Generative Decoding for Quantum Error-correcting Codes

Hanyan Cao, ^{1,2} Feng Pan, ^{1,3,4} Dongyang Feng, ^{1,2,5} Yijia Wang, ^{1,2} and Pan Zhang ^{1,5,*}

¹CAS Key Laboratory for Theoretical Physics, Institute of Theoretical Physics,

Chinese Academy of Sciences, Beijing 100190, China

²School of Physical Sciences, University of Chinese Academy of Sciences, Beijing 100049, China

³Centre for Quantum Technologies, National University of Singapore, 117543, Singapore

⁴Science, Mathematics and Technology Cluster, Singapore University

of Technology and Design, 8 Somapah Road, 487372 Singapore

⁵School of Fundamental Physics and Mathematical Sciences,

Hangzhou Institute for Advanced Study, UCAS, Hangzhou 310024, China

Efficient and accurate decoding of quantum error-correcting codes is essential for fault-tolerant quantum computation, however, it is challenging due to the degeneracy of errors, the complex code topology, and the large space for logical operators in high-rate codes. In this work, we propose a decoding algorithm utilizing generative modeling in machine learning. We employ autoregressive neural networks to learn the joint probability of logical operators and syndromes in an unsupervised manner, eliminating the need for labeled training data. The learned model can approximately perform maximum likelihood decoding by directly generating the most likely logical operators for k logical qubits with $\mathcal{O}(2k)$ computational complexity. Thus, it is particularly efficient for decoding high-rate codes with many logical qubits. The proposed approach is general and applies to a wide spectrum of quantum error-correcting codes including surface codes and quantum low-density paritycheck codes (qLDPC), under noise models ranging from code capacity noise to circuit level noise. We conducted extensive numerical experiments to demonstrate that our approach achieves significantly higher decoding accuracy compared to the minimum weight perfect matching and belief propagation with ordered statistics on the surface codes and high-rate quantum low-density parity-check codes. Our approach highlights generative artificial intelligence as a potential solution for the real-time decoding of realistic and high-rate quantum error correction codes.

I. INTRODUCTION

Quantum computers have the potential to solve practical problems that are intractable for classical computers [1–4]. However, current quantum computing devices are noisy, which restricts their capability. Quantum error correction (QEC) as a crucial step in managing these noises and achieving fault-tolerant quantum computing has become a key research frontier in both theoretical studies [5–7] and hardware developments [8–10] of quantum computation. In QEC, logical states comprising k logical qubits are encoded using n physical qubits, and continuous errors can be digitized into a finite set of discrete errors by measuring the redundant ancilla qubits, resulting in an error syndrome [11, 12]. A decoding algorithm then interprets the error information based on the syndrome and determines an appropriate operation to correct the logical error.

Decoding in the quantum error correction codes is considered more challenging than in its classical counterpart due to the presence of multiple error types and their degeneracy [13]. Consequently, the corresponding Tanner graphs for quantum codes are more complex. For instance, the Tanner graph of Calderbank-Shor-Steane codes [14–17] invariably contains loops of various sizes due to commutation relations [13]. This complexity significantly diminishes the performance of standard decoding algorithms, such as belief propagation, in quantum systems compared to their classical counterparts. Widely applied decoding algorithms include the minimum weight perfect matching (MWPM) algorithm [18–20] and Belief propagation with ordered statistics post-processing (BPOSD) [21–24]. Both methods fall within the minimum weight decoding category, which identifies the most probable error rather than the most probable logical error for given syndromes.

Another category of decoding algorithms, maximum likelihood decoding (MLD), addresses the issues by identifying the most likely logical operator among all 4^k possibilities. In principle, MLD can provide optimal decoding results [9, 25–27], however, this computational task belongs to the #P-hard complexity class and remains highly challenging [28]. Recently, efficient approximate implementations of MLD usually employ tensor networks, such as the boundary matrix product state method [25, 26] for surface code with code-capacity noise [17, 29]. However, due to its high computational cost and the limited accuracy of approximate tensor-network contractions, tensor-network-based MLD

^{*} panzhang@itp.ac.cn

poses significant challenges for quantum codes with complex topologies. For instance, this applies to three-dimensional codes [30], codes with long-range interactions such as quantum low-density parity-check (qLDPC) codes [5, 6], and codes subject to circuit-level noise. Another limitation of existing maximum-likelihood decoders is the requirement to compute probabilities for all 4^k logical operators corresponding to k logical qubits, a task that becomes computationally infeasible for large k.

In this work, we introduce a neural MLD framework that overcomes the restrictions of arbitrary topology and the 4^k computational complexity of canonical maximum likelihood decoding methods. It is inspired by unsupervised generative modeling developed recently in machine learning [31–34]. Our approach employs autoregressive neural networks to model the joint distribution of logical operators given a syndrome, using the product of conditional probabilities. The parameters of the neural network are learned from samples of the error model or experimental data without needing labels. The decoding process involves sequentially generating logical operators one by one for each logical qubit using the learned conditional probabilities, akin to generating chat text word by word in large language models [35]. We term it generative neural decoder (GND for short). Various autoregressive neural networks can be incorporated into our generative framework, including masked autoregressive network for density estimation (MADE) [32] and transformers [33, 36].

In the past few years, various neural network (NN) decoders have been introduced [37–43]. These approaches typically frame decoding problems as classification tasks, utilizing neural networks such as multilayer perceptrons, convolutional neural networks, recurrent networks, and Transformers as classifiers to address the classification task. Recently, some neural network decoders have been utilized to decode the surface code under the circuit-level noise [44, 45] and demonstrate significant advantages over classical decoding algorithms. However, classifier-based NN decoders usually work for a small number of logical qubits, particularly with k=1. Because they face a challenge that the number of classification labels scales exponentially as 4^k [43]. Some NN decoders resolve the 4^k complexity issue by treating 4^k logical classifications as 2k binary-classifications [37, 39, 44]. Effectively, this approach employs an approximation known as the mean-field approximation to approximate the joint probability distribution of the 4^k logical operators using a product distribution, which approximation introduces additional errors. In comparison, the generative decoding scheme proposed in this work naturally works for k>1 logical qubits and does not rely on the mean-field approximations, because it employs the autoregressive modeling [46] to represent the joint distribution as a product of conditional distributions. In other words, rather than simply mapping syndromes to error functions, our approach utilizes generative neural networks to model the joint distribution of logical operators given a syndrome, and sequentially generate the most probable logical operators, thereby circumventing the need for 4^k labels or 4^k computational complexity.

