

# Transformer-QEC: Quantum Error Correction Code Decoding with Transferable Transformers

Hanrui Wang<sup>1</sup>, Pengyu Liu<sup>2</sup>, Kevin Shao<sup>1</sup>, Dantong Li<sup>3</sup>,  
Jiaqi Gu<sup>4</sup>, David Z. Pan<sup>5</sup>, Yongshan Ding<sup>3</sup>, Song Han<sup>1</sup>

<sup>1</sup>Massachusetts Institute of Technology, <sup>2</sup>Carnegie Mellon University, <sup>3</sup>Yale University, <sup>4</sup>Arizona State University, <sup>5</sup>University of Texas at Austin

## Abstract

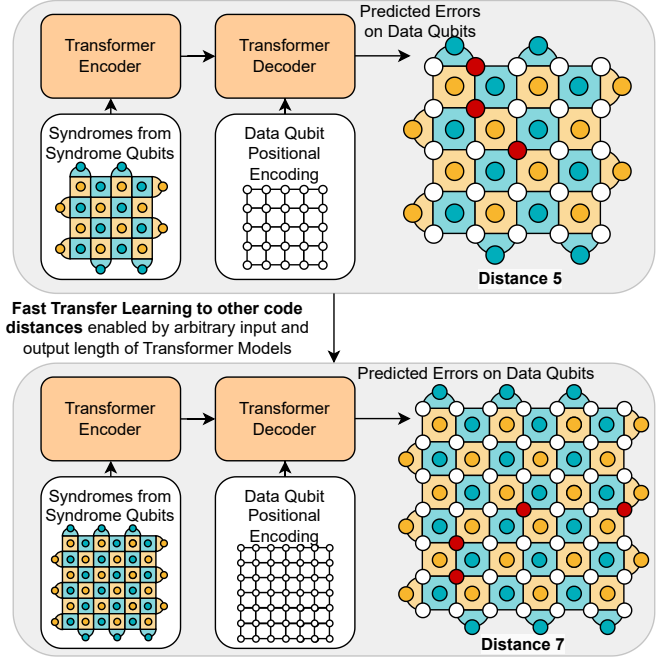
Quantum computing has the potential to solve problems that are intractable for classical systems, yet the high error rates in contemporary quantum devices often exceed tolerable limits for useful algorithm execution. Quantum Error Correction (QEC) mitigates this by employing redundancy, distributing quantum information across multiple data qubits and utilizing syndrome qubits to monitor their states for errors. The syndromes are subsequently interpreted by a decoding algorithm to identify and correct errors in the data qubits. This task is complex due to the multiplicity of error sources affecting both data and syndrome qubits as well as syndrome extraction operations. Additionally, identical syndromes can emanate from different error sources, necessitating a decoding algorithm that evaluates syndromes collectively. Although machine learning (ML) decoders such as multi-layer perceptrons (MLPs) and convolutional neural networks (CNNs) have been proposed, they often focus on *local* syndrome regions and require *retraining* when adjusting for different code distances. To address these issues, we introduce a transformer-based QEC decoder, termed Transformer-QEC, which employs self-attention to achieve a global receptive field across all input syndromes. It incorporates a *mixed loss* training approach, combining both local physical error and global parity label losses. Moreover, the transformer architecture’s inherent adaptability to variable-length inputs allows for efficient *transfer learning*, enabling the decoder to adapt to varying code distances without retraining.

Evaluation on six code distances and ten different error configurations demonstrates that our model consistently outperforms non-ML decoders, such as Union Find (UF) and Minimum Weight Perfect Matching (MWPM), and other ML decoders, thereby achieving best logical error rates. Moreover, the transfer learning can save over 10 $\times$  of training cost.

## 1 Introduction

The field of Quantum Computing (QC) has attracted significant attention in research circles as a novel computational paradigm, poised to tackle challenges previously beyond the reach of conventional computing with remarkable efficiency. QC offers promising prospects across various industries and academic fields, notably impacting areas such as cryptography [33], database searching [17], combinatorial optimization [14, 26], molecular dynamics simulations [31], and advancements in machine learning [5, 24, 25, 28, 39–42, 44, 48], among others.

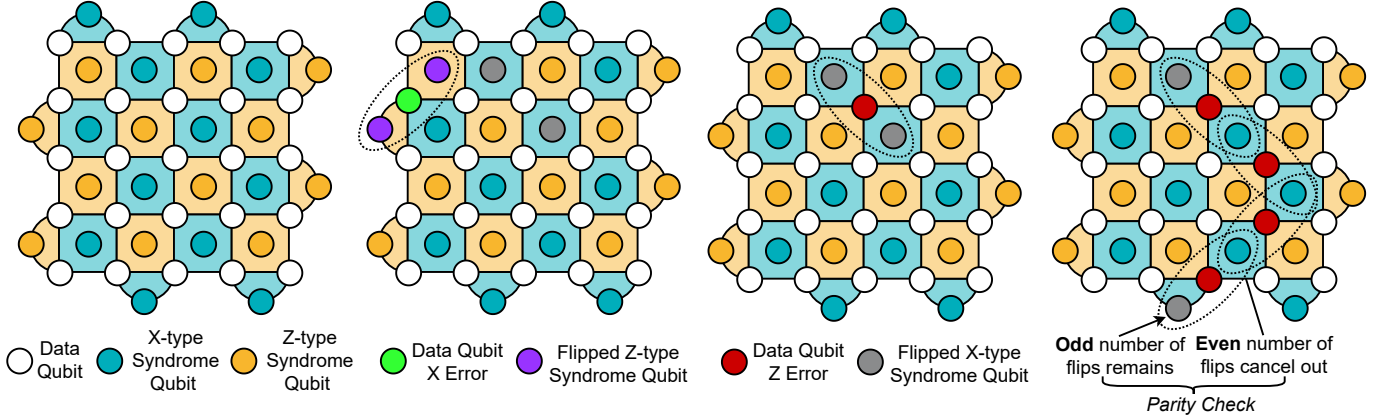
Despite significant progress in quantum computing hardware [21], current devices exhibit high error rates, ranging from 10<sup>-3</sup> to 10<sup>-2</sup>, which are orders of magnitude above the thresholds required for practical applications (below 10<sup>-10</sup>) [23]. Addressing these error rates is crucial for advancing the field. Quantum Error Correction (QEC) serves as a key technique for error mitigation by incorporating redundancy—distributing the information of a single logical qubit across multiple physical qubits. With adequate redundancy in QEC codes, the logical error rate can decrease exponentially, provided the



