

# New Methods for Error Correction Codes with Deep Learning

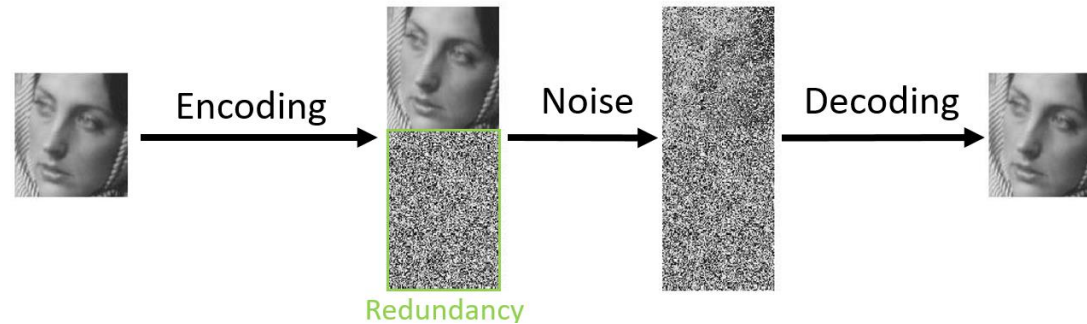
Yoni Choukroun

School of Computer Science, Tel Aviv University

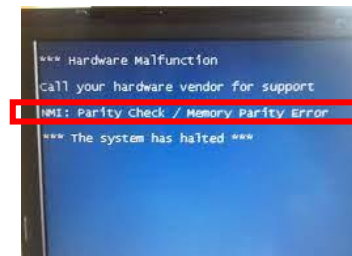
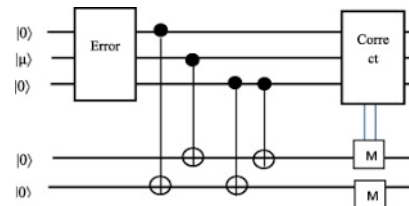
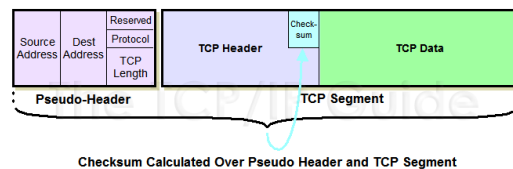
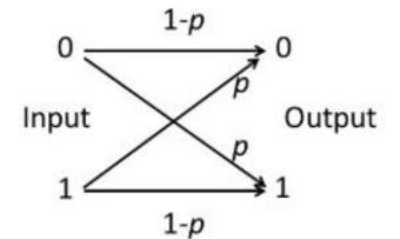


# Error Correction Code: Motivation

- Error-correcting codes (ECC) are clever ways of **representing data** so that one can **recover** the original information even if parts of it are corrupted.
- The basic idea is to **judiciously add redundancy (encoding)** so that the original information can be recovered (**decoding**) even when parts of the data have been corrupted (**noise**).



- E.g., assume we want to transmit **one single bit  $b$**  and the probability of bit flip is  $p$ .
  - Using the **3-repetition code** (sending  $[b, b, b]$ ), the probability of error is **reduced** to  $\sim 3p^2$  (but transmission rate  $1/3$ ).
- ECC is present in all types of communication (transmission over *space*) and storage (transmission over *time*)



# Error Correction Code: Setting

➤ **Goal:** allow reliable data transmission over a noisy communication channel.

❖ For a Binary Linear Block Code  $(n, k)$

• A desired binary message  $m \sim \text{Bern}^k(0.5)$  is encoded to a redundant codeword

▪  $x \equiv Gm$ . Multiplication over  $GF(2)$ , i.e.,  $x = \text{mod}(Gm, 2)$ ,  $G \in \{0, 1\}^{n \times k}$

▪ E.g., the natural **structure** of images allowing their denoising is here inserted **artificially** into the data to allow optimal decoding.

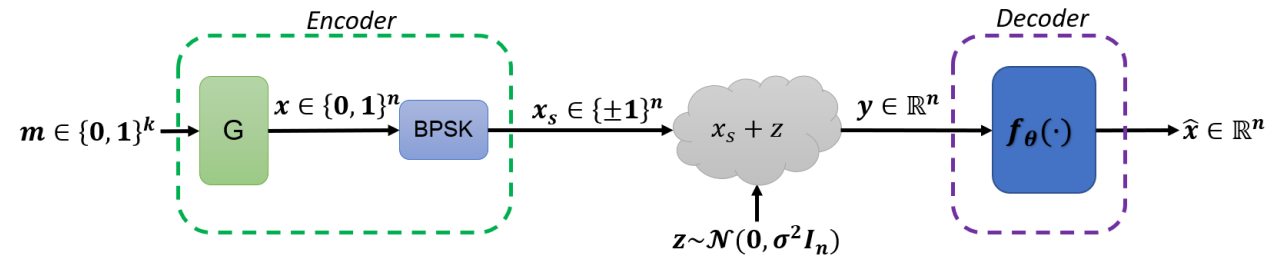
• The codeword  $x$  is modulated (via BPSK) for transmission

▪  $x_s = 1 - 2x$ .

• The channel adds (AWGN) noise  $z$

▪  $y = x_s + z$ .

• The (parameterized) decoder  $f_\theta(y)$  aims at retrieving the original codeword  $x$  from  $y$ .



• For the 3-repetition code (3,1):

1.  $m = 1 \in \{0,1\}$
2.  $G = (1, 1, 1)^T$
3.  $x = x_s = Gm = (1,1,1)$
4.  $y = x_s \oplus z = (1, -1, -1)$
5.  $\hat{x} = f(y) = \left\lfloor \frac{\sum_i y_i}{3} \right\rfloor = (1,1,1)$

• For the (7,4) Hamming code:

$$x = Gm \equiv_2 \begin{pmatrix} 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 3 \\ 1 \\ 2 \\ 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}$$

$$y = x \oplus z = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

# Error Correction Code: Setting

## ➤ How can I detect that an error occurred?

- Check if this is one of the  $2^k$  codewords
- The **parity check** matrix  $H \in \{0, 1\}^{(n-k) \times n}$  kernel defines the codewords
  - $G^T H \equiv 0 \Rightarrow Hx \equiv 0$ .

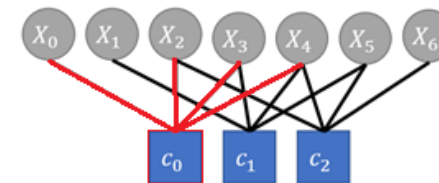
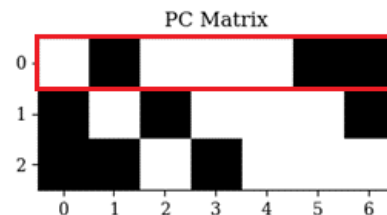
$$H = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} \quad Hx \equiv_2 \begin{pmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \\ 2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

- The parity check equations allow parity check **errors discovery**
  - $s = H(x \oplus z_b) = Hx \oplus Hz_b \equiv Hz_b \in \{0, 1\}^{n-k}$ .
  - This **binary** vector of parity check errors is called the **SYNDROME**

$$s = Hy \equiv_2 \begin{pmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 3 \\ 4 \\ 3 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

- The **Tanner graph** is the (factor) graph representation of  $H$  (edge for 1 in each column)

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



# Error Correction Code: Background

- Given a code (i.e.,  $\mathbf{H}$ ), the best possible decoder is defined by the **NP-hard maximum likelihood decoder** (search for the closest codeword among the  $2^k$ ).
- Must have been done during the past 70 years in *information theory* for the design of **provably (asymptotical) optimal codes** and their efficient ( $\leq$  polynomial time) decoding.
- Shannon's channel coding theorem (1948) shows that, given a noisy channel with capacity  $\mathbf{C}$ , if the information is transmitted at rate  $\mathbf{R} = \mathbf{k}/\mathbf{n}$  such that  $\mathbf{R} < \mathbf{C}$ , then **there exists** a coding scheme that guarantees **negligible probability** of miscommunication.
- Recently, neural networks-based decoders (aka **neural decoders**) have shown promising results in this field.

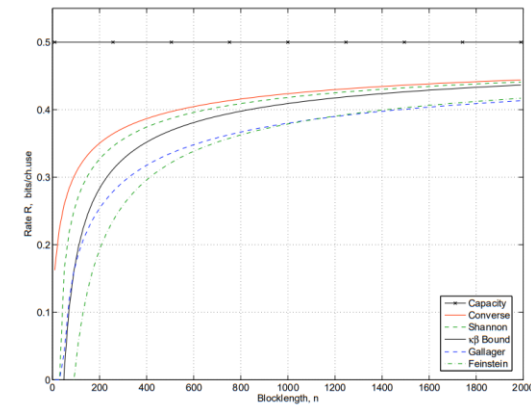


Fig. 6. Bounds for the AWGN channel,  $S/N = 0$  dB,  $\epsilon = 10^{-5}$ .

# Neural Decoders

❖ *Two main families of neural decoders:*

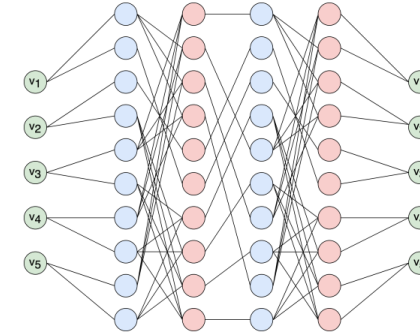
- **Model-based** decoders implement *augmented* parameterized versions of the classical *Belief Propagation* decoder built upon the **Tanner graph** (graph representation of  $H$ ).

- Pros:

- Invariant to the codeword (robust to codewords overfitting)
- Built on **iterative** (message-passing) legacy methods

- Cons:

- Suffers from **heavy and restrictive inductive bias**.
- **Improvement vanishes** as the code length and the number of iterations increase



- **Model-free** decoders employ **general** types of neural network architectures (e.g., MLP, RNN)

- Pros:

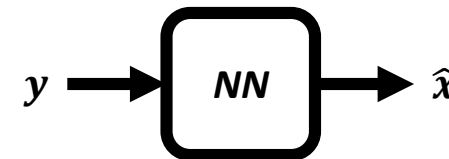
- Total **freedom** in model design

- Cons:

1. **Overfitting** (*exponential number of codewords for training*)
2. **Difficulties for the model to learn the code.**
3. **Lack Iterative formulation**

- Cons in Common :

4. **Lack Code invariance** (*one decoder for each code*)
5. **Not adapted to modern ECC settings** (*e.g., quantum computing*)
6. **Cannot design/learn the code**





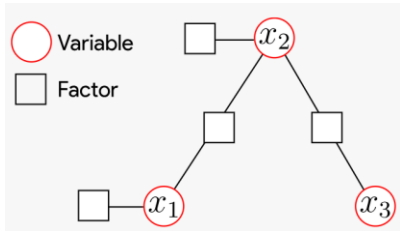
# Model-*BASED* Decoders



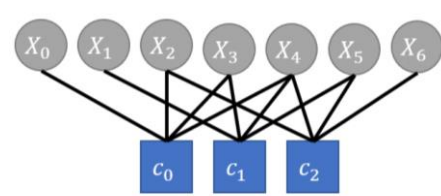


# Belief Propagation (BP)

- A **factor graph** is a (bipartite) representation of a discrete probability distribution that takes advantage of conditional independencies between variables to make the representation *more compact*
  - The Hammersley-Clifford theorem tells us that any positive joint distribution can be represented as a product of factors



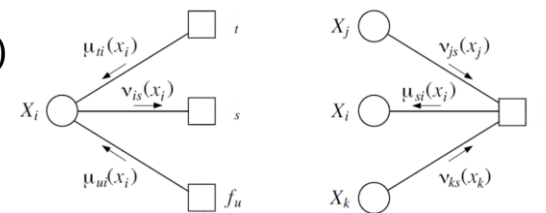
$$p(x_1, \dots, x_n) = \frac{\prod_{i=1}^{n-k} \phi_i(X_{S_i})}{Z}$$



$$p(x_1, x_2, x_3) \propto \phi_1(x_1) \phi_2(x_2) \phi_3(x_1, x_2) \phi_4(x_2, x_3)$$

- Belief-Propagation** (Sum-product algorithm) allows efficient marginal inference  $p(x_i)$  via variable elimination as *message-passing* over the factor graph.

$$v_{x_i \rightarrow \phi_k}(x_i) = \prod_{\phi_k \in N(i) \setminus \{\phi_k\}} \mu_{\phi_k \rightarrow x_i}(x_i) \quad \text{and} \quad \mu_{\phi_k \rightarrow x_i}(x_i) = \sum_{X_{S_k \setminus \{i\}}} \phi_s(X_{S_k}) \prod_{j \in S_k \setminus \{i\}} v_{X_j \rightarrow \phi_k}(x_j)$$



- For non-tree Bayesian graphs, the inference is not tractable. Thus, (**loopy**) belief-propagation is performed **iteratively** (until convergence of the beliefs to a local minimum)

$$b_i(x_i) = P(x_i) \propto \phi_i(x_i) \prod_{\phi_k \in N(i)} \mu_{\phi_k \rightarrow x_i}^T(x_i)$$

- Derivation:

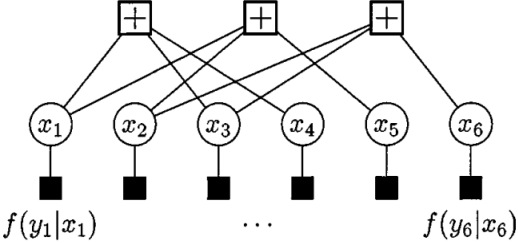
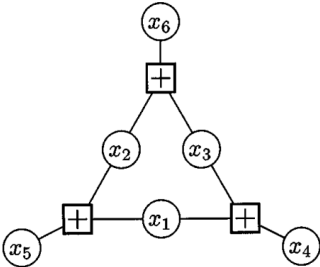
- Write the posterior factor graph  $p(X) = \prod_i \phi_i(X_{S_i})/Z$  and the **tree graph distribution**  $q(X) = \prod_i b_i(x_i)^{1-d_i} \prod_j b_{S_j}(X_{S_j})$
- Minimize with respect to the (constrained) beliefs  $b_i, b_{S_j}$ :  $KL(q||p) = H_q(X) - \sum_i \mathbb{E}_q \log(\phi_i(X_{S_i}))$

# Belief Propagation Decoding

- Early (Gallager) application of BP to ECC with a Posteriori distribution

- $P(x|y) \sim P(y|x)P(x) = \prod_i f(y_i|x_i)P(\bigcup_k \bigoplus_{k \in N(j)} x_k = 0)$
- $P(\bigcup_k \bigoplus_{k \in N(j)} x_k = 0) = \prod_k P(\bigoplus_{k \in N(j)} x_k = 0)$
- $\phi(X_{S_k}) = P(\bigoplus_{k \in N(j)} x_k = 0) = \frac{1}{2} (1 + \prod_{k \in N(j) \setminus i} 2(v_{x_k \rightarrow c_j}^T(a)))$ 
  - Since the variables are independent

$$H = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$



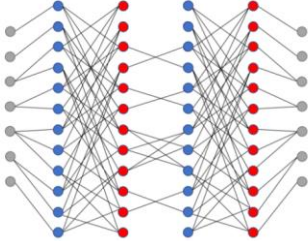
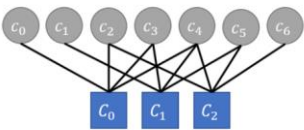
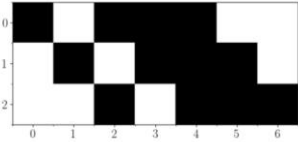
- Belief-propagation *decoding* is generally represented as a **Trellis** graph unrolling of the factor/Tanner graph (*log-likelihoods*).

$$L_v = \log \left( \frac{\Pr(c_v = 1|y_v)}{\Pr(c_v = 0|y_v)} \right)$$

$$x_e^{2k+1} = x_{(c,v)}^{2k+1} = L_v + \sum_{e' \in N(v) \setminus \{(c,v)\}} x_{e'}^{2k}$$

$$x_e^{2k} = x_{(c,v)}^{2k} = 2 \operatorname{arctanh} \left( \prod_{e' \in N(c) \setminus \{(c,v)\}} \tanh \left( \frac{x_{e'}^{2k-1}}{2} \right) \right)$$

$$o_v = L_v + \sum_{e' \in N(v)} x_{e'}^{2L}$$



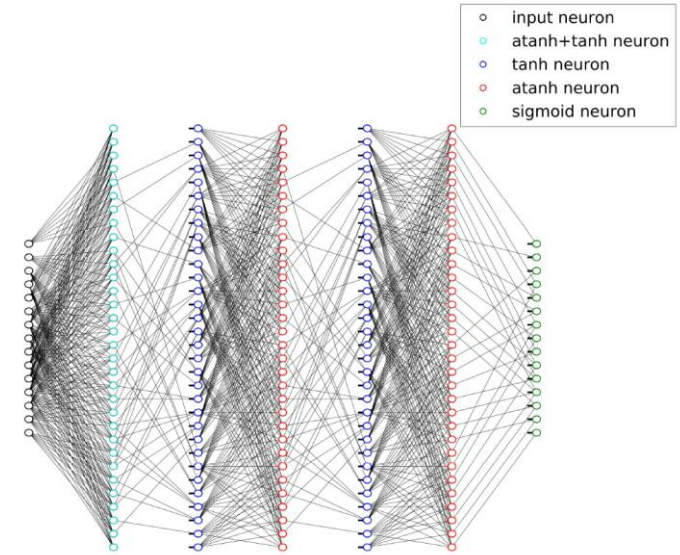
**How can we augment BP decoders?**

# Learning to Decode Linear Codes Using Deep Learning, Nachmani et al. Allerton 2016

- ▶ Input - LLR

$$l_v = \log \frac{\Pr(C_v = 1|y_v)}{\Pr(C_v = 0|y_v)}$$

$y_v$  is the channel output corresponding to the  $v$ th codebit,  $C_v$ .