In the following text, we will first introduce our generative decoding scheme, then present numerical experiments on the surface code [17, 47], Bivariate Bicycle (BB) code [7], and qLDPC codes [48] under various error models to demonstrate the superiority of our generative decoding approach.

II. RESULTS

A. The generative decoding scheme

Consider a [[n,k,d]] quantum code, where k logical qubits are encoded using n physical qubits, and the code distance is d. A state $|\psi\rangle$ in the logical space is encoded as a fixed point of operators $\{S\}$ with eigenvalue equal to 1, i.e. $S|\psi\rangle = |\psi\rangle$. These operators form an Abelian subgroup of the n-qubit Pauli group (ignoring the global phase) $\mathcal{P}_n = \{I, X, Y, Z\}^{\otimes n}$, known as the stabilizer group. Apparently, applying them to the quantum circuit does not influence the encoded states. This group $S = \langle g_1, g_2, \cdots, g_m \rangle$ has m = n - k generators, referred to as stabilizer generators in the stabilizer formalism [11, 12]. When an error occurs on the state $|\psi\rangle$, it is modeled as the application of an error operator $E \in \mathcal{P}_n$. Most errors do not commute with the stabilizers, yielding -1 eigenvalue upon measurement. As illustrated in Fig. 1, these errors can be detected by measuring the ancilla qubits associated with the stabilizer generators, producing a syndrome $\gamma(E) = \{\gamma_1(E), \gamma_2(E), \cdots, \gamma_m(E)\}$, where $\gamma_i(E) = 0$ if g_i and E commute, and $\gamma_i(E) = 1$ if they anti-commute.

Whether and how the error E affects the encoded states depends on its associated logical operation. For instance, a quantum code encoding a single logical qubit has four possible logical operators L_I, L_X, L_Z, L_Y . An error E belongs to one of the four corresponding logical sectors, determined by its commutation relations with the logical operators. By ignoring the global phase, we only need to consider L_X and L_Z , the 4 logical sectors are characterized by commutation relations with these two operators. This can be represented using 2 binary variables $\beta = \{\beta^X(E) \text{ and } \beta^Z(E)\}$, indicating whether E commute with L_X and L_Z . The formulation above can be generalized to codes with k logical qubits with 4^k logical operators, by using 2^k binary variables for logical sectors $\{\beta_1, \dots, \beta_k\} = \{\beta_1^X(E), \beta_1^Z(E), \dots, \beta_k^X(E), \beta_k^Z(E)\}$. Then we can see that the task of decoding is to categorize logi-

cal sectors β given syndrome γ , that is, model the conditional probability $p(\beta|\gamma)$ with $\beta \in \{0,1\}^{2k}$ and $\gamma \in \{0,1\}^m$. Given an error model, samples of the joint distribution $p(\beta,\gamma)$ can, in principle, be directly drawn from the model. However, accurately computing the conditional probability distribution

$$p(\beta|\gamma) = \frac{p(\beta,\gamma)}{\sum_{\beta} p(\beta,\gamma)}$$

for a large number of logical qubits is challenging, as naively computing the normalization factor (i.e., the partition function in statistical physics) has a computational complexity of 2^{2k} . In the exsiting neural network decoders[38, 45, 49–51], the number of logical qubits k is usually considered to be small (particularly k=1 in the surface code) and thus this challenge is often ignored.

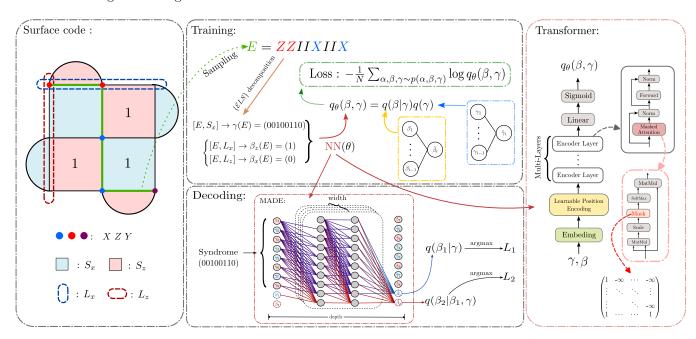


FIG. 1. Illustration of GND's training and decoding processes, using a distance-3 surface code as an example. The red and blue plaquettes represent the X and Z stabilizer generators, respectively, acting on the qubits located at the lattice vertices they cover. The training process begins by sampling errors from the error model, which are then digitalized using the \mathcal{ELS} decomposition to obtain the corresponding configurations. These configurations are fed as input to the autoregressive neural network, and the output of the network is used to compute the negative log-likelihood (NLL), followed by a gradient descent optimization step. The observed syndrome is fed into the neural network during decoding to obtain the first logical variable. This variable is then fed back into the network along with the syndrome until the final logical variable is determined. The figure also illustrates the specific structures of the two autoregressive neural networks employed in this work, MADE and causal Transformer.

In this work, we propose to approach $p(\beta, \gamma)$ by learning a structured variational distribution known as *autoregressive* model. The model factorizes the variational joint distribution $q_{\theta}(\beta, \gamma)$ into the product of conditional probabilities.

$$q_{\theta}(\beta, \gamma) = \prod_{i=1}^{2k} q(\beta_i | \beta_{j < i}, \gamma) \cdot \prod_{i=1}^{m} q(\gamma_i | \gamma_{j < i}) . \tag{1}$$

Here, the subscript j < i denotes variables in front of variable i given an order of variables. The i'th logical variable β_i depends solely on the syndrome γ and logical variables before it $\beta_{j < i}$. The order of variables is given naturally in this case and can be considered as a predefined arrow of time. In this sense, the variables in front of i are the history of i, and variables behind of i are i's future. So the states of i only depend on its history j < i, not on its future j > i. This restriction is essential in autoregressive modeling and is sometimes called "casual". Autoregressive models encompass various renowned implementations, including recurrent neural networks (RNNs) [31], long short-term memory (LSTM) [34], and causal transformers [33, 36], all of which can be integrated into our GND framework. In this work, we consider two concrete examples of autoregressive networks: the masked autoregressive network for density estimation (MADE) [32], which has applications in statistical mechanics [52], and the Transformer [33, 36],

a standard model in natural language processing for generating subsequent words based on a given sentence. Their structures are illustrated in Fig.1.