**Figure 1: Transformer for error correction decoding overview. The Transformer takes syndrome inputs and processes them through both the transformer model encoder and decoder. The output of this process consists of error predictions. One notable advantage is its ability to be seamlessly applied to different code distances due to the transformer model’s flexible input and output size.**

physical error rate  $p$  remains below a specific threshold. Thus, judiciously managed redundancy via QEC can enable quantum systems to achieve the low error rates necessary for practical computation. In QEC, logical qubits encode across multiple data qubits, with error detection aided by parity qubits, as shown in Figure 1. Periodic syndrome extraction localizes errors into detectable syndromes over several cycles. A classical decoder then interprets these syndromes to correct errors, influenced by physical error rates and decoder performance. Decoders often require a broad syndrome field due to potential syndrome overlap from different errors.

The rotated surface code [20] is a promising candidate for realizing fault tolerance. Well-established algorithms for surface code include Minimum Weight Perfect Matching (MWPM) and Union Find (UF). Recently, machine learning (ML), especially neural network (NN) based decoders have gained attention due to a few desirable characteristics. First, they generally run in constant time, which is necessary to prevent a backlog of syndrome outcomes. Second, unlike MWPM, they are capable of learning both correlations between physical errors (such as the correlation between X and Z error



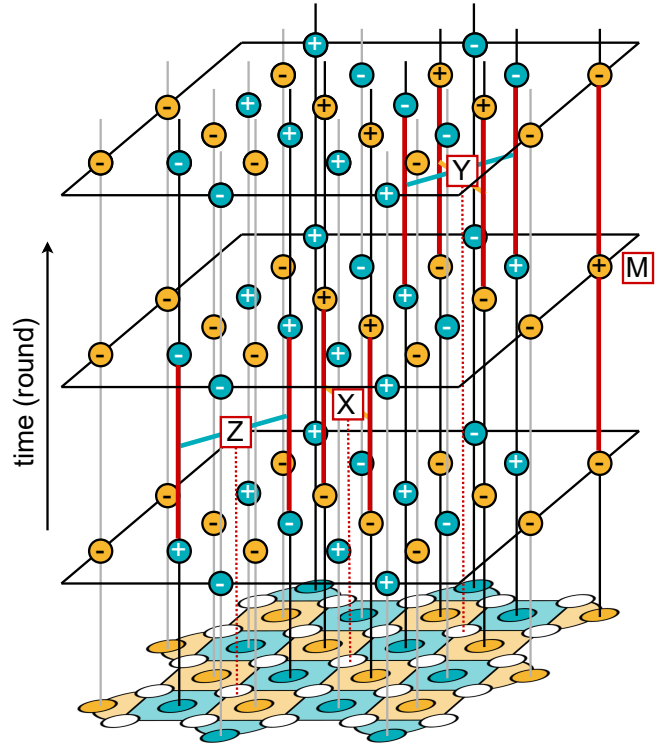
**Figure 2: The surface code contains data qubits and two kinds of syndrome qubits. X-type syndrome qubits in green checks Z errors while Z-type syndrome qubits in yellow check X errors. When error occurs on data qubits, the nearby syndromes may be flipped depending on the parity of data qubits. When multiple error occurs, the syndrome patterns will be more difficult to decode.**

in depolarizing errors) as well as learning hidden and potentially changing underlying physical error distributions.

However, ML based decoders also bring significant challenges. First, several models, such as those grounded in convolutional neural networks (CNN), are limited by a small receptive field. This constraint may hinder their ability to accurately pinpoint long error chains. Second, many models, like the multi-layer perceptron (MLP), has a fixed size for both input and output. As a result, changes in code distance would mandate retraining of an entirely new model, leading to considerable overhead.

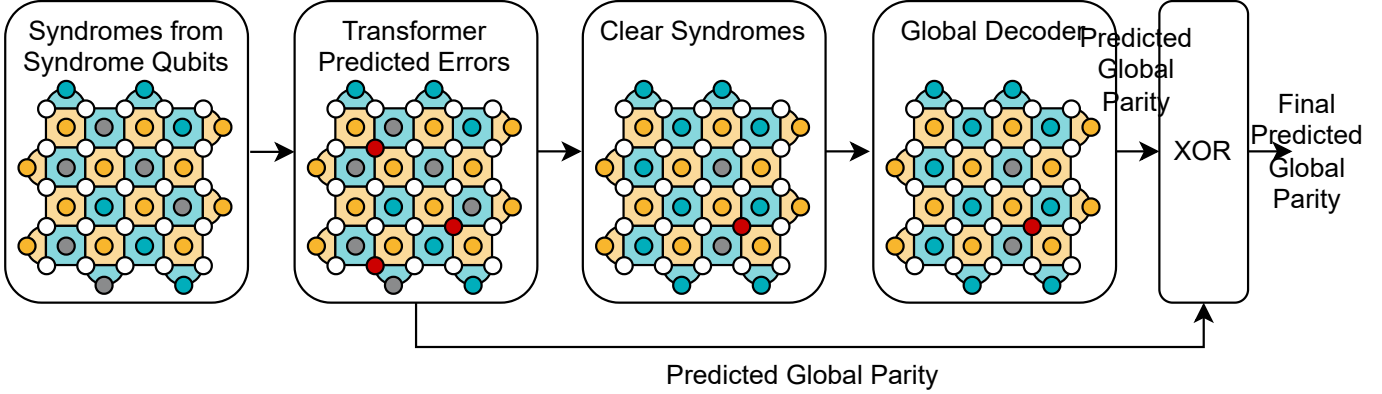
Therefore, to solve these challenges, we propose Transformer-QEC, a transferable transformer model designed for accurate and efficient decoding of surface code, as illustrated in Fig. 1. For the sake of simplicity, the figure only depicts two dimensions, but in reality, the input syndromes include an additional temporal dimension – *round* as in Fig. 2. Our proposed model employs a transformer structure, incorporating both an encoder and a decoder to process the syndromes. Binary features on each syndrome qubit are projected to token embeddings and augmented with a 3D sinusoidal positional encoding, informing the model about the location of each qubit. The embeddings of the 3D inputs are then flattened to 1D input sequence and processed by the transformer encoder layer. Thanks to the global interaction capability brought by attention layer, all input syndromes can be considered holistically which boosts accuracy. The decoder then uses the positional encoding of the data qubits to predict the X or Z errors on each of them. Moreover, we propose a *mixed loss* that combines the loss from the local physical error of each qubit with the loss from global parity prediction.