- ▶ For odd  $i$  and  $e = (v, c)$  -

$$x_{i,e=(v,c)} = l_v + \sum_{e'=(v,c'), c' \neq c} x_{i-1,e'}$$

- ▶ For even  $i$  and  $e = (v, c)$  -

$$x_{i,e=(v,c)} = 2 \tanh^{-1} \left( \prod_{e'=(v',c), v' \neq v} \tanh \left( \frac{x_{i-1,e'}}{2} \right) \right)$$

- ▶ The final  $v$ th output -

$$o_v = l_v + \sum_{e'=(v,c')} x_{2L,e'}$$

- ▶ For odd  $i$  and  $e = (v, c)$  -

$$x_{i,e=(v,c)} = 2 \tanh^{-1} \left( \prod_{e'=(v',c), v' \neq v} x_{i-1,e'} \right)$$

- ▶ For even  $i$  and  $e = (v, c)$  -

$$x_{i,e=(v,c)} = \tanh \left( \frac{1}{2} \left( w_{i,v} l_v + \sum_{e'=(v,c'), c' \neq c} w_{i,e,e'} x_{i-1,e'} \right) \right)$$

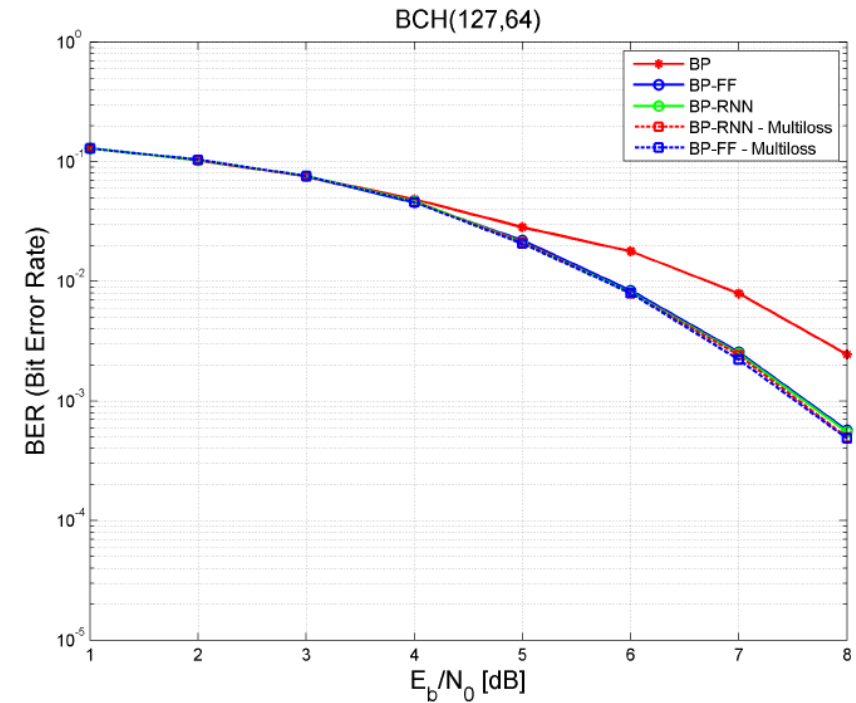
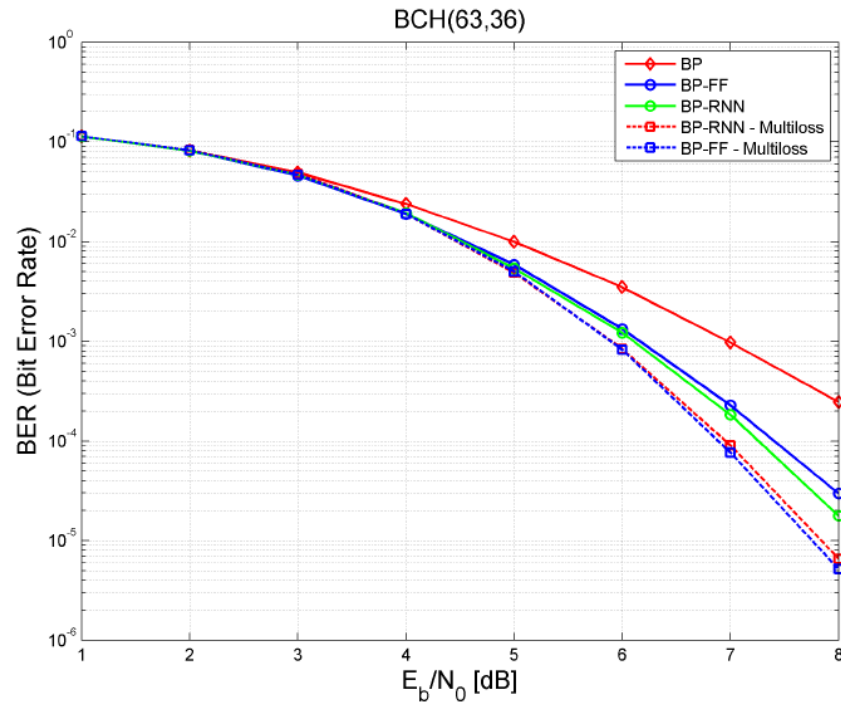
- ▶ The final  $v$ th output -

$$o_v = \sigma \left( w_{2L+1,v} l_v + \sum_{e'=(v,c')} w_{2L+1,v,e'} x_{2L,e'} \right)$$

where  $\sigma(x) \equiv (1 + e^{-x})^{-1}$

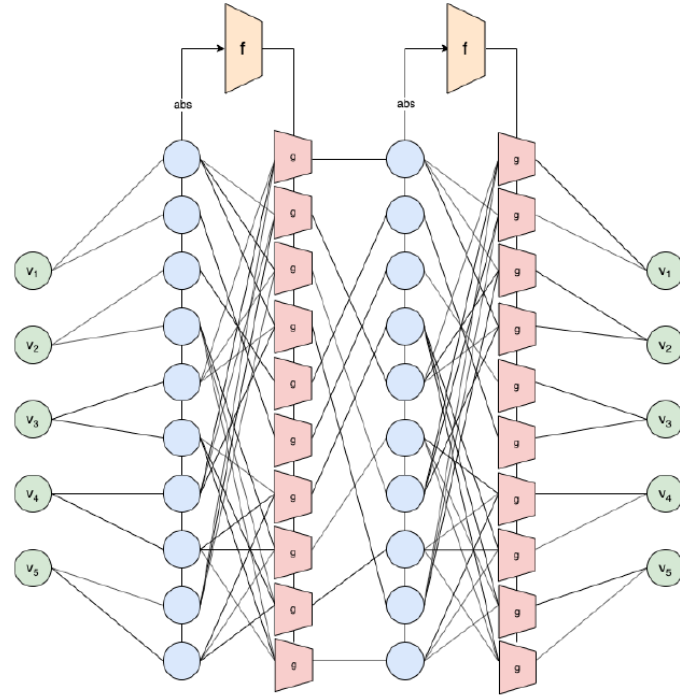
$$L(o, y) = -\frac{1}{N} \sum_{v=1}^N y_v \log(o_v) + (1 - y_v) \log(1 - o_v)$$

# Learning to Decode Linear Codes Using Deep Learning, Nachmani et al.



- RNN == Shared parameters
- Multi loss  $L(o, y) = -\frac{1}{N} \sum_{t=1}^T \sum_{v=1}^N y_v \log(o_{v,t}) + (1 - y_v) \log(1 - o_{v,t})$  for better gradient update (vanishing).
- Ensemble can also be applied

# Hyper-graph-network decoders for block codes, Nachmani and Wolf. Neurips19

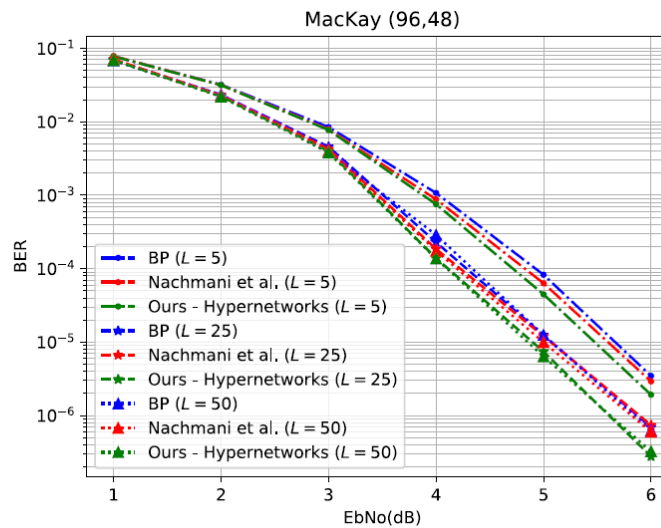
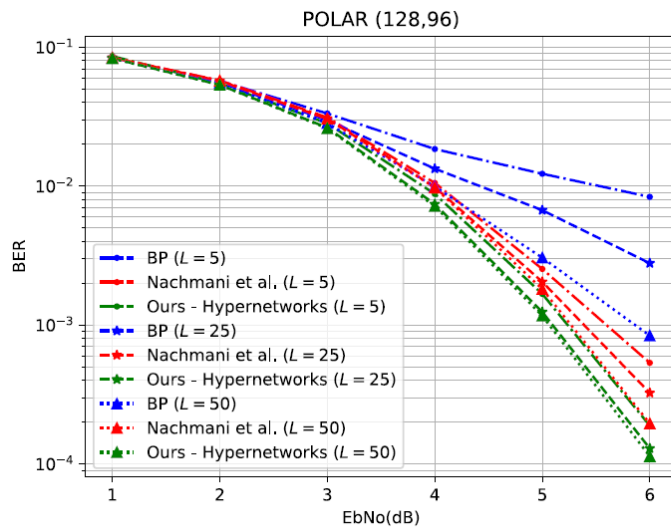
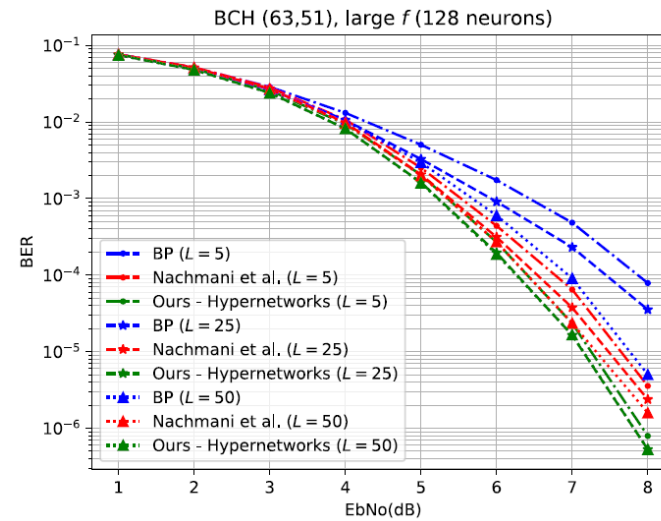
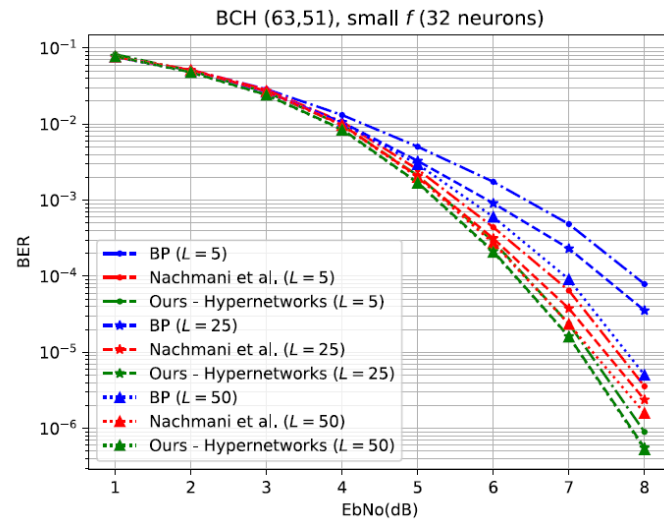


$$\theta_g^j = f(|x^{j-1}|, \theta_f)$$

$$x_e^j = x_{(c,v)}^j = g(l_v, x_{N(v,c)}^{j-1}, \theta_g^j)$$

- Replace arctanh with first order approximation for stable training
- The decoder maintains symmetry condition for zero codeword training (no overfitting)

# Hyper-graph-network decoders for block codes, Nachmani and Wolf. Neurips19



# Model-FREE Decoders



**How can we remain robust to overfitting?**

# Solving *Model Free* Overfitting

- We want  $h(\cdot)$  such that  $P(h(y)|x) = P(h(y))$

- Multiplicative noise**

- The multiplicative noise  $\tilde{z}$  is an **equivalent** statistical model to the true physical additive one  $z$

$$y = x_s + z = x_s \odot \tilde{z} \Rightarrow \tilde{z} = y \odot x_s$$

- Preprocessing for invariance**

- We can remain **invariant to the transmitted codeword** by processing other measures of  $y$

- The **Magnitude**:  $|y|$

- $P(|y||x) = P(|x_s \tilde{z}| |x) = P(|x_s \tilde{z}| |x) = P(|x_s| |\tilde{z}| |x) = P(|\tilde{z}| |x) = P(|\tilde{z}|)$

- The **Syndrome**:  $s(y) = H \text{bin}(\text{sign}(y))$

- $P(s(y)|x) = P(Hx + H z|x) = P(H z)$

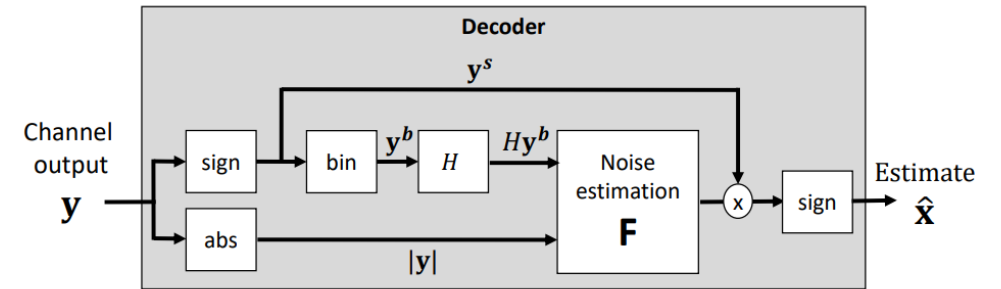
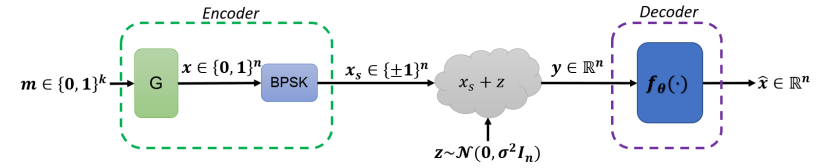
- Extends classical *syndrome decoding*

- Result**

- Even with loss of information

- Does not involve any intrinsic performance penalty (in terms of BER and MSE)

- Guarantees the generalization of performance obtained during training.

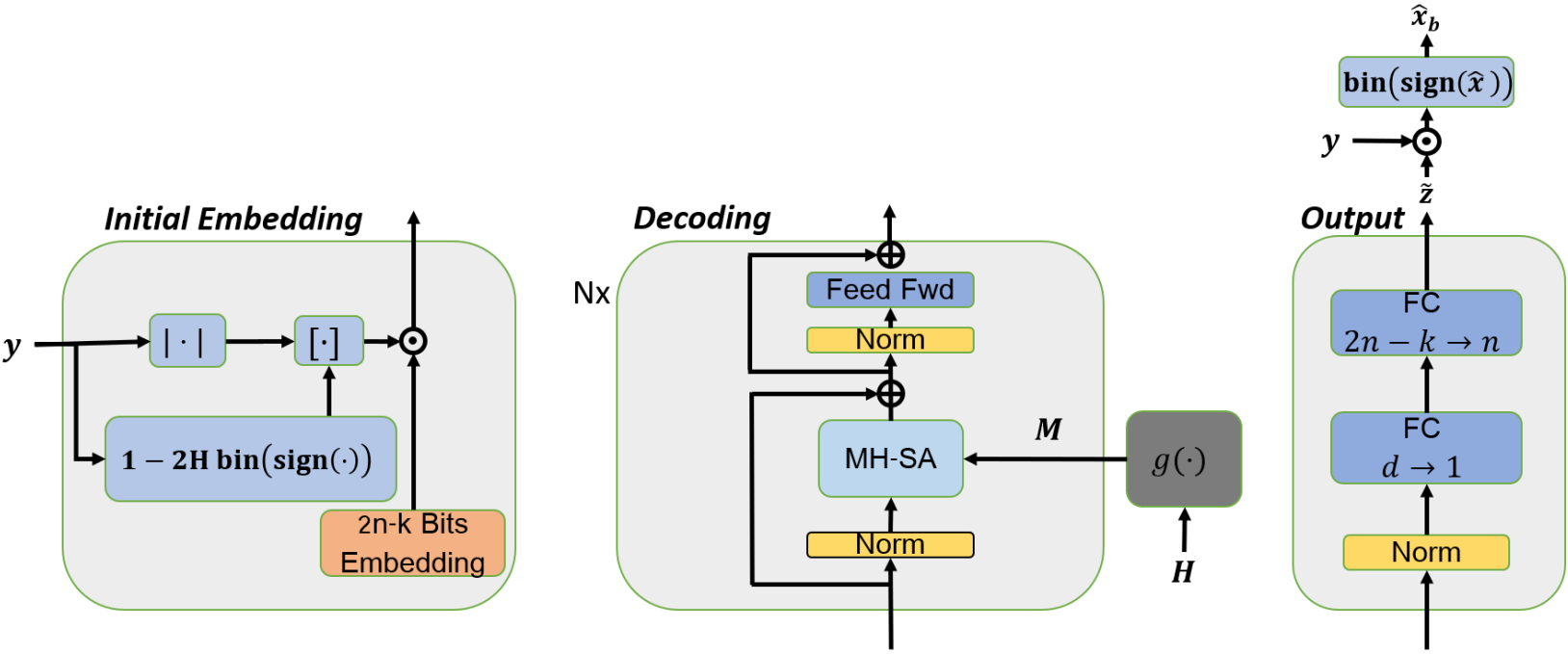


❖ The prediction  $\hat{x} = f_\theta(y)$  is now defined by  $f_\theta: \mathbb{R}_+^n \times \{0, 1\}^{n-k} \rightarrow \mathbb{R}^n$  such that  $f_\theta([|y|, s(y)]) \approx \tilde{z}$



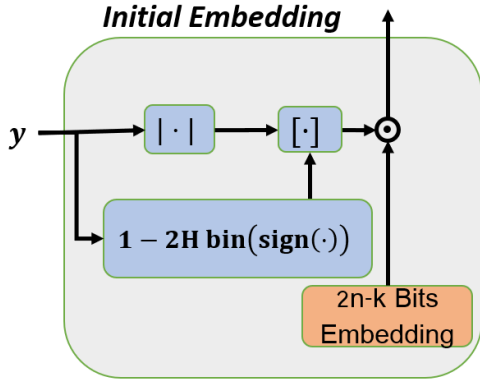
***How can we insert information about the code  
into the model free solutions?***

# Error Correction Code *Transformer* (ECCT) Neurips22



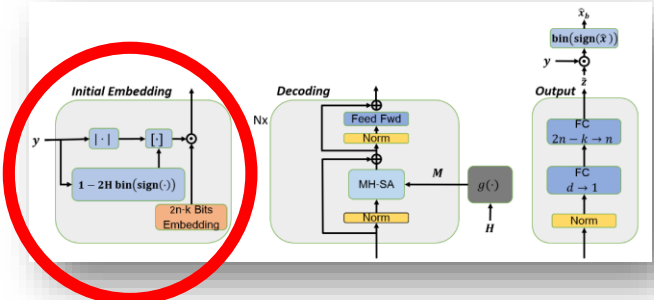
# Positional Reliability Encoding

- Regular architectures (FC/RNN) lose every initial/positional information along the layers.
- The one-hot **high dimensional** embedding is **modulated** by the **code invariant** magnitude and syndrome values.  
 $\Rightarrow$  *Less reliable elements* (i.e., low magnitude) *collapse to the origin* ( $\leq 0$ ).
- Now the input is of dimension  $(n + (n - k)) \times d = (2n - k) \times d$



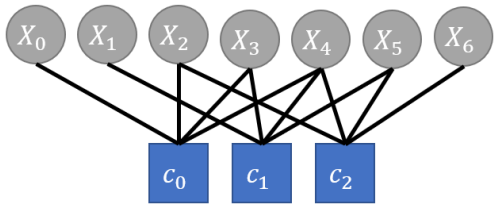
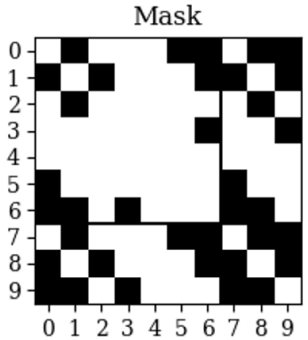
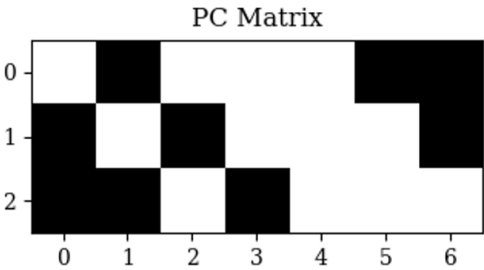
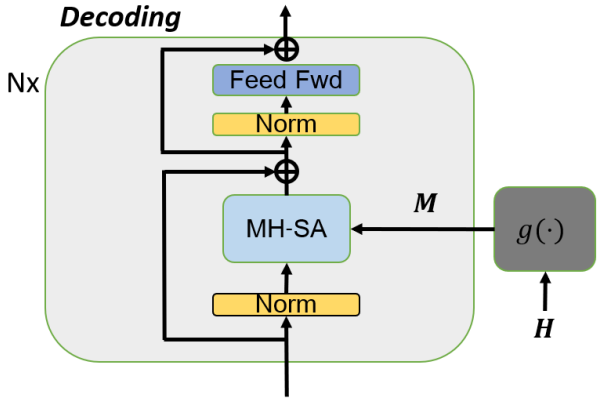
$$\phi_i = \begin{cases} |y_i|W_i, & \text{if } i \leq n \\ (1 - 2(s(y))_{i-n+1})W_i, & \text{otherwise} \end{cases}$$

$$\Phi = (h(y) \cdot \mathbf{1}_d^T) \odot W$$



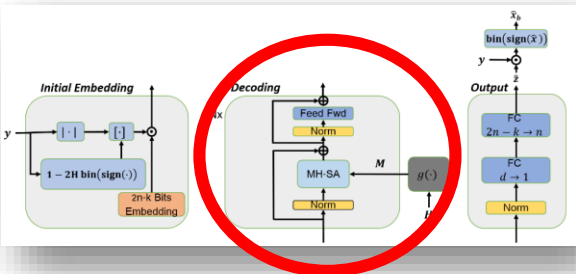
# Code Aware Self-Attention

- Decoding requires the cross analysis between elements.
- We propose to **insert the code** via an **adapted masked self-attention**.
- The proposed mask can be seen as the **sparse adjacency matrix of the Tanner graph** extended to a **two rings connectivity** for simultaneous cross analysis.