The autoregressive model q_{θ} is trained to approximate the data distribution $p(\beta, \gamma)$. This is achieved by learning the parameters θ to minimize the distance between the $p(\beta, \gamma)$ and $q_{\theta}(\beta, \gamma)$. We assume that the samples of the true distribution $p(\beta, \gamma)$ can be obtained efficiently, either using an error model or from experiments. In this work, we choose the forward Kullback-Leibler divergence

$$D_{KL}(p||q_{\theta}) = \sum_{\beta,\gamma} p(\beta,\gamma) \log \left[\frac{p(\beta,\gamma)}{q_{\theta}(\beta,\gamma)} \right], \tag{2}$$

as the distance measure, this yields a negative log-likelihood loss function

$$\mathcal{L} = \underset{\theta}{\operatorname{arg\,min}} - \sum_{\alpha,\beta,\gamma} p(\alpha,\beta,\gamma) \log q_{\theta}(\beta,\gamma). \tag{3}$$

Here, α denotes the configuration of stabilizer generators, as detailed in Section IV A. Directly optimizing \mathcal{L} using Eq.2 requires samples drawn from $p(\beta, \gamma)$, which is challenging since $p(\beta, \gamma) = \sum_{\alpha} p(\alpha, \beta, \gamma)$. Fortunately, as Eq.3 indicates, we can directly sample from the error model $p(\alpha, \beta, \gamma)$ to train the autoregressive model, eliminating the need to prepare datasets of (β, γ) .

The training process is depicted in Fig. 1. First, errors are sampled from the error model, and the $\{\mathcal{E}, \mathcal{L}, \mathcal{S}\}$ decomposition is applied to each error to obtain the (β, γ) configuration, which serves as training data for the autoregressive neural network. Two specific autoregressive neural networks, the masked autoregressive network for density estimation (MADE) [32] and the Transformer [33, 36], are illustrated. To facilitate learning, we iteratively minimize the loss function in Eq. (3) to reduce the discrepancy between the autoregressive variational distribution and the error distribution. After training, the autoregressive network learns all conditional probabilities as specified in Eq. (1). Decoding proceeds by sequentially generating a logical operator configuration $\beta = \{0,1\}^{2k}$ using the learned conditional probabilities conditioned on the syndrome.

$$\hat{\beta}_i = \underset{\beta_i}{\arg\max} \, q_{\theta}(\beta_i | \beta_1, \cdots, \beta_{i-1}, \gamma_1, \cdots, \gamma_m) \,, \tag{4}$$

This process is analogous to text generation by ChatGPT [35, 53, 54], where text is generated word-by-word given a prompt (which is analogous to syndromes in our setting). In more detail, decoding the *i*th logical variable involves setting the syndrome γ and the logical variable $\beta_{< i}$ as input (other variables are set to zero to maintain consistent input length) and running the forward pass of the neural network to obtain the output, i.e. the conditional probability $p(\beta_i|\gamma,\beta_{j< i})$, then decide the logical operator for *i* using Eq. (4). So for each logical operator, it takes a forward pass of the neural network, so in total requires 2k neural network passes for decoding *k* logical qubits, contrasting sharply with the $\mathcal{O}(4^k)$ complexity of the conventional tensor network-based MLD method. Also, our scheme is capable of handling arbitrary code topologies, easily accommodating any stabilizer generator connectivity.

B. Numerical experiments on high-rate quantum codes

To evaluate our generative decoding approach, we conducted extensive numerical experiments on several quantum codes with k=1 and k>1 logical qubits. We tested both MADE and Transformer as autoregressive networks within our GND framework. Although Transformer has demonstrated superior performance in language modeling [35], it did not exhibit greater capability than MADE in our framework. Consequently, in the main text, we report only the results obtained with MADE. A detailed comparison between MADE and Transformer is provided in the Sec IV C.

We first evaluate GND on the standard rotated surface code (with k=1) under the depolarizing noises. The depolarising noise only considers the equal probability errors on the qubit and simplifies the decoding problem, so we can test GND with a distance up to 13. The comparisons between GND and the MWPM algorithms are shown in Fig. 2. We can see that at any physical error rate, GND gives a lower logical error rate than MWPM, and the improvements are quite significant at low physical error rates.

Next, we examine the codes with k > 1 logical qubits. We first choose the defected surface code under the circuit-level noise. As a variation of the surface code, the defected surface code removes several stabilizers to achieve a higher code rate [55]. This code retains the two-dimensional structure of the surface code and is well-suited for widely used superconducting quantum chips. Experiments on the defected surface code can be conveniently conducted with only minor modifications to the conventional surface code setup, thus serving as an effective method

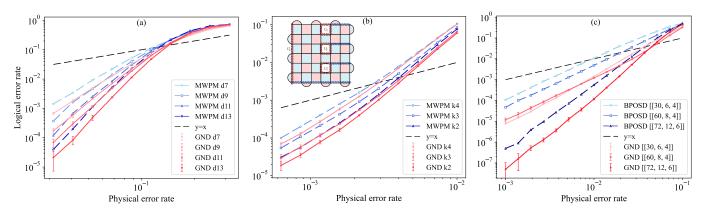


FIG. 2. Performance of our generative neural decoder (GND) compared with the minimum weight perfect matching (MWPM) and the BPOSD algorithm. (a) Rotated surface code with distances up to 13 under depolarizing noise. (b) Defected surface code originally of distance 7. We introduce 3 defects by removing 3 Z-stabilizer generators. Each defect provides an additional logical qubit, for which we establish a valid logical Z operator base represented by red dashed lines. The minimum weight of a logical Z operator is 4, corresponding to the Z-distance. The blue dashed lines represent the corresponding X-operators L_X . (c) High-rate [[60, 8, 4]], [[72, 12, 6]] BB code and [[30, 6, 4]] qLDPC codes under code capacity depolarizing noise model. We used the PYMATCHING [19] for MWPM decoding and the LDPC package [23] for BPOSD decoding (using min-sum BP, the max number of iterations is 1000, and the post-processing method is OSD-CS). For GND, the neural networks are trained at a single physical error rate equal (or close) to the threshold and used for decoding at all physical error rates. For rotated surface codes and [[30, 6, 4]] ldpc code, we used 2×10^9 samples to train 2×10^5 epochs using 10000 batch; for [[60, 8, 4]] BB codes and defected surface codes under curcuit level noise, we used 5×10^9 data to train 5×10^5 epochs using 10000 batch; for [[72, 12, 6]] BB code, we used 1×10^{10} data to train 5×10^5 epochs using 20000 batch. Each data point in (a) and (b) is averaged over 100 decoding tasks with each task comprising 10^5 error instances, and in (c) is averaged over 10 decoding tasks, with each task involving 10^7 error instances.