To optimize computational efficiency in the context of varying code distances for quantum error correction, we advocate for a transfer learning paradigm. Taking advantage of the transformer model’s capacity for arbitrary input lengths, we repurpose pretrained models for different code distances by merely modifying the input sequence. For example, a model trained initially for a code distance  $d = 5$  can be efficiently fine-tuned for alternate distances ( $d = 7$  or  $d = 9$ ), thereby achieving enhanced performance metrics. Our approach yields a tenfold reduction in computational costs as compared to training from scratch, as verified through our experiments.

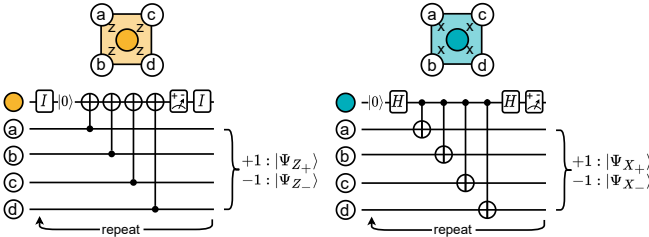


**Figure 3: Multiple rounds of surface code measurement. The progression of time is depicted by moving upwards from the array at the base, with each horizontal plane representing a step in the measurement process. In reality, errors will also occur in the syndrome extraction circuit and syndrome qubits, necessitating the need to repeat multiple rounds for decoding. On the right side, the measurement error on the syndrome qubit will also flip the syndrome.**

We extensively evaluate Transformer-QEC across six code distances, 3, 5, 7, 9, 11 and 13 and compare it with MWPM, UF and MLP baselines under 10 different error rates. Our results demonstrate that



**Figure 4: Overall workflow of Transformer-QEC.** The syndromes are firstly processed by the transformer model to predict the errors. Since the errors may not fully clear all syndromes, we will pass the cleared syndromes to a global decoder to predict a global parity. The final global parity is the XOR of the global parity from transformer predicted physical error and that predicted by global decoder.



**Figure 5: Syndrome extraction circuit. Top: Z-type syndrome qubits. Bottom: X-type syndrome qubits.**

Transformer-QEC consistently surpasses these baselines, achieving the lowest logical error rates. In summary, Transformer-QEC makes four key contributions:

- **A novel transformer-based model** for surface code decoding which uses the syndrome with positional encoding as inputs and predict errors.
- **A mixed loss** approach combined loss from local physical error prediction and global parity prediction improves the model’s trainability and performance.
- **Transfer learning across different code distances.** For the first time, we propose to transfer the knowledge learn on one distance to another, thus reducing costs.
- **Extensive evaluations** on different physical error rates and distances demonstrates that our model consistently outperforms baselines such as MWPM, UF and MLP.

## 2 Background and Related Works

**Quantum error correction.** Quantum Error Correction (QEC) enhances logical qubit fidelity by expanding into multiple physical qubits. The method involves syndrome extraction and error decoding steps. Using auxiliary syndrome qubits, errors are detected and classified into discrete Pauli categories [29]. If the physical error rate is below a certain threshold, QEC reduces the logical error rate, requiring more physical qubits. Quantum techniques often have added complexity compared to classical ones, with specific schemes viable for NISQ devices. This work focuses on the rotated surface code.

**Surface code.** The surface code is a leading QEC scheme that employs a two-dimensional lattice of interlaced data and syndrome qubits to encode a logical qubit. Notable for its high error threshold and need for only nearest-neighbor connectivity, it offers practicality for real-world quantum systems. The code distance  $D$  dictates the lattice size and correlates with error resilience. Adjacent parity qubits detect errors on data qubits via a stabilizer circuit, capable of identifying X, Z, or Y errors, as depicted in Fig. 5. The code can correct error chains up to a length of  $\lfloor \frac{D-1}{2} \rfloor$ . A simplified variant, the ‘rotated’ surface code (Fig. 2), reduces qubit and gate overhead and is often favored in practice. Characterized as a  $[[D^2, 1, D]]$  stabilizer code, it stands as a viable option for near-term fault-tolerant quantum computation. Its design allows transversal single-qubit operations and enables two-qubit CNOT gates via lattice surgery [7], thereby enhancing its feasibility for physical implementations. Decoders analyze syndromes from ancilla measurements to correct data qubit errors. They independently address X-type and Z-type errors, implicitly fixing Y-type errors. For large-scale FTQCs, decoders must be accurate, fast, and scalable [8]. Accuracy indicates reliable error detection, latency mandates cycle-limited operation, and scalability ensures efficient resource use. Increased accuracy may prolong operation time.

