and Quantum States

Definition. A **qubit** is a quantum two-level system. Its state is written using **ket notation**:

$$|\psi\rangle = \alpha |0\rangle + \beta |1\rangle \quad (64)$$

where:

- $|0\rangle = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and $|1\rangle = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ are the **computational basis states**
- α, β are complex numbers with $|\alpha|^2 + |\beta|^2 = 1$
- The ket symbol $|\cdot\rangle$ is standard notation for quantum states

Common quantum states include:

- $|0\rangle, |1\rangle$ = computational basis
- $|+\rangle = \frac{1}{\sqrt{2}}(|0\rangle + |1\rangle)$ = superposition (plus state)
- $|-\rangle = \frac{1}{\sqrt{2}}(|0\rangle - |1\rangle)$ = superposition (minus state)

7.2 Pauli Operators

Definition. The **Pauli operators** are the fundamental single-qubit error operations:

Symbol	Matrix	Binary repr.	Effect on states
$\mathbb{1}$ (Identity)	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$	(0, 0)	No change
X (bit flip)	$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$	(1, 0)	$ 0\rangle \leftrightarrow 1\rangle$
Z (phase flip)	$\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$	(0, 1)	$ +\rangle \leftrightarrow -\rangle$
$Y = iXZ$	$\begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}$	(1, 1)	Both flips

Table 9: Pauli operators

Key Point. Quantum errors are modeled as random Pauli operators:

- **X errors** = bit flips (like classical errors)
- **Z errors** = phase flips (uniquely quantum, no classical analogue)
- **Y errors** = both (can be written as $Y = iXZ$)

7.3 Binary Representation of Pauli Errors

Definition. An n -qubit Pauli error E can be written in **binary representation**:

$$E \mapsto e_Q = (\mathbf{x}, \mathbf{z}) \quad (65)$$

where:

- $\mathbf{x} = (x_1, \dots, x_n)$ indicates X components ($x_j = 1$ means X error on qubit j)
- $\mathbf{z} = (z_1, \dots, z_n)$ indicates Z components ($z_j = 1$ means Z error on qubit j)

For example, the error $E = X_1 Z_3$ on 3 qubits (X on qubit 1, Z on qubit 3) has binary representation:

$$e_Q = (\mathbf{x}, \mathbf{z}) = ((1, 0, 0), (0, 0, 1)) \quad (66)$$

7.4 CSS Codes

Definition. A **CSS code** (Calderbank-Shor-Steane code) is a quantum error-correcting code with a structure that allows X and Z errors to be corrected independently.

A CSS code is defined by two classical parity check matrices H_X and H_Z satisfying:

$$H_X \cdot H_Z^T = \mathbf{0} \quad (\text{orthogonality constraint}) \quad (67)$$

The combined quantum parity check matrix is:

$$H_{\text{CSS}} = \begin{pmatrix} H_Z & \mathbf{0} \\ \mathbf{0} & H_X \end{pmatrix} \quad (68)$$

Key Point. The orthogonality constraint $H_X \cdot H_Z^T = \mathbf{0}$ ensures that the quantum stabilizers **commute** (a necessary condition for valid quantum codes).

7.5 Syndrome Measurement in CSS Codes

For a CSS code with error $E \mapsto e_Q = (x, z)$:

Definition. The **quantum syndrome** is:

$$s_Q = (s_x, s_z) = (H_Z \cdot x, H_X \cdot z) \quad (69)$$

- $s_x = H_Z \cdot x$ detects X (bit-flip) errors
- $s_z = H_X \cdot z$ detects Z (phase-flip) errors



Two independent classical problems!

Figure 19: CSS codes allow independent X and Z decoding

Key Point. CSS codes allow **independent decoding**:

- Decode X errors using matrix H_Z and syndrome s_x
- Decode Z errors using matrix H_X and syndrome s_z

Each is a classical syndrome decoding problem — so BP can be applied!