to evaluate the efficiency of high-rate codes in near-term applications. For the noise, we consider the circuit-level noise model, which is more experimentally realistic than the code capacity noise considered in the last example. The circuit-level noise accounts for errors in actual quantum circuits, including measurements, gate operations, idling, and reset [9, 10, 56, 57]. To correct such noise, repetitive syndrome measurements are required [58], which often introduces redundant information. A more compact representation is provided by the detector error model (DEM) [59], where a detector is defined as the linear exclusive-or (XOR) of multiple measurements. DEM captures discrepancies in consecutive measurements across circuit layers, offering a more concise representation than full syndrome measurement results. In this framework, detector states serve as syndromes γ , while final logical errors correspond to β in our method.

As illustrated in the inset of Fig. 2, we consider a distance-7 surface code with three stabilizers removed. The defected surface code features one Z-operator L_Z^0 of length 7 and three Z-operators L_Z^1, L_Z^2, L_Z^3 of length 4. Consequently, it encodes 4 logical qubits and has a Z-distance of 4. To evaluate decoding performance across different k values, we sequentially removed the stabilizers corresponding to L_Z^1, L_Z^2, L_Z^3 , resulting in three defected surface codes with k=2,3,4. In the numerical experiments, we initialized the circuit in the Z-basis and performed four rounds of syndrome measurements. We utilize the Stim [60] library to construct the DEM, simulate the circuit, and perform sampling. The results, shown in Fig. 2(b), compare the logical error rates of our method GND and MWPM for each k value. The figures demonstrate that GND achieves significantly lower logical error rates than MWPM for all k values across all physical error rates.

Although the surface code has a two-dimensional structure that is convenient to the existing superconducting quantum hardware, it has the disadvantage of la ow code rate R=k/n, the ratio between the number of logical qubits and the number of physical qubits. So it is inefficient to encode a large number of logical qubits. Recently, researchers investigated quantum codes that have a much larger logical rate than the surface code; here we term them as "high-rate codes". Notable examples include Bivariate Bicycle (BB) code [7], and qLDPC codes [48]. The BB codes family, as a specialized class of qLDPC codes, exhibits several advantageous properties. First, compared to commonly used k=1 surface codes, it significantly reduces physical qubit resources while encoding the same number of logical qubits with equivalent accuracy. Additionally, it features a well-designed syndrome measurement circuit with only a few layers of parallel two-qubit gates, maintaining circuit distance and preventing hook errors. Furthermore, due to its specially designed topology, it can be embedded into a two-layer quantum hardware architecture with efficient long-range connectivity, making it a promising candidate for future fault-tolerant quantum computing. To further demonstrate the universality of our GND method for arbitrary topology high-rate codes, we also evaluated small

qLDPC codes. The codes are chosen after a random search for those lacking a regular topological structure, making them ideal candidates for our evaluation. In our numerical experiments, we take the [[60, 8, 4]] and [[72, 12, 6]] BB codes, along with a [[30, 6, 4]] qLDPC code under code capacity depolarizing noise model as illustrative examples to demonstrate the performance of GND.

We note that in exchange for higher encoding rates, the BB codes and qLDPC codes exhibit more complex structures, with each qubit checked by multiple stabilizers. Consequently, when mapping the decoding problem to a matching problem, hyperedges emerge, rendering MWPM inapplicable. Thus, we compare our algorithm with the only currently available decoding algorithm, BPOSD [21–24]. The comparison results are shown in Fig. 2(c). The results demonstrate that our method surpasses BPOSD in logical error rates for both BB codes and the selected qLDPC code. With lower error rates, the gap between GND and BPOSD is larger. Notably, with a physical error rate below 10^{-2} , GND is about 10 times more accurate than BPOSD.

III. DISCUSSIONS

We have presented a neural network decoding framework for quantum codes based on generative modeling in machine learning. In this framework, an autoregressive generative model is employed to approximate the conditional probability of the target logical sector given the syndrome (detector), and trained by minimizing the discrepancy (measured by KL divergence) between the empirical error distribution and the distribution of the autoregressive model. After training, it can generate logical operators sequentially using the learned conditional probabilities, analogous to generating the text in the large language models. Our approach shares advantages common to other NN decoders, such as applicability to arbitrary code topologies, arbitrary error models, and fast neural network computation. The significant advantage of our approach is to avoid computing 4^k probabilities for logical operators. This is because our model factorizes the variational joint distribution encompassing 4^k configurations into a product of 2k conditional probabilities. This greatly reduces the generation of logical operators for k > 1 logical qubits, and the number of neural network computations required for decoding scales linearly with the number of logical qubits. Furthermore, compared to existing NN decoding schemes, our approach avoids using exponentially large classification labels or using product-of-marginal approximation. We have carried out extensive numerical experiments using surface codes, the defected surface code, BB code, and qLDPC codes under both code capacity noise and circuit-level noise. By comparison to the standard decoding algorithms, including MWPM and BPOSD, we have demonstrated the superiority of our approach in accuracy and efficiency.

Within the generative decoding framework, we have experimented with two types of autoregressive neural networks, MADE and Transformer, and found that MADE with fewer parameters performs better. We suspect that we did not invest enough engineering effort to fully unlock the potential of the Transformer, which is considered the state-of-the-art autoregressive model and has demonstrated its power in large language models.

Regarding the decoding speed, the time for decoding in our current implementation is in the order of 10^{-2} seconds. We attribute this to the relatively low GPU efficiency in our implementation, given that the number of parameters, ranging from 10^6 to 10^7 , is still comparatively small when compared with large-scale neural network models. As a consequence, the decoding speed is considerably slower than the operation time of the superconducting qubits, which is about 10^{-6} seconds per round. So it remains a challenge for GND to improve the efficiency through engineering efforts. In future work, we aim to enhance the performance and scalability of our generative decoding framework by engineering efforts such as leveraging multiple GPUs, decreasing the latency using FPGAs, employing lower-precision floating-point numbers in decoding, and reducing model size using pruning, distillation, or quantization.

