

Incremental Data-Uploading for Full-Quantum Classification

Maniraman Periyasamy, Nico Meyer, Christian Ufrecht, Daniel D. Scherer, Axel Plinge, and Christopher Mutschler
Fraunhofer IIS, Fraunhofer Institute for Integrated Circuits IIS, Nuremberg, Germany

Abstract—The data representation in a machine-learning model strongly influences its performance. This becomes even more important for quantum machine learning models implemented on noisy intermediate scale quantum (NISQ) devices. Encoding high dimensional data into a quantum circuit for a NISQ device without any loss of information is not trivial and brings a lot of challenges. While simple encoding schemes (like single qubit rotational gates to encode high dimensional data) often lead to information loss within the circuit, complex encoding schemes with entanglement and data re-uploading lead to an increase in the encoding gate count. This is not well-suited for NISQ devices. This work proposes ‘incremental data-uploading’, a novel encoding pattern for high dimensional data that tackles these challenges. We spread the encoding gates for the feature vector of a given data point throughout the quantum circuit with parameterized gates in between them. This encoding pattern results in a better representation of data in the quantum circuit with a minimal pre-processing requirement. We show the efficiency of our encoding pattern on a classification task using the MNIST and Fashion-MNIST datasets, and compare different encoding methods via classification accuracy and the effective dimension of the model.

Index Terms—image classification, variational quantum computing, data uploading.

I. INTRODUCTION

THE field of quantum computing has witnessed a surge of interest from different fields in the scientific community recently. With the realization of quantum devices with increasing qubits and fidelity by several manufacturers like IBM [1] and Google [2], research in quantum computing started shifting from theoretical research towards applied research [3]. Quantum devices employed for information processing and computational purposes are currently being explored to overcome many of the limitations posed by classical hardware in different industry segments such as finance, cybersecurity, and chemical industry [3]. Research groups from diverse fields of science and technology are studying numerous heuristic and non-heuristic quantum computing methods to solve a given problem and achieve the so-called quantum supremacy [4].

However, the current limitation in the number of qubits and low gate fidelity makes non-heuristic approaches impractical.

The research is supported by the Bavarian Ministry of Economic Affairs, Regional Development and Energy with funds from the Hightech Agenda Bayern.

email address for correspondence:

maniraman.periyasamy@iis.fraunhofer.de, axel.plinge@iis.fraunhofer.de

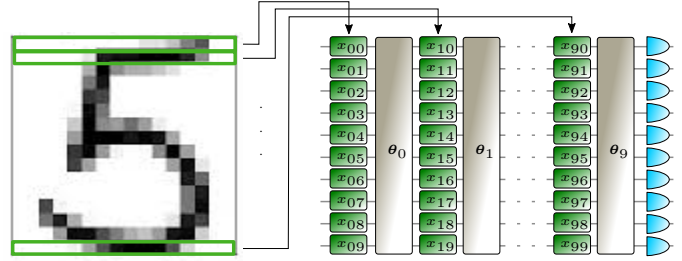


Figure 1: Proposed method: The image is downscaled to 10×10 , then fed row by row to the 10 bit quantum circuit, each encoding layer (green) is followed by a variational layer (gray), then the 10 qubits are measured (blue).

Hence heuristic approaches, especially in the domain of machine learning (ML), are deemed to be one of the prime candidates for practical quantum computing in the NISQ era. Though classical ML is well matured and has decades of domain-specific enhancements, its quantum counterpart is still a lively research field with many loose ends and uncertainties [5]. Adding in the factors from NISQ devices like low gate fidelity and qubit connectivity limitations to this mixture of uncertainties, learning-based approaches like quantum machine learning (QML) and quantum reinforcement learning (QRL) become more complex.

Of all the open questions yet to be answered in the field of gate based quantum computing and quantum machine learning, the question of selecting the optimal variational quantum circuit (VQC) for a given problem is of utmost importance and significance [6]. The gates in a quantum circuit (QC) used for QML are grouped into three categories, namely, encoding gates, decoding gates and the variational gates. These encoding gates, decoding gates and variational gates are further grouped into encoding layer, decoding layer and variational layer in many QML works, though these layers are just a visual representation and do not reflect the theory behind layers in classical ML. The encoding gates are selected based on the encoding method chosen and the number of input features. The optimal selection of the encoding method is pivotal to successful learning of a QML model as this represents the classical information to be fed into the circuit [7]–[10]. While the effect of different encoding methods in data representation and expressivity of a VQC have been studied before [10]–[14],

the encoding pattern we propose (see Section II) has not been discussed in the literature previously to the best of our knowledge.