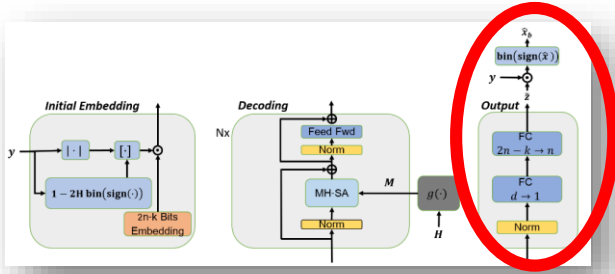
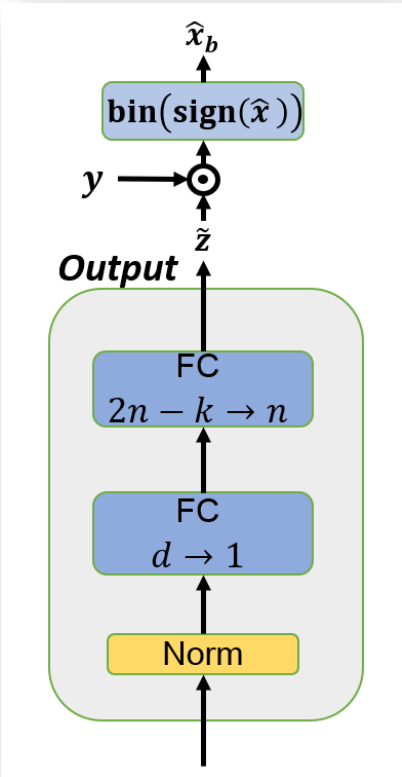
$$A_H(Q, K, V) = \text{Softmax}\left(\frac{QK^T + g(H)}{\sqrt{d}}\right)V.$$

$$g(H) : \{0, 1\}^{(n-k) \times k} \rightarrow \{-\infty, 0\}^{2n-k \times 2n-k}$$



# Noise Prediction Module

In order to **predict the noise**, the embeddings are shrunk to a **one-dimensional scalar** representation and further reduced to **the code length**.

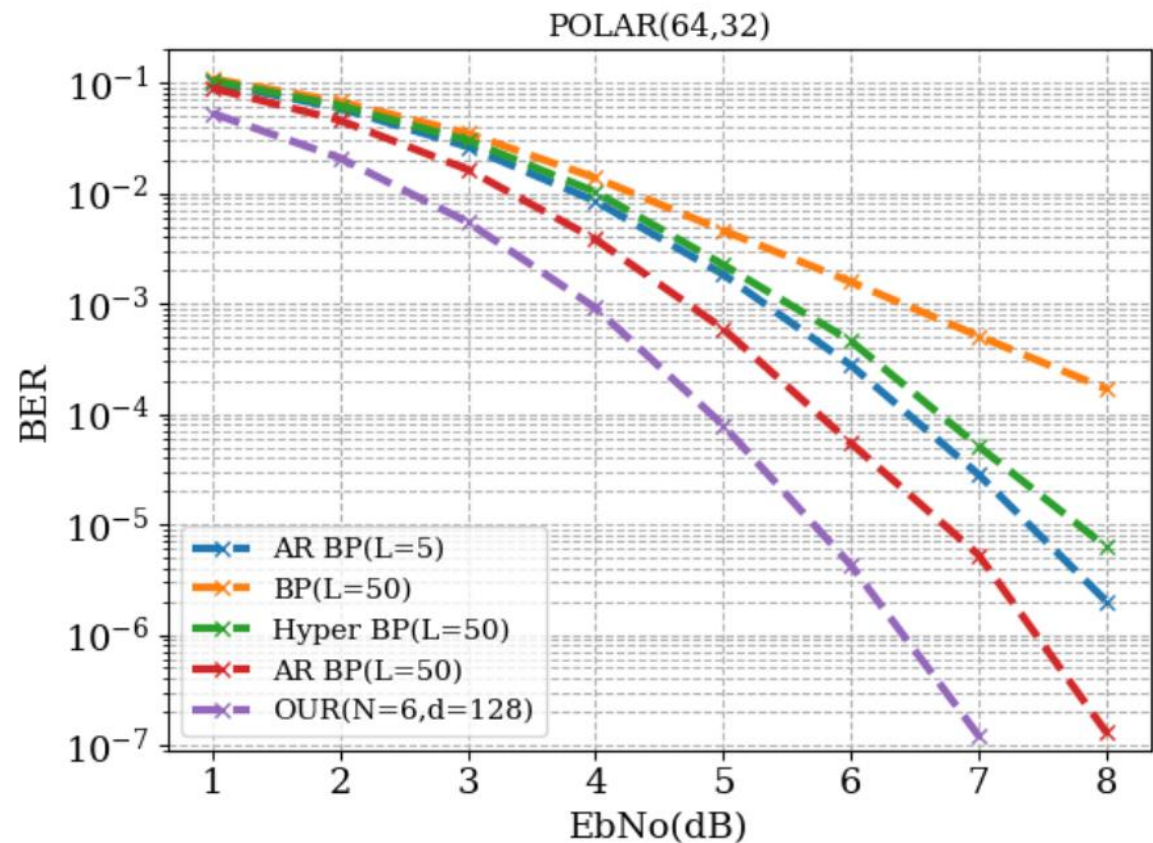




# Experiments

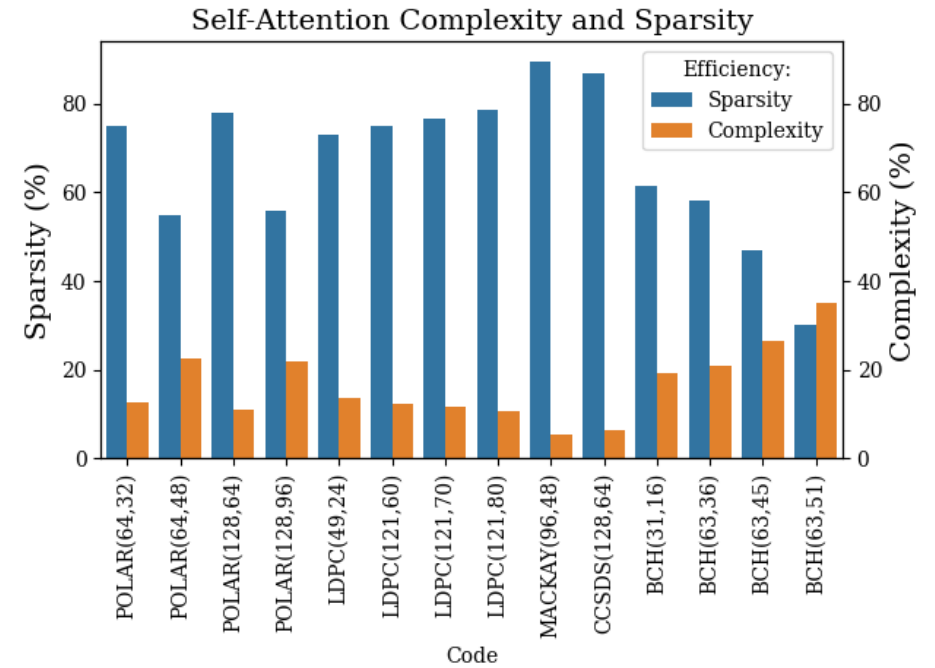
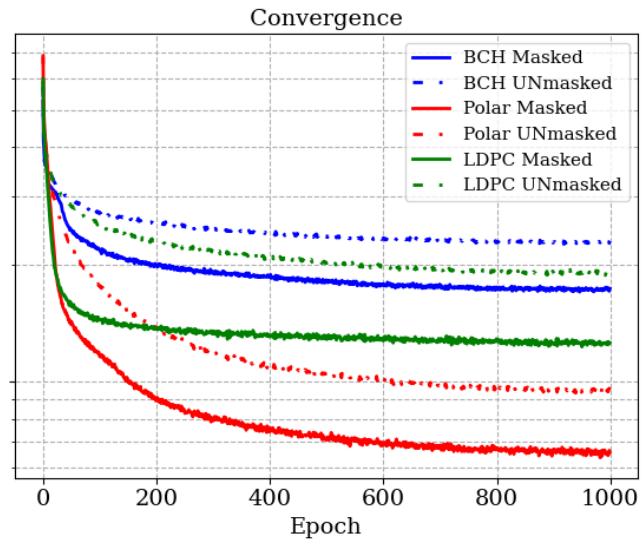
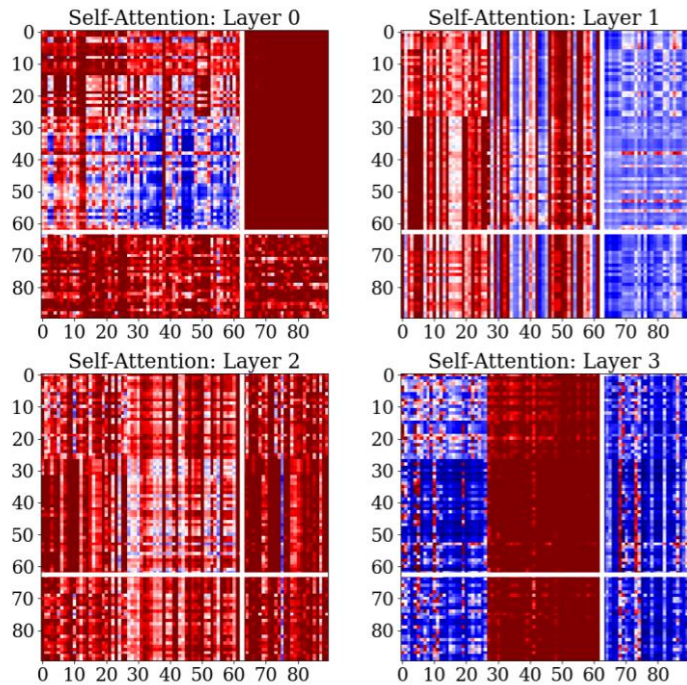
Method	BP[27]			Hyp BP[24]			AR BP[25]			Ours N=2			Ours N=6			Ours N=10														
	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6												
Polar(64,32)	3.52	4.04	4.48	4.25	5.49	7.02	4.77	6.30	8.19	4.27	5.44	6.95	5.71	7.63	9.94	4.57	5.86	7.50	6.48	8.60	11.43	4.87	6.2	7.93	6.99	9.44	12.32			
Polar(64,48)	4.15	4.68	5.31	4.91	6.48	8.41	5.25	6.96	9.00	4.92	6.46	8.41	5.82	7.81	10.24	5.14	6.78	8.9	6.15	8.20	10.86	5.36	7.12	9.39	6.36	8.46	11.09			
Polar(128,64)	3.38	3.80	4.15	3.89	5.18	6.94	4.02	5.48	7.55	3.51	4.52	5.93	4.47	6.34	8.89	3.83	5.16	7.04	5.12	7.36	10.48	4.04	5.52	7.62	5.92	8.64	12.18			
Polar(128,86)	3.80	4.19	4.62	4.57	6.18	8.27	4.81	6.57	9.04	4.30	5.58	7.34	5.36	7.45	10.22	4.49	5.90	7.75	5.75	8.16	11.29	4.75	6.25	8.29	6.31	9.01	12.45			
Polar(128,96)	3.99	4.41	4.78	4.73	6.39	8.57	4.92	6.73	9.30	4.56	5.98	7.93	5.39	7.62	10.45	4.69	6.20	8.30	5.88	8.33	11.49	4.88	6.58	8.93	6.31	9.12	12.47			
LDPC(49,24)	5.30	7.28	9.88	5.76	7.90	11.17	6.05	8.13	11.68	4.51	6.07	8.11	5.74	8.13	11.30	4.58	6.18	8.46	5.91	8.42	11.90	4.71	6.38	8.73	6.13	8.71	12.10	6.68	9.53	13.30
LDPC(121,60)	4.82	7.21	10.87	5.22	8.29	13.00	5.22	8.31	13.07	3.88	5.51	8.06	4.98	7.91	12.70	3.89	5.55	8.16	5.02	7.94	12.72	3.93	5.66	8.51	5.17	8.31	13.30	5.73	9.35	15.01

- $E_b/N_0 \sim SNR$
- Code( $n, k$ )



- Our method **surpasses** every existing neural decoder by very large margins (even with shallow ECCT) and at a **fraction of the complexity** of the previous SOTA method.

# Analysis



- **Self-attention maps:** first layers have higher self-attention values in the **syndrome** part for the parity-check analysis to finally focus on the **information** bits part.

- **Impact of masking:** Masking improves the performance by orders of magnitude demonstrating the importance of **integrating code** information.

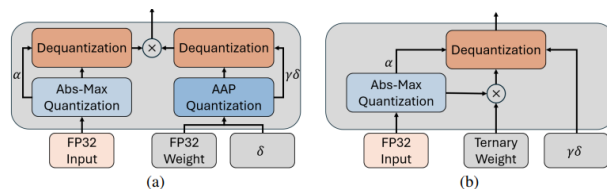
- **Sparsity and complexity ratio:** Sparsity and self-attention complexity can reach up to 80% and as low as 5, respectively.

“But it’s not very efficient...”

# Accelerating Error Correction Code Transformers

- Accelerate and improve decoding via

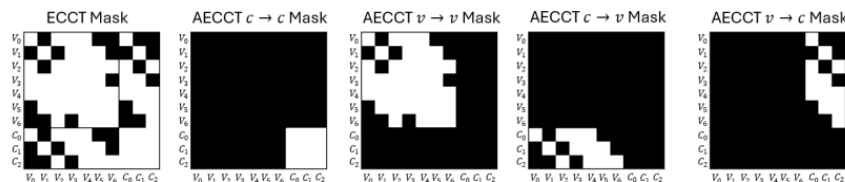
- Adaptive absolute percentile **ternary (0,1,-1) quantization** (90% compression and >224x less energy)



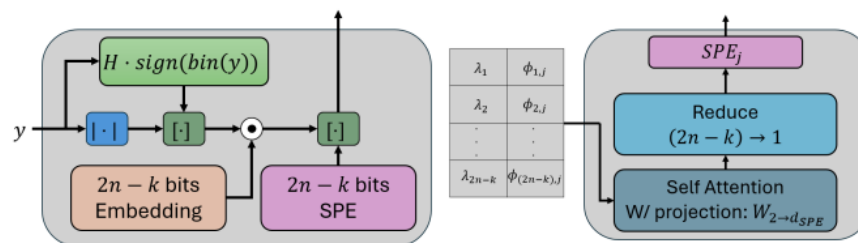
$$\text{AAPLinear}(x, W, b) = \text{Quant}(x) \cdot \text{Ternary}(W) \cdot \frac{\gamma\delta\alpha}{Q_b} + b$$

Figure 1: AAP Linear Layer: (a) QAT: Training with quantization noise; (b) Inference: Matrix multiplication using only integer additions with fixed ternary weights and fixed weight scale.

- Head partitioning self attention

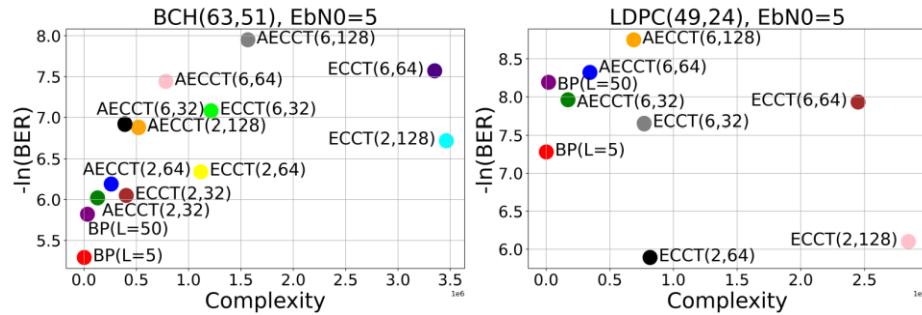


- Positional encoding of the tanner graph

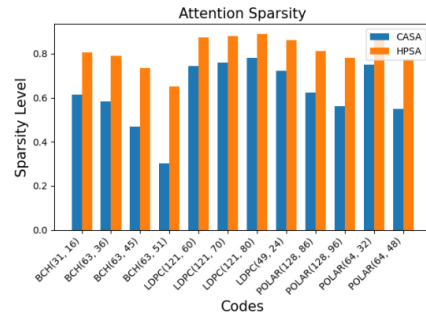


# Accelerating Error Correction Code Transformers

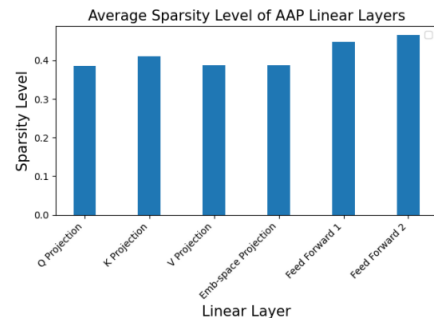
- Matches/improve over ECCT
- Close to BP's complexity



- Increase sparsity masking



- Linear Layers are extremely sparse



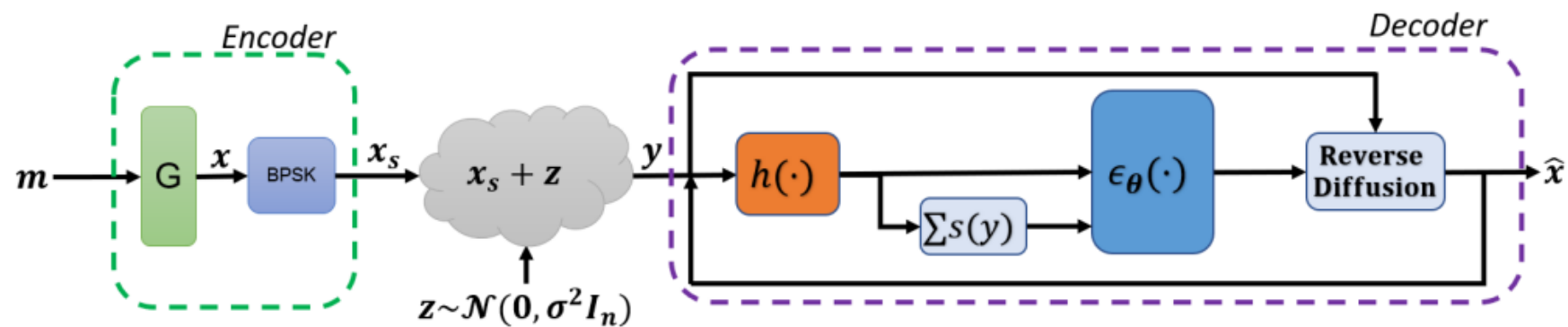
Method	BP			ECCT $N = 6$			AECCT $N = 6$			ECCT $N = 10$			AECCT $N = 10$		
	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
Polar(64,48)	3.52	4.04	4.48	6.36	8.46	11.09	<b>6.43</b>	<b>8.54</b>	<b>11.12</b>	6.43	8.40	11.00	<b>6.54</b>	<b>8.51</b>	<b>11.12</b>
Polar(128,86)	3.80	4.19	4.62	<b>6.31</b>	<b>9.01</b>	<b>12.45</b>	6.04	8.56	11.81	7.26	10.60	<b>14.80</b>	<b>7.28</b>	10.60	14.59
Polar(128,96)	3.99	4.41	4.78	<b>6.31</b>	<b>9.12</b>	<b>12.47</b>	6.11	8.81	12.15	<b>6.85</b>	<b>9.78</b>	12.90	6.79	9.68	<b>12.93</b>
LDPC(49,24)	5.30	7.28	9.88	5.79	8.13	11.40	<b>6.10</b>	<b>8.65</b>	<b>12.34</b>	6.35	9.01	12.43	<b>6.67</b>	<b>9.35</b>	<b>13.56</b>
LDPC(121,60)	4.82	7.21	10.87	5.01	7.99	12.78	<b>5.17</b>	<b>8.32</b>	<b>13.40</b>	5.51	8.89	14.51	<b>5.71</b>	<b>9.31</b>	<b>14.90</b>
LDPC(121,70)	5.88	8.76	13.04	6.19	9.89	15.58	<b>6.38</b>	<b>10.1</b>	<b>16.01</b>	6.86	11.02	16.85	<b>7.05</b>	<b>11.40</b>	<b>17.30</b>
LDPC(121,80)	6.66	9.82	13.98	7.07	10.96	16.25	<b>7.27</b>	<b>11.50</b>	<b>16.90</b>	7.76	12.30	17.82	<b>7.98</b>	<b>12.60</b>	<b>18.10</b>
BCH(31,16)	4.63	5.88	7.60	6.39	8.29	10.66	<b>7.01</b>	<b>9.33</b>	<b>12.27</b>	6.41	8.30	10.77	<b>7.21</b>	<b>9.47</b>	<b>12.45</b>
BCH(63,36)	3.72	4.65	5.66	4.68	6.65	9.10	<b>5.19</b>	<b>6.95</b>	<b>9.33</b>	<b>5.09</b>	<b>6.96</b>	<b>9.43</b>	4.90	6.64	9.19
BCH(63,45)	4.08	4.96	6.07	5.60	7.79	10.93	<b>5.90</b>	<b>8.24</b>	<b>11.46</b>	5.72	7.99	11.21	<b>5.83</b>	<b>8.15</b>	<b>11.52</b>
BCH(63,51)	4.34	5.29	6.35	5.66	7.89	11.01	<b>5.72</b>	<b>8.01</b>	<b>11.24</b>	5.38	7.40	10.50	<b>5.68</b>	<b>7.88</b>	<b>11.04</b>



Both highly and slightly corrupted codewords go through the same computationally demanding *neural* decoding procedure.