Decoders interpret syndromes, outcomes of ancilla measurements as shown in Fig.3, to identify and correct errors in data qubits. X-type and Z-type errors are treated independently, which inherently rectifies Y-type errors. Effective decoders for large-scale fault-tolerant quantum computers (FTQCs) must satisfy three primary criteria: accuracy, latency, and scalability[8]. Accuracy denotes the decoder’s reliability in error identification. Latency requires the decoder to function within a single syndrome extraction cycle. Scalability entails efficient resource utilization, crucial for hardware-constrained environments. Notably, increased accuracy often comes at the cost of extended operational time.

**ML based decoder.** The landscape of quantum computing has been significantly enriched by ML-based decoders, notably through Neural Network (NN) and Reinforcement Learning (RL) paradigms. In the realm of NNs, Boltzmann machines initiated ML-based decoding in toric codes [35], later broadened by Multi-Layer Perceptrons [4] and Long Short-Term Memory networks for surface codes [2]. On

the RL front, advancements include transforming decoding into an RL environment [32] and optimizing RL decoders for specific error types [1]. Efforts for scalable decoding have ranged from low-depth Convolutional Neural Networks [3] to multilevel architectures [37]. Guided by these contributions, our research focuses on employing Transformer-based NNs for local decoding, aiming to develop novel architectures and enhance scalability through a two-level local-global design. There exist numerous hardware accelerators [22, 27, 38, 45, 46] that can be used to improve the efficiency and speed of NN-based decoders.

**Non-ML based decoder.** Various decoding algorithms for QEC offer distinct advantages. Minimum Weight Perfect Matching (MWPM) uses Edmonds’ method for optimal error syndrome pairing in topological codes [15, 47]. The Union Find (UF) decoder features linear-time complexity for toric and surface codes [10]. Lookup Table (LUT) decoders rely on predetermined error patterns [34], while Tensor Network (TN) decoders utilize tensor structures for high error thresholds in topological codes [6]. MWPM decoders [11, 13, 15, 19] are considered to have a good balance between the decoding accuracy and speed. It has almost linear time complexity [16, 18, 47], is practical for hardware implementation for real-time decoding [36] and is more accurate than LUT and Union Find decoders. [43] proposes to tackle the noise drift by updating the MWPM decoding graph weights.

### 3 Methodology

In this section, we delineate the error correction pipeline of our Transformer-QEC framework, followed by detailed discussions on the incorporated transformer model and transfer learning.

**Overall workflow.** The iterative syndrome extraction procedure yields syndrome outcomes at discrete rounds, serving as the basis for the decoder’s error predictions. Purely ML-based decoders, while effective, may yield predictions misaligned with the original syndromes. Therefore, an auxiliary, non-ML decoder is employed for syndrome clearance, as depicted in Figure 4. Post-ML decoding, the predicted errors are utilized to reconcile syndromes and ascertain their global parity. Any residual errors are addressed by a global decoder like MWPM, ensuring complete syndrome clearance and yielding an additional global parity. The decoding process culminates in the XOR of the two global parities obtained.

**Transformer model.** Given that the computational complexity of global decoders such as MWPM is often proportional to the number of non-zero syndromes, the efficiency of a preceding ML-based decoder in clearing syndromes can be highly advantageous. To this end, we introduce a novel Transformer-based decoder, schematically represented in Figure 6. To encode input syndromes, we employ a cubic grid framework. For a surface code characterized by a distance  $D$ , a  $D + 1$  square grid is utilized to ensure that each syndrome qubit resides at a grid intersection. Typically, the number of iterative rounds for syndrome extraction is set equal to the code distance. In our architecture, we incorporate an additional layer dedicated to the final syndrome measurement, thereby extending the round dimension to  $D + 1$ . Consequently, the feature space manifests as a  $D + 1$  cubic grid. Each cell within this grid encapsulates a six-dimensional feature vector. The initial two dimensions specifically indicate the positions of the X-check and Z-check syndrome qubits, as further illustrated in Fig. 6.

Location encoding for the X and Z check syndrome qubits:

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

Subsequent to the positional dimensions, the next two channels in the six-dimensional feature vector encode the syndromes, varying across iterative rounds. For ultimate error correction, data qubits are measured in a designated basis. To enable network generalization across different syndrome measurement rounds, we define temporal lattice boundaries. Specifically, the fifth channel is initialized to 1 for the inaugural round and 0 for later rounds. Analogously, the sixth channel is set to 1 solely for the terminal round and 0 otherwise.

The feature vectors undergo dimensionality elevation via a learnable embedding layer. Subsequently, 3D sinusoidal positional encodings are integrated to convey qubit locations. This enhanced representation is then reshaped from 3D to 1D before forwarding to the Transformer encoder. The encoder comprises multiple layers, each featuring one multi-head self-attention (MHSA) module and one feed-forward network (FFN). The MHSA facilitates long-range contextual awareness by allowing each syndrome to attend to any other within the 3D grid. The FFN consists of two fully-connected layers, serving to further elevate the feature dimensionality, apply an activation function, and project back.

To predict physical errors, we employ Transformer decoder layers, using the positional encodings of data qubit locations as inputs. Each decoder layer comprises a self-attention module and a cross-attention module interfacing with the encoder. Queries for cross-attention originate from the decoder inputs, while keys and values are sourced from the encoder, thereby enabling the decoder to incorporate preceding syndrome information. Subsequently, a FFN layer and a prediction layer output the logits. Owing to the higher cost of false positives relative to false negatives, a confidence threshold is used for error prediction. Specifically, post-Sigmoid confidence  $> 0.95$  triggers a positive prediction.

**Mixed loss:** We introduce a composite loss function during training, integrating contributions from two distinct aspects. The first component arises from the prediction of local physical errors, while the second pertains to global parity prediction. The latter is computed via global average pooling of encoder output embeddings, subsequently processed through a prediction layer, as depicted in Figure 6 (top right). This global parity serves as auxiliary information, enhancing the model’s generalization across diverse syndrome patterns.