The encoding gate set positioning becomes a salient factor for VQC-based QML models that learn from high dimensional data due to the limitations in the size of the quantum device in terms of the number of qubits and the depth of the QC that can be executed both in a simulation and on a real device. This limitation brings in a trade-off between the number of encoding gates and the number of parameterized gates for a fixed circuit depth and gate count. The larger the dimension of the input, the larger the number of encoding gates required and the larger the number of encoding gates used, the fewer the number of parameterized gates that can be used. The reduction in the parameterized gate count reduces the expressivity of the model. This reduction of the dimension of the input results in information loss. Naively encoding the data at the start of the circuit or encoding patterns like data re-uploading is not very promising for high dimensional data as it results in an increase in the QC depth and gate count or the information in data becomes less accessible by the model.

This paper investigates the impact of encoding gate set positioning on the trainability of a QML model and proposes an encoding pattern for high dimensional inputs, see Fig. 1. The key concept behind our method is incremental uploading.

We evaluate our approach on an image classification task (i.e., MNIST [15] and Fashion MNIST [16]) as image inputs are among the most common high dimensional inputs with various practical significance. While classifying MNIST with QML is not new, i.e., early work uses heavily downsampled 4×4 images for binary classification [17] or quantum techniques for dimensionality reduction and classification [18], a truncation of the dataset or extreme reduction in dimensions using a dimensionality reduction will only work for simple classification tasks which tolerate these information losses. One other approach which handles high dimensional image data effectively is the quantum convolutional neural network [19]. However, this is not a pure quantum approach, and the input image is broken into smaller pieces and fed into the circuit sequentially, resulting in a longer runtime, however, with the prospect for parallelization. To show how to handle more complex classification tasks, we use a high dimensional representation of the full MNIST dataset.

II. METHOD

Quantum gates in a VQC can be grouped into three categories: encoding layers, decoding layers and variational layers. Due to the limitation in circuit depth and to avoid possible barren plateau effects [20], the QC has to be designed in such a way that it allows for maximum classification performance and expressivity for a given gate count and circuit depth. To this end, we studied the effect of encoding layer positioning on the performance of a QML model by decomposing them into smaller encoding layers and progressively increasing the number of variational parameters between the layers in the VQC. Also, we would like to introduce the nomenclature used

for grouping the encoding gates throughout this paper. From here on, the collection of all encoding gates is to be called an encoding block and the encoding block split in to a smaller group of encoding gates are to be called as encoding layer.

An obvious decomposition of encoding block for image data splits and groups the gates used to encode raw features from each row of the input image. These row-wise grouped encoding gates (hereafter referred to as *encoding layers*) can per design choice be freely moved across the VQC, though each move results in a different architecture with an impact on the performance of the model. Hence, we designed five different encoding block split patterns with incremental number parameterized gates between them. These circuits are as follows: IDU_1, here, there is no variational layer between the encoding layers. All encoding gates are placed at the start followed by all variational layers. IDU_2, IDU_4, IDU_8, IDU_10 represent the circuits where the encoding block is split into 2, 4, 8, and 10 parts respectively. Between each split, there is a variational layer and the remaining variational layers are appended at the end. The number of variational layers between any two encoding layers is restricted to one as we wanted to analyze the performance boost attained by introducing a minimal and constant number of variational layers in between them. This restriction is only a design choice and other design choices are of course possible. The overall working of this proposed encoding pattern is shown in Fig. 2. We call this encoding pattern incremental data-uploading (IDU). The performance of this pattern is compared against the data re-uploading (DRU) encoding pattern [21] which is deemed to be the state-of-the-art following the evaluation metrics of Skolik et al. [22] and theoretical support from Schuld et al. [12]. However, to have a fair comparison, the number of parameters of the QML model is kept constant. Hence in the DRU architecture, the entire image information is encoded into the circuit followed by a variational layer and this is repeated until the number of variational parameters matches the number of parameters used in the incremental data-uploading experiment.

III. EVALUATION

A. Data Statistics

Our experiments have been conducted using the MNIST [15] and Fashion-MNIST [16] datasets. MNIST is a handwritten digit dataset consisting of 70,000 images representing the digits 0-9 with a size of 28×28 pixels each. Each digit class contains roughly between 5,400 - 6,750 images. The dataset of 70,000 grayscale images has been randomly grouped into 48,000 images for training, 12,000 images for validation, and 10,000 images for evaluation. As our quantum experiments use Tensorflow Quantum and the Cirq simulator on classical hardware, deeper and larger circuits become computationally intractable. To reduce computation time we reduce the size of the images from 28×28 to 10×10 using a bilinear filter so that the encoding gate count required to encode the data is small. Similarly, Fashion MNIST is also made of 70,000 grayscale images of size 28×28 representing 10 categories of clothes. As