***How can we develop an adaptive and iterative decoding scheme?***

# Denoising Diffusion Error Correction Codes ICLR23





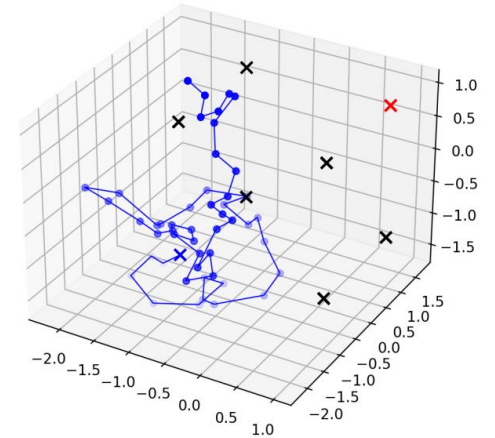
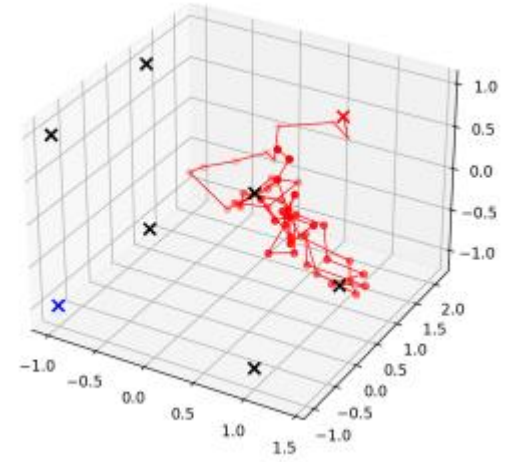
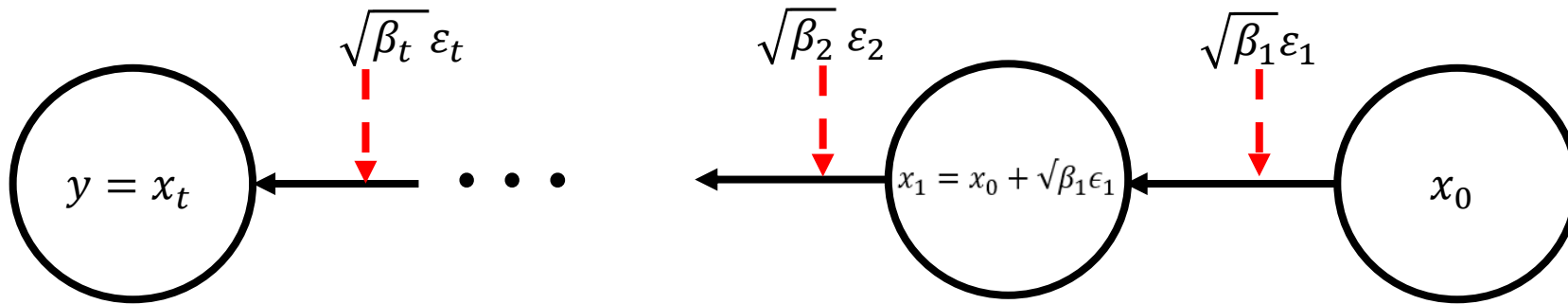
# Channel Corruption as Diffusion

- We formulate the channel corruption process as an **iterative forward diffusion process**

$$y := x_t = x_0 + \sigma \varepsilon, \varepsilon \sim \mathcal{N}(0, I)$$

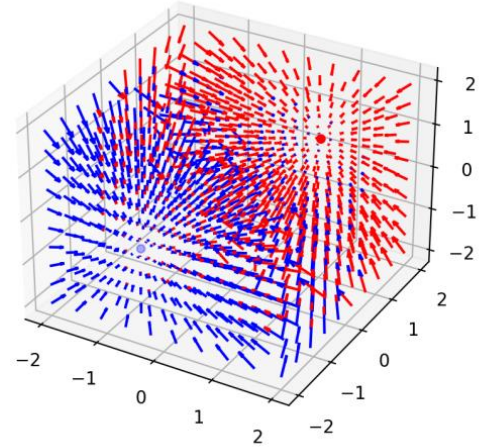
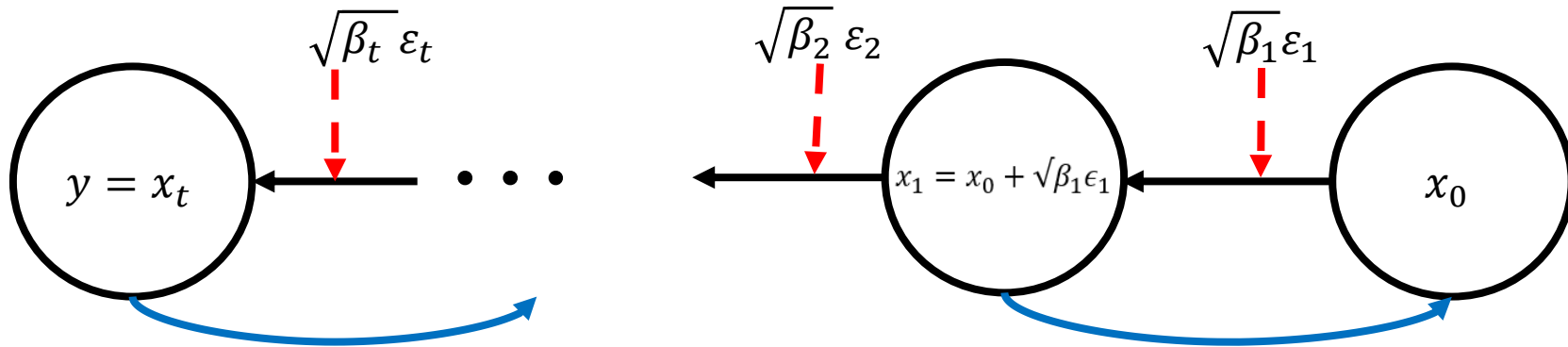
$$= x_0 + \sqrt{\beta_1} \varepsilon_1 + \dots + \sqrt{\beta_t} \varepsilon_t = x_0 + \sum_{i=1}^t \sqrt{\beta_i} \varepsilon_i,$$

$$= x_0 + \sqrt{\bar{\beta}_t} \varepsilon \sim \mathcal{N}(x_t; x_0, \bar{\beta}_t I) \quad \bar{\beta}_t = \sum_{i=1}^t \beta_i$$



# Decoding as Denoising Diffusion

- As in traditional *reverse* diffusion process, we are interested in learning to **iteratively denoise** the corrupted codeword  $y$ .



- The original **scaled** setting is adapted to the ECC setting

$$\begin{aligned}
 q(x_{t-1}|x_t, x_0) &= q(x_t|x_{t-1}, x_0) \frac{q(x_{t-1}|x_0)}{q(x_t|x_0)} \\
 &\propto \exp\left(-\frac{1}{2}\left(\frac{(x_t - x_{t-1})^2}{\beta_t} + \frac{(x_{t-1} - x_0)^2}{1 - \bar{\beta}_t} - \frac{(x_t - x_0)^2}{1 - \bar{\beta}_t}\right)\right) \\
 &= \exp\left(-\frac{1}{2}\left(\left(\frac{1}{\beta_t} + \frac{1}{\bar{\beta}_t}\right)x_{t-1}^2 - \left(\frac{2}{\beta_t}x_t - \frac{2}{\bar{\beta}_t}x_0\right)x_{t-1} + C(x_t, x_0)\right)\right)
 \end{aligned}$$

$$\tilde{\mu}_t(x_t, x_0) = \frac{\bar{\beta}_t}{\bar{\beta}_t + \beta_t} x_t + \frac{\beta_t}{\bar{\beta}_t + \beta_t} x_0 = x_t - \frac{\sqrt{\bar{\beta}_t \beta_t}}{\bar{\beta}_t + \beta_t} \varepsilon, \text{ and } \tilde{\beta}_t = \frac{\bar{\beta}_t \beta_t}{\bar{\beta}_t + \beta_t}.$$

# Parity Check Conditioning and Reverse Diffusion

- The reverse denoising process of *traditional* DDPM is conditioned by the *time step*.
- Since we are not interested in generative models, we suggest **conditioning** the diffusion decoder according to the number of **parity check errors** which conveys information about the level of corruption.

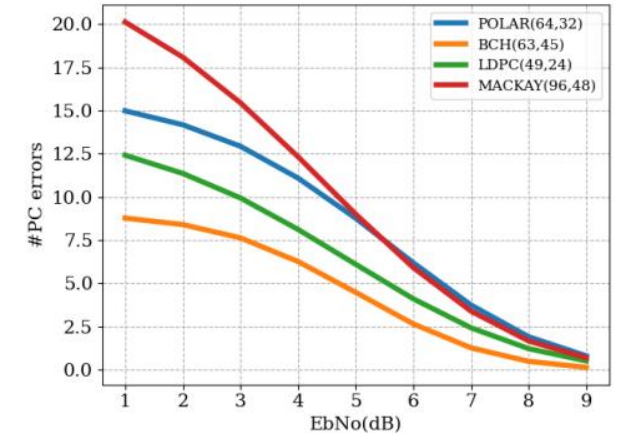
$$\mathcal{L}(\theta) = -\mathbb{E}_{t,x_0,\varepsilon} \log \left( \epsilon_\theta(x_0 + \bar{\beta}_t^{1/2} \varepsilon, \mathbf{e}_t), \tilde{\varepsilon}_b \right)$$

- The diffusion model is trained to predict the binary (sign of the ) **multiplicative** noise
- We obtain the traditional additive noise by *subtracting the predicted codeword*

$$\hat{\varepsilon} = y - \text{sign}(\hat{x}) = y - \text{sign}(\hat{\varepsilon}y).$$

- The final *reverse* diffusion process is given by

$$x_{t-1} = x_t - \frac{\sqrt{\bar{\beta}_t} \beta_t}{\bar{\beta}_t + \beta_t} (x_t - \text{sign}(x_t \epsilon_\theta(x_t, e_t)))$$

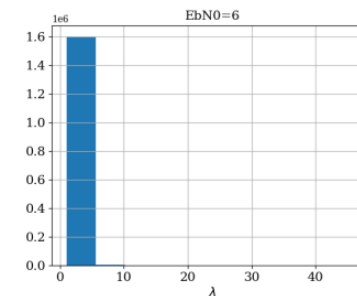
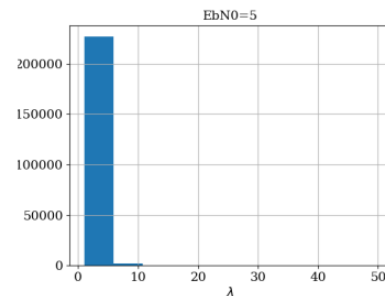
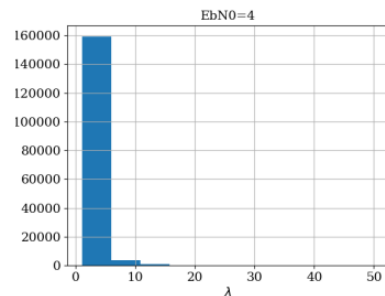
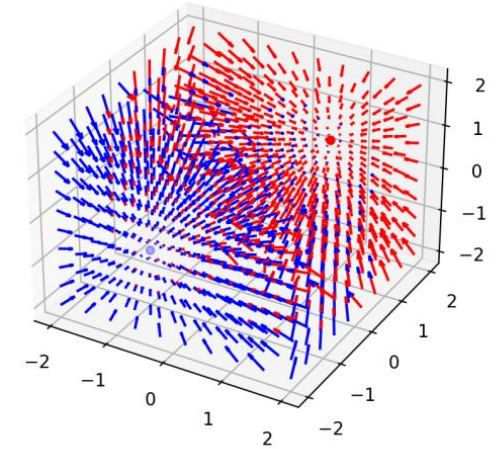


# Optimal Diffusion Step Size

- One major limitation of the generative neural diffusion process is **the large number of steps** required - generally a thousand - to generate high-quality samples.
- The number of parity check errors conveys information about the level of corruption
  - Given the reverse diffusion direction/vector, we define the optimal step size as the one which **minimizes the number of parity check errors**.
- We propose to find the optimal step size  $\lambda$  by solving the minimum number of parity errors.

$$\lambda^* = \arg \min_{\lambda \in \mathbb{R}^+} \left\| s \left( x_t - \lambda \frac{\sqrt{\bar{\beta}_t} \beta_t}{\bar{\beta}_t + \beta_t} \hat{\varepsilon} \right) \right\|_1,$$

- Since the objective is highly non-differentiable and non-convex (sign, modulo 2), we suggest using a grid **line-search** approach



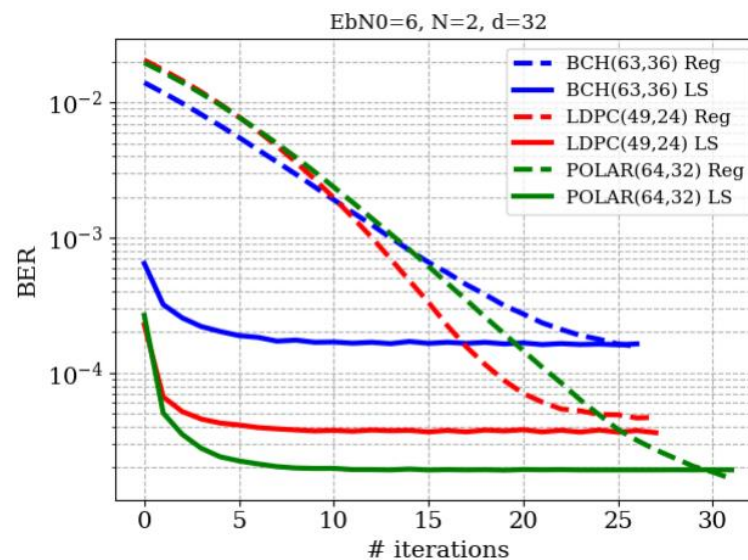
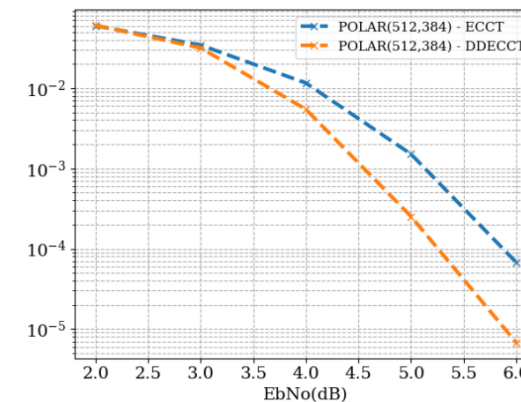
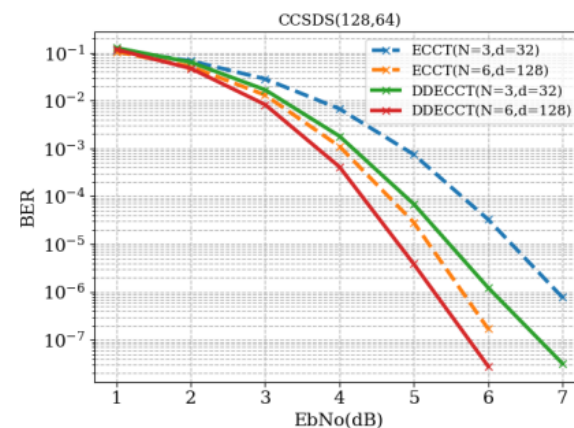
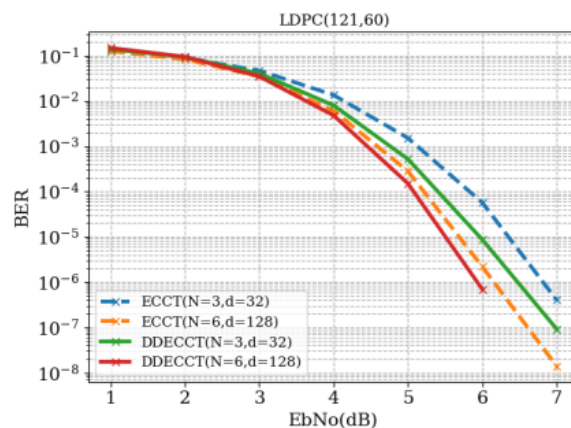
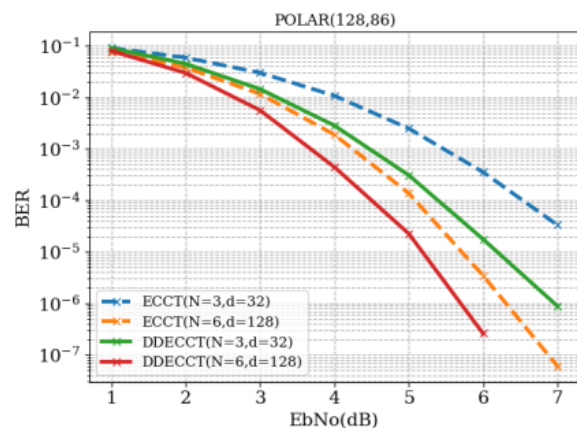
# Experiments

- The base noise estimator is an adapted (“time” conditioned) *ECCT*
- Our approach **outperforms** the current SOTA results (obtained by ECCT) by extremely large margins on several families of codes of different lengths and rates, at a **fraction** of the capacity.
- Especially for shallow models, the difference can be **orders of magnitude**.

Method $E_b/N_0$	BP			ARBP			ECCT N=2			ECCT N=6			Ours N=2			Ours N=6		
	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
Polar(64,32)	3.52	4.04	4.48	4.77	6.30	8.19	4.27	5.44	6.95	5.71	7.63	9.94	5.99	8.16	10.90	6.76	9.14	12.31
	4.26	5.38	6.50	5.57	7.43	9.82	4.57	5.86	7.50	6.48	8.60	11.43	6.23	8.52	11.23	6.90	9.43	12.85
Polar(64,48)	4.15	4.68	5.31	5.25	6.96	9.00	4.92	6.46	8.41	5.82	7.81	10.24	5.55	7.67	10.08	5.98	8.02	10.94
	4.74	5.94	7.42	5.41	7.19	9.30	5.14	6.78	8.9	6.15	8.20	10.86	5.74	7.85	10.40	5.98	8.26	11.13
Polar(128,64)	3.38	3.80	4.15	4.02	5.48	7.55	3.51	4.52	5.93	4.47	6.34	8.89	5.37	7.75	10.51	6.34	9.26	12.77
	4.10	5.11	6.15	4.84	6.78	9.30	3.83	5.16	7.04	5.12	7.36	10.48	5.97	8.52	11.72	7.24	10.70	14.56
Polar(128,86)	3.80	4.19	4.62	4.81	6.57	9.04	4.30	5.58	7.34	5.36	7.45	10.22	5.61	7.76	10.42	6.52	9.21	12.64
	4.49	5.65	6.97	5.39	7.37	10.13	4.49	5.90	7.75	5.75	8.16	11.29	5.99	8.19	11.00	7.09	10.20	13.84
Polar(128,96)	3.99	4.41	4.78	4.92	6.73	9.30	4.56	5.98	7.93	5.39	7.62	10.45	5.60	7.83	10.56	6.46	9.41	12.52
	4.61	5.79	7.08	5.27	7.44	10.2	4.69	6.20	8.30	5.88	8.33	11.49	5.95	8.42	11.38	6.83	9.99	13.36
LDPC(49,24)	5.30	7.28	9.88	6.05	8.13	11.68	4.51	6.07	8.11	5.74	8.13	11.30	5.27	7.38	10.23	5.87	8.22	11.56
	6.23	8.19	11.72	6.58	9.39	12.39	4.58	6.18	8.46	5.91	8.42	11.90	5.31	7.35	10.40	5.84	8.29	11.85
LDPC(121,60)	4.82	7.21	10.87	5.22	8.31	13.07	3.88	5.51	8.06	4.98	7.91	12.70	4.48	6.95	10.65	5.25	8.43	13.80
	-	-	-	-	-	-	3.89	5.55	8.16	5.02	7.94	12.72	4.56	7.02	10.64	5.32	8.69	13.82
LDPC(121,70)	5.88	8.76	13.04	6.45	10.01	14.77	4.63	6.68	9.73	6.11	9.62	15.10	5.41	8.22	12.22	6.49	10.39	15.43
	-	-	-	-	-	-	4.64	6.71	9.77	6.28	10.12	15.57	5.52	8.47	12.63	6.64	10.65	16.21
LDPC(121,80)	6.66	9.82	13.98	7.22	11.03	15.90	5.27	7.59	10.08	6.92	10.74	15.10	6.12	9.38	13.25	7.68	12.19	17.83
	-	-	-	-	-	-	5.29	7.63	10.90	7.17	11.21	16.31	6.26	9.41	13.41	7.39	11.46	17.65
MacKay(96,48)	6.84	9.40	12.57	7.43	10.65	14.65	4.95	6.67	8.94	6.88	9.86	13.40	6.18	8.63	11.53	7.86	11.61	15.51
	-	-	-	-	-	-	5.04	6.80	9.23	7.10	10.12	14.21	6.28	8.8	11.78	7.93	11.65	15.51
CCSDS(128,64)	6.55	9.65	13.78	7.25	10.99	16.36	4.35	6.01	8.30	6.34	9.80	14.40	5.79	8.48	12.24	7.28	11.66	17.02
	-	-	-	-	-	-	4.41	6.09	8.49	6.65	10.40	15.46	5.81	8.79	12.29	7.55	12.01	17.62
BCH(63,36)	3.72	4.65	5.66	4.33	5.94	8.21	3.79	4.87	6.35	4.42	5.91	8.01	4.71	6.45	8.72	5.01	6.84	9.30
	4.03	5.42	7.26	4.57	6.39	8.92	4.05	5.28	7.01	4.62	6.24	8.44	4.84	6.65	9.01	5.07	7.02	9.85
BCH(63,45)	4.08	4.96	6.07	4.80	6.43	8.69	4.47	5.88	7.81	5.16	7.02	9.75	5.12	7.16	9.95	5.49	7.71	10.86
	4.36	5.55	7.26	4.97	6.90	9.41	4.66	6.16	8.17	5.41	7.49	10.25	5.33	7.49	10.18	5.60	8.02	11.05
BCH(63,51)	4.34	5.29	6.35	4.95	6.69	9.18	4.60	6.05	8.05	5.20	7.08	9.65	5.09	7.08	9.87	5.35	7.49	10.38
	4.5	5.82	7.42	5.17	7.16	9.53	4.78	6.34	8.49	5.46	7.57	10.51	5.19	7.23	10.20	5.39	7.48	10.53
							5.01	6.72	9.03	5.66	7.89	11.01	5.21	7.29	10.13	5.26	7.40	10.49

