

## Pruning neural belief propagation decoders

↳ tailor an overcomplete parity-check matrix to BP decoding using machine learning.

① weight = indication of importance of the connected check nodes (CNs) to decoding.

② use them to prune unimportant CNs.

↳ iteration w/  $H$  &  $z$  parity-check matrix  $arg$ .

ML (maximum likelihood decoding)

⇒ achieving near ML performance = computationally complex.

↓  
sol: BP decoding.

LDPC + sparse parity check matrix.

⇒ BP decoding ≈ near optimal performance.

WBP (weighted belief propagation) =

① In each iteration, the VNs and CNs are now referred to as VN layers and CN layers, respectively.

② weights can be introduced at the edges which then are optimized using SGD.

③ each edge has a different weight

\* while WBP decoding improves upon conventional BP decoding, its performance is still limited by the underlying parity-check matrix.

↳ we introduce a pruning-based approach to selecting the best

parity-check equations for each iteration of the BP decoder for short linear block codes.

① starts with a large overcomplete parity-check matrix under WBP decoding.

② consider the weights in the Tanner graph, the magnitude of the weights gives an indication of the importance of the edge in the decoding process.

↳ A magnitude close to zero

indicate that edge has low importance.

• By trying the weights for each CN, reinforcing that the weights of all incoming edges to a single CN are equal, the weights can be interpreted as an indication of importance of the CN in the decoding process.

CNs with connected low-weight edges do not play an important role in the decoding process and can be removed.

↓

by allowing pruning of different CNs in each iteration, the optimization results in a different parity-check equation for each iteration.

각 iteration마다 다른 CN을 pruning을 적용함으로써, 각 iteration마다 다른 parity check equation의 최적화된 결과를 얻는다.

★ The weights in the corresponding Tanner graph can be tied, leading to the largest complexity.

$$\lambda_{CN}^{(l)} = 2 \cdot w_c^{(l)} \cdot \tanh^{-1} \left( \prod_{\tilde{v} \in N(c)} \tanh \left( \frac{1}{2} \lambda_{\tilde{v} \rightarrow c}^{(l)} \right) \right)$$

↓

modified WBP: weights are tied at the CNs [same weight  $w_c^{(l)}$ ]

•  $w_c^{(l)} \in [0, 1]$  = how much the CN contributes to decoding.

• A large magnitude indicates high importance whereas a magnitude of zero indicates that the CN is irrelevant to the decoding process.

1. Initialize

↑ iteration of PCM

$$H = \{H_1, H_2, \dots, H_L\}, w_c = \{1, 1, \dots, 1\}$$

weight

2. Run one batch of gradient descent to optimize  $w$ .

Not

3. converged or max. number of batches?

↓ Yes

4. Find the min. check node weight and remove the check node

↓

5. Target complexity or performance reached? → No: (2)

↓ Yes.

decoder  $D_1 = N_{opt}, W_{opt}$

$D_2 = N_{opt}$ , all weight = 1

$D_3 = N_{opt}$ , untied optimized weights over all iterations.

The weights in  $w$  are then optimized using the Adam optimizer within the TensorFlow programming framework.

After the optimization has converged, we find the index & iteration of lowest CN weight  $w_c^{(l)}$  and set it to zero.

As this may change the optimal value for the remaining weights, we rerun the training.

### neural decoding with optimization of node activations.

#### III. Method

A) the Architecture

B) The knowledge distillation loss term  
C) The sparse node activation Loss term.

↳ Abstract = two novel loss terms on

the node's activations. (노드 활성도에 대한)

새로운 loss terms를 제시되었다.

#### IV. Results

① first loss term = sparse constraint,  
on the node's activations.

Pruning neural belief decoders

↳

remove irrelevant check nodes.

② second loss term tried to mimic the node's activations from a teacher decoder, which has better performance,

③ finding better sparse parity check matrix.

↓

#### 1. introduction

The neural belief propagation (NBP)

was the first deep neural decoder that

↳ NBP > Vanilla BP decoder, provide an improvement over the vanilla belief propagation decoder.

① 이전 연구 = to find better neural architecture to improve the decoder's performance.

② tried to change the loss function to improve the performance.

⇒ learning from the syndrome =

penalizes the neural decoder for producing outputs that do not correspond to

valid codewords, the loss function

↳  $\sum_{i=1}^n \frac{1}{2} \sum_{j=1}^n \hat{v}_{i,j}$  at the output layer, BER  $\frac{1}{2}$ .



① list-based decimation stage.

② learned decimation stage.

weights are tied over iterations.

$$w_{ch,v_i}^{(l)} = w_{ch,v_i}, w_{v_i \rightarrow c_j}^{(l)} = w_{v_i \rightarrow c_j}$$

$$NBP-D(l_{max}, n_D, n_{LP})$$

iterate between decimation process and conventional NBP list-based decimation.

Identify  $VN'$  with the lowest absolute a posteriori LLR  $\mu_{v_i}^{(l_{max})}$

$\therefore$  least reliable LLR  $\frac{n_D}{2}$

$\pm \infty$ : decimate.  
 $= \mu_{ch,v}$

Sign  $\frac{n_D}{2} \times \Rightarrow$  tree build decoding.

$\begin{bmatrix} +\infty \\ -\infty \end{bmatrix}$  two new graph.

② learned decimation stage.

$\hookrightarrow \Theta = \text{sign} \uparrow, \frac{y}{2} \text{ 정도}$

$\Rightarrow$  use an NN to decide to which value each VN should be decimated.

a posteriori LLR in  $l$ th iteration

$$\mu_{v_i}^{(l)} = \mu_{ch,v_i} + \sum_{c \in N(v_i)} \mu_{c \rightarrow v_i}^{(l)}$$

$$\mu_{ch,v} = \mu_{ch,v} + \text{sign}(\mu_v^{(l_{max})})$$

$$f_{NN}(\mu_{ch,v}, \{\mu_{c \rightarrow v}^{(l_{max})} | c \in N(v)\}, \Theta)$$

All trainable parameters are  $VN$  of the NN.

$\hookrightarrow$  같은 NN이 적용되고 weights는 shared.

two-stage decimation

① we build a list by iterating between a conventional NBP decoder

가장 신뢰성 낮은 bits를 추출하고, 기존의

NBP decoder 사이에서 lists를 만들고 and guessing the least reliable bit.

② learned decimation = 어떤 bits는 decimation할지 결정.

$\hookrightarrow$  The second stage iterates between a conventional NBP decoder and learned decimation, where we use a neural network to decide the decimation value for each bit.

$\star$  a posteriori LLR.