**Transfer learning:** Each code type encompasses a family of codes with varying distances, which necessitates different logical error rates based on the quantum algorithm deployed. Traditional approaches typically require retraining for each new code distance, incurring substantial computational and temporal overhead. To mitigate this, Transformer-QEC employs a transfer learning strategy, capitalizing on the inherent similarities between codes of different distances. For instance, handling syndromes in a distance-5 code bears resemblance to managing a sub-block in a distance-7 code. Utilizing a pre-trained model, we fine-tune it on the new distance’s dataset. This flexibility is enabled by the Transformer architecture’s inherent ability to accommodate sequences of arbitrary lengths. The sole aspect warranting careful adjustment is the positional encoding for the new distance.

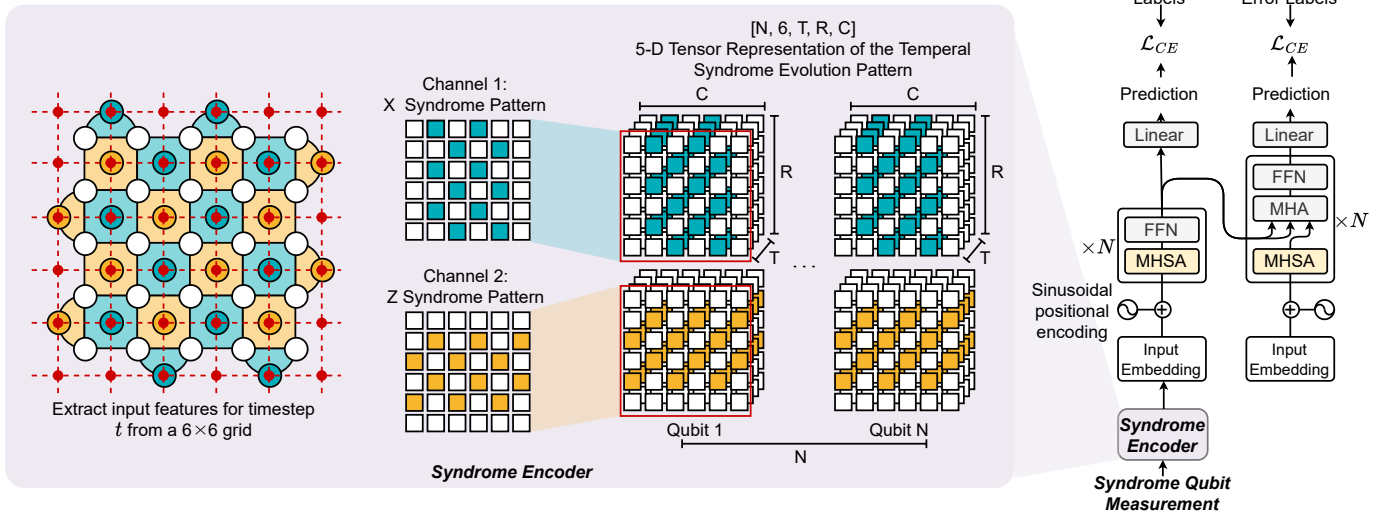


Figure 6: Transformer model architecture. The input of the syndromes will be encoded by a  $(D + 1)$  cubic grid. The input will go through the transformer encoder with self attention and FFN layers. Then the transformer decoder will produce the physical error predictions by processing the positional encoding of data qubits with size  $D$  cubic.

Table 1: Comparison of logical error rates under different code distance and physical error rates.

Distance	Phys. Err. Rate	Logical Error Rate ↓			
		UF	MWPM	MLP	Transformer-QEC
3	0.0500	0.16745	0.14063	0.14794	<b>0.13005</b>
	0.0100	0.01039	0.00800	0.00903	<b>0.00784</b>
5	0.0500	0.24120	0.17279	0.20888	<b>0.17232</b>
	0.0100	0.00406	0.00268	0.00443	<b>0.00254</b>
7	0.0500	0.29813	0.20178	0.28454	<b>0.20590</b>
	0.0100	0.00113	0.00064	0.00197	<b>0.00059</b>
9	0.0500	0.35250	0.23161	0.32770	<b>0.23144</b>
	0.0100	0.00028	0.00002	0.00017	<b>0.00001</b>

## 4 Evaluation

### 4.1 Evaluation Methodology

**Benchmarks.** We have selected the rotated surface code with distances of 3, 5, 7, 9. The round is set to be the same as the distance. The phenomenological error model [12] we use encompasses errors on syndrome measurement and data qubits. Each syndrome qubit experiences a measurement error with a probability  $p$ . The errors on data qubits are depolarizing errors, which occur with a probability  $p$ , causing Pauli X, Y, or Z errors with equal probability. As assumed in previous work [9], the error probabilities of these two types are considered to be equal. We choose values of  $p$  from the set 0.05, 0.04, 0.03, 0.025, 0.02, 0.015, 0.01, 0.0075, 0.005, 0.0025. The Google Stim package is used to construct the circuit and perform stabilizer simulations.

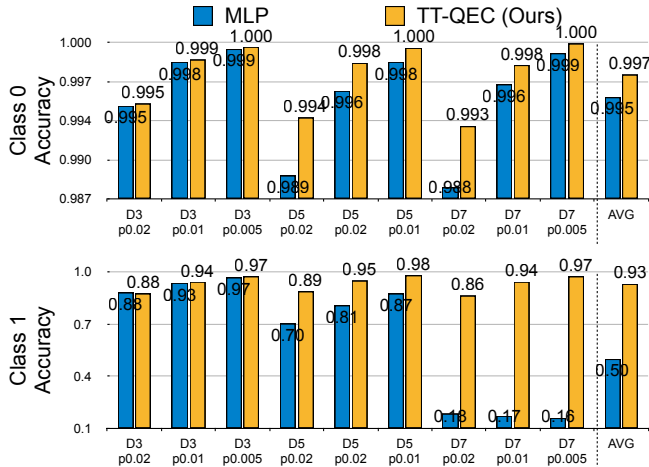
**Baselines.** Our three baselines include the Union Find decoder, the MWPM decoder as implemented in [47], and a MLP architecture. Following [30], our MLP architecture has two hidden layers, with the dimensions of these layers set empirically. As the MLP requires fixed-size inputs and generates fixed-size outputs, it does not facilitate transfer learning like the Transformer does.