IV. METHOD

A. The ELS decomposition and the maximum likelihood decoding

The Pauli group of logical qubits P_L is generated by $\mathcal{L} = \langle l_1^x, l_1^z, l_2^x, l_2^z, \cdots, l_k^x, l_k^z \rangle$, where l_i^x and l_i^z denote the logical X and Z operators of the *i*'th logical qubits, respectively. This group commutes with the stabilizer group. Beyond the 2^m stabilizer operators and 4^k logical operators, there exist another 2^m operators that do not commute with the stabilizers; these belong to the pure error subgroup \mathcal{E} , which is Abelian and satisfies the commutation relation $e_i g_j = (-1)^{\delta_{ij}} g_j e_i$. The three subgroups introduced here elucidate the structure of the Pauli group, as manifested in the decomposition $\mathcal{P}_n = \mathcal{E} \otimes \mathcal{L} \otimes \mathcal{S}$ [11, 57], and is termed as $\{\mathcal{E}, \mathcal{L}, \mathcal{S}\}$ decomposition.

The $\{\mathcal{E}, \mathcal{L}, \mathcal{S}\}$ decomposition of the Pauli group tells us that an error E can be mapped to a configuration $\{\alpha, \beta, \gamma\}$ and the error distribution (from the error model) can be equivalently expressed as $p(E) = p(\alpha, \beta, \gamma)$. Here $\alpha \in \{0, 1\}^m$ represents the configuration of m stabilizers, with each value α_i determined by the commutation relation between

E and the pure error generator E_i . The configuration $\beta \in \{0,1\}^{2k}$ denotes the state of the logical X and logical Z operators, while $\gamma \in \{0,1\}^m$ represents the configuration of the m pure error generators, with each value γ_i determined by the commutation relation between E and the stabilizers. So the configuration γ is synonymous with the syndrome. Utilizing this mapping, we can observe the degeneracy of errors; that is, for a given logical configuration β , there are 2^m assignments of $\{\alpha\}$ that yield the same syndrome. MLD computes the likelihood of a logical operator configuration β , which necessitates considering all α configurations

Consider a [[n, k, d]] quantum code with n physical qubits and distance d. A logical state $|\phi\rangle$ with k logical qubits is encoded using a codeword $|\psi\rangle$.

$$\arg\max_{\beta} p(\beta|\gamma) = \arg\max_{\beta} \sum_{\alpha} p(\alpha, \beta|\gamma) . \tag{5}$$

Thus, in the language of the $\{\mathcal{E}, \mathcal{L}, \mathcal{S}\}$ decomposition, MLD considers the total probability of a coset of the stabilizer sub-group \mathcal{S} , rather than the probability of a single error as in minimum-weight decoding [18, 19]. This problem generally belongs to the computational class of #P problems [28], meaning that no algorithms can exactly solve MLD in polynomial time. Specifically, the computation poses two challenges: summing over an exponential number of α configurations and performing this summation for the 4^k possible β configurations [25].

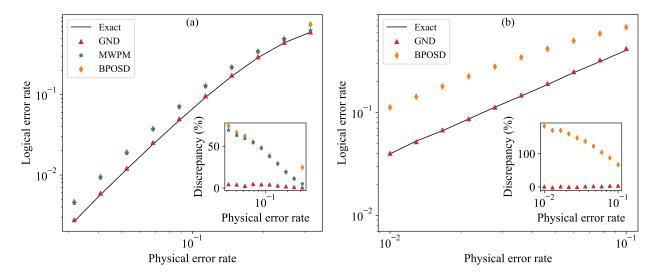


FIG. 3. Logical error rates of our algorithm GND with different physical error rates. (a) Compared with MWPM, the exact MLD algorithm, and BPOSD algorithm on the [[25, 1, 5]] surface code. (b) Compared with an exact MLD algorithm, and BPOSD algorithm on [[18, 4, 4]] BB code. The error model is the depolarizing error model. Each data point in the figures is averaged over 1000000 error instances. The black lines are the optimal maximum-likelihood decoding algorithm, which exactly sums all stabilizer configurations. The insets show the discrepancy between the approximate algorithms and the exact algorithm.

B. Comparing with existing neural network decoder

Recently, numerous neural network-based decoders (NN decoders) have been proposed alongside the rapid development of artificial intelligence. We can mainly group the NN decoders into two categories, the neural network minimum-weight decoders (NNMWD) and the neural network maximum likelihood decoder (NNMLD). The NNMWD infers the most likely error for given syndromes by modeling the error probabilities of each data qubits [38, 49–51]. One issue for NNMWD is that the output errors are not always consistent with the syndrome [49–51]. This is very similar to the issue of belief propagation which usually does not give an error that satisfies the syndrome in quantum codes, hence requiring an additional operation like OSD to make the syndrome satisfied [21–24]. Compared with the existing NNMWD, our algorithm does not have the syndrome-consistency issue, and the performance of our algorithm (which belongs to MLD) is theoretically better than that of NNMWD.

The second category, NNMLD, learns the joint distribution of logical sectors given syndrome $q(\beta|\gamma)$ using neural networks. For k=1, when the number of logical sectors is small, the problem can be converted to a classification

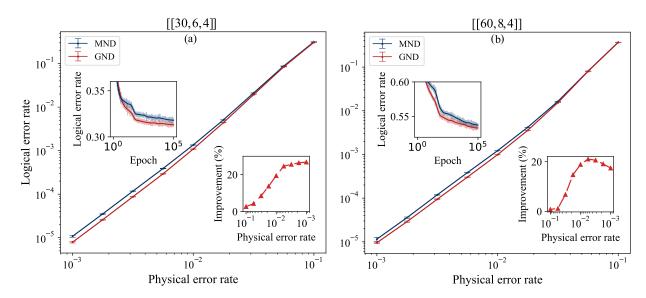


FIG. 4. Comparing GND and MND. The sub-figures (a) and (b) show the decoding performance of GND and MND on [[30, 6, 4]] and [[60, 8, 4]] codes, respectively. The lower right insert figures show the improvement of GND relative to MND. And the upper left insert figures show the variation of the logical error rate with training epoch. The red and blue solid lines represent the average logical error rate of GND and MND, respectively, and shades of the same color indicate variance. MADE and MLP are implemented to construct the decoders. The number of parameters used in GND is 1980972 and 11666828 for [[30, 6, 4]] and [[60, 8, 4]] codes separately. These numbers are 1993600 and 11557728 of MND. The logical error rate with epoch of the [[30, 6, 4]] code is shown in the inset. The decoders are trained at p = 0.1 for [[30, 6, 4]] code and p = 0.12 for [[60, 8, 4]]. Each data point averages over 10 decoding tasks and each task has 10^7 shots.

problem and solved using a neural network classifier. A successful example is the recently proposed AlphaQubit [45] which uses a recurrent neural network and Transformer as a binary classifier (a classifier with two distinct classes for two logical sectors) [45], which successfully solvers the decoding problem of the surface code up to distance 7 under circuit level noise [10]. In this way, the classifier directly models the joint distribution of all possible logical sectors (e.g., using the softmax function) given syndrome. Then we can see that a significant issue of training a classifier is that the number of classes increases with the number of logical sectors 4^k [61] and the computational cost for molding the joint distribution of all 4^k logical sectors is not scalable with the number of qubits k.