- $E_b/N_0 \sim SNR$
- Code( $n, k$ )

# Experiments and Analysis



**LS:** Our line search approach allows convergence within very few iterations.



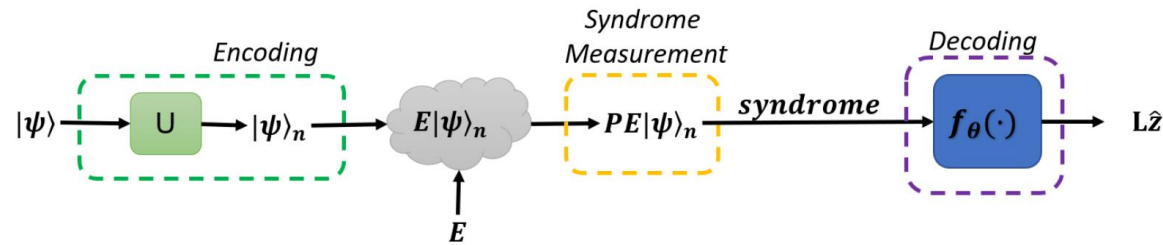
***How to adapt neural decoders to Quantum error correction?***



# Deep Quantum Error Correction AAI24

- **Goal:** allow the protection of quantum information from quantum noise (e.g., quantum gates, decoherence).
- A quantum bit (*qubit*) is defined as the superposition of two states

$$|\psi\rangle = \alpha|0\rangle + \beta|1\rangle, \quad \alpha, \beta \in \mathbb{C}, \quad |\alpha|^2 + |\beta|^2 = 1$$



- There are three major differences/challenges with classical error correction
  1. **Syndrome Decoding:** There is no arbitrary access to the current state (due to quantum wave measurement collapse) such that only partial information defined by the syndrome is available.  
*It requires an adaptation of the existing neural decoders to syndrome decoding.*
  2. **Logical Decoding:** We are interested in the logical qubits only, meaning we wish to predict the codeword **up to** the logical operators mapping (i.e.,  $\mathbb{L}\hat{Z}$  instead of  $\hat{Z}$ ).  
*However, this mapping is defined over the highly non-differentiable  $GF(2)$  (i.e., XOR).*
  3. **Noisy Syndrome measurement:** The syndrome measurement itself being noisy, the decoding must be performed based on multiple noisy measurements of the syndrome is obtained upon **multiple** noisy syndrome observations.  
*Efficient decoding methods must be developed.*
- *These QECC challenges are at the core of our contributions.*

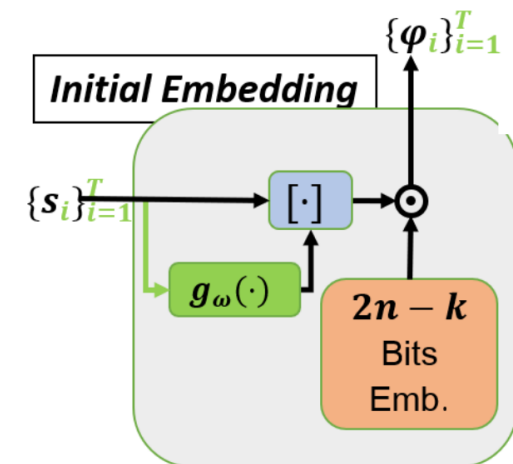
# Overcoming Measurement Collapse by Prediction

- In the QECC setting, only the code **syndrome** is available, since classical measurements are not allowed due to the *wave function collapse phenomenon*. Thus, it's not possible to arbitrarily access/measure the quantum state.
- We thus propose to *extend the ECCT*, by replacing the magnitude of the channel output (i.e.,  $|\mathbf{y}|$ ) with an **initial estimate of the noise** to be further refined by the code-aware network.
  - Reminiscent of MCMC methods.
- Given  $g_\omega : \{0, 1\}^{n_s} \rightarrow \mathbb{R}^n$  the **initial** parameterized **noise estimator** from a given syndrome

$$\hat{z} = f_\theta(h_q(s)) = f_\theta([g_\omega(s), s])$$

- The initial estimator is trained to predict the noise from the syndrome (i.e., *syndrome decoding*)

$$\mathcal{L}_g = \text{BCE}(g_\omega(s), \varepsilon)$$



# Logical Decoding

- The logical error rate (LER) metric provides valuable information on the **practical** decoding performance.
- We wish to minimize the following LER objective

$$\mathcal{L}_{\text{LER}} = \text{BCE}(\mathbb{L}f_{\theta}(s), \mathbb{L}\varepsilon) \quad \mathbb{L} \in \{0, 1\}^{k \times n}.$$

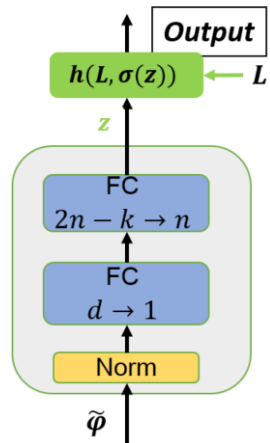
where the multiplications are performed over *the highly **nondifferentiable** GF(2)*.

- We propose to optimize the objective using a differentiable equivalence mapping of the XOR operator via **bipolar mapping**
  - $\phi(u) = 1 - 2u$  induces  $\phi(u \oplus v) = \phi(u)\phi(v)$  such that we have

$$\left( \Lambda(\mathbb{L}, x) \right)_i := \mathbb{L}_i \oplus x = \phi^{-1} \left( \prod_j \phi((\mathbb{L})_{ij} \cdot x_j) \right)$$

- The LER training objective becomes

$$\mathcal{L}_{\text{LER}} = \text{BCE} \left( \Lambda(\mathbb{L}, \text{bin}(f_{\theta}(s))), \mathbb{L}\varepsilon \right)$$



# Noisy Syndrome Measurement

- In the presence of measurement errors, each syndrome measurement is repeated several times, and efficient noisy measurement decoding is required.

- At each time sample we have the syndrome defined as

$$s_t = \left( H(x \oplus \varepsilon_1 \oplus \cdots \oplus \varepsilon_t) \right) \oplus \tilde{\varepsilon}_t$$

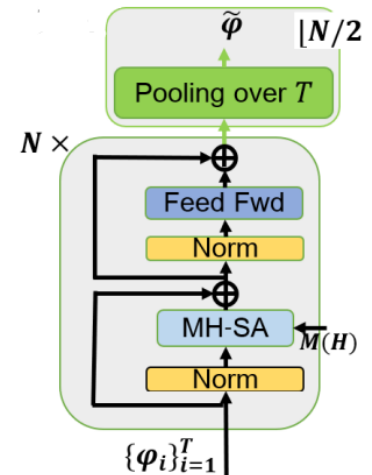
- We first analyze each measurement separately and then perform global decoding **at the embedding level** by applying a symmetric pooling function, e.g. an average, in the **middle** of the neural network.

- Given a neural decoder with  $N$  layers and the activations  $\varphi \in \mathbb{R}^{T \times n \times d}$ , the pooled embedding is given by  $\tilde{\varphi} = \frac{1}{T} \sum_t \varphi_t$  at layer  $l = \lfloor N/2 \rfloor$

- **Final Objective:** Given  $\mathcal{L}_{\text{BER}} = \text{BCE}(f_\theta(s), \varepsilon)$  the overall objective is

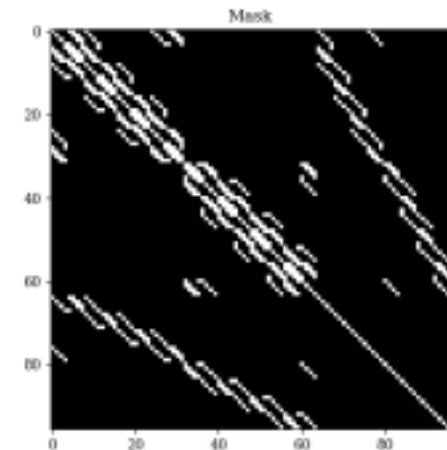
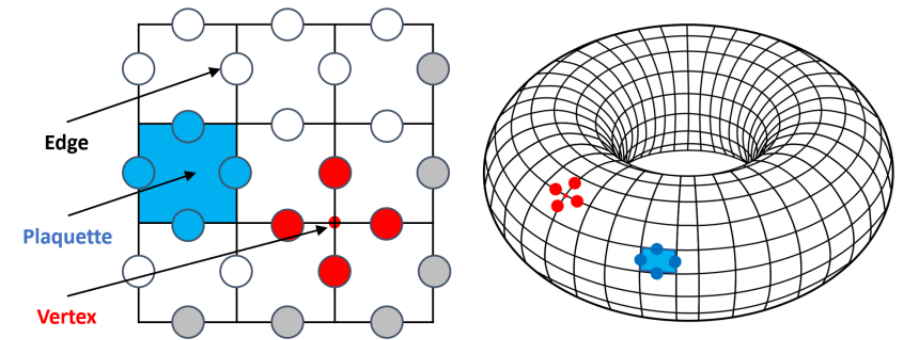
$$\mathcal{L} = \lambda_{\text{BER}} \mathcal{L}_{\text{BER}} + \lambda_{\text{LER}} \mathcal{L}_{\text{LER}} + \lambda_g \mathcal{L}_g$$

- The **BER regularization is important** since the  $GF2$  optimization induces severe *saddle point optimization*.

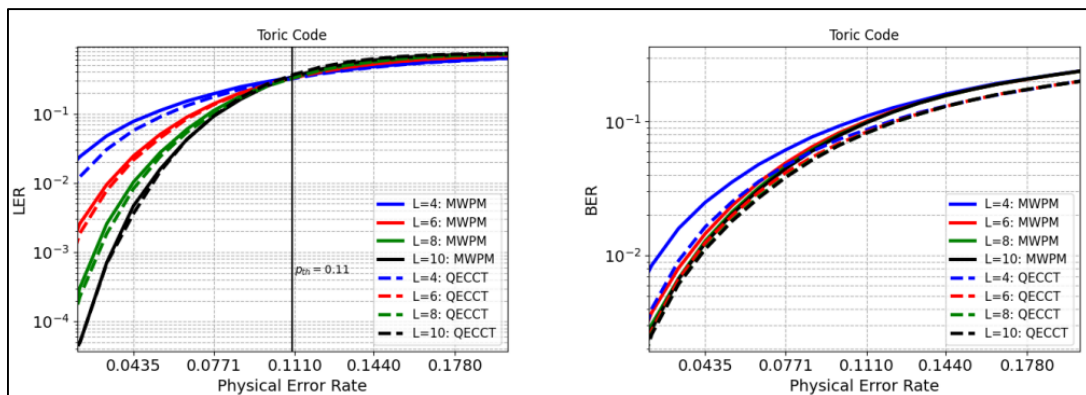


# Experiments

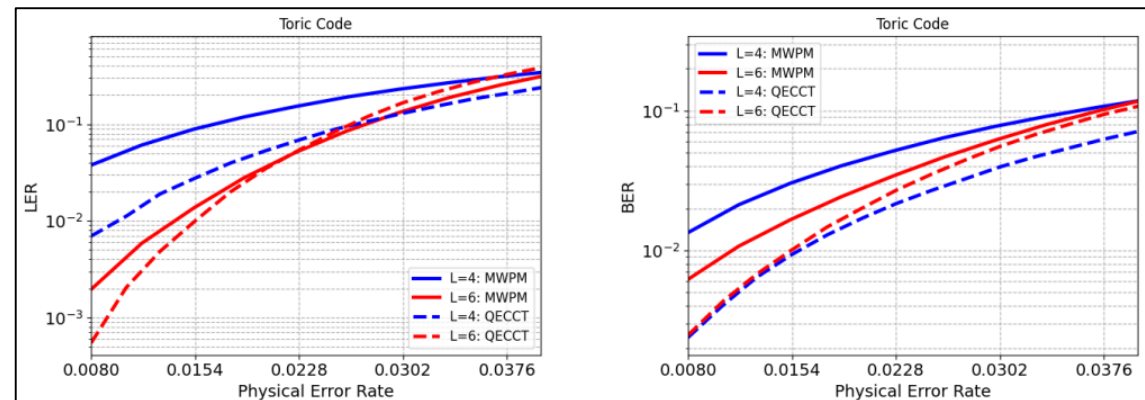
- Experiments are performed on the popular **Toric and Surface** codes.
  - The performances are reported for lattice of up to length 10 (i.e., hundreds of qubits)
- Several noise settings are experimented
  - **Independent** noise
  - **Depolarization** noise
  - **Circuit** noise
  - With **faulty syndrome measurements**
- The model is based on the *refined* ECCT with a *shallow* architecture of 6 layers with highly **sparse** self-attention.
- The baseline is the very popular Minimum Weight Perfect Matching (**MWPM**) algorithm ( $O((n^3 + n^2) \log(n))$  to  $O(n^2)$ )