7.6 Quantum Code Parameters

Quantum codes use double-bracket notation $[[n, k, d]]$:

- n = number of physical qubits
- k = number of logical qubits encoded
- d = code distance (minimum weight of undetectable errors)

Compare to classical $[n, k, d]$ notation (single brackets).

Definition. A **quantum LDPC (QLDPC) code** is a CSS code where H_{CSS} is sparse.

An (l_Q, q_Q) -QLDPC code has:

- Each column of H_{CSS} has at most l_Q ones
- Each row of H_{CSS} has at most q_Q ones

7.7 The Hypergraph Product Construction

Definition. The **hypergraph product** constructs a quantum CSS code from a classical code.

Given classical code with $m \times n$ parity check matrix H :

$$H_X = (H \otimes \mathbb{1}_n \quad \mathbb{1}_m \otimes H^T) \quad (70)$$

$$H_Z = (\mathbb{1}_n \otimes H \quad H^T \otimes \mathbb{1}_m) \quad (71)$$

Where:

- \otimes = **Kronecker product** (tensor product of matrices)
- $\mathbb{1}_n = n \times n$ identity matrix
- H^T = transpose of H

A well-known example is the **Toric Code**, which is the hypergraph product of the ring code (cyclic repetition code). From a classical $[n, 1, n]$ ring code, we obtain a quantum $[[2n^2, 2, n]]$ Toric code. Its properties include:

- (4, 4)-QLDPC: each stabilizer involves at most 4 qubits
- High threshold (10.3% with optimal decoder)
- Rate $R = \frac{2}{2n^2} \rightarrow 0$ as $n \rightarrow \infty$

8 Results and Performance

8.1 Error Threshold

Definition. The **threshold** p_{th} is the maximum error rate below which the logical error rate decreases with increasing code distance.

- If $p < p_{\text{th}}$: Larger codes \rightarrow exponentially better protection
- If $p > p_{\text{th}}$: Larger codes \rightarrow worse protection (error correction fails)

8.2 Experimental Results

Code Family	BP Only	BP+OSD-0	BP+OSD-CS
Toric	N/A (fails)	$9.2 \pm 0.2\%$	$9.9 \pm 0.2\%$
Semi-topological	N/A (fails)	$9.1 \pm 0.2\%$	$9.7 \pm 0.2\%$
Random QLDPC	$6.5 \pm 0.1\%$	$6.7 \pm 0.1\%$	$7.1 \pm 0.1\%$

Table 10: Observed thresholds from the paper

Key Results for Toric Code

- **BP alone:** Complete failure due to degeneracy (no threshold)
- **BP+OSD-CS:** 9.9% threshold (optimal decoder achieves 10.3%)
- **Improvement:** Combination sweep gains 0.7% over OSD-0
- **Low-error regime:** Exponential suppression of logical errors

8.3 Complexity

Component	Complexity	Notes
BP (per iteration)	$O(n)$	Linear in block length
OSD-0	$O(n^3)$	Dominated by matrix inversion
Combination sweep	$O(\lambda^2)$	$\lambda = 60 \rightarrow 1830$ trials
Total	$O(n^3)$	Practical for moderate n

Table 11: Complexity analysis

9 Summary

9.1 Key Takeaways

1. **Classical BP** computes marginal probabilities via message passing on factor graphs
2. **Quantum codes suffer from degeneracy**: multiple errors can produce the same syndrome, causing BP to output invalid solutions (split beliefs)
3. **OSD resolves degeneracy** by selecting a basis guided by BP soft decisions, then solving via matrix inversion to get a unique valid solution
4. **Combination sweep** efficiently improves OSD-0 by testing low-weight configurations of the remainder bits
5. **BP+OSD is general**: works for Toric codes, semi-topological codes, and random QLDPC codes, achieving near-optimal thresholds