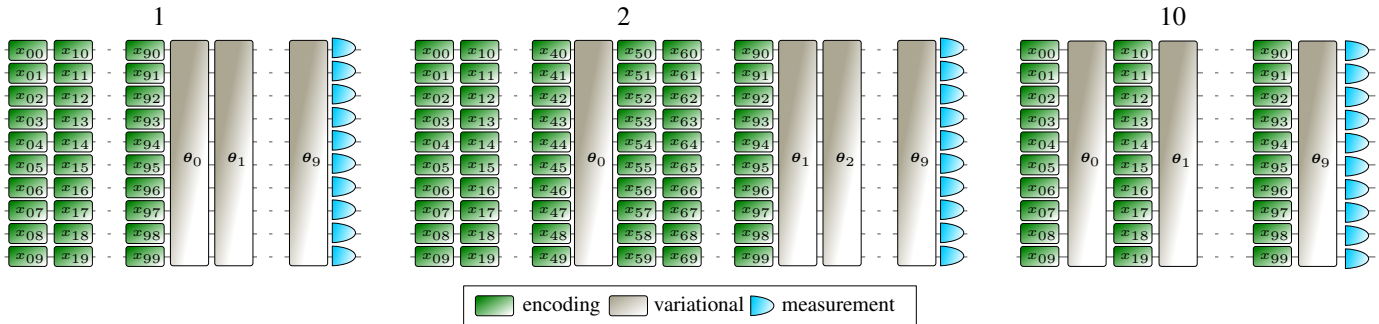


Figure 2: Different splits of interleaved layers. For the 1-split, there are 10 encoding layers followed by 10 variational layers, for the 2-split, there are 5 encoding layers, followed by 1 variational layer, then 5 encoding layers, followed by the remaining 9 variational layers. The proposed 10-split is interleaving encoding and variational layers.

for MNIST we reduce the image size and hence the overall computational time.

B. Quantum Encoding, Variational and Decoding Layer

All our datasets consist of 10 classes. Hence, we designed a ten qubit quantum circuit along with the softmax function to learn a mapping function $f(\cdot) : X \rightarrow R^o$, where X is the dataset with n data points and R^o is the probability of a data point belonging to each class. The quantum circuit consists of multiple encoding and variational layers as explained below.

The process of embedding a classical data point $x \in X$ into a quantum circuit is commonly known as data encoding, sometimes also referred to as data uploading [21]. In practice, one of the most common ways to encode a data point into a quantum circuit is via a state preparation circuit acting on state $|0\rangle^{\otimes n}$ in computational basis [9]. A state preparation circuit often consists of single qubit rotational gates matching the dimension of x with or without entangling gates so that each raw feature of the x can be scaled between $[0, \pi]$ or $[0, 2\pi]$ and used as the rotational angles for one gate in the state preparation circuit [9]. We choose a total of 100 single qubit rotational gates R_x that match the feature dimension of each data point x as the encoding layer(s) for all our experiments. The R_x gates are split into groups of ten where each group acts on one qubit. Each pixel value in x , ranging between $[0, 255]$ is scaled to $[0, \pi]$ and are fed as the rotational angle for the encoding gates.

The variational layers hold the learnable parameters that are optimized using gradient descent to approximate the mapping function f . In quantum circuits, the variational layers are again realized using single-qubit rotational and multi-qubit entangling gates where the rotational angles of the rotational gates act as the learnable parameters. Our variational layers consist of single-qubit R_y and R_z rotational gates with nearest neighbour controlled- R_z entanglements. The complete quantum circuit with encoding and variational layers is shown in Fig. 3. The circuit is measured in the computational basis, and the expectation values of the individual qubit along with the softmax function are used for class prediction.

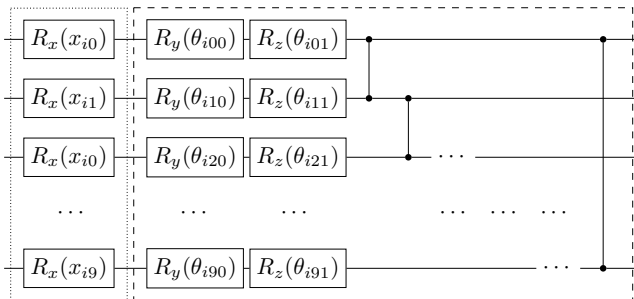


Figure 3: Single encoding block followed by a variational layer. The encoding block and the variational block are repeated 10 times with different intervals in between to form different architectures.

C. Incremental Data Uploading

We implement our interval uploading method with each of the architectures representing a different split in the encoding layer. We train each architecture on all datasets with a learning rate of 0.001 using the ADAM optimizer [23] for 25 epochs.

From Fig. 4, we can infer that there is a direct correlation between the split in the encoding layers and the performance of the model. The more the number of variational layers between the encoding blocks, the better the approximation of the mapping function and accurate the classification in both the training and testing phase. The architecture with ten encoding blocks yields the highest accuracy of around 60%.¹ We have observed the same pattern on Fashion MNIST.

To validate the argument that the interval uploading method is not data-dependent and is expected to work on arbitrary classification tasks, we shuffled every pixel value within each image in the MNIST dataset with a fixed permutation chosen randomly. The models trained on this shuffled MNIST dataset also displayed the same pattern as in the other datasets. The test accuracy of these models is shown in Table I.

¹We acknowledge that the classification accuracy is not competitive for a simple dataset such as MNIST. However, the goal of our work is not the accurate classification of the MNIST dataset but to study the impact of the encoding pattern on the trainability of the model by comparing the relative change in classification accuracy without increasing the number of parameters. Hence, we did not optimize the model for an increased classification accuracy.