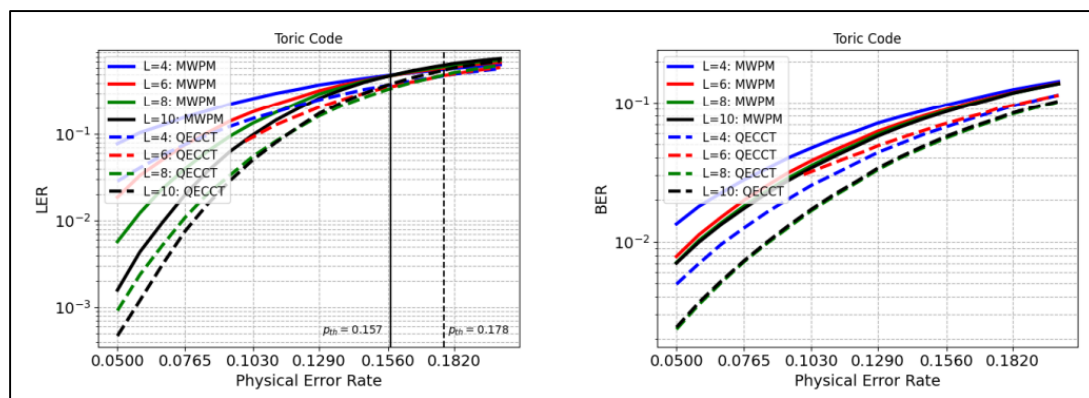
# Toric Code Results



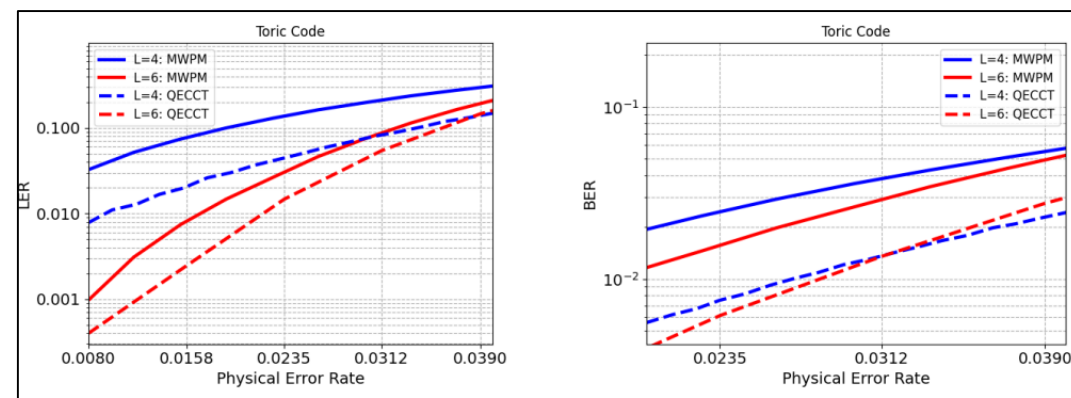
Independent noise



Independent noise  
with faulty syndrome measurements

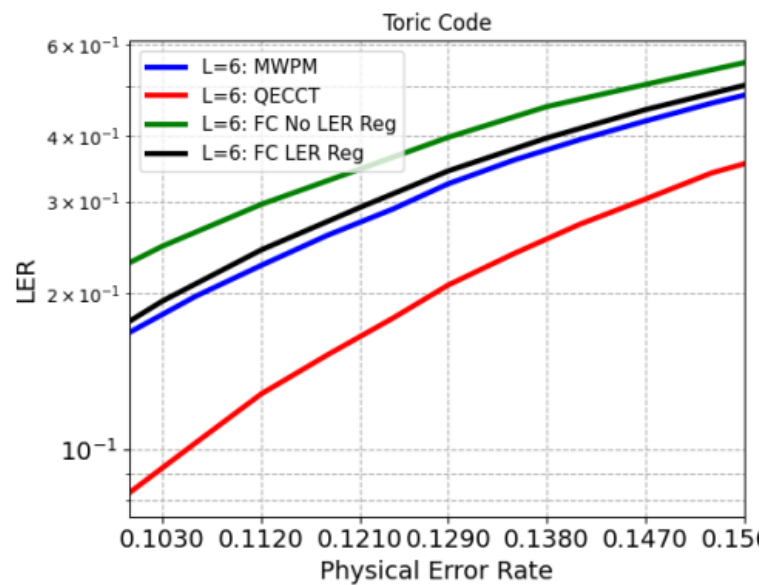


Depolarization noise

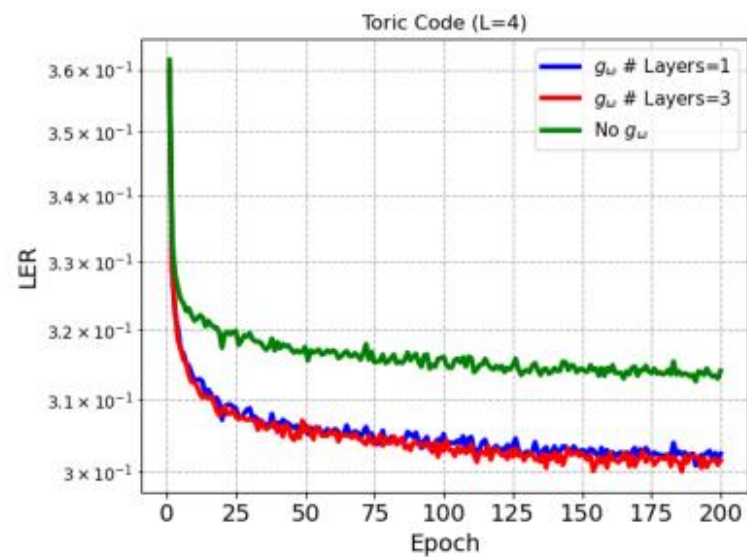


Depolarization noise  
with faulty syndrome measurements

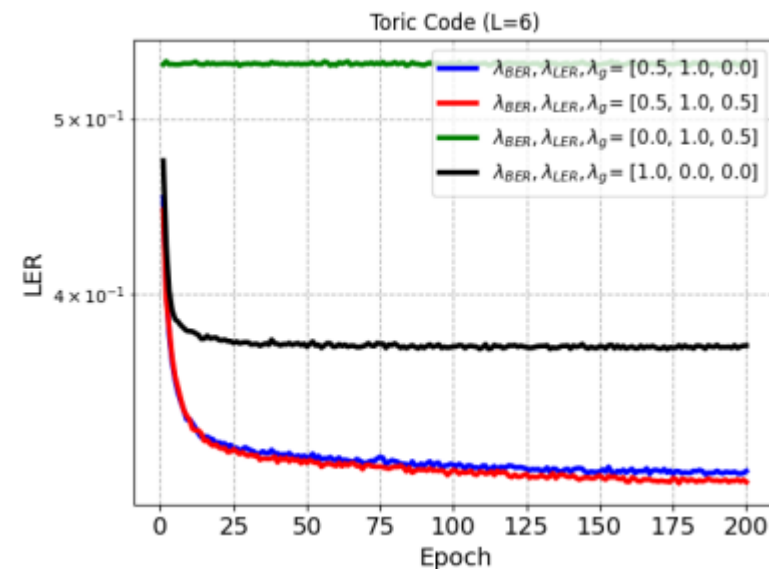
# Analysis



Impact of the architecture and of the LER optimization



Impact of the initial noise estimator



Impact of the BER regularization

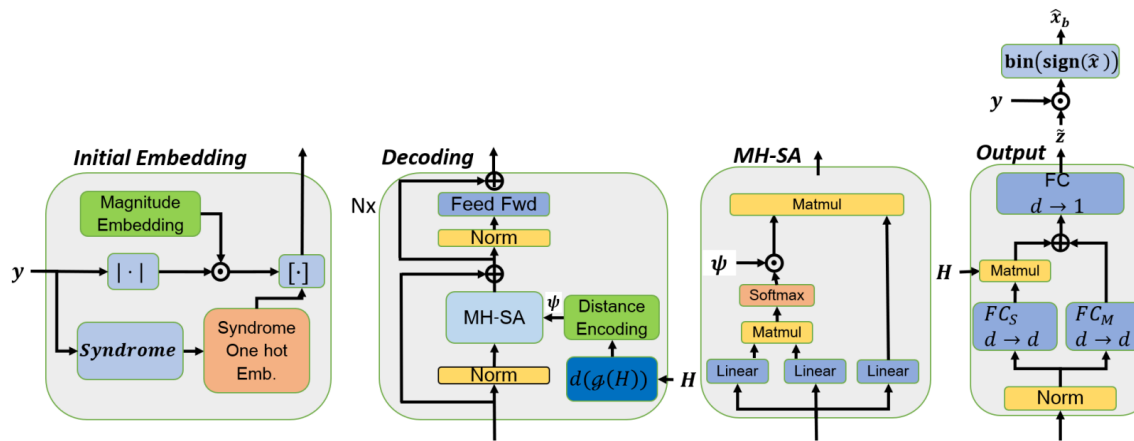
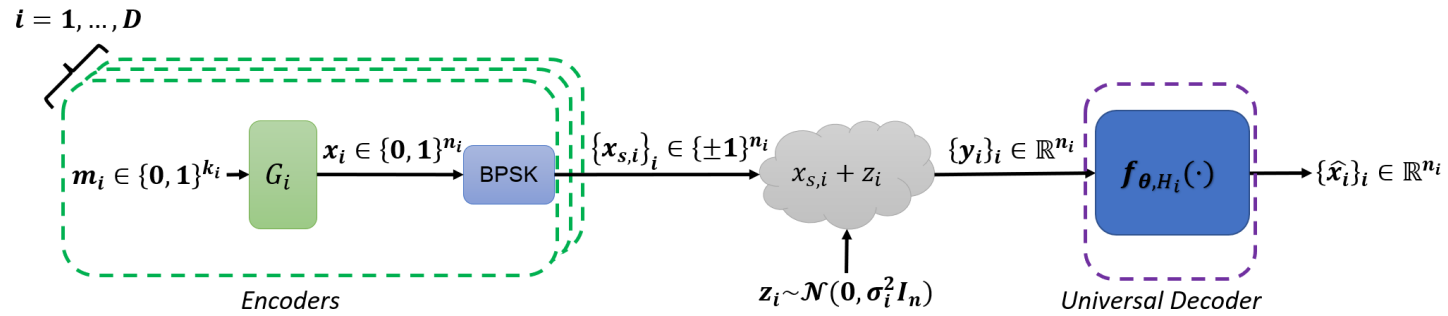




One needs to develop, train, and deploy one (neural) decoder for each family of code, length, and rate.

***How can we develop a single universal neural decoder  
which is code/length/rate invariant?***

# A Foundation Model for Error Correction Codes ICLR24



# Code-Invariant Initial Embedding

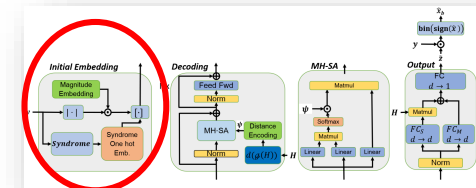
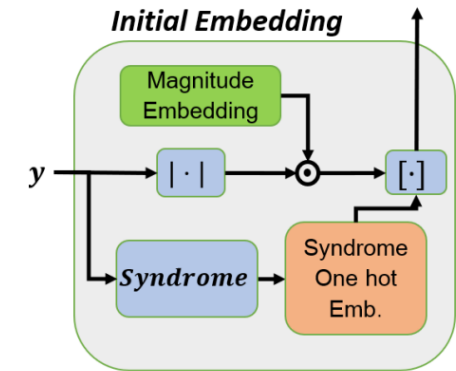
- In ECCT, a unique model is crafted for every code and length where the initial embedding is designed such that **each input bit possesses its distinct embedding vector**, providing, as a byproduct, a **learned positional encoding**.  $\Phi = (h(y) \cdot 1_d^T) \odot W$

- In our **length-invariant** model (FECCT), we propose a new code-invariant embedding, where a **single embedding** is given for all **magnitude** elements, and **two embeddings** are given for every element of the binary **syndrome**.

$$\phi_i = \begin{cases} |y_i| W_M, & \text{if } i \leq n \\ W_{(s(y))_{i-n+1}}^S, & \text{otherwise} \end{cases}$$

With  $\{W_M, W_0^S, W_1^S\} \in \mathbb{R}^d$ .

- This new **length/rate-invariant initial encoding** requires **three** embedding vectors compared to the  $2n - k$  vectors of the ECCT.
- In contrast to ECCT which captures the bit position with learned embedding, our method **lacks positional information**.



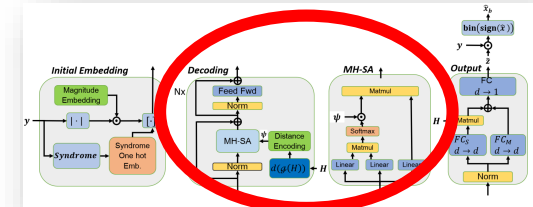
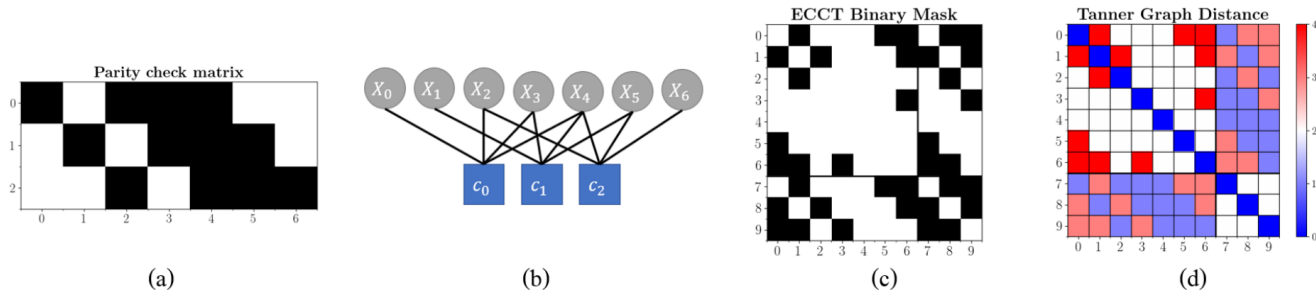
# Tanner Graph Distance Masking as Code and Positional encoding

- FECCT's SA masking serves two purposes.
  - Similar to ECCT, it **integrates the code structure** into the transformer.
  - Adds the **relative position information** to the processed elements.
- The Tanner graph captures the **relations** between every two bits in the code (**relative positional encoding**).
- We consider the distance matrix  $\mathcal{G}(H) \in \mathbb{N}^{(2n-k) \times (2n-k)}$ , induced by the code (Tanner graph).
  - Each element  $(i, j)$  in this matrix is defined as the *length of the shortest path in the Tanner graph* between node  $i$  and node  $j$ .

- We learn a *parameterized mapping*  $\psi: \mathbb{N} \rightarrow \mathbb{R}$  of the *distance matrix*, incorporated into the self-attention

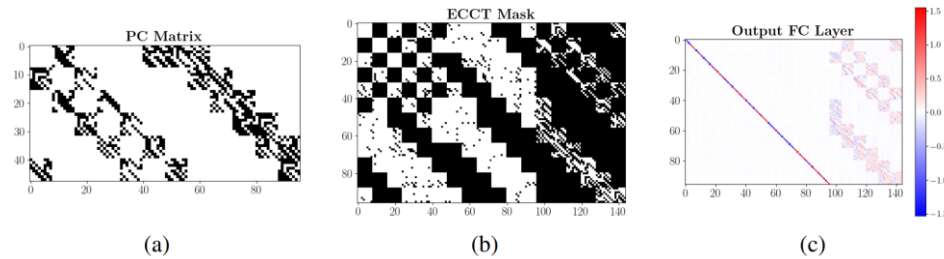
$$A_H(Q, K, V) = \left( \text{Softmax} \left( \frac{QK^T}{\sqrt{d}} \right) \odot \psi(\mathcal{G}(H)) \right) V.$$

- This attention mechanism **generalizes** the ECCT which captures **only up to two rings** connectivity information.



# Parity-Check Aware Prediction

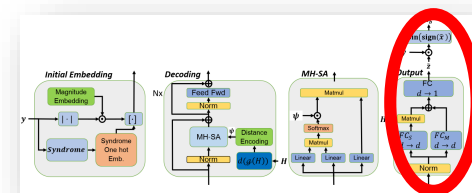
- ECCT makes use of two *fully connected layers* (least length invariant modules) for the final prediction ( $(2n - k) \rightarrow n$ )
- ECCT's learned output layer is (surprisingly) greatly **induced** by the code/parity check matrix.



- Motivated by this phenomenon, we explicitly **enforce** a similar dependency structure.
- By *splitting* the syndrome and the channel output elements we integrate the remaining syndrome information by **aggregation** according to the parity check matrix connectivity

$$\hat{\varepsilon} = (\phi_{o,M} W_M + H^T (\phi_{o,S} W_S)) W_{d \rightarrow 1}$$

- This way, the final prediction is **code-aware** but also **code/length invariant**.
- Finally, the FECCT being invariant, its *number of parameters* is **independent of the code**.



# Experiments

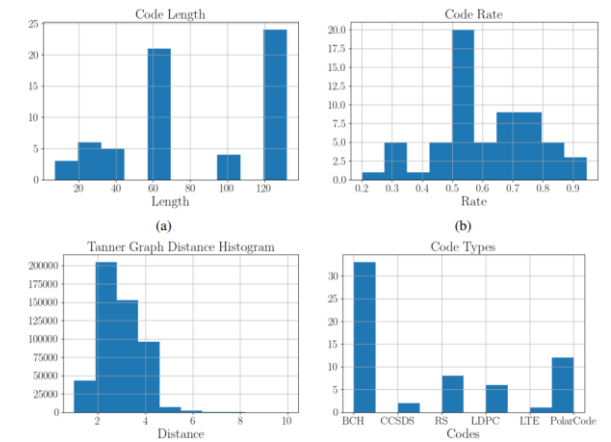
- Trained on multiple codes, **our single decoder (with smaller capacity)** can match and even outperform **other methods designed and trained separately on each code**, in multiple scenarios
  - Pretrained codes
  - Zero-shot codes
  - Fine-tuned codes

Supervision	Unlearned			Fully supervised						Zero-Shot					
	BP			Hyp BP			ARBP			ECCT			Ours		
Method															
$E_b/N_0$	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
BCH(63,45)	4.08	4.96	6.07	4.48	6.07	8.45	4.80	6.43	8.69	5.18	7.24	10.20	5.18	7.32	10.31
	4.36	5.55	7.26	4.64	6.27	8.51	4.97	6.90	9.41						
BCH(63,51)	4.34	5.29	6.35	4.64	6.08	8.16	4.95	6.69	9.18	5.63	7.96	11.22	5.71	8.07	11.31
	4.50	5.82	7.42	4.80	6.44	8.58	5.17	7.16	9.53						
BCH(127,92)	NA	NA	NA	NA	NA	NA	NA	NA	NA	4.10	5.71	8.38	4.11	5.84	8.79
BCH(255,163)	NA	NA	NA	NA	NA	NA	NA	NA	NA	3.34	4.13	5.80	3.34	4.13	5.76
CCSDS(128,64)	6.55	9.65	13.78	6.99	10.57	15.27	7.25	10.99	16.36	6.77	10.51	15.90	6.52	9.67	15.01
CCSDS(32,16)	NA	NA	NA	NA	NA	NA	NA	NA	NA	5.93	7.77	10.02	5.23	7.00	9.21
POLAR(128,86)	3.80	4.19	4.62	4.57	6.18	8.27	4.81	6.57	9.04	6.39	9.08	12.70	5.53	7.90	11.29
	4.49	5.65	6.97	4.95	6.84	9.28	5.39	7.37	10.13						
POLAR(64,32)	3.52	4.04	4.48	4.25	5.49	7.02	4.77	6.30	8.19	6.91	9.18	12.34	5.88	7.91	10.76
	4.26	5.38	6.50	4.59	6.10	7.69	5.57	7.43	9.82						

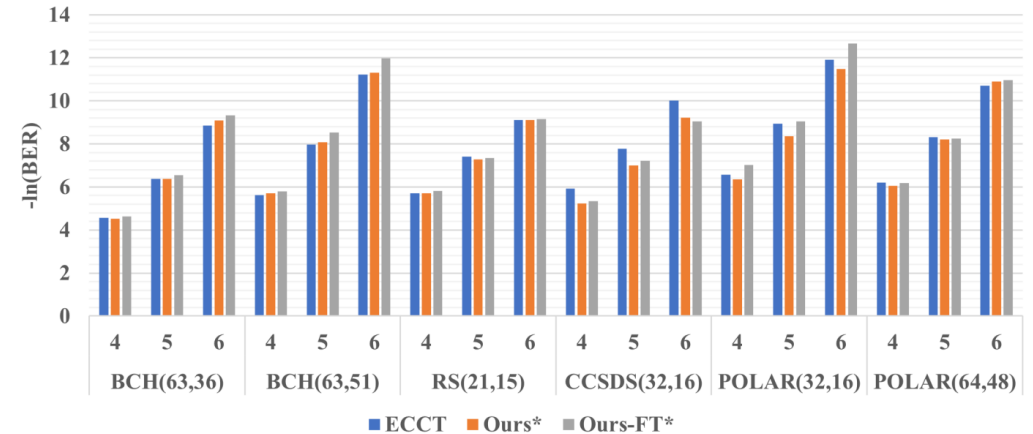
Zero-Shot Codes

- $E_b/N_0 \sim SNR$
- Code( $n, k$ )

Method	BP			Hyp BP			ARBP			ECCT			Ours		
	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
BCH(63,36)	3.72	4.65	5.66	3.96	5.35	7.20	4.33	5.94	8.21	4.56	6.37	8.85	4.53	6.38	9.10
BCH(127,120)	NA	NA	NA	NA	NA	NA	NA	NA	NA	4.70	6.37	8.95	4.62	6.33	8.95
Reed Solomon(21,15)	NA	NA	NA	NA	NA	NA	NA	NA	NA	5.71	7.42	9.11	5.71	7.28	9.12
Reed Solomon(60,52)	NA	NA	NA	NA	NA	NA	NA	NA	NA	5.53	7.54	9.98	5.47	7.59	10.21
POLAR(32,16)	NA	NA	NA	NA	NA	NA	NA	NA	NA	6.57	8.94	11.91	6.36	8.36	11.49
POLAR(64,48)	3.52	4.04	4.48	4.25	5.49	7.02	4.77	6.30	8.19	6.21	8.32	10.71	6.06	8.21	10.90
	4.26	5.38	6.50	4.59	6.10	7.69	5.57	7.43	9.82						

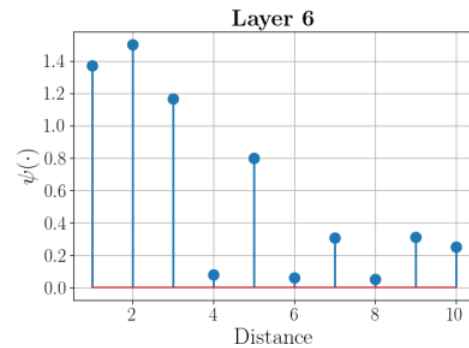
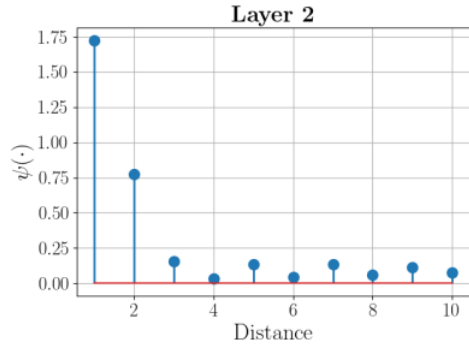


Pretrained Codes



Fine-Tuned Codes

# Analysis



Method	POLAR(64,32)			BCH(63,45)		
	4	5	6	4	5	6
<b>ECCT</b>	4.12	5.22	6.67	4.45	5.81	7.65
ECCT + <b>II</b>	4.27	5.54	7.14	4.52	5.98	7.92
ECCT + <b>IO</b>	4.44	5.73	7.40	4.41	5.76	7.62
ECCT + <b>II + IO</b>	4.09	5.26	6.80	4.31	5.62	7.41
ECCT + <b>DM</b>	4.44	5.73	7.37	4.74	6.34	8.53
ECCT + <b>DM + II</b>	4.44	5.73	7.37	5.17	7.07	9.59
ECCT + <b>DM + IO</b>	4.36	5.64	7.32	4.53	6.01	8.03
<b>FECCT: ECCT + DM + II + IO</b>	4.36	5.64	7.32	4.52	5.98	8.05

Method	FECCT - single			FECCT		
	4	5	6	4	5	6
POLAR(64,48)	6.35	8.50	11.12	6.06	8.21	10.96
POLAR(128,86)*	3.90	5.36	7.57	5.53	7.90	11.29
BCH(63,36)	4.01	5.42	7.30	4.53	6.38	9.10
BCH(63,51)*	4.65	6.35	8.73	5.71	8.07	11.31
Reed Solomon(21,15)	4.25	4.62	4.97	4.56	6.83	10.51
Reed Solomon(60,52)	3.68	3.81	3.77	5.47	7.49	10.24
CCSDS(128,64)*	2.90	3.42	4.30	6.52	9.67	15.01
CCSDS(32,16)*	4.10	4.54	4.43	5.23	7.00	9.21

- **Learned Distance Mapping:** FECCT seems to assign the **most** impactful mapping for the elements distanced by one and two nodes, **remarkably matching the ECCT's two-ring heuristic.**

- **Architectural Ablation**  
The ablation demonstrate the beneficial impact of **each of the contributions** on the obtained accuracy compared to SOTA

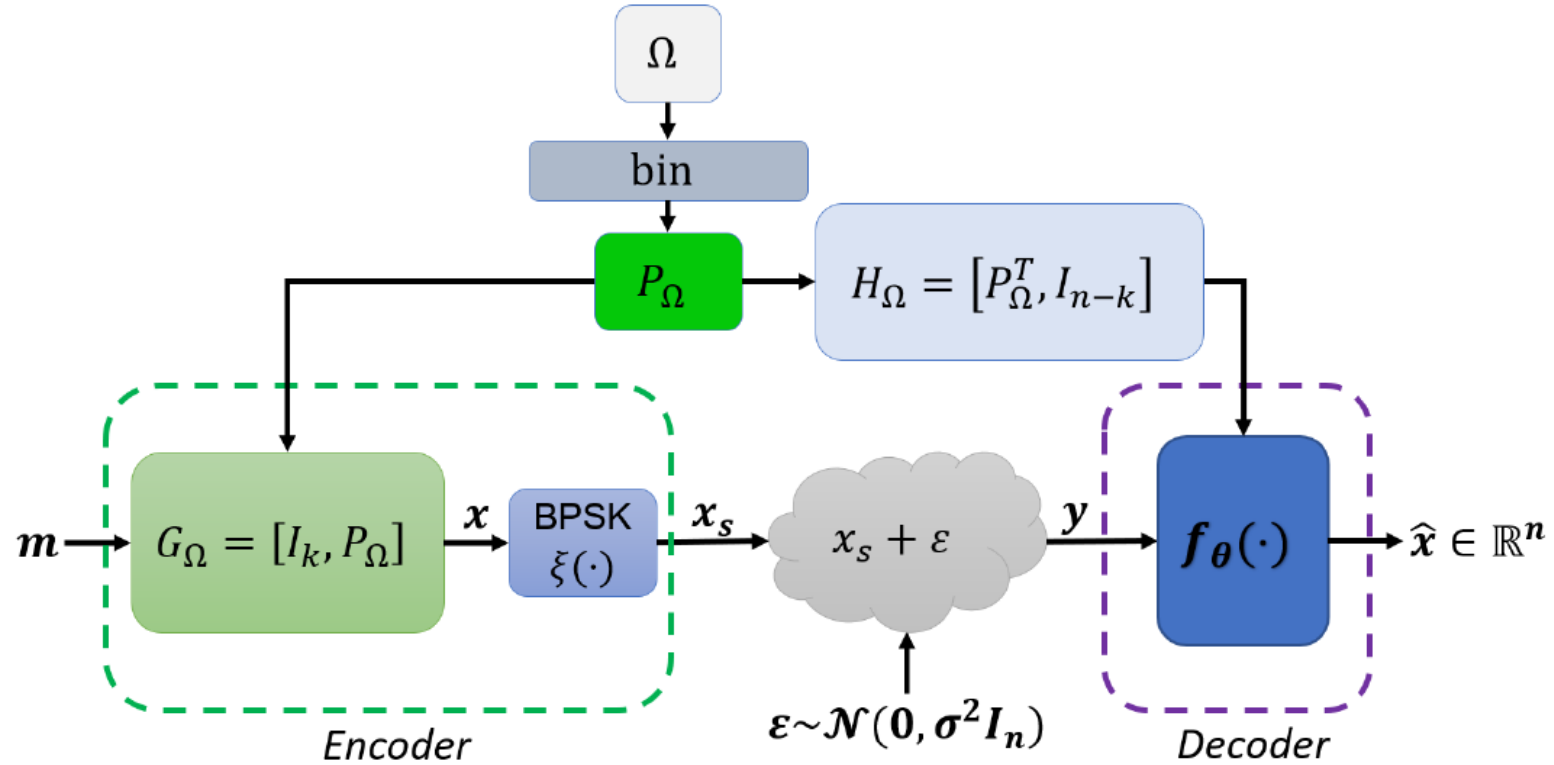
- **Generalization:**  
To show the importance of dataset diversity, we show that training FECCT on one **single** code is slightly better on the trained code but totally lacks generalization