**Training settings.** Our main model is a Transformer with 6 layers, an embedding dimension of 256, 8 heads, and a FFN hidden dimension of 512. This model contains 7.9 million parameters. We also have a smaller model with 6 layers, an embedding dimension of 64, 2 heads, and an FFN hidden dimension of 128, which includes 0.5 million parameters. For training, we collect a dataset of 1,000,000 samples with a 1% error rate. We use a learning rate of 0.001 with linear warmup and cosine decay, a weight decay with lambda 0.0001, and we train for 100 epochs. We utilize the Adam optimizer with a weighted binary cross-entropy loss for local physical errors, and a normal binary cross-entropy loss for global parity errors. For the MLP model, we use a physical error rate of 1% for  $d = 3, 5$  and 2.5% for  $d = 7, 9$ . Like the initial Transformer model, we train for 100 epochs with the Adam optimizer. Training is conducted on a single NVIDIA A6000 GPU.

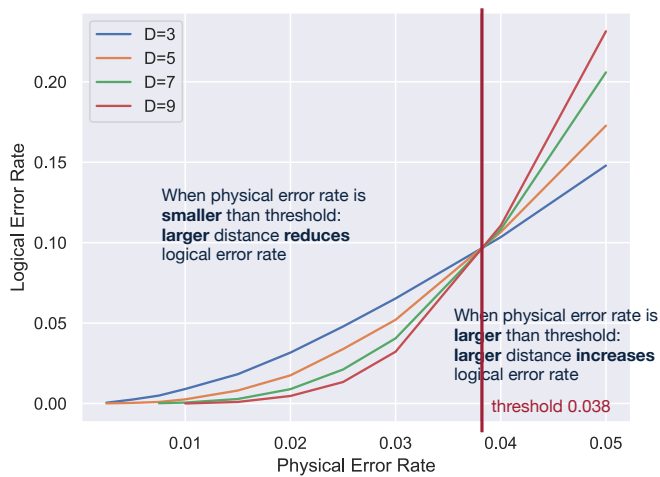
**Transfer learning settings.** The distance 5 model, trained from scratch, is used as the source model for transfer learning. For all other distances, we use a constant learning rate of 0.0005 and train for 10 epochs. All other settings remain identical to the settings of training from scratch.

### 4.2 Experiment Results

**Main results.** Table 1 presents our primary results for varying code distances and physical error rates. Notice that all the models are transferred from the distance 5 model. In general, Transformer-QEC achieves a lower logical error rate for all benchmarks. The improvements over Union Find and MLP decoders are considerably more significant than the MWPM. This is likely because the global decoder in Transformer-QEC’s framework also employs an MWPM. The MLP model can surpass the Union Find but is generally not as good as MWPM and Transformer-QEC, even though the MLP models are trained individually for each code distance. This result highlights the effectiveness of our proposed transfer learning techniques. Furthermore, in Figure 7, we show the physical error prediction accuracy for baseline MLP and our Transformer-QEC. The class 0 accuracy means



**Figure 7: Accuracy comparison between the Transformer-QEC and an MLP baseline. Class 0 accuracy is the accuracy of correctly identify a no error data qubit as no error (True Negative). Class 1 accuracy is the accuracy of correctly identify an error when the data qubit has error (True Positive).**



**Figure 8: Threshold of transformer based decoder. The threshold indicated the largest acceptable physical error rate for which using QEC can reduce error rate. Transformer obtains about 0.038 threshold.**

the ratio of predicted 0 when the ground truth is 0 (true negative). The class 1 accuracy means the ratio of predicted 1 when the ground truth is 1 (true positive). We can see that the accuracy for class 0 is in general much higher than class 1 because of the imbalance of training dataset. Moreover, our Transformer-QEC can achieve 43% higher accuracy on the class 1 which means the Transformer-QEC model can identify errors with much higher reliability. That is beneficial when we desire the low level decoder to clear as many as syndromes as possible and speedup the end-to-end process.

**Evaluation of the threshold.** Figure 8 shows the threshold evaluation of Transformer-QEC, with the X-axis as physical and Y-axis as logical error rates. The curves of different distances intersect at a point where the physical error rate is 0.0038 and the logical error rate is around 0.09. When  $p$  is smaller than the threshold, larger

**Table 2: Comparison of logical error rate with global loss.**

Error Rate	0.0500	0.0400	0.0300	0.0250	0.0200
Local loss	0.17276	<b>0.10659</b>	0.05207	<b>0.03384</b>	0.01751
<b>+ Global loss</b>	<b>0.17232</b>	<b>0.10659</b>	<b>0.05196</b>	<b>0.03384</b>	<b>0.01744</b>

Error Rate	0.0150	0.0100	0.0075	0.0050	0.0025
Local loss	0.00808	0.00259	<b>0.00097</b>	0.00039	0.00007
<b>+ Global loss</b>	<b>0.00802</b>	<b>0.00254</b>	0.00103	<b>0.00035</b>	<b>0.00005</b>

**Table 3: Comparison of logical error rates under different model size.**

Error Rate	0.0200	0.0150	0.0100	0.0075	0.0050
503K Params	0.01812	0.00860	0.00290	0.00127	0.00045
7,911K Params	<b>0.01744</b>	<b>0.00802</b>	<b>0.00254</b>	<b>0.00103</b>	<b>0.00035</b>

code distances reduce the logical error rate. However, when  $p$  is larger than the threshold, larger distances do not help. Instead, we observe larger logical error rates. This trend can be attributed to the increased error introduced by larger system sizes, which eclipses the benefits of greater redundancy with more qubits.

**Effectiveness of mixed loss.** To evaluate the mixed loss function, we perform an ablation study on the distance 5 code, as shown in Table 2. Each column shows the comparison of the logical error rate under a specific physical error rate. The performance with both local and global loss can achieve better or equivalent performance for nine out of ten cases. This demonstrates that the global parity loss provides valuable guidance during the model’s training process.