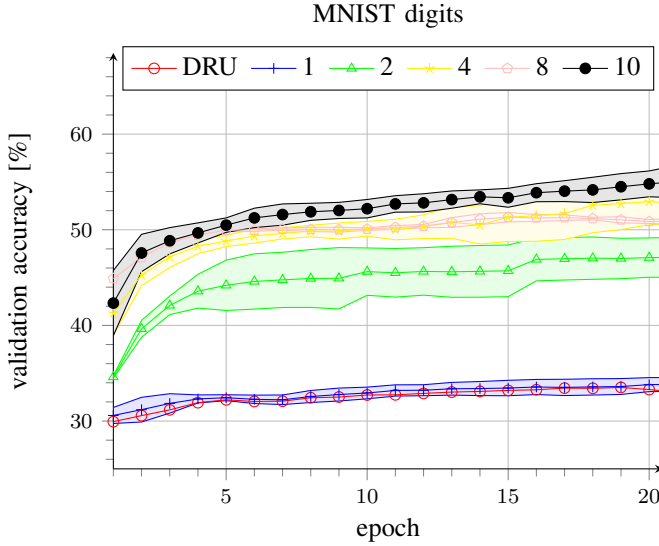


Figure 4: Average validation accuracy and its standard deviation over 5 training runs for quantum circuits with different splits in the encoding layer. DRU stands for the Data Re-uploading method, and the numbers 1, 2, 4, 8, 10 represents the quantum circuits with 1 whole encoding layer, encoding layer split into 2, 4, 8, 10 blocks respectively.

Table I: Interval uploading performance on test sets.

Dataset	DRU	IDU				
		1	2	4	8	10
MNIST	33.2±0.01	34.0±0.01	47.3±0.03	50.9±0.01	51.5±0.00	56.7±0.02
shuffled	32.2±0.00	47.1±0.01	52.2±0.01	53.8±0.01	56.1±0.01	58.6±0.01
Fashion	43.5±0.17	43.8±0.01	48.3±0.01	52.5±0.01	53.6±0.03	56.9±0.03

D. IDU in a "Deeper" circuit

From Fig. 4 and Table I it becomes clear that the quantum architecture with data re-uploading type encoding exhibits similar or lower performance than the least performing IDU architecture. Intuitively, this is an expected result as the data used for the data re-uploading architecture is a reduced dataset where each image is summed over its columns. This summation results in information loss, hence the loss in performance by the model. However, the architecture with a single encoding layer performs the same summation over the image columns (as only R_x gates are used for encoding) and performs slightly better than the DRU architecture. The poor performance of the model with DRU encoding can be correlated to the low expressive power of the variational layers. A explanation of this hypothesis is given in Section IV-B.

To further validate this hypothesis, we increased the number of parameters in the variational layer from 20 to 60 to increase the trainability of the model, see Fig. 5. This in turn increased the performance of the DRU architecture. The accuracy of DRU with a higher number of parameters is better than the architecture with a single encoding layer for the same number of parameters. However, the DRU still demonstrate a significantly low performance compared to all IDU architec-

Table II: Performance on a "deeper" architecture.

Dataset	DRU	IDU				
		1	2	4	8	10
MNIST	42.7±0.00	41.5±0.01	54.0±0.01	57.9±0.02	62.6±0.01	63.9±0.01

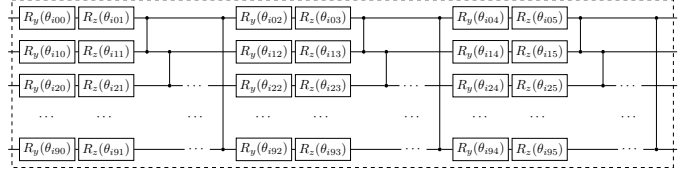


Figure 5: Variational layer with 60 parameters used in deeper models.

tures with split greater than two. Table II shows the results of different architectures with a higher number of parameters.

E. IDU for Advanced Encoding schemes

Using single-qubit R_x type encoding gates sequentially results in partial or complete summation of the input image data along the column resulting in some information loss. To further validate the effect of IDU-type encoding pattern without the influence of the summation effect, we tested two other encoding methods: 1) a R_x - R_y encoding, where we used a sequence of alternating R_x and R_y rotational gates for each row of the image instead of just R_x gates, see Fig. 6, and 2) a R_x - CR_z - R_y encoding, which is similar to R_x - R_y encoding but that uses a CR_z gate in between R_x and R_y gates, see Fig. 6.