*How can we co-learn binary linear block codes  
along with the neural decoder?*

# Learning Linear Block Error Correction Codes ICML24



# End to End Optimization

- We assume the **standard** (/canonical/systematic) **form** of the code for differentiable and fast optimization

$$G = [I_k, P], \quad P \in \{0,1\}^{k \times (n-k)}$$
$$\Rightarrow H = [P^T, I_{n-k}]$$

- Using a *general form matrix* form is preferable but requires **fast and differentiable** computation of its inverse at each iteration.
- We want **trainable parameterization** of the code  $\Omega \in \mathbb{R}^{k \times (n-k)}$  such that, with  $\text{bin}(\cdot)$  a binarization function, we have

$$P = P_\Omega = \text{bin}(\Omega),.$$

- The **end-to-end optimization objective** is now defined by

$$\mathcal{L}(\Omega, \theta) = \mathbb{E}_{m \sim \text{Bern}^k(1/2), \varepsilon \sim \mathcal{Z}} \text{BCE}(f_\theta(h_\Omega(y_\Omega)), \text{bin}(\tilde{\varepsilon}))$$

$$y_\Omega = \xi(\phi(m, G_\Omega)) + \varepsilon$$

$$h_\Omega(y_\Omega) = [|y_\Omega|, H_\Omega \text{bin}(y_\Omega)]$$

# Optimization over $GF(2)$

❖ **Two main non-differentiable** modules in the pipeline

- **Binarization**

- Performed via the Straight-Through-Estimator (STE)

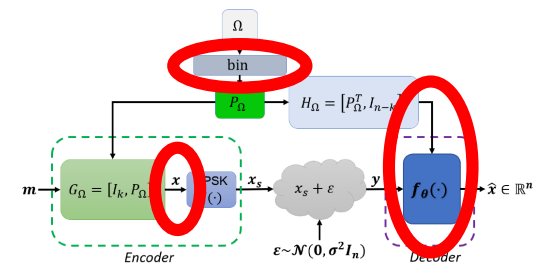
$$\begin{cases} \text{bin}(u) = \xi^{-1}(\text{sign}(u)) \\ \frac{\partial \text{bin}(u)}{\partial u} = -\frac{\mathbb{1}_{|u| \leq \tau}}{2} \end{cases}$$

- **Dot products over  $GF(2)$**

- Performed via polar transform [1]  $\xi(u) = 1 - 2u, u \in \{0, 1\} \rightarrow \xi(u \oplus v) = \xi(u)\xi(v), \forall u, v \in \{0, 1\}$

$$(\phi(m, G_\Omega))_i := G_{\Omega_i} \oplus m = \xi^{-1}\left(\prod_{j=1}^k \xi((G_\Omega)_{ij} \cdot m_j)\right)$$

- The new form defines a **multilinear polynomial** (potentially *inducing saddle-point optimization*) and the gradient can now be computed in a **differentiable** manner.



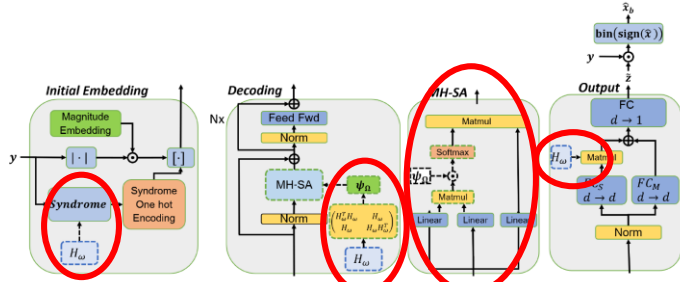
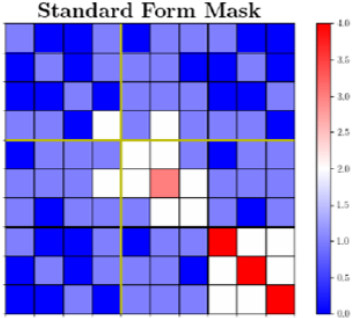
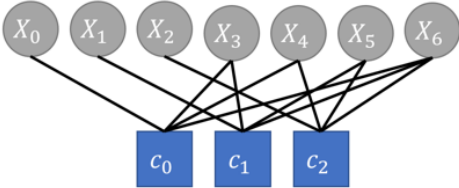
# Differentiable Masking

- We suggest to **extend FECCT** to support end-to-end training
- Existing masking methods induced from the code are **extracted once in a non-differentiable fashion**,
  - No information can be **backpropagated** during the optimization **from the mask to the code**
- We learn a parameterized mapping  $\psi_\gamma: \mathbb{N} \rightarrow \mathbb{R}$  of the elements constituting the mask, which is **derived by the parity-check matrix**, such that

$$A_H(Q, K, V) = \text{Softmax}\left(\frac{QK^T + \psi_\gamma(g(H_\Omega))}{\sqrt{d}}\right)V,$$

$$g(H_\Omega) = \begin{pmatrix} H_\Omega^T H_\Omega & H_\Omega \\ H_\Omega^T & H_\Omega H_\Omega^T \end{pmatrix}$$

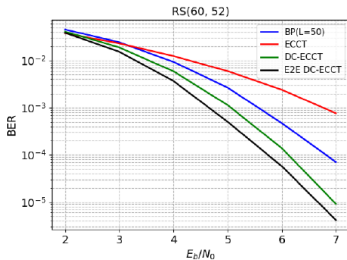
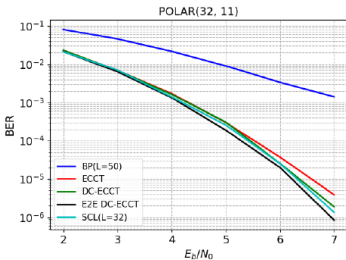
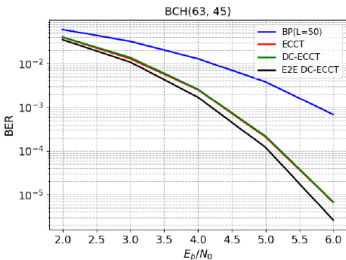
- Represents the **two-step transition matrix** between every two nodes



# Experiments

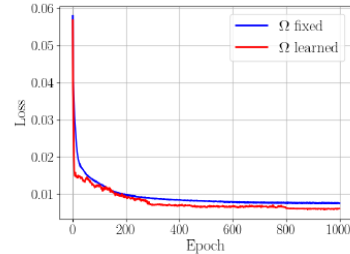
- Even for fixed (not trained) code the **proposed neural decoder** outperforms the state-of-the-art neural decoder.
- The **end-to-end optimization** of the code improves the performance by very large margins.

Method	BP			SCL			ECCT			DC-ECCT			E2E DC-ECCT			
	$E_b/N_0$	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
POLAR(32,11)		3.29	3.77	4.21	6.22	8.06	10.28	4.46	5.57	7.01	4.53	5.69	7.08	4.65	5.81	7.28
		3.84	4.71	5.70	6.45	8.37	10.60	6.37	8.12	10.19	6.41	8.09	10.57	<b>6.62</b>	<b>8.57</b>	<b>11.71</b>
POLAR(64,32)		3.53	4.02	4.45	7.24	9.74	12.91	4.56	5.93	7.75	4.40	5.76	7.67	4.72	6.22	8.13
		4.29	5.35	6.45	<b>8.16</b>	<b>10.73</b>	<b>13.98</b>	7.19	9.70	13.33	7.45	10.49	13.74	7.59	10.38	13.09
BCH(31,16)		4.59	5.87	7.57	NA	NA	NA	4.61	5.97	7.69	4.97	6.56	8.54	5.30	6.89	9.09
		5.12	6.87	9.27				6.37	8.32	10.63	7.07	<b>9.69</b>	12.54	<b>7.19</b>	9.08	<b>12.84</b>
BCH(63,45)		4.07	4.92	6.03	NA	NA	NA	4.59	6.07	8.13	4.54	6.07	8.14	4.98	6.69	8.97
		4.35	5.60	7.24				6.25	8.78	12.45	6.08	8.64	12.41	<b>6.37</b>	<b>9.09</b>	<b>13.12</b>
LDPC(49,24)		5.25	7.15	9.86	NA	NA	NA	4.21	5.32	6.56	4.08	5.29	6.55	4.95	6.46	8.33
		6.09	8.75	11.91				5.34	6.43	7.21	5.57	6.55	7.21	<b>6.28</b>	<b>8.77</b>	<b>12.33</b>
RS(60,52)		4.43	5.32	6.43	NA	NA	NA	4.37	5.11	6.03	5.04	6.68	8.82	5.12	6.80	9.02
		4.69	6.43	7.56				4.37	5.13	6.04	<b>5.61</b>	<b>7.59</b>	9.82	<b>5.61</b>	7.56	<b>9.90</b>

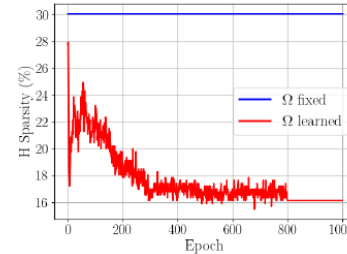


# Analysis

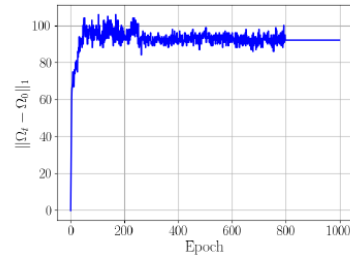
Code	Baseline			Standard form			Random $H$			Random $\Omega$			E2E $\Omega$			
	$E_b/N_0$	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
(31,16)		4.59	5.87	7.57	3.97	4.67	5.57	3.23	3.72	4.36	4.22	5.01	5.86	<b>6.13</b>	<b>7.95</b>	<b>9.90</b>
		5.12	6.87	9.27	4.64	5.89	7.37	3.42	4.18	5.33	4.89	6.14	7.38	<b>6.42</b>	<b>8.31</b>	<b>10.24</b>
(63,45)		4.07	4.92	6.03	4.37	5.33	6.42	3.72	4.26	4.93	3.96	4.61	5.41	<b>5.55</b>	<b>7.33</b>	<b>9.23</b>
		4.35	5.60	7.24	4.78	6.30	8.22	3.93	4.79	5.93	4.27	5.29	6.69	<b>6.15</b>	<b>8.60</b>	<b>11.32</b>
(60,52)		4.43	5.32	6.43	4.60	5.65	6.94	4.41	5.32	6.45	4.60	5.63	6.93	<b>4.73</b>	<b>5.79</b>	<b>7.01</b>
		4.69	5.95	7.56	4.83	6.21	8.00	4.65	5.90	7.53	4.83	6.19	7.97	<b>4.99</b>	<b>6.45</b>	<b>8.35</b>
(64,32)		3.53	4.02	4.45	3.82	4.37	5.03	2.92	3.37	3.73	3.34	3.90	4.64	<b>6.93</b>	<b>9.49</b>	<b>12.51</b>
		4.29	5.35	6.45	4.81	5.74	6.71	2.96	3.54	4.17	3.59	4.53	5.87	<b>7.69</b>	<b>10.04</b>	<b>12.94</b>



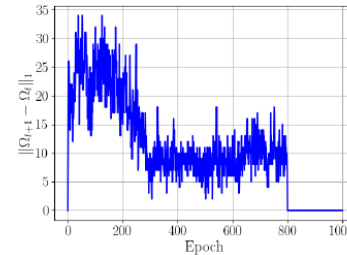
(a)



(b)



(c)



(d)

Code	(31,16)			(32,11)			
	$E_b/N_0$	4	5	6	4	5	6
Our		5.16	6.51	8.19	4.40	5.54	7.04
Mask V2		4.83	6.18	7.85	4.04	5.05	6.28
STE		5.11	6.53	8.20	4.63	5.89	7.39
Mask S.G.		4.54	5.81	7.50	4.26	5.39	6.79
Random $m$		5.04	6.44	8.13	4.38	5.50	6.84

- **BP on learned codes:**

The learned codes outperform other codes under the BP decoder by very large margins, *even if the code is presented in standard form.*

Our method seems to provide **good codes** in a broader sense.

- **Training Dynamics**

The encoder-decoder models enable faster and better training.

Fast change during first stages.

The learned code is *generally sparse*.

- **Ablation Study:**

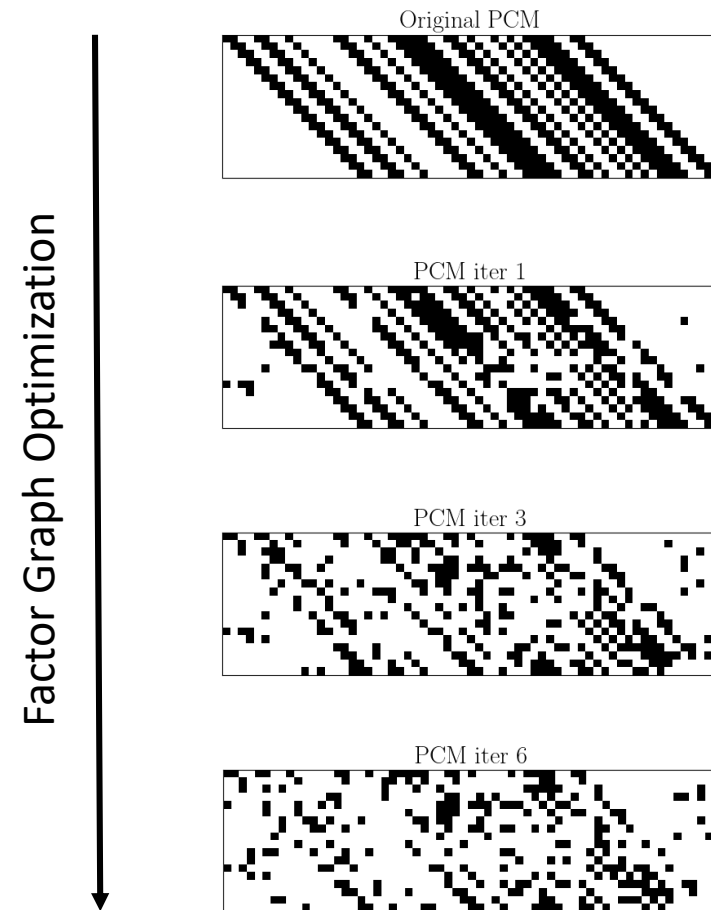
We show that the different elements of the **architecture** and **training** are crucial for performant results.





***How can we design optimal codes according to  
Belief Propagation decoders' inductive bias?***

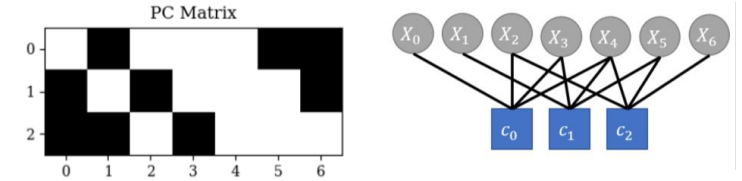
# Factor Graph Optimization of Error-Correcting Codes for Belief Propagation Decoding *sub to Neurips24*



# Belief Propagation Decoding

- A **factor graph** is a representation of a discrete probability distribution that takes advantage of conditional independencies between variables to make the representation *more compact*

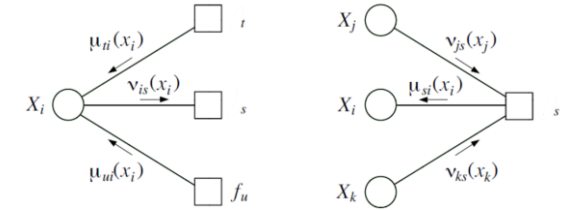
$$p(x_1, \dots, x_n) = \frac{\prod_{i=1}^{n-k} \phi_i(x_{S_i})}{Z}$$



- Belief-Propagation** (Sum-product algorithm) allows efficient marginal inference  $p(x_i)$  via variable elimination as *message-passing* over the factor graph.

$$v_{X_i \rightarrow \phi_t}(x_i) = \prod_{\phi_k \in N(i) \setminus \{\phi_t\}} \mu_{\phi_k \rightarrow X_i}(x_i)$$

$$\mu_{\phi_t \rightarrow X_i}(x_i) = \sum_{x_{S_t \setminus \{i\}}} \phi_s(x_{S_t}) \prod_{j \in S_t \setminus \{i\}} v_{X_j \rightarrow \phi_t}(x_j)$$

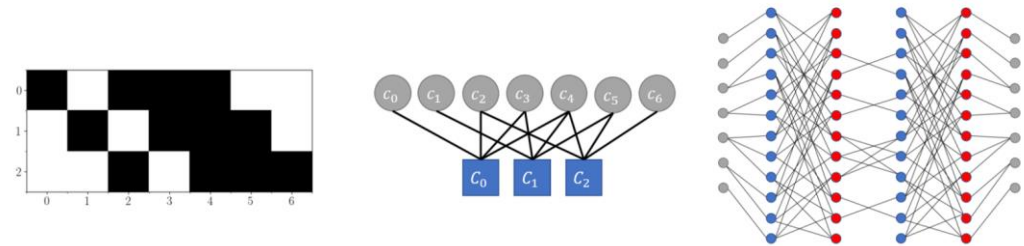


- For non-tree Bayesian graphs, the inference is not tractable. Thus, (**loopy**) belief-propagation is performed **iteratively** (until convergence)

$$P(x_i) = \phi_i(x_i) \prod_{\phi_k \in N(i)} \mu_{\phi_k \rightarrow X_i}^T(x_i)$$

- Belief-propagation *decoding* is generally represented as a **Trellis** graph unrolling of the factor/**Tanner** graph (*log-likelihoods*).

- $L_v = \log \left( \frac{\Pr(c_v = 1|y_v)}{\Pr(c_v = 0|y_v)} \right)$
- $x_e^{2k+1} = x_{(c,v)}^{2k+1} = L_v + \sum_{e' \in N(v) \setminus \{(c,v)\}} x_{e'}^{2k}$
- $x_e^{2k} = x_{(c,v)}^{2k} = 2 \operatorname{arctanh} \left( \prod_{e' \in N(c) \setminus \{(c,v)\}} \tanh \left( \frac{x_{e'}^{2k-1}}{2} \right) \right)$
- $o_v = L_v + \sum_{e' \in N(v)} x_{e'}^{2L}$



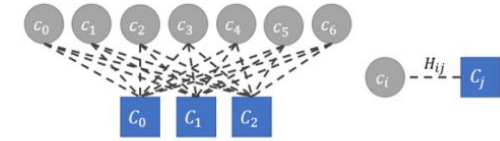
# Belief Propagation Codes

- We wish to obtain BP-optimized codes by solving the following **factor graph learning** (*structure learning*) objective

$$H^* = \arg \min_{H \in \{0,1\}^{(n-k) \times n}} \mathbb{E}_{m \sim \text{Bern}^k(1/2), \varepsilon \sim \mathcal{Z}, T \in \mathbb{N}_+} \mathcal{D} \left( f_{H,T}(\phi(G(H)m) + \varepsilon), m \right) + \mathcal{R}(H)$$

- $f_{H,T}$  is the **BP decoder** built upon  $H$  with  $T$  iterations
  - $\phi(\cdot)$  denotes the codeword **modulation**
  - $\mathcal{D}$  is the **metric** of interest and  $\mathcal{R}$  denotes a potential **regularization** of interest (e.g., sparsity or code structure)
- 
- Multiple challenges arise
    1. The optimization is highly **non-differentiable** (NP-hard binary non-linear integer programming)
    2.  $x = G(H)m = Gm$  is both **highly non-differentiable** (matrix-vector multiplication over  $GF(2)$ ) and **computationally expensive** (inverse via Gaussian elimination of  $H$ )
    3. The modulation  $\phi(\cdot)$  can be non-differentiable
    4. BP assumes a **fixed code** (i.e., the factor graph edges) upon which the decoder is implemented.