**Ablation on model size.** We evaluate two models with different sizes but the same training setting in Table 3 under code distance 5 and varying physical error rates. It is evident that the larger model, with approximately 8 million parameters, outperforms the smaller model with 500 thousand parameters. The larger model is not overfitted to the training set and performs poorly on testing. This outcome is mainly due to the large size of the training set.

## 5 Conclusion

In conclusion, our study presents a robust QEC decoder for rotated surface codes using machine learning and transformer models. It outperforms existing benchmarks across code distances and facilitates quick transfer learning. Key factors for this success include global decoding and large Transformer models. This work advances ML-based Transformer decoders, enhancing accuracy and speed in error correction of the surface codes and beyond.

## References

- [1] P. Andreassen et al. 2019. Quantum error correction for the toric code using deep reinforcement learning. *Quantum* 3 (Sept. 2019), 183.
- [2] P. Baireuther et al. 2018. Machine-learning-assisted correction of correlated qubit errors in a topological code. *Quantum* 2 (Jan. 2018), 48.
- [3] Nikolas P. Breuckmann and Xiaotong Ni. 2018. Scalable Neural Network Decoders for Higher Dimensional Quantum Codes. *Quantum* 2 (May 2018), 68.
- [4] Christopher Chamberland and Pooya Ronagh. 2018. Deep neural decoders for near term fault-tolerant experiments. *Quantum Science and Technology* 3, 4 (jul 2018).
- [5] Jinglei Cheng, Hanrui Wang, Zhiding Liang, Yiyu Shi, Song Han, and Xuehai Qian. 2022. TopGen: Topology-Aware Bottom-Up Generator for Variational Quantum Circuits. *arXiv preprint arXiv:2210.08190* (2022).
- [6] Christopher T. Chubb. 2021. General tensor network decoding of 2D Pauli codes. *arXiv:2101.04125* [quant-ph]