The classification results of MNIST dataset using these two encoding methods are given in Table III. We see that the effect of incremental data-uploading type encoding pattern is more general and not restricted to single-qubit R_x type encoding. Please note that the encoding methods are simple design choices where the encoding gates do not commute. These encoding methods are not optimized toward the MNIST dataset as the intent behind the experiment was to study the effect

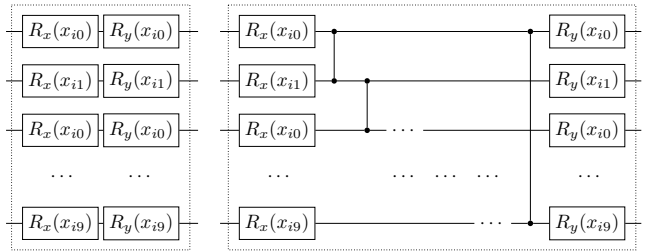


Figure 6: Single encoding block for R_x - R_y type encoding (left) and R_x - CR_z - R_y type encoding (right).

Table III: Incremental data-uploading performance on advanced encoding methods.

Dataset	DRU	IDU				
		1	2	4	8	10
R_x - R_y	0.23	0.29	0.34	0.37	0.43	0.45
R_x - CR_z - R_y	0.19	0.34	0.34	0.36	0.37	0.41

of incremental data-uploading on different encoding methods and not the effect of encoding method on MNIST dataset in itself.

IV. THEORETICAL CONSIDERATIONS

A. Quantification of Trainability and Expressibility

Two important properties of an ML model are its expressibility and trainability. Abbas et al. [11] generalizes tools for quantitative analysis to the quantum realm. Both concepts are based on the Fisher information matrix (FIM) [24] associated with the statistical model [25] $p_\theta(x, y)$ implemented by the VQCs. In practice, we use the empirical FIM defined as

$$\tilde{F}_k(\theta) = \frac{1}{k} \sum_{j=1}^k \frac{\partial}{\partial \theta} \ln p_\theta(x^{(j)}, y^{(j)}) \frac{\partial}{\partial \theta} \ln p_\theta(x^{(j)}, y^{(j)})^t. \quad (1)$$

Here, $(x^{(j)}, y^{(j)})_{j=1}^k$ are i.i.d. drawn from the joint distribution $p_\theta(x, y) = p_\theta(y | x)p(x)$. For the MNIST dataset one has inputs $x \in \mathbb{R}^{10 \times 10}$ and labels $y \in \{0, \dots, 9\}$. However, the following consideration generalize to data of any finite dimensionality.

The FIM captures the geometry of the parameter space, which has a crucial influence on the trainability of a model. To assess this, the spectrum of the positive semidefinite matrix, i.e., the distribution of its eigenvalues, is considered. A degenerate spectrum indicates a distorted parameter space, which is disadvantageous for any gradient-based optimization technique. Furthermore, an increasing accumulation of eigenvalues around zero for growing model size (i.e., qubit number) indicates the presence of barren plateaus [11].

The effective dimension [26] is a tool to capture the expressibility or capacity of a ML model. It is based upon the (empirical) FIM, and therefore can be estimated relatively straightforward by sampling. The effective dimension of a statistical model \mathcal{M}_Θ is defined as

$$ed_n(\mathcal{M}_\Theta) := 2 \frac{\ln \left(\frac{1}{V_\Theta} \int_{\Theta} \sqrt{\det \left(I_d + c_n \hat{F}(\theta) \right)} d\theta \right)}{\ln(c_n)}, \quad (2)$$

where $d = |\theta|$ is the number of parameters, $V_\Theta := \int_{\Theta} d\theta$ is the volume of the parameter space, and $\hat{F}(\theta) \in \mathbb{R}^{d \times d}$ is a normalized version of the (empirical) FIM. The parameter n captures the effective resolution of the parameter space (i.e. is related to the data availability). It enters the definition in the normalization factor $c_n = \frac{n}{2\pi \ln n}$. Under certain conditions the effective dimension provides an upper bound to the generalization error [11]. In more plain words, the measure quantifies the range of different functions, that a given model can approximate. In order to compare different models, a normalized version of the effective dimension is preferable. A division by d restricts the measure to the range $[0, 1]$, where higher values indicate a more expressible model.

The empirical Fisher information matrix was estimated using 200 random samples x^k from the MNIST data set and 100 random parameter sets θ of 200 parameters each drawn from a uniform distribution with range $[0, \pi]$. The eigenvalue spectra