9.2 The BP+OSD Recipe

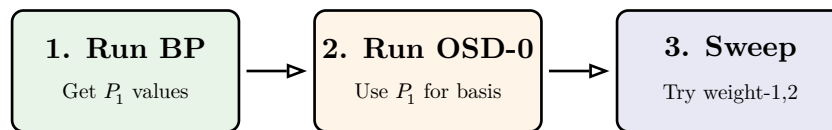


Figure 20: BP+OSD in three steps

10 References

Bibliography

- [1] D. J. MacKay, *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003.
- [2] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, 1988.
- [3] T. Richardson and R. Urbanke, “The Capacity of Low-Density Parity-Check Codes Under Message-Passing Decoding,” *IEEE Transactions on Information Theory*, vol. 47, no. 2, pp. 599–618, 2001.
- [4] J. S. Yedidia, W. T. Freeman, and Y. Weiss, “Understanding Belief Propagation and Its Generalizations,” *Exploring Artificial Intelligence in the New Millennium*. pp. 239–269, 2003.
- [5] L. Dolecek, Z. Zhang, V. Anantharam, M. J. Wainwright, and B. Nikolić, “Analysis of Absorbing Sets and Fully Absorbing Sets of Array-Based LDPC Codes,” *IEEE Transactions on Information Theory*, vol. 56, no. 1, pp. 181–201, 2010.
- [6] J. Chen, A. Dholakia, E. Eleftheriou, M. P. Fossorier, and X.-Y. Hu, “Reduced-Complexity Decoding of LDPC Codes,” *IEEE Transactions on Communications*, vol. 53, no. 8, pp. 1288–1299, 2005.
- [7] M. P. Fossorier and S. Lin, “Soft-Decision Decoding of Linear Block Codes Based on Ordered Statistics,” *IEEE Transactions on Information Theory*, vol. 41, no. 5, pp. 1379–1396, 1995.
- [8] M. P. Fossorier and S. Lin, “Computationally Efficient Soft-Decision Decoding of Linear Block Codes Based on Ordered Statistics,” *IEEE Transactions on Information Theory*, vol. 42, no. 3, pp. 738–750, 1996.
- [9] J. Roffe, D. R. White, S. Burton, and E. Campbell, “Decoding Across the Quantum Low-Density Parity-Check Code Landscape,” *Physical Review Research*, vol. 2, no. 4, p. 43423, 2020.
- [10] R. W. Hamming, “Error Detecting and Error Correcting Codes,” *The Bell System Technical Journal*, vol. 29, no. 2, pp. 147–160, 1950.
- [11] G. D. Forney, “Generalized Minimum Distance Decoding,” *IEEE Transactions on Information Theory*, vol. 12, no. 2, pp. 125–131, 1966.
- [12] J. Hagenauer, E. Offer, and L. Papke, “Iterative Decoding of Binary Block and Convolutional Codes,” *IEEE Transactions on Information Theory*, vol. 42, no. 2, pp. 429–445, 1996.
- [13] C. Yue, M. Shirvanimoghaddam, Y. Li, and B. Vucetic, “A Revisit to Ordered Statistics Decoding: Distance Distribution and Decoding Rules,” *IEEE Transactions on Information Theory*, vol. 67, no. 7, pp. 4288–4337, 2021.
- [14] A. Montanari, “Belief Propagation.” [Online]. Available: <https://web.stanford.edu/~montanar/RESEARCH/BOOK/partD.pdf>
- [15] T. Richardson and R. Urbanke, *Modern Coding Theory*. Cambridge University Press, 2008.

- [16] S. Tatikonda and M. I. Jordan, “Loopy Belief Propagation and Gibbs Measures,” in *Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence*, 2002, pp. 493–500.
- [17] A. T. Ihler, J. W. Fisher III, and A. S. Willsky, “Loopy Belief Propagation: Convergence and Effects of Message Errors,” *Journal of Machine Learning Research*, vol. 6, pp. 905–936, 2005.

End of Lecture Note