To resolve this issue, several works proposed to simplify the joint distribution of 4^k logical sectors using an approximation of the product of marginal distributions [37, 39, 44]

$$p(\beta_1, \dots, \beta_{2k}|\gamma) \approx p(\beta_1|\gamma) \dots p(\beta_{2k}|\gamma)$$
 (6)

In other words, this approach approximates a multi-classification task with 4^k classes as 2k independent binary-classification tasks. We term this approach marginal neural decoding (MND). In statistical physics, the ansatz of using the product of marginals as an approximation to the joint distribution is known as naïve mean-field approximation [62, 63], and it has been explicitly shown that the autoregressive distribution (as in our generative approach) can provide achieves much better performance in solving statistical mechanics problems [63]. In machine learning, it is well known that the product ansatz can not provide representative representations for complex distributions such as distributions of words in natural languages [46], and modern large language models [54] adopt autoregressive distributions (such as in causal Transformers [33]) rather than product distributions.

In the context of the neural network decoder, we explicitly compared the performance of our generative neural decoding GND approach with the neural network decoder using the product of marginals MND. In detail, we performed numerical experiments on [[30, 6, 4]] and [[60, 8, 4]] high-rate codes and compared the logical error rates given by GND and MND. In both GND and MND, we explored the same neural network architecture, using MADE [32] and Multi-Layer Perceptron (MLP) [37] so the parameters of them are quite similar. The results are shown in Fig. 4. From the figure, we can see that GND demonstrates significant advantages over MND, indicating that the principle, the autoregressive modeling, is significantly superior to the simple product of marginals as used in previous MND studies.

We observe that another notable distinction between our generative decoder and existing neural network decoders is that our generative decoder models the joint probability $p(\beta, \gamma)$, as illustrated in Eq. (3). This implies that the generative decoder not only models $p(\beta \mid \gamma)$ but also $p(\gamma)$. Consequently, in addition to generating logical sectors, it

also has the capability to generate syndromes.

C. Comparing MADE with Transformer in decoding

In this section, we evaluated two different autoregressive models in our GND framework, MADE [32] and causal Transformer [36]. The architectures of these models are illustrated in Fig. 1. MADE is a lightweight autoregressive model composed exclusively of masked linear layers and activation functions, enforcing the autoregressive property through binary connectivity masks. Its simplicity and efficiency make it a popular choice for probabilistic modeling in statistical physics, particularly for partition function estimation [52]. In contrast, the causal transformer leverages the decoder module of the Transformer architecture [33], which has achieved state-of-the-art performance in language modeling tasks [35]. Despite its success in other domains, the transformer underperformed MADE in our quantum decoding framework. This may be due to the large number of parameters of the transformer that are difficult to optimize, as well as the self-attention mechanism, which dynamically models dependencies between variables, may struggle with problems where variable relationships are unknown or poorly structured, as noted in previous work [64]. To quantify this disparity, Fig. 5 compares the decoding accuracy of MADE and Transformer on increasing size rotated surface codes. While both models perform similarly for small code distances, the Transformer's performance degrades significantly as the system scales. In addition, its memory footprint exceeds the capacity of a single NVIDIA A100 GPU at larger scales, when the code distance is up to 11. We conducted extensive studies to isolate the cause of the Transformer's shortcomings and tried to remove modules that might affect performance, such as dropout and layer norm, among others. All these attempts failed to improve the performance of the Transformer any further.

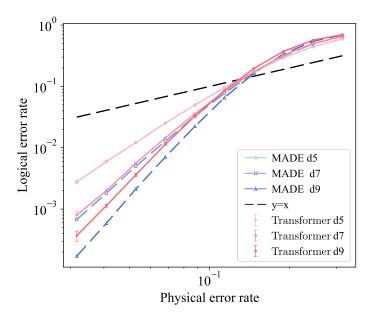


FIG. 5. Comparing the decoding performance of MADE with Transformer on d = [3, 5, 7] Rotated Surface Codes. In this figure, we use the hyperparameters of Transformer as $d_{model} = 256$, $N_{head} = 4$, $d_{ff} = 256$, $N_{layers} = 2$. And the hyperparameters of MADE can be found in Tab. I. Each data point has 10 decoding tasks, and each task has 10^5 shots.

D. Parameters of the neural networks

In the numerical experiments, all neural networks were trained on a single NVIDIA A100 GPU. The hyperparameters for training the neural networks are detailed in Tab. I. The table provides the number of inputs N_{in} , Depth, Width and number of parameters N_{para} of neural networks, across multiple error-correcting codes, including RSC (Rotated Rurface codes), LDPC (low-density parity-check codes), BB (Bivariate Bicycle codes) and DSC^{cir} (Defected Surface codes with circuit level noise). P_t is the physical error rate, T_{decoding} is the decoding time in seconds.

Code	N_{in}	Depth	Width	P_t	$N_{ m para}$	$T_{ m decoding}$
RSC7	50	4	30	0.189	3522050	8.5×10^{-4}
RSC9	82	4	30	0.189	9397282	8.5×10^{-4}
RSC11	122	4	20	0.189	9308722	1.1×10^{-3}
RSC13	170	4	25	0.189	27988545	9.2×10^{-4}
LDPC30	36	4	20	0.1	1980972	4×10^{-3}
BB60	68	5	30	0.12	11666828	6.4×10^{-3}
BB72	84	6	30	0.12	16289364	9.3×10^{-3}
$\mathrm{DSC}^{\mathrm{cir}}_{k2}$	191	4	30	0.015	66789071	1.2×10^{-3}
$\mathrm{DSC}_{k3}^{\mathrm{cir}}$	189	4	30	0.015	65397969	1.5×10^{-3}
$\mathrm{DSC}_{k4}^{\mathrm{cir}}$	187	4	30	0.015	64021507	1.9×10^{-3}

TABLE I. Hyperparameters of GND. Where N_{in} is the number of inputs of the model, Depth is the number of hidden layers, width is the number of channels, P_t is the physical error rate for training, N_{para} is the number of parameters and T_{decoding} is the decoding time in seconds.