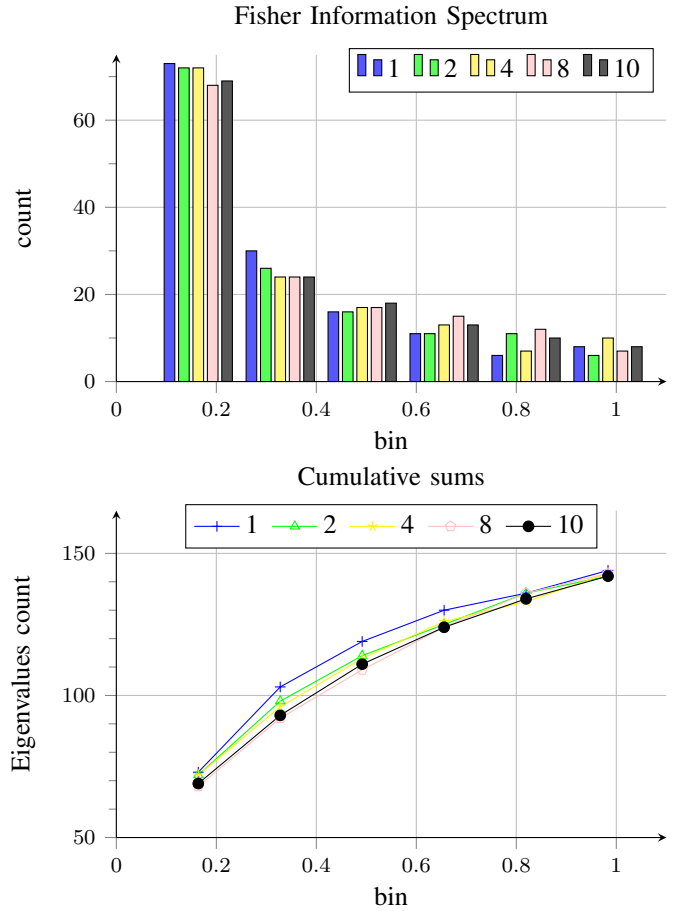


Figure 7: Fisher information

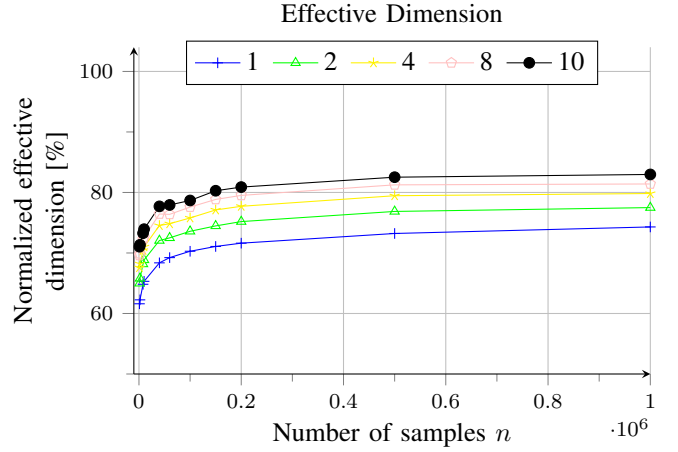


Figure 8: Effective dimension for different IDU architectures.

of the FIMs over 4000 samples for different IDU architectures are shown in Fig. 7. For reasons of presentation the histograms cut values larger than one, which anyhow do not change the overall picture. The normalized effective dimension for different IDU architectures for sample sizes ranging from 10^3 to 10^6 is shown in Fig. 8. Results presented in Fig. 7 depict

that the eigenvalue spectrum becomes more uniform for IDU architectures with a higher number of splits. This normalizing effect becomes more obvious when considering the cumulative sum plot shown in Fig. 7, bottom.

Although the peculiarity of the normalization effect is quite small in the considered instances, it indicates an improvement in trainability when employing the proposed approach. As the spectrum is more uniform with fewer eigenvalues close to zero, the parameter space is less distorted, which is beneficial for optimization methods. The difference in terms of the effective dimension is more distinct, i.e. it clearly increases when using more IDU layers. This indicates an increase in model expressibility, while the number of parameters stays the same. In all instances the normalized effective dimension grows with larger resolution of the parameter space, which is a reasonable behaviour for machine learning models.

B. Frequency spectrum

To gain more insight into the performance differences observed in the previous sections, in the following we investigate the function class represented by the different architectures. Slightly generalizing the setting, we consider $x \in \mathbb{R}^{N \times M}$ in the following and denote the j th rows of the matrix by the column vector \mathbf{x}_j , that is $(\mathbf{x}_j)_k = x_{jk}$ for $k = 0, \dots, M - 1$ and $j = 0, \dots, N - 1$. The vectors \mathbf{x}_j therefore correspond to the data fed in the j th encoding layer in Fig. 1. It was shown by Schuld et al. in Ref. [12] that the functions f_θ represented by VQCs are Fourier sums when each encoding layer is given by single-qubit rotations about a given axis for each qubit. In particular, the variational layers determine the amplitudes and the frequency spectrum is fixed by the data-encoding layers. Following Ref. [12], we find

$$f_\theta(x) = \sum_{\omega_0, \dots, \omega_{N-1} \in \Omega} c_\omega(\theta) \exp \left\{ i \sum_{j=0}^{N-1} \omega_j \mathbf{x}_j \right\}, \quad (3)$$