- [7] H. Clare et al. 2012. Surface code quantum computing by lattice surgery. *New Journal of Physics* 14, 12 (dec 2012), 123011.
- [8] P. Das et al. 2022. AFS: Accurate, Fast, and Scalable Error-Decoding for Fault-Tolerant Quantum Computers. In *HPCA*. IEEE, 259–273.
- [9] Poulami Das, Aditya Locharla, and Cody Jones. 2021. LILLIPUT: A Lightweight Low-Latency Lookup-Table Based Decoder for Near-term Quantum Error Correction. arXiv:2108.06569 [quant-ph]
- [10] N. Delfosse and G. Zémor. 2017. Linear-Time Maximum Likelihood Decoding of Surface Codes Over the Quantum Erasure Channel. *Quantum Information & Computation* (2017).
- [11] Eric Dennis, Alexei Kitaev, Andrew Landahl, and John Preskill. 2002. Topological quantum memory. *J. Math. Phys.* 43, 9 (2002), 4452–4505.
- [12] Eric Dennis, Alexei Kitaev, Andrew Landahl, and John Preskill. 2002. Topological quantum memory. *J. Math. Phys.* 43, 9 (2002), 4452–4505.
- [13] Jack Edmonds. 1965. Maximum matching and a polyhedron with 0, 1-vertices. *Journal of research of the National Bureau of Standards B* 69, 125–130 (1965), 55–56.
- [14] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. 2014. A quantum approximate optimization algorithm. *arXiv preprint arXiv:1411.4028* (2014).
- [15] A. Fowler et al. 2012. Surface codes: Towards practical large-scale quantum computation. *Physical Review A* 86, 3 (2012), 032324.
- [16] Austin G Fowler, Adam C Whiteside, and Lloyd CL Hollenberg. 2012. Towards practical classical processing for the surface code: timing analysis. *Physical Review A* 86, 4 (2012), 042313.
- [17] Lov K Grover. 1996. A fast quantum mechanical algorithm for database search. In *STOC*. 212–219.
- [18] Oscar Higgott and Craig Gidney. 2023. Sparse Blossom: correcting a million errors per core second with minimum-weight matching. *arXiv preprint arXiv:2303.15933* (2023).
- [19] Adam Holmes, Mohammad Reza Jokar, Ghasem Pasandi, Yongshan Ding, Massoud Pedram, and Frederic T Chong. 2020. NISQ+: Boosting quantum computing power by approximating quantum error correction. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 556–569.
- [20] C. Horsman et al. 2012. Surface code quantum computing by lattice surgery. *New Journal of Physics* 14, 12 (2012), 123011.
- [21] IBM. [n. d.]. IBM Unveils Breakthrough 127-Qubit Quantum Processor.
- [22] Zexi Ji\*, Hanrui Wang\*, Miaorong Wang, Win-San Khwa, Meng-Fan Chang, Song Han, and Anantha P. Chandrakasan. 2023. SpAtten-Chip: A Fully-Integrated Energy-Scalable Transformer Accelerator Supporting Adaptive Model Configuration and Token Pruning for Language Understanding on Edge Devices. In *2023 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*. 1–6. <https://doi.org/10.1109/ISLPED58423.2023.10244459>
- [23] J. Lee et al. 2021. Even more efficient quantum computations of chemistry through tensor hypercontraction. *PRX Quantum* 2, 3 (2021), 030305.
- [24] Z. Liang et al. 2021. Can noise on qubits be learned in quantum neural network? a case study on quantumflow. In *ICCAD*. IEEE, 1–7.
- [25] Zhiding Liang, Jinglei Cheng, Hang Ren, Hanrui Wang, Fei Hua, Yongshan Ding, Fred Chong, Song Han, Yiyu Shi, and Xuehai Qian. 2022. Pan: Pulse ansatz on nisq machines. *arXiv preprint arXiv:2208.01215* (2022).
- [26] Zhiding Liang, Zhixin Song, Jinglei Cheng, Zichang He, Ji Liu, Hanrui Wang, Ruiyang Qin, Yiru Wang, Song Han, Xuehai Qian, and Yiyu Shi. 2023. Hybrid Gate-Pulse Model for Variational Quantum Algorithms. In *2023 60th ACM/IEEE Design Automation Conference (DAC)*. 1–6. <https://doi.org/10.1109/DAC56929.2023.10247923>
- [27] Yujun Lin, Zhekai Zhang, Haotian Tang, Hanrui Wang, and Song Han. 2021. PointAcc: Efficient Point Cloud Accelerator. In *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture (Virtual Event, Greece) (MICRO '21)*. Association for Computing Machinery, New York, NY, USA, 449–461. <https://doi.org/10.1145/3466752.3480084>
- [28] S. Lloyd et al. 2013. Quantum algorithms for supervised and unsupervised machine learning. *arXiv preprint arXiv:1307.0411* (2013).
- [29] Michael A Nielsen and Isaac Chuang. 2002. Quantum computation and quantum information.
- [30] Ramon W. J. Overwater, Masoud Babaie, and Fabio Sebastiano. 2022. Neural-Network Decoders for Quantum Error Correction Using Surface Codes: A Space Exploration of the Hardware Cost-Performance Tradeoffs. *IEEE TQC* (2022).
- [31] A. Peruzzo et al. 2014. A variational eigenvalue solver on a photonic quantum processor. *Nature communications* 5, 1 (2014), 1–7.
- [32] S. Ryan et al. 2020. Reinforcement learning decoders for fault-tolerant quantum computation. *Machine Learning: Science and Technology* 2, 2 (dec 2020), 025005.
- [33] Peter W Shor. 1999. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. *SIAM review* 41, 2 (1999), 303–332.
- [34] Yu Tomita and Krysta M. Svore. 2014. Low-distance surface codes under realistic quantum noise. *Phys. Rev. A* 90 (Dec 2014), 062320. Issue 6.
- [35] Giacomo Torlai and Roger G. Melko. 2017. Neural Decoder for Topological Codes. *Phys. Rev. Lett.* 119 (Jul 2017), 030501. Issue 3. <https://doi.org/10.1103/PhysRevLett.119.030501>
- [36] Suhas Vittal, Poulami Das, and Moinuddin Qureshi. 2023. Astrea: Accurate Quantum Error-Decoding via Practical Minimum-Weight Perfect-Matching. In *Proceedings of the 50th Annual International Symposium on Computer Architecture (Orlando, FL, USA) (ISCA '23)*. Association for Computing Machinery, New York, NY, USA, Article 2, 16 pages. <https://doi.org/10.1145/3579371.3589037>
- [37] Thomas Wagner, Hermann Kampermann, and Dagmar Bruß. 2020. Symmetries for a high-level neural decoder on the toric code. *Phys. Rev. A* 102 (Oct 2020), 042411. Issue 4. <https://doi.org/10.1103/PhysRevA.102.042411>
- [38] Hanrui Wang et al. 2020. *Efficient algorithms and hardware for natural language processing*. Ph.D. Dissertation. Massachusetts Institute of Technology.
- [39] Hanrui Wang, Yongshan Ding, Jiaqi Gu, Yujun Lin, David Z Pan, Frederic T Chong, and Song Han. 2022. QuantumNAS: Noise-adaptive search for robust quantum circuits. In *2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 692–708.
- [40] Hanrui Wang, Jiaqi Gu, Yongshan Ding, Zirui Li, Frederic T Chong, David Z Pan, and Song Han. 2022. QuantumNAT: quantum noise-aware training with noise injection, quantization and normalization. In *Proceedings of the 59th ACM/IEEE Design Automation Conference*. 1–6.
- [41] Hanrui Wang, Zirui Li, Jiaqi Gu, Yongshan Ding, David Z Pan, and Song Han. 2022. QOC: quantum on-chip training with parameter shift and gradient pruning. In *Proceedings of the 59th ACM/IEEE Design Automation Conference*. 655–660.
- [42] Hanrui Wang, Pengyu Liu, Jinglei Cheng, Zhiding Liang, Jiaqi Gu, Zirui Li, Yongshan Ding, Weiwen Jiang, Yiyu Shi, Xuehai Qian, et al. 2022. Quest: Graph transformer for quantum circuit reliability estimation. In *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*.
- [43] Hanrui Wang, Pengyu Liu, Yilian Liu, Jiaqi Gu, Jonathan Baker, Frederic T. Chong, and Song Han. 2023. DGR: Tackling Drifted and Correlated Noise in Quantum Error Correction via Decoding Graph Re-weighting. In *arXiv preprint*.
- [44] Hanrui Wang, Yilian Liu, Pengyu Liu, Jiaqi Gu, Zirui Li, Zhiding Liang, Jinglei Cheng, Yongshan Ding, Xuehai Qian, Yiyu Shi, David Z. Pan, Frederic T. Chong, and Song Han. 2023. RobustState: Boosting Fidelity of Quantum State Preparation via Noise-Aware Variational Training. *arXiv preprint* (2023).
- [45] Hanrui Wang, Zhekai Zhang, and Song Han. 2021. Spatten: Efficient sparse attention architecture with cascade token and head pruning. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 97–110.
- [46] Hanrui Wang\*, Zhekai Zhang\*, Song Han, and William J Dally. 2020. Sparch: Efficient architecture for sparse matrix multiplication. In *2020 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 261–274.
- [47] Yue Wu and Lin Zhong. 2023. Fusion Blossom: Fast MWPM Decoders for QEC. *arXiv preprint arXiv:2305.08307* (2023).
- [48] Han Zheng, Gokul Subramanian Ravi, Hanrui Wang, Kanav Setia, Frederic T Chong, and Junyu Liu. 2023. SnCQA: A hardware-efficient equivariant quantum convolutional circuit architecture. In *2023 IEEE International Conference on Quantum Computing and Engineering (QCE)*.