ACKNOWLEDGMENTS

A PYTHON implementation and a Jupyter Notebook tutorial of our algorithm are available at [65]. We thank Weilei Zeng and Ying Li for their helpful discussions. We are also grateful to the anonymous reviewers for their valuable suggestions during the early stages of this paper. This work is supported by Project 12325501, 12047503, and 12247104 of the National Natural Science Foundation of China and Project ZDRW-XX-2022-3-02 of the Chinese Academy of Sciences. P. Z. is partially supported by the Innovation Program for Quantum Science and Technology project 2021ZD0301900.

AUTHOR CONTRIBUTIONS

H.C., F.P., and P.Z. had the original idea for this work. H.C., F.P., D.F., and Y. W. performed the study, and all authors contributed to the preparation of the manuscript.

COMPETING INTERESTS

The authors declare no competing interests.

ADDITIONAL INFORMATION

Correspondence and requests for materials should be addressed to P.Z. Reprints and permission information is available online at [URL will be inserted by publisher].

^[1] P. W. Shor, Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer, SIAM Journal on Computing 26, 1484 (1997), arXiv:quant-ph/9508027.

^[2] C. Gidney and M. Ekerå, How to factor 2048 bit rsa integers in 8 hours using 20 million noisy qubits, Quantum 5, 433 (2021), arXiv:1905.09749 [quant-ph].

^[3] A. Aspuru-Guzik, A. D. Dutoi, P. J. Love, and M. Head-Gordon, Simulated quantum computation of molecular energies, Science 309, 1704 (2005).

^[4] W. van Dam, H. Liu, G. H. Low, A. Paetznick, A. Paz, M. Silva, A. Sundaram, K. Svore, and M. Troyer, End-to-end quantum simulation of a chemical system (2024), arXiv:2409.05835 [quant-ph].

^[5] P. Panteleev and G. Kalachev, Quantum LDPC codes with almost linear minimum distance, IEEE Transactions on Information Theory 68, 213 (2022).

^[6] P. Panteleev and G. Kalachev, Asymptotically good quantum and locally testable classical ldpc codes, in *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2022 (Association for Computing Machinery, New York, NY, USA, 2022) p. 375-388

- [7] S. Bravyi, A. W. Cross, J. M. Gambetta, D. Maslov, P. Rall, and T. J. Yoder, High-threshold and low-overhead fault-tolerant quantum memory, Nature 627, 778-782 (2024)
- [8] Google Quantum AI, Exponential suppression of bit or phase errors with cyclic error correction, Nature 595, 383 (2021).
- [9] Google Quantum AI, Suppressing quantum errors by scaling a surface code logical qubit, Nature 614, 676 (2023).
- [10] Google Quantum AI, Quantum error correction below the surface code threshold (2024), arXiv:2408.13687 [quant-ph].
- [11] D. Gottesman, Stabilizer codes and quantum error correction (California Institute of Technology, 1997).
- [12] M. A. Nielsen and I. L. Chuang, Quantum Computation and Quantum Information: 10th Anniversary Edition (Cambridge University Press, 2010).
- [13] Z. Babar, P. Botsinis, D. Alanis, S. X. Ng, and L. Hanzo, Fifteen years of quantum ldpc coding and improved decoding strategies, IEEE Access 3, 2492 (2015).
- [14] P. W. Shor, Scheme for reducing decoherence in quantum computer memory, Phys. Rev. A 52, R2493 (1995).
- [15] A. R. Calderbank and P. W. Shor, Good quantum error-correcting codes exist, Physical Review A 54, 1098-1105 (1996)
- [16] A. M. Steane, Error correcting codes in quantum theory, Phys. Rev. Lett. 77, 793 (1996).
- [17] S. B. Bravyi and A. Y. Kitaev, Quantum codes on a lattice with boundary, arXiv preprint quant-ph/9811052 (1998).
- [18] E. Dennis, A. Kitaev, A. Landahl, and J. Preskill, Topological quantum memory, Tech. Rep. 9 (2002).
- [19] O. Higgott, Pymatching: A python package for decoding quantum codes with minimum-weight perfect matching (2021), arXiv:2105.13082 [quant-ph].
- [20] Y. Wu and L. Zhong, Fusion blossom: Fast mwpm decoders for qec, arXiv preprint arXiv:2305.08307 (2023).
- [21] M. Fossorier and S. Lin, Soft-decision decoding of linear block codes based on ordered statistics, IEEE Transactions on Information Theory 41, 1379 (1995).
- [22] M. Fossorier, Iterative reliability-based decoding of low-density parity check codes, IEEE Journal on Selected Areas in Communications 19, 908 (2001).
- [23] J. Roffe, D. R. White, S. Burton, and E. Campbell, Decoding across the quantum low-density parity-check code landscape, Physical Review Research 2, 10.1103/PhysRevResearch.2.043423 (2020).
- [24] P. Panteleev and G. Kalachev, Degenerate Quantum LDPC Codes With Good Finite Length Performance, Quantum 5, 585 (2021).
- [25] S. Bravyi, M. Suchara, and A. Vargo, Efficient algorithms for maximum likelihood decoding in the surface code, Physical Review A 90, 032326 (2014).
- [26] C. T. Chubb, General tensor network decoding of 2D Pauli codes (2021), arxiv:2101.04125 [quant-ph].
- [27] C. Piveteau, C. T. Chubb, and J. M. Renes, Tensor Network Decoding Beyond 2D (2023), arxiv:2310.10722 [quant-ph].
- [28] C. T. Chubb and S. T. Flammia, Statistical mechanical models for quantum codes with correlated noise, Annales de l'Institut Henri Poincaré D 8, 269-321 (2021)
- [29] H. Bombin, R. S. Andrist, M. Ohzeki, H. G. Katzgraber, and M. A. Martin-Delgado, Strong resilience of topological codes to depolarization, Physical Review X 2, 10.1103/physrevx.2.021004 (2012).
- [30] M. Vasmer and D. E. Browne, Three-dimensional surface codes: Transversal gates and fault-tolerant architectures, Phys. Rev. A 100, 012312 (2019).
- [31] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, Empirical evaluation of gated recurrent neural networks on sequence modeling (2014), arXiv:1412.3555 [cs.NE].
- [32] M. Germain, K. Gregor, I. Murray, and H. Larochelle, Made: Masked autoencoder for distribution estimation (2015), arXiv:1502.03509 [cs.LG].
- [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, Attention is all you need, Advances in neural information processing systems 30 (2017).
- [34] S. Hochreiter and J. Schmidhuber, Long short-term memory, Neural Computation 9, 1735 (1997).
- [35] https://openai.com/chatgpt.
- [36] R. Fakoor, P. Chaudhari, J. Mueller, and A. J. Smola, Trade: Transformers for density estimation (2020), arXiv:2004.02441 [cs.LG].
- [37] S. Varsamopoulos, B. Criger, and K. Bertels, Decoding small surface codes with feedforward neural networks, Quantum Science and Technology 3 (2017).
- [38] S. Krastanov and L. Jiang, Deep Neural Network Probabilistic Decoder for Stabilizer Codes, Scientific Reports 7 (2017).
- [39] P. Baireuther, T. E. O'Brien, B. Tarasinski, and C. W. Beenakker, Machine-learning-assisted correction of correlated qubit errors in a topological code, Quantum 2, 48 (2018).
- [40] R. W. Overwater, M. Babaie, and F. Sebastiano, Neural-Network Decoders for Quantum Error Correction Using Surface Codes: A Space Exploration of the Hardware Cost-Performance Tradeoffs, IEEE Transactions on Quantum Engineering 3, 1 (2022).
- [41] A. Davaasuren, Y. Suzuki, K. Fujii, and M. Koashi, General framework for constructing fast and near-optimal machine-learning-based decoder of the topological stabilizer codes, Physical Review Research 2 (2020).
- [42] S. Gicev, L. C. Hollenberg, and M. Usman, A scalable and fast artificial neural network syndrome decoder for surface codes, arXiv preprint arXiv:2110.05854 (2021).
- [43] K. Meinerz, C.-Y. Park, and S. Trebst, Scalable neural decoder for topological surface codes, Phys. Rev. Lett. **128**, 080505 (2022).
- [44] M. Lange, P. Havström, B. Srivastava, V. Bergentall, K. Hammar, O. Heuts, E. van Nieuwenburg, and M. Granath, Data-driven decoding of quantum error correcting codes using graph neural networks, arXiv preprint arXiv:2307.01241 10.48550/arXiv.2307.01241 (2023).