where $\Omega = \{-1, 0, 1\}^M$ is the frequency spectrum and ω the matrix containing ω_j as j th row. Since f_θ is real valued, we find $c_{-\omega} = c_\omega^*$. More intuitively, the functions represent NM dimensional Fourier sums with frequencies ± 1 and 0 . Note that the coefficients $c_\omega(\theta)$ are only independent and can be chosen freely if the variational layers are universal, i.e. can represent any unitary matrix. In practice, the expressivity of the circuit might be severely limited by the number of variational parameters, indeed a general n -qubit unitary requires exponentially many parameters in the numbers of qubits. Nevertheless, equation (Eq. (3)) qualitatively explains the behaviour shown in Table I where an increase of performance with increasing number of interleaved variational layers is observed. Since the data encoding is based on R_x rotations only, the setup in the left subfigure of Fig. 2 is equivalent to summing the vectors and feeding the result into the circuit by only one encoding layer. As a result, the input dimension in equation (Eq. (3)) decreases from NM to M so that the model loses access to much of the information present in the data x , explaining the poor

performance in the left column of Table I. As the number of variational layers increase, f_θ gains access to more information as only some of the rows in the data are summed, finally reaching optimal performance for the fully interleaved setup shown in the right subfigure of Fig. 2. It is worthwhile noting that decreasing the number of interleaved layers while keeping the number of variational layers constant, increases the expressibility of the final variational layers but it seems conceivable that this increase cannot compensate for the information loss by partially or fully summing the rows in the data x . The same argument applies to the interpretation of the DRU column in Table I and Table II. Here, in the data re-uploading setting the rows of the image are first summed and then repeatedly encoded into the circuit with variational layers in between. While the frequency spectrum of the Fourier sum now contains all integers between $-N$ and N [12], as can be seen from equation (Eq. (3)) by replacing \mathbf{x}_j by the sum over the rows for all j , again the model seems unable to compensate for the information which is lost in summing the rows of the image. In case of more general encoding schemes such as alternating R_x - R_y gates for subsequent encoding layers, intuitively, this effect is less dramatic due to the non-commutativity of the encoding gates. However, Table III indicates that also in case of non-commuting encoding gates the information in the data is much better accessible by the model in which the more variational layers are interleaved with encoding layers. A more rigorous discussion of this situation will be at the focus of further work.

V. CONCLUSION

This paper proposes an encoding pattern called incremental data-uploading for high dimensional data. It acts as a guideline for positioning the encoding layers in a variational quantum circuit. Here, the encoding and variational layers alternate one after the other so that the data is fed incrementally into the circuit and becomes more accessible to the model. IDU with maximum variational layers in between them showed a performance boost of 15 - 25 percentage points in image classification tasks. The effective dimension and Fisher information results also support our claim that the IDU pattern increases the trainability and expressivity of the QML model without increasing the number of its parameters. In addition, expressing the quantum model as a partial Fourier sum, we were able to connect its performance to the range of accessible frequencies.

Our experiments also showed that an encoding pattern like data re-uploading exhibits low accuracy when dealing with high dimensional data and fewer parameters in a QML model even though it increases the overall QC depth and number of quantum gates. Hence, we conclude that the presented data encoding pattern shows an improvement in the performance of a QML model with high dimensional data, shallow circuit depth and a given encoding method. However, finding an optimal encoding method in the IDU framework for a given dataset is left for future work.