# Structure Learning via *Tensor Belief Backpropagation*



- BP decoders are generally implemented using *sparse graphs* via the *Trellis graph* depiction.
- We propose a *Tanner graph learning* approach, where the *bipartite graph* is assumed as **complete with binary weighted edges**.
- The two alternating stages of BP can be represented in a **differentiable matrix form** rather than its *static graph* formulation

$$- Q_{ij} = L + \sum_{j' \in C_i \setminus j} R_{j'i} \equiv L + \sum_{j'} R_{j'i} H_{j'i} - R_{ji}$$

$$- R_{ji} = 2 \operatorname{arctanh} \left( \prod_{i' \in V_j \setminus i} \tanh \left( \frac{Q_{i'j}}{2} \right) \right)$$

$$= 2 \operatorname{arctanh} \left( \prod_{i'} \left( \tanh \left( \frac{Q_{i'j} H_{j'i'}}{2} \right) + (1 - H_{j'i'}) \right) / \left( \tanh \left( \frac{Q_{ij}}{2} \right) + \epsilon \right) \right)$$

- $(1 - H) \in \{0, 1\}^{(n-k) \times n}$  represents the identity element of multiplication

- BP remains differentiable with respect to  $\mathbf{H}$  as a composition of differentiable functions

---

## Algorithm 1: Tensor Belief Propagation

---

```

1 function BP(llr, H, iters, eps=1e-7)
  // llr is the batched LLR matrix, (B, n)
  // H is the binary parity-check matrix, (n-k,n)
  // iters is the number of BP iterations
2 H = H.unsqueeze(dim=0).T
3 C = llr.unsqueeze(dim=-1)
4 for t in range(iters) do
5     Q = C if t == 0 else C + sum(R*H,dim=-1).unsqueeze(dim=-1) - R
6     tmp = tanh(0.5*Q)
7     R = 2*atanh( prod(tmp*H+(1-H),dim=1)/(tmp+eps) )
8 return C.squeeze()+sum(R*H,dim=-1)

```

---

# Belief Propagation Codes Optimization

1. Highly non-differentiable optimization
2.  $x = G(H)m = Gm$
3. ~~Non-differentiable modulation  $\phi(\cdot)$~~
4. BP assumes a fixed code

- The tensor reformulation solves the major *graph learning challenge* (challenge 4).
- For any given  $H$  the conditional independence of error probability under symmetry [1] is satisfied for message passing algorithms
  - It is enough to optimize the **zero codeword only** i.e.,  $c = Gm = 0$  (challenge 2)
  - As a by-product, the optimization problem is **invariant to the choice of modulation** (challenge 3)

- To optimize  $H$  (challenge 1) we relax the NP-hard binary programming problem to an unconstrained parameterized problem
  - Defining  $\Omega \in \mathbb{R}^{(n-k) \times n}$  we have  $H := H(\Omega) = \text{bin}(\Omega)$
  - Implemented via shifted *straight-through estimator* (STE)

$$\text{bin}(u) = (1 - \text{sign}(u))/2, \quad \partial \text{bin}(u)/\partial u = -0.5\mathbb{1}_{|u| \leq 1}$$

- The final objective is given by the empirical risk objective (with  $c_s = \phi(c_0)$  denotes the *modulated zero* codeword)

$$\mathcal{L}(\Omega) = \sum_{t=1}^T \sum_{i=1}^n \text{BCE} \left( f_{\text{bin}(\Omega), t}(c_s + \varepsilon_i), c \right) + \mathcal{R}(\text{bin}(\Omega))$$

- While *highly non-convex*, the objective is (sub)differentiable and thus optimizable via classical **first-order methods**
  - Since  $H$  is binary, only changes in the **sign** of  $\Omega$  are relevant for the optimization.
  - Given the gradient  $\nabla_{\Omega} \mathcal{L}$  computed on sufficient statistic we propose a binary programming aware **line-search procedure** limiting the number of steps to up to  $n(n-k) \ll 2^{n(n-k)}$

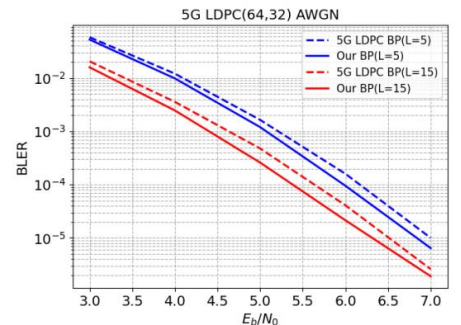
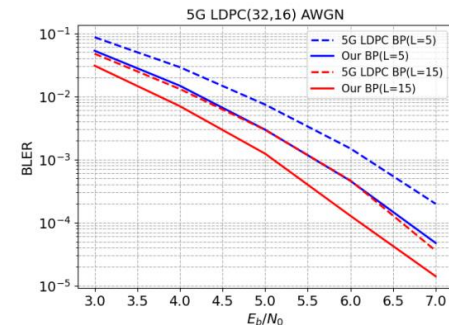
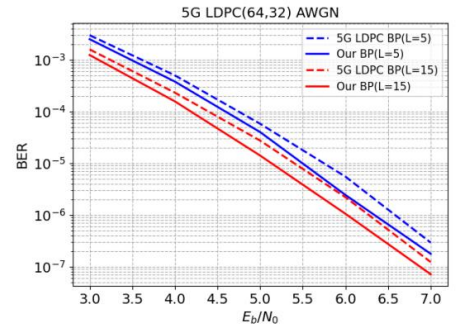
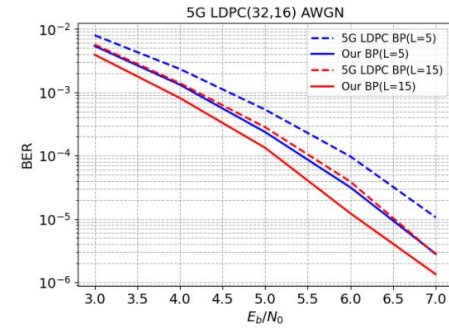
$$\lambda^* = \arg \min_{\lambda \in \mathcal{I}_{\Omega}} \mathcal{L}(\Omega - \lambda \nabla_{\Omega} \mathcal{L}), \quad \mathcal{I}_{\Omega} = \left\{ s_i = \left( \frac{\Omega}{\nabla_{\Omega} \mathcal{L}} \right)_i \mid s_i > 0 \right\}$$

# Experiments

- Our method improves by *large margins* all code families on the three different **channel noise** scenarios and with various number of **decoding iterations** (L=5,15; first and second row).

Channel Method $E_b/N_0$	AWGN						Fading						Bursting					
	BP			Our			BP			Our			BP			Our		
	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
BCH(63,45)	4.06	4.91	6.04	<b>5.44</b>	<b>6.93</b>	<b>8.60</b>	3.09	3.46	3.90	<b>3.96</b>	<b>4.58</b>	<b>5.27</b>	3.60	4.32	5.19	<b>4.05</b>	<b>5.07</b>	<b>6.27</b>
	4.21	5.24	6.59	<b>5.70</b>	<b>7.35</b>	<b>9.16</b>	3.13	3.55	4.04	<b>4.10</b>	<b>4.80</b>	<b>5.56</b>	3.67	4.52	5.59	<b>4.21</b>	<b>5.40</b>	<b>6.85</b>
CCSDS(128,64)	6.46	9.61	13.99	<b>7.34</b>	<b>10.48</b>	<b>14.37</b>	5.72	7.42	9.47	<b>6.73</b>	<b>8.45</b>	<b>10.45</b>	5.29	7.81	11.25	<b>6.23</b>	<b>8.80</b>	<b>11.90</b>
	7.32	10.83	15.43	<b>8.61</b>	<b>12.26</b>	<b>16.00</b>	6.43	8.29	10.28	<b>8.05</b>	<b>10.07</b>	<b>12.37</b>	5.98	8.85	12.53	<b>7.39</b>	<b>10.43</b>	<b>13.28</b>
LDPC(121,60)	4.81	7.17	10.75	<b>7.70</b>	<b>10.87</b>	<b>14.25</b>	4.10	5.23	6.68	<b>6.68</b>	<b>8.47</b>	<b>10.50</b>	3.97	5.75	8.40	<b>6.23</b>	<b>8.89</b>	<b>11.98</b>
	5.31	7.96	11.85	<b>8.86</b>	<b>11.91</b>	<b>14.41</b>	4.42	5.61	7.04	<b>7.71</b>	<b>9.67</b>	<b>11.76</b>	4.31	6.37	9.25	<b>7.26</b>	<b>10.03</b>	<b>12.88</b>
LDPC(121,80)	6.59	9.68	13.43	<b>7.77</b>	<b>11.21</b>	<b>15.06</b>	4.60	5.80	7.22	<b>5.55</b>	<b>6.90</b>	<b>8.36</b>	5.30	7.60	10.66	<b>6.23</b>	<b>8.87</b>	<b>12.19</b>
	7.35	10.94	15.46	<b>8.75</b>	<b>12.45</b>	<b>15.67</b>	4.97	6.29	7.82	<b>6.25</b>	<b>7.80</b>	<b>9.47</b>	5.81	8.50	12.15	<b>6.99</b>	<b>10.09</b>	<b>13.74</b>
LDPC(128,64)	3.66	4.65	5.80	<b>5.54</b>	<b>7.37</b>	<b>9.44</b>	3.22	3.80	4.44	<b>4.86</b>	<b>5.94</b>	<b>7.15</b>	3.23	4.08	5.09	<b>3.72</b>	<b>5.00</b>	<b>6.54</b>
	4.00	5.16	6.42	<b>6.56</b>	<b>8.70</b>	<b>10.81</b>	3.51	4.18	4.84	<b>5.64</b>	<b>6.85</b>	<b>8.14</b>	3.48	4.51	5.66	<b>4.13</b>	<b>5.72</b>	<b>7.66</b>
LDPC(32,16)	4.36	5.59	7.18	<b>5.48</b>	<b>7.02</b>	<b>8.92</b>	4.03	4.70	5.47	<b>5.26</b>	<b>6.02</b>	<b>6.82</b>	3.88	4.89	6.18	<b>4.77</b>	<b>6.02</b>	<b>7.52</b>
	4.64	6.07	7.94	<b>5.76</b>	<b>7.44</b>	<b>9.41</b>	4.29	5.06	5.90	<b>5.43</b>	<b>6.23</b>	<b>6.97</b>	4.09	5.26	6.76	<b>5.01</b>	<b>6.35</b>	<b>7.96</b>
LDPC(96,48)	6.73	9.48	12.98	<b>7.22</b>	<b>9.96</b>	<b>13.37</b>	3.83	4.57	5.35	<b>5.37</b>	<b>6.51</b>	<b>7.71</b>	5.68	7.94	10.90	<b>5.90</b>	<b>8.19</b>	<b>10.91</b>
	7.50	10.61	<b>14.26</b>	<b>8.29</b>	<b>11.12</b>	14.06	4.17	4.94	5.73	<b>6.14</b>	<b>7.38</b>	<b>8.65</b>	6.33	8.91	<b>11.99</b>	<b>6.71</b>	<b>9.28</b>	11.75
LTE(132,40)	2.94	3.32	3.57	<b>3.25</b>	<b>3.71</b>	<b>4.04</b>	3.17	3.45	3.67	<b>4.49</b>	<b>4.99</b>	<b>5.47</b>	2.75	3.17	3.47	<b>2.99</b>	<b>3.44</b>	<b>3.78</b>
	3.37	3.79	4.09	<b>3.93</b>	<b>4.49</b>	<b>4.89</b>	3.60	3.82	4.01	<b>5.32</b>	<b>5.81</b>	<b>6.31</b>	3.17	3.62	3.96	<b>3.53</b>	<b>4.03</b>	<b>4.41</b>
MACKAY(96,48)	6.75	9.45	<b>12.85</b>	<b>7.03</b>	<b>9.63</b>	12.78	6.28	7.86	9.55	<b>6.53</b>	<b>8.06</b>	<b>9.77</b>	5.72	7.97	10.81	<b>5.95</b>	<b>8.23</b>	<b>10.91</b>
	7.59	10.52	<b>14.09</b>	<b>7.99</b>	<b>10.97</b>	14.05	7.04	8.76	10.64	<b>7.47</b>	<b>9.32</b>	<b>11.19</b>	6.39	8.90	11.91	<b>6.82</b>	<b>9.41</b>	<b>12.71</b>
POLAR(128,86)	3.76	4.17	4.58	<b>4.83</b>	<b>5.87</b>	<b>6.58</b>	3.15	3.53	3.91	<b>3.64</b>	<b>4.28</b>	<b>4.94</b>	3.48	3.96	4.37	<b>3.69</b>	<b>4.51</b>	<b>5.18</b>
	4.02	4.67	5.38	<b>5.37</b>	<b>6.88</b>	<b>8.10</b>	3.28	3.73	4.18	<b>3.92</b>	<b>4.70</b>	<b>5.52</b>	3.65	4.31	4.97	<b>3.87</b>	<b>4.91</b>	<b>5.91</b>
RS(60,52)	4.41	5.32	6.41	<b>5.02</b>	<b>6.38</b>	<b>7.99</b>	3.11	3.41	3.77	<b>3.37</b>	<b>3.73</b>	<b>4.12</b>	3.85	4.58	5.44	<b>4.17</b>	<b>5.18</b>	<b>6.40</b>
	4.54	5.52	6.64	<b>5.07</b>	<b>6.47</b>	<b>8.12</b>	3.13	3.43	3.81	<b>3.38</b>	<b>3.75</b>	<b>4.15</b>	3.91	4.72	5.67	<b>4.21</b>	<b>5.27</b>	<b>6.56</b>
LDPC PGE2(64,32)	4.38	5.12	6.04	<b>4.45</b>	<b>5.19</b>	<b>6.10</b>	4.08	4.44	4.81	<b>4.10</b>	<b>4.46</b>	<b>4.85</b>	4.07	4.69	5.43	<b>4.07</b>	<b>4.70</b>	<b>5.43</b>
	4.38	5.13	6.04	<b>4.44</b>	<b>5.19</b>	<b>6.10</b>	4.08	4.44	4.81	<b>4.10</b>	<b>4.47</b>	<b>4.85</b>	4.06	4.69	5.43	<b>4.06</b>	<b>4.69</b>	<b>5.44</b>
LDPC PGE5(64,32)	6.02	8.20	10.95	<b>6.53</b>	<b>8.73</b>	<b>11.56</b>	5.63	6.86	8.31	<b>6.22</b>	<b>7.48</b>	<b>8.82</b>	5.18	6.97	9.34	<b>5.59</b>	<b>7.41</b>	<b>9.51</b>
	6.63	9.06	<b>12.30</b>	<b>7.13</b>	<b>9.48</b>	12.20	6.19	7.52	9.02	<b>6.96</b>	<b>8.34</b>	<b>9.85</b>	5.68	7.75	<b>10.19</b>	<b>6.12</b>	<b>8.06</b>	10.13
LDPC PGE10(64,32)	3.98	5.17	6.70	<b>5.56</b>	<b>7.22</b>	<b>9.13</b>	3.52	4.18	4.95	<b>5.02</b>	<b>6.00</b>	<b>7.11</b>	3.48	4.47	5.75	<b>4.26</b>	<b>5.50</b>	<b>7.01</b>
	4.27	5.77	7.67	<b>6.25</b>	<b>8.28</b>	<b>10.59</b>	3.71	4.47	5.30	<b>5.60</b>	<b>6.72</b>	<b>7.90</b>	3.67	4.90	6.46	<b>4.73</b>	<b>6.22</b>	<b>8.01</b>

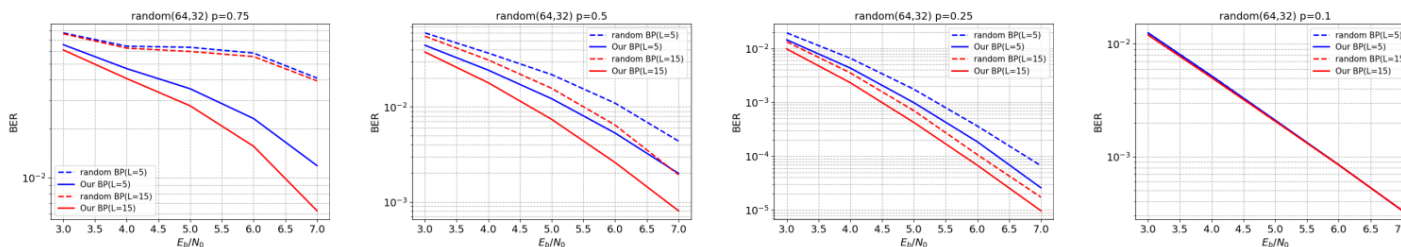
Method $E_b/N_0$	SCL			BP (Polar)			Our (Polar)			BP (5G LDPC)			Our (5G LDPC)		
	4	5	6	4	5	6	4	5	6	4	5	6	4	5	6
(32,16)	6.22	8.06	10.28	4.36	5.59	7.18	5.48	7.02	8.92	6.06	7.53	9.23	<b>6.63</b>	<b>8.53</b>	<b>10.36</b>
	6.45	8.37	10.60	4.64	6.07	7.94	5.76	7.44	9.41	6.57	8.17	10.16	<b>7.11</b>	<b>8.92</b>	<b>11.31</b>
(64,32)	7.36	9.82	12.98	-	-	-	-	-	-	7.59	9.75	12.10	<b>7.87</b>	<b>10.12</b>	<b>12.92</b>
	8.10	10.73	14.00	-	-	-	-	-	-	8.36	10.50	13.02	<b>8.75</b>	<b>11.17</b>	<b>13.76</b>
(128,64)	8.49	11.46	16.16	-	-	-	-	-	-	9.90	13.20	16.73	<b>9.98</b>	<b>13.27</b>	<b>17.02</b>
	9.59	13.12	17.48	-	-	-	-	-	-	12.31	15.98	18.06	<b>12.04</b>	<b>16.18</b>	<b>18.66</b>



# Analysis

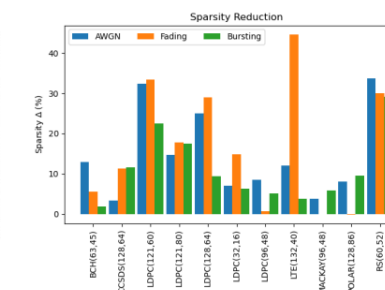
- **Impact of initialization**

- Initialization has a **large impact** the performance (local convergence)



- **Learned regularization**

- Regularization enforces structure and can improve performance. **Sparsity** constraint has **low influence** since it coincides with *BP's inductive bias*.



- **Convergence**

- Optimal step sizes are close to the *working point vicinity*.

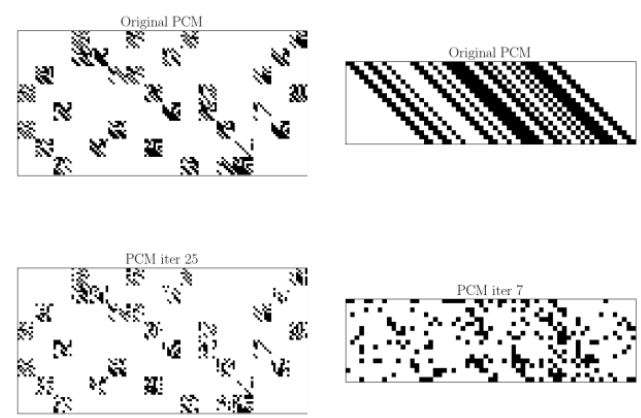
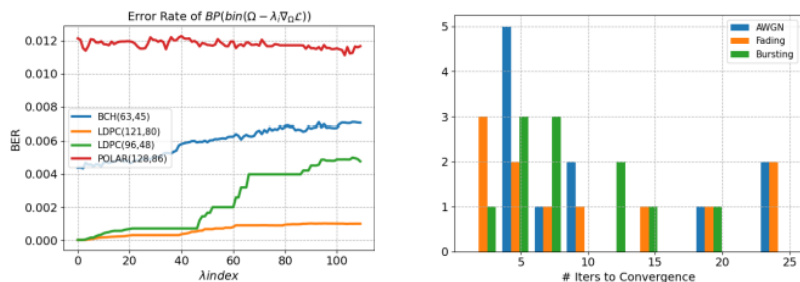


Figure 5: Sparsity reduction of the proposed codes.



# Links

- **Papers:**

- <https://yonilc.github.io/>

- **Code:**

- <https://github.com/yonilc>

- **Correspondence:**

- [choukroun.yoni@gmail.com](mailto:choukroun.yoni@gmail.com)