- [45] J. Bausch, A. W. Senior, F. J. H. Heras, T. Edlich, A. Davies, M. Newman, C. Jones, K. Satzinger, M. Y. Niu, S. Blackwell, et al., Learning high-accuracy error decoding for quantum processors, Nature (2024).
- [46] C. M. Bishop and H. Bishop, Deep learning: Foundations and concepts (Springer Nature, 2023).
- [47] H. Bombin and M. A. Martin-Delgado, Optimal resources for topological two-dimensional stabilizer codes: Comparative study, Physical Review A **76** (2007).
- [48] https://github.com/WeileiZeng/CSS-Code-Database/blob/main/dataI0.ipynb.
- [49] S. Varsamopoulos, K. Bertels, and C. G. Almudever, Comparing neural network based decoders for the surface code, IEEE Transactions on Computers 69, 300 (2020).
- [50] G. Torlai and R. G. Melko, Neural decoder for topological codes, Phys. Rev. Lett. 119, 030501 (2017).
- [51] Y.-H. Liu and D. Poulin, Neural belief-propagation decoders for quantum error-correcting codes, Phys. Rev. Lett. 122, 200501 (2019).
- [52] D. Wu, L. Wang, and P. Zhang, Solving statistical mechanics using variational autoregressive networks, Phys. Rev. Lett. 122, 080602 (2019).
- [53] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, Improving language understanding by generative pre-training (OpenAI, 2018).
- [54] OpenAI, Gpt-4 technical report (2023), arXiv:2303.08774 [cs.CL].
- [55] A. G. Fowler, M. Mariantoni, J. M. Martinis, and A. N. Cleland, Surface codes: Towards practical large-scale quantum computation, Phys. Rev. A 86, 032324 (2012).
- [56] Y. Zhao, Y. Ye, H.-L. Huang, Y. Zhang, D. Wu, H. Guan, Q. Zhu, Z. Wei, T. He, S. Cao, et al., Realization of an error-correcting surface code with superconducting qubits, Physical Review Letters 129, 030501 (2022).
- [57] T. C. Bohdanowicz, E. Crosson, C. Nirkhe, and H. Yuen, Good approximate quantum ldpc codes from spacetime circuit hamiltonians, in *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing* (ACM, Phoenix AZ USA, 2019) pp. 481–490.
- [58] T. C. Bohdanowicz, Quantum Constructions on Hamiltonians, Codes, and Circuits, Ph.D. thesis, California Institute of Technology (2022).
- [59] M. McEwen, D. Bacon, and C. Gidney, Relaxing Hardware Requirements for Surface Code Circuits using Time-dynamics, Quantum 7, 1172 (2023).
- [60] C. Gidney, Stim: A fast stabilizer circuit simulator, Quantum 5, 497 (2021), 2103.02202 [quant-ph].
- [61] A. Davaasuren, Y. Suzuki, K. Fujii, and M. Koashi, General framework for constructing fast and near-optimal machine-learning-based decoder of the topological stabilizer codes, Phys. Rev. Res. 2, 033399 (2020).
- [62] J. S. Yedidia, W. T. Freeman, Y. Weiss, et al., Understanding belief propagation and its generalizations, Exploring artificial intelligence in the new millennium 8, 0018 (2003).
- [63] D. Wu, L. Wang, and P. Zhang, Solving statistical mechanics using variational autoregressive networks, Physical Review Letters 122, 10.1103/physrevlett.122.080602 (2019).
- [64] M. Hahn. Theoretical limitations of self-attention inneural sequence models. Transactions the Association for Computational Linguistics 8, 156 (2020).https://direct.mit.edu/tacl/articlepdf/doi/10.1162/tacl_a_00306/1923102/tacl_a_00306.pdf.
- [65] https://github.com/CHY-i/GND.