REFERENCES

- [1] “Ibm quantum,” <https://quantum-computing.ibm.com/>, 2022.
- [2] F. Arute, K. Arya, R. Babbush, D. Bacon, J. C. Bardin, R. Barends, R. Biswas, S. Boixo, F. G. S. L. Brandao, D. A. Buell, B. Burkett, Y. Chen, Z. Chen, B. Chiaro, R. Collins, W. Courtney, A. Dunsworth, E. Farhi, B. Foxen, A. Fowler, C. Gidney, M. Giustina, R. Graff, K. Guerín, S. Habegger, M. P. Harrigan, M. J. Hartmann, A. Ho, M. Hoffmann, T. Huang, T. S. Humble, S. V. Isakov, E. Jeffrey, Z. Jiang, D. Kafri, K. Kechedzhi, J. Kelly, P. V. Klimov, S. Knysh, A. Korotkov, F. Kostrița, D. Landhuis, M. Lindmark, E. Lucero, D. Lyakh, S. Mandrà, J. R. McClean, M. McEwen, A. Megrant, X. Mi, K. Michielsen, M. Mohseni, J. Mutus, O. Naaman, M. Neeley, C. Neill, M. Y. Niu, E. Ostby, A. Petukhov, J. C. Platt, C. Quintana, E. G. Rieffel, P. Roushan, N. C. Rubin, D. Sank, K. J. Satzinger, V. Smelyanskiy, K. J. Sung, M. D. Trevithick, A. Vainsencher, B. Villalonga, T. White, Z. J. Yao, P. Yeh, A. Zalcman, H. Neven, and J. M. Martinis, “Quantum supremacy using a programmable superconducting processor,” *Nature*, vol. 574, no. 7779, pp. 505–510, Oct 2019.
- [3] F. Bova, A. Goldfarb, and R. G. Melko, “Commercial applications of quantum computing,” *EPJ Quantum Technology*, vol. 8, no. 1, p. 2, Jan 2021.
- [4] C. Moussa, H. Calandra, and V. Dunjko, “To quantum or not to quantum: Towards algorithm selection in near-term quantum optimization,” *Quantum Science and Technology*, vol. 5, p. 044009, 10 2020.
- [5] N. Mishra, M. Kapil, H. Rakesh, A. Anand, N. Mishra, A. Warke, S. Sarkar, S. Dutta, S. Gupta, A. Prasad Dash, R. Gharat, Y. Chatterjee, S. Roy, S. Raj, V. Kumar Jain, S. Bagaria, S. Chaudhary, V. Singh, R. Maji, P. Dalei, B. K. Behera, S. Mukhopadhyay, and P. K. Panigrahi, “Quantum machine learning: A review and current status,” in *Data Management, Analytics and Innovation*, N. Sharma, A. Chakrabarti, V. E. Balas, and J. Martinovic, Eds., vol. 1175. Singapore: Springer Singapore, 2021, pp. 101–145.
- [6] H. Watanabe, R. Raymond, Y.-Y. Ohnishi, E. Kaminishi, and M. Sugawara, “Optimizing parameterized quantum circuits with free-axis selection,” 10 2021, pp. 100–111.
- [7] H. Yano, Y. Suzuki, K. M. Itoh, R. Raymond, and N. Yamamoto, “Efficient discrete feature encoding for variational quantum classifier,” *IEEE Trans. Quantum Engineering*, vol. 2, 2021.
- [8] M. C. Caro, E. Gil-Fuster, J. J. Meyer, J. Eisert, and R. Sweke, “Encoding-dependent generalization bounds for parametrized quantum circuits,” *Quantum*, vol. 5, p. 582, Nov. 2021.
- [9] R. LaRose and B. Coyle, “Robust data encodings for quantum classifiers,” *Phys. Rev. A*, vol. 102, p. 032420, Sep 2020.
- [10] L. Banchi, J. Pereira, and S. Pirandola, “Generalization in quantum machine learning: A quantum information standpoint,” *PRX Quantum*, vol. 2, 11 2021.
- [11] A. Abbas, D. Sutter, C. Zoufal, A. Lucchi, A. Figalli, and S. Woerner, “The power of quantum neural networks,” *Nature Computational Science*, vol. 1, no. 6, pp. 403–409, jun 2021.
- [12] M. Schuld, R. Sweke, and J. J. Meyer, “Effect of data encoding on the expressive power of variational quantum-machine-learning models,” *Phys. Rev. A*, vol. 103, p. 032430, 2021.
- [13] M. C. Caro, E. Gil-Fuster, J. Meyer, J. Eisert, and R. Sweke, “Encoding-dependent generalization bounds for parametrized quantum circuits,” vol. 5, p. 582, 2021.
- [14] M. Franz, L. Wolf, M. Periyasamy, C. Ufrecht, D. D. Scherer, A. Plinge, C. Mutschler, and W. Mauerer, “Uncovering instabilities in variational-quantum Deep Q-Networks,” *arXiv preprint arXiv:2202.05195*, 2022.
- [15] Y. LeCun, C. Cortes, and C. J. Burges, “The MNIST database of handwritten digits,” 1998. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [16] H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms,” *arXiv preprint arXiv:1708.07747*, 2017.
- [17] E. Farhi and H. Neven, “Classification with quantum neural networks on near term processors,” *arXiv preprint arXiv:1802.06002*, 2018.
- [18] I. Kerenidis and A. Luongo, “Classification of the MNIST data set with quantum slow feature analysis,” *Physical Review A*, vol. 101, no. 6, Jun. 2020.
- [19] A. Matic, M. Monnet, J. M. Lorenz, B. Schachtner, and T. Messerer, “Quantum-classical convolutional neural networks in radiological image classification,” 2022. [Online]. Available: <https://arxiv.org/abs/2204.12390>
- [20] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Babbush, and H. Neven, “Barren plateaus in quantum neural network training landscapes,” *Nature Communications*, vol. 9, no. 1, p. 4812, 2018.
- [21] A. Pérez-Salinas, A. Cervera Lierta, E. Gil-Fuster, and J. Latorre, “Data re-uploading for a universal quantum classifier,” *Quantum*, vol. 4, p. 226, 02 2020.
- [22] A. Skolik, S. Jerbi, and V. Dunjko, “Quantum agents in the gym: a variational quantum algorithm for deep q-learning,” *arXiv preprint arXiv:2103.15084*, 2021.
- [23] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *International Conference on Learning Representations*, 12 2014.
- [24] M. Thomas and A. T. Joy, *Elements of information theory*. Wiley-Interscience, 2006.
- [25] J. Rissanen, “Fisher information and stochastic complexity,” *IEEE Transactions on Information Theory*, vol. 42, no. 1, pp. 40–47, 1996.
- [26] O. Berezniuk, A. Figalli, R. Ghigliazza, and K. M. M. Musaelian, “A scale-dependent notion of effective dimension,” *arXiv preprint arXiv:2001.10872*, 2020.